



CA' FOSCARI UNIVERSITY OF VENICE

AND

MASARYK UNIVERSITY

DOCTOR OF PHILOSOPHY PROGRAMME  
IN COMPUTER SCIENCE

CYCLE XXXIII

FINAL THESIS

**String Analysis  
for  
Software Verification**

SSD: INF/01

**PROGRAMME COORDINATOR**

Prof. Agostino CORTESI

**SUPERVISOR (CA' FOSCARI UNIVERSITY)**

Prof. Agostino CORTESI

**SUPERVISOR (MASARYK UNIVERSITY)**

Prof. Vashek MATYAS

**GRADUATE STUDENT**

Martina OLLIARO

Matriculation Number 834397

UČO 47767



*To the human revolution*



# *Abstract*

This thesis aims to investigate string manipulation with security implications in different programming languages and to improve the state-of-the-art by applying the Abstract Interpretation theory to string analysis. Erroneous string manipulation is a challenging problem in software verification. In fact, it is one of the major causes of program vulnerabilities that can be exploited by malicious users, leading to severe consequences for the affected systems. By string analysis, we mean statically computing the set of string values that are possibly assigned to a variable. As for other analysis issues, this is undecidable. Thus a certain degree of approximation is necessary to find evidence of bugs and vulnerabilities in string manipulating code. We take advantage of the Abstract Interpretation theory, *id est*, a powerful mathematical theory that enables us to define and prove the soundness of approximations. The five main contributions of this thesis are:

*Abstracting shape and content of strings.* We introduce a new sophisticated string abstract domain for the C language. The domain (called M-String) is parametrized on an index (bound) domain and a character domain. Picking different constituent domains, i.e., both shape information on the array structure and value information on the contained characters, allows M-String to be tailored for specific verification tasks (e.g., detection of buffer overflows), balancing precision against complexity. We describe the concrete and the abstract semantics of basic string operations and prove their soundness formally. In addition to a selection of string functions from the standard C library, we provide semantics for character access and update, enabling automatic lifting of arbitrary string-manipulating code into the domain. Furthermore, we provide an executable implementation of abstract operations. Using a tool that automatically lifts existing programs into the M-String domain along with an explicit-state model checker, we evaluate the accuracy of the proposed domain experimentally on real-case test programs.

*Combining string domains.* We combine abstract domains resulting from the reduced product between string shape abstraction and string content abstraction, in order to better detect inconsistent states leading to program errors without a major impact on efficiency. In particular, the combinations involve some string abstract domains introduced in the literature with the segmentation domain that we instantiate for string analysis.

*Completeness of string domains.* In Abstract Interpretation, completeness ensures that the analysis does not lose information with respect to the property of interest. We provide a systematic and constructive approach for generating the completion of string domains for dynamic languages, and we apply it to the refinement of existing string abstractions. Indeed, for dynamic languages, lack of completeness is a key security issue, as poorly managed string manipulation code may easily lead to significant security flaws. We also provide an effective procedure to measure the precision improvement obtained when lifting the analysis to complete domains.

*Relational abstract domains for string analysis.* Almost all the existing string abstract domains track information of single variables in a program (e.g., if a string contains a specific character) without inspecting their relationship with other values, causing loss of relevant knowledge. Thus, we introduce a generic framework that allows formalizing relational string abstract domains based on ordering relationships. We instantiate this framework to several domains built upon different well-known string orders (e.g., substring relationships). We implemented the domain based on substring ordering, and we provide an experimental evaluation of its effectiveness on some case studies.

*String manipulation in watermarking scenarios.* We manipulate string values in the context of relational database watermarking. We propose a semantic-driven watermarking approach of relational textual databases, which marks multi-word textual attributes, exploiting the synonym substitution technique for text watermarking together with notions in semantic similarity analysis, and dealing with the semantic perturbations provoked by the watermark embedding. We show the effectiveness of our approach through an experimental evaluation, highlighting the resulting capacity, robustness, and imperceptibility watermarking requirements. We also prove the resilience of our approach with respect to the random synonym substitution attack.

# *Acknowledgements*

Doctoral studies represent a crucial point in my personal and educational growth. Among all the people who walked with me during this journey, four are of particular importance.

*To my Mother.* Since the moment she gave me life, she never abandoned me, fully accomplishing her duties. She is kindness and resilience made person. Her support pushed me towards places and thoughts which I would not have had the courage to explore. I owe her everything.

*To my Sisters.* Diletta and Carola are the light of my days. They relieve my bad moods and take care of me, making me always feel safe and protect. I am never alone with them.

*To my Mentor.* Tino believed in me. There is nothing more precious than finding someone who guides you and teaches you. I am incredibly grateful for the time he dedicated to me.

Thanks to all my beloved, in particular to my colleagues and friends Gianluca, Vincenzo, and Maikel, and to my best friend Elena.

*Martina Olliaro*

*Treviso, 10/02/2021*





# CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Methodology . . . . .	3
1.2.1	Static Analysis . . . . .	4
1.2.2	Abstract Interpretation . . . . .	4
1.3	Contribution . . . . .	6
1.4	Thesis Structure . . . . .	7
<b>2</b>	<b>Preliminaries</b>	<b>9</b>
2.1	Sets and Sequences . . . . .	9
2.2	Order Relations . . . . .	10
2.3	Functions . . . . .	12
2.4	Fixpoints . . . . .	12
2.5	Traces . . . . .	13
2.6	Abstract Interpretation . . . . .	14
2.6.1	Galois Connection . . . . .	14
2.6.2	Soundness and Completeness . . . . .	16
2.6.3	Fixpoints Approximation . . . . .	18
2.6.4	Product Operators . . . . .	18
<b>3</b>	<b>String Analysis for C</b>	<b>21</b>
3.1	Introduction . . . . .	21
3.2	FunArray . . . . .	23

3.2.1	Array Concrete Representation . . . . .	23
3.2.2	Array Abstract Domain Functor . . . . .	25
3.3	Syntax . . . . .	26
3.4	Concrete Domain and Semantics . . . . .	27
3.4.1	Character Array Concrete Representation . . . . .	28
3.4.2	Concrete Domain . . . . .	29
3.4.3	Concrete Semantics . . . . .	29
3.5	M-String . . . . .	35
3.5.1	Character Array Abstract Domain Functor . . . . .	35
3.5.2	Abstract Semantics . . . . .	41
3.5.3	Soundness . . . . .	45
3.6	Program Abstraction . . . . .	48
3.6.1	Compilation-Based Approach . . . . .	49
3.6.2	Syntactic Abstraction . . . . .	49
3.6.3	Aggregate Domains . . . . .	51
3.6.4	Semantic Abstraction . . . . .	53
3.6.5	Abstract Operations . . . . .	53
3.7	Instantiating M-String . . . . .	54
3.7.1	Symbolic Scalar Values . . . . .	55
3.7.2	Concrete Characters, Symbolic Bounds . . . . .	55
3.7.3	Symbolic Characters, Symbolic Bounds . . . . .	57
3.7.4	Implementation . . . . .	57
3.8	Experimental Evaluation . . . . .	58
3.9	Discussion . . . . .	62
<b>4</b>	<b>Combining String Domains</b>	<b>65</b>
4.1	Introduction . . . . .	65
4.2	Syntax . . . . .	68
4.3	Concrete Domain and Semantics . . . . .	69
4.3.1	Concrete Domain . . . . .	69
4.3.2	Concrete Semantics . . . . .	69
4.4	String Abstract Domains . . . . .	71

4.4.1	String Length . . . . .	71
4.4.2	Character Inclusion . . . . .	72
4.4.3	Prefix and Suffix . . . . .	73
4.5	Segmentation Abstract Domain . . . . .	74
4.5.1	String Concrete Representation . . . . .	74
4.5.2	Abstract Domain . . . . .	75
4.5.3	Abstract Semantics . . . . .	81
4.5.4	Soundness . . . . .	83
4.6	Refined String Abstract Domains . . . . .	84
4.6.1	Meaning of Refinement . . . . .	85
4.6.2	Combining Segmentation and String Length Domains . . . . .	85
4.6.3	Combining Segmentation and Character Inclusion Domains . . . . .	89
4.6.4	Combining Segmentation and Prefix Domains . . . . .	94
4.7	Discussion . . . . .	98
<b>5</b>	<b>Completeness of String Domains</b>	<b>101</b>
5.1	Introduction . . . . .	101
5.2	Making Abstract Interpretation Complete . . . . .	102
5.2.1	Complete Shell vs Complete Core . . . . .	103
5.2.2	Domain Completion Procedure . . . . .	105
5.2.3	Motivating Example . . . . .	106
5.3	Core Language . . . . .	108
5.3.1	Syntax . . . . .	108
5.3.2	Concrete Semantics . . . . .	109
5.4	Making JavaScript String Abstract Domains Complete . . . . .	109
5.4.1	Completing TAJIS String Abstract Domain . . . . .	110
5.4.2	Completing SAFE String Abstract Domain . . . . .	113
5.5	Benefits of Adopting Complete String Abstractions . . . . .	117
5.5.1	Precision . . . . .	117
5.5.2	Qualitative Evaluation of Complete Shells . . . . .	118
5.5.3	False Positives Reduction . . . . .	119
5.6	Relative Precision . . . . .	120

5.6.1	Abstract Domains Precision: an Overview . . . . .	120
5.6.2	Measuring Precision Gained by Complete Shells . . . . .	121
5.6.3	Experimental Evaluation . . . . .	124
5.7	Discussion . . . . .	125
<b>6</b>	<b>Relational String Abstract Domains</b>	<b>127</b>
6.1	Introduction . . . . .	127
6.2	Core Language . . . . .	129
6.2.1	Syntax . . . . .	129
6.2.2	Concrete Semantics . . . . .	130
6.3	A Suite of String Relational Abstract Domains . . . . .	132
6.3.1	General Relational Framework . . . . .	132
6.3.2	String Length Relational Abstract Domain . . . . .	134
6.3.3	Character Inclusion Relational Abstract Domain . . . . .	135
6.3.4	Substring Relational Abstract Domain . . . . .	136
6.3.5	Extension to String Expressions . . . . .	137
6.4	Experimental Evaluation . . . . .	140
6.4.1	Test Cases . . . . .	140
6.4.2	Evaluating a Real World Sample . . . . .	142
6.5	Discussion . . . . .	143
<b>7</b>	<b>String Manipulation in Watermarking Scenarios</b>	<b>145</b>
7.1	Introduction . . . . .	146
7.2	Motivating Examples . . . . .	149
7.3	Preliminaries . . . . .	151
7.3.1	Semantic Similarity Theory . . . . .	151
7.3.2	Text Watermarking . . . . .	151
7.4	Semantic-based Watermarking Approach . . . . .	153
7.4.1	Architecture of the Proposal . . . . .	153
7.4.2	Watermarking Procedure . . . . .	157
7.4.3	Analysis of the Watermark Capacity . . . . .	159
7.4.4	Considerations for the Adversary Model . . . . .	161
7.5	Experimental Results . . . . .	163

7.5.1	Improvement of the Watermark Capacity . . . . .	165
7.5.2	Detectability Analysis . . . . .	166
7.5.3	Watermark Imperceptibility . . . . .	169
7.5.4	Technique's Robustness . . . . .	170
7.5.5	Scalability and Complexity . . . . .	171
7.6	Discussion . . . . .	173
<b>8</b>	<b>Conclusion and Future Work</b>	<b>175</b>
<b>A</b>	<b>Unification Algorithm</b>	<b>177</b>
<b>B</b>	<b>String Abstract Domains</b>	<b>179</b>
B.1	String Length . . . . .	179
B.2	Character Inclusion . . . . .	183
B.3	Prefix . . . . .	186
<b>C</b>	<b>Relational String Abstract Domains</b>	<b>189</b>
C.1	Abstract Semantics of <code>Len</code> . . . . .	189
C.2	Abstract Semantics of <code>Char</code> . . . . .	190
C.3	Abstract Semantics of <code>Sub</code> . . . . .	191
C.4	Abstract Semantics of <code>Sub<sup>*</sup></code> . . . . .	193
C.5	<code>Len<sup>*</sup></code> Relational Abstract Domain . . . . .	195
C.6	<code>Char<sup>*</sup></code> Relational Abstract Domain . . . . .	197
	<b>Bibliography</b>	<b>199</b>



## LIST OF FIGURES

2.1	Hasse Diagram . . . . .	11
3.1	Interpretation-based vs Compilation-based . . . . .	50
3.2	Syntactic Abstraction . . . . .	50
3.3	M-String with Symbolic Bounds . . . . .	56
5.1	Coalesced Sum Abstract Domain . . . . .	107
5.2	$\mu$ Dyn Syntax . . . . .	108
5.3	$\mu$ Dyn Semantics . . . . .	109
5.4	SAFE and TAJs Domains for $\mu$ Dyn . . . . .	110
5.5	TAJs Complete Shell . . . . .	112
5.6	SAFE Concatenation Abstract Semantics . . . . .	114
5.7	SAFE Absolute Complete Shell . . . . .	116
5.8	$\mu$ Dyn Program Example . . . . .	122
6.1	SECNAME Sample . . . . .	128
6.2	IMP Syntax . . . . .	130
6.3	Control Flow Graph Generation . . . . .	130
6.4	Len Lattice Operations . . . . .	135
6.5	Char Lattice Operations . . . . .	136
6.6	Sub Lattice Operations . . . . .	137
6.7	A* Lattice Operations . . . . .	138
6.8	Sub* Lattice Operations . . . . .	139
6.9	NOTCON and REP Programs . . . . .	140

7.1	Embedding Process Architecture . . . . .	154
7.2	Synonyms Sets . . . . .	156
7.3	Evaluating Watermark Capacity . . . . .	161
7.4	WM Source Samples . . . . .	163
7.5	Analysis Correlation . . . . .	172



## LIST OF TABLES

3.1	C Syntax . . . . .	27
3.2	M-String Measurements . . . . .	59
3.3	Benchmark Abstraction . . . . .	59
3.4	PDCLib Verification Results . . . . .	60
3.5	Variabs Overflow Benchmarks Results . . . . .	61
3.6	Measurements for Automatically Generated Parsers Analysis . . . . .	62
3.6	Measurements for Automatically Generated Parsers Analysis: <i>cont</i> . . . . .	62
4.1	Introductory Example . . . . .	68
4.2	Shortcuts of String Constants . . . . .	71
4.3	Program Analysis with $\overline{\mathbf{SL}}$ . . . . .	72
4.4	Program Analysis with $\overline{\mathbf{CI}}$ . . . . .	73
4.5	Program Analysis with $\overline{\mathbf{PR}}$ . . . . .	74
5.1	Completing SAFE Domain . . . . .	115
5.2	TAJS Domain vs TAJS Complete Shell Domain . . . . .	124
6.1	Analyses Results of NOTCON . . . . .	141
6.2	Analyses Results of REP . . . . .	141
6.3	Analyses Results of SECNAME . . . . .	143
7.1	AHK Approach Notation . . . . .	147
7.2	Motivating Example . . . . .	149
7.3	Structure of the Dataset “Amazon Fine Food Reviews” . . . . .	164

7.4	Binary Capacity Values . . . . .	166
7.5	Weight-based Capacity Values . . . . .	166
7.6	WSD Precision During WM Detection . . . . .	167
7.7	Dected WM Quality . . . . .	168
7.8	WM Detectability . . . . .	168
7.9	Fixed Word Rate Values . . . . .	169
7.10	Similarities for WM UTM . . . . .	170
7.11	Detected WMs After Pseudo-random Tuple Deletion Attacks . . . . .	171
7.12	Detected WMs After Pseudo-Random Update Attacks . . . . .	171
7.13	WM Embedding Required Time . . . . .	172
B.1	$\overline{\mathbf{SL}}$ Abstract Semantics . . . . .	181
B.2	$\overline{\mathbf{CI}}$ Abstract Semantics . . . . .	185
B.3	$\overline{\mathbf{PR}}$ Abstract Semantics . . . . .	187

# Chapter 1

---

## INTRODUCTION

The goal of this thesis is to improve the state-of-the-art of string analysis in software verification. We take into account different scenarios of string manipulation, and by applying the Abstract Interpretation theory, we introduce new abstract domains for string analysis. We discuss their ability to enhance the precision when combining them with other existing strings domains. As an additional contribution, we apply a semantic approach to string manipulation in the context of watermarking for relational databases.

### 1.1 Motivation

Programming with strings is often prone to errors. Indeed the different ways strings are treated by the programming languages and the lack of support for string manipulation often lead to run-time errors. Depending on the language we are using, strings can be handled as arrays of characters, primitive data types, or objects. Due to their different implementation, strings suffer from manipulation errors causing vulnerabilities that result in irreversible damages [30].

For instance, in the `C` programming language, strings are not a built-in type. As a matter of fact, they are represented as null terminating arrays of characters. Due to how arrays are managed in `C`, it is common for those that manipulate strings to suffer from errors such as: unbounded string copies, off-by-one errors, null-termination errors, and string truncation [154]. Some of those errors lead to buffer overflows that can cause the program to behave unpredictably, crash or be subject to malicious code injection.

Injection attacks also target programs written in a memory safe language. In the `Java` programming language, strings are objects corresponding to the instances of the `String` class, and even though the `Java` language is memory safe, programmers can make logical programming errors leading to unexpected consequences. The same happens in software written in the `JavaScript` language, where strings are a primitive data type and language features like, e.g., dynamic typing, implicit type conversion, and reflection, compromise

code readability, data consistency and make the code prone to bugs or to vulnerability exploitations [12, 145].

Nevertheless, strings play a crucial role in programming. In particular, in web-based software, strings are used any time the web application needs to communicate with end-users and to store information, as well as in generating dynamic code.

For the reasons above string analysis techniques, i.e., a static analysis approach that, given a particular point of a program, aims to determine the literals in that point held by a string variable, are recognised to be an effective approach to detect string manipulation errors and prevent catastrophic events. The relevance of string analysis may be better understood by the following two security scenarios.

### **Buffer Overflow Vulnerability**

According to the IEEE and TIOBE rankings, C is currently one of the top ten programming languages in use [31, 168], and many software systems of critical importance are written in the C programming language.

A well-known security issue that affects C code is the buffer overflow vulnerability. A buffer overflow consists of access to a buffer outside its assigned bounds. It occurs because the C programming language does not provide any built-in protection to over-reading or over-writing data in any parts of the memory. For example, there is no automatic check if the insertion of data into an array (the built-in buffer type) is confined within its boundaries. Static methods that automatic detect buffer overflows in C programs have been widely studied in the literature (with an emphasis also on strings) [62, 63, 129, 161, 173, 175, 181]. Those methods exploit different inference techniques based, among others, on constraint-based techniques, tainted data-flow analysis, string pattern matching analysis or annotation analysis [156].

### **Injection**

Web applications are often affected by problems related to how strings are (improperly) manipulated [30, 91, 95, 164], leading to severe consequences, e.g., deleting content from a database or executing malicious scripts in the end user's browser. According to the OWASP Top Ten [141], two of the most common vulnerabilities in web applications that arise from string manipulation errors are: Cross Site Scripting (XSS) and SQL injection (SQLi). Both vulnerabilities are classified as injection attacks, i.e., when untrusted inputs are supplied to a program, and are caused by inadequate string input validation and sanitization, specifically:

- “A SQL Injection attack consists of insertion (or injection) of a SQL query via the input data from the client to the application” [163].

- “A Cross Site Scripting attack consists of injection of malicious scripts into otherwise benign and trusted websites” [182].

In the last decade, the number of web applications is increasing. Consequently, the amount of new XSS and SQLi vulnerabilities is advancing despite both the suggestions provided by secure coding practices and the available techniques of input validation and sanitization.

To this end, precise static analysis techniques can prevent string manipulation errors that lead to security attacks in web-based software. In the literature, several frameworks that use static analysis to counter injection attacks have been designed [71, 94, 100, 109, 124, 174, 177]. Furthermore, different string analysis techniques have been developed to tackle the problems mentioned above [114, 133, 167, 170, 176, 178, 186, 188]. For instance, Minimade [133] proposed to use a context free grammar to approximate strings to validate and guarantee the security of web pages dynamically generated by a server side program. This technique has been extended by Wasserman and Su to find SQLi [176] and XSS vulnerabilities [177]. In [188], Yu et al. applied string analysis, by means of automaton abstraction, to check the correctness of sanitization operations and later to automatically repair faulty string manipulation code [184]. Detection and verification of sanitizers have been carried out also by Tateishi et al. [167], where constraints over program variables and string operations are represented with monadic second-order logic. Finally, string analysis by means of Abstract Interpretation has been used by Tripp et al. [170] to detect JavaScript security vulnerabilities in the client side of web applications. Other survey contributions in this field can be found in [90, 95, 116, 157].

## 1.2 Methodology

As mentioned above, string analysis is a static analysis technique that, given a particular program point where a string variable occurs, aims to determine the literals possibly assigned to that variable. In general, this is an undecidable problem, i.e., no algorithm can solve it for all programs and inputs. However, approximations can be used to prove anyway the presence or the absence of software bugs and vulnerabilities in the string manipulating code.

In the recent literature, different approximation techniques for string analysis have been developed, such as [30]: automata-based [15, 34, 187, 188], abstraction-based [8, 9, 10, 11, 49], constraint-based [1, 123, 150, 152, 172], and grammar-based [133, 177]. For instance, in [15], Arceri and Mastroeni defined a new automaton-based semantics for string analysis to handle dynamic languages string features such as dynamic typing and implicit type conversion. Amadini et al., in [8], approximate strings as a dashed string, namely a sequence of concatenated blocks that specify the number of times the character they contain must appear and the number of times the latter may appear. Samirni et al.

[150] repaired HTML generation errors in PHP programs by solving a system of string constraints.

We choose to adopt the Abstract Interpretation theory [51, 56] to approach the string analysis problem as it provides a rigorous mathematical framework within which we can design sound approximations varying the degrees of abstractions.

### 1.2.1 Static Analysis

Static analysis aims to automatic reasoning about the behaviour of computer programs without executing them [136]. Firstly used in compiler optimization (e.g., to produce efficient code), it is now widely applied in software verification. In particular, the verification process is about providing guarantees that a computer program complies with its requirement (e.g., correctness) [101] and, as a matter of fact, static program analysis has been shown useful to find bugs [19, 27]. That is extremely important when code bugs may be exploited to perform security attacks.

It is important to note that from Rice's theorem [148], we know that all non trivial properties about the behaviour of programs, written in a Turing complete programming language, are undecidable. To overcome this problem, static analysis provides *sound* approximations. By sound approximation, we mean both the approximation of concrete values handled by a program by means of abstract values and an approximation of operations by corresponding operations on abstract values that preserves the semantics with respect to the observed property kept by the abstraction.

In the context of software verification, a sound verification analysis must detect all the errors, possibly producing also false positive due to the approximation.

Different static analysis techniques have been developed since the 1970s, such as: data flow analysis [113], control flow analysis [7], type and effect analysis [138], Abstract Interpretation [51], model checking [22] and symbolic execution [117]. As the years go by, for some of those techniques there have been defined several approaches. For example, the data flow analysis can be equational-based or constraint-based [139] and the model checking technique can be abstract-based [35] or symbolic-based [36].

### 1.2.2 Abstract Interpretation

Abstract Interpretation [51, 56] has been proposed by P. Cousot and R. Cousot in the 1970s as a general theory of sound abstraction of the uncomputable concrete semantics of computer programs, which is able to cover all the aforementioned static analysis techniques. Now it is widely integrated in software verification tools and used to rigorous mathematical proofs of approximations correctness. Static analysis by Abstract Interpretation is sound by construction, but in general, it is not complete, as some loss of information may occur during the abstract computation. Concrete values managed by a program are represented by elements belonging to algebraic structures called abstract domains. Abstract domains

are designed taking into account the properties to infer, e.g., numerical or lexical, relational or non-relational, and they can be both language-specific or general purpose. Moreover, Abstract Interpretations allows the combination of abstract domains, through suitable operators, i.e., refinement and product operators, to enhance the accuracy of the abstract computation.

### **Numerical Abstract Domains**

A numerical abstract domain is an abstract domain that approximates sets of numerical values [126]. Different numerical abstract domains have been designed, such as: the Sign and the Parity abstract domains [50, 51] that respectively approximate integers by their sign or by their “be even or odd”, the Interval abstract domain [51] that takes into account the minimum and the maximum value that an integer variable can be assigned to. The domains just mentioned do not track relationships between variables. Conversely, the following ones are relational: the Pentagon abstract domain [126] used to prove the safety of array accesses in byte-code and intermediate languages, the Octagon abstract domain [135] which allows to prove program invariants and the Trapezoid Step Function abstract domain [43] which approximates continuous functions. Recently, Bautista et al. [25] defined the Numeric Path Relations abstract domain, which denotes relations over structured data containing scalar values. Moreover, combinations of numerical abstract domains have been extensively used to improve the precision of the analysis and implemented in various tools [59, 64].

### **String Abstract Domains**

A string abstract domain approximates sets of string values. In the recent literature, different abstract domains for string analysis have been defined [11, 10, 34, 49, 57, 115, 127, 185] and different tools have been designed for their usage and combination [105, 112, 122, 162]. For instance, a basic, well-known domain is the String Set domain, which simply keeps track of a set of strings and it is a specific instance of the general (bounded) set domain. There exist general-purpose abstract domains that approximate shape and/or content information of strings, such as the Character Inclusion and the Prefix domains [49] that track the characters possibly and certainly contained in a string and its prefix, respectively. The String Length abstract domain approximates the minimum and maximum length of a string, similar to the interval domain for numerical values. Another general-purpose string domain is the String Hash domain proposed in [127], based on a distributive hash function. Some string abstract domains focus on specific programming languages such as `JavaScript` [16, 105] and `C` [108, 134], others target specific properties and rely on different approaches like automata [188], regular expressions [142], and grammatical summaries [115].

### 1.3 Contribution

The main goal of this thesis is to investigate existing string abstractions, understand the nature of string problems, design novel string abstract domains tailored for software security and verification, and refine the existing ones. Our contribution can be summarized as follows:

- We introduce M-String, a sophisticated C-specific string abstract domain. It fully approximates both the content and the shape of strings by representing them as abstractions of consecutive, non-overlapping, possibly empty segments. A segment represents a sequence of characters that share the same property and is surrounded by the segment bounds that abstracts concrete array indexes. By using M-String we may detect buffer overflows, among other verification tasks.
- We instantiate the Segmentation domain [57] for string analysis. We combine it with some existing general-purpose string abstract domains, to improve the accuracy of the analysis, also highlighting inconsistent states which may lead to program errors.
- We show how to compute the completion of string abstract domains to obtain an optimal abstract analysis, i.e., both sound and complete. In particular, we compute the complete version of two JavaScript-specific string abstract domains, i.e., TAJIS [105] and SAFE [122], with respect to some operations of interest. Moreover, we quantify the precision gained using a complete domain with respect to its uncomplete version.
- We design a framework that allows the definition of relational string abstract domains based on string orders. We make the analysis able to capture relations between single string variables and then between string expressions and string variables. Moreover, we instantiate our framework over three well-known string orders, i.e., the length inequality, char inclusion, and substring relation.
- As the last contribution, not strictly related to the previous ones, we introduce a new watermarking technique for relational textual databases based on synonyms substitution, which allows the preservation of the semantic consistency between the original database and the watermarked one. Our technique achieves a high level of capacity and robustness against the synonyms substitution attack, among others, while not damaging its imperceptibility requirement.

Note that most of the results we achieved are equipped with an experimental evaluation.



## 1.4 Thesis Structure

Chapter 2 introduces the technical notation used in the following chapters.

Chapter 3 introduces the definition of the M-String abstract domain for the analysis of `C` programs handling strings. M-String detects the presence of common strings management errors that may lead to undefined behaviour, e.g., buffer overflows. The contribution of this chapter has been published in [47, 48, 121].

Chapter 4 presents the Segmentation string abstract domain, an instantiation of the FunArray domain functor by Cousot et al. [57], for the analysis of arrays content, to string analysis. We discuss the limitations of basic existing string domains, and we combine them with the Segmentation to improve the accuracy of the analysis result. Combinations are carried out by exploiting the notions of reduced product [56] and Granger methodology [87]. Moreover, we study inconsistencies between abstract states as they may lead to program errors. The contribution of this chapter is submitted for publication.

Chapter 5 defines a way to proceed to compute the completion of string abstract domains [80] and applies the completion procedure to two incomplete `JavaScript` string abstract domains to enhance their precision. Namely, we complete TAJIS [105] and SAFE [122] abstract domains with respect to two operations of interest. Moreover, we provide a measure of the analysis precision increment gained using those complete domains with respect to the original ones. The contribution of this chapter has been published in [18].

Chapter 6 defines a general framework to design relational string abstract domains based on orders between string elements. The framework we present is first defined as tracking relations between single string variables and then enhanced to capture relations between string variables and string expressions. We instantiate it to a suit of string order relations of interest, i.e., length inequality, character inclusion, and substring relation. In particular, we evaluate the latter substring relational abstract domain comparing it with state-of-the-art general purpose string abstract domains. The contribution of this chapter is submitted for publication.

Chapter 7 presents a new watermarking technique to protect ownership of textual relational databases, which takes care of semantic consistency between the original database and the watermarked one. We design it carefully balancing robustness and imperceptibility of the watermark, and with the aim of increasing both its capacity and resilience to attacks, e.g., synonyms substitution attack. Moreover, we formalize a new approach to measure watermark capacity, and we prove the effectiveness of our technique through an experimental evaluation. The contribution of this chapter has been published in [85, 84].

Chapter 8 concludes and presents perspectives for future work. Appendix B recalls formal details of existing string abstract domains.



## Chapter 2

---

### PRELIMINARIES

In this chapter, we introduce the mathematical background and the theoretical results that will be used throughout the thesis. In particular, we recall some basic notions on set and order theory, functions, definitions on fixpoints and traces, and the basics of Abstract Interpretation, i.e., the framework we use to conduct string analysis.

#### Chapter Structure

Section 2.1 presents some notation on sets and sequences. Section 2.2 and Section 2.3 recall some definitions on order theory and functions, respectively. Section 2.4 and Section 2.5 provide notions on fixpoints and traces, respectively. Section 2.6 introduces results on Abstract Interpretation.

### 2.1 Sets and Sequences

#### Sets

A set is a, possibly finite, collection of distinct elements. Let  $A$  and  $B$  be sets. The set membership is denoted by  $x \in A$ , i.e., it says that  $x$  is an element of  $A$ . Let  $\phi$  be a formula. The set of all elements of  $A$  which satisfy the formula  $\phi$  is  $\{x \in A \mid \phi(x)\}$ . The inclusion relation is denoted by  $A \subseteq B$ , meaning that all the elements of  $A$  are also elements of  $B$  or, in other words, that  $A$  is a subset of  $B$ . The powerset of  $A$ , i.e.,  $\mathcal{P}(A)$ , is the set of all subsets of  $A$ . The difference between two sets is denoted by  $A \setminus B$  representing the set of all elements in  $A$  but not in  $B$ , i.e.,  $\{x \in A \mid x \notin B\}$ . The intersection between two sets is denoted by  $A \cap B$  and it represents the set of elements that belong to both  $A$  and  $B$ , i.e.,  $\{x \mid x \in A \wedge x \in B\}$ . The union between two sets is written  $A \cup B$  and it denotes the set of elements that belong to either  $A$  and  $B$ , i.e.,  $\{x \mid x \in A \vee x \in B\}$ . The Cartesian product between  $A$  and  $B$  is  $A \times B$ , i.e.,  $\{(x, y) \mid x \in A \wedge y \in B\}$ . We denote by  $A^n$  the  $n$ -ary Cartesian product of a set  $A$ , with  $n \geq 2$  and by  $|A|$  the cardinality of  $A$ . The emptyset is denoted by  $\emptyset$ . We denote by  $\exists^n$  the existential quantifier with cardinality  $n$ , which is read as “there exist exactly  $n$  objects”.

As usual in mathematics,  $\mathbb{N}$ ,  $\mathbb{Z}$ ,  $\mathbb{Q}$  and  $\mathbb{R}$  denote the set of all natural, integer, rational, and real numbers, respectively.

### Sequences

Let  $c \in \Sigma$  be a finite set of characters, i.e., an alphabet. We denote by  $\Sigma^*$  the set of all possible finite sequences of characters. Let  $\sigma \in \Sigma^*$ , we denote by  $|\sigma|$  the length of  $\sigma$ , by  $\text{char}(\sigma)$  the set of characters that occurs in the string  $\sigma$ . The empty string is denoted by  $\epsilon$ . A string  $\sigma$  may be also denoted by  $\sigma_1 \dots \sigma_n$  in order to refer to its characters.  $\sigma + \sigma'$  denotes the string concatenation, and  $\sigma[x/i]$  is  $\sigma$  where the character occurring in position  $i$  has been substituted with the character  $x$ . Given a string  $\sigma$ , we denote by  $\text{char}(\sigma)$  the set of of the characters occurring in  $\sigma$ .

## 2.2 Order Relations

A binary relation  $\sqsubseteq_L$  over a set  $L$  is a subset of  $L \times L$ , i.e.,  $\sqsubseteq_L \subseteq L \times L$ . A binary relation  $\sqsubseteq_L$  over a set  $L$  is a pre-order if and only if it is reflexive, i.e.,  $\forall x \in L : x \sqsubseteq_L x$ , and transitive, i.e.,  $\forall x_1, x_2, x_3 \in L : x_1 \sqsubseteq_L x_2 \wedge x_2 \sqsubseteq_L x_3 \Rightarrow x_1 \sqsubseteq_L x_3$ . A binary relation  $\sqsubseteq_L$  over a set  $L$  is a partial order if and only if it is a pre-order, and anti-symmetric, i.e.,  $\forall x_1, x_2 \in L : x_1 \sqsubseteq_L x_2 \wedge x_2 \sqsubseteq_L x_1 \Rightarrow x_1 = x_2$ . A binary relation  $\sqsubseteq_L$  over a set  $L$  is a total order if and only if it is a partial order, and  $\forall x_1, x_2 \in L : x_1 \sqsubseteq_L x_2 \vee x_2 \sqsubseteq_L x_1$ .

A set  $L$  equipped with an order relation is denoted by  $(L, \sqsubseteq_L)$ . A partially ordered set is also called poset. Given a poset  $(L, \sqsubseteq_L)$ , the infimum  $\perp_L$ , when it exists, is the element in  $L$  smaller than any other element of  $L$  with respect to  $\sqsubseteq_L$ . The supremum  $\top_L$ , when it exists, is the element in  $L$  greater than any other element of  $L$  with respect to  $\sqsubseteq_L$ . When both  $\perp_L$  and  $\top_L$  exist then  $(L, \sqsubseteq_L)$  is said to be bounded. Given a poset  $(L, \sqsubseteq_L)$  and a  $X \subseteq L$ , an upper bound of  $X$  is an element  $y \in L$  such that  $x \sqsubseteq_L y$ ,  $\forall x \in X$ . An upper bound  $y$  for  $X$  is the least upper bound (lub) of  $X$  if and only if for every other upper bound  $y'$  of  $X$  it holds that  $y \sqsubseteq_L y'$ . The least upper bound, when it exists, it is unique. Dually, a lower bound of  $X$  is an element  $y \in L$  such that  $y \sqsubseteq_L x$ ,  $\forall x \in X$ . A lower bound  $y$  of  $X$  is the greatest lower bound (glb) of  $X$  if and only if for every other lower bound  $y'$  of  $X$  it holds that  $y' \sqsubseteq_L y$ . Given a poset  $(L, \sqsubseteq_L)$  and  $X \subseteq L$ , we denote by  $\max^{\sqsubseteq_L}(X) = \{x \in X \mid \forall y \in X. x \sqsubseteq_L y \Rightarrow x = y\}$  the set of the maximal elements of  $X$  in  $L$  with respect to  $\sqsubseteq_L$ .

A subset  $X$  of a poset  $(L, \sqsubseteq_L)$  is a chain if it is totally ordered with respect to  $\sqsubseteq_L$ , i.e.,  $\forall x_1, x_2 \in X : x_1 \sqsubseteq_L x_2 \vee x_2 \sqsubseteq_L x_1$ .  $X$  is an ascending chain if it is ordered and each element is greater or equal to the previous one with respect to  $\sqsubseteq_L$ .  $X$  is an descending chain if it is ordered and each element is smaller or equal to the previous one with respect to  $\sqsubseteq_L$ . A poset  $(L, \sqsubseteq_L)$  satisfies the ascending chain condition (ACC) if every ascending chain  $x_1 \sqsubseteq_L x_2 \sqsubseteq_L \dots$  of elements in  $L$  is eventually stationary, i.e.,  $\exists n \in \mathbb{N} : \forall m > n : x_m = x_n$ .

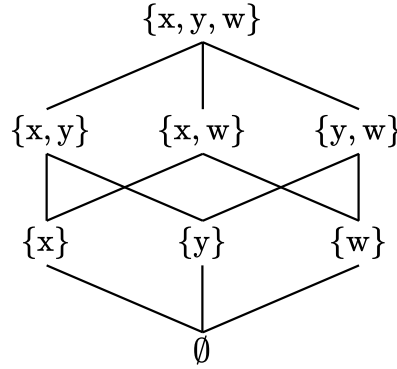


FIGURE 2.1: The Hasse diagram of the powerset of a 3-element set ordered by inclusion, i.e.,  $(\mathcal{P}(\{x, y, w\}), \subseteq)$ .

Similarly,  $(L, \sqsubseteq_L)$  satisfies the descending chain condition (DCC) if there is no infinite descending chain. Every finite poset satisfies the ACC and DCC.

Note that a poset can be graphically represented by the so-called Hasse diagram (cf. Figure 2.1). Given a poset  $(L, \sqsubseteq_L)$ , its corresponding Hasse diagram has a node for each element of  $L$  and an edge links two elements  $x_1, x_2 \in L$  if  $x_1 \sqsubseteq_L x_2 \wedge \nexists y \in L : x_1 \sqsubseteq_L y \wedge y \sqsubseteq_L x_2$ .

A poset  $(L, \sqsubseteq_L)$  is a directed set if  $\forall X \subseteq L$  such that  $X \neq \emptyset$  and  $X$  is finite, then  $X$  has a least upper bound in  $L$ . Dually, a poset  $(L, \sqsubseteq_L)$  is a co-directed set if  $\forall X \subseteq L$  such that  $X \neq \emptyset$  and  $X$  is finite, then  $X$  has a greatest lower bound in  $L$ . A poset  $(L, \sqsubseteq_L)$  is a complete partial order (cpo) on directed sets if it has an infimum and for each directed set  $D \subseteq L$ ,  $D$  has a least upper bound in  $L$ .

### Lattices

A poset  $(L, \sqsubseteq_L)$  is a join semi lattice if for each  $x_1, x_2 \in L$  there exists a least upper bound  $x_1 \sqcup_L x_2$  and it belongs to  $L$ . A poset  $(L, \sqsubseteq_L)$  is a meet semi lattice if for each  $x_1, x_2 \in L$  there exists a greatest lower bound  $x_1 \sqcap_L x_2$  and it belongs to  $L$ . A poset  $(L, \sqsubseteq_L)$  is a lattice if it has an infimum, a supremum and, for each  $x_1, x_2 \in L$ , both the greatest lower bound  $x_1 \sqcap_L x_2$  and the least upper bound  $x_1 \sqcup_L x_2$  exist and belong to  $L$ . A lattice is denoted by  $(L, \sqsubseteq_L, \perp_L, \top_L, \sqcap_L, \sqcup_L)$ . A complete lattice is a lattice where each subset of  $L$  has both a least upper bound and a greatest lower bound, i.e.,  $\forall X \subseteq L : \prod_L X, \bigsqcup_L X \in L$ . Any finite lattice is complete. A complete lattice is always a cpo.

Let  $L$  be a complete lattice.  $X \subseteq L$  is a Moore family of  $L$  if  $X = \mathcal{M}(X) = \{\prod_L S \mid S \subseteq X\}$ , where  $\prod_L \emptyset = \top_L \in \mathcal{M}(X)$ . We denote by  $\mathcal{M}(X)$  the Moore closure of  $X \subseteq L$ , that is the smallest subset of  $L$  which contains  $X$  and it is a Moore family of  $L$ .

## 2.3 Functions

A function is a mathematical object that relates two sets by associating each element of the first set (domain) to an element of the second set (codomain). Let  $f$  be a function with domain  $A$  and codomain  $B$ . We denote by  $f : A \rightarrow B$  a function from elements of  $A$  to elements of  $B$ . Given  $f : A \rightarrow B$  and  $g : B \rightarrow C$  we denote with  $g \circ f : A \rightarrow C$  their composition, i.e.,  $g \circ f = \lambda x.g(f(x))$ . When  $f$  is a function of arity  $n$ , i.e.,  $f : A^n \rightarrow B$ , with  $p \in A^n$  and  $i \in [0, n)$ ,  $f_p^i = \lambda x.f(p[x/i]) : A \rightarrow B$  is the same function where all the parameters but the  $i$ -th are fixed by  $p$ .

Given two posets  $(A, \sqsubseteq_A, \sqcap_A, \sqcup_A)$  and  $(B, \sqsubseteq_B, \sqcap_B, \sqcup_B)$ , a function  $f : A \rightarrow B$  is:

- monotone

$$\text{if } \forall x_1, x_2 \in A : x_1 \sqsubseteq_A x_2 \Rightarrow f(x_1) \sqsubseteq_B f(x_2)$$

- continuous

$$\text{if } \forall X \subseteq A : X \text{ is a chain} \wedge f(\bigsqcup_A X) = \bigsqcup_B \{f(x) \mid x \in X\}$$

- join preserving

$$\text{if } \forall x_1, x_2 \in A : f(x_1 \sqcup_A x_2) = f(x_1) \sqcup_B f(x_2)$$

- complete join preserving

$$\text{if } \forall X \subseteq A : \bigsqcup_A X \text{ exists} \wedge f(\bigsqcup_A X) = \bigsqcup_B f(X)$$

- meet preserving

$$\text{if } \forall x_1, x_2 \in A : f(x_1 \sqcap_A x_2) = f(x_1) \sqcap_B f(x_2)$$

- complete meet preserving

$$\text{if } \forall X \subseteq A : \sqcap_A X \text{ exists} \wedge f(\sqcap_A X) = \sqcap_B f(X)$$

The notion of function can be generalized by the one of partial function which allows a function to map some elements of its domain to an element of its codomain. Formally, a partial function  $f : A \rightarrow B$  is a function  $f : X \rightarrow B$  for some  $X \subseteq A$ .

An upper closure operator on a poset  $(A, \sqsubseteq_A)$  is an operator  $\rho : A \rightarrow A$  which is monotone, idempotent (i.e.,  $\forall x \in A : \rho(\rho(x)) = \rho(x)$ ), and extensive (i.e.,  $\forall x \in A : x \sqsubseteq_A \rho(x)$ ). The set of all closure operators on a poset  $A$  is denoted by  $uco(A)$  and it is a complete lattice.

## 2.4 Fixpoints

Given a poset  $(L, \sqsubseteq_L)$  and an operator  $f : L \rightarrow L$ , a fixpoint of  $f$  is an element  $x \in L$  such that  $f(x) = x$ . The set of all fixpoints of  $f$  is  $\text{Fix}(f) = \{x \in L \mid f(x) = x\}$ . Moreover, a

pre-fixpoint of  $f$  is an element  $x \in L$  such that  $x \sqsubseteq_L f(x)$  and a post-fixpoint of  $f$  is an element  $x \in L$  such that  $f(x) \sqsubseteq_L x$ . It follows that the sets of all pre- and post-fixpoints of  $f$  are  $\text{preFix}(f) = \{x \in L \mid x \sqsubseteq_L f(x)\}$  and  $\text{postFix}(f) = \{x \in L \mid f(x) \sqsubseteq_L x\}$ . A fixpoint  $x \in \text{Fix}(f)$  is a least fixpoint of  $f$  if  $\forall y \in \text{Fix} : x \sqsubseteq_L y$ . A fixpoint  $x \in \text{Fix}(f)$  is a greatest fixpoint of  $f$  if  $\forall y \in \text{Fix} : y \sqsubseteq_L x$ . When the least and greatest fixpoints of  $f$  exist they are unique and they are denoted by  $\text{lfp}_x^{\sqsubseteq_L}(f)$  and  $\text{gfp}_x^{\sqsubseteq_L}(f)$ , respectively. Precisely,  $\text{lfp}_x^{\sqsubseteq_L}(f)$  denotes the least fixpoint of  $f$  greater than  $x$  with respect to the order  $\sqsubseteq_L$ , and  $\text{gfp}_x^{\sqsubseteq_L}(f)$  denotes the greatest fixpoint of  $f$  smaller than  $x$  with respect to the order  $\sqsubseteq_L$ .

The following theorems guarantee the existence of the least and greatest fixpoints of a monotonic operator defined over a complete lattice.

**Theorem 2.1** (Tarski's Theorem [74]). Let  $(L, \sqsubseteq_L, \perp_L, \top_L, \sqcap_L, \sqcup_L)$  be a complete lattice and let  $f : L \rightarrow L$  be a monotonic operator on this lattice. The set of fixpoints is a non-empty complete lattice, and:

$$\text{lfp}_{\perp}^{\sqsubseteq_L}(f) = \sqcap_L \{x \in L \mid f(x) \sqsubseteq_L x\}$$

$$\text{gfp}_{\perp}^{\sqsubseteq_L}(f) = \sqcup_L \{x \in L \mid x \sqsubseteq_L f(x)\}$$

□

**Theorem 2.2** (Constructive version of Tarski's Theorem [55]). Let  $(L, \sqsubseteq_L, \perp_L, \top_L, \sqcap_L, \sqcup_L)$  be a complete lattice and let  $f : L \rightarrow L$  be a monotonic operator on this lattice. Define the following sequence:

$$f^0 = \perp_L$$

$$f^\delta = f(f^{\delta-1}) \quad \text{for every successor ordinal } \delta$$

$$f^\delta = \bigsqcup_{\alpha < \delta} f^\alpha \quad \text{for every limit ordinal } \delta$$

Then the ascending chain  $\{f^i \mid 0 \leq i \leq \delta\}$ , where  $\delta$  is an ordinal, is ultimately stationary for some  $\rho \in \mathbb{N}$  that is  $f^\rho = \text{lfp}_{\perp}^{\sqsubseteq_L}(f)$ .

□

## 2.5 Traces

Given a set  $X$ , a trace  $\tau : \mathbb{N} \rightarrow X$  is a partial function such that:

$$\forall i \in \mathbb{N} : i \notin \text{dom}(\tau) \Rightarrow \forall j > i : j \notin \text{dom}(\tau)$$

where  $\text{dom}(\tau)$  denotes the domain of a trace  $\tau$ . It follows that the domain of all non-empty traces is a segment of  $\mathbb{N}$ . The empty trace is denoted by  $\varepsilon_\tau$ , i.e., the trace  $\tau$  such that  $\text{dom}(\tau) = \emptyset$ . Let  $X$  be a generic set of elements, we denote by  $X^+$  the set of all finite

traces composed by elements in  $X$ . Let  $\text{len}_{\text{tr}} : X^+ \rightarrow \mathbb{N}$  be the function that, given a trace, returns its length, then:  $\text{len}_{\text{tr}}(\tau) = i + 1$  such that  $i \in \text{dom}(\tau) \wedge i + 1 \notin \text{dom}(\tau)$ . Note that if  $\tau = \varepsilon_\tau$  then  $\text{len}_{\text{tr}}(\tau) = 0$ . A trace can be represented as a sequence of states, i.e.,  $s_0 \rightarrow s_1 \rightarrow \dots$ , that corresponds to the trace  $\{(0, s_0), (1, s_1), \dots\}$ . We denote by  $X_{\overset{\text{T}}{\rightarrow}}^+$  the set of traces in  $X^+$  ending with a final state with respect to the transition  $\overset{\text{T}}{\rightarrow}$ , i.e.,  $X_{\overset{\text{T}}{\rightarrow}}^+ = \{s_0 \rightarrow \dots \rightarrow s_i \mid s_0 \rightarrow \dots \rightarrow s_i \in X^+ \wedge \nexists s_j \in X \text{ such that } s_i \overset{\text{T}}{\rightarrow} s_j\}$ . Given a set of initial elements  $X_0$ , a set of states  $S$ , and a transition relations  $\overset{\text{T}}{\rightarrow} \subseteq S \times S$ , the partial trace semantics [56] builds up all the traces that can be obtained by starting from traces containing only a single element from  $X_0$  and then iteratively applying the transition relation until a fixpoint is reached.

**Definition 2.1** (Partial trace semantics [56]). Let  $S$  be a set of states,  $X_0 \subseteq S$  be a set of initial elements, and  $\overset{\text{T}}{\rightarrow} \subseteq S \times S$  be a transition relation. Let  $f : \mathcal{P}(S) \rightarrow (S^+ \rightarrow S^+)$  be the function defined as:

$$F(X_0) = \lambda T. \{s_0 \mid s_0 \in X_0\} \cup \\ \{s_0 \rightarrow \dots \rightarrow s_{i-1} \rightarrow s_i \mid s_0 \rightarrow \dots \rightarrow s_{i-1} \in T \wedge s_{i-1} \overset{\text{T}}{\rightarrow} s_i\}$$

The partial trace semantics is defined as:

$$\mathbb{PT}\llbracket X_0 \rrbracket = \text{lfp}_{\emptyset}^{\sqsubseteq} F(X_0)$$

△

## 2.6 Abstract Interpretation

Abstract Interpretation [51, 56] is a theoretical framework for sound reasoning about program semantic properties of interest. Precisely it formalizes the (concrete) semantics of a program as a fixpoint of all its execution traces, and it allows to define and prove the soundness of (possible computable) semantics at different degrees of abstraction.

### 2.6.1 Galois Connection

Formally, the concrete and abstract semantics are defined on a concrete  $\mathbf{D}$  and abstract  $\overline{\mathbf{D}}$  sets, respectively. The concrete and the abstract domains are modelled as posets, i.e.,  $(\mathbf{D}, \sqsubseteq_{\mathbf{D}})$  and  $(\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$ , respectively. The concrete and the abstract domains are related by a pair of monotone functions: the concretization  $\gamma_{\overline{\mathbf{D}}} : \overline{\mathbf{D}} \rightarrow \mathbf{D}$  and the abstraction  $\alpha_{\overline{\mathbf{D}}} : \mathbf{D} \rightarrow \overline{\mathbf{D}}$  functions.

**Definition 2.2** (Galois Connection [51]). Let  $(\mathbf{D}, \sqsubseteq_{\mathbf{D}})$  and  $(\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$  be two partial orders. Two functions  $\alpha_{\overline{\mathbf{D}}} : \mathbf{D} \rightarrow \overline{\mathbf{D}}$  and  $\gamma_{\overline{\mathbf{D}}} : \overline{\mathbf{D}} \rightarrow \mathbf{D}$  form a Galois Connection if and only if

$$\forall d \in \mathbf{D} : \exists \overline{d} \in \overline{\mathbf{D}} : \alpha_{\overline{\mathbf{D}}}(d) \sqsubseteq_{\overline{\mathbf{D}}} \overline{d} \Rightarrow d \sqsubseteq_{\mathbf{D}} \gamma_{\overline{\mathbf{D}}}(\overline{d})$$



We denote this fact by writing

$$(\mathbf{D}, \sqsubseteq_{\mathbf{D}}) \xleftrightarrow[\alpha_{\overline{\mathbf{D}}}]^{\gamma_{\overline{\mathbf{D}}}} (\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$$

△

**Theorem 2.3** ([51]). Let  $(\mathbf{D}, \sqsubseteq_{\mathbf{D}})$  and  $(\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$  be two partial orders and let  $\alpha_{\overline{\mathbf{D}}} : \mathbf{D} \rightarrow \overline{\mathbf{D}}$  and  $\gamma_{\overline{\mathbf{D}}} : \overline{\mathbf{D}} \rightarrow \mathbf{D}$  be two maps such that:

- $\alpha_{\overline{\mathbf{D}}}$  and  $\gamma_{\overline{\mathbf{D}}}$  are monotone
- $\alpha_{\overline{\mathbf{D}}} \circ \gamma_{\overline{\mathbf{D}}}$  is reductive, i.e.,  $\forall \overline{\mathbf{d}} \in \overline{\mathbf{D}} : \alpha_{\overline{\mathbf{D}}} \circ \gamma_{\overline{\mathbf{D}}}(\overline{\mathbf{d}}) \sqsubseteq_{\overline{\mathbf{D}}} \overline{\mathbf{d}}$
- $\alpha_{\overline{\mathbf{D}}} \circ \gamma_{\overline{\mathbf{D}}}$  is extensive, i.e.,  $\forall \mathbf{d} \in \mathbf{D} : \mathbf{d} \sqsubseteq_{\mathbf{D}} \gamma_{\overline{\mathbf{D}}} \circ \alpha_{\overline{\mathbf{D}}}(\mathbf{d})$

Then, it holds that  $(\mathbf{D}, \sqsubseteq_{\mathbf{D}}) \xleftrightarrow[\alpha_{\overline{\mathbf{D}}}]^{\gamma_{\overline{\mathbf{D}}}} (\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$

□

A Galois Connection can be induced by an abstraction function that is a complete lub  $(\sqcup_{\overline{\mathbf{D}}})$  preserving map or by a concretization function that is a complete glb  $(\sqcap_{\mathbf{D}})$  preserving map, as proved by Proposition 7 of [52].

**Theorem 2.4** (Galois Connection induced by lub preserving maps). Let  $\alpha_{\overline{\mathbf{D}}} : \mathbf{D} \rightarrow \overline{\mathbf{D}}$  be a complete join preserving maps between posets  $(\mathbf{D}, \sqsubseteq_{\mathbf{D}})$  and  $(\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$ . Define:

$$\gamma_{\overline{\mathbf{D}}} = \lambda \overline{\mathbf{d}}. \sqcup_{\mathbf{D}} \{ \mathbf{d} \mid \alpha_{\overline{\mathbf{D}}}(\mathbf{d}) \sqsubseteq_{\overline{\mathbf{D}}} \overline{\mathbf{d}} \}$$

If  $\gamma_{\overline{\mathbf{D}}}$  is well define, then:

$$(\mathbf{D}, \sqsubseteq_{\mathbf{D}}) \xleftrightarrow[\alpha_{\overline{\mathbf{D}}}]^{\gamma_{\overline{\mathbf{D}}}} (\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$$

□

**Theorem 2.5** (Galois Connection induced by glb preserving maps). Let  $\gamma_{\overline{\mathbf{D}}} : \overline{\mathbf{D}} \rightarrow \mathbf{D}$  be a complete join preserving maps between posets  $(\mathbf{D}, \sqsubseteq_{\mathbf{D}})$  and  $(\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$ . Define:

$$\alpha_{\overline{\mathbf{D}}} = \lambda \mathbf{d}. \sqcap_{\overline{\mathbf{D}}} \{ \overline{\mathbf{d}} \mid \mathbf{y} \sqsubseteq_{\mathbf{D}} \gamma_{\overline{\mathbf{D}}}(\overline{\mathbf{d}}) \}$$

If  $\alpha_{\overline{\mathbf{D}}}$  is well define, then:

$$(\mathbf{D}, \sqsubseteq_{\mathbf{D}}) \xleftrightarrow[\alpha_{\overline{\mathbf{D}}}]^{\gamma_{\overline{\mathbf{D}}}} (\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$$

□

Galois Connections are compositional, i.e., the composition of two Galois Connections is still a Galois Connection.

**Theorem 2.6** (Composition of Galois Connections). Suppose that:

- $(\mathbf{D}, \sqsubseteq_{\mathbf{D}}) \xleftrightarrow[\alpha_{\overline{\mathbf{D}}}]^{\gamma_{\overline{\mathbf{D}}}} (\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$  and

$$\bullet (\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}}) \xleftrightarrow[\alpha_{\overline{\mathbf{D}'}}]{\gamma_{\overline{\mathbf{D}'}}} (\overline{\mathbf{D}'}, \sqsubseteq_{\overline{\mathbf{D}'}})$$

Then,

$$(\mathbf{D}, \sqsubseteq_{\mathbf{D}}) \xleftrightarrow[\alpha_{\mathbf{D}} \circ \alpha_{\overline{\mathbf{D}'}}]{\gamma_{\overline{\mathbf{D}}} \circ \gamma_{\overline{\mathbf{D}'}}} (\overline{\mathbf{D}'}, \sqsubseteq_{\overline{\mathbf{D}'}})$$

□

Note that, as it is possible to define Galois Connections also in terms of upper closure operators, the framework of Abstract Interpretation can be equivalently formalized either as Galois Connections or closure operators on a given concrete domain, which is a complete lattice  $\mathbf{D}$  [56]. In particular, given  $(\mathbf{D}, \sqsubseteq_{\mathbf{D}}) \xleftrightarrow[\alpha_{\mathbf{D}}]{\gamma_{\overline{\mathbf{D}}}} (\overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}})$ , then  $\gamma_{\overline{\mathbf{D}}} \circ \alpha_{\overline{\mathbf{D}}} \in uco(\mathbf{D})$ . In particular, an upper closure operator (over  $\mathbf{D}$ ) uniquely identifies a Galois Connection, and vice-versa.

### Lattice of Abstract Domains

If  $\mathbf{D}$  is a complete lattice, then  $(uco(\mathbf{D}), \sqsubseteq, \lambda x.x, \lambda x.\top, \sqcap, \sqcup)$  forms a complete lattice [51], which is the set of all possible abstractions of  $\mathbf{D}$ , where  $\lambda x.x$  is the bottom element and  $\lambda x.\top$  the top element. Let  $\rho, \eta \in uco(\mathbf{D})$ ,  $\{\rho_i\}_{i \in I} \subseteq uco(\mathbf{D})$ . We say that  $\rho$  is more precise than  $\eta$ , i.e.,  $\rho \sqsubseteq \eta$  if and only if  $\forall \mathbf{d} \in \mathbf{D} : \rho(\mathbf{d}) \sqsubseteq \eta(\mathbf{d})$ . Let  $\mathbf{d} \in \mathbf{D}$ , the greatest lower bound is  $(\sqcap_{i \in I} \rho_i)(\mathbf{d}) = \sqcap_{i \in I} \rho_i(\mathbf{d})$  and the least upper bound is  $(\sqcup_{i \in I} \rho_i)(\mathbf{d}) = \mathbf{d}$  if and only if  $\forall i \in I : \rho_i(\mathbf{d}) = \mathbf{d}$ .

### 2.6.2 Soundness and Completeness

Two fundamental notions in Abstract Interpretation are those of soundness and completeness. Abstract Interpretation requires that the abstract semantics must be sound with respect to the concrete one, i.e., the properties verified by the abstraction process must be valid with respect to the concrete semantics. On the other hand, Abstract Interpretation accepts incompleteness of the abstract results, i.e., the abstract results can be less precise than the concrete ones.

#### Soundness

To accomplish soundness, the concretization of the result of the abstract computation must over-approximate the result of the concrete semantics.

**Definition 2.3** (Soundness). Let  $\mathbf{D}$  and  $\overline{\mathbf{D}}$  form a Galois Connection. Moreover, let  $\mathbb{S} : \mathbf{D} \rightarrow \mathbf{D}$  and  $\mathbb{S}_{\overline{\mathbf{D}}} : \overline{\mathbf{D}} \rightarrow \overline{\mathbf{D}}$  be the concrete and the abstract semantics, respectively. The abstract semantics is sound if and only if  $\forall \rho \in \text{preFix}(\mathbb{S}_{\overline{\mathbf{D}}}) \subseteq \overline{\mathbf{D}}$ , we have that:

$$\gamma_{\overline{\mathbf{D}}} \circ \mathbb{S}_{\overline{\mathbf{D}}}[\rho] \sqsupseteq_{\mathbf{D}} \mathbb{S}[\gamma_{\overline{\mathbf{D}}}(\rho)]$$

△

Definition 2.3 presents soundness based on the computational process of the concrete domain, but we could define it also in terms of the computational process of the abstract domain.

Among all the sound abstract functions  $\mathbb{S}_{\overline{\mathbf{D}}}$ , we aim at the best one, namely the one that loses less information approximating  $\mathbb{S}$  (cf. Definition 2.4).

**Definition 2.4** (Best correct approximation). Given  $\mathbf{D} \xleftrightarrow[\alpha]{\gamma} \overline{\mathbf{D}}$  and a concrete function  $\mathbb{S} : \mathbf{D} \rightarrow \mathbf{D}$ , the function  $\alpha_{\overline{\mathbf{D}}} \circ \mathbb{S} \circ \gamma_{\overline{\mathbf{D}}} : \overline{\mathbf{D}} \rightarrow \overline{\mathbf{D}}$  is the best correct approximation of  $\mathbb{S}$ .

△

Note that there exist different ways of proving the soundness of an abstract semantics [52]. Moreover, it is still acceptable if a Galois Connection between the concrete and the abstract domain cannot be induced, i.e., if does not exist a minimal element abstracting a concrete property of interest, as we can always choose an arbitrary abstract element among all the available abstractions of a given concrete element [52] or weaken the relation between the concrete and the abstract semantics without affecting the soundness [54].

### Completeness

In Abstract Interpretation, there exist two notions of completeness: backward and forward. Backward completeness property focuses on complete abstractions of the inputs, while forward completeness [75, 76, 77] focuses on complete abstractions of the outputs, both with respect to an operation of interest.

**Definition 2.5** (Backward completeness). Given  $\mathbf{D} \xleftrightarrow[\alpha]{\gamma} \overline{\mathbf{D}}$ , a concrete function  $\mathbb{S} : \mathbf{D} \rightarrow \mathbf{D}$  and an abstract function  $\mathbb{S}_{\overline{\mathbf{D}}} : \overline{\mathbf{D}} \rightarrow \overline{\mathbf{D}}$ ,  $\mathbb{S}_{\overline{\mathbf{D}}}$  is backward complete with respect to  $\mathbb{S}$  if:

$$\alpha_{\overline{\mathbf{D}}} \circ \mathbb{S} = \mathbb{S}_{\overline{\mathbf{D}}} \circ \alpha_{\overline{\mathbf{D}}}$$

△

**Definition 2.6** (Forward completeness). Given  $\mathbf{D} \xleftrightarrow[\alpha]{\gamma} \overline{\mathbf{D}}$ , a concrete function  $\mathbb{S} : \mathbf{D} \rightarrow \mathbf{D}$  and an abstract function  $\mathbb{S}_{\overline{\mathbf{D}}} : \overline{\mathbf{D}} \rightarrow \overline{\mathbf{D}}$ ,  $\mathbb{S}_{\overline{\mathbf{D}}}$  is forward complete with respect to  $\mathbb{S}$  if:

$$\mathbb{S} \circ \gamma_{\overline{\mathbf{D}}} = \gamma_{\overline{\mathbf{D}}} \circ \mathbb{S}_{\overline{\mathbf{D}}}$$

△

Having backward complete abstract functions means that we have the guarantee that no loss of information arises during the input abstraction process, with respect to an operation of interest. Conversely, having forward complete abstract functions means that we have the guarantee that no loss of information arises during the output abstraction process, with respect to an operation of interest. While forward completeness is less known, backward completeness is the one recognized ad standard notion of completeness in Abstract Interpretation, and in this thesis, we will focus on it.

### 2.6.3 Fixpoints Approximation

As mentioned above, the concrete and the abstract semantics are defined as the fixpoint computation of monotonic functions.

If the abstract domain respects the ACC, then the abstract computation terminates in a finite time. Otherwise a widening operator [46, 51] must be defined to force the convergence of the analysis, still maintaining the soundness.

**Definition 2.7** (Widening). Let  $(L, \sqsubseteq_L)$  be a poset. A widening operator  $\nabla_L : L \times L \rightarrow L$  satisfies the following conditions:

- $\forall x, y \in L : x \sqsubseteq_L (x \nabla_L y) \wedge y \sqsubseteq_L (x \nabla_L y)$
- for every increasing chain  $x_0 \sqsubseteq_L x_1 \sqsubseteq_L \dots$ , the chain defined as

$$\begin{aligned} w_0 &= x_0 \\ w_{n+1} &= w_n \nabla_L x_{n-1} \end{aligned}$$

is not strictly increasing.

△

The widening often lead to a dramatic loss of information. Thus, a narrowing operator is usually defined to improve the accuracy of the widening result.

**Definition 2.8** (Narrowing). Let  $(L, \sqsubseteq_L)$  be a poset. A narrowing operator  $\Delta_L : L \times L \rightarrow L$  satisfies the following conditions:

- $\forall x, y \in L : x \sqsubseteq_L (x \Delta_L y) \sqsubseteq_L y$
- for every decreasing chain  $x_0 \supseteq_L x_1 \supseteq_L \dots$ , the chain defined as

$$\begin{aligned} w_0 &= x_0 \\ w_{n+1} &= w_n \Delta_L x_{n-1} \end{aligned}$$

is not strictly decreasing.

△

### 2.6.4 Product Operators

Abstract Interpretation allows the combination, through *enhancing* operators, of different abstract domains within the same analysis, without affecting its soundness. Combinations aim to improve the accuracy of the analysis, also at the cost of increasing its complexity.

Enhancing operators are those who refine the information tracked by abstract domains. They can be formalized as lower closure operators on the set of all abstract interpretations

of a concrete domain  $\mathbf{D}$  [42, 68, 79]. The best-known enhancement operators in Abstract Interpretation theory include: reduced product, disjunctive completion, reduced cardinal power [56], tensor product [137], open product, pattern completion [40, 41], functional dependencies [78], complete shell [80], and logical product [88]. The reduced product is a refinement of the Cartesian product as it enhance the accuracy of the information tracked by one domain using the information tracked by the other, and vice-versa, instead of running the two abstract analysis in parallel. The disjunctive completion [56] enhances an abstract domain by adding denotations for concrete disjunctions of its values [68]. In contrast, the reduced cardinal power [56] captures disjunctive information over abstract states [42], being suitable for relational analysis. A further Cartesian product refinement is the open product, presented by Cortesi et al. in [40, 41]. It works on open abstract interpretations, which include queries and open operations. In [80], Giacobazzi et al. presented a constructive way to obtain the so-called complete shell of an abstract interpretation, i.e., a domain transformer which includes the minimal number of abstract point to a certain abstract domain  $\overline{\mathbf{D}}$  to make it complete with respect to a certain operation of interest. The benefits of working with complete string abstract domains have been proven in Chapter 5. Finally, the logical product [88], is more precise than the reduced product, and it combines lattices which are defined over convex, stably infinite and disjoint theories.

In the following, we recall the Cartesian, reduced, and Granger products.

### Cartesian Product

Let  $\overline{\mathbf{D}}$  and  $\overline{\mathbf{D}'}$  be two abstract domains representing sound approximations of the same concrete domain  $\mathbf{D}$  and let  $\mathfrak{C} = \overline{\mathbf{D}} \times \overline{\mathbf{D}'}$  be their Cartesian product. The abstract elements in  $\mathfrak{C}$  are pairs  $(\overline{\mathbf{d}}, \overline{\mathbf{d}'})$  such that  $\overline{\mathbf{d}} \in \overline{\mathbf{D}}$  and  $\overline{\mathbf{d}'} \in \overline{\mathbf{D}'}$ . The partial order, upper and lower bound and widening operators are defined as the component-wise application of the corresponding operators of the two domains. The Cartesian product is a lattice. The abstraction function on  $\mathfrak{C}$  maps a concrete element  $\mathbf{d} \in \mathbf{D}$  to the pair  $(\alpha_{\overline{\mathbf{D}}}(\mathbf{d}), \alpha_{\overline{\mathbf{D}'}}(\mathbf{d}))$ , while the concretization function on  $\mathfrak{C}$  maps an abstract element to  $\gamma_{\overline{\mathbf{D}}}(\overline{\mathbf{d}}) \sqcap_{\mathbf{D}} \gamma_{\overline{\mathbf{D}'}}(\overline{\mathbf{d}'})$ . Then, the Cartesian product forms a Galois Connection with the concrete domain. An abstract domain functor is a function from the parameter abstract domains  $\overline{\mathbf{D}}_1, \overline{\mathbf{D}}_2, \dots, \overline{\mathbf{D}}_n$  to a new abstract domain  $\overline{\mathbf{D}}(\overline{\mathbf{D}}_1, \overline{\mathbf{D}}_2, \dots, \overline{\mathbf{D}}_n)$ . The abstract domain functor  $\overline{\mathbf{D}}(\overline{\mathbf{D}}_1, \overline{\mathbf{D}}_2, \dots, \overline{\mathbf{D}}_n)$  composes abstract domain properties of the parameter abstract domains to build a new class of abstract properties and operations [57].

### Reduced Product

The reduced product [56] combines abstract domains, mutually refining them. Informally speaking, it improves the precision of the information tracked by one domain exploiting the information tracked by the other, and vice-versa [42]. Let  $(\overline{\mathbf{d}}_1, \overline{\mathbf{d}}_2) \in \mathfrak{C}$  be the Cartesian product of  $\overline{\mathbf{D}}_1$  and  $\overline{\mathbf{D}}_2$ . The reduced product is a Cartesian product equipped with a

reduction operator. In particular, the reduced product search for the smallest element  $(\bar{\mathbf{d}}'_1, \bar{\mathbf{d}}'_2)$  such that:  $\gamma_{\bar{\mathbf{D}}_1}(\bar{\mathbf{d}}'_1) \subseteq \gamma_{\bar{\mathbf{D}}_1}(\bar{\mathbf{d}}_1)$ ,  $\gamma_{\bar{\mathbf{D}}_2}(\bar{\mathbf{d}}'_2) \subseteq \gamma_{\bar{\mathbf{D}}_2}(\bar{\mathbf{d}}_2)$ , and  $\gamma_{\bar{\mathbf{D}}_1}(\bar{\mathbf{d}}'_1) \cap \gamma_{\bar{\mathbf{D}}_2}(\bar{\mathbf{d}}'_2) = \gamma_{\bar{\mathbf{D}}_1}(\bar{\mathbf{d}}_1) \cap \gamma_{\bar{\mathbf{D}}_2}(\bar{\mathbf{d}}_2)$ . Formally, the reduction operator,  $\rho : \mathfrak{C} \rightarrow \mathfrak{C}$ , is defined as follows: let  $\mathfrak{c}$  denote the pair  $(\bar{\mathbf{d}}_1, \bar{\mathbf{d}}_2)$  and  $\mathfrak{c}^*$  denote the pair  $(\bar{\mathbf{d}}'_1, \bar{\mathbf{d}}'_2)$  then,  $\rho(\mathfrak{c}) = \prod_{\mathfrak{c}} \{ \mathfrak{c}^* \in \mathfrak{C} : \gamma_{\mathfrak{c}}(\mathfrak{c}) \sqsubseteq_{\mathfrak{c}} \gamma_{\mathfrak{c}}(\mathfrak{c}^*) \}$ , where  $(\bar{\mathbf{d}}_1, \bar{\mathbf{d}}_2) \sqsubseteq_{\mathfrak{c}} (\bar{\mathbf{d}}'_1, \bar{\mathbf{d}}'_2) \Leftrightarrow \bar{\mathbf{d}}_1 \sqsubseteq_{\bar{\mathbf{D}}_1} \bar{\mathbf{d}}'_1$  and  $\bar{\mathbf{d}}_2 \sqsubseteq_{\bar{\mathbf{D}}_2} \bar{\mathbf{d}}'_2$ . In the following we will refer to the reduced product between two abstract elements  $\bar{\mathbf{d}}_1$  and  $\bar{\mathbf{d}}_2$  as  $\bar{\mathbf{d}}_1 \otimes \bar{\mathbf{d}}_2$ .

### Granger Product

The Granger product [87] is an approximation of the reduction operator. It is based on two refinement operators iterated until a fixpoint is reached (or, in other words, when the smallest reduction is obtained). Each of these operators takes advantage of the information of one of the two domains involved in the product and it improves the information of the other using the information of the first one [42]. Let  $\bar{\mathbf{D}}_1$  and  $\bar{\mathbf{D}}_2$  be two abstract domains,  $\bar{\mathbf{d}}_1$  and  $\bar{\mathbf{d}}_2$  be abstract elements belonging to  $\bar{\mathbf{D}}_1$  and  $\bar{\mathbf{D}}_2$  respectively, and  $\mathfrak{C}$  be their Cartesian product. The Granger operators are defined as follows:  $\rho_1 : \mathfrak{C} \rightarrow \bar{\mathbf{D}}_1$  and  $\rho_2 : \mathfrak{C} \rightarrow \bar{\mathbf{D}}_2$ . In order to have a sound over-approximation of the reduction operator,  $\rho_1$  and  $\rho_2$  have to satisfy the following conditions: let  $\mathfrak{c}$  denote the pair  $(\bar{\mathbf{d}}_1, \bar{\mathbf{d}}_2) \in \mathfrak{C}$  then,  $\rho_1(\mathfrak{c}) \sqsubseteq_{\bar{\mathbf{D}}_1} \bar{\mathbf{d}}_1 \wedge \gamma_{\mathfrak{c}}((\rho_1(\mathfrak{c}), \bar{\mathbf{d}}_2)) = \gamma_{\mathfrak{c}}(\mathfrak{c})$  and  $\rho_2(\mathfrak{c}) \sqsubseteq_{\bar{\mathbf{D}}_2} \bar{\mathbf{d}}_2 \wedge \gamma_{\mathfrak{c}}((\bar{\mathbf{d}}_1, \rho_2(\mathfrak{c}))) = \gamma_{\mathfrak{c}}(\mathfrak{c})$ .

## Chapter 3

---

# STRING ANALYSIS FOR C

In this chapter, we present the M-String abstract domain for strings manipulating C code. In particular, M-String is a refinement of the segmentation approach to array content representation proposed by Cousot et al. [57], and it is tailored to detect the presence of common C strings managements errors. An experimental evaluation is also given to prove its effectiveness.

The content of this chapter reports contributions published in [47, 48, 121].

### Chapter Structure

Section 3.1 introduces the problem of handling C strings and explains our contribution. Section 3.2 recalls the array segmentation abstract domain [57] on which M-String is based. Section 3.3 gives the syntax of some operations of interest. Section 3.4 defines the concrete domain and semantics. Section 3.5 presents the M-String abstract domain for C character arrays and its semantics, whose soundness is formally proved. Section 3.6 presents a general approach to abstraction as a program transformation and extends it to abstraction of program strings. Sections 3.7 and 3.8 present implementation and evaluation details of M-String abstraction. Section 3.9 concludes.

### 3.1 Introduction

C is still one of the mainstream programming languages [31], and a large portion of systems of critical relevance is written in this language, e.g., embedded systems. Unfortunately, C programs suffer from bugs due to the way they are laid out in memory, which malicious parties may exploit to drive security attacks. Ensuring the correctness of such software is of great concern. Our main interest is to guarantee the correctness of C programs that manage strings because the incorrect string manipulation may lead to several catastrophic events, ranging from loss or exposure of sensitive data to crashes in critical software components.

Strings in C are not a basic data type. As a matter of fact, strings in C are represented by zero-terminated arrays of characters and there are libraries that provide functions which

allow operating on them [30]. C programs that manipulate strings can suffer from buffer overflows and related issues due to the possible discrepancy between the size of the string and the size of the array (buffer). A buffer overflow is a bug that affects C code when a buffer is accessed out of its bounds. In particular, an out-of-bounds write is a particular (and very dangerous) case of buffer overflow. Out-of-bounds read is less critical as a bug. It is important to design methods supporting the automatic correctness verification of string management in C programs for the previously mentioned reasons and also because buffer overflows are usually exploitable and can easily lead to arbitrary code execution [140]. Existing bugs can be identified by enhancing tools for code analysis, which can also reduce the risk of introducing new bugs and limiting the occurrence of costly security incidents.

### State of the Art

Static methods tailored to identify buffer overflows automatically have been extensively studied in the literature, and several inference techniques were proposed and implemented: tainted data-flow analysis, constraint solvers for various theories (including string theories), and techniques based on them (e.g., symbolic execution), annotation analysis or string pattern matching analysis [156]. Furthermore, the above mentioned techniques and a large number of bug hunting tools based on static analysis had been implemented [62, 63, 97, 129, 175, 181].

For instance, Jones and Kelly [107] introduced a backward compatible method of bounds checking of C programs, which leaves the representation of pointers unchanged, allowing inter-operation between checked and unchecked code, with recompilation confined to the modules where problems might occur. The just mentioned feature differentiates the proposed schema from previously existing techniques. Dor et al. [62] introduced CSSV, the static verifier of C strings. Contracts are supplied to the tool, which acts in 4 stages, reducing the problem of checking code that manipulates string to checking code that manipulates integers. Finally, Splat, described in [183], is a tool that automatically generates test inputs, symbolically reasoning about lengths of input buffers.

Most of the existing string abstract domains are general-purpose domains, focusing on the generic aspects of strings, without accounting for the specifics of string handling by the different programming languages. However, it is often beneficial to consider specific aspects of string representation when designing abstract domains for program analysis. Referring to the C programming language, Journault et al. proposed an abstract domain for C strings which tracks both their length and the buffer allocated size into which they are contained [108]. Combining it with the cell abstraction [134], such domain can describe relations between length of variables and offsets of pointers. Amadini et al. [10] have evaluated several abstract string domains (and their combinations) for analysis of JavaScript programs. In [142] the simplified regular expression domain for JavaScript analysis was defined too. In addition to theoretical work, a number of tools based on the



above mentioned abstract domains and their combinations have been designed and implemented [105, 112, 142, 162]. While dynamic languages heavily rely on strings and their analysis benefits much from tailored abstract domains, the specifics of the `C` approach to strings also earns attention.

## Contribution

We introduce M-String, a new abstract domain tailored for the analysis of strings in `C`. This domain approximates sets of `C` character arrays, allows the abstraction of both shape information on the array structure and value information on the contained characters, and it highlights the presence of well-formed strings in the approximated character arrays.

M-String refines the segmentation approach to array representation introduced in [57]. M-String’s goal is to detect the presence of common string management errors that may lead to undefined behaviours or, more specifically, which may result in buffer overflows. Moreover, keeping track of the content of the characters occurring after the first null character, we reduce the number of false positives. In fact, rewriting the first null character in the array is not always an error, as further occurrences of the null character may follow. Just as the array segmentation-based representation introduced in [57], M-String is parametric in two ways: with respect to both the representation of the indices of the array and with respect to the abstraction of the element values.

To provide evidence of the effectiveness of M-String, we extend LART [120], a tool which performs automatic abstraction on programs, making it supporting also sophisticated (non-scalar) domains such as M-String.

We extend LART along with DIVINE 4 [23], an explicit state model checker based on LLVM. This way, we can verify the correctness of operations on strings in `C` programs automatically. The experimental evaluation is performed by analyzing several `C` programs, ranging from quite simple to moderately complex code, including parsers generated by `bison`, a tool which translates context-free grammars into `C` parsers. The results show the actual impact of an ad-hoc segmentation-based abstract domain on model checking of `C` programs.

## 3.2 FunArray

In the following, we recall the array segmentation analysis presented in [57]. Notice that we slightly modify the notation to be consistent with the whole chapter. For more details, we refer the reader to the original paper.

### 3.2.1 Array Concrete Representation

Below, we recall the definition of the array concrete representation.

**Definition 3.1.** Let  $R_a$  be the set of concrete array environments. A concrete array environment  $\theta \in R_a$  maps array variables  $a \in \mathbf{A}$  to their values  $\theta(a) \in \mathbf{A}$ , such that:

- $\theta(a) = (\rho, l_a, h_a, A_a)$  and,
- $\theta(a) \in \mathbf{A} = R_v \times \mathbf{E} \times \mathbf{E} \times (\mathbb{Z} \rightarrow (\mathbb{Z} \times \mathbf{V}))$

where

1.  $R_v$  is the set of concrete variable environments. A concrete variable environment  $\rho \in R_v$  maps variables (of basic types)  $x \in \mathbf{X}$  to their values  $\rho(x) \in \mathbf{V}$ .
2.  $\mathbf{E}$  is the set of program expressions made up of constants, variables, mathematical unary and binary operators;  $l_a, h_a \in \mathbf{E}$  are expressions whose value,<sup>1</sup> given by  $\llbracket l_a \rrbracket \rho$  and  $\llbracket h_a \rrbracket \rho$ , respectively represents the lower bound and the upper bound of an array  $a$ , i.e., the lower and the upper bound of its indexes range.<sup>2</sup> Note that the value of an upper bound of an array concrete value corresponds to the index immediately after the one that points to the last memory block allocated to the array when it has been initialized. As usual, array indexes are 0-based.
3.  $\mathbb{Z}$  is the set of integer numbers and  $\mathbf{V}$  is the set of values. Let  $\text{idx}(a)$  be the set of indexes  $i$  of an array  $a$ , i.e.,  $\text{idx}(a) = \{i \mid i \in [\llbracket l_a \rrbracket \rho, \llbracket h_a \rrbracket \rho]\} \subseteq \mathbb{Z}$  and let  $\text{pair}(a)$  be the set of pairs  $(i, v)$  such that  $v$  is the value of the element indexed by  $i$  in an array  $a$ , i.e.,  $\text{pair}(a) = \{(i, v) \mid i \in \text{idx}(a) \wedge \llbracket a[i] \rrbracket \rho = v \in \mathbf{V}\} \subseteq \mathbb{Z} \times \mathbf{V}$ .<sup>3</sup> Thus,  $A_a : \text{idx}(a) \rightarrow \text{pair}(a)$  is a function mapping the indexes of an array  $a$  to their corresponding pairs (index, indexed array value).

△

**Example 3.1.** Let  $a$  be a C integer array initialized as follows:  $a[5] = \{5, 7, 9, 11, 13\}$ . The concrete value of  $a$  is given by the tuple  $\theta(a) = (\rho, 0, 5, A_a)$ , where the value of the lower and the upper bound of  $a$  are clear from the context and the codomain of the function  $A_a$  is the set  $\text{pair}(a) = \{(0, 5), (1, 7), (2, 9), (3, 11), (4, 13)\}$ . Moreover, let  $b$  denote the sub-array of  $a$  from position 2 to 3 included, its concrete value is given by  $\theta(b) = (\rho, 2, 4, A_b)$  such that  $\text{pair}(b) = \{(2, 9), (3, 11)\}$ .

◻

---

<sup>1</sup>For simplicity, in the following, expressions are evaluated to integers.

<sup>2</sup>According to the denotational semantics approach [153], in [57] the value of an arithmetic expression  $e$  is denoted by  $\llbracket e \rrbracket \rho$ , where: (1) the double square brackets notation denotes the semantic evaluation function and, (2)  $\rho$  is a variable environment. Typically,  $\llbracket x \rrbracket \rho$  is equivalent to  $\rho(x)$ , with  $x \in \mathbf{X}$ , and  $\llbracket n \rrbracket \rho$ , where  $n$  is a constant, is equivalent to  $n$  itself. Thus, for example, if  $e$  is the expression  $x - 1$ , its semantics  $\llbracket x - 1 \rrbracket \rho$  is defined as  $\llbracket x \rrbracket \rho - \llbracket 1 \rrbracket \rho$ , which corresponds to  $\rho(x) - 1$ .

<sup>3</sup>Note that, in some cases, we use the interval notation to denote numeric sets, e.g., let  $x, y \in \mathbb{Z}$ , the interval between  $x$  and  $y$  (not included) is denoted by  $[x, y)$  which corresponds to the set  $\{w \in \mathbb{Z} \mid x \leq w < y\}$ .

Observe that this array representation allows reasoning about the correspondence between shape components of an array and actual values of the array elements.

### 3.2.2 Array Abstract Domain Functor

According to [57], the FunArray abstract domain  $\overline{\mathbf{S}}$  (shortcut for  $\overline{\mathbf{S}}(\overline{\mathbf{B}}, \overline{\mathbf{A}}, \overline{\mathbf{R}})$ ) allows representing a sequence of consecutive, non-overlapping and possibly empty segments that over-approximate a set of concrete array values in  $\mathcal{P}(\mathbf{A})$ , i.e., the powerset of  $\mathbf{A}$ . Each segment represents a sub-array whose elements share the same property (e.g., being positive integer values) and is surrounded by the so-called segment bounds, i.e., abstractions on its lower and upper bound.

**Example 3.2.** Consider the integer array  $\mathbf{a}[5] = \{5, 7, 9, 10, 12\}$ . As an abstraction of  $\mathbf{a}$  we may consider  $\{0\}$  odd  $\{3\}$  even  $\{5\}$  saying that the array contains odd numbers in the first three elements (indexed from 0 to 2) and two even elements (indexed from 3 to 4).

◻

The elements of FunArray belong to the set  $\overline{\mathbf{S}} = \{(\overline{\mathbf{B}} \times \overline{\mathbf{A}}) \times (\overline{\mathbf{B}} \times \overline{\mathbf{A}} \times \{\_, ?\})^k \times (\overline{\mathbf{B}} \times \{\_, ?\}) \mid k \geq 0\} \cup \{\perp_{\overline{\mathbf{S}}}\}$ , and have the form  $\overline{\mathbf{b}}_1 \overline{\mathbf{p}}_1 \overline{\mathbf{b}}_2 [?_2] \overline{\mathbf{p}}_2 \dots \overline{\mathbf{p}}_{n-1} \overline{\mathbf{b}}_n [?_n]$  where

1.  $\overline{\mathbf{B}}$  is the segment bound abstract domain, approximating array indexes, with abstract properties  $\overline{\mathbf{b}}_i \in \overline{\mathbf{B}}$  such that  $i \in [1, n]$  and  $n > 1$ .

We denote by  $\overline{\mathbf{E}}$  the set of expressions of canonical form  $\overline{\mathbf{x}} + k$ , where  $\overline{\mathbf{x}} \in \overline{\mathbf{X}}$  and  $k \in \mathbb{Z}$ . The segment bounds  $\overline{\mathbf{b}}_i$  are sets of expressions  $\{\overline{\mathbf{e}}_i^1, \dots, \overline{\mathbf{e}}_i^m\}$ , such that  $\overline{\mathbf{e}}_i^j \in \overline{\mathbf{E}}$ . The variable abstract domain  $\overline{\mathbf{X}}$  encodes program variables, i.e.,  $\overline{\mathbf{X}} = \mathbf{X} \cup \{v_0\}$ , where  $v_0$  is a special variable whose value is assumed to be zero. Moreover,  $\overline{\mathbf{b}}_i = \emptyset$  denotes unreachability; if  $\overline{\mathbf{b}}_i \neq \emptyset$ , the expressions appearing in a segment bound are all equivalent symbolic denotations of some concrete value (generally unknown in the abstract representation except when one of the  $\overline{\mathbf{e}}_i^j$  is a constant). Thus,  $\overline{\mathbf{B}}$  depends on the expression abstract domain  $\overline{\mathbf{E}}$ , which, in turn, depends on the variable abstract domain  $\overline{\mathbf{X}}$ .

2.  $\overline{\mathbf{A}}$  is the array element abstract domain, with abstract properties  $\overline{\mathbf{p}}_i \in \overline{\mathbf{A}}$ . It denotes possible values of pairs (index, indexed array element) in a segment, for relational abstractions, array elements otherwise.
3.  $\overline{\mathbf{R}}$  is the variable environment abstract domain, which depends on the variable abstract domain  $\overline{\mathbf{X}}$ , with abstract properties  $\overline{\mathbf{p}} \in \overline{\mathbf{R}}$ .
4. the question mark, if present, expresses that the segment that precedes it may be empty. The question mark can never precede  $\overline{\mathbf{b}}_1$ . The space symbol  $\_$  in  $\{\_, ?\}$  represents a non-empty segment.

**Example 3.3.** Let  $\overline{\mathbf{A}}$  be the classical sign abstraction of numerical values. The segmentation abstract predicate  $\{0\} \text{ pos } \{3\}? \text{ neg } \{5\}$  represents arrays of length 5, with either 0 or 3 positive elements followed by either 5 or 2 negative elements, respectively. For instance, it represents:  $[7, 9, 10, -11, -9]$ ,  $[6, 8, 5, -4, -2]$  and  $[-2, -6, -3, -1, -4, -8]$ . Please note that in the last case, the lack of positive values is justified by the presence of the question mark that says that the first segment is optional.

◻

The *unification* algorithm [57] (cf. Appendix A) modifies two compatible segmentations in order to align them with respect to the same list of bounds.<sup>4</sup> The unification algorithm does not guarantee the maximality of the result, but it is always well-defined, it does terminate, and it is deterministic. The partial order  $\sqsubseteq_{\overline{\mathbf{S}}}$  over  $\overline{\mathbf{S}}$  is defined over unified segmentations as well as the join  $\sqcup_{\overline{\mathbf{S}}}$  and the meet  $\sqcap_{\overline{\mathbf{S}}}$  operators. Please note that  $\overline{\mathbf{S}}$  is not necessarily a lattice [73]. Moreover,  $\overline{\mathbf{S}}$  does not respect the ACC, therefore, in order to ensure the convergence of the analysis, it is equipped with a widening operator  $\nabla_{\overline{\mathbf{S}}}$ . A narrowing operator  $\Delta_{\overline{\mathbf{S}}}$ , which improves the precision of the widening result, is also defined.  $\nabla_{\overline{\mathbf{S}}}$  and  $\Delta_{\overline{\mathbf{S}}}$  operate on unified segmentations.

This abstract array representation is effective for analyzing the content of arrays. In the case of the C programming language, where a string is defined as a null-terminating character array, it is however, not powerful enough to detect common string manipulation errors.

### 3.3 Syntax

Strings in the programming language C are arrays of characters, whose length is determined by the position of a terminating null character '\0'. Thus, for example, the string literal "bee" has four characters: 'b', 'e', 'e', '\0'. Moreover, C supports several string handling functions defined in the standard library `string.h`.

We focus on the most significant functions in the `string.h` header (see Table 3.1), manipulating null-terminated sequences of characters, plus the array elements access and update operations. Recall that `char`, `int` and `size_t` are data types in C, `const` is a qualifier applied to the declaration of any variable which specifies the immutability of its value, and `*str` denotes that `str` is a pointer variable.

- `strcat` appends the null-terminated string pointed to by `str2` to the null-terminated string pointed to by `str1`. The first character of `str2` overwrites the null-terminator of `str1` and `str2` should not overlap `str1`. The string concatenation returns the pointer `str1`.

---

<sup>4</sup>Two segmentations,  $\overline{\mathbf{b}}_1^1 \dots \overline{\mathbf{b}}_n^1[?_n^1]$  and  $\overline{\mathbf{b}}_1^2 \dots \overline{\mathbf{b}}_n^2[?_n^2]$ , are compatible if  $\overline{\mathbf{b}}_1^1 \cap \overline{\mathbf{b}}_1^2 \neq \emptyset$  and  $\overline{\mathbf{b}}_n^1 \cap \overline{\mathbf{b}}_n^2 \neq \emptyset$ .

<code>char *strcat(char *str1, const char *str2)</code>
<code>char *strchr(char *str, int c)</code>
<code>int strcmp(const char *str1, const char *str2)</code>
<code>char *strcpy(char *str1, const char *str2)</code>
<code>size_t strlen(const char *str)</code>

TABLE 3.1: String functions syntax in C.

- `strchr` locates the first occurrence of `c` (converted to a `char`) in the string pointed to by `str`. The terminating null character is considered to be part of the string. The string character function returns a pointer to the located character, or a null pointer if the character does not occur in the string.
- `strcmp` lexicographically compares the string pointed to by `str1` to the string pointed to by `str2`. The string compare function returns an integer greater than (equal to, or less than) zero when the string pointed to by `str1` is greater than (equal to, or less than) the string pointed to by `str2`.
- `strcpy` copies the null-terminated string pointed to by `str2` to the memory pointed to by `str1`. `str2` should not overlap `str1`. The string copy function returns the pointer `str1`.
- `strlen` computes the number of bytes in the string to which `str` points, not including the terminating null byte. The string length function returns the length of `str`.

Accessing an array element is by indexing the array name. Let `i` be an index, the `i`-th element of the character array `str` is accessed by `str[i]`. Instead, a character array element is updated (or an assignment is performed to a character array element) by `str[i] = 'x'`, where `'x'` denotes a character literal.

As mentioned in Section 3.1, C does not guarantee bounds checking on array accesses and, in case of strings, the language does not ensure that the latter are null-terminated. As a consequence, improper string manipulation leads to several vulnerabilities and exploits [154]. For instance, if non null-terminated strings are passed to the functions above, the latter may return misleading results or read out of the array bound. Moreover, since `strcat` and `strcpy` do not allow the size of the destination array `str1` to be specified, they are frequent sources of buffer overflows.

### 3.4 Concrete Domain and Semantics

We aim to capture the presence of well-formed strings in C character arrays, to avoid undesired execution behaviours that may be security relevant. To reach our goal, we

propose a character array concrete value that highlights the occurrence of null characters in it, and we introduce the notion of *string of interest* of an array of chars. The concrete semantics of the operations presented in Section 3.3 is also given.

### 3.4.1 Character Array Concrete Representation

Let  $\mathbf{C}$  be the finite set of characters representable by the encoding in use, equipped with a top element  $\top_{\mathbf{C}}$  representing an unknown value; and let  $\mathbf{M}$  be the set of character array variables, such that  $\mathbf{M} \subseteq \mathbf{A}$  (with  $\mathbf{A}$  being the set of array variables - of any type - presented in Section 3.2.2). Then, the operational semantics of character array variables are concrete array environments  $\mu \in \mathbf{R}_{\mathbf{m}}$ . The semantics maps character arrays  $\mathbf{m} \in \mathbf{M}$  to their values  $\mu(\mathbf{m})$ . Precisely:

- $\mu(\mathbf{m}) = (\rho, l_{\mathbf{m}}, h_{\mathbf{m}}, M_{\mathbf{m}}, N_{\mathbf{m}})$  and,
- $\mu(\mathbf{m}) \in \mathbf{M} = \mathbf{R}_v \times \mathbf{E} \times \mathbf{E} \times (\mathbb{Z} \rightarrow (\mathbb{Z} \times \mathbf{C})) \times \mathcal{P}(\mathbb{Z})$

so that  $\mathbf{R}_{\mathbf{m}}$  is a map from  $\mathbf{M}$  to  $\mathbf{M}$ , where  $\mathbf{R}_v$  and  $\mathbf{E}$  are the concrete variable environment and the expression domain of Definition 3.1 respectively,  $\mathbb{Z}$  is the integer domain and  $\mathbf{C}$  is the character set introduced above. Note that with respect to the concrete array environment  $\theta$  introduced in Definition 3.1, the function  $\mu$  returns as a last component the set of indexes which map to the string terminating characters  $N_{\mathbf{m}} = \{i \mid i \in \text{idx}(\mathbf{m}) \wedge M_{\mathbf{m}}(i) = (i, '\backslash 0')\}$ , with  $\text{idx}(\mathbf{m})$  being the domain of the function  $M_{\mathbf{m}}$ . Instead,  $M_{\mathbf{m}}$  behaves exactly as  $A_{\mathbf{a}}$  in  $\theta(\mathbf{a})$ , mapping each index  $i$  of the considered array to the pair of the index  $i$  and the indexed array element  $v$ .

Thus,  $\mathbf{M}$  extends  $\mathbf{A}$  (cf. Section 3.2.1) by adding a parameter that takes into account the presence of null characters in a character array. For well-formed strings,  $N_{\mathbf{m}}$  is not empty. Moreover, character array elements that have not been initialized are mapped to the top value  $\top_{\mathbf{C}}$  as they may be values already present in the memory assigned to the locations array itself.

**Example 3.4.** Let  $\mathbf{m}$  be a  $\mathbf{C}$  character array initialized as follows:  $\mathbf{m}[6] = \{'b', 'e', 'e', '\backslash 0', 'b'\}$ . The concrete value of  $\mathbf{m}$  is given by the tuple  $\mu(\mathbf{m}) = (\rho, 0, 6, M_{\mathbf{m}}, N_{\mathbf{m}})$ , where the codomain of the function  $M_{\mathbf{m}}$  is the set  $\text{pair}(\mathbf{m}) = \{(0, 'b'), (1, 'e'), (2, 'e'), (3, '\backslash 0'), (4, 'b'), (5, \top_{\mathbf{C}})\}$  and  $N_{\mathbf{m}}$  is the singleton  $\{3\}$ , being the array cell of index 3 the only one certainly containing a null character.

◻

### String of Interest

We formally define the string of interest of a character array as the sequence of its elements up to the first terminating one (included).

**Definition 3.2** (String of interest). Let  $\mathbf{m} \in \mathbf{M}$  be an array of characters with concrete value  $\mu(\mathbf{m}) = (\rho, l_{\mathbf{m}}, h_{\mathbf{m}}, M_{\mathbf{m}}, N_{\mathbf{m}})$ . The string of interest of the character array described by  $\mu(\mathbf{m})$  is defined as follows:

$$\text{strInt}(\mu(\mathbf{m})) = \begin{cases} \langle v_i : i \in [l_{\mathbf{m}}\rho, \min(N_{\mathbf{m}})] \wedge M_{\mathbf{m}}(i) = (i, v) \rangle & \text{if } N_{\mathbf{m}} \neq \emptyset \\ \text{undef} & \text{otherwise} \end{cases}$$

with  $v_i$  denoting the character value which occurs in the pair  $(i, v)$ , and with  $\min(N_{\mathbf{m}})$  denoting the minimum element of  $N_{\mathbf{m}}$ .

△

**Example 3.5.** Consider the concrete character array value of Example 3.4. Its string of interest is the sequence of characters "bee\0".

◻

Our definition of string of interest of character arrays allows us to distinguish well-formed strings and avoid bad usage of arrays of characters. If the null character appears at the first index of a character array, then we refer to its string of interest as null (`null`). In general, we refer to character arrays which contain a well-defined or null string of interest as character arrays that contain a *well-formed string*.

Moreover, when allocated memory capacity is not sufficient for a declared character array, the system writes a null character outside the array, occupying memory that is not destined for it and causing a buffer overflow. We do not represent this system behaviour since it leads to an undefined one, so we consider the string of interest of such character arrays as undefined (`undef`).

### 3.4.2 Concrete Domain

As a concrete domain for array of characters we refer to the complete lattice  $\mathcal{P}(\mathbf{M})$  defined as  $(\mathcal{P}(\mathbf{M}), \subseteq_{\mathcal{P}(\mathbf{M})}, \perp_{\mathcal{P}(\mathbf{M})}, \top_{\mathcal{P}(\mathbf{M})}, \cup_{\mathcal{P}(\mathbf{M})}, \cap_{\mathcal{P}(\mathbf{M})})$  where:  $\mathcal{P}(\mathbf{M})$  is the powerset of concrete character array values, the set inclusion  $\subseteq_{\mathcal{P}(\mathbf{M})}$  corresponds to the partial order, the bottom element  $\perp_{\mathcal{P}(\mathbf{M})}$  is the emptyset  $\emptyset$ , the top element  $\top_{\mathcal{P}(\mathbf{M})}$  is the superset of any subset of  $\mathbf{M}$  (i.e.,  $\mathbf{M}$  itself), the set union  $\cup_{\mathcal{P}(\mathbf{M})}$  denotes the least upper bound and, the set intersection  $\cap_{\mathcal{P}(\mathbf{M})}$  denotes the greatest lower bound.

We stress that our concrete domain is used as a framework to help us create the abstract representation, and it is not how the (concrete) character array values are actually represented in `C` programs.

### 3.4.3 Concrete Semantics

To formalize the concrete semantics of the `C` standard library functions from `string.h` introduced in Section 3.3, we define the following auxiliary functions embedding (`emb`),

extraction (**ext**), comparison (**cmp**) and substitution (**sub**) over single concrete character array values.

**Definition 3.3** (Embedding). Let  $\mu(m_1), \mu(m_2) \in \mathbf{M}$  be two concrete character array values and let  $[l_1, u_1] \subseteq [\llbracket l_{m_1} \rrbracket \rho, \llbracket h_{m_1} \rrbracket \rho]$ ,  $[l_2, u_2] \subseteq [\llbracket l_{m_2} \rrbracket \rho, \llbracket h_{m_2} \rrbracket \rho]$  be two indexes ranges of the same length. The function  $\text{emb}(\mu(m_1), [l_1, u_1], \mu(m_2), [l_2, u_2])$  embeds the sequence of characters of  $\mu(m_2)$  which occurs from the index  $l_2$  to the index  $u_2$  into  $\mu(m_1)$  from the index  $l_1$  to the index  $u_1$ . Formally,  $\text{emb}(\mu(m_1), [l_1, u_1], \mu(m_2), [l_2, u_2]) = \mu(m_1)'$  such that:

- $\llbracket l_{m_1}' \rrbracket \rho = \llbracket l_{m_1} \rrbracket \rho$  and  $\llbracket h_{m_1}' \rrbracket \rho = \llbracket h_{m_1} \rrbracket \rho$
- $M_{m_1}' :$ 
  - $\forall i \in [\llbracket l_{m_1}' \rrbracket \rho, l_1]: M_{m_1}'(i) = (i, v)$  such that  $k = i \wedge M_{m_1}(k) = (k, v)$
  - $\forall i \in [l_1, u_1]: M_{m_1}'(i) = (i, v)$  such that  $k = l_2 + (i - l_1) \wedge M_{m_2}(k) = (k, v)$
  - $\forall i \in (l_1, \llbracket h_{m_1}' \rrbracket \rho): M_{m_1}'(i) = (i, v)$  such that  $k = i \wedge M_{m_1}(k) = (k, v)$
- $N_{m_1}' = (N_{m_1} \setminus \{i \mid i \in [l_1, u_1]\}) \cup \{i \mid i \in [l_1, u_1] \wedge k = l_2 + (i - l_1) \wedge M_{m_2}(k) = (k, '\backslash 0')\}$

△

**Example 3.6.** Let  $\mu(m_1) = (\rho, 0, 7, M_{m_1}, N_{m_1})$  and  $\mu(m_2) = (\rho, 0, 6, M_{m_2}, N_{m_2})$  be two concrete character array values such that:

- $\text{pair}(m_1) = \{(0, 'a'), (1, 'a'), (2, 'a'), (3, '\backslash 0'), (4, 'a'), (5, 'a'), (6, 'a')\}$
- $\text{pair}(m_2) = \{(0, 'b'), (1, 'b'), (2, 'b'), (3, 'b'), (4, 'b'), (5, '\backslash 0')\}$

Moreover, consider the intervals of equal length:

- $[2, 4]_{m_1} \subseteq [\llbracket l_{m_1} \rrbracket \rho, \llbracket h_{m_1} \rrbracket \rho]$
- $[3, 5]_{m_2} \subseteq [\llbracket l_{m_2} \rrbracket \rho, \llbracket h_{m_2} \rrbracket \rho]$

The function  $\text{emb}(\mu(m_1), [2, 4]_{m_1}, \mu(m_2), [3, 5]_{m_2}) = \mu(m_1)'$  where:

- $\llbracket l_{m_1}' \rrbracket \rho = \llbracket l_{m_1} \rrbracket \rho = 0$  and  $\llbracket h_{m_1}' \rrbracket \rho = \llbracket h_{m_1} \rrbracket \rho = 7$
- $\text{pair}(m_1') = \{(0, 'a'), (1, 'a'), (2, 'b'), (3, 'b'), (4, '\backslash 0'), (5, 'a'), (6, 'a')\}$
- $N_{m_1}' = \{4\}$

◻

**Definition 3.4** (Extraction). Let  $\mu(m) \in \mathbf{M}$  be a concrete character array value and let  $[l, u] \subseteq [\llbracket l_m \rrbracket \rho, \llbracket h_m \rrbracket \rho]$  be an indexes range. The function  $\text{ext}(\mu(m), [l, u])$  extracts the sequence of characters which occurs in  $\mu(m)$  from the index  $l$  to the index  $u$ . Formally,  $\text{ext}(\mu(m), [l, u]) = \mu(m)'$  such that:



- $\llbracket l_{m'} \rrbracket \rho = 1$  and  $\llbracket h_{m'} \rrbracket \rho = u + 1$
- $M_{m'} : \forall i \in [\llbracket l_{m'} \rrbracket \rho, \llbracket h_{m'} \rrbracket \rho): M_{m'}(i) = (i, v)$  such that  $k = i \wedge M_m(k) = (k, v)$
- $N_{m'} = N_m \setminus \{i \mid i \notin [l, u]\}$

△

**Example 3.7.** Let  $\mu(m_1)$  be the character array concrete value of Example 3.6 and  $[1, 3]_{m_1} \subseteq [\llbracket l_{m_1} \rrbracket \rho, \llbracket h_{m_1} \rrbracket \rho)$  be an indexes range of  $\mu(m_1)$ . The function  $\text{ext}(\mu(m_1), [1, 3]_{m_1}) = \mu(m_1)'$  such that:

- $\llbracket l_{m_1'} \rrbracket \rho = 1$  and  $\llbracket h_{m_1'} \rrbracket \rho = 4$
- $\text{pair}(m_1') = \{(1, 'a'), (2, 'a'), (3, '\0')\}$
- $N_{m_1'} = \{3\}$

◇

**Definition 3.5** (Comparison). Let  $\mu(m_1), \mu(m_2) \in \mathbf{M}$  be two concrete character array values which contain a fully initialized well-formed string of interest, i.e., no  $\top_{\mathbf{C}}$  occurs. The function  $\text{cmp}(\mu(m_1), \mu(m_2))$  (cf. Algorithm 1) lexicographically compares the strings of interest of  $\mu(m_1)$  and  $\mu(m_2)$  and it returns an integer value  $n$  which denotes the lexicographic distance between them.

△

**Example 3.8.** Let  $\mu(m_1)$  and  $\mu(m_2)$  be the character array concrete values of Example 3.6. Both of them contain a fully initialized well-formed string of interest and the function  $\text{cmp}(\mu(m_1), \mu(m_2))$  computes the lexicographic distance between them. Precisely, the procedure stops after the first iteration of the for loop (cf. Algorithm 1) and, assuming ASCII as the character encoding set, it returns the value  $-1$ , i.e.,  $n = 97 - 98$ , which means that  $\text{strInt}(\mu(m_1))$  lexicographically precedes  $\text{strInt}(\mu(m_2))$ .

◇

**Definition 3.6** (Substitution). Let  $\mu(m) \in \mathbf{M}$  be a concrete character array value, let  $z \in [\llbracket l_m \rrbracket \rho, \llbracket h_m \rrbracket \rho)$  be an index, and let  $c \in \mathbf{C}$  be a character. The function  $\text{sub}(\mu(m), z, c)$  substitutes the character which appears in  $\mu(m)$  at the index  $z$  with the character  $c$ . Formally,  $\text{sub}(\mu(m), z, c) = \mu(m)'$  such that:

- $\llbracket l_{m'} \rrbracket \rho = \llbracket l_m \rrbracket \rho$  and  $\llbracket h_{m'} \rrbracket \rho = \llbracket h_m \rrbracket \rho$
- $M_{m'} :$ 
  - $\forall i \in [\llbracket l_{m'} \rrbracket \rho, z): M_{m'}(i) = (i, v)$  such that  $k = i \wedge M_m(k) = (k, v)$

---

**Algorithm 1** Lexicographic comparison of concrete character array values.

---

**Input:** two concrete character array values  $\mu(m_1), \mu(m_2) \in \mathbf{M}$  such that:

- both  $N_{m_1}$  and  $N_{m_2}$  are different from the emptyset and,
- for  $i_1 \in \llbracket l_{m_1} \rrbracket \rho, \min(N_{m_1}), i_2 \in \llbracket l_{m_2} \rrbracket \rho, \min(N_{m_2})$ :
  - $M(i_1) \neq (i_1, \top_{\mathbf{C}})$
  - $M(i_2) \neq (i_2, \top_{\mathbf{C}})$

**Output:** an integer value  $n$ .

```

1:  $n = 0, i_1 = \llbracket l_{m_1} \rrbracket \rho, i_2 = \llbracket l_{m_2} \rrbracket \rho$ 
2: while  $i_1 \in \llbracket l_{m_1} \rrbracket \rho, \min(N_{m_1}) \wedge i_2 \in \llbracket l_{m_2} \rrbracket \rho, \min(N_{m_2})$  do
3:    $n = v_{i_1} - v_{i_2}$ 
4:   if  $n \neq 0$  then
5:     return  $n$ 
6:   else
7:      $i_1 = i_1 + 1, i_2 = i_2 + 1$ 
8: return  $n$ 

```

---

– for  $i = z$ :  $M_{m'}(z) = (z, c)$

–  $\forall i \in (z, \llbracket h_{m'} \rrbracket \rho)$ :  $M_{m'}(i) = (i, v)$  such that  $k = i \wedge M_m(k) = (k, v)$

$$\bullet N_{m'} = \begin{cases} N_m & \text{if } (z \in N_m \wedge c \text{ is null}) \vee (z \notin N_m \wedge c \text{ is not null}) \\ N_m \setminus \{z\} & \text{if } z \in N_m \wedge c \text{ is not null} \\ N_m \cup \{z\} & \text{otherwise} \end{cases}$$

△

**Example 3.9.** Let  $\mu(m_1)$  be the character array concrete value of Example 3.6, the index  $z$  be equal to 4 and the character  $c$  be the null termination '\0'. The function  $\text{sub}(\mu(m_1), 4, '\0') = \mu(m_1)'$  such that:

- $\llbracket l_{m_1'} \rrbracket \rho = 0$  and  $\llbracket h_{m_1'} \rrbracket \rho = 7$
- $\text{pair}(m_1') = \{(0, 'a'), (1, 'a'), (2, 'a'), (3, '\0'), (4, '\0'), (5, 'a'), (6, 'a')\}$
- $N_{m_1'} = \{3, 4\}$

◻

### Array Access

The semantic operator  $\mathfrak{A}$ , given the statement `accessj` and a set of concrete character array values  $X$  in  $\mathcal{P}(\mathbf{M})$  as parameter, returns a value in  $\mathbf{C}$ . In particular, `accessj(X)` returns the character  $v$  which occurs at position  $j$  if all the character array values in  $X$

contain  $v$  at index  $j$  and the latter is well-defined (i.e., it ranges in the array bounds) for all the character array values in  $X$ ; otherwise it returns  $\top_{\mathbf{C}}$ . Formally,

$$\mathfrak{M}\llbracket\text{access}_j\rrbracket(X) = \begin{cases} v & \text{if } \forall \mu(\mathbf{m}) \in X : j \in \llbracket\mathbf{l}_m\rrbracket\rho, \llbracket\mathbf{h}_m\rrbracket\rho \text{ and } M_m(j) = (j, v) \\ \top_{\mathbf{C}} & \text{otherwise} \end{cases}$$

### String Concatenation

The semantic operator  $\mathfrak{M}$ , given a statement and some sets of concrete character array values in  $\mathcal{P}(\mathbf{M})$  as parameters, returns a set of concrete character array values. When applied to  $\text{strcat}(X_1, X_2)$ , it returns all the possible embeddings in  $X_1$  of a string of interest taken from  $X_2$  if all the character array values (which belong to both  $X_1$  and  $X_2$ ) contain a well-formed string and the condition on the size of the destination character array values is fulfilled; otherwise it returns  $\top_{\mathcal{P}(\mathbf{M})}$ . Please note that the size condition ensures to perform the string concatenation only if the destination character array value is big enough to contain the string of interest of the source character array value, thus preventing undefined behaviours. Formally,

$$\mathfrak{M}\llbracket\text{strcat}\rrbracket(X_1, X_2) = \begin{cases} X_1' & \text{if } \forall \mu(\mathbf{m}_1) \in X_1 : \forall \mu(\mathbf{m}_2) \in X_2 : \\ & \text{strInt}(\mu(\mathbf{m}_1)) \neq \text{undef} \neq \text{strInt}(\mu(\mathbf{m}_2)) \wedge \\ & \text{size.condition is true} \\ \top_{\mathcal{P}(\mathbf{M})} & \text{otherwise} \end{cases}$$

The `size.condition is true` if:

$$(\llbracket\mathbf{h}_{m_1}\rrbracket\rho - \llbracket\mathbf{l}_{m_1}\rrbracket\rho) \geq [(\min(N_{m_1}) - \llbracket\mathbf{l}_{m_1}\rrbracket\rho - 1) + (\min(N_{m_2}) - \llbracket\mathbf{l}_{m_2}\rrbracket\rho)]$$

Moreover,  $X_1'$  is the set of  $\text{emb}(\mu(\mathbf{m}_1), [l_1, u_1], \mu(\mathbf{m}_2), [l_2, u_2])$  (cf. Definition 3.3), such that:

- $\mu(\mathbf{m}_1) \in X_1$ ,  $l_1 = \min(N_{m_1})$  and  $u_1 = l_1 + (\min(N_{m_2}) - \llbracket\mathbf{l}_{m_2}\rrbracket\rho)$
- $\mu(\mathbf{m}_2) \in X_2$ ,  $l_2 = \llbracket\mathbf{l}_{m_2}\rrbracket\rho$  and  $u_2 = \min(N_{m_2})$

### String Character

The semantic operator  $\mathfrak{M}$ , when applied to  $\text{strchr}_v(X)$ , returns the set of string of interest suffixes in  $X$  from the index corresponding to the first occurrence of the character  $v$  if all the character array values in  $X$  contain a well-formed string containing  $v$ . Otherwise, if all the character array values in  $X$  contain a well-formed string in which does not occur the character  $v$ , it returns the emptyset (denoted by  $\perp_{\mathcal{P}(\mathbf{M})}$ ); otherwise it returns  $\top_{\mathcal{P}(\mathbf{M})}$ . Formally,

$$\mathfrak{M}\llbracket\text{strchr}_v\rrbracket(X) = \begin{cases} X' & \text{if } \forall \mu(\mathbf{m}) \in X : \text{strInt}(\mu(\mathbf{m})) \neq \text{undef} \text{ and } v \in \text{strInt}(\mu(\mathbf{m})) \\ \perp_{\mathcal{P}(\mathbf{M})} & \text{if } \forall \mu(\mathbf{m}) \in X : \text{strInt}(\mu(\mathbf{m})) \neq \text{undef} \text{ and } v \notin \text{strInt}(\mu(\mathbf{m})) \\ \top_{\mathcal{P}(\mathbf{M})} & \text{otherwise} \end{cases}$$

In particular,  $X'$  is the set of  $\text{ext}(\mu(\mathbf{m}), [l, u])$  (cf. Definition 3.4), such that:

- $\mu(\mathbf{m}) \in X$ ,  $l = \min(\{i : i \in \llbracket l_{\mathbf{m}} \rrbracket \rho, \min(N_{\mathbf{m}}) \wedge M_{\mathbf{m}}(i) = (i, v)\})$  and  $u = \min(N_{\mathbf{m}})$

### String Compare

The semantic operator  $\mathfrak{P}$ , given the statement `strcmp` and two sets of concrete character array values  $X_1, X_2$  in  $\mathcal{P}(\mathbf{M})$  as parameters, returns a value in the set of integers equipped with a top element, i.e.,  $\mathbb{Z} \cup \top_{\mathbb{Z}}$ . In particular, `strcmp`( $X_1, X_2$ ) returns an integer value  $n$  which denotes the lexicographic distance between strings of interest in  $X_1$  and  $X_2$  if for all  $\mu(\mathbf{m}_1) \in X_1$  and  $\mu(\mathbf{m}_2) \in X_2$  the procedure `cmp`( $\mu(\mathbf{m}_1), \mu(\mathbf{m}_2)$ ) (cf. Definition 3.5) returns  $n$ ; otherwise it returns  $\top_{\mathbb{Z}}$ .

Note that if  $n$  is negative, the string of interest of  $\mathbf{m}_1$  precedes the string of interest of  $\mathbf{m}_2$  in lexicographic order. Conversely, if  $n$  is positive, the string of interest of  $\mathbf{m}_1$  follows the string of interest of  $\mathbf{m}_2$  in lexicographic order, and if  $n$  is equal to zero they are lexicographically equal. Formally,

$$\mathfrak{P}\llbracket\text{strcmp}\rrbracket(X_1, X_2) = \begin{cases} n & \text{if } \forall \mu(\mathbf{m}_1) \in X_1 : \forall \mu(\mathbf{m}_2) \in X_2 : \text{cmp}(\mu(\mathbf{m}_1), \mu(\mathbf{m}_2)) = n \\ \top_{\mathbb{Z}} & \text{otherwise} \end{cases}$$

### String Copy

The semantic operator  $\mathfrak{M}$ , when applied to `strcpy`( $X_1, X_2$ ), behaves similarly to the string concatenation function above. Formally,

$$\mathfrak{M}\llbracket\text{strcpy}\rrbracket(X_1, X_2) = \begin{cases} X_1' & \text{if } \forall \mu(\mathbf{m}_1) \in X_1 : \forall \mu(\mathbf{m}_2) \in X_2 : \\ & \text{strInt}(\mu(\mathbf{m}_1)) \neq \text{undef} \neq \text{strInt}(\mu(\mathbf{m}_2)) \wedge \\ & \text{size.condition is true} \\ \top_{\mathcal{P}(\mathbf{M})} & \text{otherwise} \end{cases}$$

The `size.condition` is true if:

$$(\llbracket l_{\mathbf{m}_1} \rrbracket \rho - \llbracket l_{\mathbf{m}_1} \rrbracket \rho) \geq (\min(N_{\mathbf{m}_2}) - \llbracket l_{\mathbf{m}_2} \rrbracket \rho)$$

Moreover,  $X_1'$  is the set of  $\text{emb}(\mu(\mathbf{m}_1), [l_1, u_1], \mu(\mathbf{m}_2), [l_2, u_2])$ , such that:

- $\mu(\mathbf{m}_1) \in X_1$ ,  $l_1 = \llbracket l_{\mathbf{m}_1} \rrbracket \rho$  and  $u_1 = l_1 + (\min(N_{\mathbf{m}_2}) - \llbracket l_{\mathbf{m}_2} \rrbracket \rho)$
- $\mu(\mathbf{m}_2) \in X_2$ ,  $l_2 = \llbracket l_{\mathbf{m}_2} \rrbracket \rho$  and  $u_2 = \min(N_{\mathbf{m}_2})$

### String Length

The semantic operator  $\mathfrak{L}$ , given the statement `strlen` and a set of concrete character array values  $X$  in  $\mathcal{P}(\mathbf{M})$  as parameter, returns a value in the set of integers equipped with a top element, i.e.,  $\mathbb{Z} \cup \top_{\mathbb{Z}}$ . In particular, `strlen`( $X$ ) returns an integer value  $n$  which corresponds to the length of the sequence of characters before the first null one of the character arrays values in  $X$  if all the character array values in  $X$  contain a well-formed string of the same length; otherwise it returns  $\top_{\mathbb{Z}}$ . Formally,

$$\mathfrak{L}\llbracket\text{strlen}\rrbracket(X) = \begin{cases} n & \text{if } \forall \mu(\mathbf{m}) \in X : \text{strInt}(\mu(\mathbf{m})) \neq \text{undef} \wedge (\min(N_{\mathbf{m}}) - \llbracket\mathbf{l}_m\rrbracket\rho) = n \\ \top_{\mathbb{Z}} & \text{otherwise} \end{cases}$$

### Array Update

The semantic operator  $\mathfrak{M}$ , when applied to `updatej,v`( $X$ ), returns the set of character array values in  $X$  where the character that occurs at position  $j$  has been substituted with the character  $v$  if the index  $j$  is well-defined for all the character array values in  $X$ ; otherwise it returns  $\top_{\mathcal{P}(\mathbf{M})}$ . Formally,

$$\mathfrak{M}\llbracket\text{update}_{j,v}\rrbracket(X) = \begin{cases} X' & \text{if } \forall \mu(\mathbf{m}) \in X : j \in \llbracket\mathbf{l}_m\rrbracket\rho, \llbracket\mathbf{h}_m\rrbracket\rho \\ \top_{\mathcal{P}(\mathbf{M})} & \text{otherwise} \end{cases}$$

In particular,  $X'$  is the set of `sub`( $\mu(\mathbf{m}), j, v$ ) (cf. Definition 3.6).

## 3.5 M-String

In the previous section, we defined the concrete value of a character array, highlighting the presence of a well-formed string in it. Moreover, we presented our concrete domain  $\mathcal{P}(\mathbf{M})$ , made of sets of character array values, and its concrete semantics of some operations of interest. In the following, we formalize the M-String abstract domain, which approximates elements in  $\mathcal{P}(\mathbf{M})$ , and its semantics for which soundness is proved.

### 3.5.1 Character Array Abstract Domain Functor

The M-String ( $\overline{\mathbf{M}}$ ) abstract domain approximates sets of concrete character array values with a pair of segmentations that highlights the nature of their strings of interest. The elements of the domain are split segmentation abstract predicates. As for `FunArray` (recalled in Section 3.2.2), segments represent sequences of characters which share the same property and are delimited by the so-called segment bounds. More precisely, the M-String abstract domain is a functor given by  $\overline{\mathbf{M}}(\overline{\mathbf{B}}, \overline{\mathbf{C}}, \overline{\mathbf{R}})$  where:

1.  $\overline{\mathbf{B}}$  denotes the abstraction of segment bounds, equipped with the addition ( $\oplus_{\overline{\mathbf{B}}}$ ) and subtraction ( $\ominus_{\overline{\mathbf{B}}}$ ) operations.
2.  $\overline{\mathbf{C}}$  is the abstraction of the character array elements, it is signed, it contains the value 0, and it is equipped with `isNull`, a special monotonic function lifting abstract elements in  $\overline{\mathbf{C}}$  to a value in the set `{true, false, maybe}` and with subtraction ( $\ominus_{\overline{\mathbf{C}}}$ ).
3.  $\overline{\mathbf{R}}$  denotes the abstraction of scalar variable environments (cf. Section 3.2.2). Namely, the constant propagation domain on the set of variables  $\mathbf{X}$ .

The elements of M-String belong to the set  $\overline{\mathbf{M}} \triangleq (\overline{\mathbf{M}}_s, \overline{\mathbf{M}}_{ns}) \cup \{\perp_{\overline{\mathbf{M}}}, \top_{\overline{\mathbf{M}}}\}$  where:

- $\overline{\mathbf{M}}_s$  is  $\{\{\overline{\mathbf{B}} \times \overline{\mathbf{C}}\} \times \{\overline{\mathbf{B}} \times \overline{\mathbf{C}} \times \{-, ?\}\}^k \times \{\overline{\mathbf{B}} \times \{-, ?\}\} \mid k \geq 0\} \cup \{\overline{\mathbf{B}}\} \cup \{\emptyset\}$  and it represents the segmentation of the strings of interest of a set of character arrays.
- $\overline{\mathbf{M}}_{ns}$  is  $\{\{\overline{\mathbf{B}} \times \overline{\mathbf{C}}\} \times \{\overline{\mathbf{B}} \times \overline{\mathbf{C}} \times \{-, ?\}\}^k \times \{\overline{\mathbf{B}} \times \{-, ?\}\} \mid k \geq 0\} \cup \{\emptyset\}$  and it represents the segmentation of the content of character arrays after their string of interests, or character arrays that do not contain the null terminating character.
- $\perp_{\overline{\mathbf{M}}}, \top_{\overline{\mathbf{M}}}$  are special elements denoting the bottom/top element of  $\overline{\mathbf{M}}$ .

The elements in  $\overline{\mathbf{M}}$  are split segmentation abstract predicates of the form  $\overline{\mathbf{m}} = (s, ns)$ . For instance, when  $\overline{\mathbf{m}}$  is equal to  $(\overline{\mathbf{b}}_1, \emptyset)$ , it abstracts concrete character array values of length 1 and containing a `null` string of interest (cf. Section 3.4.1). Instead, when  $\overline{\mathbf{m}}$  is equal to  $(\overline{\mathbf{b}}_1, \overline{\mathbf{b}}_2 \overline{\mathbf{p}}_2 \overline{\mathbf{b}}_3[?_3] \dots \overline{\mathbf{b}}_n[?_n])$ , it approximates concrete character array values of length greater than or equal to 1 containing a `null` string of interest. In particular:

1.  $\overline{\mathbf{b}}_i \in \overline{\mathbf{B}}$  denotes the segment bounds, chosen in abstract domain  $\overline{\mathbf{B}}$ , such that  $i \in [1, n]$  and  $n > 1$ . A segment bound approximates a set of indexes (i.e., positive integers  $\mathbb{Z}^+$ ), but contrary to what defined for the `FunArray` abstraction, the choice of  $\overline{\mathbf{B}}$  is unconstrained.

For the sake of readability, we apply arithmetic operators on  $\overline{\mathbf{b}}_i$  directly. For instance,  $\overline{\mathbf{b}} \oplus_{\overline{\mathbf{B}}} 1$  should be read as  $\alpha_{\overline{\mathbf{B}}}(\{i + 1 \mid i \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}})\})$  or  $\overline{\mathbf{b}}_1 \oplus_{\overline{\mathbf{B}}} \overline{\mathbf{b}}_2$  as  $\alpha_{\overline{\mathbf{B}}}(\{i_1 + i_2 \mid i_1 \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_1) \wedge i_2 \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_2)\})$ , where  $\alpha_{\overline{\mathbf{B}}}$  and  $\gamma_{\overline{\mathbf{B}}}$  are respectively the abstraction and concretization functions over the bounds abstract domain.

Please note that  $\overline{\mathbf{b}}_1$  and  $\overline{\mathbf{b}}_n$ , respectively, represent the segmentation lower and upper bound and in the case in which  $\overline{\mathbf{m}}$  corresponds to the split segmentation  $(\overline{\mathbf{b}}_1 \overline{\mathbf{p}}_1 \overline{\mathbf{b}}_2[?_2] \dots \overline{\mathbf{b}}_{n-1}[?_{n-1}], \emptyset)$  the segmentation upper bound is hidden, due to a representative choice, and equal to  $\overline{\mathbf{b}}_{n-1} \oplus_{\overline{\mathbf{B}}} 1$ . Moreover, in a segmentation  $\dots \overline{\mathbf{b}}_i[?_i] \overline{\mathbf{p}}_{i+1} \overline{\mathbf{b}}_{i+1}[?_{i+1}] \dots$  we always assume that  $\min(\gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_{i+1})) > \max(\gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_i))$ .

2.  $\overline{\mathbf{p}}_i \in \overline{\mathbf{C}}$  are abstract predicates, chosen in an abstract domain  $\overline{\mathbf{C}}$ , denoting possible values of pairs (index, character array element value) in a segment, for relational abstraction, character array elements otherwise.

3. the question mark ?, if present, indicates that the preceding segment might be empty, while  $\_$  indicates a non-empty segment, and, as for [57], non-empty segments are not marked.

**Example 3.10.** Consider the split segmentation abstract predicate  $\bar{\mathbf{m}} = ([0, 0] \text{'a'} [2, 5], \emptyset)$  where  $\bar{\mathbf{C}}$  is the constant propagation domain for characters and  $\bar{\mathbf{B}}$  the interval domain.  $\bar{\mathbf{m}}$  approximates character arrays certainly containing a string of interest, which is actually a sequence of 'a', whose length goes from 2 to 5, followed by a null character, e.g., "aa\0" and "aaaa\0".

In the rest of the chapter, we will refer to the  $s$  and the  $ns$  parameters of a given split segmentation abstract predicate  $\bar{\mathbf{m}}$  by  $\bar{\mathbf{m}}.s$  and  $\bar{\mathbf{m}}.ns$  respectively.

M-String, like FunArray, is equipped with partial order  $\sqsubseteq_{\bar{\mathbf{M}}}$ , join  $\sqcup_{\bar{\mathbf{M}}}$ , meet  $\sqcap_{\bar{\mathbf{M}}}$ , widening  $\nabla_{\bar{\mathbf{M}}}$  and narrowing  $\Delta_{\bar{\mathbf{M}}}$  operators (cf. Section 3.2.2). We highlight the fact that the choice of  $\bar{\mathbf{B}}$  is let free, so the segmentation unification algorithm presented in [57] needs to be accordingly modified, while preserving its original requirements. The unify procedure behaves as follows: given  $\bar{\mathbf{m}}_1, \bar{\mathbf{m}}_2 \in \bar{\mathbf{M}}$ ,  $\text{unify}(\bar{\mathbf{m}}_1, \bar{\mathbf{m}}_2)$  results into the pair  $\text{unify}(\bar{\mathbf{m}}_1.s, \bar{\mathbf{m}}_2.s)$  and  $\text{unify}(\bar{\mathbf{m}}_1.ns, \bar{\mathbf{m}}_2.ns)$ , where  $\bar{\mathbf{m}}_1.s$  and  $\bar{\mathbf{m}}_2.s$  ( $\bar{\mathbf{m}}_1.ns$  and  $\bar{\mathbf{m}}_2.ns$  resp.) are compatible, leading to two abstract predicates  $(\bar{\mathbf{m}}'_1.s, \bar{\mathbf{m}}'_1.ns)$  and  $(\bar{\mathbf{m}}'_2.s, \bar{\mathbf{m}}'_2.ns)$ , respectively. Given two split segmentations  $\bar{\mathbf{m}}_1$  and  $\bar{\mathbf{m}}_2$ , let  $l_{\bar{\mathbf{m}}_1.s}$  and  $h_{\bar{\mathbf{m}}_1.s}$  ( $l_{\bar{\mathbf{m}}_2.s}$  and  $h_{\bar{\mathbf{m}}_2.s}$  resp.) denote the lower and upper bounds of  $\bar{\mathbf{m}}_1.s$  ( $\bar{\mathbf{m}}_2.s$  resp.).  $\bar{\mathbf{m}}_1.s$  and  $\bar{\mathbf{m}}_2.s$  are compatible if  $l_{\bar{\mathbf{m}}_1.s} \sqcap_{\bar{\mathbf{B}}} l_{\bar{\mathbf{m}}_2.s} \neq \perp_{\bar{\mathbf{B}}}$  and  $h_{\bar{\mathbf{m}}_1.s} \sqcap_{\bar{\mathbf{B}}} h_{\bar{\mathbf{m}}_2.s} \neq \perp_{\bar{\mathbf{B}}}$ . The same apply to  $\bar{\mathbf{m}}_1.ns$  and  $\bar{\mathbf{m}}_2.ns$ . Definitions 3.7 and 3.8 present how the join and the meet operators over  $\bar{\mathbf{M}}$  are computed. The widening and narrowing can be easily derived.

**Example 3.11.** Consider following segmentations:  $\bar{\mathbf{m}}_1 = ([0, 0] \text{ odd } [2, 4] \text{ even } [7, 7], \emptyset)$  and  $\bar{\mathbf{m}}_2 = ([0, 0] \text{ odd } [1, 2] \top_{\text{par}} [3, 6] \text{ even } [7, 7], \emptyset)$ , where  $\text{even}$ ,  $\text{odd}$  and  $\top_{\text{par}}$  are elements of the parity domain [56]. Their unification leads to the abstract elements  $\bar{\mathbf{m}}'_1 = ([0, 0] \text{ odd } [2, 4] \text{ even } [7, 7], \emptyset)$  and  $\bar{\mathbf{m}}'_2 = ([0, 0] \text{ odd } [2, 2] \top_{\text{par}} [7, 7], \emptyset)$ . Please observe that the unify yields to a pair of segmentations with the same number of segments and is not always optimal.

◻

**Definition 3.7** (M-String join).  $\sqcup_{\bar{\mathbf{M}}}$  represents the join operator that defines a minimal upper bound between two abstract elements. Let  $\text{unify}(\bar{\mathbf{m}}_1, \bar{\mathbf{m}}_2) = \bar{\mathbf{m}}'_1, \bar{\mathbf{m}}'_2$ , then  $\bar{\mathbf{m}}'_1 \sqcup_{\bar{\mathbf{M}}} \bar{\mathbf{m}}'_2 = (\bar{\mathbf{m}}'_1.s \sqcup_{\bar{\mathbf{M}}} \bar{\mathbf{m}}'_2.s, \bar{\mathbf{m}}'_1.ns \sqcup_{\bar{\mathbf{M}}} \bar{\mathbf{m}}'_2.ns)$  where:

- $\bar{\mathbf{m}}'_1.s \sqcup_{\bar{\mathbf{M}}} \bar{\mathbf{m}}'_2.s = \bar{\mathbf{b}}_1^1 \sqcup_{\bar{\mathbf{B}}} \bar{\mathbf{b}}_1^2 \bar{\mathbf{p}}_1^1 \sqcup_{\bar{\mathbf{C}}} \bar{\mathbf{p}}_1^2 \dots \bar{\mathbf{b}}_k^1 \sqcup_{\bar{\mathbf{B}}} \bar{\mathbf{b}}_k^2 [?^1_k] \gamma [?^2_k]$
- $\bar{\mathbf{m}}'_1.ns \sqcup_{\bar{\mathbf{M}}} \bar{\mathbf{m}}'_2.ns = \bar{\mathbf{b}}_{k+1}^1 \sqcup_{\bar{\mathbf{B}}} \bar{\mathbf{b}}_{k+1}^2 \bar{\mathbf{p}}_{k+1}^1 \sqcup_{\bar{\mathbf{C}}} \bar{\mathbf{p}}_{k+1}^2 \dots \bar{\mathbf{b}}_n^1 \sqcup_{\bar{\mathbf{B}}} \bar{\mathbf{b}}_n^2 [?^1_n] \gamma [?^2_n]$

if  $\bar{\mathbf{m}}'_1.s$  and  $\bar{\mathbf{m}}'_2.s$  ( $\bar{\mathbf{m}}'_1.ns$  and  $\bar{\mathbf{m}}'_2.ns$  resp.) are compatible;  $\top_{\bar{\mathbf{M}}}$  otherwise.<sup>5</sup>

△

**Definition 3.8** (M-String meet).  $\sqcap_{\overline{\mathbf{M}}}$  represents the meet operator that defines a maximal lower bound between two abstract elements. Let  $\text{unify}(\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2) = \overline{\mathbf{m}}'_1, \overline{\mathbf{m}}'_2$ , then  $\overline{\mathbf{m}}'_1 \sqcap_{\overline{\mathbf{M}}} \overline{\mathbf{m}}'_2 = (\overline{\mathbf{m}}'_1.s \sqcap_{\overline{\mathbf{M}}} \overline{\mathbf{m}}'_2.s, \overline{\mathbf{m}}'_1.ns \sqcap_{\overline{\mathbf{M}}} \overline{\mathbf{m}}'_2.ns)$  where:

- $\overline{\mathbf{m}}'_1.s \sqcap_{\overline{\mathbf{M}}} \overline{\mathbf{m}}'_2.s = \overline{\mathbf{b}}_1^1 \sqcap_{\overline{\mathbf{B}}} \overline{\mathbf{b}}_1^2 \overline{\mathbf{p}}_1^1 \sqcap_{\overline{\mathbf{C}}} \overline{\mathbf{p}}_1^2 \dots \overline{\mathbf{b}}_k^1 \sqcap_{\overline{\mathbf{B}}} \overline{\mathbf{b}}_k^2 \wedge [?_k^1] \wedge [?_k^2]$
- $\overline{\mathbf{m}}'_1.ns \sqcap_{\overline{\mathbf{M}}} \overline{\mathbf{m}}'_2.ns = \overline{\mathbf{b}}_{k+1}^1 \sqcap_{\overline{\mathbf{B}}} \overline{\mathbf{b}}_{k+1}^2 \overline{\mathbf{p}}_{k+1}^1 \sqcap_{\overline{\mathbf{C}}} \overline{\mathbf{p}}_{k+1}^2 \dots \overline{\mathbf{b}}_n^1 \sqcap_{\overline{\mathbf{B}}} \overline{\mathbf{b}}_n^2 \wedge [?_n^1] \wedge [?_n^2]$

if  $\overline{\mathbf{m}}'_1.s$  and  $\overline{\mathbf{m}}'_2.s$  ( $\overline{\mathbf{m}}'_1.ns$  and  $\overline{\mathbf{m}}'_2.ns$  resp.) are compatible;  $\perp_{\overline{\mathbf{M}}}$  otherwise.<sup>6</sup>

△

### Concretization

The concretization function on the M-String abstract domain  $\gamma_{\overline{\mathbf{M}}}$  maps an abstract element to a set of concrete character array values as follows:  $\gamma_{\overline{\mathbf{M}}}(\perp_{\overline{\mathbf{M}}}) = \emptyset$ , otherwise  $\gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}})$  is the set of all possible character array values represented by a split segmentation abstract predicate  $\overline{\mathbf{m}}$ .

Formally, we firstly define the concretization function of a generic segment  $(\overline{\mathbf{b}}\overline{\mathbf{p}}\overline{\mathbf{b}}'[?])$  (regardless of what part of the split it is part of)  $\gamma_{\overline{\mathbf{M}}}^*$ , following [57], which corresponds to the set of character array values whose elements in the segment  $[\overline{\mathbf{b}}, \overline{\mathbf{b}}'[?])$  satisfy the predicate  $\overline{\mathbf{p}}$ .

$$\begin{aligned} \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{b}}\overline{\mathbf{p}}\overline{\mathbf{b}}'[?])\overline{\rho} = \{ & (\rho, l, h, M, N) \mid \rho \in \gamma_{\overline{\mathbf{R}}}(\overline{\rho}) \wedge \forall b, b' : b \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}), b' \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}') \wedge \\ & \llbracket l \rrbracket \rho \leq b \leq b' \leq \llbracket h \rrbracket \rho \wedge \forall i \in [b, b'] : M(i) \in \gamma_{\overline{\mathbf{C}}}(\overline{\mathbf{p}}) \wedge \\ & N = \{i \mid M(i) = (i, '\backslash 0')\} \} \end{aligned}$$

where  $\gamma_{\overline{\mathbf{R}}} \in \overline{\mathbf{R}} \rightarrow \mathcal{P}(\mathbf{R}_v)$  is the concretization function for the variable environment abstract domain,  $\gamma_{\overline{\mathbf{B}}} \in \overline{\mathbf{B}} \rightarrow \mathcal{P}(\mathbb{Z}^+)$  is the concretization function for the segment bounds abstract domain, and  $\gamma_{\overline{\mathbf{C}}} \in \overline{\mathbf{C}} \rightarrow \mathcal{P}(\mathbb{Z} \times \mathbf{C})$  is the concretization function for the array characters abstract domain.

We remind that the upper bound of  $\overline{\mathbf{m}}.s$  is not followed by a segment abstract predicate. Let  $\overline{\mathbf{b}}$  be the upper bound of  $\overline{\mathbf{m}}.s$  (which may coincide with the lower bound of  $\overline{\mathbf{m}}.s$  in the case in which  $\overline{\mathbf{m}}$  approximates characters arrays containing null strings of interest).  $\overline{\mathbf{b}}$  is equivalent to the segment  $\overline{\mathbf{b}}\overline{\mathbf{p}}\overline{\mathbf{b}}'$  such that  $\overline{\mathbf{b}}' = \overline{\mathbf{b}} \oplus_{\overline{\mathbf{B}}} 1$  and  $\overline{\mathbf{p}}$  is null.

An abstract element in the M-String domain is a pair of segmentations. Thus, we define the concretization function of the possible  $\overline{\mathbf{m}}.s$  and  $\overline{\mathbf{m}}.ns$  belonging to a character array abstract predicate  $\overline{\mathbf{m}}$ , i.e.,  $\gamma_{\overline{\mathbf{M}}}^* \in \overline{\mathbf{M}} \rightarrow \overline{\mathbf{R}} \rightarrow \mathcal{P}(\mathbf{M})$ . Let  $\upharpoonright_{\overline{\mathbf{M}}}$  denote the concatenation of several concrete values.

<sup>5</sup>Note that  $\sqcup_{\overline{\mathbf{B}}}$ ,  $\sqcup_{\overline{\mathbf{C}}}$  and  $\vee$  denote the join operator of  $\overline{\mathbf{B}}$ ,  $\overline{\mathbf{C}}$  and  $\{?, ?\}$ , respectively.

<sup>6</sup>Note that  $\sqcap_{\overline{\mathbf{B}}}$ ,  $\sqcap_{\overline{\mathbf{C}}}$  and  $\wedge$  denote the meet operator of  $\overline{\mathbf{B}}$ ,  $\overline{\mathbf{C}}$  and  $\{?, ?\}$ , respectively.



$$\begin{aligned}
\gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}}.s)\overline{\rho} &= \left\{ (\rho, l, h, M, N) \in \prod_{i=1}^k \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{b}}_i \overline{\mathbf{p}}_i \overline{\mathbf{b}}_{i+1} [?_{i+1}]) \overline{\rho} \mid \forall b_1, b_k : b_1 \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_1), \right. \\
&\quad b_k \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_k) \wedge b_1 = \llbracket l \rrbracket \rho \wedge b_k + 1 \leq \llbracket h \rrbracket \rho \} \\
&\quad \text{if } \overline{\mathbf{m}}.s = \overline{\mathbf{b}}_1 \overline{\mathbf{p}}_1 \overline{\mathbf{b}}_2 [?_2] \dots \overline{\mathbf{b}}_{k-1} [?_{k-1}] \overline{\mathbf{p}}_{k-1} \overline{\mathbf{b}}_k [?_k] \\
&= \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{b}}_1) \overline{\rho} \\
&\quad \text{if } \overline{\mathbf{m}}.s = \overline{\mathbf{b}}_1 \\
&= \emptyset \\
&\quad \text{otherwise} \\
\gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}}.ns)\overline{\rho} &= \left\{ (\rho, l, h, M, N) \in \prod_{i=1}^{n-1} \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{b}}_i \overline{\mathbf{p}}_i \overline{\mathbf{b}}_{i+1} [?_{i+1}]) \overline{\rho} \mid \forall b_1, b_n : \right. \\
&\quad b_1 \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_1), b_n \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_n) \wedge b_1 = \llbracket l \rrbracket \rho \wedge b_n = \llbracket h \rrbracket \rho \} \\
&\quad \text{if } \overline{\mathbf{m}}.ns = \overline{\mathbf{b}}_1 \overline{\mathbf{p}}_1 \overline{\mathbf{b}}_2 [?_2] \dots \overline{\mathbf{b}}_n [?_n] \\
&= \left\{ (\rho, l, h, M, N) \in \prod_{i=k+1}^{n-1} \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{b}}_i \overline{\mathbf{p}}_i \overline{\mathbf{b}}_{i+1} [?_{i+1}]) \overline{\rho} \mid \forall b_{k+1}, b_n : \right. \\
&\quad b_{k+1} \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_{k+1}), b_n \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_n) \wedge \llbracket l \rrbracket \rho < b_{k+1} \wedge b_n = \llbracket h \rrbracket \rho \} \\
&\quad \text{if } \overline{\mathbf{m}}.ns = \overline{\mathbf{b}}_{k+1} \overline{\mathbf{p}}_{k+1} \overline{\mathbf{b}}_{k+2} [?_{k+2}] \dots \overline{\mathbf{b}}_n [?_n] \\
&= \emptyset \\
&\quad \text{otherwise}
\end{aligned}$$

Finally, the concretization function of a split segmentation abstract predicate  $\overline{\mathbf{m}}$  is as follows:

$$\gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}})\overline{\rho} \triangleq \left\{ (\rho, l, h, M, N) \in \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}}.s)\overline{\rho} +_{\mathbf{M}} \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}}.ns)\overline{\rho} \mid \forall b_1, b_n : b_1 \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_1), \right. \\
\left. b_n \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_n) \wedge b_1 = \llbracket l \rrbracket \rho \wedge b_n = \llbracket h \rrbracket \rho \right\}$$

where  $+_{\mathbf{M}}$  returns all the possible concatenations between a concrete array value taken from  $\gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}}.s)$ , and a concrete array value taken from  $\gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}}.ns)$ .

**Definition 3.9** (Invalid segment). Given a generic segment  $\overline{\mathbf{b}}\overline{\mathbf{p}}\overline{\mathbf{b}}' [?]$ , it is considered invalid if its segment abstract predicate  $\overline{\mathbf{p}}$  is equal to  $\perp_{\overline{\mathcal{C}}}$  and its upper bound  $\overline{\mathbf{b}}'$  is not followed by a question mark.

△

**Theorem 3.1.** Let  $\overline{\mathbf{X}} \subseteq \overline{\mathbf{M}}$  such that all elements in  $\overline{\mathbf{X}}$  are compatible and their meet does not result in split segmentation abstract predicates which contain invalid abstract elements. Then, it holds that:

$$\gamma_{\overline{\mathbf{M}}} \left( \prod_{\overline{\mathbf{m}} \in \overline{\mathbf{X}}} \overline{\mathbf{m}} \right) = \bigcap_{\overline{\mathbf{m}} \in \overline{\mathbf{X}}} \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}})$$

*Proof.* The following inference chain holds:

$$\begin{aligned}
& \gamma_{\overline{\mathbf{M}}} \left( \prod_{\overline{\mathbf{m}} \in \overline{\mathbf{X}}} \overline{\mathbf{m}} \right) \\
&= \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}'}) \quad \text{where } \overline{\mathbf{m}'} \text{ is the result of the meet operation over } \overline{\mathbf{X}} \text{ (cf. Definition 3.8)} \\
&= \{(\rho, l, h, M, N) \in \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}'}.s)\overline{\rho} +_{\mathbf{M}} \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}'}.ns)\overline{\rho} \mid \forall b_1, b_n : b_1 \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_1), \\
&\quad b_n \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_n) \wedge b_1 = \llbracket l \rrbracket \rho \wedge b_n = \llbracket h \rrbracket \rho\} \quad \text{by definition of } \gamma_{\overline{\mathbf{M}}} \\
&= \bigcap_{\overline{\mathbf{m}} \in \overline{\mathbf{X}}} \{(\rho, l, h, M, N) \in \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}}.s)\overline{\rho} +_{\mathbf{M}} \gamma_{\overline{\mathbf{M}}}^*(\overline{\mathbf{m}}.ns)\overline{\rho} \mid \forall b_1, b_n : b_1 \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_1), \\
&\quad b_n \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_n) \wedge b_1 = \llbracket l \rrbracket \rho \wedge b_n = \llbracket h \rrbracket \rho\} \\
&= \bigcap_{\overline{\mathbf{m}} \in \overline{\mathbf{X}}} \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}) \quad \text{by definition of } \gamma_{\overline{\mathbf{M}}} \\
& \quad \square
\end{aligned}$$

Observe that if the hypotheses of Theorem 3.1 are not satisfied, i.e., if either the abstract predicates in  $\overline{\mathbf{X}}$  are not compatible or their meet leads to invalid segmentations, then  $\gamma_{\overline{\mathbf{M}}} \left( \prod_{\overline{\mathbf{m}} \in \overline{\mathbf{X}}} \overline{\mathbf{m}} \right) = \gamma_{\overline{\mathbf{M}}}(\perp_{\overline{\mathbf{M}}}) = \emptyset$ , and  $\bigcap_{\overline{\mathbf{m}} \in \overline{\mathbf{X}}} \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}) = \emptyset$ .

In the implementation, we will use the following two functions: *lift* and *lower* that relate single strings to their abstraction in M-String.

**Definition 3.10** (Lift). Let  $\mu(\mathbf{m}) \subseteq \mathbf{M}$  be a concrete character array value. Given the abstraction function  $\alpha_{\overline{\mathbf{M}}}$  on M-String, we define the lift operation of  $X$  as follows:

$$\text{lift}(\mu(\mathbf{m})) = \alpha_{\overline{\mathbf{M}}}(\mu(\mathbf{m})).$$

△

**Definition 3.11** (Lower). Let  $\overline{\mathbf{m}}$  denote  $\text{lift}(\mu(\mathbf{m}))$  (cf. Definition 3.10). Given the concretization function  $\gamma_{\overline{\mathbf{M}}}$  on M-String, we define the lower operation of  $\overline{\mathbf{m}}$  as follows:

$$\text{lower}(\overline{\mathbf{m}}) = \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}})$$

△

### Abstraction

Let  $X$  be a set of concrete character array values. The abstraction function on the M-String abstract domain  $\alpha_{\overline{\mathbf{M}}}$  maps  $X$  to  $\perp_{\overline{\mathbf{M}}}$  in the case in which  $X$  is empty. Otherwise,  $\alpha_{\overline{\mathbf{M}}}$  maps  $X$  to an element of M-String that over-approximates values in  $X$ .

### 3.5.2 Abstract Semantics

We now formalize the abstract semantics of the concrete operations defined in Section 3.4.3, over the M-String domain. In doing so, we will take advantage of the auxiliary function  $\text{minLen}_{\overline{\mathbf{M}}}$  which computes the minimum length of an element  $\overline{\mathbf{m}} \in \overline{\mathbf{M}}$ , as the upper bound of a split segmentation is possibly followed by a question mark.

**Definition 3.12** ( $\overline{\mathbf{m}}$  minimum length). Let  $\overline{\mathbf{m}} \in \overline{\mathbf{M}}$  be different from  $\perp_{\overline{\mathbf{M}}}$  and let  $l_{\overline{\mathbf{m}}}, h_{\overline{\mathbf{m}}} \in \overline{\mathbf{B}}$  denote the lower and the upper bound of  $\overline{\mathbf{m}}$ , respectively. We define the minimum length of a split segmentation abstract predicate  $\overline{\mathbf{m}}$ , denoted by  $\text{minLen}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}})$ , as follows:

$$\text{minLen}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}) = \begin{cases} \overline{\mathbf{b}}_k \ominus_{\overline{\mathbf{B}}} l_{\overline{\mathbf{m}}} & \text{if } \overline{\mathbf{m}}.\text{ns} \neq \emptyset \wedge h_{\overline{\mathbf{m}}} \text{ is followed by } ? \wedge \exists k \in \overline{\mathbf{m}}.\text{ns} : \\ & k = \max(\{i \in \overline{\mathbf{m}}.\text{ns} \mid \overline{\mathbf{b}}_i \text{ is not followed by } ?\}) \\ \overline{\mathbf{h}}_{\overline{\mathbf{m}}} \ominus_{\overline{\mathbf{B}}} l_{\overline{\mathbf{m}}} & \text{otherwise} \end{cases}$$

△

Please note that in the second case of Definition 3.12, the minimum length of a split segmentation corresponds to its length, denoted by  $\text{len}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}})$ . This operation can be also applied to the parameters of  $\overline{\mathbf{m}}$  themselves, when they are different from the emptyset and their upper bound is not question marked, which is always the case with  $\overline{\mathbf{m}}.\text{s}$ .

**Example 3.12.** Consider the split segmentation abstract predicate  $\overline{\mathbf{m}} = ([0, 0] \text{'a'} [2, 5], [3, 6] \text{'b'} [7, 7] \text{'c'} [8, 8]?)$  where  $\overline{\mathbf{C}}$  is the constant propagation domain for characters and  $\overline{\mathbf{B}}$  the interval domain. The minimum length of  $\overline{\mathbf{m}}$  is given by  $\text{minLen}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}) = [7, 7] \ominus_{\overline{\mathbf{B}}} [0, 0] = [7, 7]$  as its upper bound is followed by a question mark. Logically speaking, the maximum length of  $\overline{\mathbf{m}}$  is  $[8, 8] \ominus_{\overline{\mathbf{B}}} [0, 0] = [8, 8]$ . The length of  $\overline{\mathbf{m}}.\text{s}$  is given by  $\text{len}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}.\text{s}) = [2, 5] \ominus_{\overline{\mathbf{B}}} [0, 0] = [2, 5]$ .

○

#### Abstract Array Access

The semantic operator  $\mathfrak{A}_{\overline{\mathbf{M}}}$  is the abstract counterpart of  $\mathfrak{A}$  (cf. Section 3.4.3). In particular,  $\text{access}_{\bar{j}}(\overline{\mathbf{m}})$  returns, if  $\bar{j}$  is valid for  $\overline{\mathbf{m}}$  (i.e., there exist, and it is unique, a segment bounds interval  $[\overline{\mathbf{b}}_i[?^i], \overline{\mathbf{b}}_{i+1})$  in  $\overline{\mathbf{m}}$  to which  $\bar{j}$  belongs), the segment abstract predicate  $\overline{\mathbf{p}}_i$ ; otherwise it returns  $\top_{\overline{\mathbf{C}}}$ . Formally,

$$\mathfrak{A}_{\overline{\mathbf{M}}}[\text{access}_{\bar{j}}](\overline{\mathbf{m}}) = \begin{cases} \overline{\mathbf{p}}_i & \text{if } \exists! i \in \overline{\mathbf{m}} : \bar{j} \in [\overline{\mathbf{b}}_i[?^i], \overline{\mathbf{b}}_{i+1}) \\ \top_{\overline{\mathbf{C}}} & \text{otherwise} \end{cases}$$

where  $\bar{j} \in [\overline{\mathbf{b}}_i[?^i], \overline{\mathbf{b}}_{i+1})$  if  $\forall j \in \gamma_{\overline{\mathbf{B}}}(\bar{j}) : \forall [l, u) \in \{[l, u) \mid l \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_i) \wedge u \in \gamma_{\overline{\mathbf{B}}}(\overline{\mathbf{b}}_{i+1})\} : j \in [l, u)$ .

### Abstract String Concatenation

The semantic operator  $\mathfrak{M}_{\overline{\mathbf{M}}}$  is the abstract counterpart of  $\mathfrak{M}$ . When applied to  $\text{strcat}(\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2)$ , it returns  $\overline{\mathbf{m}}_1'$  that is  $\overline{\mathbf{m}}_1$  into which  $\overline{\mathbf{m}}_2.s$  has been embedded starting from the upper bound of  $\overline{\mathbf{m}}_1.s$ , if both the input split segmentations approximate character arrays which contain a well-formed string and the condition on the size of the destination split segmentation is fulfilled; otherwise it returns  $\top_{\overline{\mathbf{M}}}$ . Formally,

$$\mathfrak{M}_{\overline{\mathbf{M}}}[\text{strcat}](\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2) = \begin{cases} \overline{\mathbf{m}}_1' & \text{if } \overline{\mathbf{m}}_1.s \neq \emptyset \neq \overline{\mathbf{m}}_2.s \wedge \text{size.condition is true} \\ \top_{\overline{\mathbf{M}}} & \text{otherwise} \end{cases}$$

The `size.condition` is true if  $\text{minLen}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_1) \geq_{\overline{\mathbf{B}}} (\text{len}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_1.s) \oplus_{\overline{\mathbf{B}}} \text{len}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_2.s) \oplus_{\overline{\mathbf{B}}} 1)$ . Let:

- $\overline{\mathbf{m}}_1 = (\overline{\mathbf{b}}_1 \overline{\mathbf{p}}_1^1 \overline{\mathbf{b}}_2^{[?1]_2} \dots \overline{\mathbf{p}}_{k-1}^1 \overline{\mathbf{b}}_k^{[?1]_k}, \overline{\mathbf{b}}_{k+1}^1 \overline{\mathbf{p}}_{k+1}^1 \overline{\mathbf{b}}_{k+2}^{[?1]_{k+2}} \dots \overline{\mathbf{b}}_n^{[?1]_n})$
- $\overline{\mathbf{m}}_2 = (\overline{\mathbf{b}}_1^2 \overline{\mathbf{p}}_1^2 \overline{\mathbf{b}}_2^{[?2]_2} \dots \overline{\mathbf{p}}_{k-1}^2 \overline{\mathbf{b}}_k^{[?2]_k}, \text{ns})$

Then,  $\overline{\mathbf{m}}_1'.s = \overline{\mathbf{b}}_1 \overline{\mathbf{p}}_1^1 \overline{\mathbf{b}}_2^{[?1]_2} \dots \overline{\mathbf{p}}_{k-1}^1 \overline{\mathbf{b}}_k^{[?1]_k} \overline{\mathbf{p}}_1^2 (\overline{\mathbf{b}}_k \oplus_{\overline{\mathbf{B}}} (\overline{\mathbf{b}}_2 \ominus_{\overline{\mathbf{B}}} \overline{\mathbf{b}}_1))^{[?2]_2} \dots \overline{\mathbf{p}}_{k-1}^2 (\overline{\mathbf{b}}' \oplus_{\overline{\mathbf{B}}} (\overline{\mathbf{b}}_k \ominus_{\overline{\mathbf{B}}} \overline{\mathbf{b}}_{k-1}^2))^{[?2]_2}$  such that  $\overline{\mathbf{b}}'$  denotes the immediately preceding adapted segment bound. On the other hand,  $\overline{\mathbf{m}}_1'.\text{ns}$  is the result of removing from  $\overline{\mathbf{m}}_1.\text{ns}$  the sub-segmentation that goes from the lower bound of  $\overline{\mathbf{m}}_1.\text{ns}$  to the upper bound of  $\overline{\mathbf{m}}_1'.s$  included.

**Example 3.13.** Let  $([0, 0] \text{ a}^* [5, 7], [6, 8] \text{ br}^* [13, 14])$  and  $([0, 0] \text{ a}^* [3, 3], \emptyset)$  be two abstract elements in  $\overline{\mathbf{M}}$ , such that  $\overline{\mathbf{B}}$  is the interval domain over array indexes and  $\overline{\mathbf{C}}$  is the prefix domain over string values. Precisely,  $([0, 0] \text{ a}^* [5, 7], [6, 8] \text{ br}^* [13, 14])$  approximates all the characters arrays having as string of interest any string starting with the character 'a' whose length goes from 5 to 7, followed by the null character and any string starting with "br" whose length goes from 5 to 8. On the other hand,  $([0, 0] \text{ a}^* [3, 3], \emptyset)$  abstracts all the array of chars with string of interest equal to a string, of length 3, starting with a . Consider now the concatenation between them,

$$\mathfrak{M}_{\overline{\mathbf{M}}}[\text{strcat}](( [0, 0] \text{ a}^* [5, 7], [6, 8] \text{ br}^* [13, 14]), ([0, 0] \text{ a}^* [3, 3], \emptyset))$$

The size condition is satisfied: the minimum length of the destination split segmentation is equal to 13, which is strictly greater than  $7 + 3 + 1$ , i.e., the maximum length of the destination abstract array plus the maximum length of the source segmentation plus one (the null character). Their concatenation results in the following abstract element:

$$([0, 0] \text{ a}^* [5, 7] \text{ a}^* [8, 10], [9, 11] \text{ br}^* [13, 14])$$

which is equivalent to  $([0, 0] \text{ a}^* [8, 10], [9, 11] \text{ br}^* [13, 14])$ .

◻

### Abstract String Character

The semantic operator  $\mathfrak{M}_{\overline{M}}$ , when applied to  $\text{strchr}_{\overline{v}}(\overline{m})$ , returns a split segmentation abstract predicate  $\overline{s}$  with the left hand side parameter equal to the suffix segmentation of the input  $\overline{m}.s$  from the first segment in which  $\overline{v}$  certainly occurs and the right hand side parameter equal to the emptyset, if  $\overline{m}$  approximates character arrays which contain a well-formed string and the character  $\overline{v}$  appears in at least one segment whose bounds are not question marked. Otherwise, if  $\overline{m}$  approximates character arrays which contain a well-formed string of interest and the abstract character  $\overline{v}$  does not occur in  $\overline{m}.s$ , it returns  $\perp_{\overline{M}}$ ; otherwise it returns  $\top_{\overline{M}}$ . Formally,

$$\mathfrak{M}_{\overline{M}}[\text{strchr}_{\overline{v}}](\overline{m}) = \begin{cases} \overline{s} & \text{if } \overline{m}.s \neq \emptyset \wedge \exists i \in \overline{m}.s : \overline{p}_i = \overline{v} \wedge \\ & \overline{b}_i, \overline{b}_{i+1} \text{ are not followed by ?} \\ \perp_{\overline{M}} & \text{if } \overline{m}.s \neq \emptyset \wedge \nexists i \in \overline{m}.s : \overline{p}_i = \overline{v} \\ \top_{\overline{M}} & \text{otherwise} \end{cases}$$

where  $\overline{s} = (\overline{b}_j \overline{p}_j \overline{b}_{j+1} \dots \overline{b}_k [?_k], \emptyset)$  such that

$$j = \min(\{i \in \overline{m}.s \mid \overline{p}_i = \overline{v} \wedge \overline{b}_i, \overline{b}_{i+1} \text{ are not question marked}\})$$

### Abstract String Compare

The semantic  $\mathfrak{P}_{\overline{M}}$  is the abstract counterpart of  $\mathfrak{P}$ . In particular,  $\text{strcmp}(\overline{m}_1, \overline{m}_2)$  returns a value  $\overline{n}$  denoting the lexicographic distance between  $\overline{m}_1.s$  and  $\overline{m}_2.s$  if both the input split segmentations approximate character arrays which contain a well-formed string and they can be unified; otherwise it returns  $\top_{\mathbb{Z}}$ .

Note that if  $\overline{n}$  is negative, the strings of interest approximated by  $\overline{m}_1$  precede those represented by  $\overline{m}_2$  in lexicographic order. Conversely, if  $\overline{n}$  is positive, the strings of interest approximated by  $\overline{m}_1$  follow those represented by  $\overline{m}_2$  in lexicographic order, and if  $\overline{n}$  is equal to zero they are lexicographically equal. Formally,

$$\mathfrak{P}_{\overline{M}}[\text{strcmp}](\overline{m}_1, \overline{m}_2) = \begin{cases} \overline{n} & \text{if } \overline{m}_1.s \neq \emptyset \neq \overline{m}_2.s \\ \top_{\mathbb{Z}} & \text{otherwise} \end{cases}$$

where  $\overline{n} = \text{cmp}_{\overline{M}}(\overline{m}_1, \overline{m}_2)$  (cf. Algorithm 2).

### Abstract String Copy

The semantic  $\mathfrak{M}_{\overline{M}}$ , when applied to  $\text{strcpy}(\overline{m}_1, \overline{m}_2)$ , returns  $\overline{m}_1'$  that is  $\overline{m}_1$  into which  $\overline{m}_2.s$  has been embedded starting from the lower bound of  $\overline{m}_1$ , if both the input split segmentations approximate character arrays which contain a well-formed string and the con-

---

**Algorithm 2** Lexicographic comparison of split segmentation abstract predicates.

---

**Input:** two compatible split segmentation abstract predicates  $\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2 \in \overline{\mathbf{M}}$ .

**Output:** an integer value  $\overline{n}$ .

```

1:  $\overline{n} = 0, i = 1$ 
2:  $\text{unify}(\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2) = \overline{\mathbf{m}}'_1, \overline{\mathbf{m}}'_2$ 
3: if  $\overline{\mathbf{m}}'_1.s = \overline{\mathbf{b}}_1^1 \wedge \overline{\mathbf{m}}'_2.s = \overline{\mathbf{b}}_1^2$  then
4:   return  $\overline{n}$ 
5: else if  $\overline{\mathbf{m}}'_1.s = \overline{\mathbf{b}}_1^1 \wedge \overline{\mathbf{m}}'_2.s \neq \overline{\mathbf{b}}_1^2$  then
6:    $\overline{n} = \overline{n} \ominus_{\mathbb{C}} \overline{\mathbf{p}}_1^2$ 
7:   return  $\overline{n}$ 
8: else if  $\overline{\mathbf{m}}'_1.s \neq \overline{\mathbf{b}}_1^1 \wedge \overline{\mathbf{m}}'_2.s = \overline{\mathbf{b}}_1^2$  then
9:    $\overline{n} = \overline{\mathbf{p}}_1^1 \ominus_{\mathbb{C}} \overline{n}$ 
10:  return  $\overline{n}$ 
11: else
12:   while  $i \in \overline{\mathbf{m}}'_1.s \wedge i \in \overline{\mathbf{m}}'_2.s$  do
13:      $\overline{n} = \overline{\mathbf{p}}_i^1 \ominus_{\mathbb{C}} \overline{\mathbf{p}}_i^2$ 
14:     if  $\overline{n} \neq 0$  then
15:       return  $\overline{n}$ 
16:     else
17:        $i = i + 1$ 
18: return  $\overline{n}$ 

```

---

dition on the size of the destination split segmentation is fulfilled; otherwise it returns  $\top_{\overline{\mathbf{M}}}$ . Formally,

$$\mathfrak{M}_{\overline{\mathbf{M}}}[\text{strcpy}](\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2) = \begin{cases} \overline{\mathbf{m}}'_1 & \text{if } \overline{\mathbf{m}}_1.s \neq \emptyset \neq \overline{\mathbf{m}}_2.s \wedge \text{size.condition is true} \\ \top_{\overline{\mathbf{M}}} & \text{otherwise} \end{cases}$$

The `size.condition` is true if  $\text{minLen}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_1) \geq_{\mathbb{B}} \text{len}_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_2.s) \oplus_{\mathbb{B}} 1$ . Let:

- $\overline{\mathbf{m}}_1 = (\overline{\mathbf{b}}_1^1 \overline{\mathbf{p}}_1^1 \overline{\mathbf{b}}_2^1 [?_2^1] \dots \overline{\mathbf{p}}_{k-1}^1 \overline{\mathbf{b}}_k^1 [?_k^1], \overline{\mathbf{b}}_{k+1}^1 \overline{\mathbf{p}}_{k+1}^1 \overline{\mathbf{b}}_{k+2}^1 [?_{k+2}^1] \dots \overline{\mathbf{b}}_n^1 [?_n^1])$
- $\overline{\mathbf{m}}_2 = (\overline{\mathbf{b}}_1^2 \overline{\mathbf{p}}_1^2 \overline{\mathbf{b}}_2^2 [?_2^2] \dots \overline{\mathbf{p}}_{k-1}^2 \overline{\mathbf{b}}_k^2 [?_k^2], \text{ns})$

Then,  $\overline{\mathbf{m}}'_1.s = \overline{\mathbf{b}}_1^1 \overline{\mathbf{p}}_1^2 (\overline{\mathbf{b}}_1^1 \oplus_{\mathbb{B}} (\overline{\mathbf{b}}_2^2 \ominus_{\mathbb{B}} \overline{\mathbf{b}}_1^2)) [?_2^2] \dots \overline{\mathbf{p}}_{k-1}^2 (\overline{\mathbf{b}}^1 \oplus_{\mathbb{B}} (\overline{\mathbf{b}}_k^2 \ominus_{\mathbb{B}} \overline{\mathbf{b}}_{k-1}^2)) [?_k^2]$  where  $\overline{\mathbf{b}}^1$  denotes the immediately preceding adapted segment bound. Instead,  $\overline{\mathbf{m}}'_1.\text{ns}$  is the subsegmentation of  $\overline{\mathbf{m}}_1$  that goes from the upper bound of  $\overline{\mathbf{m}}'_1.s$  plus one to the upper bound of  $\overline{\mathbf{m}}_1$ .

### Abstract String Length

The semantic  $\mathfrak{L}_{\overline{\mathbf{M}}}$  is the abstract counterpart of  $\mathfrak{L}$ . In particular, `strlen` returns a value  $\overline{n}$ , if  $\overline{\mathbf{m}}$  approximates character arrays which contain a well-formed string, the upper bound

of  $\bar{\mathbf{m}}.s$  is not followed by a question mark and in  $\bar{\mathbf{m}}.s$  do not occur possibly null segment abstract predicates; otherwise it returns  $\top_{\bar{\mathbf{B}}}$ . Formally,

$$\mathfrak{L}_{\bar{\mathbf{M}}}\llbracket\text{strlen}\rrbracket(\bar{\mathbf{m}}) = \begin{cases} \bar{n} & \text{if } \bar{\mathbf{m}}.s = \bar{\mathbf{b}}_1\bar{\mathbf{p}}_1\bar{\mathbf{b}}_2[?_2] \dots \bar{\mathbf{b}}_k \wedge \\ & \nexists i \in \bar{\mathbf{m}}.s : \text{isNull}(\bar{\mathbf{p}}_i) = \text{maybe} \\ \top_{\bar{\mathbf{B}}} & \text{otherwise} \end{cases}$$

where  $\bar{n} = \bar{\mathbf{b}}_k \ominus_{\bar{\mathbf{B}}} \bar{\mathbf{b}}_1$ .

### Abstract Array Update

The semantic  $\mathfrak{M}_{\bar{\mathbf{M}}}$ , when applied to  $\text{update}_{\bar{j}, \bar{v}}(\bar{\mathbf{m}})$ , returns, if  $\gamma_{\bar{\mathbf{B}}}(\bar{j})$  corresponds to the singleton  $\{j\}$  and  $\bar{j}$  is valid for  $\bar{\mathbf{m}}$  (i.e., there exists - and it is unique - a segment bounds interval  $[\bar{\mathbf{b}}_i[?_i], \bar{\mathbf{b}}_{i+1})$  in  $\bar{\mathbf{m}}$  to which  $\bar{j}$  belongs),  $\bar{\mathbf{m}}'$  that is  $\bar{\mathbf{m}}$  where the segment  $\bar{\mathbf{b}}_i[?_i]\bar{\mathbf{p}}_i\bar{\mathbf{b}}_{i+1}$  is split so that the segment abstract predicate at position  $\bar{j}$  is substituted with  $\bar{v}$ ; otherwise it returns  $\top_{\bar{\mathbf{M}}}$ . Formally,

$$\mathfrak{M}_{\bar{\mathbf{M}}}\llbracket\text{update}_{\bar{j}, \bar{v}}\rrbracket(\bar{\mathbf{m}}) = \begin{cases} \bar{\mathbf{m}}' & \text{if } \gamma_{\bar{\mathbf{B}}}(\bar{j}) = \{j\} \wedge \exists! i \in \bar{\mathbf{m}} : \bar{j} \in [\bar{\mathbf{b}}_i[?_i], \bar{\mathbf{b}}_{i+1}) \\ \top_{\bar{\mathbf{M}}} & \text{otherwise} \end{cases}$$

### 3.5.3 Soundness

**Theorem 3.2.**  $\mathfrak{A}_{\bar{\mathbf{M}}}$  is a sound over-approximation of  $\mathfrak{A}$ . Formally,

$$\gamma_{\bar{\mathbf{C}}}(\mathfrak{A}_{\bar{\mathbf{M}}}\llbracket\text{stm}\rrbracket(\bar{\mathbf{m}})) \supseteq \{\mathfrak{A}\llbracket\text{stm}\rrbracket(\mu(\mathbf{m})) : \mu(\mathbf{m}) \in \gamma_{\bar{\mathbf{M}}}(\bar{\mathbf{m}})\}$$

*Proof.*

Consider the unary operator  $\text{access}_{\bar{j}}$  and let  $\bar{\mathbf{m}}$  be a split segmentation abstract predicate. We have to prove that:

$$\gamma_{\bar{\mathbf{C}}}(\mathfrak{A}_{\bar{\mathbf{M}}}\llbracket\text{access}_{\bar{j}}\rrbracket(\bar{\mathbf{m}})) \supseteq \{\mathfrak{A}\llbracket\text{access}_{\bar{j}}\rrbracket(\mu(\mathbf{m})) : \mu(\mathbf{m}) \in \gamma_{\bar{\mathbf{M}}}(\bar{\mathbf{m}})\}$$

$\text{access}_{\bar{j}}$  of  $\mu(\mathbf{m})$  returns, by definition of  $\mathfrak{A}$ , the character array value  $v$  that occurs at position  $j$ , if  $j$  belongs to  $[[l_{\mathbf{m}}]\rho, [h_{\mathbf{m}}]\rho)$ ;  $\top_{\bar{\mathbf{C}}}$  otherwise. Let  $\alpha_{\bar{\mathbf{B}}}(\bar{j}) = \bar{j}$ . Then,  $v$  belongs to  $\gamma_{\bar{\mathbf{C}}}(\mathfrak{A}_{\bar{\mathbf{M}}}\llbracket\text{access}_{\bar{j}}\rrbracket(\bar{\mathbf{m}}))$  because  $\text{access}_{\bar{j}}$  of  $\bar{\mathbf{m}}$ , by definition of  $\mathfrak{A}_{\bar{\mathbf{M}}}$ , is equal to the segment abstract predicate  $\bar{\mathbf{p}}_i$ , if there exists - and it unique - a segment bounds interval  $[\bar{\mathbf{b}}_i[?_i], \bar{\mathbf{b}}_{i+1})$  to which  $\bar{j}$  belongs;  $\top_{\bar{\mathbf{C}}}$  otherwise. □

**Theorem 3.3.**  $\mathfrak{M}_{\bar{\mathbf{M}}}$  is a sound over-approximation of  $\mathfrak{M}$ . Formally,

$$\gamma_{\bar{\mathbf{M}}}(\mathfrak{M}_{\bar{\mathbf{M}}}\llbracket\text{stm}\rrbracket(\bar{\mathbf{m}})) \supseteq \{\mathfrak{M}\llbracket\text{stm}\rrbracket(\mu(\mathbf{m})) : \mu(\mathbf{m}) \in \gamma_{\bar{\mathbf{M}}}(\bar{\mathbf{m}})\}$$

*Proof.*

- Consider the binary operator `strcat` and let  $\overline{\mathbf{m}}_1$  and  $\overline{\mathbf{m}}_2$  be two split segmentation abstract predicates. We have to prove that:

$$\begin{aligned} & \gamma_{\overline{\mathbf{M}}}(\mathfrak{M}_{\overline{\mathbf{M}}}\llbracket\text{strcat}\rrbracket(\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2)) \\ & \supseteq \\ & \{\mathfrak{M}\llbracket\text{strcat}\rrbracket(\mu(\mathbf{m}_1), \mu(\mathbf{m}_2)) : \mu(\mathbf{m}_1) \in \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_1) \wedge \mu(\mathbf{m}_2) \in \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_2)\} \end{aligned}$$

by definition of  $\mathfrak{M}$ , `strcat` of  $\mu(\mathbf{m}_1)$  and  $\mu(\mathbf{m}_2)$  returns  $\mu(\mathbf{m}_1)'$  where the first null-terminating memory block of  $\mu(\mathbf{m}_2)$  (including the null terminator), i.e., its string of interest, is embedded into  $\mu(\mathbf{m}_1)$  starting from the index to which the first null character in  $\mu(\mathbf{m}_1)$  occurs, if both  $\mu(\mathbf{m}_1)$  and  $\mu(\mathbf{m}_2)$  contain a well-formed string and the size condition on the destination character array value is fulfilled;  $\top_{\overline{\mathbf{M}}}$  otherwise. Then,  $\mu(\mathbf{m}_1)'$  belongs to  $\gamma_{\overline{\mathbf{M}}}(\mathfrak{M}_{\overline{\mathbf{M}}}\llbracket\text{strcat}\rrbracket(\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2))$  because `strcat` of  $\overline{\mathbf{m}}_1$  and  $\overline{\mathbf{m}}_2$ , by definition of  $\mathfrak{M}_{\overline{\mathbf{M}}}$ , is equal to  $\overline{\mathbf{m}}_1'$  that is  $\overline{\mathbf{m}}_1$  into which  $\overline{\mathbf{m}}_2$ 's has been embedded starting from the upper bound of  $\overline{\mathbf{m}}_1$ 's, if both  $\overline{\mathbf{m}}_1$  and  $\overline{\mathbf{m}}_2$  approximate character arrays which contain a well-formed string and the size condition on the destination segmentation abstract predicate is fulfilled;  $\top_{\overline{\mathbf{M}}}$  otherwise.

- Consider the unary operator `strchr $\overline{v}$` , and let  $\overline{\mathbf{m}}$  be a split segmentation abstract predicate. We have to prove that:

$$\gamma_{\overline{\mathbf{M}}}(\mathfrak{M}_{\overline{\mathbf{M}}}\llbracket\text{strchr}_{\overline{v}}\rrbracket(\overline{\mathbf{m}})) \supseteq \{\mathfrak{M}\llbracket\text{strchr}_{\overline{v}}\rrbracket(\mu(\mathbf{m})) : \mu(\mathbf{m}) \in \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}})\}$$

`strchr $\overline{v}$`  of  $\mu(\mathbf{m})$  returns, by definition of  $\mathfrak{M}$ ,  $\mu(\mathbf{s})$  that corresponds to the suffix of the string of interest of  $\mu(\mathbf{m})$  starting from the index to which appears the first occurrence of  $\overline{v}$ , if  $\mu(\mathbf{m})$  contains a well-formed string and  $\overline{v}$  occurs in  $\mu(\mathbf{m})$ ; the emptyset (i.e.,  $\perp_{\overline{\mathbf{M}}}$ ), if  $\mu(\mathbf{m})$  contains a well-formed string and  $\overline{v}$  does not occur in  $\mu(\mathbf{m})$ ;  $\top_{\overline{\mathbf{M}}}$  otherwise. Let  $\alpha_{\overline{\mathbf{C}}}(\overline{v}) = \overline{v}$ . Then,  $\mu(\mathbf{s})$  belongs to  $\gamma_{\overline{\mathbf{M}}}(\mathfrak{M}_{\overline{\mathbf{M}}}\llbracket\text{strchr}_{\overline{v}}\rrbracket(\overline{\mathbf{m}}))$  because `strchr $\overline{v}$`  of  $\overline{\mathbf{m}}$ , by definition of  $\mathfrak{M}_{\overline{\mathbf{M}}}$ , is equal to  $\overline{\mathbf{s}}$  that is the split segmentation abstract predicate with  $\overline{\mathbf{s}}$ 's equal to the sub-segmentation of  $\overline{\mathbf{m}}$ 's starting from the first segment to which  $\overline{v}$  certainly occurs and  $\overline{\mathbf{s}}$ 'ns equal to the emptyset if  $\overline{\mathbf{m}}$  approximates character arrays which contain a well-defined string and  $\overline{v}$  appears in at least one segment whose bounds are not question marked;  $\perp_{\overline{\mathbf{M}}}$  if  $\overline{\mathbf{m}}$  approximates character arrays which contain a well-formed string and  $\overline{v}$  does not appear in  $\overline{\mathbf{m}}$ 's;  $\top_{\overline{\mathbf{M}}}$  otherwise.

- Consider the binary operator `strcpy` and let  $\overline{\mathbf{m}}_1$  and  $\overline{\mathbf{m}}_2$  be two split segmentation abstract predicates. We have to prove that:

$$\gamma_{\overline{\mathbf{M}}}(\mathfrak{M}_{\overline{\mathbf{M}}}\llbracket\text{strcpy}\rrbracket(\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2))$$



$$\begin{aligned} &\supseteq \\ &\{\mathfrak{M}[\text{strcpy}](\mu(\mathbf{m}_1), \mu(\mathbf{m}_2)) : \mu(\mathbf{m}_1) \in \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_1) \wedge \mu(\mathbf{m}_2) \in \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_2)\} \end{aligned}$$

`strcpy` of  $\mu(\mathbf{m}_1)$  and  $\mu(\mathbf{m}_2)$  returns, by definition of  $\mathfrak{M}$ ,  $\mu(\mathbf{m}_1)'$  where the first null-terminating memory block of  $\mu(\mathbf{m}_2)$  (including the null terminator), i.e., its string of interest, is embedded into  $\mu(\mathbf{m}_1)$  starting from the lower bound of  $\mu(\mathbf{m}_1)$ , if both  $\mu(\mathbf{m}_1)$  and  $\mu(\mathbf{m}_2)$  contain a well-formed string and the size condition on the destination character array value is fulfilled,  $\top_{\overline{\mathbf{M}}}$  otherwise. Then,  $\mu(\mathbf{m}_1)'$  belongs to  $\gamma_{\overline{\mathbf{M}}}(\mathfrak{M}_{\overline{\mathbf{M}}}[\text{strcpy}](\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2))$  because `strcpy` of  $\overline{\mathbf{m}}_1$  and  $\overline{\mathbf{m}}_2$ , by definition of  $\mathfrak{M}_{\overline{\mathbf{M}}}$ , is equal to  $\overline{\mathbf{m}}_1'$  that is  $\overline{\mathbf{m}}_1$  into which  $\overline{\mathbf{m}}_2$ 's has been embedded starting from the lower bound of  $\overline{\mathbf{m}}$ , if both  $\overline{\mathbf{m}}_1$  and  $\overline{\mathbf{m}}_2$  approximate character arrays which contain a well-formed string and the size condition on the destination segmentation abstract predicate is fulfilled;  $\top_{\overline{\mathbf{M}}}$  otherwise.

- Consider the unary operator  $\text{update}_{\bar{j}, \bar{v}}$  and let  $\overline{\mathbf{m}}$  be a split segmentation abstract predicate. We have to prove that:

$$\gamma_{\overline{\mathbf{M}}}(\mathfrak{M}_{\overline{\mathbf{M}}}[\text{update}_{\bar{j}, \bar{v}}](\overline{\mathbf{m}})) \supseteq \{\mathfrak{M}[\text{update}_{\bar{j}, \bar{v}}](\mu(\mathbf{m})) : \mu(\mathbf{m}) \in \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}})\}$$

$\text{update}_{\bar{j}, \bar{v}}$  of  $\mu(\mathbf{m})$  returns, by definition of  $\mathfrak{M}$ ,  $\mu(\mathbf{m})'$  that is  $\mu(\mathbf{m})$  where the character at position  $\bar{j}$  has been substituted with the character  $\bar{v}$ , if  $\bar{j}$  is a valid index for  $\mu(\mathbf{m})$ ;  $\top_{\overline{\mathbf{M}}}$  otherwise. Let  $\alpha_{\overline{\mathbf{B}}}(\bar{j}) = \bar{j}$  and  $\alpha_{\overline{\mathbf{C}}}(\bar{v}) = \bar{v}$ . Then,  $\mu(\mathbf{m})'$  belongs to  $\gamma_{\overline{\mathbf{M}}}(\mathfrak{M}_{\overline{\mathbf{M}}}[\text{update}_{\bar{j}, \bar{v}}](\overline{\mathbf{m}}))$  because  $\text{update}_{\bar{j}, \bar{v}}$  of  $\overline{\mathbf{m}}$ , by definition of  $\mathfrak{M}_{\overline{\mathbf{M}}}$ , is equal to  $\overline{\mathbf{m}}'$  that is  $\overline{\mathbf{m}}$  where the segment that is valid for  $\bar{j}$  is split so that the segment abstract predicate which occurs at position  $\bar{j}$  is substituted with  $\bar{v}$ , if  $\gamma_{\overline{\mathbf{B}}}(\bar{j})$  is equal to the singleton  $\{\bar{j}\}$  and  $\bar{j}$  is valid for  $\overline{\mathbf{m}}$ ;  $\top_{\overline{\mathbf{M}}}$  otherwise. □

**Theorem 3.4.**  $\mathfrak{P}_{\overline{\mathbf{M}}}$  is a sound over-approximation of  $\mathfrak{P}$ . Formally,

$$\gamma_{\overline{\mathbf{B}}}(\mathfrak{P}_{\overline{\mathbf{M}}}[\text{stm}](\overline{\mathbf{m}})) \supseteq \{\mathfrak{P}[\text{stm}](\mu(\mathbf{m})) : \mu(\mathbf{m}) \in \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}})\}$$

*Proof.*

Consider the binary operator `strcmp` and let  $\overline{\mathbf{m}}_1$  and  $\overline{\mathbf{m}}_2$  be two split segmentation abstract predicates. We have to prove that:

$$\begin{aligned} &\gamma_{\overline{\mathbf{C}}}(\mathfrak{P}_{\overline{\mathbf{M}}}[\text{strcmp}](\overline{\mathbf{m}}_1, \overline{\mathbf{m}}_2)) \\ &\supseteq \\ &\{\mathfrak{P}[\text{strcmp}](\mu(\mathbf{m}_1), \mu(\mathbf{m}_2)) : \mu(\mathbf{m}_1) \in \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_1) \wedge \mu(\mathbf{m}_2) \in \gamma_{\overline{\mathbf{M}}}(\overline{\mathbf{m}}_2)\} \end{aligned}$$

`strcmp` of  $\mu(\mathbf{m}_1)$  and  $\mu(\mathbf{m}_2)$  returns an integer value  $n$ , resulting from the difference between corresponding character array elements, denoting the lexicographic distance between

the strings of interest of  $\mu(m_1)$  and  $\mu(m_2)$ , if both contain a well-formed string,  $\top_Z$  otherwise, by definition of  $\mathfrak{P}$ . Then  $n$  belongs to  $\gamma_{\overline{C}}(\mathfrak{P}_{\overline{M}}[\text{strcmp}]((\overline{m}_1, \overline{m}_2)))$  because  $\text{strcmp}$  of  $\overline{m}_1$  and  $\overline{m}_2$ , by definition of  $\mathfrak{P}_{\overline{M}}$ , is equal to  $\bar{n}$  that is the difference between corresponding segment abstract predicates, denoting the lexicographic distance between  $\overline{m}_1.s$  and  $\overline{m}_2.s$ , if  $\overline{m}_1$  and  $\overline{m}_2$  are comparable, both approximate character arrays which contain a well-formed string where  $\top_{\overline{C}}$  does not occur;  $\top_Z$  otherwise.

□

**Theorem 3.5.**  $\mathcal{L}_{\overline{M}}$  is a sound over-approximation of  $\mathcal{L}$ . Formally,

$$\gamma_{\overline{B}}(\mathcal{L}_{\overline{M}}[\text{stm}]((\overline{m}))) \supseteq \{\mathcal{L}[\text{stm}](\mu(m)) : \mu(m) \in \gamma_{\overline{M}}(\overline{m})\}$$

*Proof.*

Consider the unary operator  $\text{strlen}$  and let  $\overline{m}$  be a split segmentation abstract predicate. We have to prove that:

$$\gamma_{\overline{B}}(\mathcal{L}_{\overline{M}}[\text{strlen}]((\overline{m}))) \supseteq \{\mathcal{L}[\text{strlen}](\mu(m)) : \mu(m) \in \gamma_{\overline{M}}(\overline{m})\}$$

$\text{strlen}$  of  $\mu(m)$  returns, by definition of  $\mathcal{L}$ , an integer value  $n$  which denotes the length of the sequence of character before the first null one in  $\mu(m)$ , if  $\mu(m)$  contains a well-formed string;  $\top_Z$  otherwise. Then  $n$  belongs to  $\gamma_{\overline{B}}(\mathcal{L}_{\overline{M}}[\text{strlen}]((\overline{m})))$  because  $\text{strlen}$  of  $(\overline{m})$ , by definition of  $\mathcal{L}_{\overline{M}}$  is equal to the difference between the lower and the upper bound of  $\overline{m}.s$  if  $\overline{m}$  approximates character arrays which contain a well-formed string of interest;  $\top_{\overline{B}}$  otherwise.

□

### 3.6 Program Abstraction

Adapting M-String to the analysis of real-world C programs requires, first of all, a procedure that identifies string operations automatically. A subset of these operations must be performed using abstract functions, carried out on a suitable abstract representation. The technique that captures this approach is known as abstract interpretation. A typical implementation is based on an interpreter in the programming language sense: it executes the program by directly performing the operations written down in the source code. However, rather than using concrete values and concrete operations on those values, part (or the entirety) of the computation is performed in an *abstract domain*, which over approximates the semantics of the concrete program.

In this chapter, we mainly focus on string abstraction. Therefore we will interpret the program's portions that do not make use of strings without abstracting values. We only apply abstraction to strings that, within the program, are manipulated by string operations. When the program deals with string variables that exhibit minimal variation,

e.g., string literals, the M-String representation will provide no benefit. Instead, it could either hurt performance or it may introduce spurious counterexamples.

Based on the considerations above, it is clear that it is beneficial to reuse and refactor existing tools that implement abstract verification in a modular way on explicit programs. A compilation-based abstraction design that follows this approach was introduced and implemented in [120]. However, this tool is designed to abstract scalar values only. This is why we need to extend it to operate with more sophisticated domains that represent more complex objects, such as strings.

In the rest of this section, we will first summarize the general approach to abstraction as a program transformation. In Section 3.6.3, we explore the implications of aggregate (as opposed to scalar) domains within this framework. Sections 3.6.4 and 3.6.5 discuss the semantic (run-time) aspects of the abstraction and which operations we consider as primitives of the abstraction.

### 3.6.1 Compilation-Based Approach

Instead of (re-)interpreting instructions abstractly, in a compilation-based approach, abstract instructions are transformed into an equivalent explicit code that implements the abstract computation. The transformation takes place before the analysis of the program (e.g., model checking) during the compilation process.

Consequently, the analysis processes the program without needing special knowledge of the abstract domains in use, as the abstraction is encoded directly in the program. Figure 3.1 depicts a comparison of the compilation-based approach with respect to the interpretation-based approach adopted by more conventional abstract interpreters.

In a compilation-based approach, two different abstraction perspectives are considered:

1. *static*, referring to the syntax and the type system,
2. *dynamic*, or semantic, referring to execution and values.

The LART tool performs syntactic (*static*) abstraction on LLVM bitcode [119]. Syntactic abstraction replaces some of the LLVM instructions that occur in the program with their abstract counterparts, as depicted in Figure 3.2.

### 3.6.2 Syntactic Abstraction

The first step of program abstraction performed by LART is a syntactic abstraction. Syntactic abstraction replaces LLVM instructions or whole functions with their abstract counterparts. Since we do not want to perform all operations abstractly, we need to classify only those that might be applied to abstract values. The abstract values emerge in the program as input values. From these values, LART computes all the operations that might deal with abstract values through a combination of data flow and alias analysis. Finally, as a

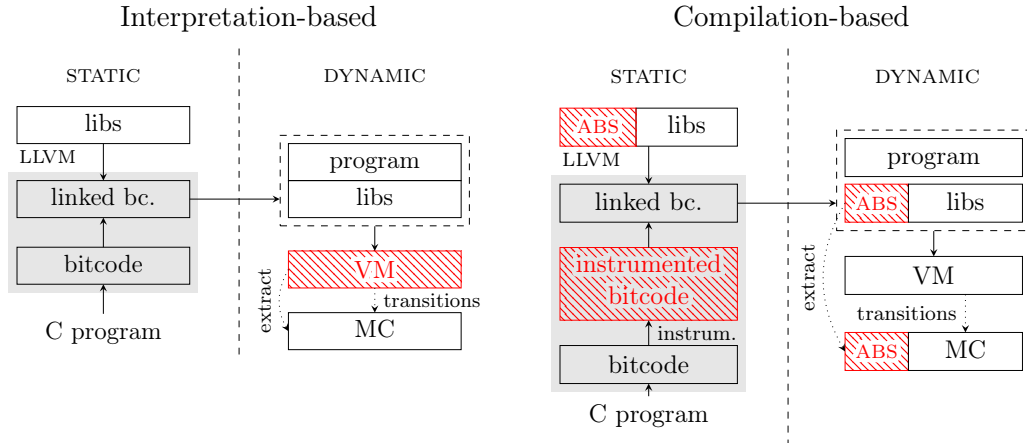


FIGURE 3.1: Comparison of an abstract interpretation and a compilation-based approach. In the interpretation-based approach, the whole abstract interpretation is performed at runtime. The bitcode operations are interpreted abstractly by the virtual machine (VM) that maintains an abstract state. In this way, an abstract state-space is generated for a model-checking algorithm (MC). The compilation-based approach is different. The abstract operations are instrumented into the compiled program and their implementation is provided as a library. Then, the virtual machine executes the instrumented program as a regular bitcode [120].

```

char *a = input_string();    abstract a = abstract_string();
char *b = string();         char *b = string();
char *c = strcat(a, b);     abstract c = abstract_strcat(a, lift(b));
int l = strlen(c);         abstract l = abstract_strlen(c);

```

FIGURE 3.2: Syntactic abstraction.

result of the analysis, LART obtains a set of possibly abstract operations that are replaced by their abstract equivalents, e.g., `strcat`, `strlen` are replaced by `abstract_strcat` and `abstract_strlen`. Abstract operations then implement the manipulation of abstract values, in our case with M-Strings as described in Section 3.4. In other words the specific meaning of abstract instructions and abstract values then defines the semantic abstraction.

For the precise formulation of syntactic abstraction, we take advantage of the static type system of LLVM. We leverage the fact that we can assign to each variable its type, which is either concrete or abstract. In this way, we can precisely set a boundary between concrete and abstract values.

We consider a simplified version of LLVM. It defines a set of *concrete scalar types*  $S$ . The set of all possible types is given by a map  $\Gamma$  that inductively defines all finite (non-recursive) algebraic types over the set of given scalars. To be precise, the set of types  $\Gamma(T)$

derived from a set of scalars  $T$  is as follows:

1.  $T \subseteq \Gamma(T)$ , meaning each scalar type is included in  $\Gamma(T)$ ,
2. if  $t_1, \dots, t_n \in \Gamma(T)$  then also the *product type* is in  $\Gamma(T)$ :  $(t_1, \dots, t_n) \in \Gamma(T)$ ,  $n \in \mathbb{N}$ ,
3. if  $t_1, \dots, t_n \in \Gamma(T)$  then also *disjoint union* is in  $\Gamma(T)$ :  $t_1 \mid t_2 \mid \dots \mid t_n \in \Gamma(T)$ ,  $n \in \mathbb{N}$ ,
4. if  $t \in \Gamma(T)$  then  $t^* \in \Gamma(T)$ , where  $t^*$  denotes pointer type.

In a concrete LLVM program, the set of admissible types comprise those derived from concrete scalars  $S$ , i.e.,  $\Gamma(S)$ . In syntactic abstraction, we need to extend admissible types by abstract types. From these, we generate all possible types using  $\Gamma$ . Depending on the type of abstraction, we use a different set of basic abstract types. In the case of scalar abstraction, a set of basic abstract types contains abstract scalar types  $\bar{S}$ . Correspondence between abstract and concrete scalars is given by a bijective map  $\Lambda : S \rightarrow \bar{S}$ . Finally, each value, which exists in the abstracted program, has an assigned type of  $\Gamma(S \cup \bar{S})$ . Specifically, this implies that the abstraction works with *mixed types*—products and unions might contain both concrete and abstract fields. Moreover, it is possible to create pointers to both abstract or mixed values.

### 3.6.3 Aggregate Domains

In addition to scalar values that cannot be further decomposed, programs typically operate with more complex data, which can be seen as compositions—aggregates—of multiple scalar values. Depending on aggregates' nature, we can classify them as aggregates that contain a variable number of items (arrays), records that contain a fixed number of items in a fixed layout, where each of these can be of a different type. The items in such aggregates can be (and often are) scalars. However, more complex aggregates are also possible: arrays of records, records which in turn contain other records, and so on.

While scalar domains only dealt with simple values, in aggregate abstraction, we consider composite data in the spirit of the above definition. Similar to scalar domains, abstract aggregate domains approximate concrete aggregate values by describing a particular set of aggregate properties. For example, we can describe a set of aggregates by their length or a set of values that appear in the aggregate. In the M-String, we keep properties in the form of segmentation, where segments are further abstracted by bounds and characters. Values in an aggregate domain then keep the representation of chosen properties and operations updates them. For instance, consider an array length property domain—the domain operations in such a case operate only with lengths of arrays, e.g., *abstract concat* of arrays adds together lengths of its arguments (abstract arrays).

In general, aggregate domains can provide arbitrary operations. However, two operations are, in some sense, universal, being elementary memory manipulation operations,

namely: byte-wise **access** and **update** of the aggregate. The universality of these operations originates from the fact that all aggregate operations can be represented as accesses and updates. In a low-level representation of a program (assembly), they usually are presented in this form. LLVM allows a slightly higher level of manipulation to access and update individual scalars present in the aggregates (as opposed to bytes). For M-String, though, this distinction is not essential because the scalars stored in C strings are individual bytes (characters). All other operations are present in the form of sequences of elementary instructions—possibly encapsulated in functions. Moreover, as in concrete programs, the **access** and **update** represent an interface between scalars and memory, in the abstraction they form an interface between scalar and aggregate domains (even in the case of byte-oriented access since bytes are also scalars). We refer the reader to the Section 3.5.2 for abstract semantics of **access** and **update**.

In comparison to scalar abstraction, the syntactic abstraction of aggregates does not operate directly on aggregate types. In LLVM, aggregate values are usually represented by a pointer to the underlying aggregate type. Therefore all the accesses and updates are made through the pointers to the aggregates. For instance, strings are represented as a pointer to a character array. We need to take this fact into account when we perform the syntactic abstraction. In the analysis, we consider the pointers to aggregates as base types for the abstraction. In the case of arrays, the base types are concrete pointers to those arrays: we call them  $P^*$ , where  $P^* \subseteq \Gamma(S)$ . A set of abstract pointers types  $P^*$  then describes types of abstracted aggregates (arrays). As for scalar domains, we define a natural correspondence between pointers to concrete values and abstract aggregates as a bijective map  $\Lambda : P^* \rightarrow \overline{P^*}$ . For instance, in the case of M-String abstraction, the map  $\Lambda$  assigns to **char\*** a type of M-String value. Finally, we allow all the mixed types generated from scalars and abstract aggregates:  $\Gamma(S \cup \overline{P^*})$ .

Observe that pointers, in general, also in LLVM maintain two pieces of information about memory location: they represent both the memory *object* and an *offset* into that object. In particular, our implementation treats the first 32-bits of the pointer as an object identifier and the last 32-bits as its offset. This distinction is not very relevant in explicit programs because those two components are represented uniformly in a single value and often, they cannot be distinguished at all. However, the distinction becomes relevant when dealing with abstract aggregate values. In fact, in this case, the *object* component of the pointer is concrete as it determines a single specific abstract object. Instead, the *offset* component may or may not be concrete. The choice depends on the specific abstract aggregate domain: it may be more advantageous representing the offset in an abstract way, i.e., by a 32-bit abstract scalar value. Observe that a memory access through such a pointer needs to be treated in both cases as an abstract **access** or **update** operation.

The two basic memory access operations In LLVM, **load** and **store**, correspond to

the `access` and `update` operations. It is important to note that memory access is always explicit: memory is never directly used in a computation. This observation is used in the design of aggregate abstraction, where we can assume that the access to the content of an aggregate will always go through a pointer associated with the abstract object.

### 3.6.4 Semantic Abstraction

In syntactic abstraction, we dealt with operations' syntax, their types, and the types of values and variables. It described how LART performs a source-to-source transformation. In contrast, semantic abstraction concerns the values computed at runtime by a program. It defines how abstract operations modify values and how to transfer between concrete and abstract values. Therefore, similarly to syntactic abstraction that defined the maps  $\Lambda$  and  $\Lambda^{-1}$  to transfer between concrete and abstract *types*, the semantic abstraction makes use of `lift` and `lower` (cf. Definitions 3.10 and 3.11): operations (instructions) converting values between their concrete and abstract representations. They realize a runtime implementation of domain functions: abstraction ( $\alpha_{\overline{M}}$  in the case of M-String) and concretization ( $\gamma_{\overline{M}}$ ).

The `lift` operation implements abstraction of concrete values by a single over-approximating abstract value. For example, in Figure 3.2, on line 3 of the abstracted program, a concrete string `b` is lifted to the abstract domain. This allows performing `abstract_strcat` in a single abstract domain. In other words, operations do not need to consider concrete values because all their arguments are lifted to the abstract domain. This simplifies the implementation of a domain and reduces the number of possible domain interactions. While  $\Lambda$ , which was a purely syntactic construct, `lift` and `lower` accomplish the actual conversion of values between domains during program runtime. During program execution, lowering an abstract value into multiple concrete values can be seen as a nondeterministic branching in the program (and the `lower` operator is indeed based on a non-deterministic choice operator).<sup>7</sup> For further details of the program transformation performed by LART, we kindly refer the reader to [120].

### 3.6.5 Abstract Operations

As a result of syntactic abstraction, we obtain a program that temporarily contains abstract operations. These operations take abstract values as operands and return abstract values as a result. Though, after the program transformation, the resulting program is required to be a semantically valid LLVM bitcode. Therefore, we demand that each abstract operation can be realized as a sequence of concrete instructions. This allows us to obtain an abstract

---

<sup>7</sup>In a model checker, the non-deterministic choice would be typically implemented as branching in the state space (and the consequences of all possible outcomes would be explored). In a testing context, however, the choice might be implemented by choosing one particular path at random.

program that does not contain any abstract operations and executes it using standard (concrete, explicit) methods.

Thoroughly, syntactic abstraction substitutes concrete operations with their abstract counterparts: an operation with type  $(t_1, \dots, t_n) \rightarrow t_r$  is substituted by an abstract operation of type  $(\Lambda(t_1), \dots, \Lambda(t_n)) \rightarrow \Lambda(t_r)$ . Furthermore, transformation inserts lift and lower operations as needed, typically, in places where concrete values are operands of abstract operations. The implementation is free to select the operations to be abstracted and where value lifting and lowering be inserted, so long type constraints are satisfied. However, it tends to minimize the number of abstracted operations.

In addition to LLVM instructions, the M-String abstraction requires the transformation to abstract function calls to standard library functions, e.g., `strcmp`, `strcat`. From the perspective of syntactic abstraction, we can treat function calls as single atomic operations that take abstract values and produce abstract results. Hence, the transformation substitutes them in the same way as instructions: for instance `strcmp` operation of type  $(m, m) \rightarrow s$  is replaced by `abstract_strcmp` of type  $(\Lambda(m), \Lambda(m)) \rightarrow \Lambda(s)$  where  $m$  is a concrete character array and  $s$  is a concrete scalar result of the string comparison. Afterwards, all abstract operations are implemented by using concrete subroutines (implementation of abstract semantics). For details, see [120].

Observe that, as an alternative approach, the standard library functions `strcat`, `strcmp`, etc. could have been transformed instruction by instruction, by using abstract access and update of a content only. However, the price to pay would have been losing a certain degree of accuracy in the abstraction, the exact amount depending on the single operation.

### 3.7 Instantiating M-String

As an aggregate domain, M-String is parametrizable by scalar domains of characters and indices (bounds). This allows us to tailor the abstraction to the needs of the analysis of string values. Depending on the precision of chosen domains, the instance of the M-String domain will inherit their properties. With more precise domains, the M-String values will maintain a higher granularity of segmentation. But, simpler character representation will decrease the segmentation granularity for the cost of a higher rate of false alarms.

A particular instance of M-String is automatically derived from a parametric description given in Section 3.5, provided a suitable scalar domain  $\overline{\mathbf{C}}$  for characters and scalar domain  $\overline{\mathbf{B}}$  to represent segment bounds. The instantiation demands that both scalar domains  $\overline{\mathbf{C}}$  and  $\overline{\mathbf{B}}$  are equipped with operations that appear in the operations with the segmentation. These are mainly elementary arithmetic and relational operations. In the implementation, we provide an M-String domain template that automatically derives all the operations from provided scalar domains.



### 3.7.1 Symbolic Scalar Values

In program verification, it is common practice to represent certain values symbolically (for instance, inputs from the environment). The symbolic representation allows the verifier to consider all admissible values with a reasonably small overhead. In DIVINE, symbolic verification is implemented using a similar abstraction to one described in the previous section: symbolic scalar values represent their content by SMT formula expressions (terms) in the form of abstract syntax trees. The input values are represented as unconstrained variables in the bit vector logic. Operations then build formulae trees from their arguments. In addition to these so-called data definitions, symbolic representation also maintains one global formula of constraints (path-condition), which is derived from the control flow of the program. A more detailed description of this symbolic representation is presented in [120].

The domain of symbolic values (we call it a term domain) requires DIVINE to be augmented with an SMT solver from a suitable theory.<sup>8</sup> DIVINE uses the solver to detect computations that have reached the bottom of the term domain (those are the infeasible paths through the program). Furthermore, as a model checker, it needs to identify equal states or whether the state subsumes another one. This is achieved by the equivalence check of corresponding formulae. With these prerequisites, the symbolic representation in joint with the bit-vector theory is a precise abstraction (i.e., it is not an approximation but models the program state faithfully).

### 3.7.2 Concrete Characters, Symbolic Bounds

In the evaluation, we instantiate the M-String domain in two ways. The first simpler instantiation sets the domain of characters  $\overline{\mathbf{C}}$  to be the concrete domain (i.e., we let the characters be represented by themselves). We let the domain of segment bounds  $\overline{\mathbf{B}}$  to be a symbolic 32b integers. This instantiation balances between simplicity on the one hand (both domains we used for parameters were already present in DIVINE) and the ability to describe strings with undetermined length and structure.

At the implementation level (described in more detail in the following section), the domain remains generic: the particular domains we picked can be easily substituted by other domains. Compared to the theoretical description of M-String, the implementation uses a slightly simplified representation of segmentation made of a pair of arrays (cf. Figure 3.3). The elements of these arrays are characters and bounds, whose type is derived from parametrization, i.e., from the scalar domains  $\overline{\mathbf{C}}$  and  $\overline{\mathbf{B}}$ . The modification of the representation is just for optimizing the implementation and does not affect the operations' semantics. The analysis with this representation is presented in Example 3.14.

---

<sup>8</sup>For scalars in C programs, we use the bitvector theory.

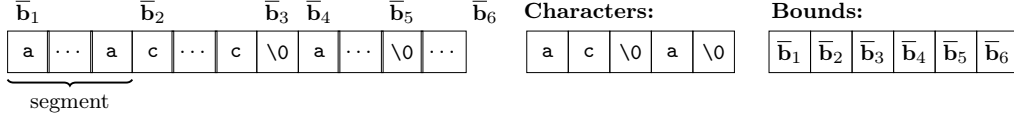


FIGURE 3.3: M-String value with symbolic bounds, where string of interest is from  $\bar{b}_1$  to  $\bar{b}_3$ .

This instantiation of M-String is particularly suitable for representing strings with sequences of a single character of variable length, i.e., the strings of the form  $a^k b^l c^m \dots$  where relationships between  $k, l, m, \dots$  can be specified using standard arithmetic and relational operators and each of  $a, b, c$  is a concrete letter. In turn, this allows M-String to be used for the analysis of program behavior on broad classes of input strings described this way. A more detailed description of this approach can be found in Section 3.8.

**Example 3.14.** Simple program analysis with symbolic bounds and concrete characters:

```

mstring str = abstract_string(x,  $\bar{b}_1$ , \0,  $\bar{b}_2$ , y,  $\bar{b}_3$ , \0,  $\bar{b}_4$ );
symbolic idx = abstract_int();
if (idx <  $\bar{b}_4$ ) {
    str[idx] = 'y';
    symbolic len = abstract_strlen(str);
}

```

Imagine we are given symbolic bounds  $\bar{b}_1 < \bar{b}_2 < \bar{b}_3 < \bar{b}_4$ , then the first line of the transformed program creates `mstring` value with characters  $[x, \0, y, \0]$  and bounds  $[0, \bar{b}_1, \bar{b}_2, \bar{b}_3, \bar{b}_4]$ . In the following, we describe `mstring` values as pairs of these two arrays. The second line creates a symbolic index of arbitrary value. On line 3, the program constraints the index to be smaller than the `mstring` maximal length. Otherwise, the update on the next line would yield an error. Next the program assigns to the position of abstract index a character  $y$ . The assignment is implemented as update operation on `mstring` value. Depending on the value of the `idx`, the operations results in the following strings  $\overline{str}_x$ , as result we join all possibilities:

1. if  $idx < \bar{b}_1$  :  $idx$  falls in the first segment:  $\overline{str}_1 = ([x, y, x, \0, y, \0], [0, idx, idx + 1, \bar{b}_1, \bar{b}_2, \bar{b}_3, \bar{b}_4])$  and creates a new segment between  $idx$  and  $idx + 1$  containing character  $y$ . Notice that if  $idx = 0$  the first segment is empty, similarly the third segment for  $idx + 1 = \bar{b}_1$ . The string of interest for  $\overline{str}_1$  is of form  $x^{idx} y^1 x^{\bar{b}_1 - idx - 1}$ .
2. if  $\bar{b}_1 \leq idx < \bar{b}_2$  : than  $\overline{str}_2 = ([x, \0, y, \0, y, \0], [0, \bar{b}_1, idx, idx + 1, \bar{b}_2, \bar{b}_3, \bar{b}_4])$ , with string of interest as join of following forms:
  - if the update is performed right after the first segment, i.e.,  $idx = \bar{b}_1$ :
    - if and  $|\bar{b}_1 - \bar{b}_2| > 1$ , i.e., the segment of zeros contains more elements, then the string has form  $x^{\bar{b}_1} y$ ,

- otherwise the update overwrites the single zero character, hence extends the string of interest by segment of  $y$  characters:  $x^{\bar{b}_1}y^{\bar{b}_3-\bar{b}_1}$ .
  - otherwise between first segment and  $\text{idx}$  is a terminating zero, hence the string of interest remains unchanged:  $x^{\bar{b}_1}$ .
3. if  $\bar{b}_2 \leq \text{idx} < \bar{b}_3$  : then  $\overline{\text{str}}_3 = \overline{\text{str}}$ , because update stores the same character as is already present in the segment.
  4. if  $\bar{b}_3 \leq \text{idx} < \bar{b}_4$  : then update creates a new segment inside of sequence of last zeros:  $\overline{\text{str}}_4 = ([x, \backslash 0, y, \backslash 0, y, \backslash 0], [0, \bar{b}_1, \bar{b}_2, \bar{b}_3, \text{idx}, \text{idx} + 1, \bar{b}_4])$ .

Consequently, the `abstract_strlen` operation on the last line of the program computes the join of all possible lengths of strings of interest, i.e.,  $\bar{b}_1 \sqcup_{\mathbb{B}} \bar{b}_3$ .

◻

### 3.7.3 Symbolic Characters, Symbolic Bounds

The second instantiation is used in benchmarks, where the computation with M-String values encountered abstract scalars (characters). This occurs when the program obtains some character as input from the environment and tries to store it into the M-String value. Therefore, we instantiated the M-String domain with an abstract representation of characters by setting the domain  $\overline{\mathbb{C}}$  to be the term domain, which keeps track of symbolic 8b bitvectors (characters in  $\mathbb{C}$  language). In this way, we do not need to lower abstract characters before storing them in the M-Strings, which was needed for the concrete domain used in the previous instantiation. However, we pay the price for more expensive computation with symbolic characters.

### 3.7.4 Implementation

Finally, we implemented the M-String abstraction as a LART domain.<sup>9</sup> The LART domain is a C++ library that implements abstract semantics of M-String operations presented in Section 3.5. This library is then linked to the transformed program allowing the program to perform abstract analysis with model-checker DIVINE. An abstract domain definition in LART consists of a C++ class that describes both the representation (in terms of data) and the operations (in terms of code) of the abstract domain.

In the case of M-String domain, this class contains 2 attributes: an array of *bounds* and an array of *characters*, as outlined in Section 3.7.2 and depicted in Figure 3.3. The class has two type parameters: the domain to use for representing segment bounds and the domain to represent individual characters (i.e., the content of segments). A specific instantiation

---

<sup>9</sup>The implementation with examples and documentation of domain usage can be found online on the supplementary page: <https://divine.fi.muni.cz/2020/mstring>

is then automatically derived by the C++ compiler from the classes representing the type parameters and the parametric class, representing M-String values.

As a minimal set of operations, the M-String domain implements all requisite aggregate operations: these are `lift`, `update` and `access`. Furthermore, the implementation provides an optimized version of string operations described in Sections 3.5: `strlen`, `strcpy`, `strcat`, `strcmp` and `strchr`. These operations reduce the loss of abstraction precision that would arise if only the abstraction of accesses and updates from strings were used.

Since C strings are stored as shared, mutable character arrays, the implementation of the M-String domain reflects the sharing semantics of these arrays. If multiple pointers exist into the same abstract string, modifications through one of them must also be visible when the string is accessed through another pointer. Moreover, the pointers do not have to be equal: they may point to different suffixes of the same string. Therefore, the representation of pointers to abstract strings must treat the *object* and the *offset* components separately (see also Section 3.6.3), and the representation of the *offset* component must be compatible with the bound domain  $\bar{\mathbf{B}}$ .

## 3.8 Experimental Evaluation

In the evaluation, we chose a few scenarios to demonstrate the properties of the abstraction. In the first scenario, we show that using abstract versions of standard functions is more efficient than if concrete versions were transformed using only abstract string accesses and updates. The second scenario investigates several implementations of standard library functions: we transform them automatically in the means of accesses and updates, and we show that their results agree with results generated by M-String library operations. In the third scenario, we evaluate M-String instantiation with symbolic characters on the set of benchmarks from real software that contain buffer-overflow errors. Here we show that M-String can efficiently detect real-world bugs as well as prove that a program does not contain them after they are fixed. The last benchmark shows the use of abstractions on more complex C programs. As an example, we analyze automatically generated parsers from `bison` and `flex` tools on abstract (M-String) inputs. The resource limits for all scenarios were the same: each verification run was limited to 4 processing units (cores), 80 GB of memory, and 1 hour of CPU time. The processor used to run benchmarks was AMD EPYC 7371 clocked at 2.60GHz.

### M-String Operations

The first group of benchmarks focuses on the use of resources by abstraction. Benchmarks compare the effectiveness of abstract domain operations with the automatically abstracted implementation of standard library functions from `PDCLib`, a public-domain `libc` implementation, using only essential abstract operations: `lift`, `update` and `access`. The results

	Word						Sequence					
	Verification(s)						Verification(s)					
	States	8	64	1024	4096	LART(s)	States	8	64	1024	4096	LART(s)
<code>strcmp</code>	3562	480	498	472	481	1.70	70	0.26	0.24	0.21	0.25	1.76
<code>strcpy</code>	368	9.8	9.1	9.3	9.4	1.70	48	0.20	0.20	0.21	0.20	1.71
<code>strcat</code>	7398	898	873	865	843	1.72	105	0.51	0.52	0.53	0.51	1.72
<code>strchr</code>	49	0.3	0.4	0.3	0.3	1.71	15	0.04	0.04	0.03	0.04	1.70
<code>strlen</code>	78	1.1	1.2	1.0	1.3	1.70	16	0.05	0.04	0.05	0.06	1.81

TABLE 3.2: Measurements of M-String operations on two types of inputs: *Word* and *Sequence* described in Section 3.8. Each benchmark measures a size of state space and verification time for input M-Strings of a given length. Lastly, the table shows an average transformation time (LART). All measurements of time are in seconds. The size of state space does not change for different lengths of input—for more details, see discussion in Section 3.8.

	8		64		1024	
	Time(s)	States	Time(s)	States	Time(s)	States
<code>strcmp</code>	1.24	197	260	1597	T	–
<code>strcpy</code>	0.7	122	61.5	962	T	–
<code>strcat</code>	15.8	1102	T	–	T	–
<code>strchr</code>	0.04	16	0.05	16	0.05	16
<code>strlen</code>	0.19	46	9.57	326	T	–

TABLE 3.3: Benchmark of standard library functions abstracted using only the M-String definitions of *access* and *update* operations for *Sequence* inputs of size 8, 64 and 1024 characters. Verification for *Word* strings times out in most of the instances.

depicted in Table 3.2 were measured with parametrized M-String inputs of two kinds ( $l$  is a parametric length of the input):

- *Word*  $w$  is a string of the form:  $w = c_1^{i_1} \cdot c_2^{i_2} \cdot \dots \cdot c_l^{i_l}$  where  $\sum_{k=1}^l i_k \leq l$  and  $c_x$  is an arbitrary character from domain  $\overline{\mathbf{C}}$ .
- *Sequence*  $w$  is a string of the form  $w = c^i$ , where  $i \leq l$  and  $c$  is a character from domain  $\overline{\mathbf{C}}$ .

For each standard library function and input type, we created an isolated benchmark in two variants: one using an abstract semantics of M-String operations (see Table 3.2) and the other variant (Table 3.3) only with an automatic abstraction of essential aggregate operations.

The first notable difference between automatically abstracted implementations of library functions and M-String operations is that the analysis of the former timeouts for input strings longer than 64 characters. The main cause of the lifted implementation’s inefficiency is that it iterates over all characters, while M-String operations leverage iteration

	Word						Sequence					
	4		8		16		4		8		16	
	Time(s)	States	Time(s)	States	Time(s)	States	Time(s)	States	Time(s)	States	Time(s)	States
<code>strcmp</code>	14.3	1005	105	2989	1350	9741	2.17	204	5.09	376	16.5	720
<code>strcpy</code>	5.15	515	57.4	1823	912	6935	0.83	183	2.49	347	9.14	675
<code>strcat</code>	468	5748	T	–	T	–	8.56	751	113	2535	1940	9463
<code>strchr</code>	0.08	22	0.08	22	0.08	22	0.3	17	0.3	17	0.4	17
<code>strlen</code>	0.66	91	4.13	259	68.8	883	0.15	34	0.28	54	0.65	94

TABLE 3.4: Verification results of functions from `PDCLib` with time-out of 1 hour. Measurements show the size of state space and verification time for the parametric length of the input.

over larger segments. This difference also causes a blow-up of the model checker’s state space for the lifted implementations while the state space size does not change for M-String operations. The reason for this is that the number of segments does not change with the length of the input. Therefore M-String operations always perform the same computation independently of the M-String length.

### C Standard Libraries

In the second set of benchmarks, we investigate whether the implementation from several standard libraries matches the expected abstract implementation results. In other words, we perform an equivalence check of results obtained from M-String operations with the results of the automatically abstracted (originally concrete) standard library functions. We expect that both give the same results. For the evaluation, we picked three open-source libraries: `PDCLib`, `musl-libc` and `μCLibc`. Since the results for the libraries are rather similar, we present here only an evaluation of `PDCLib` functions. The remaining results are provided in the Supplementary Material. All benchmarks showed that our implementation matches the standard one.

Similarly, as in the previous case, these benchmarks suffer from the state space blow up caused by an exponential number of possible character combinations. For this reason, we decreased the size of the input strings. In addition to large state space, many string accesses and updates of concrete implementations result in a large SMT formulae, causing a long time spent in solvers.

Furthermore, notice that the computation analysis with *Word* input, which has more segments, results in longer execution times than the analysis with *Sequence*. The reason is that the more segments naturally also causes overhead for the analyses. For example, the M-String needs to consider cases when some segments have zero length: this causes a hard SMT queries because, in the worst case, it needs to check all possible strings for given segment bounds and characters.

### Veriabs Overflow Benchmarks

In this scenario, we show that the domain is capable of efficient overflow bugs finding. Veriabs benchmarks exhibit overflow errors and fixed variants of real-world software. To soundly prove correctness of these benchmarks, we instantiate M-string with term domain also for characters. Hence we can reason about arbitrary strings of a symbolic length. However, a drawback of this instantiation is that whenever the length of the string bounds a loop, we might have to unroll the loop infinitely in the analysis—these cases timeouts in the correct benchmarks.

	Correct		Error Found		Timeout
	Tests	Time(s)	Tests	Time(s)	
apache	0	–	26	384.26	24
openser	43	234.13	45	105.93	6
wu-ftp	8	35.78	14	2461.27	19
libgd	4	9.01	4	1.85	0
madwifi	5	0.51	5	0.55	0
gxine	1	0.53	1	0.25	0

TABLE 3.5: Veriabs overflow benchmarks depict a few categories of programs exhibiting an overflow error and their fixed variants. The table shows the number of solved benchmarks (*tests*) and accumulated time for each category. For each category, we correctly depict verified benchmarks, benchmarks where the verifier finds an error, and the number of timeouts.

### Parsers

Lastly, we evaluate our implementation on more complex programs: automatically generated parsers. For the generation, we use a tool `Bison`. It reads a language specification in the form of context-free grammar and produces a C parser that accepts the language. In the benchmarks, we generate two such parsers. The first one accepts a language of numerical expressions (mathematical expressions that consist of numbers and binary operators). The second parser is for a simple programming language with variables and branching. We present an evaluation for both parsers in Table 3.6. As with the previous benchmark sets, the M-String inputs with a smaller number of segments outperformed other analyses.

Numeric Expressions Grammar						
	10		20		35	
	Time(s)	States	Time(s)	States	Time(s)	States
add	40.2	416	319	3548	T	–
ones	5.54	62	8.12	196	189	2186
alter	708	105	1582	11k	T	–

TABLE 3.6: Measurements of time and size of state space for analyses of automatically generated parsers. In these benchmarks, we use specifically hand-crafted M-String inputs for parsers. For parsing of mathematical expressions, it was: **addition** input had a form of two arbitrary numbers with a plus sign between them, **ones** was a simple input of a single digit sequence, and lastly, **alternation** was input that produced complicated M-Strings by alternating digits inside of expressions. The other parser of simple programming language was evaluated on: **value** was in input that created a variable and assigned a constant to it, **loop** was a short program with some control flow and **wrong** was a program that contained a syntax error.

Simple Programming Language Grammar						
	10		100		1000	
	Time(s)	States	Time(s)	States	Time(s)	States
value	6.58	38	90.4	488	1100	4988
loop	1.53	23	4.88	23	33.3	23
wrong	7.34	82	67.7	892	311	8992

TABLE 3.6: *cont*

### 3.9 Discussion

A new segmentation-based abstract domain for approximating C strings has been introduced, whose main novelty lies in abstracting both index bounds and substrings while managing strings as a pair of two string buffers: the string of interest itself, and a tail of allocated and possibly initialized but unused memory.

The presented approach enables more precise modelling of the functions in the standard C library for strings, also considering the known weaknesses for the management of terminating null characters and buffer bounds. The M-string domain results effective for identifying security leaks caused by string manipulation errors, e.g., buffer overflows.



After theoretically describing the domain and the basic operations on strings, we have implemented (using C++ language) the abstract semantics, combining them with a tool that starting from string-manipulating C programs lifts them to the M-String domain.

Our experimental results have also focused on tuning the parameters of M-String (the domains for both segment content and segment bounds) by instantiating them by both concrete and symbolic characters and by symbolic (bitvector) bounds.



## Chapter 4

---

# COMBINING STRING DOMAINS

In this chapter, we specialize the segmentation domain `FunArray`, introduced in Section 3.2, for array content analysis [57], and we combine it, through reduced product, with less sophisticated existing string abstract domains to improve the precision of the analysis. Since the M-String abstract domain (cf. Section 3) is inspired by `FunArray`, and it is tailored for string analysis of `C` programs, in the following, we use `FunArray` “as is”, and we combine it with general-purpose existing string abstract domains.

The contribution of this chapter is submitted for publication.

### Chapter Structure

Section 4.1 highlights some limits of existing string abstract domains and explains our contribution. Section 4.2 introduces the syntax of some string operations of interest. Section 4.3 defines the concrete domain and semantics. Section 4.4 introduces the basic abstract domains for string analysis. Section 4.5 specializes the Segmentation domain [57], summarised in Section 3.2, to the approximation of strings. Section 4.6 defines our refined string abstract domains. Section 4.7 concludes.

### 4.1 Introduction

Existing string abstract domains partially track the content and/or the shape information of string values. We now consider some basic domains for strings [49]. The String Length domain approximates the length of a string through an interval (e.g., the interval  $[2, 5]$  represents all the possible strings whose length ranges from 2 to 5). The Character Inclusion domain approximates strings through a pair of sets highlighting the characters definitely and possibly contained (e.g., the pair of sets  $(\{a, b\}, \{a, b, c\})$  represents all the possible strings definitely containing the characters ‘a’ and ‘b’ and possibly containing the character ‘c’). The String Length and the Character Inclusion domains approximate string shape and string content information, respectively. Instead, the Prefix domain partially detects the shape and content information of the strings it approximates. For example, the prefix `ab*`

represents all the possible strings sharing the starting sequence “ab”, with minimal length 2. The Segmentation domain [57], summarised in Section 3.2, is a more sophisticated abstract domain that was originally proposed to abstract arrays content by bounded, consecutive, non-overlapping, and possibly empty segments. It allows both the abstraction of array element values and relational abstractions between arrays content and their indexes.

Each abstract domain represents a trade-off between the precision and the efficiency of the analysis, since usually more precise analyses require less efficient abstract domains, slowing down the convergence of the computation of the fixpoint algorithm. However, different abstract domains can be composed in various ways (such as the Cartesian and the reduced product) within the Abstract Interpretation framework [38, 42, 56].

### State of the Art

While the compositional operators introduced in Section 2.6.4 have been extensively used to improve the precision of numerical analyses and implemented in various tools [59, 64], when we talk about string analysis, their usages and instantiations are short.

In particular, Yu et al. [185] conducted an analysis of strings on Web applications, to verify properties on strings that are relevant for security by combining (using the Cartesian product) the relation and alphabet abstractions. Elements in the product lattice are the so-called *abstraction classes* which are chosen heuristically. Amadini et al. [10] designed a framework that allows a flexible combination of several string abstract domains for JavaScript analysis. This combination was implemented as an extension of the SAFE tool [122]. Since this approach relied on the Cartesian product, the resulting analysis exposed more precise results than the single abstract domains, but such information could have been reconstructed from the individual execution of the analyses. A further improvement of such approach [11] introduced a general and modular framework to combine several string abstract domains through the reduced product by introducing the concept of *reference domain*. Let  $\mathbf{D}$  be a concrete domain, and  $\overline{\mathbf{D}}_1, \dots, \overline{\mathbf{D}}_n$  be abstract domains. A reference domain soundly approximates  $\mathbf{D}$  and captures any information expressible in each of the abstract domains  $\overline{\mathbf{D}}_1, \dots, \overline{\mathbf{D}}_n$ . The reference domain can be applied as a medium for systematically transferring information from all these abstract domains. The information is exchanged through a *strengthening function*, i.e., a closure operator on  $\overline{\mathbf{D}}_1, \dots, \overline{\mathbf{D}}_n$ . Depending on the considered domains, the strengthening function may be computationally expensive. Therefore, the authors introduced a *weak strengthening function* that is less precise but more efficient.

### Contribution

In this chapter, we combine through the reduced product some well-known string abstract domains (such as the String Length [11], the Character Inclusion and the Prefix [49] domains) with the Segmentation domain proposed by Cousot et al. in [57] to get more

sophisticated and precise abstractions. On the one hand, string domains track information about the content of string values. On the other hand, the Segmentation domain behaves like a functor that allows lifting basic domains through structural information. Since it represents the string values by distinct segments, it allows string domains to keep track of information only on a part of the string.

The main contribution of the chapter can be summarized as follows:

- We highlight the limits of the existing basic string abstract domains that deal with string shape and content properties separately or partially integrated, and the need to overcome these limits by providing a systematic construction of domains where string shape information is balanced with respect to string content information.

- We specialize the Segmentation domain for string analysis.

For instance, the string segmentation 'a' [2, 4] 'b' [0, 2] has two consecutive segments (i.e., 'a' [2, 4] and 'b' [0, 2]), each of which is composed by the character representation followed by a numerical interval. The latter indicates how many times the character approximation that precedes it is repeated. Thus, the abstract string 'a' [2, 4] 'b' [0, 2] represents all the possible strings having the character ‘a’ repeated from 2 to 4 times followed by the character ‘b’ repeated from 0 to 2 times.

- We combine basic string abstract domains with the Segmentation domain, where the tracked string information is homogeneous by exploiting the notions of Granger methodology and refinement operators [87], and we show how to actually compute these reduced products.

Consider the segmentation  $a^* [2, 4] b^* [0, 2]$  where character values are abstracted by the Prefix domain. This segmentation represents all the strings of length ranging from 2 to 4 and beginning with the character ‘a’ concatenated with strings of length ranging from 0 to 2 and beginning with the character ‘b’. Note that this combination of domains is in a position to fully cover both string shape and content information.

- We show how inconsistency conditions can be associated with abstract values, enabling the analysis to detect potential vulnerabilities within the source code.

To better understand the motivation and the results of this work, consider the example below.

**Example 4.1.** Let  $\overline{\mathbf{SL}}$ ,  $\overline{\mathbf{CI}}$ ,  $\overline{\mathbf{PR}}$  and  $\overline{\mathbf{S}}$  be the String Length, Character Inclusion, Prefix and Segmentation string abstract domains respectively. Table 4.1 depicts possible elements of these abstract domains. Note that all the abstractions in Table 4.1 approximate a set of string values that contains, for example, the literal “aadd00xyyy88”.

Consider in particular the following abstract values:

- |   |                            |
|---|----------------------------|
| (i) aadd*   | $(\overline{\mathbf{PR}})$ |
| (ii) $[\mathbf{a}, \mathbf{d}][5, 5] [\mathbf{0}, \mathbf{0}][2, 2] [\mathbf{x}, \mathbf{y}][5, 5] [\mathbf{7}, \mathbf{10}][2, 2]$ | $(\overline{\mathbf{S}})$  |

String domain	Abstract elements
$\overline{\text{CI}}$	$(\{a,d,0,x,y,8\}, \{a,d,0,x,y,8,4\}) (\{a,d\}, \{a\dots z,0\dots 9\}) (\emptyset, a,d,x,y,z,0\dots 9) \dots$
$\overline{\text{PR}}$	$a^*, aad^*, aadd^* \dots$
$\overline{\text{SL}}$	$[14, 14], [5, 16], [0, 20] \dots$
$\overline{\text{S}}$	'a' [2, 2] 'd' [3, 3] '0' [2, 2] 'x' [2, 2] 'y' [3, 3] '8' [2, 2] literal [5, 5], cypher [2, 2] literal [5, 5] cypher [2, 2] <b>[a,d]</b> [5, 5] <b>[0,0]</b> [2, 2] <b>[x,y]</b> [5, 5] <b>[7,10]</b> [2, 2]

TABLE 4.1: Introductory example.

Prefix (i) approximates all the strings starting with the sequence of characters “aadd”, while segmentation (ii) approximates all the strings starting with a character from ‘a’ to ‘d’ repeated 5 times concatenated to the cypher ‘0’ repeated two times and so on. Note that we write on boldface the segment character representations to distinguish them from the segment bounds. The reduced product between these two values leads to an improvement of the information precision tracked by each of the abstract elements. In particular, we can safely add to the end of the prefix  $aadd^*$  the sequence of cyphers 00 as the segmentation approximates strings where the integer 0 occurs at position 5 and 6, obtaining  $aadd00^*$ . Instead, the prefix clearly abstracts strings starting with a sequence of two ‘a’ followed by a sequence of three ‘d’. Thus the first segment of the segmentation value can be split accordingly, yielding to:

$$\langle aadd00^*, [\mathbf{a,a}][2, 2] [\mathbf{d,d}][3, 3] [\mathbf{0,0}][2, 2] [\mathbf{x,y}][5, 5][\mathbf{7,10}][2, 2] \rangle (\overline{\text{PR}} \otimes \overline{\text{S}})$$

Observe that in this way, both abstract components have been lifted to a more accurate representation in the resulting reduced product.

◻

## 4.2 Syntax

We briefly recall the string operators defined in [49], together with their intuitive semantics. While mainstream programming languages support a wider set of operators, we will focus our discussion and approach on this minimal set since they support the most important computations over string values.

In particular, let  $\text{str}$  be a sequence of characters.  $\text{new String}(\text{str})$  is the operator that generates a new constant string. Also, let  $\text{s1}$  and  $\text{s2}$  be two strings;  $\text{concat}(\text{s1}, \text{s2})$  concatenates  $\text{s1}$  and  $\text{s2}$ . Then, let  $\text{s}$  be a string and, let  $\text{b}$  and  $\text{e}$  be two integer values.  $\text{substring}_b^e(\text{s})$  returns the substring of  $\text{s}$  from the index  $\text{b}$  to the index  $\text{e}$ . Finally, let  $\text{s}$  be a string and let  $\text{c}$  be a character.  $\text{contains}_c(\text{s})$  returns true if and only if the character  $\text{c}$  appears in  $\text{s}$ .

## 4.3 Concrete Domain and Semantics

In the following, we recall the concrete domain and semantics of [49]. We appropriately modify the original notation to be consistent with the rest of this chapter.

### 4.3.1 Concrete Domain

Formally, let  $\Sigma^*$  be the set of all possible finite sequences of characters introduced in Section 2.1, a concrete domain is the complete lattice  $(\mathcal{P}(\Sigma^*), \subseteq, \emptyset, \Sigma^*, \cap, \cup)$  where:  $\mathcal{P}(\Sigma^*)$  denotes the powerset of  $\Sigma^*$  (that is, the set of all string values),  $\subseteq$  denotes the set inclusion (i.e., the partial order between elements in  $\mathcal{P}(\Sigma^*)$ ), the emptyset  $\emptyset$  is the bottom element of the lattice  $\mathcal{P}(\Sigma^*)$ ,  $\Sigma^*$  is the top element of the lattice  $\mathcal{P}(\Sigma^*)$ , the set intersection  $\cap$  is the greatest lower bound operand in  $\mathcal{P}(\Sigma^*)$ , and the set union  $\cup$  is the least upper bound operand in  $\mathcal{P}(\Sigma^*)$ .

### 4.3.2 Concrete Semantics

We now define the concrete semantics  $\mathfrak{S}$  and  $\mathfrak{B}$ . Formally:

$$\mathfrak{S} : \mathbf{Stm} \times \Sigma^* \cup [\mathcal{P}(\Sigma^*)]^k \rightarrow \mathcal{P}(\Sigma^*)$$

$$\mathfrak{B} : \mathbf{Stm} \times \mathcal{P}(\Sigma^*) \rightarrow \{\mathbf{true}, \mathbf{false}, \top_{\mathbf{B}}\}$$

where  $\mathbf{Stm}$  denotes a generic string operators. In particular, the semantics  $\mathfrak{S}$  applies to **new String**, **concat** and **substring**. Note that we deal with unary and binary operations then  $k$  might be only 1 or 2. On the other hand, the semantics  $\mathfrak{B}$  applies to **contains**. The concrete semantics is defined on the string operators introduced in Section 4.2 as follows:

- New String

$\mathfrak{S}[\mathbf{new\ String}(\sigma)]() = \{\sigma\}$ , that is, the semantics  $\mathfrak{S}$ , when applied to **new String**( $\sigma$ ), returns the singleton  $\{\sigma\}$ .

- String concatenation

$\mathfrak{S}[\mathbf{concat}](S_1, S_2) = \{\sigma_1 + \sigma_2 \mid \sigma_1 \in S_1 \wedge \sigma_2 \in S_2\}$ . In this case,  $S_1, S_2 \in \mathcal{P}(\Sigma^*)$  and let  $\sigma_1 + \sigma_2$  denote the concatenation between the strings  $\sigma_1$  and  $\sigma_2$ . The semantics  $\mathfrak{S}$  returns a set containing all the possible concatenations between the strings which belong to  $S_1$  and  $S_2$ .

- Substring

$\mathfrak{S}[\mathbf{substring}_b^e](S) = \{\sigma_b \dots \sigma_e \mid \sigma_1 \dots \sigma_n \in S \wedge n \geq e \wedge b \leq e\}$  where  $S \in \mathcal{P}(\Sigma^*)$ . In this case, the semantics  $\mathfrak{S}$  returns a set containing all the substrings from the  $b$ -th to the  $e$ -th character ( $\sigma_b$  and  $\sigma_e$  respectively) of the strings which belong to  $S$ . If a string is too short, the resulting set will not contain any element related to it.

- Is contained

$$\mathfrak{B}[\text{contains}_c](S) = \begin{cases} \text{true} & \text{if } \forall \sigma \in S : c \in \text{char}(\sigma) \\ \text{false} & \text{if } \forall \sigma \in S : c \notin \text{char}(\sigma) \\ \top_{\mathbf{B}} & \text{otherwise} \end{cases}$$

Let  $S \in \mathcal{P}(\Sigma^*)$ . In this case,  $\mathfrak{B}$  returns **true** if all the strings in  $S$  contain the character  $c$ , **false** if  $c$  is not contained by any of the strings in  $S$ , otherwise it returns  $\top_{\mathbf{B}}$ .

### Example

```

1 ResultSet getPerishablePrices(String lowerBound) {
2   String query = "SELECT '$' || (RETAIL/100) FROM INVENTORY WHERE ";
3   if (lowerBound != null)
4     query += "WHOLESALE > " + lowerBound + " AND ";
5
6   query += "TYPE IN (" + getPerishableTypeCode() + ")";
7   return statement.executeQuery(query);
8 }
9 String getPerishableTypeCode() {
10  return "SELECT TYPECODE, TYPEDESC FROM TYPES WHERE NAME = 'fish' OR NAME
        = 'meat'";
11 }

```

LISTING 4.1: Java code building up and executing an SQL query

Listing 4.1 reports the source code of the example that will be used to explain how the analysis with the basic string abstract domains works. This example is taken from Gould et al. [86]. In particular, method `getPerishablePrices` aims at executing a SQL query that selects, given a string representing a `lowerBound`, all the items from an inventory that are either `fish` or `meat` and whose `WHOLESALE` is greater than the `lowerBound` if such parameter is not null. This SQL query contains several errors:

1. `'$' || (RETAIL/100)` concatenates the character `$` with the numeric expression `RETAIL/100`,
2. `lowerBound` is an arbitrary string that might contain non-numeric characters, and therefore the comparison between `WHOLESALE` and its value might cause a runtime error, and
3. the subquery returned by `getPerishableTypeCode` returns two columns instead of one, and this could cause an error.

For the sake of readability we assign a shortcut to the string constants appearing in Listing 4.1 (cf. Table 4.2).



Name	String constant
$\sigma_1$	"SELECT '\$'    (RETAIL/100) FROM INVENTORY WHERE "
$\sigma_2$	"WHOLESALE > "
$\sigma_3$	" AND "
$\sigma_4$	"TYPE IN ("
$\sigma_5$	"SELECT TYPECODE, TYPEDESC FROM TYPES WHERE NAME = 'fish' OR NAME = 'meat'"
$\sigma_6$	");"

TABLE 4.2: Shortcuts of string constants in Listing 4.1.

## 4.4 String Abstract Domains

In this section, we provide an overview of a suite of basic string abstract domains [11, 49] (cf. Appendix B for more details), which collect the set of possible string values. In particular, we highlight and compare how the following string abstract domains track content and shape information of strings.

### 4.4.1 String Length

The first domain we consider is the String Length abstract domain  $\overline{\mathbf{SL}}$ , as presented in [11, 127]. This domain tracks, through a numerical interval  $[m, M]$  (such that  $m \in \mathbb{N}$  and  $M \in \mathbb{N} \cup \{\infty\}$ ), the minimum (i.e.,  $m$ ) and the maximum (i.e.,  $M$ ) length of the concrete strings it represents.

The String Length abstract domain detects information about the shape of the concrete strings it represents, i.e., their lengths. In particular, it precisely approximates only the empty string. In all the other cases, the abstraction totally loses the information about the content of the strings (e.g., the knowledge about characters order, repetitions, etc.). Any abstract element different from the top element ( $\top_{\overline{\mathbf{SL}}} = [0, \infty]$ ) leads to a finite set of concrete strings whose cardinality depends on the cardinality of the alphabet  $\Sigma$  and on the width of the abstract interval.  $\overline{\mathbf{SL}}$  can be implemented in a simple and efficient way as operations on it run in constant time.  $\overline{\mathbf{SL}}$  is an infinite lattice, and it does not respect the ACC, which is why it has been equipped with a widening operator.

Appendix B.1 reports the complete formalization and proof of soundness of this domain and its semantics.

**Example 4.2.** The result of the analysis of the program in Listing 4.1, using  $\overline{\mathbf{SL}}$ , is in Table 4.3. At pp.2, the variable `query` is associated to a state containing the abstraction of  $\sigma_1$ . The latter is approximated by the length interval  $[|\sigma_1|, |\sigma_1|]$ . The input variable `lowerBound` appears at pp.3 and, as it is unknown, it is abstracted by the top element of the  $\overline{\mathbf{SL}}$  lattice, i.e.,  $\top_{\overline{\mathbf{SL}}} = [0, \infty]$ . At pp.4, the variable `query` is associated to a

Program point	Variable	$\overline{\mathbf{SL}}$
pp.2	query	$[ \sigma_1 ,  \sigma_1 ]$
pp.3	lowerBound	$[0, \infty]$
pp.4	query	$[ \sigma_1  +  \sigma_2  +  \sigma_3 , \infty]$
pp.5	query	$[ \sigma_3 , \infty]$
pp.6	query	$[ \sigma_3  +  \sigma_4  +  \sigma_5  +  \sigma_6 , \infty]$

TABLE 4.3: Program analysis with  $\overline{\mathbf{SL}}$ .

state containing the concatenation of the abstractions of  $\sigma_1$ ,  $\sigma_2$ , `lowerBound` and  $\sigma_3$ , i.e.,  $[|\sigma_1| + |\sigma_2| + 0 + |\sigma_3|, |\sigma_1| + |\sigma_2| + \infty + |\sigma_3|] = [|\sigma_1| + |\sigma_2| + |\sigma_3|, \infty]$ . Then, at pp.5, the least upper bound ( $\perp_{\overline{\mathbf{SL}}}$ ) between the abstract value of `query` after pp.2 and after pp.4 is computed, i.e.,  $[|\sigma_1|, |\sigma_1|] \perp_{\overline{\mathbf{SL}}} [|\sigma_1| + |\sigma_2| + |\sigma_3|, \infty] = [\min(|\sigma_1|, |\sigma_2|, |\sigma_3|), \max(|\sigma_1|, \infty)] = [|\sigma_3|, \infty]$ . Finally, at pp.6, `query` is associated to a state containing the concatenation of the abstractions of itself after pp.5 and the strings  $\sigma_4$ ,  $\sigma_5$  and  $\sigma_6$ . Thus, at the end, `query` will have a length between  $|\sigma_3| + |\sigma_4| + |\sigma_5| + |\sigma_6|$  and  $\infty$ .  $\square$

#### 4.4.2 Character Inclusion

The Character Inclusion abstract domain  $\overline{\mathbf{CI}}$ , as defined in [49], uses a *pair of sets*  $(C, MC)$ , to track the characters that are certainly contained (i.e.,  $C$ ) and those that might be contained (i.e.,  $MC$ ) by the concrete strings it represents.

The Character Inclusion abstract domain detects information about the content of the concrete strings it represents. Like the String Length abstract domain, also  $\overline{\mathbf{CI}}$  can precisely approximate the empty string only, and it makes it possible to infer the minimum length of the strings it abstracts, which corresponds to the cardinality of the set  $C$ . In all the other cases, the abstraction loses the information about concrete strings shape. Indeed, similarly to  $\overline{\mathbf{SL}}$ ,  $\overline{\mathbf{CI}}$  does not preserve the information about the order of appearance of the characters, characters repetitions, and other relevant string properties. Any abstract element different from the bottom element ( $\perp_{\overline{\mathbf{CI}}}$ ) leads to an infinite set of concrete strings.  $\overline{\mathbf{CI}}$  is not computationally expensive, and it has finite height; consequently, the termination of the analysis is guaranteed by its least upper bound (used as widening operator).

Appendix B.2 reports the complete formalization and proof of soundness of this domain and its semantics.

**Example 4.3.** The result of the analysis of the program in Listing 4.1, using  $\overline{\mathbf{CI}}$ , is in Table 4.4. At pp.2, the variable `query` is associated with a state containing the abstraction of  $\sigma_1$ . The latter is approximated by a pair of set of characters, i.e.,  $(C_1, MC_1)$  with  $C_1 = MC_1$ . The input variable `lowerBound` is abstracted by the top element of the  $\overline{\mathbf{CI}}$  lattice, i.e.,  $\top_{\overline{\mathbf{CI}}} = (\emptyset, \Sigma)$ . At pp.4, the variable `query` is associated with a state containing the concatenation of the abstractions of  $\sigma_1$ ,  $\sigma_2$ , `lowerBound` and  $\sigma_3$ , i.e.,  $(C_1 \cup C_2 \cup \emptyset \cup$

Program point	Variable	$\overline{\mathbf{CI}}$
pp.2	query	$(C_1, MC_1)$
pp.3	lowerBound	$(\emptyset, \Sigma)$
pp.4	query	$(C_1 \cup C_2 \cup C_3, \Sigma)$
pp.5	query	$(C_1, \Sigma)$
pp.6	query	$(C_1 \cup C_4 \cup C_5 \cup C_6, \Sigma)$

TABLE 4.4: Program analysis with  $\overline{\mathbf{CI}}$ .

$C_3, MC_1 \cup MC_2 \cup \Sigma \cup MC_3) = (C_1 \cup C_2 \cup C_3, \Sigma)$ . Then, at pp.5, the least upper bound ( $\sqcup_{\overline{\mathbf{CI}}}$ ) between the abstract value of `query` after pp.2 and after pp.4 is computed, i.e.,  $(C_1, MC_1) \sqcup_{\overline{\mathbf{CI}}} (C_1 \cup C_2 \cup C_3, \Sigma) = (C_1 \cap (C_1 \cup C_2 \cup C_3), MC_1 \cup \Sigma) = (C_1, \Sigma)$ . Finally, at pp.6, `query` is associated with a state containing the concatenation of the abstractions of itself after pp.5 and the strings  $\sigma_4$ ,  $\sigma_5$  and  $\sigma_6$ . Thus, at the end, `query` will surely contain the characters in  $\sigma_1$ ,  $\sigma_4$ ,  $\sigma_5$  and  $\sigma_6$  and it will probably contain any character.  $\square$

#### 4.4.3 Prefix and Suffix

The Prefix abstract domain  $\overline{\mathbf{PR}}$ , presented in [49], approximates a set of concrete strings through a sequence of characters whose last element is  $*$ , that denotes any possible suffix string (the empty string  $\varepsilon$  is included). Instead, the Suffix abstract domain  $\overline{\mathbf{SU}}$  for string values mirrors the Prefix domain, and its notation and all its operators are dual to those of  $\overline{\mathbf{PR}}$ . The suffix domain abstracts strings through their suffix preceded by  $*$ , which denotes any possible prefix string,  $\varepsilon$  included. Moreover, Amadini et al. [11] discussed the Prefix-Suffix abstract domain  $\overline{\mathbf{PS}}$ , which approximates string values by their prefix and suffix simultaneously. Again, the notation and the operators of this domain can be easily induced by the Prefix and the Suffix domains. Indeed, the Prefix-Suffix domain abstracts strings through a pair of strings  $(p, s)$  which concretizes to the set of all possible strings having  $p$  as prefix and  $s$  as suffix (note that here  $*$  is not included in the definition).

The domains discussed above partially detect both content and shape of the concrete string they represent. Indeed,  $\overline{\mathbf{PR}}$ ,  $\overline{\mathbf{SU}}$ , and  $\overline{\mathbf{PS}}$  can track part of the strings structure, such as the initial part, the ending one or both of them, the minimum strings length, the characters surely contained, etc. Again, any abstract element different from the bottom element ( $\perp_{\overline{\mathbf{PR}}}$ ,  $\perp_{\overline{\mathbf{SU}}}$  and  $\perp_{\overline{\mathbf{PS}}}$  respectively) represents an infinite set of strings. Even though operations on these domains can be computed in linear time, they suffer from having an infinite height, e.g., given any prefix we can always add a character at the end of it, obtaining a new prefix. However, the domains respect the ACC, and the termination of the analysis is ensured.

Appendix B.3 contains the complete formalization and proof of soundness of the Prefix domain. Costantini et al. [49] provides this on the Suffix domain, while the one for  $\overline{\mathbf{PS}}$  can

Program point	Variable	$\overline{\mathbf{PR}}$
pp.2	query	$\sigma_1*$
pp.3	lowerBound	*
pp.4	query	$\sigma_1*$
pp.5	query	$\sigma_1*$
pp.6	query	$\sigma_1*$

TABLE 4.5: Program analysis with  $\overline{\mathbf{PR}}$ .

be obtained by the pointwise application of the operators and semantics on the Cartesian product of  $\overline{\mathbf{PR}}$  and  $\overline{\mathbf{SU}}$ .

**Example 4.4.** The result of the analysis of the program in Listing 4.1, using  $\overline{\mathbf{PR}}$ , is in Table 4.5. At pp.2, the variable `query` is associated with a state containing the abstraction of  $\sigma_1$ . The latter is approximated by the prefix  $\sigma_1*$ . The input variable `lowerBound` is abstracted by the top element of the  $\overline{\mathbf{PR}}$  lattice, i.e.,  $\top_{\overline{\mathbf{PR}}} = *$ . At pp.4, the variable `query` is associated with a state containing the concatenation of the abstractions of  $\sigma_1$ ,  $\sigma_2$ , `lowerBound` and  $\sigma_3$ , i.e.,  $\sigma_1* +_{\overline{\mathbf{PR}}} \sigma_2* +_{\overline{\mathbf{PR}}} * +_{\overline{\mathbf{PR}}} \sigma_3* = \sigma_1*$ . Then, at pp.5, the least upper bound ( $\sqcup_{\overline{\mathbf{PR}}}$ ) between the abstract value of `query` after pp.2 and after pp.4 is computed, i.e.,  $\sigma_1* \sqcup_{\overline{\mathbf{PR}}} \sigma_1* = \sigma_1*$ . Finally, at pp.6, `query` is associated with a state containing the concatenation of the abstractions of itself after pp.5 and the strings  $\sigma_4$ ,  $\sigma_5$  and  $\sigma_6$ . Thus, at the end, `query` will for sure begin with  $\sigma_1$  followed by any possible suffix string  $*$ .

◻

## 4.5 Segmentation Abstract Domain

In Section 3.2 we recalled `FunArray`, i.e., the array segmentation abstract domain functor by Cousot et al. [57]. Since the order of characters in strings is fundamental to track precise information on these values, we instantiate the `FunArray` abstract domain for string analysis. In the following, we slightly modify the notation introduced in Section 3.2 to highlight the fact that we are instantiating `FunArray` over strings.

### 4.5.1 String Concrete Representation

Let  $\varsigma \in R_s = \mathbf{S} \rightarrow \mathbf{S}$  be concrete string environments mapping string variables  $s \in \mathbf{S}$  to their instrumented values  $\varsigma(s) \in \mathbf{S} \triangleq R_v \times \mathbf{E} \times \mathbf{E} \times (\mathbb{Z} \rightarrow (\mathbb{Z} \times \Sigma))$ . Thus, a string variable  $s$  is represented by a quadruple  $\varsigma(s) = (\rho, l_s, h_s, A_s) \in \mathbf{S}$ , such that:

- $\rho \in R_v \triangleq \mathbf{X} \rightarrow \mathbf{V}$  are concrete scalar variable environments mapping variables  $x \in \mathbf{X}$  to their values  $\llbracket x \rrbracket \rho \in \mathbf{V}$ .

- $\mathbf{E}$  is the expressions domain, built from constants and scalar variables, through mathematical unary and binary operators, and  $l_s, h_s \in \mathbf{E}$ . The values of  $l_s$  and  $h_s$  ( $\llbracket l_s \rrbracket \rho$  and  $\llbracket h_s \rrbracket \rho$ ) denote respectively the lower and the upper limit of the string variable  $s$ .
- $A_s$  is a function mapping an index  $i \in \llbracket l_s \rrbracket \rho, \llbracket h_s \rrbracket \rho$  to the pair  $\langle i, c \rangle$  of the index  $i$  and the corresponding string character  $c \in \Sigma$ .

**Example 4.5.** Let  $s$  be a string variable initialized to the value "bunny". The concrete value of  $s$  is given by the tuple  $\varsigma(s) = (\rho, 0, 5, A_s)$ , where the value of the lower and the upper bound  $s$  is inferred from the context and the function  $A_s$  maps an index  $i \in [0, 5)$  to the pair (index, indexed character value). Thus, the codomain of  $A_s$  is the set  $\{(0, 'b'), (1, 'u'), (2, 'n'), (3, 'n'), (4, 'y')\}$ .

◻

### 4.5.2 Abstract Domain

The Segmentation abstract domain functor  $\bar{\mathbf{S}}$  for strings is a function from the parameter abstract domains  $\bar{\mathbf{B}}$ ,  $\bar{\mathbf{C}}$  and  $\bar{\mathbf{R}}$ , where  $\bar{\mathbf{B}}$  and  $\bar{\mathbf{R}}$  are, in turn, abstract domain functors. The segment bound abstract domain functor  $\bar{\mathbf{B}}$  is a function of the expression abstract domain  $\bar{\mathbf{E}}(\bar{\mathbf{X}})$  which depends on the variable abstract domain  $\bar{\mathbf{X}}$ , leading to the instantiated segment bound abstract domain  $\bar{\mathbf{B}}(\bar{\mathbf{E}}(\bar{\mathbf{X}}))$ .  $\bar{\mathbf{C}}$  is the string element abstract domain. Finally,  $\bar{\mathbf{R}}$  is the variable environment abstract domain functor which depends on  $\bar{\mathbf{X}}$  too, leading to the variable environment abstract domain  $\bar{\mathbf{R}}(\bar{\mathbf{X}})$ . Precisely:

- The variable abstract domain  $\bar{\mathbf{X}}$  encodes program variables.
- The variable environment abstract domain functor  $\bar{\mathbf{R}}$  depends on  $\bar{\mathbf{X}}$ , leading to the variable environment abstract domain  $\bar{\mathbf{R}}(\bar{\mathbf{X}})$ . Elements in  $\bar{\mathbf{R}}$  (shorthand for  $\bar{\mathbf{R}}(\bar{\mathbf{X}})$ ) are abstract variable environments  $\bar{\rho} \in \bar{\mathbf{R}} = \mathbf{X} \rightarrow \bar{\mathbf{V}}$ , where the value abstract domain  $\bar{\mathbf{V}}$  approximates properties of values in  $\mathbf{V}$ .  $\bar{\mathbf{R}}$  approximates sets of concrete variable environments. Formally, the concretization is  $\gamma_{\bar{\mathbf{R}}} : \bar{\mathbf{R}} \rightarrow \mathcal{P}(R_v)$ , where  $R_v = \mathbf{X} \rightarrow \mathbf{V}$  (cf. Section 4.5.1).
- The expression abstract domain functor  $\bar{\mathbf{E}}$  depends on  $\bar{\mathbf{X}}$ , leading to the expression abstract domain  $\bar{\mathbf{E}}(\bar{\mathbf{X}})$ . Elements in  $\bar{\mathbf{E}}$  (shorthand for  $\bar{\mathbf{E}}(\bar{\mathbf{X}})$ ) are symbolic expressions  $\bar{e} \in \bar{\mathbf{E}}(\bar{\mathbf{X}})$ , restricted to a canonical normal form, which depend on variables in  $\bar{\mathbf{X}}$  (notice that  $\perp_{\bar{\mathbf{E}}}, \top_{\bar{\mathbf{E}}} \in \bar{\mathbf{E}}$ ).  $\bar{\mathbf{E}}$  approximates program expressions. Formally, the concretization is  $\gamma_{\bar{\mathbf{E}}} : \bar{\mathbf{E}} \rightarrow \bar{\mathbf{R}} \rightarrow \mathcal{P}(\mathbf{V})$ . Moreover,  $\bar{\mathbf{E}}$  is equipped with sum ( $\oplus_{\bar{\mathbf{E}}}$ ), subtraction ( $\ominus_{\bar{\mathbf{E}}}$ ) and comparison ( $\leq_{\bar{\mathbf{E}}}$ ) operations. The comparison result depends on the normal form of the expressions. In general, two expressions are said to be *comparable* if and only if their comparison returns **true**.

The choice of the expression canonical form is let free.

---

**Algorithm 3** align procedure.

---

**Input:**  $\bar{s}_1, \bar{s}_2$

**Output:**  $\bar{s}_1, \bar{s}_2$  (possibly modified) or error

```

1:  $i \leftarrow 1$ 
2:  $k_1 \leftarrow 0; k_2 \leftarrow 0$ 
3: if  $\bar{s}_1 \neq \perp_{\bar{s}} \wedge \bar{s}_2 \neq \perp_{\bar{s}}$  then
4:   while  $i \leq \text{numSeg}(\bar{s}_1) + k_1 \wedge i \leq \text{numSeg}(\bar{s}_2) + k_2$  do
5:     if  $\bar{e}'.\bar{b}.\bar{s}_1[i] =_{\bar{E}} \bar{e}'.\bar{b}.\bar{s}_2[i]$  then
6:        $i++$ 
7:     else if  $\bar{e}'.\bar{b}.\bar{s}_2[i] \leq_{\bar{E}} \bar{e}'.\bar{b}.\bar{s}_1[i]$  then
8:        $\bar{e}'.\bar{b}.\bar{s}_1[i] \leftarrow \bar{e}'.\bar{b}.\bar{s}_2[i]$ 
9:        $\bar{s}_1[i].\text{append}(\bar{c}.\bar{s}_1[i][0], \bar{e}'.\bar{b}.\bar{s}_1[i] \ominus_{\bar{E}} \bar{e}'.\bar{b}.\bar{s}_2[i])$ 
10:       $k_1++, i++$ 
11:    else if  $\bar{e}'.\bar{b}.\bar{s}_1[i] \leq_{\bar{E}} \bar{e}'.\bar{b}.\bar{s}_2[i]$  then
12:       $\bar{e}'.\bar{b}.\bar{s}_2[i] \leftarrow \bar{e}'.\bar{b}.\bar{s}_1[i]$ 
13:       $\bar{s}_2[j].\text{append}(\bar{c}.\bar{s}_2[i][0], \bar{e}'.\bar{b}.\bar{s}_2[i] \ominus_{\bar{E}} \bar{e}'.\bar{b}.\bar{s}_1[i])$ 
14:       $k_2++, i++$ 
15:    else
16:      return error ( $\bar{s}_1$  and  $\bar{s}_2$  can not be aligned)
17: else
18:   return error ( $\bar{s}_1$  and  $\bar{s}_2$  can not be aligned)
19: return  $\bar{s}_1, \bar{s}_2$ 

```

---

Note that in the following examples we will use the expression normal form  $\bar{x} + k$  where  $k \in \mathbb{Z}$  and  $\bar{x} \in \bar{\mathbf{X}} \cup \{v_0\}$ , with  $v_0$  being a special variable whose value is always 0 and we omit it.

- The string element abstract domain  $\bar{\mathbf{C}}$  approximates sets of pairs (index, indexed string element), where  $\bar{c} \in \bar{\mathbf{C}}$ . Formally, the abstraction is  $\alpha_{\bar{\mathbf{C}}} : \mathcal{P}(\mathbb{Z} \times \Sigma) \rightarrow \mathcal{P}(\Sigma) \rightarrow \bar{\mathbf{C}}$ , i.e., elements in  $\mathcal{P}(\mathbb{Z} \times \Sigma)$  may be first abstracted to  $\mathcal{P}(\Sigma)$  so to perform a non-relational analysis. The concretization is  $\gamma_{\bar{\mathbf{C}}} : \bar{\mathbf{C}} \rightarrow \mathcal{P}(\mathbb{Z} \times \Sigma)$ .

The choice of  $\bar{\mathbf{C}}$  is let free.

- The segment bound abstract domain functor  $\bar{\mathbf{B}}$  is a function of the expression abstract domain  $\bar{\mathbf{E}}$ , leading to the instantiated segment bound abstract domain  $\bar{\mathbf{B}}(\bar{\mathbf{E}})$ . Elements in  $\bar{\mathbf{B}}$  (shorthand for  $\bar{\mathbf{B}}(\bar{\mathbf{E}})$ ) are symbolic intervals  $\bar{b} \in \bar{\mathbf{B}}$ , such that  $\bar{b} = [\bar{e}, \bar{e}']$ , where  $\bar{e}, \bar{e}' \in \{\bar{\mathbf{E}} \setminus \{\perp_{\bar{\mathbf{E}}}, \top_{\bar{\mathbf{E}}}\}\}$  and  $\bar{e} \leq_{\bar{\mathbf{E}}} \bar{e}'$ . Formally, the concretization is  $\gamma_{\bar{\mathbf{B}}} : \bar{\mathbf{B}} \rightarrow \mathcal{P}(R_v)$ .

Elements of  $\bar{\mathbf{S}}$  belong to the set  $\bar{\mathbf{S}} = \{(\bar{\mathbf{C}} \times \bar{\mathbf{B}})^k \mid k \geq 1\} \cup \{\perp_{\bar{\mathbf{S}}}, \top_{\bar{\mathbf{S}}}\}$ , where  $\perp_{\bar{\mathbf{S}}}$  and  $\top_{\bar{\mathbf{S}}}$  are special elements denoting the bottom/top element of  $\bar{\mathbf{S}}$ . In particular, elements of  $\bar{\mathbf{S}}$  are in the form of  $\bar{s} = \bar{c}_1 \bar{b}_1 \dots \bar{c}_n \bar{b}_n$ , where a segment  $\bar{c}_i \bar{b}_i$  (with  $\bar{b}_i = [\bar{e}_i, \bar{e}'_i]$ ) abstracts a sequence of equal characters whose length goes from  $\bar{e}_i$  to  $\bar{e}'_i$ .

**Example 4.6.** Consider the string variable of the Example 4.5. Its value in the Segmentation abstract domain  $\bar{\mathbf{S}}$  is 'b' [1, 1] 'u' [1, 1] 'n' [2, 2] 'y' [1, 1], where  $\bar{\mathbf{C}}$  is the constant propagation domain for characters.

◻

Before defining the join and the meet operators of the Segmentation domain we present the following helping procedures: **align** (cf. Algorithm 3) and **fold** (cf. Algorithm 4).

Algorithm 3 aligns two segmentations if they are both different from the bottom element of  $\bar{\mathbf{S}}$  and comparable. The comparability is actually restricted to the segment bound upper limits of each pair of corresponding segments under consideration during the alignment procedure. Let  $\bar{s}_1, \bar{s}_2 \in \bar{\mathbf{S}}$  be different from  $\perp_{\bar{\mathbf{S}}}$  (line 3), the alignment procedure starts analysing the two segmentations from their leftmost segments, i.e.,  $\bar{s}_1[1]$  and  $\bar{s}_2[1]$ , and continues along their number of segments (line 4), i.e.,  $\text{numSeg}(\bar{s}_1)$  and  $\text{numSeg}(\bar{s}_2)$ , which may change during the procedure, if all the corresponding segment bound upper limits under consideration are comparable. If this is the case and the segment bound upper limits are equal, then the procedure moves to the next segments; otherwise, if one of the two segment bound upper limits is greater or equal than the other, the first is modified according to the latter segment bound upper limit (line 10 o line 14), and the “remaining part” is appended to the previously modified segment (line 11 or line 15). In the case in which one or both the input segmentations are equals to  $\perp_{\bar{\mathbf{S}}}$  or if during the alignment procedure two corresponding segment bound upper limits are not comparable, the algorithm stops, and it returns an error message.

**Example 4.7.** Consider the segmentations:

- $\bar{s}_1 = \text{'a' [0, 4] 'b' [1, 1] 'c' [0, x]}$
- $\bar{s}_2 = \text{'a' [2, 5 + x]}$

where the segment predicate abstract domain  $\bar{\mathbf{C}}$  is the constant propagation domain for characters and  $x$  is an integer variable whose value is greater than or equal to 0. Algorithm 3 on  $\bar{s}_1$  and  $\bar{s}_2$  is applied as follows: starting from  $i$  equal to 1, we enter into the loop. The first segment bound upper limit of  $\bar{s}_1$  is strictly smaller than the first segment bound upper limit of  $\bar{s}_2$  (for any value of  $x$ ), i.e.,  $\bar{e}'\bar{b}.\bar{s}_1[1] <_{\bar{\mathbf{E}}} \bar{e}'\bar{b}.\bar{s}_2[1]$  (line 11), then the value of  $\bar{e}'\bar{b}.\bar{s}_1[1]$  is assigned to  $\bar{e}'\bar{b}.\bar{s}_2[1]$  (line 12), and the segment  $\bar{c}.\bar{s}_2[1][0, \bar{e}'\bar{b}.\bar{s}_2[1] \ominus_{\bar{\mathbf{E}}} \bar{e}'\bar{b}.\bar{s}_1[1]]$  is appended to  $\bar{s}_2[1]$  (line 13). Thus, after the first iteration of the for loop,  $\bar{s}_1$  is unchanged and  $\bar{s}_2$  is updated. We obtain:

- $\bar{s}_1 = \text{'a' [0, 4] 'b' [1, 1] 'c' [0, x]}$
- $\bar{s}_2 = \text{'a' [2, 4] 'a' [0, 1 + x]}$

---

**Algorithm 4** fold procedure.

---

**Input:**  $\bar{s}$

**Output:**  $\bar{s}$  (possibly modified)

```

1:  $i \leftarrow 1$ 
2: while  $i \leq \text{numSeg}(\bar{s})$  do
3:   if  $\bar{c}.\bar{s}[i] =_{\bar{c}} \bar{c}.\bar{s}[i+1] \wedge \bar{e}.\bar{b}.\bar{s}[i+1] =_{\bar{e}} 0$  then
4:      $\bar{b}.\bar{s}[i] \leftarrow \bar{b}.\bar{s}[i] + \bar{b}.\bar{s}[i+1]$ 
5:      $\bar{s}.\text{remove}(\bar{s}[i+1])$ 
6:   else
7:      $i++$ 
8: return  $\bar{s}$ 

```

---

The counter  $i$  is increased by 1 and, since  $\bar{e}.\bar{b}.\bar{s}_1[2] \leq_{\bar{e}} \bar{e}.\bar{b}.\bar{s}_2[2]$  (line 11), to the second segment bound upper limit of  $\bar{s}_2$  is assigned the value of  $\bar{e}.\bar{b}.\bar{s}_1[2]$  (line 12) and the segment  $\bar{c}.\bar{s}_2[2][0, \bar{e}.\bar{b}.\bar{s}_2[2] \ominus_{\bar{e}} \bar{e}.\bar{b}.\bar{s}_1[2]]$  is appended to  $\bar{s}_2[2]$  (line 13). After the second iteration of the for loop,  $\bar{s}_1$  is still unchanged and  $\bar{s}_2$  is updated as follows:

- $\bar{s}_1 = \text{'a' } [0, 4] \text{ 'b' } [1, 1] \text{ 'c' } [0, x]$
- $\bar{s}_2 = \text{'a' } [2, 4] \text{ 'a' } [0, 1] \text{ 'a' } [0, x]$

Again the counter  $i$  is increased by 1. The third segment bound upper limits of  $\bar{s}_1$  and  $\bar{s}_2$  are equal (line 5), i.e.,  $\bar{e}.\bar{b}.\bar{s}_1[3] =_{\bar{e}} \bar{e}.\bar{b}.\bar{s}_2[3]$ ,  $i$  is increased by 1, but, now, they exceed the number of segments of both  $\bar{s}_1$  and  $\bar{s}_2$ . Thus the final result is the pair of segmentations obtained after the second iteration of the while loop.

◻

The alignment procedure (cf. Algorithm 3) is used to facilitate the computation of the join  $\sqcup_{\bar{g}}$  and the meet  $\sqcap_{\bar{g}}$ . The latter may lead to segmentations that need to be taken back to their normal form, i.e., where there are no consecutive segments with the same segment abstract predicate  $\bar{c}$ . Thus, Algorithm 4 brings a segmentation to its normal form and returns it, merging consecutive segments sharing the same segment abstract predicate (line 4). The join and meet operators and the partial order of  $\bar{\mathbf{S}}$  are defined as follows:

- $\sqcap_{\bar{g}}$  represents the meet operator between two string segmentations. The meet between two segmentations  $\bar{s}_1$  and  $\bar{s}_2$  is computed on their alignment if *i*)  $\text{align}(\bar{s}_1, \bar{s}_2)$  raises no error, and *ii*) one of the two aligned segmentations has more segments than the other, the exceeding segments have to be possibly empty, i.e., their lower limit has to be equal to 0; otherwise  $\perp_{\bar{g}}$  is returned. Formally, let  $\bar{s}'_1$  and  $\bar{s}'_2$  be the results of  $\text{align}(\bar{s}_1, \bar{s}_2)$  then:

$$\bar{s}_1 \sqcap_{\bar{g}} \bar{s}_2 = \begin{cases} \text{fold}(\bar{s}'_1 \sqcap_{\bar{g}} \bar{s}'_2) & \text{see (1)} \\ \perp_{\bar{g}} & \text{otherwise} \end{cases}$$



(1)  $\bar{s}_1 \sqcap_{\bar{S}} \bar{s}_2 = \text{fold}(\bar{s}'_1 \sqcap_{\bar{S}} \bar{s}'_2)$  if the following conditions hold:

\*  $\text{align}(\bar{s}_1, \bar{s}_2) \neq \text{error}$

\* given  $j, k \in [1, 2]$  with  $j \neq k$ , if  $\text{numSeg}(\bar{s}'_j) > \text{numSeg}(\bar{s}'_k)$  then,

$$\forall i \in [\text{numSeg}(\bar{s}'_k) + 1, \text{numSeg}(\bar{s}'_j)] \text{ of } \bar{s}'_j : \bar{e}_i =_{\bar{E}} 0$$

The meet between two aligned segmentations is performed so as the string character abstract domain meet ( $\sqcap_{\bar{C}}$ ) and the bound abstract domain meet ( $\sqcap_{\bar{B}}$ ) are applied segment-wise. Notice that if the number of segments of  $\bar{s}'_1$  ( $\text{numSeg}(\bar{s}'_1)$ ) is strictly greater than the number of segments of  $\bar{s}'_2$  ( $\text{numSeg}(\bar{s}'_2)$ ) then all the exceeding segments of  $\bar{s}'_1$  are not preserved by  $\sqcap_{\bar{S}}$  (the vice-versa is similar). Finally, we compute the fold of the meet between  $\bar{s}'_1$  and  $\bar{s}'_2$ .

- $\sqcup_{\bar{S}}$  represents the join operator between two string segmentations. The join between two segmentations  $\bar{s}_1$  and  $\bar{s}_2$  is computed on their alignment if  $\text{align}(\bar{s}_1, \bar{s}_2)$  does not raise an error; if only one of the two segmentations is equal to  $\perp_{\bar{S}}$  then their join returns the one which is different from the bottom element; if both  $\bar{s}_1$  and  $\bar{s}_2$  are the bottom element then their join returns  $\perp_{\bar{S}}$ ; otherwise  $\top_{\bar{S}}$  is returned. Formally, let  $\bar{s}'_1$  and  $\bar{s}'_2$  be the results of  $\text{align}(\bar{s}_1, \bar{s}_2)$  (cf. Algorithm 3) then:

$$\bar{s}_1 \sqcup_{\bar{S}} \bar{s}_2 = \begin{cases} \text{fold}(\bar{s}'_1 \sqcup_{\bar{S}} \bar{s}'_2) & \text{if } \text{align}(\bar{s}_1, \bar{s}_2) \neq \text{error} \\ \bar{s}_1 & \text{if } \bar{s}_2 = \perp_{\bar{S}} \wedge \bar{s}_1 \neq \perp_{\bar{S}} \\ \bar{s}_2 & \text{if } \bar{s}_1 = \perp_{\bar{S}} \wedge \bar{s}_2 \neq \perp_{\bar{S}} \\ \perp_{\bar{S}} & \text{if } \bar{s}_1 = \perp_{\bar{S}} \wedge \bar{s}_2 = \perp_{\bar{S}} \\ \top_{\bar{S}} & \text{otherwise} \end{cases}$$

The join between two aligned segmentations is performed so as the string character abstract domain join ( $\sqcup_{\bar{C}}$ ) and the bound abstract domain join ( $\sqcup_{\bar{B}}$ ) are applied segment-wise. If the number of segments of  $\bar{s}'_1$  is strictly greater than the number of segments of  $\bar{s}'_2$  then all the exceeding segments of  $\bar{s}'_1$  are preserved by  $\sqcup_{\bar{S}}$ , but their segment bound lower limit is set to 0 (the vice-versa is similar).

- Let  $\bar{s}_1$  and  $\bar{s}_2$  be two abstract values in the Segmentation domain. The partial order on  $\bar{S}$  is defined as follows:  $\forall \bar{s} \in \bar{S} : \perp_{\bar{S}} \sqsubseteq_{\bar{S}} \bar{s} \wedge \bar{s} \sqsubseteq_{\bar{S}} \top_{\bar{S}}$ . Otherwise, if both  $\bar{s}_1$  and  $\bar{s}_2$  are different from  $\perp_{\bar{S}}$  and  $\top_{\bar{S}}$  then,  $\bar{s}_1 \sqsubseteq_{\bar{S}} \bar{s}_2 \Leftrightarrow \bar{s}_1 \sqcup_{\bar{S}} \bar{s}_2 = \bar{s}_2$

### Concretization

The concretization function on the segmentation abstract domain  $\gamma_{\bar{S}} : \bar{S} \rightarrow \bar{\mathbf{R}} \rightarrow \mathcal{P}(\mathbf{S})$  maps an abstract element to a set of strings as follows:  $\gamma_{\bar{S}}(\perp_{\bar{S}}) = \emptyset$ ,  $\gamma_{\bar{S}}(\top_{\bar{S}}) = \mathbf{S}$ , while

it is the set of all possible sequences of characters derivable from a segmentation abstract predicate. The formalization below follows the one defined in Section 11.3 of [57]. Let  $\gamma_{\overline{\mathbf{R}}} : \overline{\mathbf{R}} \rightarrow \mathcal{P}(\mathbf{R}_v)$  be the concretization function for the variable abstract domain and let  $\gamma_{\overline{\mathbf{C}}} : \overline{\mathbf{C}} \rightarrow \mathcal{P}(\mathbb{Z} \times \Sigma)$  be the concretization function for the string elements abstract domain (cf. Section 4.5.2). Then  $\gamma_{\overline{\mathbf{S}}}^*$  denotes the concretization of a generic segment  $\overline{\mathbf{c}\mathbf{b}}$  where  $\overline{\mathbf{b}} = [\overline{\mathbf{e}}, \overline{\mathbf{e}'}]$ , formally:

$$\gamma_{\overline{\mathbf{S}}}^*(\overline{\mathbf{c}\mathbf{b}})\overline{\rho} = \{(\rho, l, h, A) \mid \rho \in \gamma_{\overline{\mathbf{R}}}(\overline{\rho}) \wedge \llbracket l \rrbracket \rho = 0 \wedge \exists e \in \llbracket \overline{\mathbf{e}} \rrbracket \rho, \llbracket \overline{\mathbf{e}'} \rrbracket \rho : \llbracket h \rrbracket \rho = e \wedge \forall i \in [0, e-1] : A(i) \in \gamma_{\overline{\mathbf{C}}}(\overline{\mathbf{c}})\}$$

where  $\overline{\rho} \in \overline{\mathbf{R}}$ . Then, the concretization function of a string segmentation is as follows:

$$\gamma_{\overline{\mathbf{S}}}(\overline{\mathbf{c}_1\mathbf{b}_1} \dots \overline{\mathbf{c}_n\mathbf{b}_n}) = \{(\rho, l, h, A) \in \prod_{i=1}^n \gamma_{\overline{\mathbf{S}}}^*(\overline{\mathbf{c}_i\mathbf{b}_i}) \mid \llbracket l \rrbracket \rho = 0 \wedge \llbracket h \rrbracket \rho = \sum_{i=1}^n e_i\}$$

and  $\gamma_{\overline{\mathbf{S}}}(\perp_{\overline{\mathbf{S}}}) = \emptyset$

Note that a segmentation abstract predicate is *valid* if the upper bounds of segments that contain the bottom element in an abstract domain  $\overline{\mathbf{C}}$  are possibly empty; otherwise, a string segmentation is *invalid*. The concretization function of an *invalid* segmentation maps the latter abstract value to the empty-set.

**Theorem 4.1.** Let  $\overline{\mathbf{X}} \subseteq \overline{\mathbf{S}}$  such that all elements in  $\overline{\mathbf{X}}$  can be aligned and their meet does not result in an invalid segmentation. Then, it holds that:

$$\gamma_{\overline{\mathbf{S}}}\left(\prod_{\overline{\mathbf{s}} \in \overline{\mathbf{X}}} \overline{\mathbf{s}}\right) = \bigcap_{\overline{\mathbf{s}} \in \overline{\mathbf{X}}} \gamma_{\overline{\mathbf{S}}}(\overline{\mathbf{s}})$$

*Proof.* The following inference chain holds:

$$\begin{aligned} \gamma_{\overline{\mathbf{S}}}\left(\prod_{\overline{\mathbf{s}} \in \overline{\mathbf{X}}} \overline{\mathbf{s}}\right) &= \gamma_{\overline{\mathbf{S}}}(\overline{\mathbf{s}}^*) \quad \text{by definition of } \prod_{\overline{\mathbf{S}}} \\ &\quad \text{where } \overline{\mathbf{s}}^* \text{ denotes the result of the meet between the segmentations } \overline{\mathbf{s}} \text{ in } \overline{\mathbf{X}} \\ &= \{(\rho, l, h, A) \in \prod_{i=1}^n \gamma_{\overline{\mathbf{S}}}^*(\overline{\mathbf{c}_i\mathbf{b}_i}) \mid \llbracket l \rrbracket \rho = 0 \wedge \llbracket h \rrbracket \rho = \sum_{i=1}^n e_i\} \quad \text{by definition of } \gamma_{\overline{\mathbf{S}}} \\ &= \bigcap_{\overline{\mathbf{s}} \in \overline{\mathbf{X}}} \{(\rho, l, h, m) \in \prod_{i=1}^n \gamma_{\overline{\mathbf{S}}}^*(\overline{\mathbf{c}_i\mathbf{b}_i}) \mid \llbracket l \rrbracket \rho = 0 \wedge \llbracket h \rrbracket \rho = \sum_{i=1}^n e_i\} \\ &= \bigcap_{\overline{\mathbf{s}} \in \overline{\mathbf{X}}} \gamma_{\overline{\mathbf{S}}}(\overline{\mathbf{s}}) \quad \text{by definition of } \gamma_{\overline{\mathbf{S}}} \end{aligned}$$

□

Observe that if the hypotheses of Theorem 4.1 are not satisfied, i.e., if either the abstract elements in  $\overline{\mathbf{X}}$  can not be aligned or their meet leads to an invalid segmentation, then  $\gamma_{\overline{\mathbf{S}}}\left(\prod_{\overline{\mathbf{s}} \in \overline{\mathbf{X}}} \overline{\mathbf{s}}\right) = \gamma_{\overline{\mathbf{S}}}(\perp_{\overline{\mathbf{S}}}) = \emptyset$ , and  $\bigcap_{\overline{\mathbf{s}} \in \overline{\mathbf{X}}} \gamma_{\overline{\mathbf{S}}}(\overline{\mathbf{s}}) = \emptyset$

### Abstraction

Let  $X \in \mathcal{P}(\mathbf{S})$  be a set of concrete string values. The abstraction function on the segmentation abstract domain  $\alpha_{\bar{\mathbf{S}}}$  maps  $X$  to  $\perp_{\bar{\mathbf{S}}}$  in the case in which  $X$  is equal to the empty set, otherwise to the segmentation that over-approximates values in  $X$ .

To summarize,  $\bar{\mathbf{S}}$  abstracts strings by a sequence of segments in an abstract domain  $\bar{\mathbf{C}}$  and a segment bound in an abstract domain  $\bar{\mathbf{B}}$ . The segments are computed according to how the string content is manipulated. The value of a string segmentation abstract predicate depends on three parameters: the abstract domain representing symbolic segment bound expressions, the domain abstracting the pairs (index, character), and the abstract domain assigning values to segment bound expressions [57]. Thus,  $\bar{\mathbf{S}}$  can be instantiated with different abstract domains, achieving different levels of precision and costs of the analysis. Note that the segmentation abstract domain has infinite height as a string may have infinitely symbolic segments, and a segment might take successive strictly increasing abstract values. Therefore, to guarantee the convergence of the analysis, we define a widening similarly to the one presented by Cousot et al. in [57]. Informally, if the number of segments exceeds a certain threshold during the alignment procedure between two segmentations, the widening transforms those segmentations into  $\top_{\bar{\mathbf{S}}}$ . Note that a preciser widening operator may be defined, possibly together with a narrowing operator, to improve the accuracy of the analysis.

#### 4.5.3 Abstract Semantics

In Section 4.2 we recalled the syntax of some operations of interest (i.e., `newString`, `concat`, `substring`, and `contains`) presented in [49]. Moreover, in Section 4.3.2 we reported their concrete semantics. Below, we formally define their approximation in the string segmentation abstract domain  $\bar{\mathbf{S}}$ . The semantics operators  $\mathfrak{S}_{\bar{\mathbf{S}}}$  and  $\mathfrak{B}_{\bar{\mathbf{S}}}$  are the abstract counterparts of  $\mathfrak{S}$  and  $\mathfrak{B}$  respectively (cf. Section 4.3.2).

- New string

Let  $\sigma$  be a string value. The semantics  $\mathfrak{S}_{\bar{\mathbf{S}}}$ , when applied to `new String( $\sigma$ )`, returns the segmentation of the string constant  $\sigma$ . Formally,

$$\mathfrak{S}_{\bar{\mathbf{S}}}[\text{new String}(\sigma)]() = \bar{\mathbf{s}}$$

$\bar{\mathbf{s}}$  abstracts, in an abstract domain  $\bar{\mathbf{C}}$ , the sequences of equal characters in  $\sigma$ . Notice that each segment bound of  $\bar{\mathbf{s}}$  will have equal lower and upper limits in an abstract domain  $\bar{\mathbf{E}}$ .

- String concatenation

Let  $\bar{s}_1$  (i.e.,  $\bar{c}_{1,1}\bar{b}_{1,1}\dots\bar{c}_{1,n}\bar{b}_{1,n}$ ) and  $\bar{s}_2$  (i.e.,  $\bar{c}_{2,1}\bar{b}_{2,1}\dots\bar{c}_{2,n}\bar{b}_{2,n}$ ) be two abstract values in  $\bar{\mathbf{S}}$ . The semantics  $\mathfrak{S}_{\bar{\mathbf{S}}}$ , applied to  $\text{concat}(\bar{s}_1, \bar{s}_2)$  returns the concatenation of the two input segmentations. Formally,

$$\mathfrak{S}_{\bar{\mathbf{S}}}[\text{concat}](\bar{s}_1, \bar{s}_2) = \bar{s}_1 +_{\bar{\mathbf{S}}} \bar{s}_2$$

Note that if the last segment of  $\bar{s}_1$  and the first segment of  $\bar{s}_2$  share the same abstract character then these segments are unified. More precisely, given  $\bar{c}_{1,1}\bar{b}_{1,1}\dots\bar{c}_{1,n}\bar{b}_{1,n}$  and  $\bar{c}_{2,1}\bar{b}_{2,1}\dots\bar{c}_{2,n}\bar{b}_{2,n}$ , if  $\bar{c}_{1,n} =_{\bar{\mathbf{C}}} \bar{c}_{2,1}$  then their concatenation is equal to

$$\bar{c}_{1,1}\bar{b}_{1,1}\dots\bar{c}[\bar{e}_{1,n} \oplus_{\bar{\mathbf{E}}} \bar{e}_{2,1}, \bar{e}'_{1,n} \oplus_{\bar{\mathbf{E}}} \bar{e}'_{2,1}]\bar{c}_{2,2}\bar{b}_{2,2}\dots\bar{c}_{2,n}\bar{b}_{2,n}$$

where  $\bar{c}$  represents the character contained in  $\bar{c}_{1,n}$  and  $\bar{c}_{2,1}$ .

- Substring

Let  $\bar{s}$  (i.e.,  $\bar{c}_1\bar{b}_1\dots\bar{c}_n\bar{b}_n$ ) be an abstract value in  $\bar{\mathbf{S}}$ . The semantics  $\mathfrak{S}_{\bar{\mathbf{S}}}$ , applied to  $\text{substring}_{\bar{\mathbf{B}}}^{\bar{\mathbf{e}}}(\bar{s})$  returns the subsegmentation of  $\bar{s}$  from the segment whose associated interval on the indexes that may refer to it contains  $\mathbf{b}$  to the segment whose associated interval on the indexes that may refer to it contains  $\mathbf{e}$ , otherwise it returns  $\top_{\bar{\mathbf{S}}}$ . Formally,

$$\mathfrak{S}_{\bar{\mathbf{S}}}[\text{substring}_{\bar{\mathbf{B}}}^{\bar{\mathbf{e}}}](\bar{s}) =$$

$$\left\{ \begin{array}{ll} \bar{c}_k\bar{b}_k\dots\bar{c}_j\bar{b}_j & \text{if } \exists \langle \bar{c}_k\bar{b}_k, [i, i']_k \rangle, \langle \bar{c}_j\bar{b}_j, [i, i']_j \rangle \in \text{index}(\bar{s}) : \\ & k = \min(\{i \mid \langle \bar{c}_i\bar{b}_i, [i, i']_i \rangle \in \text{index}(\bar{s}) \wedge \mathbf{b} \in [i, i']_i\}) \wedge \\ & j = \max(\{i \mid \langle \bar{c}_i\bar{b}_i, [i, i']_i \rangle \in \text{index}(\bar{s}) \wedge \mathbf{e} \in [i, i']_i\}) \\ \top_{\bar{\mathbf{S}}} & \text{otherwise} \end{array} \right.$$

where  $\text{index}(\bar{s})$  is the function which associates each segment of  $\bar{s}$  to the interval of the index abstract values that may refer to it, i.e.,  $\text{index}(\bar{s}) = \{\langle \bar{c}_i\bar{b}_i, [i, i']_i \rangle \mid i \in [1, \text{numSeg}(\bar{s})]\}$ . For instance, consider the segmentation  $\bar{s} = \text{'a' } [0, 2] \text{'b' } [4, 4] \text{'c' } [0, 2]$ ,  $\text{index}(\bar{s}) = \{\langle \text{'a' } [0, 2], [0, 1] \rangle, \langle \text{'b' } [4, 4], [0, 5] \rangle, \langle \text{'c' } [0, 2], [4, 7] \rangle\}$  and the subsegmentation of  $\bar{s}$  from the index 2 to the index 5 is  $\text{'b' } [4, 4] \text{'c' } [0, 2]$ .

- Is contained

Let  $\bar{s}$  (i.e.,  $\bar{c}_1\bar{b}_1\dots\bar{c}_n\bar{b}_n$ ) be an abstract value in  $\bar{\mathbf{S}}$ . The semantics  $\mathfrak{B}_{\bar{\mathbf{S}}}$ , applied to  $\text{contains}_{\bar{\mathbf{C}}}(\bar{s})$  returns (i) **true** if there exists a segment abstract predicate in  $\bar{s}$  which approximates only the character  $\mathbf{c}$  and its segment bound lower limit is different from zero, (ii) **false** if does not exist a segment abstract predicate in  $\bar{s}$  which approximates only the character  $\mathbf{c}$  and its segment bound lower limit is different from zero, otherwise (iii)  $\top_{\bar{\mathbf{B}}}$ . Formally,

$$\mathfrak{B}_{\bar{\mathbf{S}}}\llbracket \text{contains}_c \rrbracket(\bar{\mathbf{s}}) = \begin{cases} \text{true} & \text{if } \exists \bar{c}_i[\bar{e}_i, \bar{e}'_i] \in \bar{\mathbf{s}} : \text{char}_{\bar{c}}(\bar{c}_i) = \{c\} \wedge \bar{e}_i \neq_{\bar{\mathbf{E}}} 0 \\ \text{false} & \text{if } \nexists \bar{c}_i[\bar{e}_i, \bar{e}'_i] \in \bar{\mathbf{s}} : \text{char}_{\bar{c}}(\bar{c}_i) = \{c\} \wedge \bar{e}_i \neq_{\bar{\mathbf{E}}} 0 \\ \top_{\mathbf{B}} & \text{otherwise} \end{cases}$$

where  $\text{char}_{\bar{c}}(\bar{c}) = \{c : \langle i, c \rangle \in \gamma_{\bar{c}}(\bar{c})\}$ .

#### 4.5.4 Soundness

We prove the soundness of the Segmentation abstract semantics defined above.

**Theorem 4.2.**  $\mathfrak{G}_{\bar{\mathbf{S}}}$  is a sound over-approximations of  $\mathfrak{G}$ . Formally,

$$\gamma_{\bar{\mathbf{S}}}(\mathfrak{G}_{\bar{\mathbf{S}}}\llbracket \text{Stm} \rrbracket(\bar{\mathbf{s}})) \supseteq \{\mathfrak{G}\llbracket \text{Stm} \rrbracket(\sigma) \mid \sigma \in \gamma_{\bar{\mathbf{S}}}(\bar{\mathbf{s}})\}$$

*Proof.*

We prove the soundness separately for each operator.

- Consider the **new String** operator and let  $\sigma$  be a sequence of characters. We have to prove that,

$$\gamma_{\bar{\mathbf{S}}}(\mathfrak{G}_{\bar{\mathbf{S}}}\llbracket \text{new String}(\sigma) \rrbracket()) \supseteq \{\mathfrak{G}\llbracket \text{new String}(\sigma) \rrbracket()\}$$

This holds by the definitions of  $\mathfrak{G}_{\bar{\mathbf{S}}}$  and  $\gamma_{\bar{\mathbf{S}}}$ .

- Consider the **concat** operation and let  $\bar{\mathbf{s}}_1, \bar{\mathbf{s}}_2 \in \bar{\mathbf{S}}$ . We have to prove that,

$$\gamma_{\bar{\mathbf{S}}}(\mathfrak{G}_{\bar{\mathbf{S}}}\llbracket \text{concat} \rrbracket(\bar{\mathbf{s}}_1, \bar{\mathbf{s}}_2)) \supseteq \{\mathfrak{G}\llbracket \text{concat} \rrbracket(\sigma_1, \sigma_2) \mid \sigma_1 \in \gamma_{\bar{\mathbf{S}}}(\bar{\mathbf{s}}_1) \wedge \sigma_2 \in \gamma_{\bar{\mathbf{S}}}(\bar{\mathbf{s}}_2)\}$$

Let  $\sigma$  be a generic element in  $\{\mathfrak{G}\llbracket \text{concat} \rrbracket(\sigma_1, \sigma_2) \mid \sigma_1 \in \gamma_{\bar{\mathbf{S}}}(\bar{\mathbf{s}}_1) \wedge \sigma_2 \in \gamma_{\bar{\mathbf{S}}}(\bar{\mathbf{s}}_2)\}$ , that represents any possible concatenation between a string taken from  $\gamma_{\bar{\mathbf{S}}}(\bar{\mathbf{s}}_1)$  and a string taken from  $\gamma_{\bar{\mathbf{S}}}(\bar{\mathbf{s}}_2)$ . Therefore  $\sigma$  belongs to  $\gamma_{\bar{\mathbf{S}}}(\mathfrak{G}_{\bar{\mathbf{S}}}\llbracket \text{concat} \rrbracket(\bar{\mathbf{s}}_1, \bar{\mathbf{s}}_2))$  because  $\mathfrak{G}_{\bar{\mathbf{S}}}\llbracket \text{concat} \rrbracket(\bar{\mathbf{s}}_1, \bar{\mathbf{s}}_2) = \bar{\mathbf{s}}_1 +_{\bar{\mathbf{S}}} \bar{\mathbf{s}}_2$ , by definition of  $\mathfrak{G}_{\bar{\mathbf{S}}}$ . Indeed  $\gamma_{\bar{\mathbf{S}}}(\bar{\mathbf{s}}_1 +_{\bar{\mathbf{S}}} \bar{\mathbf{s}}_2)$  contains all the strings which come from the concatenation of  $\bar{\mathbf{s}}_1$  and  $\bar{\mathbf{s}}_2$ , by definition of  $\gamma_{\bar{\mathbf{S}}}$ .

- Consider the **substring** operation and let  $\bar{\mathbf{s}} \in \bar{\mathbf{S}}$ . we have to prove that,

$$\gamma_{\bar{\mathbf{S}}}(\mathfrak{G}_{\bar{\mathbf{S}}}\llbracket \text{substring}_b^e \rrbracket(\bar{\mathbf{s}})) \supseteq \{\mathfrak{G}\llbracket \text{substring}_b^e \rrbracket(\sigma) \mid \sigma \in \gamma_{\bar{\mathbf{S}}}(\bar{\mathbf{s}})\}$$

We have two possible cases:

- If both **b** and **e** belong to an interval of the indexes associated with each segment of a segmentation  $\bar{\mathbf{s}}$ , then  $\mathfrak{G}_{\bar{\mathbf{S}}}\llbracket \text{substring}_b^e \rrbracket(\bar{\mathbf{s}})$  returns the subsegmentation of  $\bar{\mathbf{s}}$  from the leftmost segment whose associated indexes interval contains **b** to the

rightmost segment whose associated indexes interval contains  $e$ , by definition of  $\mathfrak{S}_{\bar{s}}$ . This means that any of the concrete strings in  $\gamma_{\bar{s}}(\bar{s})$  contain the substring from the index  $b$  to the index  $e$ , by definition of  $\mathfrak{S}$  and it follows that any of the substrings in  $\{\mathfrak{S}\llbracket\text{substring}_b^e\rrbracket(\sigma) \mid \sigma \in \gamma_{\bar{s}}(\bar{s})\}$  is also contained in  $\gamma_{\bar{s}}(\mathfrak{S}\llbracket\text{substring}_b^e\rrbracket(\bar{s}))$ .

- Otherwise  $\mathfrak{S}\llbracket\text{substring}_b^e\rrbracket(\bar{s})$  returns  $\top_{\bar{s}}$  that over-approximates any possible results of the concrete semantics.

□

**Theorem 4.3.**  $\mathfrak{B}_{\bar{s}}$  is a sound over-approximations of  $\mathfrak{B}$ . Formally,

$$\gamma_{\bar{s}}(\mathfrak{B}_{\bar{s}}\llbracket\text{Stm}\rrbracket(\bar{s})) \supseteq \{\mathfrak{B}\llbracket\text{Stm}\rrbracket(\sigma) \mid \sigma \in \gamma_{\bar{s}}(\bar{s})\}$$

*Proof.*

Consider the `contains` operator and let  $\bar{s} \in \bar{\mathbf{S}}$ . We have to prove that,

$$\gamma_{\bar{s}}(\mathfrak{B}_{\bar{s}}\llbracket\text{contains}_c\rrbracket(\bar{s})) \supseteq \{\mathfrak{B}\llbracket\text{contains}_c\rrbracket(\sigma) \mid \sigma \in \gamma_{\bar{s}}(\bar{s})\}$$

We have three possible cases:

- If any concrete string in  $\gamma_{\bar{s}}(\bar{s})$  contains the character  $c$  it means that  $c$  is propagated in a segmentation  $\bar{s}$  by an abstract domain  $\bar{\mathbf{C}}$  and it appears in a segment whose segment bound lower limit is different from zero, by definition of  $\mathfrak{B}_{\bar{s}}$ . Then both the concrete and the abstract semantics return `true`.
- If any concrete string in  $\gamma_{\bar{s}}(\bar{s})$  does not contain the character  $c$  it means that  $c$  is propagated in a segmentation  $\bar{s}$  by an abstract domain  $\bar{\mathbf{C}}$  and it appears in a segment whose segment bound lower limit is equal to zero, or it is not propagated at all, by definition of  $\mathfrak{B}_{\bar{s}}$ . Then both the concrete and the abstract semantics return `false`.
- Otherwise, if not all the concrete strings in  $\alpha_{\bar{s}}(\bar{s})$  contain the character  $c$  it means that  $c$  is over-approximated in a segmentation  $\bar{s}$  by an abstract domain  $\bar{\mathbf{C}}$ , by definition of  $\mathfrak{B}_{\bar{s}}$ . Then both the concrete and the abstract semantics return  $\top_{\mathbf{B}}$ .

□

## 4.6 Refined String Abstract Domains

To obtain our refined string abstract domains, we exploit the notion of Granger product [87]. The Granger product is a binary operator between abstract domains, based on two refinement operators (cf. 2.6.4). Let  $\bar{\mathbf{D}}_1$  and  $\bar{\mathbf{D}}_2$  be two abstract domains. The Granger operators iteratively lift the information of  $\bar{\mathbf{D}}_1$  using the information of  $\bar{\mathbf{D}}_2$ , and vice-versa, until the smallest reduction is obtained. We will show that the reductions we will provide

in the next sections can be achieved without iterating the refinement operators. Formally, our Granger operators are as follows:

- $\text{refine}_{\overline{\mathbf{D}}_2}: \overline{\mathbf{D}}_1 \rightarrow \overline{\mathbf{D}}_1$
- $\text{refine}_{\overline{\mathbf{D}}_1}: \overline{\mathbf{D}}_2 \rightarrow \overline{\mathbf{D}}_2$

Given two abstract elements  $\overline{\mathbf{d}}_1$  and  $\overline{\mathbf{d}}_2$ , which belong to  $\overline{\mathbf{D}}_1$  and  $\overline{\mathbf{D}}_2$  respectively, a precondition to compute the reduced product between  $\overline{\mathbf{d}}_1$  and  $\overline{\mathbf{d}}_2$  (i.e.,  $\overline{\mathbf{d}}_1 \otimes \overline{\mathbf{d}}_2$ ) is that there must not be inconsistency (**Inc**) between them, i.e.,  $\gamma_{\overline{\mathbf{D}}_1}(\overline{\mathbf{d}}_1) \cap \gamma_{\overline{\mathbf{D}}_2}(\overline{\mathbf{d}}_2) \neq \emptyset$ . In case of inconsistency, the reduced product leads to the pair of bottom elements of their respective abstract domains, formally:

$$\overline{\mathbf{d}}_1 \otimes \overline{\mathbf{d}}_2 = \begin{cases} (\perp_{\overline{\mathbf{D}}_1}, \perp_{\overline{\mathbf{D}}_2}) & \text{if } \mathbf{Inc}(\overline{\mathbf{d}}_1, \overline{\mathbf{d}}_2) \\ (\text{refine}_{\overline{\mathbf{d}}_2}(\overline{\mathbf{d}}_1), \text{refine}_{\overline{\mathbf{d}}_1}(\overline{\mathbf{d}}_2)) & \text{otherwise} \end{cases}$$

#### 4.6.1 Meaning of Refinement

Before presenting our refined string abstract domains, we intuitively sketch what refinement means for any of the domains involved in the reductions. An improvement of the precision of an abstract element in (i) String Length domain means decreasing its range, (ii) Character Inclusion domain means either increasing the cardinality of the set of certainly contained characters and/or decreasing the cardinality of the set of maybe contained characters, (iii) Prefix domain means increasing the length of its sequence of characters, concatenating one or more characters to its end, and (iv) Segmentation domain means either increasing the precision of its segment abstract predicates and/or decreasing the range of the segment bounds.

#### 4.6.2 Combining Segmentation and String Length Domains

The first combination involves the abstract elements  $\overline{\mathbf{s}}$  and  $\overline{\mathbf{n}}$ , which belong to the Segmentation and the String Length abstract domains, respectively, where  $\overline{\mathbf{n}}$  denotes the interval  $[m, M]$  introduced in Section 4.4.1 (cf. Appendix B.1 for more details). As mentioned above, we define the reduced product between  $\overline{\mathbf{s}}$  and  $\overline{\mathbf{n}}$  in terms of two refinement operators, improving the approximation reached by one domain through the other and vice-versa. Thus, below, we introduce the notion of *inconsistency* between two abstract elements in the Segmentation and in the String Length domains and the definitions of both refinement operators involved in the reduction of the latter abstract elements. Finally, we prove the soundness of the reduction operators we present. We will follow this line-up also for the reductions we will present later in the chapter.

In the following we denote by  $\text{minLen}_{\overline{\mathbf{S}}}(\overline{\mathbf{s}})$  the minimal length of a segmentation abstract predicate  $\overline{\mathbf{s}}$ , i.e., the sum of all its segment lower bounds  $\overline{\mathbf{e}}_i$ , and by  $\text{maxLen}_{\overline{\mathbf{S}}}(\overline{\mathbf{s}})$  its maximal

length, i.e., the sum of all its segment upper bounds  $\bar{e}_i'$ . Formally, let  $\bar{s} = (\bar{c}_i[\bar{e}_i, \bar{e}_i'] )_{i=1}^n$  with  $n \geq 1$  (shortcut for  $\bar{c}_1[\bar{e}_1, \bar{e}_1'] \dots \bar{c}_n[\bar{e}_n, \bar{e}_n']$ ) then,

$$\begin{aligned} \bullet \minLen_{\bar{S}}(\bar{s}) &= \begin{cases} \sum_{i=1}^n \bar{e}_i & \text{if } \forall i \in \bar{s} : \bar{e}_i \text{ does not contain variables with constant value} \\ 0 & \text{otherwise} \end{cases} \\ \bullet \maxLen_{\bar{S}}(\bar{s}) &= \begin{cases} \sum_{i=1}^n \bar{e}_i' & \text{if } \forall i \in \bar{s} : \bar{e}_i' \text{ does not contain variables with constant value} \\ \top_{\bar{E}} & \text{otherwise} \end{cases} \end{aligned}$$

**Definition 4.1** (Inconsistency between Segmentation and String Length). Let  $\bar{s}$  and  $\bar{n}$  be two abstract elements in the Segmentation and in the String Length domains, respectively. Moreover, let  $\gamma_{\bar{S}}$  and  $\gamma_{\bar{SL}}$  be the concretization functions of the domains just mentioned. There is inconsistency between  $\bar{s}$  and  $\bar{n}$  when:

$$\gamma_{\bar{S}}(\bar{s}) \cap \gamma_{\bar{SL}}(\bar{n}) = \emptyset.$$

△

**Lemma 4.1** (Inconsistency between Segmentation and String Length). Let  $\bar{s}$  and  $\bar{n}$  be two abstract elements in the Segmentation and in the String Length domains, respectively. If one of the following conditions holds, then  $\bar{s}$  and  $\bar{n}$  are inconsistent.

1. If the intersection between the interval from the minimal to the maximal length of  $\bar{s}$  and  $\bar{n}$  is empty, then  $\bar{s}$  and  $\bar{n}$  are inconsistent. Formally:

$$[\minLen_{\bar{S}}(\bar{s}), \maxLen_{\bar{S}}(\bar{s})] \cap [m, M] = \emptyset \Rightarrow \mathbf{Inc}(\bar{s}, \bar{n})$$

2. If one of the two abstract elements is the bottom element of their lattice, then  $\bar{s}$  and  $\bar{n}$  are inconsistent. Formally:

$$\bar{s} = \perp_{\bar{S}} \text{ or } \bar{n} = \perp_{\bar{SL}} \Rightarrow \mathbf{Inc}(\bar{s}, \bar{n})$$

*Proof.*

Case 1 The segmentation  $\bar{s}$  approximates strings having length that might vary from  $\minLen_{\bar{S}}(\bar{s})$  to  $\maxLen_{\bar{S}}(\bar{s})$ . The interval  $\bar{n}$  approximates all the possible strings having length that ranges from  $m$  to  $M$ . If the intersection of the two length intervals is empty,  $\bar{s}$  and  $\bar{n}$  are abstracting strings of different length. Thus, the intersection of the concrete string sets represented by  $\bar{s}$  and  $\bar{n}$  will be empty.

Case 2 Trivial, as the concretization of bottom elements is the empty set.

□



**Definition 4.2** (Segmentation refinement through String Length). Let  $\bar{s}$  and  $\bar{n}$  be two abstract elements in the Segmentation and in the String Length domains, respectively. We define  $\text{refine}_{\bar{n}}(\bar{s})$  as the refinement operator describing how the information tracked by  $\bar{n}$  improves the one from  $\bar{s}$ . Formally,  $\text{refine}_{\bar{n}}(\bar{s}) = \bar{s}^*$  where:

- a. if  $\bar{s}$  contains only one segment and its segment bound lower limit is smaller than or equal to the one of  $\bar{n}$ , then the refinement means substituting  $\bar{e}_i$  with  $m$ . Formally:

$$\bar{s} = \bar{c}_1[\bar{e}_1, \bar{e}'_1] \wedge \bar{e}_1 \leq_{\bar{E}} m \Rightarrow \bar{s}^* = \bar{s}[\bar{e}_1/\bar{e}'_1^*]$$

where  $\bar{e}_1^* = m$ .

- b. if  $\bar{s}$  contains a single segment and its bound upper limit is greater than or equal to the one of  $\bar{n}$ , then the refinement means substituting  $\bar{e}'_i$  with  $M$ . Formally:

$$\bar{s} = \bar{c}_1[\bar{e}_1, \bar{e}'_1] \wedge M \leq_{\bar{E}} \bar{e}'_1 \Rightarrow \bar{s}^* = \bar{s}[\bar{e}'_1/\bar{e}_1^*]$$

where  $\bar{e}_1^* = M$ .

Conditions a. and b. can occur at the same time. In this case, both refinements apply. Observe that the order of application is not relevant.

If neither a. nor b. occur, then  $\bar{s}^* = \bar{s}$

△

**Example 4.8.** Consider  $\bar{s} = \text{'a' } [3, 5 + x]$  where  $\bar{C}$  is the constant propagation domain for characters and  $\bar{n} = [4, 7]$ . Moreover, assume that the value of  $x$  is greater than 4. The refinement of  $\bar{s}$  through  $\bar{n}$  is as follows:

$$\begin{aligned} \text{refine}_{\bar{n}}(\bar{s}) &= \bar{s}[\bar{e}_1/\bar{e}'_1^*] \\ \bar{b}_1^* &= 7 \\ &= \text{'a' } [3, 7] \end{aligned}$$

Indeed the segment bound of  $\bar{s}$  has the upper limit equal to  $5+x$  where  $x$  is assumed to be an integer variable whose value is greater than 4, and the maximal length of the strings approximated by  $\bar{n}$  is 7. So we could substitute the segment bound upper limit of  $\bar{s}$  with the upper limit of  $\bar{n}$ , as by case b. of Definition 4.2.

◇

**Definition 4.3** (String Length refinement through Segmentation). Let  $\bar{s}$  and  $\bar{n}$  be two abstract elements in the Segmentation and in the String Length domains, respectively. We define  $\text{refine}_{\bar{s}}(\bar{n})$  as the refinement operator describing how the information tracked by  $\bar{s}$  improves the one from  $\bar{n}$ . Formally,  $\text{refine}_{\bar{s}}(\bar{n}) = \bar{n}^*$  where:

- a. if the lower limit of  $\bar{n}$  is smaller than the the minimal length of  $\bar{s}$ , then the refinement means substituting  $m$  with  $\text{minLen}_{\bar{s}}(\bar{s})$ . Formally:

$$m <_{\bar{E}} \text{minLen}_{\bar{s}}(\bar{s}) \Rightarrow \bar{n}^* = \bar{n}[m/\text{minLen}_{\bar{s}}(\bar{s})]$$

- b. if the upper limit of  $\bar{n}$  is greater than the maximal length of  $\bar{s}$ , then the refinement means substituting  $M$  with  $\text{maxLen}_{\bar{s}}(\bar{s})$ . Formally:

$$\text{maxLen}_{\bar{s}}(\bar{s}) <_{\bar{E}} M \Rightarrow \bar{n}^* = \bar{n}[M/\text{maxLen}_{\bar{s}}(\bar{s})]$$

Conditions a. and b. can occur at the same time. In this case both refinements apply. Observe that the order of application is not relevant.

If neither a. or b. occur then  $\bar{n}^* = \bar{n}$ .

△

**Example 4.9.** Consider  $\bar{s} = \text{'a' } [3, 7] \text{'b' } [1, 5] \text{'c' } [1, 1 + x]$  (with the assumptions of the example above) and  $\bar{n} = [4, 7]$ . The refinement of  $\bar{n}$  through  $\bar{s}$  is as follows:

$$\begin{aligned} \text{refine}_{\bar{s}}(\bar{n}) &= \bar{n}[m/\text{minLen}_{\bar{s}}(\bar{s})] \\ \text{minLen}_{\bar{s}}(\bar{s}) &= 5 \\ &= [5, 7] \end{aligned}$$

Indeed, the lower limit of  $\bar{n}$  is strictly smaller than the minimum length of  $\bar{s}$ . So we could replace the lower limit of  $\bar{n}$  with  $\text{minLen}_{\bar{s}}(\bar{s})$ , as by case a. of Definition 4.3.

◇

**Definition 4.4** (Reduced product between Segmentation and String Length). Let  $\bar{s}$  and  $\bar{n}$  be two abstract elements in the Segmentation and in the String Length abstract domains, respectively. Moreover, consider the refinement operators presented in Definitions 4.2 and 4.3. The reduced product of  $\bar{s}$  and  $\bar{n}$  is obtained by applying  $\text{refine}_{\bar{n}}$  and  $\text{refine}_{\bar{s}}$  as follows:

$$\bar{s} \otimes \bar{n} = (\text{refine}_{\bar{n}}(\bar{s}), \text{refine}_{\bar{s}}(\bar{n}))$$

△

**Lemma 4.2** (Soundness of Segmentation refinement). Let  $\bar{s}$  and  $\bar{n}$  be two abstract elements in the Segmentation and in the String Length abstract domains, respectively. Moreover, let  $\gamma_{\bar{s}}$  and  $\gamma_{\bar{sL}}$  be the concretization functions of the domains just mentioned, and let  $\text{refine}_{\bar{n}}(\bar{s}) = \bar{s}^*$ . The soundness of the Segmentation refinement can be expressed by:

$$\sigma \in [\gamma_{\bar{s}}(\bar{s}) \setminus \gamma_{\bar{s}}(\bar{s}^*)] \Rightarrow \sigma \notin \gamma_{\bar{sL}}(\bar{n})$$

*Proof.*

The set  $[\gamma_{\bar{s}}(\bar{s}) \setminus \gamma_{\bar{s}}(\bar{s}^*)]$  contains representations of strings whose length does not belong to the interval  $\bar{n}$ , either in case a. and in case b. of Definition 4.2. Thus, these strings do not belong to  $\gamma_{\bar{sL}}(\bar{n})$ , by definition of  $\gamma_{\bar{sL}}$ . □

**Lemma 4.3** (Soundness of String Length refinement). Let  $\bar{s}$  and  $\bar{n}$  be two abstract elements in the Segmentation and in the String Length abstract domains, respectively. Moreover, let  $\gamma_{\bar{s}}$  and  $\gamma_{\bar{sL}}$  be the concretization functions of the domains just mentioned, and let  $\text{refine}_{\bar{s}}(\bar{n}) = \bar{n}^*$ . The soundness of the String Length refinement can be expressed by:

$$\sigma \in [\gamma_{\bar{sL}}(\bar{n}) \setminus \gamma_{\bar{sL}}(\bar{n}^*)] \Rightarrow \sigma \notin \gamma_{\bar{s}}(\bar{s})$$

*Proof.*

The set  $[\gamma_{\bar{sL}}(\bar{n}) \setminus \gamma_{\bar{sL}}(\bar{n}^*)]$  contains representations of strings whose length does not belong to the interval that goes from the minimal to the maximal length of  $\bar{s}$ , either in case a. and in case b. of Definition 4.3. Thus, these strings do not belong to  $\gamma_{\bar{s}}(\bar{s})$ , by definition of  $\gamma_{\bar{s}}$  itself. □

### 4.6.3 Combining Segmentation and Character Inclusion Domains

The second combination involves the abstract elements  $\bar{s}$  and  $\bar{r}$ , which belong to the Segmentation and the Character Inclusion abstract domains respectively, such that  $\bar{r}$  denotes the pair of sets  $(C, MC)$  introduced in Section 4.4.2 (cf. Appendix B.2 for more details).

**Definition 4.5** (Inconsistency between Segmentation and Character Inclusion). Let  $\bar{s}$  and  $\bar{r}$  be two abstract elements in the Segmentation and in the Character Inclusion domains, respectively. Moreover, let  $\gamma_{\bar{s}}$  and  $\gamma_{\bar{CI}}$  be the concretization functions of the domains just mentioned. There is inconsistency between  $\bar{s}$  and  $\bar{r}$  when:

$$\gamma_{\bar{s}}(\bar{s}) \cap \gamma_{\bar{CI}}(\bar{r}) = \emptyset.$$

△

**Lemma 4.4** (Inconsistency between Segmentation and Character Inclusion). Let  $\bar{s}$  and  $\bar{r}$  be two abstract elements in the Segmentation and in the Character Inclusion domains, respectively. If one of the following conditions holds, then  $\bar{s}$  and  $\bar{r}$  are inconsistent.

1. If for each character in  $C$ , i.e., the set of certainly contained characters of  $\bar{r}$ , does not exist at least one segment abstract predicate  $\bar{c}$  in  $\bar{s}$  that approximates it, then  $\bar{s}$  and  $\bar{r}$  are inconsistent. Formally:

$$|C| = v \wedge \forall c \in C : \nexists^v i \in \bar{s} : c \in \text{char}_{\bar{CI}}(\bar{c}_i) \Rightarrow \mathbf{Inc}(\bar{s}, \bar{r})$$

where  $\exists^v$  is the existential quantifier with cardinality  $v$  introduced in Section 2.1.

2. If there exists one segment abstract predicate  $\bar{c}$  in  $\bar{s}$  which abstracts only characters that do not occur in MC, i.e., the set of maybe contained characters of  $\bar{r}$ , and its segment bound lower limit is strictly greater than 0, then  $\bar{s}$  and  $\bar{r}$  are inconsistent. Formally:

$$\exists i \in \bar{s} : \text{char}_{\bar{c}}(\bar{c}_i) \cap \text{MC} = \emptyset \wedge \bar{e}_i >_{\mathbf{E}} 0 \Rightarrow \mathbf{Inc}(\bar{s}, \bar{r})$$

3. If the set of certainly contained C and maybe contained characters MC of  $\bar{r}$  are both equal to the empty-set (i.e.,  $\bar{r}$  abstracts an empty string) and the minimum length of  $\bar{s}$  is strictly greater than 0, then  $\bar{s}$  and  $\bar{r}$  are inconsistent. Formally:

$$C = \text{MC} = \emptyset \wedge \text{minLen}_{\bar{s}}(\bar{s}) >_{\mathbf{E}} 0 \Rightarrow \mathbf{Inc}(\bar{s}, \bar{r})$$

4. If one of the two abstract elements corresponds to the bottom element of their lattice, then  $\bar{s}$  and  $\bar{r}$  are inconsistent. Formally:

$$\bar{s} = \perp_{\bar{s}} \text{ or } \bar{r} = \perp_{\text{CI}} \Rightarrow \mathbf{Inc}(\bar{s}, \bar{r})$$

*Proof.*

Case 1 The pair  $\bar{r} = (C, \text{MC})$  approximates strings certainly containing the characters in C and possibly containing the characters in  $[\text{MC} \setminus C]$ . If for each character in C does not exist at least one segment abstract predicate  $\bar{c}$  in  $\bar{s}$  approximating it, this means that  $\bar{s}$  is abstracting strings that do not simultaneously contain all the characters belonging to the set of certainly contained characters of  $\bar{r}$ . Thus, the intersection of the concrete string sets represented by  $\bar{s}$  and  $\bar{r}$  will be empty.

Case 2 If the segmentation  $\bar{s}$  contains at least one segment abstract predicate  $\bar{c}$  approximating characters that do not belong to the set of maybe contained characters of  $\bar{r}$  and its segment bound lower limit is strictly greater than 0, this means that  $\bar{s}$  is abstracting strings in which a character not present in  $\bar{r}$  surely occurs. Thus, the intersection of the concrete string sets represented by  $\bar{s}$  and  $\bar{r}$  will be empty.

Case 3 If  $\bar{r}$  approximates exactly the empty string and the minimum length of  $\bar{s}$  is strictly greater than 0 means that the concretization of  $\bar{s}$  does not contain the empty string. Thus, the intersection of the concrete string sets represented by  $\bar{s}$  and  $\bar{r}$  will be empty.

Case 4 Trivial, as the concretization of bottom elements is an empty set.

□

**Definition 4.6** (Segmentation refinement through Character Inclusion). Let  $\bar{s}$  and  $\bar{r}$  be two abstract elements in the Segmentation and in the Character Inclusion domains, respectively. We define  $\text{refine}_{\bar{r}}(\bar{s})$  as the refinement operator describing how the information tracked by  $\bar{r}$  improves the one from  $\bar{s}$ . Formally,  $\text{refine}_{\bar{r}}(\bar{s}) = \bar{s}^*$  where:

- a. if the number of segment abstract predicates  $\bar{c}$  in  $\bar{s}$  which (also) approximate the characters belonging to  $C$ , i.e., the set of certainly contained characters of  $\bar{r}$ , is equal to the cardinality of the latter set and, their segment bound lower limit is equal to 0, then the refinement means substituting those lower limits with 1. Formally:

$$|C| = v \wedge \exists v_i \in \bar{s} : \text{char}_{\bar{C}}(\bar{c}_i) \cap C \neq \emptyset \wedge \bar{e}_i =_{\bar{E}} 0 \Rightarrow \bar{s}^* = \bar{s}[\bar{e}_i/1]$$

- b. if the number of segment abstract predicates  $\bar{c}$  in  $\bar{s}$  which (also) approximate the characters belonging to  $C$ , i.e., the set of certainly contained characters of  $\bar{r}$ , is equal to the cardinality of the latter set, then the refinement means substituting those predicates, followed by a segment upper bound equal to 1, with the abstract value, in an abstract domain  $\bar{C}$ , of the set resulting from the intersection between the set of characters represented by  $\bar{c}$  and the set of certainly contained characters of  $\bar{r}$  (if it is preciser than  $\bar{c}$ ). Formally:

$$|C| = v \wedge \exists v_i \in \bar{s} : \text{char}_{\bar{C}}(\bar{c}_i) \cap C \neq \emptyset \wedge \bar{e}_i' =_{\bar{E}} 1 \\ \alpha_{\bar{C}}(\text{char}_{\bar{C}}(\bar{c}_i) \cap C) \sqsubseteq_{\bar{C}} \bar{c}_i \Rightarrow \bar{s}^* = \bar{s}[\bar{c}_i/\alpha_{\bar{C}}(\text{char}_{\bar{C}}(\bar{c}_i) \cap C)]$$

- c. if there are segment abstract predicates  $\bar{c}$  in  $\bar{s}$  which (also) approximate characters that do not belong to  $MC$ , i.e., the set of maybe contained characters of  $\bar{r}$ , then the refinement means substituting those predicates, followed by a segment upper bound strictly greater than 1, with the abstract value, in an abstract domain  $\bar{C}$ , of the set resulting from the intersection between the set of characters represented by  $\bar{c}$  and the set of maybe contained characters of  $\bar{r}$  (if it is preciser than  $\bar{c}$ ). Formally:

$$\exists i \in \bar{s} : \text{char}_{\bar{C}}(\bar{c}_i) \setminus MC \neq \emptyset \wedge \bar{e}_i' >_{\bar{E}} 1 \\ \alpha_{\bar{C}}(\text{char}_{\bar{C}}(\bar{c}_i) \cap MC) \sqsubseteq_{\bar{C}} \bar{c}_i \Rightarrow \bar{s}^* = \bar{s}[\bar{c}_i/\alpha_{\bar{C}}(\text{char}_{\bar{C}}(\bar{c}_i) \cap MC)]$$

The conditions above can occur at the same time. In this case, all refinements apply. Observe that the order of application is not relevant.

If none of the conditions above apply, then  $\bar{s}^* = \bar{s}$ .

△

**Example 4.10.** Consider  $\bar{s} = \text{'a' } [0, 2] \top_{\bar{C}\bar{P}} [0, 1] \text{'b' } [4, 8]$ , where  $\top_{\bar{C}\bar{P}}$  is the top element in the constant propagation domains for characters, and  $\bar{r} = (\{a, b, c\}, \{a, b, c, d\})$ . The

refinement of  $\bar{s}$  through  $\bar{r}$  is as follows:

$$\begin{aligned} \text{refine}_{\bar{r}}(\bar{s}) &= \bar{s}[\bar{e}_1/1, \bar{e}_2/1, \bar{c}_2/\alpha_{\bar{C}\bar{P}}(\{c\})] \\ &= 'a' [1, 2] 'c' [1, 1] 'b' [4, 8] \end{aligned}$$

Indeed, the number of segment abstract predicates in  $\bar{s}$  approximating characters contained in the set of certainly contained characters of  $\bar{r}$  is equal to the cardinality of the latter set. So we substitute their segment bound lower limit with 1 when it was originally equal to zero, as by case a of Definition 4.6. Moreover, we substitute the second abstract predicate of  $\bar{s}$  with the value (in the constant propagation abstract domain) of the only certainly contained character of  $\bar{r}$  that is not directly present in  $\bar{s}$ , as by case b. of Definition 4.6, thus refining it.

◻

**Definition 4.7** (Character Inclusion refinement through Segmentation). Let  $\bar{s}$  and  $\bar{r}$  be two abstract elements in the Segmentation and in the Character Inclusion domains respectively. We define  $\text{refine}_{\bar{s}}(\bar{r})$  as the refinement operator describing how the information tracked by  $\bar{s}$  improves the one from  $\bar{r}$ . Formally,  $\text{refine}_{\bar{s}}(\bar{r}) = \bar{r}^*$  where:

- a. if there is a segment abstract predicate  $\bar{c}$  in  $\bar{s}$  which (also) approximates exactly one character belonging to MC, i.e., the set of maybe contained characters of  $\bar{r}$ , but that does not belong to C, i.e., the set of certainly contained characters of  $\bar{r}$ , and its segment bound lower limit is strictly greater than 0, then the refinement means adding that character to the set of certainly contained characters of  $\bar{r}$ . Formally:

$$\exists i \in \bar{s} : \text{char}_{\bar{C}}(\bar{c}_i) \cap (\text{MC} \setminus \text{C}) = \{c\} \wedge \bar{e}_i >_{\mathbb{E}} 0 \Rightarrow \bar{r}^* = \bar{r}[\text{C}/\text{C} \cup \{c\}]$$

- b. if the the set of characters represented by  $\bar{s}$  is a strict subset of MC, i.e., the set of maybe contained characters of  $\bar{r}$ , then the set of maybe contained characters of  $\bar{r}$  can be refined considering just the characters occurring in that intersection. Formally:

$$\text{char}_{\bar{C}}(\bar{s}) \subset \text{MC} \neq \text{MC} \Rightarrow \bar{r}^* = \bar{r}[\text{MC}/(\text{char}_{\bar{s}}(\bar{s}) \cap \text{MC})]$$

$$\text{where } \text{char}_{\bar{s}}(\bar{s}) = \bigcup_{\bar{c} \in \bar{s}} \text{char}_{\bar{C}}(\bar{c})$$

If none of the conditions above apply, then  $\bar{r}^* = \bar{r}$ .

△

**Example 4.11.** Consider  $\bar{s} = 'a' [2, 2] 'b' [3, 5] 'c' [4, 4] 'd' [0, 2]$  and  $\bar{r} = (\{a, b\}, \{a, b, c, d, e, f\})$ . The refinement of  $\bar{r}$  through  $\bar{s}$  is as follows:

$$\begin{aligned} \text{refine}_{\bar{s}}(\bar{r}) &= \bar{r}[C/(C \cup (\text{char}_{\bar{C}}(\bar{c}_3) \cap (MC \setminus C)), MC/\text{char}_{\bar{C}}(\bar{s})] \\ &= (\{a, b\} \cup (\{c\} \cap \{c, d, e, f\}), \{a, b, c, d\}) \\ &= (\{a, b\} \cup \{c\}, \{a, b, c, d\}) \\ &= (\{a, b, c\}, \{a, b, c, d\}) \end{aligned}$$

Indeed, the third segment abstract predicate of  $\bar{s}$  has a segment bound lower limit strictly greater than 0 and it approximates exactly the character 'c' which belongs to the set of maybe contained character of  $\bar{r}$ , but that does not belong to the set of certainly contained characters of  $\bar{r}$ . So we add the character 'c' to the set of certainly contained characters of  $\bar{r}$ , as by case a. of Definition 4.7. Moreover, we modify the set of maybe contained character of  $\bar{r}$  considering just the characters occurring in the set of characters approximated by  $\bar{s}$  ( $\text{char}_{\bar{s}}(\bar{s}) \cap MC = \{a, b, c, d\} \cap \{a, b, c, d, e, f\} = \{a, b, c, d\}$ ), as by case b. of Definition 4.7, thus refining it.

◻

**Definition 4.8** (Reduced product between Segmentation and Character Inclusion). Let  $\bar{s}$  and  $\bar{r}$  be two abstract elements in the Segmentation and in the Character Inclusion abstract domains, respectively. Moreover, consider the refinement operators presented in the Definitions 4.6 and 4.7. The reduced product of  $\bar{s}$  and  $\bar{r}$  is obtained by applying  $\text{refine}_{\bar{r}}$  and  $\text{refine}_{\bar{s}}$  as follows:

$$\bar{s} \otimes \bar{r} = (\text{refine}_{\bar{r}}(\bar{s}), \text{refine}_{\bar{s}}(\bar{r}))$$

△

**Lemma 4.5** (Soundness of Segmentation refinement). Let  $\bar{s}$  and  $\bar{r}$  be two abstract elements in the Segmentation and in the Character Inclusion abstract domains, respectively. Moreover, let  $\gamma_{\bar{s}}$  and  $\gamma_{\bar{CI}}$  be the concretization functions of the domains just mentioned, and let  $\text{refine}_{\bar{r}}(\bar{s}) = \bar{s}^*$ . The soundness of the Segmentation refinement can be expressed by:

$$\sigma \in [\gamma_{\bar{s}}(\bar{s}) \setminus \gamma_{\bar{s}}(\bar{s}^*)] \Rightarrow \sigma \notin \gamma_{\bar{CI}}(\bar{r})$$

*Proof.*

The set  $[\gamma_{\bar{s}}(\bar{s}) \setminus \gamma_{\bar{s}}(\bar{s}^*)]$  contains representations of strings in which do not occur at least one character that appears in the set of certainly contained characters of  $\bar{r}$ , either in case a. and in case b. of Definition 4.6, or at least one character that does not appear in the set of maybe contained characters of  $\bar{r}$ , in case c of Definition 4.6. Thus, those strings do not belong to  $\gamma_{\bar{CI}}(\bar{r})$ , by definition of  $\gamma_{\bar{CI}}$ .

◻

---

**Algorithm 5** isPrefix procedure.

---

**Input:**  $\bar{p}$  and  $\bar{s}$

**Output:**  $\text{res} \in \{\text{true}, \text{false}, \top_{\mathbf{B}}\}$

1:  $k \leftarrow 0, i \leftarrow 1, \text{res} \leftarrow \text{true}$

2:  $\bar{s} \leftarrow \text{unFold}(\bar{s})$

3: **while**  $k \in [0, \text{len}_{\overline{\mathbf{PR}}}(\bar{p}) - 1]$  **do**

4:     **while**  $i \in [1, \text{numSeg}(\bar{s})]$  **do**

5:         **if**  $\bar{p}[k] \in \text{char}_{\overline{\mathbf{C}}}(\bar{c}.\bar{s}[i]) \wedge \bar{b}.\bar{s}[i] \in \{[0, 1], [1, 1]\}$  **then**

6:              $k \leftarrow k + 1, i \leftarrow i + 1$

7:         **else if**  $(\bar{p}[k] \in \text{char}_{\overline{\mathbf{C}}}(\bar{c}.\bar{s}[i]) \wedge \bar{b}.\bar{s}[i] \in \{[0, x], [x, x]\}) \vee (\bar{p}[k] \notin \text{char}_{\overline{\mathbf{C}}}(\bar{c}.\bar{s}[i]) \wedge \bar{b}.\bar{s}[i] \in \{[x, x]\})$  **then**

8:             **return**  $\text{res} \leftarrow \top_{\mathbf{B}}$

9:         **else if**  $(\bar{p}[k] \notin \text{char}_{\overline{\mathbf{C}}}(\bar{c}.\bar{s}[i]) \wedge \bar{b}.\bar{s}[i] \in \{[0, 0], [0, x], [0, 1]\}) \vee (\bar{p}[k] \in \text{char}_{\overline{\mathbf{C}}}(\bar{c}.\bar{s}[i]) \wedge \bar{b}.\bar{s}[i] \in \{[0, 0]\})$  **then**

10:              $i \leftarrow i + 1$

11:         **else**

12:             **return**  $\text{res} \leftarrow \text{false}$

13: **return**  $\text{res}$

---

**Lemma 4.6** (Soundness of Character Inclusion refinement). Let  $\bar{s}$  and  $\bar{r}$  be two abstract elements in the Segmentation and in the Character Inclusion abstract domains, respectively, and let  $\gamma_{\bar{s}}$  and  $\gamma_{\overline{\mathbf{CI}}}$  be the concretization functions of the domains just mentioned, and let  $\text{refine}_{\bar{s}}(\bar{r}) = \bar{r}^*$ . The soundness of the Character Inclusion refinement can be expressed by:

$$\sigma \in [\gamma_{\overline{\mathbf{CI}}}(\bar{r}) \setminus \gamma_{\overline{\mathbf{CI}}}(\bar{r}^*)] \Rightarrow \sigma \notin \gamma_{\bar{s}}(\bar{s})$$

*Proof.*

The set  $[\gamma_{\overline{\mathbf{CI}}}(\bar{r}) \setminus \gamma_{\overline{\mathbf{CI}}}(\bar{r}^*)]$  contains representations of strings in which occurs at least one character that does not appear in the set of characters represented by  $\bar{s}$ , either in case a. and in case b. of Definition 4.7. Thus, those strings do not belong to  $\gamma_{\bar{s}}(\bar{s})$ , by definition of  $\gamma_{\bar{s}}$ . □

#### 4.6.4 Combining Segmentation and Prefix Domains

The last combination we present involves the abstract elements  $\bar{s}$  and  $\bar{p}$ , belonging to the Segmentation and the Prefix abstract domains, respectively. The reductions between the Segmentation and both the Suffix and the Prefix-Suffix abstract domains can be naturally induced by the Granger product among  $\bar{s}$  and  $\bar{p}$ . Notice that the refinement of  $\bar{s}$  is defined according to the length of the sequence of characters of  $\bar{p}$ .

Algorithm 5 determines if an abstract element in  $\overline{\mathbf{PR}}$  might be a prefix of (some of) the strings approximated by a segmentation in  $\overline{\mathbf{S}}$ . Note that  $\text{len}_{\overline{\mathbf{PR}}}(\bar{p})$  denotes the length of



the prefix  $\bar{p}$  and it is equal to the length of its sequence of characters, \* included. Once the segmentation has been unfolded (line 2),<sup>1</sup> we compare the characters in the prefix element  $\bar{p}$  and the segments in the unfolded segmentation  $\bar{s}$ . An element  $\bar{p}$  is a prefix of (some of) the strings approximated by  $\bar{s}$  if for each character of  $\bar{p}$  there exists a *corresponding* segment whose segment abstract predicate approximates (also) the character of  $\bar{p}$  under consideration, and its segment bound is of the type  $[0,1]$  or  $[1,1]$  (lines 5-6). The algorithm return  $\top_{\mathbf{B}}$  if the character of  $\bar{p}$  under consideration is approximated by a segment abstract predicate whose bound is of the type  $[0,x]$  or  $[x,x]$  (lines 7-8). It returns  $\top_{\mathbf{B}}$  also in the case in which a character of  $\bar{p}$  is not approximated by the segment abstract predicate under consideration if its segment bound is of the type  $[x,x]$ . The algorithm returns **false** if the corresponding segment of the character of  $\bar{p}$  under consideration is not empty and does not approximate it. In all the other cases, i.e., where possibly empty segments that do not approximate the character of  $\bar{p}$  under consideration or where possibly empty segments approximate the character of  $\bar{p}$  under consideration, the procedure moves the comparison to the next segment (lines 9-10).

**Definition 4.9** (Inconsistency between Segmentation and Prefix). Let  $\bar{s}$  and  $\bar{p}$  be two abstract elements in the Segmentation and in the Prefix domains respectively. Moreover, let  $\gamma_{\bar{\mathbf{S}}}$  and  $\gamma_{\bar{\mathbf{PR}}}$  be the concretization functions of the domains just mentioned. There is inconsistency between  $\bar{s}$  and  $\bar{p}$  when:

$$\gamma_{\bar{\mathbf{S}}}(\bar{s}) \cap \gamma_{\bar{\mathbf{PR}}}(\bar{p}) = \emptyset$$

△

**Lemma 4.7** (Inconsistency between Segmentation and Prefix). Let  $\bar{s}$  and  $\bar{p}$  be two abstract elements in the Segmentation and in the Prefix domains, respectively. If one of the following conditions holds, then  $\bar{s}$  and  $\bar{p}$  are inconsistent.

1. If `isPrefix` (cf. Algorithm 5) of  $\bar{p}$  and  $\bar{s}$  returns **false**, then  $\bar{s}$  and  $\bar{p}$  are inconsistent. Formally:

$$\text{isPrefix}(\bar{p}, \bar{s}) = \text{false} \Rightarrow \mathbf{Inc}(\bar{s}, \bar{p})$$

2. If one of the two abstract elements corresponds to the bottom element of their lattice, then  $\bar{s}$  and  $\bar{p}$  are inconsistent. Formally:

$$\bar{s} = \perp_{\bar{\mathbf{S}}} \text{ or } \bar{p} = \perp_{\bar{\mathbf{PR}}} \Rightarrow \mathbf{Inc}(\bar{s}, \bar{p})$$

---

<sup>1</sup>The `unFold` procedure can not be defined in general as it depends on the normal form of the expressions in  $\bar{\mathbf{E}}$ , but it is reasonable to assume that an unfolded segmentation can have a limited set of segment bounds types, that are:  $[0,0]$ ,  $[0,1]$ ,  $[1,1]$ ,  $[0,x]$ ,  $[x,x]$ , with  $x \in \bar{\mathbf{X}}$  (cf. Section 4.5) if expressions contain variables.

*Proof.*

Case 1 The prefix  $\bar{p}$  approximates all the possible strings sharing the same starting sequence of characters. The segmentation  $\bar{s}$  approximates strings highlighting sequences of equal characters. If the `isPrefix` procedure on  $\bar{p}$  and  $\bar{s}$  returns `false`, then  $\bar{p}$  and  $\bar{s}$  are abstracting strings with different prefixes. Thus, the intersection of the concrete string sets represented by  $\bar{s}$  and  $\bar{p}$  will be empty.

Case 2 Trivial, as the concretization of bottom elements is the empty set.

□

**Definition 4.10** (Segmentation refinement through Prefix). Let  $\bar{s}$  and  $\bar{p}$  be two abstract elements in the Segmentation and in the Prefix domains, respectively. We define  $\text{refine}_{\bar{p}}(\bar{s})$  as the refinement operator describing how the information tracked by  $\bar{p}$  improves the one from  $\bar{s}$ . Let  $\text{unfold}(\bar{s}) = \bar{s}'$ . Formally,  $\text{refine}_{\bar{p}}(\bar{s}) = \bar{s}^*$  where:

- a. if the abstraction of a character in  $\bar{p}$ , in an abstract domain  $\bar{C}$ , precedes the segment abstract predicate of its matching segment in  $\bar{s}'$ , then the refinement means substituting that segment abstract predicate with the abstract value of its corresponding character in  $\bar{p}$ . Formally:

$$\forall i \in \bar{s}' \forall k \in [0, \text{len}_{\bar{p}\bar{R}}(\bar{p}) - 1]: \bar{c}_i \bar{b}_i \text{ is the corresponding segment of } \bar{p}[k] \wedge \\ \alpha_{\bar{C}}(\bar{p}[k]) \sqsubseteq_{\bar{C}} \bar{c}_i \Rightarrow \text{fold}(\bar{s}^* = \bar{s}'[\bar{c}_i / \alpha_{\bar{C}}(\bar{p}[k])])$$

- b. if the matching segment in  $\bar{s}'$  of a character in  $\bar{p}$  is preceded by one or more segments possibly empty and none of those is the matching segment of the character preceding the first mentioned in  $\bar{p}$ , then the refinement means removing these segments from  $\bar{s}'$ . Formally:

$$\forall i \in \bar{s}' \forall k \in [0, \text{len}_{\bar{p}\bar{R}}(\bar{p}) - 1]: \bar{c}_i \bar{b}_i \text{ is the corresponding segment of a } \bar{p}[k] \wedge \\ \exists (\bar{c}_j \bar{b}_j)_{j=i-1}^{n \in [i-1, 1]} \in \bar{s}' : \bar{c}_j \bar{b}_j \text{ is not the corresponding segment of } \bar{p}[k-1] \wedge \\ \bar{p}[k] \notin \text{char}_{\bar{C}}(\bar{c}_j) \wedge \bar{b}_j \in \{[0, 0], [0, 1], [0, x]\} \Rightarrow \text{fold}(\bar{s}^* = \bar{s}' -_{\bar{s}} (\bar{c}_j \bar{b}_j)_{j=i-1}^{n \in [i-1, 1]})$$

The conditions above can occur at the same time. In this case, both refinements apply. Observe that the order of application is not relevant.

If none of the conditions above apply, then  $\bar{s}^* = \bar{s}$ .

△

**Example 4.12.** Consider  $\bar{s} = 'a' [2, 4] 'b' [1, 1] 'c' [0, 2]$  and  $\bar{p} = aab*$ . Let  $\bar{s}'$  denote  $\text{unFold}(\bar{s}) = 'a' [1, 1] 'a' [1, 1] 'a' [0, 1] 'a' [0, 1] 'b' [1, 1] 'c' [0, 1] 'c' [0, 1]$ . The refinement of  $\bar{s}$  through  $\bar{p}$  is as follows:

$$\begin{aligned} \text{refine}_{\bar{p}}(\bar{s}) &= \text{fold}(\bar{s}' - \bar{s} (\bar{c}_k \bar{b}_k)_{k=4}^3) \\ &= \text{fold}('a' [1, 1] 'a' [1, 1] 'b' [1, 1] 'c' [0, 1] 'c' [0, 1]) \\ &= \bar{s} = 'a' [2, 2] 'b' [1, 1] 'c' [0, 2] \end{aligned}$$

Indeed, the fifth segment of  $\bar{s}'$  is the matching segment for  $\bar{p}[2]$  and it is preceded by two possibly empty segments and none of them is the matching segment of  $\bar{p}[1]$ . So we could remove them from  $\bar{s}'$ , as by case b. of Definition 4.10, thus refining it.

◻

**Definition 4.11** (Prefix refinement through Segmentation). Let  $\bar{s}$  and  $\bar{p}$  be two abstract elements in the Segmentation and in the Prefix domains, respectively. We define  $\text{refine}_{\bar{s}}(\bar{p})$  as the refinement operator describing how the information tracked by  $\bar{s}$  improves the one from  $\bar{p}$ . Let  $\text{unFold}(\bar{s}) = \bar{s}'$ . Formally,  $\text{refine}_{\bar{s}}(\bar{p}) = \bar{p}^*$  where: if there is one or more segments in  $\bar{s}'$  after the one matching the last character of  $\bar{p}$ , such that their segment abstract element approximates exactly one character and their bound is equal to the interval  $[1, 1]$ , then the refinement means concatenating that character(s) to  $\bar{p}$ . Formally:

$$\begin{aligned} \exists i \in \bar{s}': \bar{c}_i \bar{b}_i \text{ is the corresponding segment of } \bar{p}[\text{len}_{\overline{\text{PR}}}(\bar{p}) - 1] \wedge \\ \exists (\bar{c}_j \bar{b}_j)_{j=i+1}^{n \in [i+1, \text{numSeg}(\bar{s}')] } \in \bar{s}' : \text{char}_{\bar{c}}(\bar{c}_j) = \{c\} \wedge \\ \bar{b}_j = [1, 1] \Rightarrow \bar{p}^* = \bar{p} + \langle c : \text{char}_{\bar{c}}(\bar{c}_j) = \{c\} \wedge j \in [i + 1, n] \rangle \end{aligned}$$

If the condition above does not apply, then  $\bar{p}^* = \bar{p}$ .

△

**Example 4.13.** Consider  $\bar{s} = 'a' [2, 3] 'b' [2, 2]$  and  $\bar{p} = aab*$ . Let  $\bar{s}'$  denote  $\text{unFold}(\bar{s}) = 'a' [1, 1] 'a' [1, 1] 'a' [0, 1] 'b' [1, 1] 'b' [1, 1]$ . The refinement of  $\bar{p}$  through  $\bar{s}$  is as follows:

$$\begin{aligned} \text{refine}_{\bar{s}}(\bar{p}) &= aaa + \langle b \rangle \\ &= aab + b \\ &= aabb* \end{aligned}$$

Indeed, the segment of  $\bar{s}'$  that corresponds to the last character in  $\bar{p}$  is followed by a segment whose character abstract element approximates exactly the character 'b' and its segment bound is the interval  $[1, 1]$ . So we could concatenate to the sequence of character of  $\bar{p}$ , the character approximated by  $\bar{c}_5$ , as by Definition 4.11, thus refining it.

◻

**Definition 4.12** (Reduced product between Segmentation and Prefix). Let  $\bar{s}$  and  $\bar{p}$  be two abstract elements in the Segmentation and in the Prefix abstract domains, respectively. Moreover, consider the refinement operators presented in the Definitions 4.10 and 4.11. The reduced product of  $\bar{s}$  and  $\bar{p}$  is obtained by applying  $\text{refine}_{\bar{p}}$  and  $\text{refine}_{\bar{s}}$  as follows:

$$\bar{s} \otimes \bar{p} = (\text{refine}_{\bar{p}}(\bar{s}), \text{refine}_{\bar{s}}(\bar{p}))$$

△

**Lemma 4.8** (Soundness of Segmentation refinement). Let  $\bar{s}$  and  $\bar{p}$  be two abstract elements in the Segmentation and in the Prefix abstract domains, respectively. Moreover, let  $\gamma_{\bar{s}}$  and  $\gamma_{\overline{\text{PR}}}$  be the concretization functions of the domains just mentioned, and let  $\text{refine}_{\bar{p}}(\bar{s}) = \bar{s}^*$ . The soundness of the Segmentation refinement can be expressed by:

$$\sigma \in [\gamma_{\bar{s}}(\bar{s}) \setminus \gamma_{\bar{s}}(\bar{s}^*)] \Rightarrow \sigma \notin \gamma_{\overline{\text{PR}}}(\bar{p})$$

*Proof.*

The set  $[\gamma_{\bar{s}}(\bar{s}) \setminus \gamma_{\bar{s}}(\bar{s}^*)]$  contains representations of strings which do not share the same longest common prefix with the strings approximated by  $\bar{p}$ , either in case a. and in case b. of Definition 4.10. Thus, those strings do not belong to  $\gamma_{\overline{\text{PR}}}(\bar{p})$ , by definition of  $\gamma_{\overline{\text{PR}}}$ . □

**Lemma 4.9** (Soundness of Prefix refinement). Let  $\bar{s}$  and  $\bar{p}$  be two abstract elements in the Segmentation and in the Prefix abstract domains respectively. Moreover, let  $\gamma_{\bar{s}}$  and  $\gamma_{\overline{\text{PR}}}$  be the concretization functions of the domains just mentioned, and let  $\text{refine}_{\bar{s}}(\bar{p}) = \bar{p}^*$ . The soundness of the Prefix refinement can be expressed by:

$$\sigma \in [\gamma_{\overline{\text{PR}}}(\bar{p}) \setminus \gamma_{\overline{\text{PR}}}(\bar{p}^*)] \Rightarrow \sigma \notin \gamma_{\bar{s}}(\bar{s})$$

*Proof.*

The set  $[\gamma_{\overline{\text{PR}}}(\bar{p}) \setminus \gamma_{\overline{\text{PR}}}(\bar{p}^*)]$  contains representations of strings that do not share the same longest common prefix with the strings approximated by  $\bar{s}$ , as by Definition 4.11. Thus, those strings do not belong to  $\gamma_{\bar{s}}(\bar{s})$ , by definition of  $\gamma_{\bar{s}}$ . □

## 4.7 Discussion

Strings are characterized by their content and shape. Existing string abstract domains capture information about the content or the shape of strings separately, or they can track both but partially, leading to a consistent loss of precision during the analysis. The Segmentation domain specialised for string analysis tracks information about the shape and the content of strings in their entirety. Its precision depends on the abstract domains

for characters, segment bounds, and variables it has been instantiated. More precise domains lead the segmentation analysis efficiency to decrease inevitably. Instead, less precise domains increase the efficiency of the segmentation analysis at the expense of losing important information, e.g., information useful to prevent undesired behaviours of the considered program.

By lifting basic domains for strings by segmentation abstraction, we generate more sophisticated domains where structural information and character-based information are combined, yielding more accurate representation.



## Chapter 5

---

# COMPLETENESS OF STRING DOMAINS

In this chapter, we study the problem of completeness of string abstractions, and we show how to complete two well-known string abstract domains for `JavaScript`. Moreover, we provide a procedure to measure the precision gained by the complete domain with respect to its original version.

The contribution of this chapter was published in [18].

### Chapter Structure

Section 5.2 recalls relevant concepts related to the completeness property in Abstract Interpretation. Moreover, a motivating example is given to show the importance of guaranteeing completeness in an Abstract Interpretation-based analysis with respect to strings. Section 5.3 introduces the syntax and the concrete semantics of the dynamic imperative core language used in this chapter. Section 5.4 presents the completion of two well-known string abstract domains with respect to two operations of interest. Section 5.5 highlights the strengths and usefulness of the completeness approach to static analysis of string manipulating programs. Section 5.6 defines the measurement procedure of the analysis precision increment. Section 5.7 concludes.

## 5.1 Introduction

Important features of Abstract Interpretation are *soundness* (or correctness) and *completeness* [56] (cf. Section 2.6.2). While soundness should always be guaranteed, as a basic requirement, by static analysis tools, to avoid the presence of false negatives, completeness is frequently not met.

If completeness is satisfied, the abstract computations do not lose information during the abstraction process with respect to a property of interest. The analysis can be therefore considered optimal. Guaranteeing completeness of string analysis is of particular interest, especially for dynamic languages, as poorly managed string manipulation code may easily

lead to significant security flaws. For instance, in the context of web applications, the common programming languages used for web-based software development (e.g., JavaScript) offer a wide range of dynamic features that make string manipulation dangerous.

In [80], Giacobazzi et al. highlighted the fact that completeness is an abstract domain property. In the same paper, they presented a methodology to obtain complete abstract domains with respect to operations by minimally extending or restricting the underlying domains.

## Contribution

We provide a way-to-proceed in the context of imprecise string abstractions. In particular, we exploit the theoretical framework of the complete shells (cf. Section 2.6.4), constructively showing how to improve precision of incomplete abstraction, without designing new string abstract domains.

We consider two JavaScript string abstract domains defined as part of the TAJIS [105] and SAFE [122] static analysers, focusing on their completeness with respect to the main string-manipulating operations. We compute their complete versions, and we discuss the benefits of guaranteeing completeness in the context of Abstract Interpretation-based string analysis of dynamic languages. Finally, we present an effective procedure to measure the precision increment when analysing a program with a complete abstract domain.

## 5.2 Making Abstract Interpretation Complete

In this section, we introduce the notions and methodologies that we will use through the whole chapter, as proposed in [80], to constructively build, from an initial abstract domain, a new abstract domain that is complete with respect to an operation of interest. Finally, through a motivating example, we show the usefulness of completing abstract domains for string analysis.

It is worth noting that completeness is a property related to the underlying abstract domain. Starting from this fact, Giacobazzi et al. proposed a constructive method to manipulate the underlying incomplete abstract domain to get a complete abstract domain with respect to a certain operation. In particular, given two abstract domains  $\overline{\mathbf{D}}$  and  $\overline{\mathbf{D}'}$  and an operator  $f : \overline{\mathbf{D}}^n \rightarrow \overline{\mathbf{D}'}$  with  $n \in \mathbb{N}$ , the authors gave two different notions of completion of abstract domains with respect to  $f$ . The first adds the minimal number of abstract points to the input abstract domain  $\overline{\mathbf{D}}$  and captures the notion of the complete shell of  $\overline{\mathbf{D}}$ . The second removes the minimal number of abstract points from the output abstract domain  $\overline{\mathbf{D}'}$ , and defines the complete core of  $\overline{\mathbf{D}'}$ .



### 5.2.1 Complete Shell vs Complete Core

We will focus on the construction of complete shells of string abstract domains, rather than complete cores. This choice is guided by the fact that a complete core for an operation  $f$  removes abstract points from a starting abstract domain. So, even if it is complete for  $f$ , its complete core could worsen the precision of other operations.

Conversely, complete shells augment the starting abstract domains (adding abstract points), and consequently, they cannot compromise the precision of other operations.

Below, we recall two important theorems proved in [80] that provide a constructive method to compute abstract domain complete shells, defined in terms of an upper closure operator  $\rho$  (cf. Section 2.2). Precisely, these theorems present two notions of complete shells:

- complete shells of  $\rho$  relative to  $\eta$  (where  $\eta$  is an upper closure operator), meaning that they are complete shells of operations defined on  $\rho$  that return results in  $\eta$ , and
- absolute complete shells of  $\rho$ , meaning that they are complete shells of operations defined on  $\rho$  and return results in  $\rho$ .

**Definition 5.1** (Complete shell of  $\rho$  relative to  $\eta$  [80]). Let  $\langle A, \sqsubseteq_A, \sqcup_A \rangle$  and  $\langle B, \sqsubseteq_B, \sqcup_B \rangle$  be two posets and  $f : A^n \rightarrow B$  be a continuous function. Given  $\rho \in uco(A)$  and  $\eta \in uco(B)$ ,<sup>1</sup> then let  $\mathcal{S}_f^\eta : uco(A) \rightarrow uco(A)$  be the domain transformer:

$$\mathcal{S}_f^\eta(\rho) = \bigsqcup_{uco(A)} \{ \delta \in uco(A) \mid \delta \sqsubseteq \rho, \langle \delta, \eta \rangle \in \Gamma(A, B, f) \}.$$

If  $\langle \mathcal{S}_f^\eta(\rho), \eta \rangle \in \Gamma(A, B, f)$ , then  $\mathcal{S}_f^\eta(\rho)$  is called complete shell of  $\rho$  relative to  $\eta$  with respect to an operation  $f$ .

△

As discussed in [80], Definition 5.1 does not offer a constructive methodology to compute  $\mathcal{S}_f^\eta(\rho)$ . The theorem below offers a constructive characterization of the complete shell of  $\rho$  relative to  $\eta$  with respect to  $f$ , making use of the Moore closure operator defined in Section 2.2.

**Theorem 5.1** ([80]). Let  $\langle A, \sqsubseteq_A, \sqcup_A \rangle$  and  $\langle B, \sqsubseteq_B, \sqcup_B \rangle$  be two posets and  $f : A^n \rightarrow B$  be a continuous function. Given  $\rho \in uco(A)$ ,  $\eta \in uco(B)$  and  $\mathcal{S}_f^\eta(\rho)$  as in Definition 5.1, the following equality holds:

$$\mathcal{S}_f^\eta(\rho) = \mathcal{M}(\rho \cup ( \bigsqcup_A \max^{\sqsubseteq_A} (\{z \in A \mid (f_x^i)(z) \sqsubseteq_B y\}))).$$

$i \in [1, n],$   
 $x \in A^n, y \in \eta$

---

<sup>1</sup>We recall that  $uco(A)$  denotes the set of upper closure operators of a poset  $A$ .

□

As already mentioned above, the idea behind the complete shell of  $\rho$  (input abstraction) relative to  $\eta$  (output abstraction) is to refine  $\rho$  adding the minimum number of abstract points to make  $\rho$  complete with respect to an operation  $f$ . By Theorem 5.1, this is obtained by adding to  $\rho$  the maximal elements in  $A$ , whose image under  $f$  is dominated by elements in  $\eta$ , at least in a single dimension  $i$ . Clearly, the so-obtained abstraction may not be an upper closure operator for  $A$ . Hence, the Moore closure operator is applied. Instead, absolute complete shells are involved in the case in which the operator  $f$  of interest has the same input and output abstract domain, i.e.,  $f : A^n \rightarrow A$ . In this case, given  $\rho \in uco(A)$ , absolute complete shells of  $\rho$  can be obtained as the greatest fix-point (**gfp**) of the domain transformer presented in Definition 5.1, as stated below.

**Definition 5.2** (Absolute complete shell of  $\rho$  [80]). Let  $\langle A, \sqsubseteq_A, \sqcup_A \rangle$  be a poset and let  $f : A^n \rightarrow A$  be a continuous function. Given  $\rho \in uco(A)$ , then let  $\bar{\mathcal{S}}_f : uco(A) \rightarrow uco(A)$  be the domain transformer:

$$\bar{\mathcal{S}}_f(\rho) = \bigsqcup_{uco(A)} \{ \delta \in uco(A) \mid \delta \sqsubseteq \rho, \langle \delta, \delta \rangle \in \Gamma(A, A, f) \}.$$

If  $\langle \bar{\mathcal{S}}_f(\rho), \bar{\mathcal{S}}_f(\rho) \rangle \in \Gamma(A, A, f)$ , then  $\bar{\mathcal{S}}_f(\rho)$  is called absolute complete shell of  $\rho$  with respect to an operation  $f$ .

△

**Theorem 5.2** ([80]). Let  $A$  be a poset and  $f : A^n \rightarrow A$  be a continuous function, and let  $\rho \in uco(A)$ . If  $\mathcal{S}_f^\rho(\rho)$  is the complete shell of  $\rho$  relative to  $\rho$  with respect to  $f$ , and if  $\bar{\mathcal{S}}_f(\rho)$  is the absolute complete shell of  $\rho$  with respect to  $f$ , then the following equality holds:

$$\bar{\mathcal{S}}_f(\rho) = \text{gfp}(\lambda\rho. \mathcal{S}_f^\rho(\rho)).$$

□

In [80], Giacobazzi et al. have discussed the completeness and incompleteness of the **Sign** numerical abstract domain. In particular,  $\text{Sign} = \{\text{NotPos}, \text{NotNeg}\} \cup \{\perp_{\text{Sign}}, \top_{\text{Sign}}\}$ , where: negative numbers are approximated by the abstract value **NotPos**, positive numbers by **NotNeg**, and  $\perp_{\text{Sign}}$  and  $\top_{\text{Sign}}$  are the bottom and the top elements of **Sign** respectively. **Sign** is complete for the product operation. Let  $*$  :  $\mathcal{P}(\mathbb{Z}) \times \mathcal{P}(\mathbb{Z}) \rightarrow \mathcal{P}(\mathbb{Z})$  be the concrete product operation and  $\otimes_{\text{Sign}} : \text{Sign} \times \text{Sign} \rightarrow \text{Sign}$  be the corresponding abstract product operation. Given the expression  $\mathbf{e}_1 * \mathbf{e}_2$ , the equality  $\alpha_{\text{Sign}}(\mathbf{e}_1 * \mathbf{e}_2) = \alpha_{\text{Sign}}(\mathbf{e}_1) \otimes_{\text{Sign}} \alpha_{\text{Sign}}(\mathbf{e}_2)$  holds, with  $\alpha_{\text{Sign}}$  being the abstraction function of **Sign**. As an example, consider the concrete

expression  $\{2, 5\} * \{-1, -3\}$ , then

$$\begin{aligned}
\alpha_{\text{Sign}}(\{2, 5\} * \{-1, -3\}) &= \alpha_{\text{Sign}}(\{-2, -5, -6, -15\}) \\
&= \text{NotPos} \\
&= \alpha_{\text{Sign}}(\{2, 5\}) \otimes_{\text{Sign}} \alpha_{\text{Sign}}(\{-1, -3\}) \\
&= \text{NotNeg} \otimes_{\text{Sign}} \text{NotPos} \\
&= \text{NotPos}
\end{aligned}$$

Instead, **Sign** is not complete for the sum operation. Let  $+$  :  $\mathcal{P}(\mathbb{Z}) \times \mathcal{P}(\mathbb{Z}) \rightarrow \mathcal{P}(\mathbb{Z})$  be the concrete sum operation and  $\oplus_{\text{Sign}} : \text{Sign} \times \text{Sign} \rightarrow \text{Sign}$  be the corresponding abstract sum operation. Consider the concrete expression  $\{2\} + \{-1, -2\}$ , then

$$\begin{aligned}
\alpha_{\text{Sign}}(\{2\} + \{-1, -2\}) &= \alpha_{\text{Sign}}(\{0, 1\}) \\
&= \text{NotNeg} \\
&\neq \alpha_{\text{Sign}}(\{2\}) \oplus_{\text{Sign}} \alpha_{\text{Sign}}(\{-1, -2\}) \\
&= \text{NotNeg} \oplus_{\text{Sign}} \text{NotPos} \\
&= \mathbb{Z}
\end{aligned}$$

In [80], the absolute complete shell of **Sign** with respect to the sum operation has been computed, and it corresponds to the Interval abstract domain [51].

### 5.2.2 Domain Completion Procedure

To improve the understanding about how to obtain a complete domain with respect to an operation of interest we provide a step-by-step reading of the formula in Theorem 5.1 (i.e., relative complete shell) by means of the procedure of Algorithm 6. The algorithm takes as input two posets  $\langle A, \sqsubseteq_A \rangle$ ,  $\langle B, \sqsubseteq_B \rangle$ , two closures  $\rho \in uco(A)$ ,  $\eta \in uco(B)$ , which correspond to the input and output abstraction respectively, and a continuous function  $f : A^{\text{d}} \rightarrow B$ . The procedure returns the complete shell of  $\rho$  relative to  $\eta$  with respect to  $f$ . Algorithm 6 follows Theorem 5.1 and collects in  $X$ , for each dimension  $i$  and element of  $\rho$ , any element  $z \in A$  whose image under  $f$  is dominated by elements in  $\eta$  (lines 2-4). Then, the starting input abstraction  $\rho$  is joined with the new elements collected in  $X$  and since the so obtained result may not be a closure, we apply the Moore closure operator (line 5). Finally, the complete shell of  $\rho$  relative to  $\eta$  is returned at line 6.

We reported the procedure only to improve the understanding of the complete shell since, unfortunately, this is not an effective procedure. Indeed, Algorithm 6 diverges (at lines 2-4) when  $A$  is an infinite set or  $\eta$  is not a finite closure.

---

**Algorithm 6** Relative complete shell procedure pseudo-code.

---

**Input:**

- $\langle A, \sqsubseteq_A \rangle, \langle B, \sqsubseteq_B \rangle$  (posets)
- $\rho \in uco(A), \eta \in uco(B)$  (abstractions)
- $f : A^n \rightarrow B$

**Output:**  $\mathcal{S}_f^\eta(\rho)$

- 1:  $X \leftarrow \emptyset$
  - 2: **for**  $i \in [1, n]$  and  $x \in A$  and  $y \in \eta$  **do**
  - 3:     **let**  $z \in A$  be the maximum element such that  $f_x^i(z) \sqsubseteq_B y$
  - 4:     **add**  $z$  to  $X$
  - 5: **let**  $\mathcal{S}_f^\eta(\rho)$  be the Moore closure of  $\rho \cup X$
  - 6: **return**  $\mathcal{S}_f^\eta(\rho)$
- 

### 5.2.3 Motivating Example

A common feature of dynamic languages, such as PHP or JavaScript, is to be weakly typed. Hence, in those languages, the variable type may be changed through the program execution. For example, in PHP, it is completely legal to write fragments such as `x = 1; x=true;`, where the type of the variable `x` changes from integer to boolean. The first attempt for statically reasoning about variable types was adopting the so-called Coalesced Sum abstract domain ( $\overline{\text{CS}}$ ) [13, 118], to detect whether a certain variable has constant type through the whole program execution. In Figure 5.1a, we show the Coalesced Sum abstract domain for an intra-procedural version of PHP [13], that tracks null, boolean, integer, float, and string types.<sup>2</sup> Consider the formal semantics of the sum operation in PHP [67]. When one of the operands is a string, *implicit type conversion* occurs and converts the operand string to a number since the sum operation only acts on numbers. In particular, if the prefix of the string is a number, it is converted to the maximum prefix of the string corresponding to a number, otherwise it is converted to 0. For example, the expression `e = "2.4hello" + "4"` returns 6.4. Let  $\alpha_{\overline{\text{CS}}}$  and  $\oplus_{\overline{\text{CS}}}$  be the abstraction function and the abstract sum operation on the Coalesced Sum abstract domain respectively. The type of the expression `e` is given by:  $\alpha_{\overline{\text{CS}}}(\{"2.4hello"\}) \oplus_{\overline{\text{CS}}} \alpha_{\overline{\text{CS}}}(\{"4"\}) = \text{String} \oplus_{\overline{\text{CS}}} \text{String} = \top_{\overline{\text{CS}}}$ . The static type analysis based on the Coalesced Sum abstract domain returns  $\top_{\overline{\text{CS}}}$  (i.e.,

---

<sup>2</sup>By closing the Coalesced Sum abstract domain with the powerset operation, we get a more precise domain called Union Type abstract domain [118], that tracks the set of types of a certain variable during program execution.

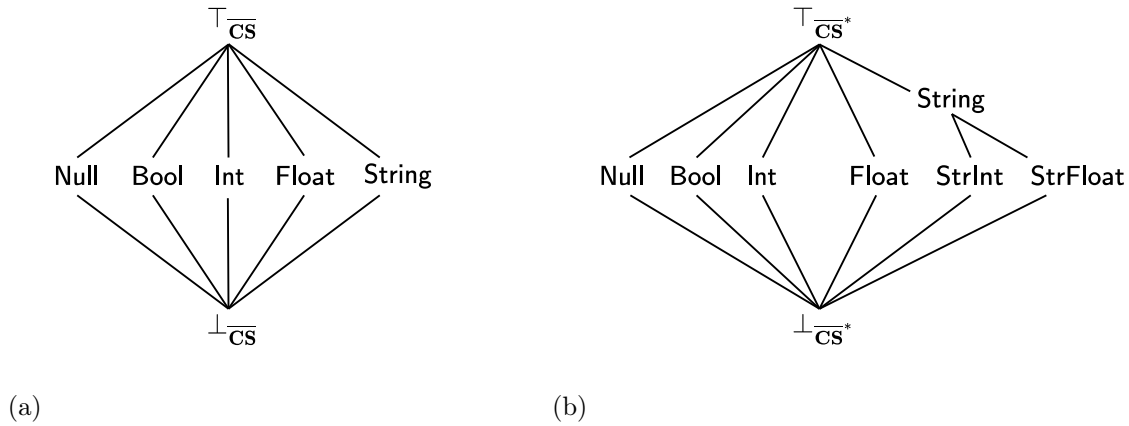


FIGURE 5.1: (a) Coalesced Sum abstract domain for PHP. (b) Complete shell of Coalesced Sum abstract domain with respect to the sum operation.

any possible value), since the sum between two strings may return either an integer or a float value. Precisely, the Coalesced Sum abstract domain is not complete with respect to the PHP sum operation, since for any string  $\sigma$  and  $\sigma'$ , it does not meet the completeness condition, i.e.,  $\alpha_{\overline{\text{CS}}}(\sigma + \sigma') = \alpha_{\overline{\text{CS}}}(\sigma) \oplus_{\overline{\text{CS}}} \alpha_{\overline{\text{CS}}}(\sigma')$ . Indeed,  $\alpha_{\overline{\text{CS}}}(\{"2.4\text{hello}" + "4"\}) = \text{Float}$  is different from  $\alpha_{\overline{\text{CS}}}(\{"2.4\text{hello}"\}) \oplus_{\overline{\text{CS}}} \alpha_{\overline{\text{CS}}}(\{"4"\}) = \top_{\overline{\text{CS}}}$ .

Intuitively, this happens because of the loss of precision that occurs during the abstraction process of the inputs, since the domain is not precise enough to distinguish between strings that may be implicitly converted to integers or floats.

Figure 5.1b shows the complete shell of the Coalesced Sum abstract domain with respect to the sum ( $\overline{\text{CS}}^*$ ). It adds two abstract values to the original domain, namely **StrFloat** and **StrInt**. They correspond to the abstractions of the strings that may be implicitly converted to floats and to integers, respectively. Notice that the type analysis on the new abstract domain is now complete with respect to the sum operation, e.g.,  $\alpha_{\overline{\text{CS}}^*}(\{"2.4\text{hello}" + "4"\}) = \text{Float}$  that is equal to  $\alpha_{\overline{\text{CS}}^*}(\{"2.4\text{hello}"\}) \oplus_{\overline{\text{CS}}^*} \alpha_{\overline{\text{CS}}^*}(\{"4"\}) = \text{StrFloat} \oplus_{\overline{\text{CS}}^*} \text{StrInt} = \text{Float}$ .

As pointed out above, guaranteeing completeness in Abstract Interpretation is a desirable property that an abstract domain should aim to, since it ensures that no loss of precision occurs during the input abstraction process of the operation of interest. It is worth noting that *guessing* a complete abstract domain for a certain operation becomes particularly hard when the operation has a tricky semantics, as in our example or, more in general, in dynamic languages operations. For this reason, complete shells become important since they can mathematically guarantee completeness for a given operation, starting from the abstract domain of interest.

```

a ::= n | x | a + a | a - a | a * a | a \ a | toNum(s) | length(s)
b ::= x | true | false | b && b | b || b | ! b
s ::= x | "σ" | concat(s1,s2)
e ::= a | b | s
st ::= x = e; | if (b) { st } else { st } | while (b) { st }
    | st1 st2
P ∈ μDyn ::= st

where x ∈ ID, σ ∈ Σ*, n ∈ INT ∪ FLOAT

```

FIGURE 5.2: μDyn syntax

### 5.3 Core Language

We define μDyn, an imperative toy language, inspired by the JavaScript programming language [128]. It is expressive enough to handle some interesting behaviors related to strings in dynamic languages, e.g., implicit type conversion.

Note that, in this chapter, the semantics notation, both concrete and abstract, is different from the one presented in the previous chapters to be consistent with the notation used in [18].

#### 5.3.1 Syntax

The syntax of μDyn is in Figure 5.2. The basic values are the set  $VAL = INT \cup FLOAT \cup BOOL \cup STR$ , where:  $INT$  is the set of signed integers  $\mathbb{Z}$ ,  $FLOAT$  is the set of signed decimal numbers,<sup>3</sup>  $BOOL$  is the set of booleans  $\{\text{true}, \text{false}\}$ , and  $STR$  is the set of strings  $\Sigma^*$  over an alphabet  $\Sigma$ . Then, let  $\Sigma^* = \Sigma_{Num}^* \cup \Sigma_{NotNum}^*$ , where:

$\Sigma_{Num}^*$  is the set of numeric strings (e.g., "42", "-7.2"), and

$\Sigma_{NotNum}^*$  is the set of non numeric strings (e.g., "foo", "-2a").

Moreover, we consider  $\Sigma_{Num}^*$  be additionally composed of four sets:

$$\Sigma_{Num}^* = \Sigma_{UInt}^* \cup \Sigma_{UFloat}^* \cup \Sigma_{SInt}^* \cup \Sigma_{SFloat}^*,$$

They correspond, from left to right, to the set of unsigned integer strings, unsigned float strings, signed integer strings and signed float strings, respectively.

---

<sup>3</sup>In general, floats are represented in programming languages in the IEEE 754 double precision format. For the sake of simplicity, we use instead decimal numbers.

$$\begin{aligned}
\llbracket x = e; \rrbracket \xi &= \xi[x \leftarrow \llbracket e \rrbracket \xi] \\
\llbracket \text{if } (b) \text{ } bl_1 \text{ else } bl_2 \rrbracket \xi &= \begin{cases} \llbracket bl_1 \rrbracket \xi & \llbracket b \rrbracket \xi = \text{true} \\ \llbracket bl_2 \rrbracket \xi & \llbracket b \rrbracket \xi = \text{false} \end{cases} \\
\llbracket \text{while } (b) \text{ } bl \rrbracket \xi &= \llbracket \text{if } (b) \{ bl \text{ while } (b) \text{ } bl \} \text{ else } \{ \} \rrbracket \xi \\
\llbracket \{ \} \rrbracket \xi &= \xi \quad \llbracket \{ st \} \rrbracket \xi = \llbracket st \rrbracket \xi \\
\llbracket st_1 st_2 \rrbracket \xi &= \llbracket st_2 \rrbracket (\llbracket st_1 \rrbracket \xi) \\
\llbracket \text{concat}(s, s') \rrbracket \xi &= \llbracket s \rrbracket \xi \cdot \llbracket s' \rrbracket \xi \quad \llbracket \text{length}(s) \rrbracket \xi = |\llbracket s \rrbracket \xi| \\
\llbracket \text{toNum}(s) \rrbracket \xi &= \begin{cases} \mathcal{N}(\llbracket s \rrbracket \xi) & \llbracket s \rrbracket \xi \in \Sigma_{\text{Num}}^* \\ 0 & \text{otherwise} \end{cases}
\end{aligned}$$

FIGURE 5.3:  $\mu$ Dyn semantics

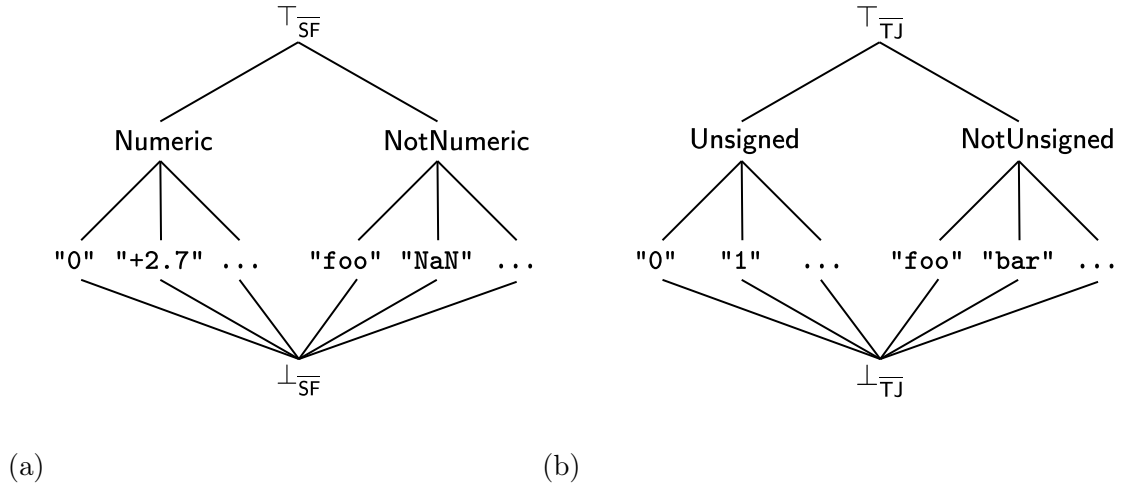
### 5.3.2 Concrete Semantics

The  $\mu$ Dyn programs are elements generated by `st` syntax rules. Program states  $\text{STATE} : \text{ID} \rightarrow \text{VAL}$ , ranged over  $\xi$ , are maps from identifiers to values. The concrete semantics of  $\mu$ Dyn statements follows [13],<sup>4</sup> and it is given by the function  $\llbracket \cdot \rrbracket : \text{STMT} \times \text{STATE} \rightarrow \text{STATE}$ , inductively defined on the structure of the statements, as in Figure 5.3. We abuse notation in defining the concrete semantics of expressions:  $\llbracket \cdot \rrbracket : \text{EXP} \times \text{STATE} \rightarrow \text{VAL}$ . Figure 5.3 shows the formal semantics of two relevant expressions involving strings we focus on: `concat`, that concatenates two strings, and string-to-number operation, namely `toNum`, that takes a string as input and returns the number that it represents if the input string corresponds to a numerical strings, 0 otherwise. Given  $\sigma \in \text{STR}$ , we denote by  $\mathcal{N}(\sigma) \in \text{INT} \cup \text{FLOAT}$  the numeric value of a given string. For example, `toNum("4.2") = 4.2` and `toNum("asd") = 0`.

## 5.4 Making JavaScript String Abstract Domains Complete

In this section, we study the completeness of two string abstract domains integrated into two state-of-the-art JavaScript static analysers based on Abstract Interpretation, namely TAJIS [105] and SAFE [122]. Both the abstract domains track important information on JavaScript strings, e.g., TAJIS can infer when a string corresponds to an unsigned integer, that may be used as an array index, and SAFE tracks numeric strings, such as "2.5" or "+5".

<sup>4</sup>Note that we only give the semantics of the operations that are of interest for this chapter. The complete concrete semantics of  $\mu$ Dyn is defined in [13].

FIGURE 5.4: (a) SAFE, (b) TAJN recasted for  $\mu\text{Dyn}$ .

For the sake of readability, we recast the original string abstract domains for  $\mu\text{Dyn}$ , following the notation adopted in [10]. Figure 5.4 depicts them. Notice that the original abstract domain part of SAFE analyser treats the string "NaN" as a numeric string. Since our core language does not provide the primitive value NaN, the corresponding string, i.e., "NaN", has no particular meaning here, and it is treated as a non-numerical string.

#### 5.4.1 Completing TAJN String Abstract Domain

Figure 5.4b depicts the string abstract domain  $\overline{\mathbf{TJ}}$ , the recasted version of the domain integrated into the TAJN static analyser [105].  $\overline{\mathbf{TJ}}$  splits the strings into **Unsigned**, that denotes the strings corresponding to unsigned numbers, and **NotUnsigned**, any other string. Hence, for example,  $\alpha_{\overline{\mathbf{TJ}}}(\{"9", "+9"\}) = \top_{\overline{\mathbf{TJ}}}$  and  $\alpha_{\overline{\mathbf{TJ}}}(\{"9.2", "foo"\}) = \text{NotUnsigned}$ . Before reaching these abstract values,  $\overline{\mathbf{TJ}}$  precisely tracks each string.

Here we focus on the `toNum` (i.e., string-to-number) operation. Since this operation clearly involves numbers, we introduce the TAJN numerical abstract domain, denoted by  $\overline{\mathbf{TJ}}_N$ . The latter domain behaves similarly to  $\overline{\mathbf{TJ}}$ , distinguishing between unsigned (**UnsignedInt**) and not unsigned integers (**NotUnsignedInt**).

Below we define the abstract semantics of the string-to-number operation for  $\overline{\mathbf{TJ}}$ . In particular, we define the function:  $\llbracket \text{toNum}(\bullet) \rrbracket^{\overline{\mathbf{TJ}}} : \overline{\mathbf{TJ}} \rightarrow \overline{\mathbf{TJ}}_N$  that takes as input a string abstract value in  $\overline{\mathbf{TJ}}$ , and returns an integer abstract value in  $\overline{\mathbf{TJ}}_N$ . Formally,

$$\llbracket \text{toNum}(s) \rrbracket^{\overline{\mathbf{TJ}}} = \begin{cases} \perp_{\overline{\mathbf{TJ}}_N} & \llbracket s \rrbracket^{\overline{\mathbf{TJ}}} = \perp_{\overline{\mathbf{TJ}}} \\ \llbracket \text{toNum}(\sigma) \rrbracket & \llbracket s \rrbracket^{\overline{\mathbf{TJ}}} = \sigma \\ \text{UnsignedInt} & \llbracket s \rrbracket^{\overline{\mathbf{TJ}}} = \text{Unsigned} \\ \top_{\overline{\mathbf{TJ}}_N} & \llbracket s \rrbracket^{\overline{\mathbf{TJ}}} = \text{NotUnsigned} \vee \llbracket s \rrbracket^{\overline{\mathbf{TJ}}} = \top_{\overline{\mathbf{TJ}}} \end{cases}$$



When the input evaluates to  $\perp_{\overline{\mathbf{TJ}}}$ , the bottom is propagated, and  $\perp_{\overline{\mathbf{TJ}}_N}$  is returned (first row). If the input evaluates to a single string value, the abstract semantics relies on its concrete one (second row), as single strings are precisely captured by  $\overline{\mathbf{TJ}}$ . When the input evaluates to the string abstract value `Unsigned`, the integer abstract value `UnsignedInt` is returned (third row). Indeed, in the concrete scenario, an unsigned and not float numeric string, which exactly represents the strings approximated by `Unsigned` in  $\overline{\mathbf{TJ}}$ , is converted into the correspondent numeric value by the string-to-number operation. Therefore, the abstraction in  $\overline{\mathbf{TJ}}_N$  of the numeric value of all the strings approximated by `Unsigned` is `UnsignedInt`. Finally, when the input evaluates to `NotUnsigned` or  $\top_{\overline{\mathbf{TJ}}}$ , the top abstract value  $\top_{\overline{\mathbf{TJ}}_N}$  is returned (fourth row). In the second case, i.e., when the input is evaluated to  $\top_{\overline{\mathbf{TJ}}}$ , it is trivial to note that the result of the abstract string-to-number operation represents the best correct approximation. However, in the first case, it is not straightforward. Note that `NotUnsigned` approximates not numeric strings and signed and/or float numeric strings. Thus the only safe abstraction of the numeric value of strings approximated by `NotUnsigned` is  $\top_{\overline{\mathbf{TJ}}_N}$ .

**Lemma 5.1.** The  $\overline{\mathbf{TJ}}$  abstract domain is not complete with respect to the `toNum` operation. In particular,<sup>5</sup>  $\exists S \in \mathcal{P}(\Sigma^*)$  such that:

$$\alpha_{\overline{\mathbf{TJ}}_N}(\llbracket \text{toNum}(S) \rrbracket) \subsetneq \llbracket \text{toNum}(\alpha_{\overline{\mathbf{TJ}}}(S)) \rrbracket^{\overline{\mathbf{TJ}}}$$

*Proof.*

As a counterexample to completeness, consider  $S = \{"2.3", "3.4"\}$ . We can show that  $\alpha_{\overline{\mathbf{TJ}}_N}(\llbracket \text{toNum}(S) \rrbracket) \neq \llbracket \text{toNum}(\alpha_{\overline{\mathbf{TJ}}}(S)) \rrbracket^{\overline{\mathbf{TJ}}}$ . Indeed,

$$\begin{aligned} \alpha_{\overline{\mathbf{TJ}}_N}(\llbracket \text{toNum}(S) \rrbracket) &= \text{NotUnsignedInt} \\ &\neq \llbracket \text{toNum}(\alpha_{\overline{\mathbf{TJ}}}(S)) \rrbracket^{\overline{\mathbf{TJ}}} \\ &= \llbracket \text{toNum}(\text{NotUnsigned}) \rrbracket^{\overline{\mathbf{TJ}}} \\ &= \top_{\overline{\mathbf{TJ}}_N} \end{aligned}$$

□

The completeness condition does not hold because the  $\overline{\mathbf{TJ}}$  string abstract domain loses too much information during the abstraction process, and it cannot be retrieved during the abstract `toNum` operation. In particular, when non-numeric strings and unsigned integer strings are converted to numbers by `toNum`, they are mapped to the same value, namely 0. Indeed,  $\overline{\mathbf{TJ}}$  does not differentiate between non-numeric and unsigned integer string values, and this is the principal cause of the  $\overline{\mathbf{TJ}}$  incompleteness with respect to

---

<sup>5</sup>We abuse notation denoting with  $\llbracket \cdot \rrbracket$  the additive lift to set of basic values of the concrete semantics, i.e., the collecting semantics.

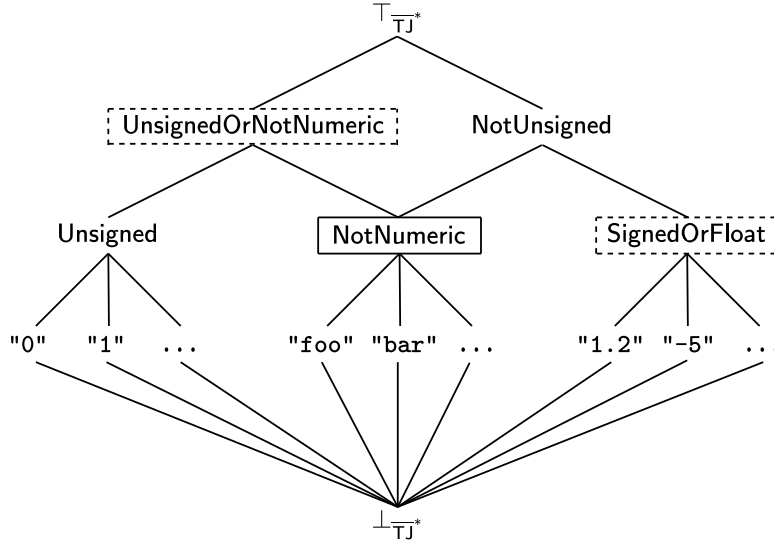


FIGURE 5.5: Complete shell of  $\rho_{\overline{\mathbf{TJ}}}$  relative to  $\rho_{\overline{\mathbf{TJ}}_N}$  w.r.t.  $\text{toNum}$ .

$\text{toNum}$ . Additionally, more precision can be obtained if we could differentiate numeric strings holding float numbers from those holding integer numbers. Thus, to make  $\overline{\mathbf{TJ}}$  complete with respect to  $\text{toNum}$ , we have to derive the complete shell of the  $\overline{\mathbf{TJ}}$  string abstract domain relative to the  $\overline{\mathbf{TJ}}_N$  numerical abstract domain, by applying Theorem 5.1.

**Definition 5.3** (Complete shell of  $\overline{\mathbf{TJ}}$ ). Let  $\rho_{\overline{\mathbf{TJ}}}$ ,  $\rho_{\overline{\mathbf{TJ}}_N}$  be the upper closure operators of  $\overline{\mathbf{TJ}}$  and  $\overline{\mathbf{TJ}}_N$  abstract domains respectively, i.e.,  $\rho_{\overline{\mathbf{TJ}}} \in \text{uco}(\mathcal{P}(\Sigma^*))$  and  $\rho_{\overline{\mathbf{TJ}}_N} \in \text{uco}(\mathcal{P}(\text{INT} \cup \text{FLOAT}))$ . Given  $\llbracket \text{toNum}(\bullet) \rrbracket^{\overline{\mathbf{TJ}}} : \overline{\mathbf{TJ}} \rightarrow \overline{\mathbf{TJ}}_N$ , we define by  $\overline{\mathbf{TJ}}^*$  the transformer  $S_{\text{toNum}}^{\rho_{\overline{\mathbf{TJ}}_N}}(\rho_{\overline{\mathbf{TJ}}})$ .

△

By Definition 5.1,  $\overline{\mathbf{TJ}}^*$  is the complete shell of  $\rho_{\overline{\mathbf{TJ}}}$  relative to  $\rho_{\overline{\mathbf{TJ}}_N}$  with respect to the  $\text{toNum}$  operation. As already argued in Section 5.4, to compute  $\overline{\mathbf{TJ}}^*$  we use the constructive characterization given by Theorem 5.1.

Figure 5.5 depicts  $\overline{\mathbf{TJ}}^*$ . In particular, the abstract points inside dashed boxes are the abstract values added during the computation of  $\overline{\mathbf{TJ}}^*$ , the point inside the solid box is instead obtained by the Moore closure of the other points of the domain, while the remaining abstract values were already in  $\overline{\mathbf{TJ}}$ . A non-intuitive point added in  $\overline{\mathbf{TJ}}^*$  is **SignedOrFloats**, namely the abstract value such that its concretization contains any float string and the signed integers. This abstract point is added during the computation of  $\overline{\mathbf{TJ}}^*$ . In particular, following Algorithm 6, instantiated with  $A = \mathcal{P}(\Sigma^*)$  and  $\eta = \rho_{\overline{\mathbf{TJ}}_N}$ , this abstract point is computed at lines 2-4 at the iteration when  $y$  is  $\gamma_{\overline{\mathbf{TJ}}_N}(\text{NotUnsignedInt})$ , i.e.,

$$\text{SignedOrFloats} \in \max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{toNum}(Z) \rrbracket \subseteq \gamma_{\overline{\mathbf{TJ}}_N}(\text{NotUnsignedInt})\})$$

Informally: which is the maximal set of strings  $Z$  such that  $\text{toNum}(Z)$  is dominated by  $\text{NotUnsignedInt}$ ? To obtain only values dominated by  $\text{NotUnsignedInt}$  from  $\text{toNum}(Z)$ , the maximal set doing so is exactly the set of the float strings and the signed strings. Other strings, such that: unsigned integer strings or not numerical strings are excluded, since they are both converted to unsigned integers, and they would violate the dominance relation. Similarly, the abstract point  $\text{UnsignedOrNotNumeric}$  is added to the complete shell  $\overline{\mathbf{TJ}}^*$ . Following again Algorithm 6, instantiated with  $A = \mathcal{P}(\Sigma^*)$  and  $\eta = \rho_{\overline{\mathbf{TJ}}_N}$ , the above abstract element is computed at lines 2-4 at the iteration when  $y$  is  $\gamma_{\overline{\mathbf{TJ}}_N}(\text{UnsignedInt})$ , i.e.,

$$\text{UnsignedOrNotNumeric} \in \max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \text{toNum}(Z) \subseteq \gamma_{\overline{\mathbf{TJ}}_N}(\text{UnsignedInt})\})$$

To obtain only values dominated by  $\text{UnsignedInt}$  from  $\text{toNum}(Z)$ , the maximal set doing so is exactly the set of the unsigned integer strings and the non-numerical strings, since the latter are converted to 0.

**Example 5.1.** Consider again the string set  $S = \{"2.3", "3.4"\}$  of Example 5.4.1. We can show that in  $\overline{\mathbf{TJ}}^*$ ,  $\alpha_{\overline{\mathbf{TJ}}_N}(\llbracket \text{toNum}(S) \rrbracket) = \llbracket \text{toNum}(\alpha_{\overline{\mathbf{TJ}}^*}(S)) \rrbracket^{\overline{\mathbf{TJ}}^*}$ . Indeed,

$$\begin{aligned} \alpha_{\overline{\mathbf{TJ}}_N}(\llbracket \text{toNum}(S) \rrbracket) &= \text{NotUnsignedInt} \\ &= \llbracket \text{toNum}(\alpha_{\overline{\mathbf{TJ}}^*}(S)) \rrbracket^{\overline{\mathbf{TJ}}^*} \\ &= \llbracket \text{toNum}(\text{SignedOrFloats}) \rrbracket^{\overline{\mathbf{TJ}}^*} \\ &= \text{NotUnsignedInt} \end{aligned}$$

◻

Note that the new abstract semantics  $\llbracket \text{toNum}(\bullet) \rrbracket^{\overline{\mathbf{TJ}}^*}$  handles the abstract points added by the complete shell. It corresponds to the best correct approximation, i.e.,  $\alpha_{\overline{\mathbf{TJ}}^*} \circ \llbracket \text{toNum}(\bullet) \rrbracket \circ \gamma_{\overline{\mathbf{TJ}}^*} : \overline{\mathbf{TJ}}^* \rightarrow \overline{\mathbf{TJ}}_N$  (cf. Definition 2.4).

### 5.4.2 Completing SAFE String Abstract Domain

Figure 5.4a depicts the string abstract domain  $\overline{\mathbf{SF}}$ , i.e., the recasted version of the domain involved into SAFE [122] static analyser. It splits strings into the abstract values:  $\text{Numeric}$  (i.e., numerical strings) and  $\text{NotNumeric}$  (i.e., all the other strings). As for  $\overline{\mathbf{TJ}}$ , before reaching these abstract values,  $\overline{\mathbf{SF}}$  precisely tracks single string values. For instance,  $\alpha_{\overline{\mathbf{SF}}}(\{"+9.6", "7"\}) = \text{Numeric}$ , and  $\alpha_{\overline{\mathbf{SF}}}(\{"+9.6", "\text{bar}"\}) = \perp_{\overline{\mathbf{SF}}}$ .

We study the completeness of  $\overline{\mathbf{SF}}$  with respect to  $\text{concat}$  operation. Figure 5.6 presents the abstract semantics of the concatenation operation for  $\overline{\mathbf{SF}}$ , that is:

$$\llbracket \text{concat}(\bullet, \bullet) \rrbracket^{\overline{\mathbf{SF}}} : \overline{\mathbf{SF}} \times \overline{\mathbf{SF}} \rightarrow \overline{\mathbf{SF}}$$

In particular, when both abstract values correspond to single strings, the standard string concatenation is applied (second row, second column). When one abstract value,

		$\llbracket \text{concat}(s_1, s_2) \rrbracket^{\overline{\mathbf{SF}}}$				
$\llbracket s_1 \rrbracket^{\overline{\mathbf{SF}}}$	$\llbracket s_2 \rrbracket^{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$	$\sigma_2 \in \Sigma^*$	Numeric	NotNumeric	$\top_{\overline{\mathbf{SF}}}$
$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$		$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$
$\sigma_1 \in \Sigma^*$	$\perp_{\overline{\mathbf{SF}}}$	$\sigma_1 \cdot \sigma_2$	$\left\{ \begin{array}{ll} \text{Numeric} & \sigma_1 = "" \text{ or} \\ & \sigma_1 \in \Sigma^*_{\text{UInt}} \\ \text{NotNumeric} & \text{otherwise} \end{array} \right.$		NotNumeric	$\top_{\overline{\mathbf{SF}}}$
Numeric	$\perp_{\overline{\mathbf{SF}}}$	$\left\{ \begin{array}{ll} \text{Numeric} & \sigma_2 = "" \text{ or} \\ & \sigma_2 \in \Sigma^*_{\text{UInt}} \\ \text{NotNumeric} & \text{otherwise} \end{array} \right.$		$\top_{\overline{\mathbf{SF}}}$	NotNumeric	$\top_{\overline{\mathbf{SF}}}$
NotNumeric	$\perp_{\overline{\mathbf{SF}}}$	NotNumeric		$\top_{\overline{\mathbf{SF}}}$	NotNumeric	$\top_{\overline{\mathbf{SF}}}$
$\top_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$	$\top_{\overline{\mathbf{SF}}}$		$\top_{\overline{\mathbf{SF}}}$	$\top_{\overline{\mathbf{SF}}}$	$\top_{\overline{\mathbf{SF}}}$

FIGURE 5.6:  $\overline{\mathbf{SF}}$  concat abstract semantics

involved in the concatenation, is a string and the other is **Numeric** (third row, second column and second row, third column) we distinguish two cases: if the string is empty or corresponds to an unsigned integer we can safely return **Numeric**, otherwise **NotNumeric** is returned. This happens because, when two float strings (hence numerical strings) are concatenated, a non-numerical string is returned (e.g.,  $\text{concat}("1.1", "2.2") = "1.12.2"$ ). For the same reason, when both input abstract values are **Numeric**, the result is not guaranteed to be numerical. Indeed,  $\llbracket \text{concat}(\text{Numeric}, \text{Numeric}) \rrbracket^{\overline{\mathbf{SF}}} = \top_{\overline{\mathbf{SF}}}$ .

**Lemma 5.2.** The  $\overline{\mathbf{SF}}$  abstract domain is not complete with respect to the `concat` operation. In particular,  $\exists S_1, S_2 \in \mathcal{P}(\Sigma^*)$  such that:

$$\alpha_{\overline{\mathbf{SF}}}(\llbracket \text{concat}(S_1, S_2) \rrbracket) \subsetneq \llbracket \text{concat}(\alpha_{\overline{\mathbf{SF}}}(S_1), \alpha_{\overline{\mathbf{SF}}}(S_2)) \rrbracket^{\overline{\mathbf{SF}}}$$

*Proof.*

As a counterexample to completeness, consider the sets  $S_1 = \{"2.2", "2.3"\}$  and  $S_2 = \{"2", "3"\}$ . Then, in  $\overline{\mathbf{SF}}$ ,  $\alpha_{\overline{\mathbf{SF}}}(\llbracket \text{concat}(S_1, S_2) \rrbracket) \neq \llbracket \text{concat}(\alpha_{\overline{\mathbf{SF}}}(S_1), \alpha_{\overline{\mathbf{SF}}}(S_2)) \rrbracket^{\overline{\mathbf{SF}}}$ . Indeed,

$$\begin{aligned} \alpha_{\overline{\mathbf{SF}}}(\llbracket \text{concat}(S_1, S_2) \rrbracket) &= \text{Numeric} \\ &\neq \llbracket \text{concat}(\alpha_{\overline{\mathbf{SF}}}(S_1), \alpha_{\overline{\mathbf{SF}}}(S_2)) \rrbracket^{\overline{\mathbf{SF}}} \\ &= \llbracket \text{concat}(\text{Numeric}, \text{Numeric}) \rrbracket^{\overline{\mathbf{SF}}} \\ &= \top_{\overline{\mathbf{SF}}} \end{aligned}$$

□

The  $\overline{\mathbf{SF}}$  abstract domain loses too much information during the abstraction process,

		$\max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(X, Z) \rrbracket \subseteq \gamma_{\overline{\mathbf{SF}}}(Y)\})$ : First Iteration				
Y		Numeric	NotNumeric	"n" $\in$ Z	"f" $\in$ F	"s" $\in$ NotNum
X	Y					
Numeric	Numeric	$\{\text{""}\} \cup \text{UInt}$	$[\text{NotNum} \setminus \{\text{""}\}] \cup \text{NotUInt} \cup \text{NotUFloat}$	$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$
NotNumeric	NotNumeric	$\perp_{\overline{\mathbf{SF}}}$	NotNumeric	$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$
"n" $\in$ Z	"n" $\in$ Z	$\{\text{""}\} \cup \text{UInt} \cup \text{UFloat}$	$[\text{NotNum} \setminus \{\text{""}\}] \cup \text{NotUInt} \cup \text{NotUFloat}$	$\perp_{\overline{\mathbf{SF}}} \vee \text{str}$	$\perp_{\overline{\mathbf{SF}}} \vee \text{str}$	$\perp_{\overline{\mathbf{SF}}} \vee \text{str}$
"f" $\in$ F	"f" $\in$ F	$\{\text{""}\} \cup \text{UInt}$	$[\text{NotNum} \setminus \{\text{""}\}] \cup \text{Float} \cup \text{NotUInt}$	$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}} \vee \text{str}$	$\perp_{\overline{\mathbf{SF}}} \vee \text{str}$
"s" $\in$ NotNum	"s" $\in$ NotNum	$\perp_{\overline{\mathbf{SF}}}$	NotNumeric	$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}}$	$\perp_{\overline{\mathbf{SF}}} \vee \text{str}$

TABLE 5.1: Completing  $\overline{\mathbf{SF}}$ .

which can not be retrieved during the abstract concatenation. Intuitively, to gain completeness with respect to `concat` operation,  $\overline{\mathbf{SF}}$  should improve the precision of the numerical strings abstraction, e.g., discriminating between float and integer strings.

**Definition 5.4** (Complete shell of  $\overline{\mathbf{SF}}$ ). Let  $\rho_{\overline{\mathbf{SF}}} \in \text{uco}(\mathcal{P}(\Sigma^*))$  be the upper closure operator related to  $\overline{\mathbf{SF}}$  abstract domain. Given  $\llbracket \text{concat}(\bullet, \bullet) \rrbracket^{\overline{\mathbf{SF}}} : \overline{\mathbf{SF}} \times \overline{\mathbf{SF}} \rightarrow \overline{\mathbf{SF}}$ , we define by  $\overline{\mathbf{SF}}^*$  the transformer  $\overline{\mathcal{S}}_{\text{concat}}(\rho_{\overline{\mathbf{SF}}})$ .

△

By Definition 5.2,  $\overline{\mathbf{SF}}^*$  is the absolute complete shell of  $\rho_{\overline{\mathbf{SF}}}$  with respect to the `concat` operation. By Theorem 5.2, the transformer  $\overline{\mathcal{S}}_{\text{concat}}(\rho_{\overline{\mathbf{SF}}})$  is equal to the Moore closure of the union between  $\overline{\mathbf{SF}}$  and the binary operator defined in Table 5.1. Table 5.1 depicts the first iteration of the fix-point computation of Theorem 5.2, where  $\perp_{\overline{\mathbf{SF}}}$  and  $\top_{\overline{\mathbf{SF}}}$  rows and columns are omitted for space limitations. In particular,  $\forall X \in \mathcal{P}(\Sigma^*)$ , when  $X \neq \perp_{\overline{\mathbf{SF}}}$ , we have that:

- $\max(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(Z, X) \rrbracket \subseteq \gamma_{\overline{\mathbf{SF}}}(\perp_{\overline{\mathbf{SF}}})\}) = \perp_{\overline{\mathbf{SF}}}$
- $\max(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(X, Z) \rrbracket \subseteq \gamma_{\overline{\mathbf{SF}}}(\perp_{\overline{\mathbf{SF}}})\}) = \perp_{\overline{\mathbf{SF}}}$
- $\max(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(X, Z) \rrbracket \subseteq \gamma_{\overline{\mathbf{SF}}}(\perp_{\overline{\mathbf{SF}}})\}) = \top_{\overline{\mathbf{SF}}}$

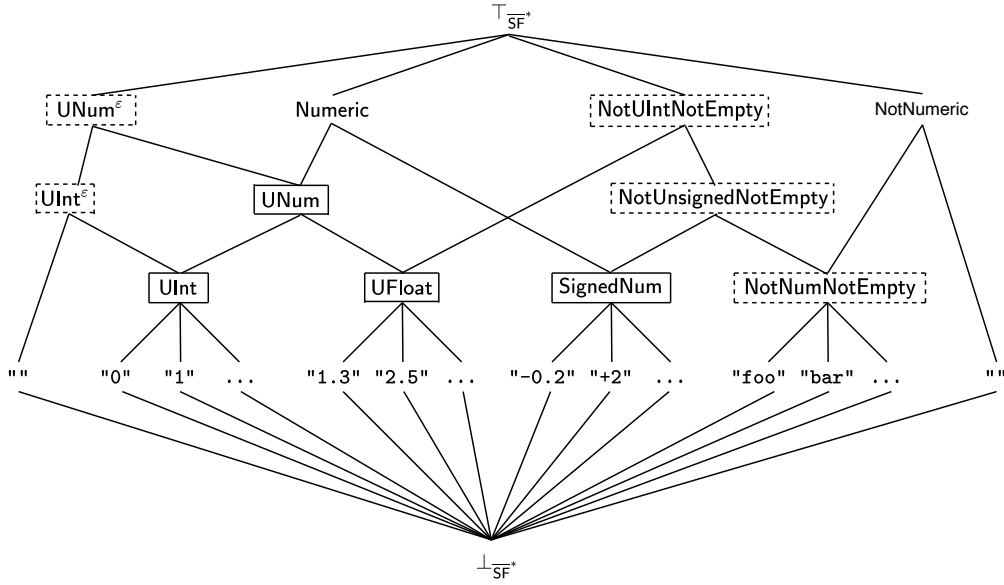


FIGURE 5.7: Absolute complete shell of  $\rho_{\overline{\mathbf{SF}}}$  w.r.t.  $\text{concat}$ .

Instead,  $\forall X \in \mathcal{P}(\Sigma^*)$ , when  $X \neq \top_{\overline{\mathbf{SF}}}$ , we have that:

$$\max(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(Z, X) \rrbracket \subseteq \gamma_{\overline{\mathbf{SF}}}(\perp_{\overline{\mathbf{SF}}})\}) = \perp_{\overline{\mathbf{SF}}}$$

The complete shell  $\overline{\mathbf{SF}}^*$  is in Figure 5.7. In particular, the points inside dashed boxes are the abstract values added during the iterative computations of  $\overline{\mathbf{SF}}^*$ , the points inside solid boxes are instead obtained by the Moore closure of the other points of the domain, while the remaining abstract values were already in  $\overline{\mathbf{SF}}$ . The meaning of abstract values in  $\overline{\mathbf{SF}}^*$  is intuitive. In order to satisfy the completeness property,  $\overline{\mathbf{SF}}^*$  splits the **Numeric** abstract value, already taken into account in  $\overline{\mathbf{SF}}$ , into all the strings corresponding to unsigned integer (**UInt**), unsigned floats (**UFloat**), and signed numbers (**SignedNum**). Moreover, particular importance is given to the empty string, since the new abstract domain specifies whether each abstract value contains the empty string  $""$ .<sup>6</sup> Indeed, the **UInt** <sup>$\epsilon$</sup>  abstract value represents the strings corresponding to unsigned integer or the empty string, and the **UNum** <sup>$\epsilon$</sup>  abstract value represents the strings corresponding to unsigned numbers or the empty string. An unexpected abstract value considered in  $\overline{\mathbf{SF}}^*$  is **NotUnsignedNotEmpty**, such that:

$$\gamma_{\overline{\mathbf{SF}}^*}(\text{NotUnsignedNotEmpty}) = \{\sigma \in \Sigma^* \mid \sigma \in \Sigma_{\text{SInt}}^* \cup \Sigma_{\text{SFloat}}^* \cup (\Sigma_{\text{NotNum}}^* \setminus \{""\})\}$$

Namely, the abstract point whose concretization corresponds to the set of any string corresponding to a signed number, and any non-numerical string, except the empty string. This abstract point has been added to  $\overline{\mathbf{SF}}^*$ , following Theorem 5.2. Since the absolute

<sup>6</sup>The empty strings ( $""$ ) at the most-left and most-right sides of Figure 5.7 are not distinct elements. They are only duplicated to declutter the figure.

complete shell is the greatest fix-point of the relative one (that in our case is reached after one iteration), the corresponding procedure is the greatest fix-point of Algorithm 6, instantiated with  $A = \mathcal{P}(\Sigma^*)$  and  $\eta = \rho_{\overline{\mathbf{SF}}}$ . In particular, the abstract element is computed at lines 2-4 at the iteration when  $y$  is  $\gamma_{\overline{\mathbf{SF}}}(\text{NotNumeric})$ ,  $x$  is  $\gamma_{\overline{\mathbf{SF}}}(\text{Numeric})$  and  $i = 1$ , namely `NotUnsignedNotEmpty` belongs to

$$\max^{\subseteq}(Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(\gamma_{\overline{\mathbf{SF}}}(\text{Numeric}), Z) \rrbracket \subseteq \gamma_{\overline{\mathbf{SF}}}(\text{NotNumeric}))$$

Informally: which is the maximal set of strings such that concatenated to any possible numerical string will produce any possible non-numerical string? Indeed, to be sure to obtain non-numerical strings, the maximal set doing so is exactly the set of any non-numerical non-empty string, and any string corresponding to a signed number, that is `NotUnsignedNotEmpty`.

**Example 5.2.** Let  $S_1 = \{"2.2", "2.3"\}$  and  $S_2 = \{"2", "3"\}$  be the string sets of Example 5.4.2. We can show that in  $\overline{\mathbf{SF}}^*$ ,  $\alpha_{\overline{\mathbf{SF}}^*}(\llbracket \text{concat}(S_1, S_2) \rrbracket) = \llbracket \text{concat}(\alpha_{\overline{\mathbf{SF}}^*}(S_1), \alpha_{\overline{\mathbf{SF}}^*}(S_2)) \rrbracket^{\overline{\mathbf{SF}}^*}$ . Indeed,

$$\begin{aligned} \alpha_{\overline{\mathbf{SF}}^*}(\llbracket \text{concat}(S_1, S_2) \rrbracket) &= \text{UFloat} \\ &= \llbracket \text{concat}(\alpha_{\overline{\mathbf{SF}}^*}(S_1), \alpha_{\overline{\mathbf{SF}}^*}(S_2)) \rrbracket^{\overline{\mathbf{SF}}^*} \\ &= \llbracket \text{concat}(\text{UFloat}, \text{UInt}) \rrbracket^{\overline{\mathbf{SF}}^*} \\ &= \text{UFloat} \end{aligned}$$

◻

As in the  $\overline{\mathbf{TJ}}^*$  case, the new abstract semantics  $\llbracket \text{concat}(\bullet, \bullet) \rrbracket^{\overline{\mathbf{SF}}^*}$  handles the abstract points added by the complete shell. It corresponds to the best correction approximation, i.e.,  $\alpha_{\overline{\mathbf{SF}}^*} \circ \llbracket \text{concat}(\bullet, \bullet) \rrbracket \circ \gamma_{\overline{\mathbf{SF}}^*} : \overline{\mathbf{SF}}^* \times \overline{\mathbf{SF}}^* \rightarrow \overline{\mathbf{SF}}^*$  (see Definition 2.4).

## 5.5 Benefits of Adopting Complete String Abstractions

Now, we discuss and evaluate the benefits gained from using the complete shells reported in Section 5.4 and, more in general, complete domains, with respect to a certain operation. In particular, we compare the  $\mu\text{Dyn}$  versions of the string abstract domains adopted by SAFE and TAJIS with their corresponding complete shells, and we discuss the complexity of the complete shells. Finally we argue how adopting complete abstract domains can be useful into static analysers.

### 5.5.1 Precision

In the previous section, we focused on the completeness of the string abstract domains integrated into SAFE and TAJIS, for  $\mu\text{Dyn}$ , with respect to two string operations, namely

`concat` and `toNum`, respectively. While string concatenation is common in any programming language, `toNum` assumes critical importance for dynamic languages, where implicit type conversion is provided. Since type conversion is often hidden to the developer, aim to completeness of the analysis increases the precision of such operations. For instance, let  $x$  be a variable, at a certain program execution point.  $x$  may take its concrete values in the set  $S = \{\text{"foo"}, \text{"bar"}\}$ . If  $S$  is abstracted into the original TAJs string abstract domain, its abstraction will correspond to `NotUnsigned`, losing the information about the fact that the concrete value of  $x$  surely does not contain numerical values. Hence, when the abstract value of  $S$  is used as input of `toNum`, the result will return  $\top_{\overline{\mathbf{TJ}}_N}$ , i.e., any possible concrete integer value. Conversely, abstracting  $S$  in  $\overline{\mathbf{TJ}}^*$  (the absolute complete shell of  $\overline{\mathbf{TJ}}$  relative to `toNum` discussed in Section 5.4.1) leads to a more precise abstraction, since  $\overline{\mathbf{TJ}}^*$  is able to differentiate between non-numerical and numerical strings. In particular, the abstract value of  $S$  in  $\overline{\mathbf{TJ}}^*$  is `NotNumeric`, and  $\llbracket \text{toNum}(\text{NotNumeric}) \rrbracket^{\overline{\mathbf{TJ}}^*}$  will precisely return 0.

Adopting a complete shell with respect to a certain operation does not compromise the precision of the others. For example, consider the domain of TAJs again and the following JavaScript fragment.

```

1  var obj = {
2      "foo" : 1,
3      "bar" : 2,
4      "1.2" : 3,
5      "2.2" : "hello"
6  }
7
8  y = obj[idx];

```

Suppose that the value of `idx` is the abstraction of the string set  $S = \{\text{"foo"}, \text{"bar"}\}$ , in the starting TAJs string abstract domain, namely the abstract value `NotUnsigned`. The variable `idx` is used to access the property of the object `obj` at line 8 and, to guarantee soundness, it accesses *all* the properties of `obj`, including the fields `"1.2"` and `"2.2"`, introducing noise in the abstract computation, since `"1.2"` and `"2.2"` are false positives values introduced by the abstraction of the values of `idx`. If we analyse the same JavaScript fragment with the absolute complete shell (with respect to `toNum` operation) of the TAJs string abstract domain defined in Section 5.4.1, we obtain more precise results. Indeed, in this case, the value of `idx` corresponds to the abstract value `NotNumeric`, and when it is used to access the object `obj` at line 8, only `"foo"` and `"bar"` are accessed, since they are the only non-numerical string properties of `obj`.

### 5.5.2 Qualitative Evaluation of Complete Shells

We evaluate the complete shells we have provided in the previous section from a qualitative point of view. As usual in static analysis, there is a trade-off between precision and



efficiency: choosing a more precise abstract domain may compromise the efficiency of the abstract computations. A representative example is in [80]: the complete shell of the Sign abstract domain with respect to addition is the interval abstract domain. Hence, starting from a finite height abstract domain (signs) we obtain an infinite height abstract domain (intervals). In particular, fix-point computations on signs converge, while on intervals it may diverge. Indeed, after the completion, the Interval abstract domain should also be equipped with a widening operator [51] to still guarantee termination. A worst-case scenario is when the complete shells with respect to a certain operation exactly corresponds to the collecting abstract domain, i.e., the concrete domain. Clearly, we cannot use the concrete domain due to undecidability reasons, so we change the starting abstract domain, since it cannot track any information related to the operation of interest. An example is the Suffix abstract domain [49] with `substring` operation: since this abstract domain tracks only the common suffix of a set of strings, it can not track the information about the indexes of the common suffix, and the complete shell of the suffix abstract domain with respect to `substring` would lead to the concrete domain. Hence, if the focus of the abstract interpreter is to improve the precision of the `substring` operation, we should change the abstract domain with a more precise one for `substring`, such as the finite state automata [16] abstract domain.

Consider now the complete shells reported in Section 5.4. The obtained complete shells still have finite height. Hence termination is still guaranteed without equipping the complete shells with widening operators. Moreover, the complexity of the string operations of interest is preserved after completion. Indeed, in both TAJIS and SAFE starting abstract domains, `concat` and `toNum` operations have constant complexity, respectively, and the same complexity is preserved in the corresponding complete shells. It is worth noting that also the complexity of the abstract domain-related operations, such as least upper bound, greatest lower bound and the ordering operator, is preserved in the complete shells. As far as the complete shells we have reported for TAJIS and SAFE are concerned, there is no worsening when we replace the original string abstract domains with the corresponding complete shells. As we have already mentioned before, this leads to completeness during the input abstraction process with respect to the relative operations, namely `concat` for SAFE and `toNum` for TAJIS.

### 5.5.3 False Positives Reduction

In static analysis, a certain degree of abstraction must be added in order to obtain decidable procedures to infer invariants on a generic program. Clearly, using less precise abstract domains lead to an increase of *false positive* values of the computed invariants. In particular, after a program is analysed, this burdens the phase of false positive detection: when a program is analysed, the following phase consists in detecting which values of the

invariants derived by the static analyser are spurious, namely those values that certainly are not computed by the concrete execution of the program of interest.

In particular, using imprecise (i.e., not complete) abstract domains clearly increases the number of false positives in the abstract computation of the static analyser, burdening the next phase of detecting the spurious values. Conversely, adopting (backward) complete abstract domains with respect to a certain operation reduces the number of false positives introduced during the abstract computations, at least in the input abstraction process. Clearly, in this way, the next phase of detecting false positives will be lighter since less noise has been introduced during the abstract computation of the invariants. Consider the JavaScript fragment of the previous paragraph again. As we already discussed before, using the starting TAJIS abstract domain to abstract the variable `idx` leads to a loss of precision, since the spurious value `"foo"` and `"bar"` are taken into account in its abstract value, namely `Unsigned`. Using the complete shell of TAJIS with respect to `toNum` instead does not add noise when `idx` is used to access `obj`.

## 5.6 Relative Precision

In this section, we recall the notion of pseudo-distance on abstract domains, firstly defined in [125], which takes into account the order relation between abstract elements (together with the fact that different abstract elements might approximate the same set of concrete values) and their possible incomparability. Then we formalise the increment of the precision obtained when analysing a program with a complete abstract domain with respect to the original version of the domain.

### 5.6.1 Abstract Domains Precision: an Overview

It is well known that abstract domains precision can be qualitatively compared by exploiting the information they can capture [44, 45]. However, quantitatively evaluating the precision of abstract domains has been proven to be quite challenging and frequently resulted in ad hoc measures.

Di Pierro and Wiklicky [144], introduced the notion of probabilistic Abstract Interpretation, used to numerically estimate the incompleteness of numerical abstract domains.

Sotin [160] presented the notion of precision of a numerical abstract value, measuring the volume it describes, to quantitatively compare the precision of numerical abstract domains. Finally, Casso et al. [32] compare the precision of different analyses on logic programs. They proposed distances in two abstract domains used in constraint logic programming, and they extended them to distances between the results of different analyses of a given program.

Logozzo et al. [125] gave a more general definition of pseudo-distance between abstract domain elements that allows to quantify the error of approximating a concrete element in an abstract domain.

**Definition 5.5** (Pseudo-distance [125]). Let  $\langle \overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}}, \perp_{\overline{\mathbf{D}}}, \top_{\overline{\mathbf{D}}}, \sqcap_{\overline{\mathbf{D}}}, \sqcup_{\overline{\mathbf{D}}} \rangle$  be an abstract domain. A function  $\delta : \overline{\mathbf{D}} \times \overline{\mathbf{D}} \rightarrow \mathbb{R} \cup \{+\infty\}$  is a *pseudo-distance* over the abstract domain  $\overline{\mathbf{D}}$  if and only if, for any  $\overline{\mathbf{d}}, \overline{\mathbf{d}}', \overline{\mathbf{d}}'' \in \overline{\mathbf{D}}$ , it satisfies the following axioms:

- *non-negativity*:  $\delta(\overline{\mathbf{d}}, \overline{\mathbf{d}}') \geq 0$
- *if-identity*:  $\overline{\mathbf{d}} =_{\overline{\mathbf{D}}} \overline{\mathbf{d}}' \Rightarrow \delta(\overline{\mathbf{d}}, \overline{\mathbf{d}}') = 0$
- *symmetry*:  $\delta(\overline{\mathbf{d}}, \overline{\mathbf{d}}') = \delta(\overline{\mathbf{d}}', \overline{\mathbf{d}})$
- *weak triangle inequality*:  $\overline{\mathbf{d}} \sqsubseteq_{\overline{\mathbf{D}}} \overline{\mathbf{d}}'' \sqsubseteq_{\overline{\mathbf{D}}} \overline{\mathbf{d}}' \Rightarrow \delta(\overline{\mathbf{d}}, \overline{\mathbf{d}}'') \leq \delta(\overline{\mathbf{d}}, \overline{\mathbf{d}}') + \delta(\overline{\mathbf{d}}', \overline{\mathbf{d}}'')$

△

We recall the path length distance definition.

**Definition 5.6** (Path length distance [125]). Let  $\langle \overline{\mathbf{D}}, \sqsubseteq_{\overline{\mathbf{D}}}, \perp_{\overline{\mathbf{D}}}, \top_{\overline{\mathbf{D}}}, \sqcap_{\overline{\mathbf{D}}}, \sqcup_{\overline{\mathbf{D}}} \rangle$  be an abstract domain. The path length distance  $\delta_{\text{plen}} : \overline{\mathbf{D}} \times \overline{\mathbf{D}} \rightarrow \mathbb{R} \cup \{+\infty\}$  is defined as

$$\delta_{\text{plen}}(\overline{\mathbf{d}}, \overline{\mathbf{d}}') = \begin{cases} \text{plen}(\overline{\mathbf{d}}, \overline{\mathbf{d}}') & \overline{\mathbf{d}} \sqsubseteq_{\overline{\mathbf{D}}} \overline{\mathbf{d}}' \\ \text{plen}(\overline{\mathbf{d}}', \overline{\mathbf{d}}) & \overline{\mathbf{d}}' \sqsubseteq_{\overline{\mathbf{D}}} \overline{\mathbf{d}} \\ +\infty & \text{otherwise} \end{cases}$$

where  $\text{plen}(\overline{\mathbf{d}}, \overline{\mathbf{d}}')$  computes the distance between two abstract values  $\overline{\mathbf{d}}, \overline{\mathbf{d}}' \in \overline{\mathbf{D}}$  and it is defined as follows.

$$\begin{aligned} \text{plen}(\overline{\mathbf{d}}, \overline{\mathbf{d}}') &= \min\{n \mid \{\overline{\mathbf{d}}_0, \overline{\mathbf{d}}_1, \dots, \overline{\mathbf{d}}_n\} \in \mathcal{P}(\overline{\mathbf{D}}), \overline{\mathbf{d}}_0 = \overline{\mathbf{d}}, \overline{\mathbf{d}}_n = \overline{\mathbf{d}}', \\ &\quad \forall i \in [0, n), \overline{\mathbf{d}}_i \sqsubseteq_{\overline{\mathbf{D}}} \overline{\mathbf{d}}_{i+1}\} \end{aligned}$$

△

### 5.6.2 Measuring Precision Gained by Complete Shells

We aim at computing the distance, and in turn, the increment of precision, between abstract points of an abstract domain  $\overline{\mathbf{D}}$  and its complete shell, with respect to a certain operation of interest  $f$ , noting that the abstract values of  $\overline{\mathbf{D}}$  are all contained in its complete shell. In order not to clutter the notation, given an abstract domain  $\overline{\mathbf{D}}$ , we denote by  $\text{Shell}_f(\overline{\mathbf{D}})$  its complete shell with respect to the function  $f$ .

Let  $P$  be a  $\mu\text{Dyn}$  program; let  $\text{Lab}(P)$  be the program points of  $P$ ; let  $\ell_i \in \text{Lab}(P)$  be  $i$ -th program point of  $P$  and let  $\text{Vars}(P)$  be the program variables. An example of  $\mu\text{Dyn}$

```

1 nums = "";
2 notnums = "";
3 i = 0;
4 while (i < length(s)) {
5   if (toNum(charAt(str, i)) == 0) {
6     notnums = concat(notnums, charAt(s, i));7
7   } else {
8     nums = concat(nums, charAt(s, i));9
9   }
10  i = i + 1;11
12 }12

```

FIGURE 5.8: Example of  $\mu$ Dyn annotated program

program annotated with its program points is in Figure 5.8. Given this program, as usual in static program analysis, the goal is to compute the abstract values associated to each variable  $x \in \text{Vars}(P)$  at each program point  $\ell \in \text{Lab}(P)$ .

We denote by  $\xi^{\overline{\mathbf{D}}} : \text{Vars}(P) \rightarrow \overline{\mathbf{D}}$  the abstract state associating each variable to an abstract value of  $\overline{\mathbf{D}}$ . When clear from the context, we denote the abstract state by  $\xi$ . Hence, given an abstract domain  $\overline{\mathbf{D}}$  and a program  $P$ , the result of the analysis of  $P$  using the abstract domain  $\overline{\mathbf{D}}$  (and the corresponding abstract semantics  $\llbracket \cdot \rrbracket^{\overline{\mathbf{D}}}$ ) is defined by

$$\text{Analysis}(P, \overline{\mathbf{D}}, \xi_\emptyset) = \{(\ell_i, \xi_i) \mid \ell_i \in \text{Lab}(P), \xi_i = \llbracket P \rrbracket^{\overline{\mathbf{D}}} \xi_\emptyset \text{ at program point } \ell_i\}$$

where  $\xi_\emptyset$  is the (initial) empty abstract state. For example, we analyse the  $\mu$ Dyn program in Figure 5.8 with the TAJ $\overline{\mathbf{J}}$  abstract domain of Figure 5.4b. The abstract state holding at program point 12 (namely the exit program point) is  $\xi_{12} = \{i \mapsto \text{UnsignedInt}, \text{nums} \mapsto \top_{\overline{\mathbf{TJ}}}, \text{notnums} \mapsto \top_{\overline{\mathbf{TJ}}}\}$ , indeed the pair  $(12, \xi_{12}) \in \text{Analysis}(P, \overline{\mathbf{TJ}}, \xi_\emptyset)$ .

We can exploit the path length distance of Definition 5.6 to define, for each program point  $\ell \in \text{Lab}(P)$ , the distance from  $\perp_{\overline{\mathbf{D}}}$  and each abstract value associated with each  $x \in \text{Vars}(P)$ .

**Definition 5.7** (Set of  $\perp$ -distances). Let  $P$  be a  $\mu$ Dyn program, let  $\overline{\mathbf{D}}$  be an abstract domain, and let  $\delta_{\text{plen}}$  be the path length distance of Definition 5.6. The set of  $\perp_{\overline{\mathbf{D}}}$ -distances for a given  $\ell \in \text{Lab}(P)$  is defined as follows.

$$\Phi_\ell(P, \overline{\mathbf{D}}) = \{(x, \delta(\perp_{\overline{\mathbf{D}}}, \xi(x))) \mid (\ell, \xi) \in \text{Analysis}(P, \overline{\mathbf{D}}, \xi_\emptyset), x \in \text{Vars}(P)\}$$

△

For example, consider again  $\xi_{12}$ , defined above, that is the abstract state holding at the program point 12 of the program in Figure 5.8. Hence, the set of the  $\perp_{\overline{\mathbf{TJ}}}$ -distances at program point 12 is  $\Phi_{12}(P, \overline{\mathbf{TJ}}) = \{(i, 2), (\text{nums}, 3), (\text{notnums}, 3)\}$ .

Note that, by the weak triangle inequality, given two abstract elements  $\bar{\mathbf{d}}, \bar{\mathbf{d}}' \in \bar{\mathbf{D}}$ , if  $\bar{\mathbf{d}} \sqsubseteq_{\bar{\mathbf{D}}} \bar{\mathbf{d}}' \Rightarrow \delta(\perp_{\bar{\mathbf{D}}}, \bar{\mathbf{d}}) \leq \delta(\perp_{\bar{\mathbf{D}}}, \bar{\mathbf{d}}')$ . Moreover, since  $\perp_{\bar{\mathbf{D}}}$  is comparable with any abstract value  $\bar{\mathbf{d}} \in \bar{\mathbf{D}}$ ,  $\delta(\perp_{\bar{\mathbf{D}}}, \bar{\mathbf{d}}) \neq \infty$ .

**Definition 5.8** (Precision entropy at a program point). Let  $P$  be a  $\mu$ Dyn program and let  $\bar{\mathbf{D}}$  be an abstract domain. Given a program point  $\ell \in \text{Lab}(P)$ , the precision entropy of  $\bar{\mathbf{D}}$  for the program  $P$  at  $\ell$  is defined as

$$\mathbb{P}_{\ell}(P, \bar{\mathbf{D}}) = \sum_{(\mathbf{x}_i, d_i) \in \Phi_{\ell}(P, \bar{\mathbf{D}})} d_i$$

where  $d_i$  is the path length distance from  $\perp_{\bar{\mathbf{D}}}$  to the abstract value of  $\mathbf{x}_i$  at the program point  $\ell$ . △

We can use Definition 5.8 to define the precision entropy of an analysis for a given program.

**Definition 5.9** (Precision entropy). Let  $P$  be a  $\mu$ Dyn program and let  $\bar{\mathbf{D}}$  be an abstract domain. The precision entropy  $\mathbb{P}(P, \bar{\mathbf{D}})$  of the abstract domain  $A$  for the program  $P$  is defined as

$$\mathbb{P}(P, \bar{\mathbf{D}}) = \sum_{\ell \in \text{Lab}(P)} \mathbb{P}_{\ell}(P, \bar{\mathbf{D}})$$

△

The precision entropy  $\mathbb{P}(P, \bar{\mathbf{D}})$  says how much information is expressed by the analysis based on the abstract domain  $\bar{\mathbf{D}}$ : the more  $\mathbb{P}(P, \bar{\mathbf{D}})$  is low, the more the analysis based on  $\bar{\mathbf{D}}$  is precise for the program  $P$ . Using this metric, we can compare the analysis results based on a certain abstract domain and its corresponding complete shell. As we have already mentioned at the beginning of Section 5.6, any abstract point of  $\bar{\mathbf{D}}$  is contained in its complete shell. Moreover, it is worth noting that, given a variable  $\mathbf{x} \in \text{Vars}(P)$  and a certain program point  $\ell \in \text{Lab}(P)$ , the abstract value computed by the analysis on  $\bar{\mathbf{D}}$  associated with  $\mathbf{x}$  is always comparable with the abstract value computed by the analysis on  $\text{Shell}_f(\bar{\mathbf{D}})$  associated with  $\mathbf{x}$ . Formally, consider  $(\ell, \xi) \in \text{Analysis}(P, \bar{\mathbf{D}}, \xi_{\emptyset})$  and  $(\ell, \xi') \in \text{Analysis}(P, \text{Shell}_f(\bar{\mathbf{D}}), \xi_{\emptyset})$ , for some program point  $\ell \in \text{Lab}(P)$ , we have that  $\forall \mathbf{x} \in \text{Vars}(P). \xi'(\mathbf{x}) \sqsubseteq_{\text{Shell}_f(\bar{\mathbf{D}})} \xi(\mathbf{x})$ . This can be expressed by the following predicate

$$\mathbb{P}(P, \text{Shell}_f(\bar{\mathbf{D}})) \leq \mathbb{P}(P, \bar{\mathbf{D}}).$$

Informally speaking, for a program point  $\ell$ , the analysis result of  $\mathbf{x}$  on  $\text{Shell}_f(\bar{\mathbf{D}})$  is always dominated (i.e., less than or equal to) by the analysis result of  $\mathbf{x}$  on  $\bar{\mathbf{D}}$  (contained in  $\text{Shell}_f(\bar{\mathbf{D}})$ ). This fact directly comes from the dominance relation involved in the definition

Variable	$\overline{\mathbf{TJ}}$	$\overline{\mathbf{TJ}}^*$	$\mathbb{P}_{12}(\mathbf{P}, \overline{\mathbf{TJ}}) - \mathbb{P}_{12}(\mathbf{P}, \overline{\mathbf{TJ}}^*)$
<code>s</code>	UnsignedStr	UnsignedStr	0
<code>i</code>	UnsignedInt	UnsignedInt	0
<code>nums</code>	$\top_{\overline{\mathbf{TJ}}}$	UnsignedOrNotNumeric	1
<code>notnums</code>	$\top_{\overline{\mathbf{TJ}}}$	UnsignedOrNotNumeric	1

TABLE 5.2: Analysis output results with  $\overline{\mathbf{TJ}}$  and  $\overline{\mathbf{TJ}}^*$ .

of complete shells given in Theorems 5.1 and 5.2. For this reason, we can always compare the analysis results produced by  $\overline{\mathbf{D}}$  and  $\text{Shell}_f(\overline{\mathbf{D}})$  for any variable and any program point.

This leads us to an automatic procedure to compute how much the analysis performed by  $\text{Shell}_f(\overline{\mathbf{D}})$  is better than the one performed by  $\overline{\mathbf{D}}$ .

### 5.6.3 Experimental Evaluation

The CLAM static analyzer for  $\mu\text{Dyn}$  programs implements the TAJIS and SAFE string abstract domains and their corresponding complete shells discussed in Section 5.4.<sup>7</sup> The abstract interpreter is parametric, as it can be tuned by selecting the string abstract domain to be used to analyse a given program. Other string abstract domains can be easily plugged into our static analyser, without re-implementing the underlying abstract interpreter. Moreover, it is possible to check the precision entropy, of Section 5.6, of an abstract domain and its complete shell. In this way, it is possible to check at which program point and for which variables the complete shell-based analysis gains precision with respect to the analysis on the original abstract domain.

For example, consider the  $\mu\text{Dyn}$  program in Figure 5.8, where the value of `s` is statically unknown. The program takes the strings `s` and puts its non-zero numerical characters into `nums` and the others into `notnums`. If the variable `s` is initialised as "24kobe8", at the program point 12 the value of `nums` is "248" and the value of `notnums` is "kobe". Consider TAJIS and its complete shell and let analyse the program with both the abstract domains. When the variable `s` is initialised with the abstract value `UnsignedStr`, the analysis output, for the exit point, is in Table 5.2, where the second and third columns correspond to the result analysis of  $\overline{\mathbf{TJ}}$  and  $\overline{\mathbf{TJ}}^*$  for the corresponding variable, respectively, and the last column is the increment precision gained by performing the analysis on  $\overline{\mathbf{TJ}}^*$ . Observe that, for the variables `s` and `i`, we have no precision improvement, as stated by the last column of Table 5.2. Concerning variables `nums` and `notnums`, the analysis on  $\overline{\mathbf{TJ}}$  returns the top values, whereas the analysis on  $\overline{\mathbf{TJ}}^*$  leads to an precision increment, since it returns the abstract value `UnsignedOrNotNumeric`, for both variables. The precision increment is in the

<sup>7</sup>Available at <https://github.com/VincenzoArceri/clam>

last column of Table 5.2 for the variables `nums` and `notnums`, since their abstract values (`UnsignedOrNotNumeric`) on the analysis on  $\overline{\mathbf{TJ}}^*$  is distant 1 from the abstract values of the analysis on  $\overline{\mathbf{TJ}}$  ( $\top_{\overline{\mathbf{TJ}}}$ ). The total precision entropy of the analysis on  $\overline{\mathbf{TJ}}$  and the analysis on  $\overline{\mathbf{TJ}}^*$  is 54 and 44, respectively, so the overall precision increment when considering the complete shell, in this case, is 10.

## 5.7 Discussion

In this chapter, we focused on backward completeness in JavaScript-specific string abstract domains, and in particular, we provided the complete shells of TAJIS and SAFE string abstract domains with respect to `toNum` and `concat` operations, together with an effective procedure to measure the precision improvement of the analysis when moving to the complete shell. Our results can also be easily applied to JSAI string abstract domain [112], as it can be seen as an extension of the SAFE domain.





## Chapter 6

---

# RELATIONAL STRING ABSTRACT DOMAINS

In this chapter, we define an Abstract Interpretation framework for relational string analysis. The framework is based on string orders and allows us to detect relations between the values of distinct string variables, in contrast to typical string abstract domains, that describe the content of each string variable independently. We instantiate the framework to some string order relations of interest. Note that the abstract domains we introduce in this chapter are not language-specific, and consequently, they may be suitable for analysing programs written in different programming languages. Moreover, we track relations between string variables and expressions and, based on the string order onto which the framework has been instantiated, we can also infer content and shape information of string values. Finally, our framework can be easily combined with basics string abstract domains and lead to preciser analyses.

The contribution of this chapter is submitted for publication.

### Chapter Structure

Section 6.1 highlights the importance of having relational string analysis and explains our contribution. Section 6.2 presents our core language. Section 6.3 formalizes the framework from which relational string abstract domains based on a given strings order can be built and defines some instances of it. Section 6.4 presents the experimental evaluation. Section 6.5 concludes.

## 6.1 Introduction

The static analysis community has spent a great effort in proposing new abstractions, to better approximate and analyse string values. Unfortunately, almost all the existing string abstract domains track information of single variables used in a program (e.g., if a string contains some characters, or it starts with a given sequence), without inspecting their relationships with other values (e.g., if a string is a substring of another one, despite their actual values are unknown).

```

String secName(String javaName, String pr1, String pr2) {
    if (javaName.startsWith(pr1)) {
        return pr2 + javaName.substring(4);
    } else if (javaName.startsWith(pr2)) {
        return pr1 + javaName.substring(4);
    } else {
        return javaName;
    }
}

```

FIGURE 6.1: SECNAME sample.

Concerning integer abstractions, advanced and sophisticated relational abstractions have been studied and improved over the years to track relations between integer variables: a representative example, among others, is the Polyhedra abstract domain [58] proposed in the late 70s, but constantly and heavily improved over the years, as reported by the more recent important works on its optimization, e.g., [21, 20]. Unfortunately, we cannot say the same for string abstract domains, where the majority of them infer only *non-relational* information.

We illustrate the problem by considering the Java function SECNAME in Figure 6.1.<sup>1</sup> SECNAME takes as input three arguments, `javaName`, `pr1` and `pr2`. Then, if `javaName` starts with `pr1`, the function returns `pr2` concatenated to the substring, from the index 4, of `javaName`. SECNAME behaves analogously when `javaName` starts with `pr2` and it concatenates `pr1` with `javaName.substring(4)`, otherwise `javaName` is returned. The relational information we want to capture here is the one relating `pr1` and `pr2` with `javaName` and the returned value. In particular, we want to infer that (i) `javaName.substring(4)` is always part of the returned value, (ii) `pr1` is returned as part of the return value if `javaName` starts with `pr2`, and (iii) `pr2` is returned as part of the returned value if `javaName` starts with `pr1`. Unfortunately, in order to catch such information, non-relational abstractions, even the more advanced and sophisticated ones, are not enough and new relational domains for strings are a necessity.

### State of the Art

Several relational abstract domains have been proposed for numerical values, such as Polyhedra [58], Octagons [135], Logahedra [98], Pentagons [126], Stripes [66], and Weighted Hexagons [72]. Some of these domains have inspired us to build the string relational domains that will be presented in Section 6.3. Indeed, consider the Octagons and the Pentagons abstract domains. Octagons track relations of the form  $\pm x \pm y \leq k$ , where

<sup>1</sup>SECNAME is the result of a slight modification made to the function available at <https://www.codota.com/code/java/classes/java.lang.String>

$k$  is a constant. Instead, Pentagons, a less precise domain than Octagons, combines the numerical properties tracked by the Interval domain (i.e.,  $x \in [n, m]$ ) and the symbolic ones captured by the Strict Upper Bound domain (i.e.,  $x < y$ ). Similar to the Strict Upper Bound domain, our framework instantiates domains that track information of the form  $x \preceq y$  where  $\preceq$  is a (general) order over string variables. Moreover, the framework extension we define to track relations between string expressions and variables, like  $x + y \preceq z$ , has been modelled similarly to Octagons.

Other abstractions have been proposed to infer information about the relations between heap-allocated data structures a program manipulates [179]. Authors defined in [88] an abstract domain that approximates must and may equalities among pointer expressions. A relational abstract domain for shape analysis has been presented in [102], built on the top of a set of logical connectives, that represents relations among memory states.

Also, for string values, a big effort has been made to improve the precision of string abstractions, but contrary to the numerical world, most of them only focus on the abstraction of a single variable (cf. Chapter 1).

### Contribution

In this chapter, we define a general framework upon which relational strings abstract domains can be instantiated, starting from a string order of interest. We introduce a suite of relational abstract domains fitting the proposed framework, based on length inequality, char inclusion, substring relations. Precisely, our framework will be firstly formalized to track relations between single string variables, and later on extended to infer relations between string expressions and variables, to appreciate the improvement of the captured information precision.

We conclude by presenting the experimental analysis results of the relational string abstract domain based on our substring relation, and we compare it with the state-of-the-art general purpose string abstract domains.

## 6.2 Core Language

As core language we use IMP, a modified version of  $\mu\text{Dyn}$  (cf. Section 5.3).

### 6.2.1 Syntax

The IMP syntax is given in Figure 6.2. Let  $P$  be an IMP program. Each IMP statement is annotated with a label  $\ell \in \text{Lab}(P)$  (not belonging to the syntax), where  $\text{Lab}(P)$  denotes the set of the  $P$  labels, i.e., its program points. We denote by  $\ell_{\perp}$  the initial program point and by  $\ell_{\dashv}$  the exit program point.

As usual in static analysis, a program can be analysed by looking at its control-flow graph (CFG for short), i.e., a directed graph that embeds the control structure of a program,

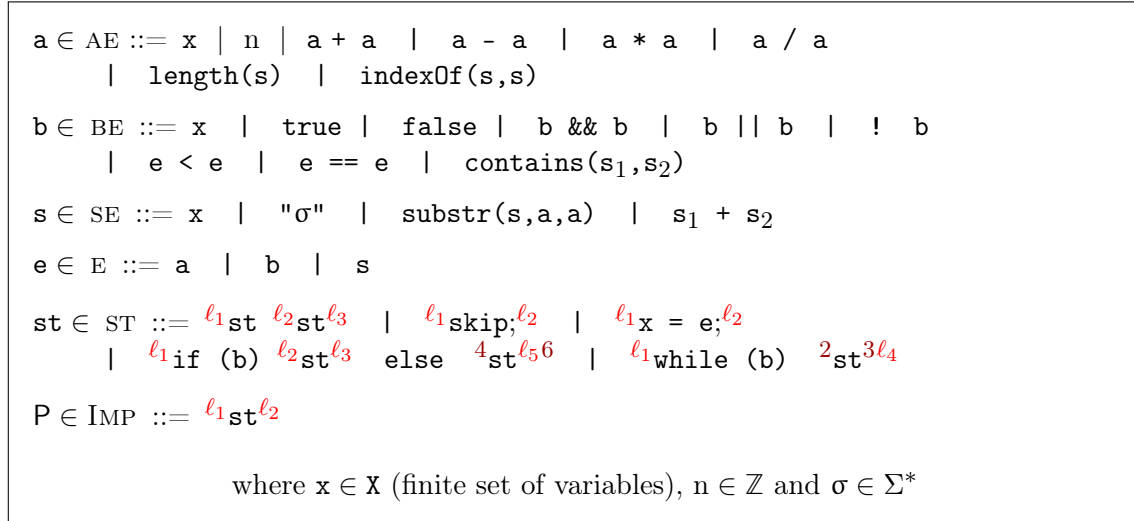


FIGURE 6.2: IMP syntax.



FIGURE 6.3: Example of CFG generation.

where the nodes are the program points and the edges express the flow paths from the entry to the exit block. Formally<sup>2</sup>, given a program  $P \in \text{IMP}$ , we define the corresponding CFG  $G_P = \langle \text{Nodes}_P, \text{Edges}_P, \text{In}_P, \text{Out}_P \rangle$  as the CFG whose nodes are the program points, i.e.,  $\text{Nodes}_P = \text{Lab}(P)$ ; the input node is the entry program point, i.e.,  $\text{In}_P = \ell_+$ ; and the output node is the last program point, i.e.,  $\text{Out}_P = \ell_-$ . The algorithm computing the CFG of a program  $P$  is standard and can be found, e.g., in [99, 14]. An example of CFG is depicted in Figure 6.3, where the label **true** is omitted.

As observed in [99], the language of the edges label is different from IMP, since the CFG embeds the control structure of the program. In particular, the grammar of the edges label is  $\text{IMP}^{\text{CFG}} ::= \text{skip} \mid x = \text{exp} \mid b$ .

## 6.2.2 Concrete Semantics

Since we aim to analyse CFGs to answer questions on the corresponding programs, we define the semantics of the edges label, elements of  $\text{IMP}^{\text{CFG}}$ , which express the effect that each edge has from its entry node to its exit node. Let  $\text{VAL} = \text{INT} \cup \text{BOOL} \cup \text{STR}$  be the set of the possible values associated with a variable, where  $\text{INT}$  denotes the set of signed integers

<sup>2</sup>Both static analysis process and notation follow the one presented in [99].

$\mathbb{Z}$ , **BOOL** denotes the set of booleans  $\{\mathbf{true}, \mathbf{false}\}$ , and **STR** denotes the set of strings  $\Sigma^*$  over an alphabet  $\Sigma$ . Let  $\mathbf{m} \in \mathbb{M} = X \rightarrow \mathbf{VAL}$  be the set of memories, where  $\mathbf{m}_\emptyset$  is the empty memory. The semantics of expressions is captured by the function  $\llbracket e \rrbracket : \mathbb{M} \rightarrow \mathbf{VAL}$ . Since the semantics of integer and Boolean expressions is standard (and not of interest of this chapter), in the following we only give the concrete semantics of string expressions.

$$\begin{aligned}
\llbracket s_1 + s_2 \rrbracket \mathbf{m} &= \llbracket s_1 \rrbracket \mathbf{m} \cdot \llbracket s_2 \rrbracket \mathbf{m} & \llbracket x \rrbracket \mathbf{m} &= \mathbf{m}(x) & \llbracket \sigma \rrbracket \mathbf{m} &= \sigma \\
\llbracket \mathbf{substr}(s, a_1, a_2) \rrbracket \mathbf{m} &= \sigma_1 \dots \sigma_{j-1} \\
&\quad \text{where } \sigma = \llbracket s \rrbracket \mathbf{m}, i = \llbracket a_1 \rrbracket \mathbf{m}, j = \llbracket a_2 \rrbracket \mathbf{m}, 0 \leq i \leq j < |\sigma| \\
\llbracket \mathbf{indexOf}(s_1, s_2) \rrbracket \mathbf{m} &= \begin{cases} \mathbf{idx}(\sigma_1, \sigma_2) & \text{if } \sigma_2 \curvearrowright \sigma_1 \\ -1 & \text{otherwise} \end{cases} \\
&\quad \text{where } \sigma_1 = \llbracket s_1 \rrbracket \mathbf{m}, \sigma_2 = \llbracket s_2 \rrbracket \mathbf{m} \\
\llbracket \mathbf{length}(s) \rrbracket \mathbf{m} &= |\llbracket s \rrbracket \mathbf{m}| \\
\llbracket \mathbf{contains}(s_1, s_2) \rrbracket \mathbf{m} &= \llbracket s_2 \rrbracket \mathbf{m} \curvearrowright \llbracket s_1 \rrbracket \mathbf{m}
\end{aligned}$$

Once we defined the expression concrete semantics, we formalize the edges label semantics. Abusing notation, we define the function  $\llbracket \mathbf{st} \rrbracket : \mathbb{M} \rightarrow \mathbb{M}$  to capture the semantics of the elements of  $\mathbf{IMP}^{\mathbf{CFG}}$ .

$$\begin{aligned}
\llbracket \mathbf{skip} \rrbracket \mathbf{m} &= \mathbf{m} & \llbracket x = e \rrbracket \mathbf{m} &= \mathbf{m}[x \leftarrow \llbracket e \rrbracket \mathbf{m}] \\
\llbracket \mathbf{b} \rrbracket \mathbf{m} &= \begin{cases} \mathbf{m} & \text{if } \llbracket \mathbf{b} \rrbracket \mathbf{m} = \mathbf{true} \\ \mathbf{m}_\emptyset & \text{if } \llbracket \mathbf{b} \rrbracket \mathbf{m} = \mathbf{false} \end{cases}
\end{aligned}$$

Finally, the store [14] is the collection of memories for each program point. The store is defined as  $\mathfrak{s} \in \mathbb{S} = \mathbf{Lab}(\mathbf{P}) \rightarrow \mathbb{M}$  and it associates a memory to each program point.

Static analysis computes (abstract) invariants for each program point; thus, we first define a collecting semantics which associates with each program point, namely to each node of a CFG, the set of its possible memories holding at that program point. This boils down to lifting the concrete semantics  $\llbracket \mathbf{st} \rrbracket : \mathbb{M} \rightarrow \mathbb{M}$  (working on single memories), to the collecting semantics  $\llbracket \mathbf{st} \rrbracket : \mathcal{P}(\mathbb{M}) \rightarrow \mathcal{P}(\mathbb{M})$  working on sets of memories. We also slightly change the notion of store, since it has to map each program point to sets of memories, namely  $\bar{\mathfrak{s}} \in \bar{\mathbb{S}} = \mathbf{Lab}(\mathbf{P}) \rightarrow \mathcal{P}(\mathbb{M})$ .

Finally, we could use standard fix-point analysis algorithms (i.e., the ones presented in [14, 99, 139]) which returns a store  $\bar{\mathfrak{s}}$  such that, for each  $\ell \in \mathbf{Lab}(\mathbf{P})$ ,  $\bar{\mathfrak{s}}(\ell)$  is the fix-point collecting semantics (i.e., a set of memories) for the program point  $\ell$ . In this way, unfortunately, the invariants for each node of a CFG are not computable because of Rice's Theorem. Therefore, we need abstraction to make static analysis decidable.

### 6.3 A Suite of String Relational Abstract Domains

In this section, we provide a suite of relational string abstract domains, based on several well-known orders over strings. We start with a general framework to build string relational abstract domains starting from a given string order. Within this framework, we present three different string relational abstract domains: length inequality, character inclusion, and substring domains, with the corresponding abstract semantics of IMP.

#### 6.3.1 General Relational Framework

We aim at capturing relations between string variables in the form  $x \preceq y$ , with respect to some (partial or pre-order) relation  $\preceq$  over strings, such as, for example, “ $y$  is a substring of  $x$ ” or “ $y$  is the prefix of  $x$ ”. As introduced in Section 6.1, in the numerical world such a relation is captured by the (strict) upper bound abstract domain [135, 126], which expresses relations in the form  $y \leq x$ , for some numerical variables  $x$  and  $y$ . In this section, we generalize the upper bound abstract domain to string variables, making it parametric with respect to a given string order.

Our starting point is an order  $\sim \subseteq \Sigma^* \times \Sigma^*$  between strings, from which, given the finite set of string variables  $X_{\text{str}}$ , we build a new order  $\preceq \subseteq X_{\text{str}} \times X_{\text{str}}$  between a pair of string variables. Upon the order  $\preceq$  we design a generalized relational string abstract domain.

**Definition 6.1** (General string relational abstract domain). Let  $\preceq \subseteq X_{\text{str}} \times X_{\text{str}}$  be an order over string variables. The general string relational abstract domain  $A$ , ranged over the meta-variable  $\mathcal{A}$ , is defined as

$$A = \mathcal{P}(\{y \preceq x \mid x, y \in X_{\text{str}}\}) \cup \{\perp_A\}$$

such that the top element  $\top_A$  corresponds to  $\emptyset$  and the bottom element is  $\perp_A$ . The least upper bound, greatest lower bound, and the partial order of  $A$  are defined as follows.

$$\begin{aligned} \mathcal{A}_1 \sqcap_A \mathcal{A}_2 &= \{y \preceq x \mid y \preceq x \in \mathcal{A}_1 \vee y \preceq x \in \mathcal{A}_2\} \\ \mathcal{A}_1 \sqcup_A \mathcal{A}_2 &= \text{Clos}(\{y \preceq x \mid y \preceq x \in \mathcal{A}_1 \wedge y \preceq x \in \mathcal{A}_2\}) \\ \mathcal{A}_1 \sqsubseteq_A \mathcal{A}_2 &\iff \mathcal{A}_1 \supseteq \mathcal{A}_2 \end{aligned}$$

where  $\text{Clos} : A \rightarrow A$  is the function performing the transitive closure of the abstract element input. △

The abstract domain  $A$  is intended to collect  $\preceq$  *must* relations, i.e., informally speaking, if a relation  $y \preceq x$  is captured in the abstract, it means that it holds in the concrete.

Note that elements of  $A$  (cf. Definition 6.1) are sets of relations  $y \preceq x$  between single string variables. Moreover, the general abstract domain  $A$  is finite, given that the set of string variables is finite and, in turn, also the number of possible relations. Thus, it is trivial to see that the domain  $(A, \sqsubseteq_A, \perp_A, \top_A, \sqcap_A, \sqcup_A)$  is a complete lattice, and that its greatest lower bound  $\sqcap_A$  and least upper bound  $\sqcup_A$  are defined as the union and the intersection between abstract elements, respectively. Concretization and abstraction functions  $\gamma_A : A \rightarrow \mathcal{P}(\mathbb{M})$  and  $\alpha_A : \mathcal{P}(\mathbb{M}) \rightarrow A$  are generally defined as

$$\gamma_A(\mathcal{A}) = \bigcap_{y \preceq x \in \mathcal{A}} \{\mathbf{m} \mid \mathbf{m}(x), \mathbf{m}(y) \in \Sigma^*, \mathbf{m}(y) \sim \mathbf{m}(x)\} \quad (6.1)$$

$$\alpha_A(\mathbb{M}) = \{y \preceq x \mid \forall \mathbf{m} \in \mathbb{M}. \mathbf{m}(y) \sim \mathbf{m}(x), x, y \in X_{\text{str}}\} \quad (6.2)$$

Also,  $(\mathcal{P}(\mathbb{M}), \alpha_A, \gamma_A A)$  is a Galois Connection, since  $\mathcal{P}(\mathbb{M})$  and  $A$  are complete lattices and  $\alpha_A$  is a join-morphism.

At this point, we define a general and parametric abstract semantics of IMP. In particular, given an abstract domain  $A$ , built upon the order  $\preceq$ , we define the function  $\llbracket \text{st} \rrbracket^A : A \rightarrow A$  capturing the  $\preceq$ -relations between variables generated by the statement  $\text{st}$ . We start by defining the parametric abstract semantics of the assignment  $x = \mathbf{s}$ . Here, the crucial point is the definition of the auxiliary function  $\text{extr} : \text{SE} \rightarrow \mathcal{P}(X_{\text{str}})$  that, given a string expression  $\mathbf{s}$ , extracts all the variables appearing in  $\mathbf{s}$  that are related with respect to  $\preceq$  with  $x$ . In particular,  $\text{extr}$  approximates the set of variables appearing in  $\mathbf{s}$  that are  $\preceq$ -related with  $x$ .

Once defined the extraction function  $\text{extr}$ , we give the general abstract semantics of the assignment, which is defined by the following steps:

- **[remove]**

$$\mathcal{A}_r = \begin{cases} \mathcal{A} \setminus \{w \preceq z \mid w = x\} & \text{if } x \text{ appears at the top-level of } \mathbf{s} \\ \mathcal{A} \setminus \{w \preceq z \mid w = x \vee z = x\} & \text{otherwise} \end{cases}$$

- **[add]**

$$\mathcal{A}_a = \mathcal{A}_r \cup \{y \preceq x \mid y \in \text{extr}(\mathbf{s})\}$$

- **[closure]**

$$\llbracket x = \mathbf{s} \rrbracket^A \mathcal{A} = \text{Clos}(\mathcal{A}_a)$$

The first phase is **[remove]**: given the input memory  $\mathcal{A}$ , it removes the relations that surely do not hold anymore after the assignment execution. In particular, we always remove the relations of the form  $x \preceq z$ , for some  $z \in X_{\text{str}}$ , since  $x$  is going to be overwritten, but we also remove any relations of the form  $w \preceq x$ , for some  $w \in X_{\text{str}}$ , if and only if  $x$  does not appear at the top-level of the expression  $\mathbf{s}$ .

Then, the **[add]** phase adds the  $\preceq$ -relations  $y \preceq x$ , for each variable  $y$  collected in  $\text{extr}(\mathbf{s})$ . Finally, **[closure]** performs the transitive closure of the abstract memory obtained from **[add]**, i.e.,  $\mathcal{A}_a$ , by means of the function **Clos**, in order to derive the implicit  $\preceq$ -relations which are not directly present in  $\mathcal{A}_a$ .

As far as Boolean expressions are concerned, the only IMP Boolean expressions that generate  $\preceq$ -relations (for the string domains presented in this chapter) are: **contains**( $\mathbf{s}_1, \mathbf{s}_2$ ),  $\mathbf{s}_1 == \mathbf{s}_2$ , conjunctive, and disjunctive expressions. Note that given the expression **contains**( $\mathbf{s}_1, \mathbf{s}_2$ ), we can infer  $\preceq$ -relations only when  $\mathbf{s}_1$  is a variable. In all the other cases, no information is gathered.

$$\llbracket \text{contains}(x, \mathbf{s}) \rrbracket^A \mathcal{A} = \text{Clos}(\mathcal{A} \cup \{y \preceq x \mid y \in \text{extr}(\mathbf{s})\})$$

Similarly, we can infer  $\preceq$ -relations in the abstract semantics of  $\mathbf{s}_1 == \mathbf{s}_2$  only when either  $\mathbf{s}_1$  or  $\mathbf{s}_2$  is a string variable.

$$\llbracket x == \mathbf{s} \rrbracket^A \mathcal{A} = \llbracket \mathbf{s} == x \rrbracket^A \mathcal{A} = \begin{cases} \text{Clos}(\mathcal{A} \cup \{y \preceq x, x \preceq y\}) & \text{if } \mathbf{s} = y \in X_{\text{str}} \\ \text{Clos}(\mathcal{A} \cup \{y \preceq x \mid y \in \text{extr}(\mathbf{s})\}) & \text{otherwise} \end{cases}$$

As far as the conjunctive and disjunctive expressions semantics, we rely on the least upper bound and greatest lower bound operators given in Definition. 6.1.

$$\begin{aligned} \llbracket e_1 \ \&\& \ e_2 \rrbracket^A \mathcal{A} &= \llbracket e_1 \rrbracket^A \mathcal{A} \sqcup_A \llbracket e_2 \rrbracket^A \mathcal{A} \\ \llbracket e_1 \ \|\ e_2 \rrbracket^A \mathcal{A} &= \llbracket e_1 \rrbracket^A \mathcal{A} \sqcap_A \llbracket e_2 \rrbracket^A \mathcal{A} \end{aligned}$$

Unlike the assignment case, the abstract semantics of Boolean expressions removes no previous substring relations, since these expressions do not affect the (concrete) memory.

Moreover, this general abstract semantics only holds for the abstract domains that we present in this chapter, and it may not hold for other abstract domains generated within the framework proposed in this section.

### 6.3.2 String Length Relational Abstract Domain

Within the formal framework presented above, we generate several relational string abstract domains. For instance, one may be interested in capturing the relations concerning the length of a string variable with respect to another, when interacting during the program execution. Formally, we are interested in identifying the relation  $\preceq_{\text{len}} \subseteq X_{\text{str}} \times X_{\text{str}}$  between string variables such that, given  $x, y \in X_{\text{str}}$ ,  $y \preceq_{\text{len}} x$  if and only if the length of  $y$  is smaller than or equal to the length of  $x$ .<sup>3</sup> We instantiate the general abstract domain of

---

<sup>3</sup> $\preceq_{\text{len}}$  is not a partial order but it is a pre-order binary relation. We will discuss later on how this affects the abstract domain gathered information.



$$\begin{aligned}
\text{Greatest lower bound: } \mathcal{L}_1 \sqcap_{\text{len}} \mathcal{L}_2 &= \{y \preceq_{\text{len}} x \mid y \preceq_{\text{len}} x \in \mathcal{L}_1 \vee y \preceq_{\text{len}} x \in \mathcal{L}_2\} \\
\text{Least upper bound: } \mathcal{L}_1 \sqcup_{\text{len}} \mathcal{L}_2 &= \text{Clos}(\{y \preceq_{\text{len}} x \mid y \preceq_{\text{len}} x \in \mathcal{L}_1 \wedge y \preceq_{\text{len}} x \in \mathcal{L}_2\}) \\
\text{Partial order: } \mathcal{L}_1 \sqsubseteq_{\text{len}} \mathcal{L}_2 &\iff \mathcal{L}_1 \supseteq \mathcal{L}_2
\end{aligned}$$

FIGURE 6.4: Lattice operations over  $\text{Len}$ .

Definition 6.1 over the order  $\preceq_{\text{len}}$ . In particular, we replace the general string order  $\preceq$  with  $\preceq_{\text{len}}$ , obtaining the relational string length abstract domain  $\text{Len} = \mathcal{P}(\{y \preceq_{\text{len}} x \mid x, y \in X_{\text{str}}\}) \cup \{\perp_{\text{len}}\}$ , ranged over the meta-variable  $\mathcal{L}$ , where the top element, denoted by  $\top_{\text{len}}$ , is  $\emptyset$ , and  $\perp_{\text{len}}$  is the bottom element. The least upper bound  $\sqcup_{\text{len}}$  and the greatest lower bound  $\sqcap_{\text{len}}$  operators and the partial order  $\sqsubseteq_{\text{len}}$  are in Figure 6.4. It is straightforward to prove that  $(\text{Len}, \sqsubseteq_{\text{len}}, \sqcup_{\text{len}}, \sqcap_{\text{len}}, \perp_{\text{len}}, \top_{\text{len}})$  is a complete lattice. Similarly, we define the concretization  $\gamma_{\text{len}} : \text{Len} \rightarrow \mathcal{P}(\mathbb{M})$  and the abstraction  $\alpha_{\text{len}} : \mathcal{P}(\mathbb{M}) \rightarrow \text{Len}$  functions of the relational string length abstract domain instantiating Equation 6.2 and 6.1 over  $\preceq_{\text{len}}$ :

$$\gamma_{\text{len}}(\mathcal{L}) = \bigcap_{y \preceq_{\text{len}} x \in \mathcal{L}} \{\mathbf{m} \mid \mathbf{m}(x), \mathbf{m}(y) \in \Sigma^*, |\mathbf{m}(y)| \leq |\mathbf{m}(x)|\}$$

$$\alpha_{\text{len}}(\mathbb{M}) = \{y \preceq_{\text{len}} x \mid \forall \mathbf{m} \in \mathbb{M}. |\mathbf{m}(y)| \leq |\mathbf{m}(x)|, x, y \in X_{\text{str}}\}$$

$(\mathcal{P}(\mathbb{M}), \alpha_{\text{len}}, \gamma_{\text{len}}, \text{Len})$  is a Galois Connection.

**Extraction function of  $\text{Len}$ .** Given  $x = \mathbf{s}$ , we can see the string expression  $\mathbf{s}$  as an ordered list of concatenated expressions  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_n$ , and the string variables that surely have length less than or equal to  $x$  are the ones at the top-level of a concatenation appearing in  $\mathbf{s}$ . For instance, consider the assignment  $x = y + z + w$ . The substring relations we aim to capture from it are  $y \preceq_{\text{len}} x$ ,  $z \preceq_{\text{len}} x$  and  $w \preceq_{\text{len}} x$ , that is  $y, z, w$  have length less than or equal to the length of  $x$ . These variables are collected by the extraction function  $\text{extr} : \text{SE} \rightarrow \mathcal{P}(X_{\text{str}})$  defined as:

$$\text{extr}(\mathbf{s}) = \begin{cases} \{y\} & \text{if } \mathbf{s} = y \in X_{\text{str}} \\ \text{extr}(\mathbf{s}_1) \cup \text{extr}(\mathbf{s}_2) & \text{if } \mathbf{s} = \mathbf{s}_1 + \mathbf{s}_2 \\ \emptyset & \text{otherwise} \end{cases}$$

The abstract semantics of  $\text{Len}$  is straightforward and reported in Appendix C.1.

### 6.3.3 Character Inclusion Relational Abstract Domain

In the following, we present the character inclusion relational abstract domain  $\text{Char}$ , tracking the characters included between a pair of string variables.

$$\begin{aligned}
\text{Greatest lower bound: } \mathcal{C}_1 \sqcap_{\text{char}} \mathcal{C}_2 &= \{y \preceq_{\text{char}} x \mid y \preceq_{\text{char}} x \in \mathcal{C}_1 \vee y \preceq_{\text{char}} x \in \mathcal{C}_2\} \\
\text{Least upper bound: } \mathcal{C}_1 \sqcup_{\text{char}} \mathcal{C}_2 &= \text{Clos}(\{y \preceq_{\text{char}} x \mid y \preceq_{\text{char}} x \in \mathcal{C}_1 \wedge y \preceq_{\text{char}} x \in \mathcal{C}_2\}) \\
\text{Partial order: } \mathcal{C}_1 \sqsubseteq_{\text{char}} \mathcal{C}_2 &\iff \mathcal{C}_1 \supseteq \mathcal{C}_2
\end{aligned}$$

FIGURE 6.5: Lattice operations over Char.

Given  $x, y \in X_{\text{str}}$ , we are also interested in capturing “if the characters which appear in  $y$  occur in  $x$  as well”. Formally, we introduce the binary relation  $\preceq_{\text{char}} \subseteq X_{\text{str}} \times X_{\text{str}}$  such that  $y \preceq_{\text{char}} x$  if and only if the set of characters of  $y$  is contained or equal to the set of characters of  $x$ . As  $\preceq_{\text{len}}, \preceq_{\text{char}}$  is not a partial order, but only a pre-order. We define the relational character inclusion string abstract domain  $\text{Char} = \mathcal{P}(\{y \preceq_{\text{char}} x \mid x, y \in X_{\text{str}}\}) \cup \{\perp_{\text{char}}\}$ , ranging over the meta-variable  $\mathcal{C}$ , where the top element, denoted by  $\top_{\text{char}}$ , is  $\emptyset$ , and  $\perp_{\text{char}}$  is the bottom element. Note that  $(\text{Char}, \sqsubseteq_{\text{char}}, \sqcup_{\text{char}}, \sqcap_{\text{char}}, \perp_{\text{char}}, \top_{\text{char}})$  is also a complete lattice and its lattice operators are presented in Figure 6.5.

The concretization  $\gamma_{\text{char}} : \text{Char} \rightarrow \mathcal{P}(M)$  and the abstraction  $\alpha_{\text{char}} : \mathcal{P}(M) \rightarrow \text{Char}$  functions are formally defined below, where  $\text{char}(\cdot)$  computes the set of characters of a string. As in the previous case,  $(\mathcal{P}(M), \alpha_{\text{char}}, \gamma_{\text{char}}, \text{Char})$  is a Galois Connection.

$$\begin{aligned}
\gamma_{\text{char}}(\mathcal{C}) &= \bigcap_{y \preceq_{\text{char}} x \in \mathcal{C}} \{m \mid m(x), m(y) \in \Sigma^*, \text{char}(m(y)) \subseteq \text{char}(m(x))\} \\
\alpha_{\text{char}}(M) &= \{y \preceq_{\text{char}} x \mid \forall m \in M. \text{char}(m(y)) \subseteq \text{char}(m(x)), x, y \in X_{\text{str}}\}
\end{aligned}$$

The abstract semantics of Char relies on the extraction function defined for Len and can be found in Appendix C.2.

### 6.3.4 Substring Relational Abstract Domain

The abstract domains Len and Char track relations about the lengths and the characters of a pair of string variables. One main limit of these domains is that they are both based on pre-orders: hence, when we have relations of the form  $x \preceq y$  and  $y \preceq x$ , we cannot assert that the variables  $x$  and  $y$  are equal. Moreover, Len loses any information about the contents of a string variable, while Char loses any information about the shape of a string variable. We propose then a strictly more precise partial order-based relational string abstract domain, still fitting in the formal framework presented in Section 6.3.1 and solving the problems of Len and Char.

Given  $x, y \in X_{\text{str}}$ , let the binary relation  $\preceq_{\text{sub}} : X_{\text{str}} \times X_{\text{str}}$  be such that  $x \preceq_{\text{sub}} y$  if and only if  $x$  is a substring of  $y$ . The relation  $\preceq_{\text{sub}}$  is a partial order, being reflexive, transitive and anti-symmetric. Thus, we define the relational string abstract domain  $\text{Sub} = \mathcal{P}(\{y \preceq_{\text{sub}} x \mid x, y \in X_{\text{str}}\}) \cup \{\perp_{\text{sub}}\}$ , ranged over the meta-variable  $\mathcal{S}$ , where the top element, denoted

$$\begin{aligned}
\text{Greatest lower bound: } \mathcal{S}_1 \sqcap_{\text{sub}} \mathcal{S}_2 &= \{x \preceq_{\text{sub}} y \mid x \preceq_{\text{sub}} y \in \mathcal{S}_1 \vee x \preceq_{\text{sub}} y \in \mathcal{S}_2\} \\
\text{Least upper bound: } \mathcal{S}_1 \sqcup_{\text{sub}} \mathcal{S}_2 &= \text{Clos}(\{x \preceq_{\text{sub}} y \mid x \preceq_{\text{sub}} y \in \mathcal{S}_1 \wedge x \preceq_{\text{sub}} y \in \mathcal{S}_2\}) \\
\text{Partial order: } \mathcal{S}_1 \sqsubseteq_{\text{sub}} \mathcal{S}_2 &\iff \mathcal{S}_1 \supseteq \mathcal{S}_2
\end{aligned}$$

FIGURE 6.6: Lattice operations over **Sub**.

by  $\top_{\text{sub}}$ , corresponds to  $\emptyset$ , and  $\perp_{\text{sub}}$  is the bottom element. **Sub** is a complete lattice, i.e.,  $(\mathbf{Sub}, \sqsubseteq_{\text{sub}}, \sqcup_{\text{sub}}, \sqcap_{\text{sub}}, \perp_{\text{sub}}, \top_{\text{sub}})$  was operator are given in Figure 6.6.

The concretization  $\gamma_{\text{sub}} : \mathbf{Sub} \rightarrow \mathcal{P}(\mathbb{M})$  and the abstraction  $\alpha_{\text{sub}} : \mathcal{P}(\mathbb{M}) \rightarrow \mathbf{Sub}$  functions follow the definitions Eq. 6.1 and 6.2 reported in Section 6.3.1 and they are defined as:

$$\begin{aligned}
\gamma_{\text{sub}}(\mathcal{S}) &= \bigcap_{y \preceq_{\text{sub}} x \in \mathcal{S}} \{m \mid m(x), m(y) \in \Sigma^*, m(y) \curvearrowright m(x)\} \\
\alpha_{\text{sub}}(\mathbb{M}) &= \{y \preceq_{\text{sub}} x \mid \forall m \in \mathbb{M}. m(y) \curvearrowright m(x), x, y \in X_{\text{str}}\}
\end{aligned}$$

$(\mathcal{P}(\mathbb{M}), \alpha_{\text{sub}}, \gamma_{\text{sub}}, \mathbf{Sub})$  is a Galois Connection.

Again, the abstract semantics of **Sub** relies on the extraction function defined for **Len** and can be found in Appendix C.3.

### 6.3.5 Extension to String Expressions

The abstract domains previously proposed are able to track only  $\preceq$ -relations between single string variables. Consider the fragment below:

$$x = y + y + w; \quad z = y + w;$$

If we analyse it with the substring abstract domain, the final abstract memory is  $\{y \preceq_{\text{sub}} x, w \preceq_{\text{sub}} x, y \preceq_{\text{sub}} z, w \preceq_{\text{sub}} z\}$ , but other substring relations may be inferred, such as  $z \preceq_{\text{sub}} x$  or  $y + w \preceq_{\text{sub}} x$ . In the following, we extend the framework introduced in Section 6.3.1 to catch also these relations.

We introduce the binary relation  $\preceq_* \subseteq \text{SE} \times X_{\text{str}}$  that relates string expressions and string variables. Note that also SE is limited to the expressions encountered in a program. Upon the binary relation  $\preceq_*$  we build the new set of abstract memories  $A^* = \mathcal{P}(\{\mathbf{s} \preceq_* x \mid \mathbf{s} \in \text{SE}, x \in X_{\text{str}}\}) \cup \{\perp_{A^*}\}$ , ranged over the meta-variable  $\mathcal{A}^*$  where the top element, denoted by  $\top_{A^*}$ , is  $\emptyset$ , and  $\perp_{A^*}$  is the bottom element. Note that,  $A^*$  is still a finite domain, since, given a program  $P \in \text{IMP}$ , both its string variables and string expressions are finite. Thus,  $(A^*, \sqsubseteq_{A^*}, \perp_{A^*}, \top_{A^*}, \sqcap_{A^*}, \sqcup_{A^*})$  is a complete lattice, whose operators and partial order are reported in Figure 6.7. The concretization  $\gamma_{A^*} : A^* \rightarrow \mathcal{P}(\mathbb{M})$  and the abstraction

$$\begin{aligned}
\text{Greatest lower bound: } \mathcal{A}_1^* \sqcap_{A^*} \mathcal{A}_2^* &= \{\mathbf{s} \preceq_* \mathbf{x} \mid \mathbf{s} \preceq_* \mathbf{x} \in \mathcal{A}_1^* \vee \mathbf{s} \preceq_* \mathbf{x} \in \mathcal{A}_2^*\} \\
\text{Least upper bound: } \mathcal{A}_1^* \sqcup_{A^*} \mathcal{A}_2^* &= \text{Clos}(\{\mathbf{s} \preceq_* \mathbf{x} \mid \mathbf{s} \preceq_* \mathbf{x} \in \mathcal{A}_1^* \wedge \mathbf{s} \preceq_* \mathbf{x} \in \mathcal{A}_2^*\}) \\
\text{Partial order: } \mathcal{A}_1^* \sqsubseteq_{A^*} \mathcal{A}_2^* &\iff \mathcal{A}_1^* \supseteq \mathcal{A}_2^*
\end{aligned}$$

FIGURE 6.7: Lattice operations over  $A^*$ .

$\alpha_{A^*} : \mathcal{P}(\mathbb{M}) \rightarrow A^*$  functions form a Galois Connection and are defined as follows:

$$\gamma_{A^*}(\mathcal{A}^*) = \bigcap_{\mathbf{s} \preceq_* \mathbf{x} \in \mathcal{A}^*} \{\mathbf{m} \mid \llbracket \mathbf{s} \rrbracket \mathbf{m}, \mathbf{m}(x) \in \Sigma^*, \llbracket \mathbf{s} \rrbracket \mathbf{m} \sim \mathbf{m}(x)\}$$

$$\alpha_{A^*}(\mathbb{M}) = \{\mathbf{s} \preceq_* \mathbf{x} \mid \forall \mathbf{m} \in \mathbb{M}. \llbracket \mathbf{s} \rrbracket \mathbf{m} \sim \mathbf{m}(x), x \in X_{\text{str}}, \mathbf{s} \in \text{SE}\}$$

At this point, we define the abstract semantics of  $A^*$ . Given an assignment  $\mathbf{x} = \mathbf{s}$ , the abstract semantics of  $A^*$  follows the one defined for  $A$ : from the expression  $\mathbf{s}$  we retrieve any sub-expression that is surely in a  $\preceq_*$ -relation with  $\mathbf{x}$ , using the auxiliary function  $\text{extr}^*$ . The function  $\text{extr}^* : \text{SE} \rightarrow \mathcal{P}(\text{SE})$  extends the previously defined function  $\text{extr}$ , extracting from a string expression  $\mathbf{s}$ , its top-level sub-expressions, in addition to its top-level variables (as  $\text{extr}$ ). For instance,  $\text{extr}^*(\mathbf{y} + \mathbf{w} + \text{"ab"}) = \{\mathbf{y}, \mathbf{w}, \mathbf{w} + \text{"ab"}, \mathbf{y} + \mathbf{w}, \mathbf{y} + \mathbf{w} + \text{"ab"}, \text{"a"}, \text{"b"}\}$ . Hence, the abstract semantics of assignment is defined by the following steps:

- **[remove]**

$$\mathcal{A}_r^* = \begin{cases} \mathcal{A}^* \setminus \{\mathbf{s}' \preceq_* \mathbf{z} \mid \mathbf{x} \text{ appears in } \mathbf{s}'\} & \text{if } \mathbf{x} \text{ appears} \\ & \text{at the top-level of } \mathbf{s} \\ \mathcal{A}^* \setminus \{\mathbf{s}' \preceq_* \mathbf{z} \mid \mathbf{z} = \mathbf{x} \vee \mathbf{x} \text{ appears in } \mathbf{s}'\} & \text{otherwise} \end{cases}$$

- **[add]**

$$\mathcal{A}_a^* = \mathcal{A}_r^* \cup \{\mathbf{s}' \preceq_* \mathbf{x} \mid \mathbf{s}' \in \text{extr}^*(\mathbf{s})\}$$

- **[inter-asg]**

$$\mathcal{A}_i^* = \mathcal{A}_a^* \cup \{\mathbf{x} \preceq_* \mathbf{y} \mid \mathcal{A}_a^*(\mathbf{x}) \subseteq \mathcal{A}_a^*(\mathbf{y})\}$$

- **[closure]**

$$\llbracket \mathbf{x} = \mathbf{s} \rrbracket^{A^*} \mathcal{A}^* = \text{Clos}(\mathcal{A}_i^*)$$

The **[remove]**, **[add]** and **[closure]** phases are similar to those of the definition of  $\llbracket \cdot \rrbracket^A$ . The intermediate phase **[inter-asg]** instead differs from the previous definition. Given  $\mathbf{z} \in X_{\text{str}}$  and a generic memory  $\mathcal{A}^* \in A^*$ , we denote by  $\mathcal{A}^*(\mathbf{z})$  the set of relations of  $\mathbf{z}$ , i.e.,  $\mathcal{A}^*(\mathbf{z}) = \{\mathbf{s}' \preceq_* \mathbf{z} \mid \mathbf{s}' \preceq_* \mathbf{z} \in \mathcal{A}^*\}$ . If from the previous steps, any relation of  $\mathbf{x}$  is

$$\begin{aligned}
\text{Greatest lower bound: } \mathcal{S}_1^* \sqcap_{\text{sub}^*} \mathcal{S}_2^* &= \{\mathbf{s} \preceq_{\text{sub}^*} \mathbf{x} \mid \mathbf{s} \preceq_{\text{sub}^*} \mathbf{x} \in \mathcal{S}_1^* \vee \mathbf{s} \preceq_{\text{sub}^*} \mathbf{x} \in \mathcal{S}_2^*\} \\
\text{Least upper bound: } \mathcal{S}_1^* \sqcup_{\text{sub}^*} \mathcal{S}_2^* &= \text{Clos}(\{\mathbf{s} \preceq_{\text{sub}^*} \mathbf{x} \mid \mathbf{s} \preceq_{\text{sub}^*} \mathbf{x} \in \mathcal{S}_1^* \wedge \mathbf{s} \preceq_{\text{sub}^*} \mathbf{x} \in \mathcal{S}_2^*\}) \\
\text{Partial order: } \mathcal{S}_1^* \sqsubseteq_{\text{sub}^*} \mathcal{S}_2^* &\iff \mathcal{S}_1^* \supseteq \mathcal{S}_2^*
\end{aligned}$$

FIGURE 6.8: Lattice operations over  $\text{Sub}^*$ .

contained in the ones of a string variable  $y$ , as checked in the **[inter-*asg*]** phase, we can safely assert that  $x$  and  $y$  are related, and we can add that relation to  $\mathcal{A}_a^*$ .

The following abstract semantics of Boolean expressions is analogous to  $\llbracket \cdot \rrbracket^A$ .

$$\begin{aligned}
\llbracket \text{contains}(x, \mathbf{s}) \rrbracket^{A^*} \mathcal{A}^* &= \text{Clos}(\mathcal{A}^* \cup \{\mathbf{s}' \preceq_* \mathbf{x} \mid \mathbf{s}' \in \text{extr}^*(\mathbf{s})\}) \\
\llbracket x == \mathbf{s} \rrbracket^{A^*} \mathcal{A}^* &= \llbracket \mathbf{s} == x \rrbracket^{A^*} \mathcal{A}^* = \text{Clos}(\mathcal{A}^* \cup \{\mathbf{s}' \preceq_* \mathbf{x} \mid \mathbf{s}' \in \text{extr}^*(\mathbf{s})\}) \\
\llbracket e_1 \ \&\& \ e_2 \rrbracket^{A^*} \mathcal{A}^* &= \llbracket e_1 \rrbracket^{A^*} \mathcal{A}^* \sqcup_{A^*} \llbracket e_2 \rrbracket^{A^*} \mathcal{A}^* \\
\llbracket e_1 \ || \ e_2 \rrbracket^{A^*} \mathcal{A}^* &= \llbracket e_1 \rrbracket^{A^*} \mathcal{A}^* \sqcap_{A^*} \llbracket e_2 \rrbracket^{A^*} \mathcal{A}^*
\end{aligned}$$

### Extended Substring Relational Abstract Domain

We extend the substring relational abstract domain following  $A^*$ . The extended version of the string length and character inclusion relational domains are in Appendix C.5 and C.6, respectively.

Let  $\preceq_{\text{sub}^*} \subseteq \text{SE} \times X_{\text{str}}$  be the binary relation such that  $\mathbf{s} \preceq_{\text{sub}^*} \mathbf{x}$  if and only if  $\mathbf{s}$  is a substring of  $\mathbf{x}$ . For instance,  $y+y \preceq_{\text{sub}^*} \mathbf{x}$  means that  $y+y$  is a substring of  $\mathbf{x}$ . Moreover, let  $\text{Sub}^* = \mathcal{P}(\{\mathbf{s} \preceq_{\text{sub}^*} \mathbf{x} \mid \mathbf{s} \in \text{SE}, \mathbf{x} \in X_{\text{str}}\}) \cup \{\perp_{\text{sub}^*}\}$  be the extended substring relational abstract domain, ranged over the meta-variable  $\mathcal{S}^*$ , where the top element, denoted by  $\top_{\text{sub}^*}$ , is  $\emptyset$ , and  $\perp_{\text{sub}^*}$  is the bottom element. The least upper bound and the greatest lower bound operators and partial order are  $\sqcup_{\text{sub}^*}$ ,  $\sqcap_{\text{sub}^*}$  and  $\sqsubseteq_{\text{sub}^*}$ , respectively, and they are in Figure 6.8. Similarly to the previous cases, also  $(\text{Sub}^*, \sqsubseteq_{\text{sub}^*}, \sqcup_{\text{sub}^*}, \sqcap_{\text{sub}^*}, \perp_{\text{sub}^*}, \top_{\text{sub}^*})$  is a complete lattice, and the concretization  $\gamma_{\text{sub}^*} : \text{Sub}^* \rightarrow \mathcal{P}(M)$  and the abstraction  $\alpha_{\text{sub}^*} : \mathcal{P}(M) \rightarrow \text{Sub}^*$  functions are defined as follows, forming again a Galois Connection.

$$\begin{aligned}
\gamma_{\text{sub}^*}(\mathcal{S}^*) &= \bigcap_{\mathbf{s} \preceq_{\text{sub}^*} \mathbf{x} \in \mathcal{S}^*} \{\mathbf{m} \mid \llbracket \mathbf{s} \rrbracket \mathbf{m}, \mathbf{m}(x) \in \Sigma^*, \llbracket \mathbf{s} \rrbracket \mathbf{m} \curvearrowright \mathbf{m}(x)\} \\
\alpha_{\text{sub}^*}(M) &= \{\mathbf{s} \preceq_{\text{sub}^*} \mathbf{x} \mid \forall \mathbf{m} \in M. \llbracket \mathbf{s} \rrbracket \mathbf{m} \curvearrowright \mathbf{m}(x), x \in X_{\text{str}}, \mathbf{s} \in \text{SE}\}
\end{aligned}$$

The abstract semantics of  $\text{Sub}$  is straightforward and is in Appendix C.4.

<pre> r = "not constant"; if ((*\$*)) r = r + " passed"; else r = r + " failed"; assert (r.contains("con")); assert (r.contains("ed")); </pre>	<pre> v = readStr(); r = "Elem: \n" + v + "\n"; while ((*\$*)) { v = readStr(); r = r + v + "\n"; } assert (r.contains("em")); assert (r.contains("\n")); assert (r.contains(v)); assert (r.contains("f")); </pre>
Program NOTCON	Program REP

FIGURE 6.9: Program samples used for domain comparison.

## 6.4 Experimental Evaluation

We have implemented  $\overline{\mathbf{RSub}}$ , a prototypical static analyser for the toy language presented in Section 6.2 using the new relational string abstract domain  $\mathbf{Sub}^*$ . Indeed it subsumes, from the precision point of view, the other abstract domains and their extensions, inferring the set of substring relations holding at each program point.  $\overline{\mathbf{RSub}}$  has been implemented and tested over several representative string manipulating fragments,<sup>4</sup> both taken from real-world software and hand-crafted. A subset of these fragments is used and discussed in this section to show limits and strengths of the aforementioned relational string domain and the corresponding analysis. In Section 6.4.1, we compare our analysis with prefix  $\overline{\mathbf{PR}}$ , suffix  $\overline{\mathbf{SU}}$ , char inclusion  $\overline{\mathbf{CI}}$  and bricks  $\overline{\mathbf{BR}}$  abstract domains [49] and with the finite state automata  $\overline{\mathbf{FA}}$  abstract domain [17] on some simple yet representative code fragments,. In Section 6.4.2 we discuss  $\overline{\mathbf{RSub}}$  on a real world example, comparing our analysis with the finite state automata abstract domain.  $\overline{\mathbf{RSub}}$  implements trace partitioning [149], which will be only used in the second part of the evaluation.

### 6.4.1 Test Cases

We consider the two code fragments manipulating strings, NOTCON and REP, in Figure 6.9. NOTCON may update the variable  $r$  either appending the string "passed" or "failed". REP iteratively appends to the variable  $r$  a read string stored in  $v$  (hence, statically unknown) concatenated with the string "\n". In both the examples, the shortcut  $?$  in the `while` and `if` statements denotes a statically unknown Boolean guard.

Table 6.1 shows the analyses results at the end of the execution of the program NOTCON. In particular, the second column reports the abstract value of  $r$  at the end of the analysis

---

<sup>4</sup>Available at <https://github.com/UniVE-SSV/rsubs/tree/main/src/test/resources/vmcai>

Domain	r abstract value	Asserts
$\overline{\text{PR}}$	not constant	$\times$
$\overline{\text{SU}}$	ed	$\times$
$\overline{\text{CI}}$	{n, o, t, c, n, s, a, e, d}, {n, o, t, c, n, s, a, p, s, e, d, f, i, l}	$\times$
$\overline{\text{BR}}$	{not constant passed, not constant failed}(1,1)	$\checkmark$
$\overline{\text{FA}}$	not constant(passed    failed)	$\checkmark$
$\overline{\text{RSub}}$	{not constant $\preceq_{\text{sub}} r$ , ed $\preceq_{\text{sub}} r$ }	$\checkmark$

TABLE 6.1: Analyses results of NOTCON.

Domain	r abstract value	Asserts
$\overline{\text{PR}}$	Elem:	$\times$
$\overline{\text{SU}}$	$\epsilon$ (unknown)	$\times$
$\overline{\text{CI}}$	{ $\circ$ }, { $\circ$ } (unknown)	$\times$
$\overline{\text{BR}}$	{ $\star$ }(0, $+\infty$ ) (unknown)	$\times$
$\overline{\text{FA}}$	Elem: $\star \backslash n (\star \backslash n)^*$	$\times$
$\overline{\text{RSub}}$	{Elem: $\preceq_{\text{sub}} r$ , $v \preceq_{\text{sub}} r$ , $r + v + \backslash n \preceq_{\text{sub}} r$ , $v + \backslash n \preceq_{\text{sub}} r$ , $\backslash n \preceq_{\text{sub}} r$ }	$\checkmark$

TABLE 6.2: Analyses results of REP.

and the third column is labeled with  $\checkmark$  if the corresponding analysis proves that all the `assert` conditions of NOTCON (lines 6 and 7) hold, or with  $\times$  otherwise. The analyses based on  $\overline{\text{PR}}$ ,  $\overline{\text{SU}}$  and  $\overline{\text{CI}}$  are not able to precisely verify all the assertions.  $\overline{\text{PR}}$  approximates the value of `r` with the common prefix, hence it asserts line 6 but not line 7. Similarly,  $\overline{\text{SU}}$  approximates the value of `r` with the common suffix, hence it asserts line 7 but not line 6.  $\overline{\text{CI}}$  neither asserts line 6 nor line 7, since the assertions check the containment of two strings (namely "con" and "ed") but  $\overline{\text{CI}}$  only tracks the *single* characters that the variable `r` may or must contain. Instead,  $\overline{\text{BR}}$  and  $\overline{\text{FA}}$  precisely track the possible string values of the variable `r` and successfully verify the assertions.  $\overline{\text{RSub}}$  also successfully verifies all the assertions, since it tracks that the strings "con" and "ed" are definitely substrings of `r`. It is worth noting that also other substring relations of `r` are tracked by  $\overline{\text{RSub}}$  (e.g., any substring of the string "not constant") but they are not interesting for the assertions that must be verified, and they have been omitted from Table 6.1 for space limitations.

Now consider REP, which differs from NOTCON since a fix-point computation is involved and deals with statically unknown strings. In particular, four assertions must be checked: the ones at lines 7-9 certainly holds, the one at line 10 *may* hold since the value of `v` is statically unknown. The analyses results are in Table 6.2, where in the second column "unknown" means that the corresponding analysis has returned the top abstract value, and the symbol  $\star$  denotes "any possible string". The third column is labelled with  $\checkmark$  if

the corresponding analysis proves all the `assert` conditions at lines 7-9 of REP or with  $\mathbf{x}$  otherwise.

$\overline{\mathbf{PR}}$  can verify the assertion at line 7 but not the ones at lines 8-9, since it loses any information on the rest of the string, except for the common prefix. In addition, it does not track the fact that the variable `v` is certainly contained in `r`.  $\overline{\mathbf{SU}}$ ,  $\overline{\mathbf{CI}}$  and  $\overline{\mathbf{BR}}$  analyses lose any information about the value of `r`, abstracting it as the corresponding top abstract value, making them unable to verify the assertions at lines 7-9.  $\overline{\mathbf{FA}}$  abstracts the value of `r` as the regular expression of Table 6.2 hence it correctly verifies the assertions at lines 7-8 but not the one at line 9, being unable to track the relation between the variables `r` and `v`. Even if statically unknown strings are manipulated by REP,  $\overline{\mathbf{RSub}}$  behaves as  $\overline{\mathbf{FA}}$  as far as assertions at lines 7-8 are concerned, since the strings `Elem:` and `\n` are definitely substrings of `r`. Also, it verifies the assertion at line 9, since it tracks that the variable `v` is a substring of the variable `r`, independently from its abstract value. The analysis results for  $\overline{\mathbf{RSub}}$  are in the last row of Table 6.2.

### 6.4.2 Evaluating a Real World Sample

We conclude this section evaluating and discussing  $\overline{\mathbf{RSub}}$  on `SECNAME`, shown in Figure 6.1 in the introduction, compared with the string analysis based on  $\overline{\mathbf{FA}}$ .

Table 6.3 has the analyses results of `SECNAME`, where we denote by `res` the abstract value returned by `SECNAME`. The second column contains the values obtained without trace partitioning, while the third column has the results obtained applying trace partitioning. We focus now on the first case. Since the function parameters are statically unknown, the Boolean guards at lines 2-4 are statically unknown, as well. Consequently,  $\overline{\mathbf{FA}}$  needs to take into account the values returned at lines 3, 5, and 7: unfortunately, being the parameter statically unknown, each line returns the top abstract value, losing any information about the value returned by `SECNAME`, as shown by the second line of Table 6.3. The same happens for  $\overline{\mathbf{RSub}}$ : each abstract value returned by each `return` statement (lines 3, 5 and 7) share no substring relation with the others, hence, lubbing the results,  $\overline{\mathbf{RSub}}$  approximates the returned values of `SECNAME` as the top abstract value.

Now, consider the same analysis with trace partitioning, whose results are in the third column of Table 6.3; each `[i]` denotes a trace and the corresponding abstract value returned at line `i`. As far as  $\overline{\mathbf{FA}}$  is concerned, the situation remains unaffected with respect to the analysis without trace partitioning. In particular, each returned abstract value still corresponds to the top abstract value.  $\overline{\mathbf{RSub}}$  precisely infers the information in the third column of Table 6.3, describing the possible substring relations at each `return` hotspot. In particular, three possible cases are detected: (i) when `pr2`  $\preceq_{\text{sub}}$  `javaName` the results definitely contains `pr1` and `javaName.substring(4)`, (ii) when `pr1`  $\preceq_{\text{sub}}$  `javaName` the results definitely contains `pr2` and `javaName.substring(4)`, (iii) in other cases, the results simply contains `javaName`.



Domain	res abstract value	res abstract value (with trace partitioning)
$\overline{\mathbf{FA}}$	$\top$	$[2] \top [5] \top [7] \top$
$\overline{\mathbf{RSub}}$	$\top$	$[2] \left\{ \begin{array}{l} \text{pr1} \preceq_{\text{sub}} \text{res} \\ \text{javaName.substring}(4) \preceq_{\text{sub}} \text{res} \end{array} \right\}$ where $\{\text{pr2} \preceq_{\text{sub}} \text{javaName}\}$ $[5] \left\{ \begin{array}{l} \text{pr2} \preceq_{\text{sub}} \text{res} \\ \text{javaName.substring}(4) \preceq_{\text{sub}} \text{res} \end{array} \right\}$ where $\{\text{pr1} \preceq_{\text{sub}} \text{javaName}\}$ $[7] \left\{ \text{javaName} \preceq_{\text{sub}} \text{res} \right\}$

TABLE 6.3: Analyses results of SECNAME.

Even if  $\overline{\mathbf{FA}}$  provides a precise string analysis, it is limited when dealing with statically unknown inputs, as highlighted in the previous example. Contrary, even if  $\overline{\mathbf{RSub}}$  does not perform a *collection* of the possible string values of a variable, tracking only the *definite* substrings of a variable, it treats variables and expressions symbolically, making itself able to infer useful information, such as in the SECNAME analysis. For instance, we can surely assert that the string `javaName.substring(4)` is a substring of the result and that the result surely includes part of the input (`javaName`).

## 6.5 Discussion

In this chapter, we have introduced a general framework to generate new relational abstract domains starting from pre- and partial orders on string values. In particular, we have introduced a new relational substring domain,  $\mathbf{Sub}^*$ , showing its impact on the accuracy of the analysis with respect to state-of-art general-purpose string abstract domains. It is worth noting that in this chapter,  $\mathbf{Sub}^*$  has been used as a *standalone* abstract domain. Of course, further accuracy can be achieved by combining  $\mathbf{Sub}^*$  with existing non-relational domains through the well-known reduced product operator. In fact, even the simple reduced product between the constant string abstract domain and  $\mathbf{Sub}^*$  would lead to more accurate analysis results, as it would be able to retrieve the exact value of some sub-expressions.



## Chapter 7

---

# STRING MANIPULATION IN WATERMARKING SCENARIOS

The last chapter of the thesis is somehow orthogonal with respect to the previous ones, as it moves to the usage of string manipulation in an applicative scenario. We investigate how a semantics-based approach to string manipulation may be effective when combined with watermarking techniques for ownership protection.

The content of this chapter reports contributions published in [85].

### Context

Information security systems protect digital assets from piracy and prying eyes. As shown in [5], those systems can be distinguished into two categories: cryptography and information hiding. Cryptography aims to hide secret information, transforming a given text into a cypher one (encryption phase), with the possibility of reversing the transformation, so restoring the original information (decryption phase). As its name suggests, information hiding looks for embedding unnoticeable data in the cover media, e.g., images, audio, video, and texts. Based on the reason why some information is hidden in some particular cover media, we can distinguish two techniques: steganography and watermarking. Steganography hides a secret message into a cover media in order to transmit it. On the contrary, watermarking hides data in cover media to protect their intellectual property from digital piracy (see [103]).

Thus, digital watermarking aims to prove ownership of digital cover data by persistently embedding, visibly or not, a watermark in the latter digital assets.<sup>1</sup> The embedding has to guarantee usability of the watermarked data, leaving them intact and recognisable. Then, a detection and extraction process of the watermark may follow, to claim the intellectual property [53, 158].

---

<sup>1</sup>Precisely, the watermark is the set of all the marks introduced into a cover data by a watermarking technique [3].

Briefly, watermarking techniques were firstly proposed to protect multimedia data [3],<sup>2</sup> exploiting the information redundancy they contain, and the limits of the human visual and auditory systems, to embed the watermark [104]. Watermarking has been subsequently extended to digital texts, due to the increased demand of digital publishing [104]. Also software have been affected by unauthorised copies and illegal reuses, then software watermarking techniques have been developed [61]. In particular, software watermarking techniques, according to the nature of the watermark extraction process, are distinguished in [53]: static, dynamic, and abstract. In [146], Dalla Preda and Pasqua defined a software watermarking system based on program semantics and abstract interpretation that defines the previous techniques as its instances.

Piracy also affected relational data, and thus watermarking database techniques have been developed. Firstly meant for numerical data, relational database watermarking techniques are based on the assumption that it is tolerable that watermarked databases contain a small amount of errors [28], despite the possibility of compromising the answers of SQL queries based on numerical conditions.

## Chapter Structure

Section 7.1 introduces the watermarking methodology for relational data to protect their integrity and intellectual property. Moreover, it highlights the problem of dealing with the semantic perturbation which may rise by watermarking relational data and explains our contribution. Section 7.2 defines the motivating examples that let clear the need of the elements we introduce in our technique. Section 7.3 gives an overview about semantic similarity theory and text watermarking techniques. Section 7.4 introduces our watermarking approach, emphasizing the definition of the elements related to the preservation of the semantic consistency, which constitutes a special feature of our work. Details related to the implementation are also given. Section 7.5 presents the validation of our approach through an experimental evaluation. Section 7.6 concludes.

## 7.1 Introduction

Relational databases are enterprise software systems, introduced in the 1970s [37], where data are stored and eventually analysed to find hidden relations among them which are useful to take strategic decisions. Data stored in relational databases can be pirated, illegally redistributed, tempered, and their ownership claimed, as happens to all the other digital assets (images, audio, video and texts). Indeed, the internet growth resulted in a multitude of web-based services through which data are continuously transmitted and easily accessible.

---

<sup>2</sup>Multimedia data refer to images, audio, and video.

Symbol	Description
SK	Secret key only known by the data owner.
PK	Primary key identifying the tuple.
$\eta$	Number of tuples in a relation R being watermarked.
$\gamma$	Fraction of tuples being watermarked $\gamma \in [1, \eta]$ . Also known as TF.
$\nu$	Number of attributes considered for the embedding.
$\xi$	Range of less significant bits ( <i>lsb</i> ) to embed the mark.
$\omega$	Number of watermarked tuples after the embedding.

TABLE 7.1: AHK approach notation.

How can we protect the integrity and the intellectual property of data? Relational database watermarking [4] has been proposed as a new field of security, to protect data property value. Its main element is the watermark, i.e., a stream of binary elements called marks, and it consists of two phases: watermark embedding and watermark extraction [3]. Generally speaking, we can classify relational database watermarking techniques into two categories [130]: those embedding the watermark in the data, causing distortions, e.g., [106, 111, 189], and those that do not cause them, e.g., [28, 29, 89]. The distortion-based techniques are mostly oriented to resist aggressive attacks as they are conceived to protect data from false ownership claims [92]. Moreover, relational databases watermarking techniques may be distinguished by: the type of watermark information to be embedded, the type of the attribute into which marks are embedded, their granularity level (bit level, character level, attribute level, tuple level), the detection and verification process in use (deterministic or probabilistic, blindly or non-blindly, publicly or privately), and their ultimate purpose [93]

The first watermarking technique for relational data was proposed by [4]. Also known as the AHK algorithm, this distortion-based approach defines the notation of the elements belonging to a relation R (see Table 7.1), embeds the marks in numerical attributes, introduces a solid criteria for the selection of the places for the watermark embedding, and it has been used by many authors [33, 65, 165] as starting point for other watermarking proposals.

Relational database watermarking techniques deal with different issues [92]. Among them, three are of particular interest: the watermark capacity (i.e., the optimal amount of marks that can be embedded in a relational database), its robustness (i.e., the ability of relational watermarking techniques to resist against malicious or unintentionally cyber incidents), and imperceptibility (i.e., the ability of distortion-based relational watermarking techniques to not affect the usability of the data). Each of these features is highly linked to the others, and this is the reason why it is necessary to consider a trade-off among them in the design of relational watermarking techniques. For example, one way to increase the robustness of a watermarking technique is by embedding more marks, which also

increases its capacity. But this may compromise its imperceptibility, affecting the data usability and giving clues to attackers of the watermark presence. Instead, the higher the imperceptibility, the higher is the robustness of the watermark, but this happens in spite of its capacity.

### State of the Art

After the approach of [4], several relational database watermarking techniques were proposed using different types of attributes to embed the watermark. Despite that, dealing with semantic perturbations provoked by the distortion is an issue ignored most of the time. Among the few works addressing this problem there are the techniques by [26], by [70] and by [69].

In [26], authors proposed a technique to protect the privacy and ownership of medical data. They worked with categorical attributes,<sup>3</sup> and they performed the watermark embedding as permutations of categorical values using a Domain Hierarchy Tree (DHT) of the selected attribute. This scheme does not consider more complex data types (e.g., multi-word textual<sup>4</sup>), and despite the use of the DHT for reducing the semantic perturbations, the changes compromise the inter-attribute semantic consistency. Furthermore, the latter approach reduces the watermark capacity by marking only one attribute per tuple, and if there was no DHT, then the watermark synchronization would be compromised, as the DHT is the structure on which the embedding and the extraction processes rely.

In [69, 70] the distortion is controlled by using ontologies, seeking for the preservation of the inter-attribute semantic consistency. In both the approaches, the requiring ontologies (defined for specific contexts) directly impact the blindness of the techniques, and make the watermark synchronization dependant on additional external information. Furthermore, just a single numerical attribute per tuple is selected to embed the watermark, which limits the application of the ontologies, considering the potential they may offer to increase the watermark capacity, among other things. In the end, these approaches depend on the PK of R to embed the watermark, which makes easy to compromise the watermark detection in scenarios where R is used separately from the database.

### Contribution

In this chapter, we propose a semantic-driven approach for watermarking multi-word textual attributes in relational databases, to protect their ownership. In order to preserve the meaning, fluency, grammaticality, writing style and value [104] of the multi-word textual attributes, the distortion is meant as a substitution of words that are strongly semantically similar in a certain context (i.e., synonyms).

---

<sup>3</sup>Categorical attributes represent discrete values which belong to a finite set of categories.

<sup>4</sup>We refer to multi-word attributes as textual attributes formed by one or more than one sentence.

Student					
IdNumber	StudentName	StudentSurname	Subject	Score	ProfessorJudgment
1001	John	Oliver	Mathematics	95	John has improved a lot.
1002	Justin	Fitzgerald	Physics	69	Justin has problems to pass Physics.
1003	Andrea	Russo	History	98	Andrea is the first in his History class.
1004	Karla	Olivare	Mathematics	100	Karla is an outstanding student.

TABLE 7.2: Motivating example.

Our aims are (i) to achieve a high degree of robustness without compromising the imperceptibility, (ii) increment the capacity taking care of the semantic of the data, and (iii) preserve the results of the queries on the watermarked data in comparison to those obtained with the same queries over the unwatermarked data. This avoids that the distortion caused by the watermark embedding from affecting the decision making of the organizations using and deploying the data.

Our approach is conceived so as to achieve a high watermark synchronization, making our scheme resilient against attacks based on the elimination of tuples and/or attributes. By involving other aspects to the watermarking process (e.g., the elements forming the relation structure, information corresponding to other data types, and the low redundancy of the stored data), our technique gets resilient against the *random synonym substitution attack*. The features of our approach increase the chaotic nature of the mark embedding places selection, making it difficult for attackers to determine their locations and to compromise the watermark detection by overwriting them.

We also introduce a novel approach to analyse watermark capacity through the calculation of the index  $c_w$ , which expresses the technique's resilience degree to malicious operations. This new measure is described in Section 7.4 highlighting its differences to traditional capacity measurement for relational data watermarking. Combining our proposal with numeric cover type relational watermarking techniques, allows us to increase the watermark capacity without compromising its imperceptibility, so also increasing the robustness, and making our technique effective for any practical scenarios.

## 7.2 Motivating Examples

Using numeric attributes to embed the watermark gives high coverage to relational watermarking techniques, making it possible the increment of the watermark capacity. Nevertheless, accomplishing the imperceptibility requirement at plain sight, numerical distortion compromises SQL query results based on numeric conditions.

**Example 7.1.** Consider the relation **Student** depicted in Table 7.2, where the attribute **IdNumber** denotes the primary key of the relation itself. According to the query below

for selecting the students who have passed a certain grade (**Score**  $\geq 70$ ), only the students {John, Andrea, Karla} are recovered.

```
SELECT StudentName
FROM Student
WHERE Score  $\geq 70$ 
```

If the attribute **Score** is selected to embed a mark, despite performing a passive distortion, e.g., by just using the two less significant bits (*lsb*), the result of the query above would be different. Indeed, Justin Fitzgerald could be listed among the students passing the grade if the value of the 2<sup>nd</sup> *lsb* is modified, changing, for example, the score from 69 to 71.

◻

The distortion caused by the change of the *lsb* of numerical values in a relation is not relevant when the values of the attribute chosen to embed the mark are not in the boundaries of some criteria for data recovery or their classification (e.g., changing the score of Andrea from 98 to 96 or 99, depending of the *lsb* selected as mark carrier, would not produce a different answer to the query above). But when this is not the case, such a distortion may lead to taking decisions based on wrong assumptions. Therefore, it follows that numerical distortions compromise the semantics of the tuples, despite the latter distortions being traditionally controlled by defining the maximum amount of tolerable error over the numerical attributes being watermarked.

Embedding the marks in textual attributes avoids compromising the results of queries based on numerical conditions, but applying the distortion over the *lsb* of a textual value compromises the watermark imperceptibility requirement.

**Example 7.2.** Given the relation defined in Table 7.2, consider the value of the attribute **ProfessorJudgment** for the tuple with **IdNumber** = 1002, i.e., “Justin has problem to pass Physics.”. Changing one of the two *lsb* of this textual value will provoke changing “Physics” to “Physicr” or “Physicq” making perceptible the distortion and creating a meaningless word.

◻

Note that even when the marks are embedded in textual attributes exploiting the limitations of the human vision for increasing the watermark capacity (e.g., by adding extra white spaces between words [6] or using invisible characters according to the database encoding [131]), is easy for the attacker to detect the position of the marks through computational techniques. Thus, for the aforementioned reasons, when the marks have to be embedded into textual attributes, a different approach is required.



## 7.3 Preliminaries

As mentioned in Section 7.1, the relational watermarking technique we propose marks multi-word textual attributes and preserves the semantic consistency between the original database and the watermarked one performing semantically similar words substitutions, and taking care of the context in which the words fall. Thus, we give below a brief overview of semantic similarity theory and semantic-based text watermarking techniques.

### 7.3.1 Semantic Similarity Theory

Semantic similarity is about computing the resemblance between the meanings of textual entities (e.g., words, sentences, texts), that are not necessarily lexically similar [24, 96]. It has application in many research fields [2, 143, 155, 166] among which: Natural Language Processing tasks (e.g., Word Sense Disambiguation and Synonym Detection), Artificial Intelligence, Cognitive Science, Psychology, Information Retrieval and Bio-Informatics.

The literature has several methods for measuring semantic similarity between textual entities [96]. They depend on one or several knowledge sources (e.g., taxonomies, thesaurus, ontologies) and rely on different theoretical properties [24]. A measure of similarity takes as input two textual entities and returns a numeric score that quantifies how much they are alike [166]. Formally, let  $e_1$  and  $e_2$  be two textual entities,  $\text{Sim}(e_1, e_2)$  denotes the semantic similarity between  $e_1$  and  $e_2$ .

Different words that have highly related meanings are called synonyms. Generally speaking, synonym words belong to the same node in a hierarchical knowledge organization scheme and their semantic similarity is maximized [159].

### 7.3.2 Text Watermarking

According to classifications of text watermarking techniques presented in [5, 104, 110, 192], we summarise them as follows:

- a) Image-based approaches, where a text document, whose content is seen as a series of text images, is used to embed the watermark bits. In general, these methods are resilient against typical image watermarking attacks and format-based attacks.
- b) Structure-based approaches, where imperceptible changes to the text structure, features and font are made, into which the watermarking information to be hidden is encoded. These techniques are more likely to be vulnerable to very simple attacks, e.g., the text retyping attack and the copy paste to notepad attack, and to the use of Optical Character Recognition (OCR) technologies.
- c) Syntactic-based approaches apply syntactic transformation to plain text document structures to embed the watermark. These schemes have been proved to be efficient

with agglutinative languages like Turkish, but in general they are not adequate for English language.

- d) Semantic-based approaches embed the watermark into the semantic structure of text documents, where text meaning analysis and text transformations are performed using natural language processing algorithms. Semantic-based text watermarking schemes are resilient against retyping attacks or to the use of OCR programs, but prone to weaknesses related to natural language processing.

Note that syntactic-based and semantic-based text (or content-based) watermarking schemes fall into the linguistic-based approach category. Therefore, they are highly dependant on the type of language in use, which represents a disadvantage in a scenario where languages rapidly evolve, and focused on do not (or minimally) alter the meaning of the cover text.

### **Synonym Substitution Approach**

Image-based and structure-based approaches are useful when the text is forced to be displayed by using specific means. In the case in which the text is stored as content (e.g., stored as multi-word textual attribute in relational databases), independent from its graphical representation, the text can actually be displayed in multiple ways, making those techniques useless. Consequently, content-based text watermarking techniques are better suited in the context of relational textual database watermarking. In particular, we focused on the synonym substitution method for watermarking textual documents, where certain words are replaced with their synonyms preserving the semantic consistency between the original cover text and the watermarked one.

Firstly exploited by steganography [180], synonym substitution technique was later used to watermark plain text document [169]. Still, this watermarking technique is limited to the English language and highly depends on the quality of the text processing tools, like the word sense disambiguator (i.e., a technique that aims to identify which is the sense of a word in a sentence). Moreover, synonym substitution watermarking techniques are vulnerable to synonym substitution attacks.

To overcome random synonym substitution attacks, Topkara et al. [169] proposed a lexical watermarking system based on substituting words with homographs<sup>5</sup> from their synonym set, and using meaning-preserving generalizing substitutions. Then, in order to guarantee the context-dependency between synonyms, they implemented a semi-automatic interactive encoding mechanism that allows a person designated to decide on the acceptability of the substitutions given by the system.

---

<sup>5</sup>Two or more words are homographs if they are spelled the same way but differ in meaning and origin, and sometimes in pronunciation [169].

## 7.4 Semantic-based Watermarking Approach

As pointed out in Example 7.2, modifying the *lsb* of a textual attribute value may, for example, introduce syntactic errors or cause semantic inconsistencies, leading to the imperceptibility requirement violation of the watermarking technique.

In this section, we present our semantic-based watermarking approach for relational data. Necessary condition for our technique to be applied is that the target relation must contain at least one multi-word textual attribute into which we will embed the watermark. The distortion is thought as synonym substitution.

Our proposal is focused on increasing the watermark capacity, without affecting its imperceptibility, achieving high degree of robustness against typical relational watermarking and textual watermarking attacks. Indeed, by considering multi-word textual attributes as cover type, multiple synonym substitutions can be performed over a single attribute value, resulting in the increment of the embedded marks, despite some sentences being composed of just few words, overcoming the downside of watermarking short documents, which reduces the watermark capacity [104]. Moreover, the capacity of the watermark increases when numerical attributes, in addition to multi-word textual attributes, are considered to embed the watermark. Also, taking care of preserving the meaning of the watermarked text (by using the synonym substitution approach with a proper word sense disambiguation), the imperceptibility remains untouched, which results in a direct increment of the technique's robustness.

Note that when dealing with multi-word textual cover type, there are malicious operations focused on compromising the watermark embedded in textual documents (e.g., random synonym substitution attack) that must be considered to prove the effective robustness of our approach.

### 7.4.1 Architecture of the Proposal

Figure 7.1 depicts the architecture of the watermark embedding process of our proposal. As usual, the watermark extraction involves the same elements of the embedding procedure but it is performed in the opposite direction.

Our watermark embedding procedure (fully described in Section 7.4.2) relies on a relational database DB storing one or several relations  $R$  with at least one multi-word textual attribute, one or several knowledge sources  $N$ , and a watermark  $WM$ . The choice of the knowledge source(s) is let free, but it has to take care of the semantic links between words. So that, on the latter, we can use a similarity engine to verify semantic consistency properties. A word sense disambiguation WSD module is needed for selecting the proper set of synonyms, depending on the context in which the words that are candidates to be replaced fall.

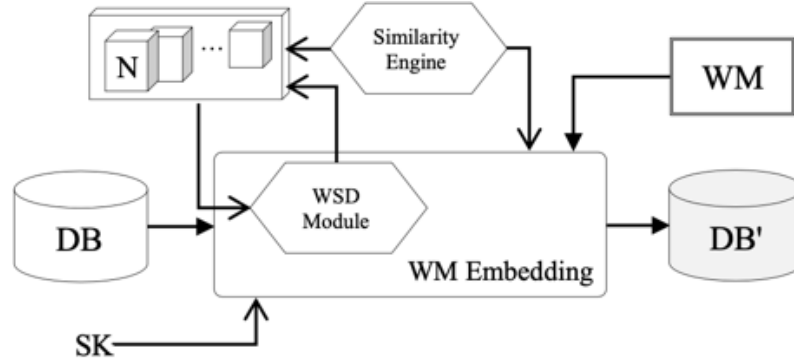


FIGURE 7.1: Embedding process architecture.

We also encourage the use of meaningful sources for the generation of the watermark considering that over this kind of signals can be applied methods to enhance the quality of the extracted watermark, contributing to its recognition despite the execution of aggressive attacks over the watermarked data. The secret key SK will be only known by the data owner and its complexity will be crucial against malicious operations (e.g., the brute force attack), as the security of watermarking techniques is based on the secrecy of the parameter values [92]. Finally, a distorted database DB' is produced. Note that our watermarking technique modifies R only for the values that are selected for the embedding. Every other value remains the same.

### Similarity Engine

Consider a relation  $R(PK, A_1, \dots, A_m)$  belonging to a database DB which stores at least one multi-word textual attribute. Let  $A_h$  (with  $h \in [1, m]$ ) be a multi-word attribute in R, and let  $r_k$  be the k-th instance of R. Then, let  $r_k.A_h$  be the value of the attribute  $A_h$  with respect to the tuple  $r_k$ . Moreover, let  $s$  be a sentence in  $r_k.A_h$ , and assume that the word  $w$  part of  $s$  has been selected by a procedure P (see Algorithm 7 in Section 7.4.2) to embed the watermark, i.e., to be replaced with its synonym  $w'$  (see Algorithm 8 in Section 7.4.2). The embedding is performed if and only if the replacement complies the *intra-attribute consistency* and the *inter-attribute consistency* semantic properties.

Note that, when it is possible, we replace only one word per sentence of a multi-word attribute value, and the semantic distortion can be embedded in more than one attribute per tuple.

**Definition 7.1** (Intra-attribute consistency). Let  $r_k.A_h \in R$  be the value of the multi-word attribute  $A_h$  for the k-th instance of R. Let  $s$  be a sentence in  $r_k.A_h$ ,  $w \in s$  be the word to replace, and let  $w'$  be the candidate substitute word. Then,  $s^*$  denotes the sentence  $s$  after replacement, i.e.,  $s[w/w']$ . Finally  $r_k.A_h[s/s^*]$  denotes the distorted result. We say that  $[w/w']$  is a substitution intra-attribute consistent if the semantic similarity score between  $r_k.A_h$  and  $r_k.A_h[s/s^*]$  is higher than or equal to a threshold  $\delta$ . Formally:

$$\text{Sim}(r_k.A_h, r_k.A_h[s/s^*]) \geq \delta$$

△

**Definition 7.2** (Inter-attribute consistency). Given  $r_k$ , i.e., the  $k$ -th instance of  $R$ , where watermark will be embedded, let  $r_k^*$  be the distorted result, i.e.,  $r_k[r_k.A_h/r_k.A_h[s/s^*]]$ . Moreover, let  $\phi$  be a function mapping tuple values to their concatenation (through “\_and\_”). Following Definition 7.1, we say that  $[w/w']$  is a substitution inter-attribute consistent if the semantic similarity score between  $\phi(r_k)$  and  $\phi(r_k^*)$  is higher than or equal to a certain threshold  $\mu$ . Formally:

$$\text{Sim}(\phi(r_k), \phi(r_k^*)) \geq \mu.$$

△

The value of the thresholds  $\delta$  and  $\mu$  depends on the chosen semantic similarity measure.

### The Word Sense Disambiguation Module

The correct functioning of the WSD module is a key element for the success of our approach. Two important issues depend on the WSD module performance: (i) maintaining the semantic value of the database being protected (ii) and the guarantee of achieving a high watermark synchronization.

WSD is considered an open research field in natural language processing. The main challenges come due to the fact that words often change meanings depending on the context they are used. For example, the noun *tree* can be used to refer to the programming data structure, but also to the living organism belonging to the vegetable kingdom. Thus, the set of synonyms allowed to replace *tree* must be shrunken according to certain context. The same happens with words used as adjectives. The adjective *hard* can be used to refer to someone with a sturdy temperament, to express determinism in business dealings, or to describe a feature of a solid object. If the ambiguity of the word is not taken away considering the context, there is a high probability the value of the text will be compromised once the word replacement is performed [104].

Let  $D$  be a function that returns the ordered set  $\mathcal{Z}$  of synonyms of a word  $w$  given a context  $\varsigma$ , denoted by  $\mathcal{Z} \leftarrow D(w, \varsigma)$ . The following rules need to be accomplished:

1.  $w \in D(w, \varsigma)$
2.  $\forall t \in D(w, \varsigma) : D(t, \varsigma) = D(w, \varsigma)$

Let  $\mathcal{Z}[t]$  be the  $t$ -th element of the set  $\mathcal{Z}$ .

For any word belonging to  $\mathcal{Z}$ , the set of synonyms for the given context must be the same to maintain the semantic value of the text from where  $w$  and  $\varsigma$  were selected. Of course, if the WSD module does not work properly, these rules can be violated considering

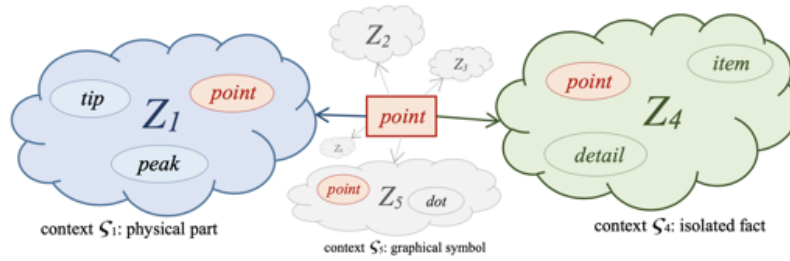


FIGURE 7.2: Synonyms sets linked to a word  $w$ .

the same word can be part of different synonym sets given by other contexts (see Fig. 7.2 where the word “point” has multiple sets of synonyms relying on different contexts in which it can be used). In general, the correct functioning of the WSD module will depend on the accuracy of the implementation of the function  $D$ , responsible for obtaining the set of synonyms of  $w$  that more fit the context  $\varsigma$ .

Synchronization is the process of aligning two signals in time or space [60]. Considering the embedded and extracted watermarks as those signals, achieving a high watermark synchronization relies on extracting the exact same marks that were embedded. If WSD fails, the rules 1. and 2. are not accomplished. Then, for example, there is the possibility that a mark embedded considering the synonyms of the set  $Z_1$  is extracted looking at synonyms of the set  $Z_4$ . Because of that, wrong mark values will be detected, compromising the quality of the extracted watermark and its synchronization. On the contrary, if previous rules are not violated, the value of the textual attribute being watermarked is preserved and the watermark synchronization is guaranteed.

### Data Quality Preservation

The keeping of watermarked data quality will mainly depend on the WSD module and the parameters defining the maximum allowable semantic distortion to perform the marks embedding. According to [169], using ambiguous words increases the resilience of the watermarking scheme against attacks, but not all documents being watermarked tolerate this kind of operation, which reduces the options to perform word replacement.

As long as rules 1. and 2. are accomplished, the word used to replace  $w$  will be obtained from  $Z$ , and the value of the mark detected during the extraction process will match the embedded one. Also, by defining the maximum tolerated semantic distortion, it is avoided the use of words belonging to  $Z$  that might involve ambiguity, increasing the probability of falling in another synonym set. Thus, the equivalence above guarantees the preservation of the data quality, no matter the nature of the text of  $R$  being protected.

---

**Algorithm 7** watermarkingEmbedding procedure
 

---

**Input:**  $R, SK, \gamma, \mathcal{A}, WM, \varphi, \chi, N, \mathcal{C}, \delta, \mu$   
**Output:**  $R'$ , i.e., the watermarked relation

- 1: **for** tuple  $r \in R$  **do**
- 2:      $k_r = \text{VPK}(SK \circ r_K)$
- 3:     **if**  $(k_r \bmod \gamma) = 0$  **then**
- 4:          $\mathcal{A}' \leftarrow \Theta(\mathcal{A}, r_K, \varphi)$
- 5:         **for**  $v \in \mathcal{A}'$  **do**
- 6:              $\mathcal{S} \leftarrow \Lambda(v, \chi)$
- 7:             **for** sentence  $s \in \mathcal{S}$  **do**
- 8:                  $k_s = H(SK \circ k_r \circ \Gamma(s))$
- 9:                  $i = k_s \bmod \Upsilon(s).\text{length}$
- 10:                  $w = \Upsilon(s)[i]$
- 11:                  $\varsigma \leftarrow \text{getSense}(w, s, N)$
- 12:                  $\mathcal{Z} \leftarrow \text{getCandidates}(\varsigma, w, N)$
- 13:                  $w' = \text{getSubstitute}(WM, k_s, \mathcal{Z}, \mathcal{C}, \vartheta)$
- 14:                  $s^* \leftarrow s[w/w']$
- 15:                  $v^* \leftarrow v[s/s^*]$
- 16:                 **if**  $\text{Sim}(v, v^*) < \delta$  **then**
- 17:                     rollback embedding
- 18:                 **else**
- 19:                      $r^* \leftarrow r[v/v^*]$
- 20:                     **if**  $\text{Sim}(\phi(r), \phi(r^*)) \geq \mu$  **then**
- 21:                          $r \leftarrow r^*$
- 22:                          $v \leftarrow v^*$
- 23:                     commit embedding
- 24:                 **else**
- 25:                     rollback embedding

---

### 7.4.2 Watermarking Procedure

Distortion-based relational watermarking techniques consist of two processes: (i) the embedding of the watermark (ii) and its extraction [92]. To achieve the watermark synchronization it is required the use of the same parameter values in both processes. Moreover, the extraction is performed when it is required to demonstrate the watermark presence in the data, as evidence in case of ownership claims, among others.

#### Watermark Embedding

Algorithm 7 presents the details of the watermark embedding procedure of our approach. Given a relation  $R$ , for each tuple  $r \in R$ , the values of the multi-word textual attributes composing the list  $\mathcal{A}$  are analyzed. Then, the virtual primary key  $k_r$  is generated using the VPK function (line 2). The input of the latter function results from the concatenation

---

**Algorithm 8** getSubstitute procedure
 

---

**Input:** WM,  $k_s$ ,  $\mathcal{Z}$ ,  $\mathcal{C}$ ,  $\vartheta$ 
**Output:**  $w'$ 

```

1:  $i = k_s \bmod \text{WM.length}$ 
2:  $m \leftarrow \text{WM}[i]$  set_order( $\mathcal{Z}$ ,  $\mathcal{C}$ )
3: if  $m = 1$  then
4:   return  $\mathcal{Z}[0]$ 
5: else
6:   return  $\mathcal{Z}[\vartheta]$ 

```

---

( $\circ$ ) between a secret key SK, and data represented by  $r_K$  identifying the tuple  $r$  (e.g., the relation's PK or other virtual primary keys generated by external schemes).

Next, in the case in which the **if-statement** condition is satisfied, a filter  $\varphi$  is applied (by the function  $\Theta$ ) to the multi-word attribute values in  $\mathcal{A}$ , in order to exclude those considering their content and links with the other attributes of the tuple (e.g., exclusion of sentences containing acronyms or abbreviations). In this way, we add an extra step to help maintain *inter-attribute consistency* while high unpredictability is incorporated into the technique<sup>6</sup>. The attributes passing the latter filter are stored in the set  $\mathcal{A}'$  (lines 3-4).

Similarly as above, for each multi-word attribute value  $v$  in  $\mathcal{A}'$ , a filter  $\chi$  is applied (by the function  $\Lambda$ ) to the sentences in  $v$ , to exclude those that do not accomplish the conditions to be considered for the embedding process (e.g., involving just sentences composed of more than certain number of words). Sentences accomplishing the conditions defined in  $\chi$  are finally considered for the embedding, and they are stored in the set  $\mathcal{S}$  (lines 5-6).

For each sentence in  $\mathcal{S}$  the key  $k_s$  is generated (from a one-way hash function  $H$  that takes as input the concatenation of SK,  $k_r$  and the elements of the sentence do not tolerating changes obtained by the  $\Gamma$  function) that identifies the sentence inside the multi-word attribute value, and using  $k_s$  the word  $w$  to be replaced is selected (lines 7-10). Note that  $\Upsilon$  is a function that returns, as an array of words, the elements of a sentence tolerating substitutions.

Then, according to the sense  $\zeta$  of the selected word in the sentence under consideration (obtained using the `getSense` method), the set  $\mathcal{Z}$  of synonyms of  $w$  is obtained (by using the `getCandidates` method) (lines 11-12). Both `getSense` and `getCandidates` functions are based on the set of rules and the definitions given by the knowledge source(s)  $N$ .

Finally, the mark to be embedded and the new word  $w'$  to replace  $w$  are selected (line 13). Algorithm 8 defines the function `getSubstitute` where the mark is selected according to the value of  $k_s$ . The set of candidate substitute words is sorted according to the criteria  $\mathcal{C}$ , and the new word is chosen depending on the value of the mark to be embedded. In lines 14-25, the replacement of the sentence in the carrier attribute and in the corresponding

---

<sup>6</sup>Each attribute's value contributes itself in varying the elements of  $R$  selected as mark carriers.



tuple is performed. The substitution will only be carried out if *intra-attribute* and *inter-attribute* consistencies properties are not violated.

If a word is not detected in the knowledge sources (i.e., *cross-linguality problem* [2]), our approach ignores the position from the marking's candidates and proceed with the rest of the data stored in the relation.

### Watermark Extraction

The watermark extraction process is similar to the embedding process but is performed in the opposite direction. Also, it must be performed using the same parameter values employed to embed the watermark to guarantee its right synchronization. In the extraction process, for the same mark position in the watermark, several elements are recovered, and before assigning the mark final value, a majority voting is performed. In this way, the effect of attacks based on low aggressive data modifications are avoided.

The extraction is performed with no need to check the semantic distortion between the words replaced to carry out the marks embedding (i.e., without considering the similarity metrics value). Nevertheless, the extracted watermark quality depends on the knowledge source(s) and on the precision of the WSD module, since words can be assigned to a set of synonyms different from those considered for the embedding, adding noise to the extracted signal. This is because, in the new set of synonyms, the original word can occupy a different position, assigning to the extracted mark a wrong value.

#### 7.4.3 Analysis of the Watermark Capacity

Since relational data have no fixed order, sequential watermarking approaches are vulnerable to subset reverse order attacks [92]. To overcome this vulnerability, techniques have been designed performing a *pseudo-random selection* of both the watermark source elements used to generate the marks and the embedding locations in R. In general, this operation has been achieved by using a specific definition of Equation 7.1. Nevertheless, besides contributing to robustness against subset reverse order attacks, *pseudo-random selection* makes it difficult for the attackers to predict embedding locations for overwriting or removing of marks.

$$\mathcal{V}(x, y) \bmod \mathcal{M}_X \quad (7.1)$$

where,  $\mathcal{V}(x, y)$  is a value generated using data from a generic position  $(x, y)$  of R, and  $\mathcal{M}_X$  is the maximum value of a given range. An example of a particular definition of this expression is in line 10 of Algorithm 7, where  $\mathcal{V}(x, y)$  is defined in line 9 as  $k_s$ .

Because of *pseudo-random selection*, during the embedding process some marks are selected more than once while others are entirely ignored. Embedding the same mark multiple times leads the technique to be resilient against update attacks, if a majority voting

is performed over each mark position in the extraction process. On the other hand, if the number of excluded marks is too high, the watermark synchronization can be compromised. Indeed, the latter process would suffer from the the same negative consequences as when aggressive attacks based on updating or deleting data are performed.

The *pseudo-random selection* downside can be reduced if the number of times the watermark source elements are considered increases. This can be achieved by marking a higher volume of data while the same watermark source is used.

In our approach, new elements to increase the watermark capacity, without affecting the imperceptibility requirement, are introduced. Moreover, we propose a new metric for measuring the watermark capacity which considers the different number of times each mark is selected during the embedding process. The latter metric allows the evaluation of the capacity in function of the technique's robustness, since it assigns to each mark a weight depending on the number of times it is embedded in R. In this way, the difficulty of compromising each mark in the watermarked data is reflected in the metric.

For techniques embedding one mark per selected tuple, the number of embedded marks E is equal to the number of marked tuples  $\omega$ . If, besides the numeric cover type, multi-word textual attributes are considered (following the approach we proposed above), the number of embedded marks increases for each tuple according to Equation (7.2), where  $\lambda_n$  represents the number of numeric attributes,  $\lambda_s$  the number of marked sentences, and  $\bar{\delta}$  the number of marks embedded in each sentence. If the watermark capacity increases due to marks embedded on multi-word textual attributes (besides those embedded on numerical attributes), using an effective WSD module, the technique becomes more resilient without compromising data usability.

$$E \approx \omega * (\lambda_n + \bar{\delta} * \lambda_s) \quad (7.2)$$

Even if no numerical attributes are considered and only one multi-word textual attribute is selected to perform the embedding, more than one mark can be embedded per tuple, according to  $\bar{\delta}$ , which still constitutes an increment of the capacity compared to other watermarking techniques.

The number of marks selected for the embedding (denoted by  $m_e$ ) using as reference the watermark size (denoted by  $n$ ) is commonly used to measure the technique's watermark capacity. We define that metric as the binary capacity, given by  $c_b$  according to Equation (7.3). The downside of  $c_b$  is that all marks present the same weight, and only their inclusion/exclusion represents information for the measurement.

$$c_b = m_e * 100/n \quad (7.3)$$

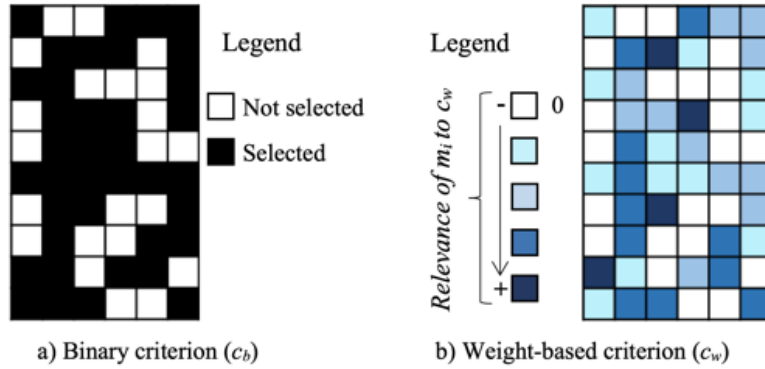


FIGURE 7.3: Criteria to evaluate the watermark capacity.

We consider the number of times each mark is embedded and we propose the weight-based capacity metric, given by  $c_w$  according to Equation (7.4).

$$c_w = \sum_{i=0}^{n-1} (\mathcal{Z}(m_i)) / n \quad (7.4)$$

where,  $\mathcal{Z}(m_i)$  represents the number of times the mark  $m_i$  was embedded. This is given since not all marks are selected the same number of times. In Figure 7.3 we used a scale of colors to illustrate with an example the differences between  $c_b$  and  $c_w$ . The value of  $c_w$  is the mean of the number of times all marks of the watermark are embedded.

The higher  $c_w$  the better, but it is also required that  $\rho_w \approx 0$ , being  $\rho_w$  the standard deviation of the number of times each mark is embedded. This tells us that each mark was selected multiple times evenly, adding higher relevance to the watermark recognition, increasing the probability of its detection despite the execution of benign updates and attacks over R.

Note that the highest value of  $c_w$  will be difficult to achieve as this would mean embedding each mark evenly the maximum possible number of times, and since the process presents *pseudo-random nature*, this is not expected.

In general,  $c_b$  does not reflect how embedding each mark multiple times contributes to obtain a different degree of resilience as  $c_w$  does. This is important to measure as when more redundancy is achieved for the embedding of a mark value, it is more difficult to compromise its value with update attacks, if a majority voting of all recovered values is performed during watermark extraction.

#### 7.4.4 Considerations for the Adversary Model

One of the major challenges textual watermarking techniques based on synonyms substitution face is the *random synonyms substitution* attacks. This vulnerability is linked to the tolerance a text has with respect to synonyms substitution on it, depending on the

context in which the data was used. Performing embedding of marks in multi-word textual attributes through semantically similar words replacement, we must consider adding resilience to this threat, as well as to the rest of the malicious operations an attacker may perform over the database relation.

Depending on the context in which data are used, we have a margin of marks allowed to be embedded. If the number of marks that can be embedded is high, an high degree of robustness can be achieved, otherwise, the technique's resilience gets compromised. Moreover, the data context also reduces the attacker's freedom to perform aggressive operations if he wants to preserve the quality of the data. Therefore, using textual watermarking to relational data increases the difficulty to perform effective attacks.

Increasing the capacity by considering both, the numerical and textual cover types, allows achieving higher robustness, making it hard to compromise the watermark detection. Several elements of the relation, e.g., attributes storing float numbers, single sentences or several paragraphs, can be selected to embed the marks, based on multiple features (e.g., numerical ranges, minimum number of nouns required per sentence, among others) which highly increase the complexity to detect the marks embedding locations. Furthermore, beyond the complexity of the *pseudo-random selection* of the marks embedding places in the relation, we take advantage of the multi-word textual data structure to add a high entropy to the embedding procedures.

[4] defined Equation (7.5) to get the probability for the attacker to successfully detect the embedding locations used by the data owner. While  $\omega$  refers to the number of tuples watermarked by the data owner,  $\gamma_A$ ,  $\nu_A$ , and  $\xi_A$  denote respectively the tuple fraction, the number of attributes, and the *lsb* considered by the attacker.

$$P\{\text{success}|\omega\} = \left(1 - \frac{1}{2\gamma_A\nu_A\xi_A}\right)^\omega \quad (7.5)$$

Considering all those elements, it is difficult to detect the marks embedding positions used by the data owner. Despite that, we increase the difficulty by adding the element  $\zeta_A$  to Equation (7.5). Equation (7.6) extends Equation (7.5) as follows:

$$P\{\text{success}|\omega\} = \left(1 - \frac{1}{2\gamma_A\nu_A(\xi_A + \zeta_A)}\right)^\omega \quad (7.6)$$

The term  $\zeta_A$  defines the complexity that derives from the consideration of marks embedding locations among multi-word textual cover types. It expresses the possibility of embedding one mark in the whole attribute value, or marking each sentence with one mark, or embedding multiple marks in each sentence. Since all embedding positions for textual attributes are also generated using *pseudo-random selection* and considering the increment of the number of elements to know by the attacker, the probability of performing successful attacks reduces. Furthermore, by accomplishing security and public system



FIGURE 7.4: Samples of the binary images used as WM sources.

requirements [3], we add secrecy to parameter values, increasing the difficulty of attackers to detect embedding locations.

## 7.5 Experimental Results

Following the recommendation given in Section 7.4.1 of considering meaningful sources to generate the watermark WM, we validated our approach by using binary images as WM sources. Besides the benefits mentioned above, this type of data allows taking the simplest pixel value for the mark generation, which contributes to perform less aggressive distortions during the watermark embedding compared to techniques generating marks from color (e.g., [190]) or gray-scale images (e.g., [191]).

To analyze the effect of the watermark length variation, images of different sizes were used. Samples of them are shown in Figure 7.4: a) the Universiti Teknologi Malaysia (UTM) logo ( $82 \times 80$  pixels), b) the logo of the World Wildlife Fund (WWF) ( $40 \times 45$  pixels), and c) the Chinese character dào ( $20 \times 21$  pixels). By convention, we used the red color to highlight the missed pixels due to watermark incomplete embedding or malicious operations by attackers.

To know the quality of the extracted WM two metrics were used: the Correction Factor (CF) and the Structural Similarity Index (SSIM). The Correction Factor, defined by Equation (7.7), is used to compare the pixels of the image generated from the embedded WM (given by  $\text{Img}_{\text{emb}}$ ) against the pixels of the image generated from the extracted WM (given by  $\text{Img}_{\text{ext}}$ ). In the equation, variables  $h$  and  $w$  represent the height and the width of the images respectively. The maximum value of CF is 100, meaning that the WMs are identical. In the case in which  $\text{CF} = 0$ , then the embedded and the extracted WMs are completely different.

$$\text{CF} = \frac{\sum_{i=1}^h \sum_{j=1}^w (\text{Img}_{\text{emb}}(i,j) \oplus \overline{\text{Img}_{\text{ext}}(i,j)})}{h \times w} \times 100 \quad (7.7)$$

The SSIM is oriented to obtain an appreciation of the image's quality closer to human perception. The index is calculated according to Equation (7.8), using multiple windows. The windows are defined by  $x$  and  $y$  and present common size  $N \times N$ . The range of possible values taken by this metric in this work is between 0 and 1, where 1 means there

Attribute	Type	Description
ProductId	String	Id. of the product
UserId	String	Id of the user
ProfileName	String	Name of the user
HelpfulnessNumerator	Numeric	Numerator of the fraction of users who found the review helpful
HelpfulnessDenominator	Numeric	Denominator of the fraction of users who found the review helpful
Score	Numeric	Rating of the product
Time	Numeric	Time of the review (unix time)
Summary	String	Review summary
Text	String	Text of the review

TABLE 7.3: Structure of the dataset “Amazon Fine Food Reviews”.

exists a perfect structural similarity between the embedded and the extracted image, and 0 indicates no structural similarity.

$$\text{SSIM}(x, y) = \frac{(2\mu_x\mu_y + C_1) + (2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \quad (7.8)$$

The symbols  $\mu_x$  and  $\mu_y$  represent the average of  $x$  and  $y$  respectively,  $\sigma_x^2$  and  $\sigma_y^2$  their respective variance, and  $\sigma_{xy}$  their covariance. The elements  $C_1$  and  $C_2$  are two stabilization constants.

The data to embed the marks was the data set *Amazon Fine Food Reviews*. The structure is depicted in Table 7.3, from where we mostly used the attribute ‘*Text*’, also storing the text with the highest length. We also used the first 30.000 tuples out of 500,000 to compare our results with previous works.

We used the relation’s PK to perform the watermark synchronization. To avoid the use of the PK we recommend the generation of virtual primary keys by using the Ext-Scheme [81] or the HQR-Scheme [83]. These schemes were originally proposed for numerical attributes, but they can be applied by combining numerical and textual cover types as well.

We used WordNet, as knowledge source, [132], which consists of a lexical database of the English language, where nouns, verbs, adjectives, and adverbs are grouped into sets of cognitive synonyms (synsets), each expressing a distinct concept [147]. So, given a word  $w$ , and the context where  $w$  is used (i.e., the meaning of the sentence to which  $w$  is part of), WordNet returns the appropriate synset of  $w$ . Notice that, in our experiment, the set of synonyms  $\mathcal{Z}$  (introduced in Section 7.4) will correspond to a synset in WordNet.

For the evaluation of our approach, we implemented a client-server architecture application using Java 1.8 programming language for the client-side, the Oracle Database 12c for the server-side. We used WordNet 3.1 database files with the Java API jwnl 1.4.1

rc2 for accessing and working with WordNet resources and ws4j 1.0.1 for using Semantic Relatedness/Similarity algorithms already developed. The WSD module was based on the Lesk algorithm [171], which compares the word definitions with the definitions of the rest of the words presented in the sentence, finding the more convenient context for its use. According to that, the most appropriate set of synonyms can be chosen. Finally, the runtime environment was a 3.60 GHz Intel i7-4790 PC with 16.0 GB of RAM running on Windows 10 OS.

### 7.5.1 Improvement of the Watermark Capacity

We compared the watermark capacity analysis of our approach with two other techniques [151], which uses only one attribute, and [82] with two attributes. Of all techniques using an image to generate the watermark (i.e., Image-Based Watermarking (IBW) [92]) these techniques constitute a representation of the ones more recent, used to mark one or several attributes per tuple. We selected only one attribute to mark with our approach to show that WM capacity improves even compared to two-attributes embedding, due to the selected cover type. The techniques we compare with performed the watermark embedding on the numerical cover type, but by involving the same number of tuples we can appreciate how much the watermark capacity increases for our approach.

Table 7.4 shows the value of  $c_b$  (see Equation (7.3) Section 7.4.3) obtained for each technique. In the table, columns titles “S & H” refers to Sardroudi and Ibrahim’s technique, “G. et al.” to Pérez Gort et al.’s and “Prop.” to our proposal. Given that, the number of marks missed by using our approach is lower than the number of marks missed by using the other techniques. Indeed, there are fewer red pixels in the images of the WMs synchronized by our proposal. The main reason for WM improvement is that for some cases the values stored in the attribute “*Text*” are composed of more than one sentence. If allowed, we only embed one mark per sentence selecting a common *nouns* from it. By more than one mark can be used as the carrier, in which case WM capacity will be even higher.

The Table 7.5 shows the value of  $c_w$  (see Equation (7.4) Section 7.4.3) with the correspondent  $\rho_w$  for each case, giving a clear idea about how we not only improve the value for  $c_b$  but for  $c_w$  as well, increasing the probability of overcoming attacks based on data updates.

The experiments to register the capacity values were applied over a set of 30000 tuples. For the case of Sardroudi & Ibrahim’s and Pérez Gort et al.’s, it was used the numerical data set *Forest Cover Type* [39] as these techniques were designed for marking numerical values. Also, the watermark embedding with Pérez Gort et al.’s technique was performed with Attribute Fraction equal to 5 in order to mark only two attributes per tuple.

Once the WM capacity increase was proven, it is critical to guarantee marks detection, otherwise, it will not be possible to recognize the WM signal in the protected data. In























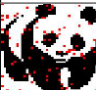






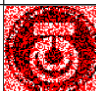

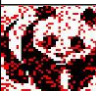




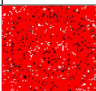
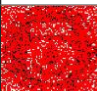
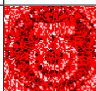






TF	UTM			WWF			Đào		
	S & H	G et al.	Prop.	S & H	G et al.	Prop.	S & H	G et al.	Prop.
2									
	88.34%	98.61%	99.91%	99.94%	99.94%	99.94%	99.76%	99.76%	99.76%
5									
	57.93%	81.46%	93.67%	95.16%	99.72%	99.89%	99.76%	99.76%	99.76%
10									
	35.64%	56.99%	75.35%	78.77%	94.88%	99.44%	99.29%	99.76%	99.76%
20									
	19.25%	34.28%	48.96%	53.88%	78.66%	91.83%	97.38%	99.52%	99.76%
40									
	9.66%	18.14%	27.99%	31.27%	51.88%	69.56%	80.95%	94.04%	98.81%

TABLE 7.4: Value of  $c_b$  for different techniques.

TF	UTM			WWF			Đào		
	S & H	G et al.	Proposal	S & H	G et al.	Proposal	S & H	G et al.	Proposal
2	2 ( $\pm 1.47$ )	4 ( $\pm 2.12$ )	6 ( $\pm 2.73$ )	7 ( $\pm 2.95$ )	15 ( $\pm 4.20$ )	24 ( $\pm 5.12$ )	33 ( $\pm 5.92$ )	67 ( $\pm 8.96$ )	105 ( $\pm 11.18$ )
5	0 ( $\pm 1.27$ )	1 ( $\pm 1.51$ )	2 ( $\pm 1.87$ )	3 ( $\pm 1.81$ )	6 ( $\pm 2.56$ )	10 ( $\pm 3.27$ )	13 ( $\pm 3.65$ )	26 ( $\pm 5.46$ )	43 ( $\pm 6.79$ )
10	0 ( $\pm 0.80$ )	0 ( $\pm 1.26$ )	1 ( $\pm 1.23$ )	1 ( $\pm 1.43$ )	3 ( $\pm 1.79$ )	5 ( $\pm 2.28$ )	6 ( $\pm 2.72$ )	13 ( $\pm 3.73$ )	21 ( $\pm 4.67$ )
20	0 ( $\pm 0.51$ )	0 ( $\pm 0.78$ )	0 ( $\pm 1.07$ )	0 ( $\pm 1.18$ )	1 ( $\pm 1.36$ )	2 ( $\pm 1.64$ )	3 ( $\pm 1.86$ )	6 ( $\pm 2.72$ )	10 ( $\pm 3.20$ )
40	0 ( $\pm 0.34$ )	0 ( $\pm 0.50$ )	0 ( $\pm 0.66$ )	0 ( $\pm 0.72$ )	0 ( $\pm 1.16$ )	1 ( $\pm 1.12$ )	1 ( $\pm 1.37$ )	3 ( $\pm 1.89$ )	5 ( $\pm 2.33$ )

TABLE 7.5:  $c_w$  value with its correspondent  $\rho_w$  for each experiment.

the following, the results focused on testing the WM detectability are shown. It is also analyzed the way the WSD module precision determines the quality of the extracted WM.

### 7.5.2 Detectability Analysis

The detectability of marks in our approach is directly linked to the precision of the WSD module. Since for WM embedding are used set of synonyms of the selected word, it is not expected data quality degradation, but if for WM extraction, the WSD module does not assign the same set of synonyms used for the embedding, then several marks will be recovered with wrong values, adding noise to the extracted WM signal. If the signal is too noisy, WM synchronization can be compromised, making impossible its identification.

The main goal of this work is not to improve the WSD algorithms, but to use those already defined that guarantee high precision for marks detection. The WSD module precision will be denoted as  $\mathcal{P}$ , which is obtained according to Equation (7.9), where  $\mathcal{W}_T$



TF	UTM	WWF	Dào
2	0.9479	0.9498	0.9518
5	0.9478	0.9500	0.9517
10	0.9518	0.9513	0.9516
20	0.9482	0.9422	0.9466
40	0.9532	0.9472	0.9463

TABLE 7.6: WSD precision during WM detection.

represents the number of tagged words (i.e., words selected for sense disambiguation) and  $\mathcal{W}_{CT}$  the number of words correctly tagged (i.e., words that during the extraction process were linked to the same synonym set used for WM embedding).

$$\mathcal{P} = \mathcal{W}_{CT}/\mathcal{W}_T \quad (7.9)$$

The WSD module we used is based on the Lesk algorithm [171]. The precision described by this module was registered through a set of experiments, whose results are shown in Table 7.6.

Note that when combining WSD along with the watermarking technique, WSD lack of precision can be compensated by the majority voting performed over WM extraction. Then, the higher  $c_w$  with  $\rho_w$  closer to zero, the stronger the effect of the majority voting to overcoming low WSD precision. The Table 7.7 shows the results experimentally supporting this point, linking the quality of the detected WM to the results of Table 7.5 despite the precision weaknesses shown in Table 7.6.

The results shown in Table 7.7 correspond to a set of experiments with different values for TF and WM sources. For each case, it is shown the embedded WM (the small image), the detected one (the big image), and the value of SSIM with a percentage corresponding to the number of pixels not matching between the two images. The low value of this metric compared to those shown in Table 7.6 directly endorses our statement.

Previous results show how wrong mark values are not reflected in red pixels, but in black and white, added to wrong regions of the image. This can be understood as a *Gaussian noise* which degree is directly linked to the precision of the WSD module. In Table 7.8 are shown other results, which describe the quality of the extracted WM in more detail. This time, the quality of the detected WM is given respect to both, the embedded WM and the original one. Of course, since depending on the parameter's values the original WM is not usually entirely embedded, the higher quality will be the one given respect to the embedded WM.

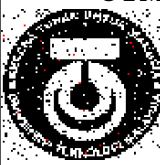









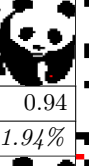


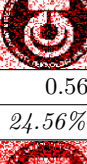


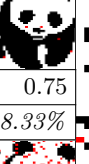
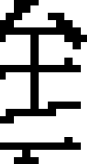
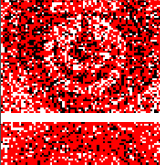
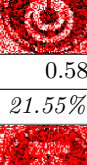


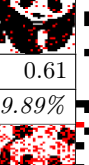


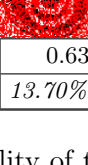


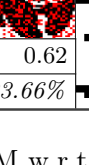

TF	UTM		WWF		Đào	
2		 0.69 4.92%		 1.00 0		 1.00 0
5		 0.55 18.00%		 0.94 1.94%		 1.00 0
10		 0.56 24.56%		 0.75 8.33%		 1.00 0
20		 0.58 21.55%		 0.61 19.89%		 1.00 0.48%
40		 0.63 13.70%		 0.62 23.66%		 0.84 7.38%

TABLE 7.7: Quality of the detected WM w.r.t. the embedded one.

TF	UTM				WWF				Đào			
	vs. Embedded		vs. Original		vs. Embedded		vs. Original		vs. Embedded		vs. Original	
	CF	SSIM	CF	SSIM	CF	SSIM	CF	SSIM	CF	SSIM	CF	SSIM
2	99.42%	0.69	95.06%	0.68	100%	1	99.94%	0.99	100%	1	99.76%	0.99
5	88.34%	0.55	80.06%	0.49	100%	0.94	98.00%	0.93	100%	1	99.76%	0.99
10	71.11%	0.56	61.32%	0.37	98.77%	0.75	91.72%	0.75	100%	1	99.76%	0.99
20	54.42%	0.58	37.85%	0.25	83.54%	0.61	76.66%	0.54	100%	1	99.28%	0.99
40	46.51%	0.63	21.9%	0.17	67.41%	0.62	55.33%	0.41	98.07%	0.84	91.9%	0.83

TABLE 7.8: Detectability achieved for each WM over different numbers of marked tuples.

Data detectability is benefited from the combination between the WSD module, which takes care of imperceptibility, and majority voting, which allows being tolerated WSD lack of precision. The positive impact of this effect increases when WM size decreases or the number of tuples being watermarked increases. The following experiments are meant to analyze how data usability and watermark imperceptibility are maintained. This two WM

<b>TF</b>	<b>UTM</b>	<b>WWF</b>	<b>Đào</b>
2	0.3752	0.3655	0.3146
5	0.3765	0.3652	0.3142
10	0.3763	0.3597	0.3153
20	0.3768	0.3591	0.3094
40	0.3961	0.3590	0.3068

TABLE 7.9: Value of  $T_w$  for previous experiments.

requirements are critical, since if are not accomplished, the technique becomes useless for practical scenarios.

### 7.5.3 Watermark Imperceptibility

As mentioned before, the embedding process is carried out through the replacement of a *pseudo-randomly* selected word by a synonym from a specific set of synonyms. The word selected from the latter set will depend on the value of the mark extracted from the binary image. Occasionally, the word selected from the set of synonyms is the same one that was selected from the sentence. That is the scenario that best contributes to WM imperceptibility since the mark is embedded without any modification in the data. We compute the rate of marks embedded without performing word replacement through the rate of fixed words given by  $T_w = \mathcal{W}_F/\mathcal{W}_T$ , where  $\mathcal{W}_F$  represents the words that do not change during the embedding process and  $\mathcal{W}_T$ , as previously defined, represents the number of tagged words. The value of  $T_w$  for each one of the experiments is shown in Table 7.9, where more or less for all cases a third of the selected words allows mark embedding without being replaced, which positively contributes to achieving WM imperceptibility.

From the knowledge source point of view, it is also possible to detect the quality of the WM imperceptibility. For this case, since we are using WordNet, we can use a set of similarity metrics that are defined to measure the relatedness or similarity between words. According to WordNet structure, and the way the metrics are defined, when two words are selected from the same synonym set, the metrics report the highest possible value between them. Table 7.10 gives the accumulated value for some metrics commonly used, for the experiments performed during mark embedding using UTM as WM source. The table’s column “Iter.” refers to iterations, meaning the number of times the measurement between words was carried out.

Thanks to the use of similarity metrics it is possible to determine and control the amount of distortion introduced during WM embedding. This contributes to maintaining data usability and WM imperceptibility, goals that highly depend on the knowledge source, the similarity engine and the WSD module. For our case, as long as the words belong to the same set of synonyms, quality results are guaranteed.

TF	Iter.	WUP	JCN	LCH	LIN	RES	PATH	LESK
2	44510	$44.64 \times 10^3$	$4.73 \times 10^{11}$	$16.40 \times 10^4$	$44.45 \times 10^3$	$38.02 \times 10^4$	$44.46 \times 10^3$	$44.22 \times 10^4$
5	18321	$18.39 \times 10^3$	$1.96 \times 10^{11}$	$67.55 \times 10^3$	$18.31 \times 10^3$	$15.71 \times 10^4$	$18.31 \times 10^3$	$18.17 \times 10^4$
10	9091	$91.20 \times 10^2$	$9.72 \times 10^{10}$	$33.52 \times 10^3$	$90.84 \times 10^2$	$78.00 \times 10^3$	$90.85 \times 10^2$	$89.75 \times 10^3$
20	4427	$44.37 \times 10^2$	$4.73 \times 10^{10}$	$16.32 \times 10^3$	$44.22 \times 10^2$	$38.00 \times 10^3$	$44.23 \times 10^2$	$43.05 \times 10^3$
40	2156	$21.63 \times 10^2$	$2.31 \times 10^{10}$	$79.53 \times 10^2$	$21.55 \times 10^2$	$18.65 \times 10^3$	$21.56 \times 10^2$	$21.45 \times 10^3$

TABLE 7.10: Similarities metrics for WM UTM.

Even though, since our technique is meant to be used for relational data copyright protection, we have to guarantee robustness against malicious operations. The core of our approach has been proved to be resilient against common malicious operations oriented to compromise relational data watermarking techniques [4, 82, 151]. Nevertheless, we need to consider another threat more focused on compromising WM detection over text documents. The next section is oriented to analyze the resilience of our technique.

#### 7.5.4 Technique’s Robustness

The major threat our technique faces is linked to *random synonym substitution* attacks from watermarking techniques created for document protection. In a similar way that WSD and majority voting combined allow overcoming WSD lack of precision, using textual as WM cover type in relational data reduces the probability of performing a successful *random synonym substitution* attack. This is because to successfully compromise the mark value, the right words need to be selected. But first, it is required to detect the right tuple, and the attribute selected for marking inside the tuple. The high number of parameters involved in the technique makes that very difficult to achieve.

This is one of the benefits of combining both cover types. The probability of successfully overwriting marks decreases if, besides relational elements, textual’s are considered. To that, once the position is correctly detected in the relation, it is necessary to know the type of word selected for the embedding (e.g., noun, adverbs, adjectives, etc.), the sentence itself, and break the secrecy of  $k_s$  (see Algorithm 7). As we mentioned in Section 7.4.4, this is ruled by the adversary model obtained as a result of extending Equation (7.5) to Equation (7.6).

To study our approach’s resilience to *random updates* we performed two types of experiments, random tuple deletion and random actualization of words stored in the same attribute we use for marking. In Table 7.11 is depicted the degree of damage the WM gets while the number of pseudo-randomly deleted tuples increases. This experiment was performed attacking the relation marked with the WM generated from the image WWF with parameters  $TF = 2$ , detected with quality of  $SSIM = 1$  with no pixels in contradiction with respect to the embedded WM. (see Table 7.7).










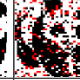
Datum	Percentage of tuples deleted (attack degree).									
	0	10%	20%	30%	40%	50%	60%	70%	80%	90%
Image										
SSIM	1	0.99	1	0.99	0.94	0.93	0.90	0.82	0.74	0.55
CF	100%	100%	100%	100%	100%	100%	100%	99.88%	97.10%	78.65%

TABLE 7.11: WMs detected after pseudo-random tuple deletion attacks.











Datum	Percentage of the attributes updated (attack degree)									
	0	10%	20%	30%	40%	50%	60%	70%	80%	90%
Image										
SSIM	1	1	0.99	0.99	0.97	0.92	0.89	0.89	0.78	0.78
CF	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%

TABLE 7.12: WMs detected after pseudo-random update attacks.

The second robustness experiment was performed under the same conditions of the previous one, but updating the value of the attribute, randomly selecting the tuple according to the attack’s percentage. In this case, the update operation is based on selecting the same attribute value but replacing a word from the sentence for a synonym. The selection of the tuple, the word being replaced, and the synonym was made pseudo-randomly. This experiment was performed this way to simulate random synonym substitution attacks focused on compromising textual watermarking. If the mark value during the WM detection is given by exclusion (assigning 1 if the word is the first one in the synonym set and 0 if it is not), the probability of success for this type of attack decreases considerably.

From the experiments performed in this section, it is clear our technique describes a resilience that guarantees WM preservation despite data degradation due to malicious operations (see Tables 7.11 and 7.12). Considering attackers are also interested in maintaining data quality, it is not expected they will exceed the degree of damage caused to the data in the experiments we performed. Then, we claim our technique is resilient, being recommended for practical scenarios to guarantee copyright protection.

### 7.5.5 Scalability and Complexity

Since WM embedding and extraction processes are similar in complexity, in this section we report the time required for WM embedding. We carried out a set of experiments, performing WM embedding multiple times involving a different number of tuples. Figure 7.5 depicts the linear correlation between the time required by WM embedding and the

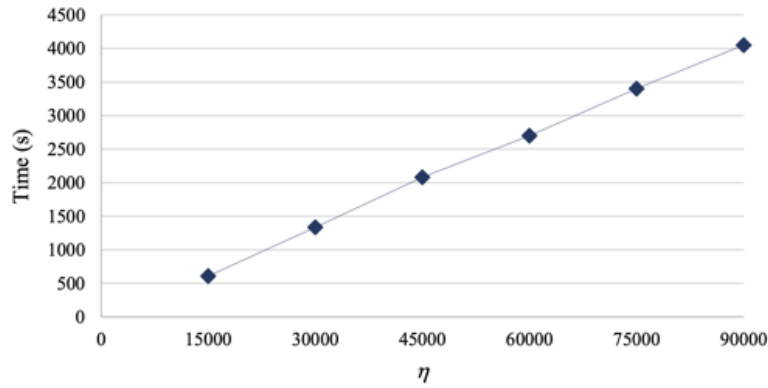


FIGURE 7.5: Correlation described by the times required for watermarking different amount of tuples.

$\eta$	Time (s)	
	Average	Proportion
15000	609.75 ( $\pm 2.70$ )	none
30000	1336.00 ( $\pm 7.35$ )	2.19
45000	2081.57 ( $\pm 3.40$ )	1.55
60000	2699.07 ( $\pm 7.90$ )	1.30
75000	3401.25 ( $\pm 14.49$ )	1.26
90000	4050.67 ( $\pm 27.41$ )	1.19

TABLE 7.13: Time required for WM embedding involving different tuples number.

number of tuples in R. This experiment was performed using  $\gamma = 5$ , the WM source WWF and the attribute “Text”, increasing the value of  $\eta$  of 15000 units each time.

For each amount of tuples, the same experiment was performed several times until reaching a standard deviation of the time required as close as zero as possible. Table 7.13 shows along with the standard deviation, the mean of the time recorded. According to column “Proportion”, which compares the average of the time required for marking the tuples with respect to the average of time required for the previous row, our approach describes a linear behavior.

The approach’s complexity will depend on the Similarity Engine, the WSD module, and the number of sentences stored in the attribute being watermarked. Despite all these factors, for the conditions given by the experimental set up to validate our work, it is recorded a linear behavior. Then, it can be established that the overall time complexity corresponds to  $O(n)$ , being reliable its application to different sizes of data stored in R.

## 7.6 Discussion

In this chapter, we proposed a watermarking technique for relational data that uses multi-word textual attributes as cover type. The embedding of the marks in our approach is performed by substitutions of synonym words in sentences, guaranteeing the semantic preservation of the data, and the total imperceptibility of the watermark. Despite multiple attributes can be considered for each tuple, when paragraphs are stored, the selection of one word per sentence allows the increment of the watermark capacity with respect to previous techniques, and with it, its robustness.

This technique works in combination with a WSD module, a semantic similarity engine, and one or several knowledge sources, linking its complexity and precision to the behavior of those external elements. For the experimental validation we used WordNet as knowledge source and the Lesk algorithm for WSD. The results show that how our technique guarantees the watermark embedding, its detection, and robustness against subset attacks and random synonym substitution attacks, making it a valuable tool for ownership protection, and the data integrity validation. For the case of random synonym substitution attacks, which constitute a serious threat for techniques focused on watermarking textual documents, the combination of the relational data structure and the multi-word textual data type guarantees the watermark persistence independently the attack's degree. The preservation of the semantic was defined as the main priority for this approach, adding as a feature the tolerance to the watermark embedding in a way that, to the best of our knowledge, no other technique did before.





## Chapter 8

---

# CONCLUSION AND FUTURE WORK

In this thesis we have studied the problem of string analysis for software verification, which is also relevant from a security point of view. Indeed, more precise string approximations better detect program errors that may lead to significant exploitations (voluntary or otherwise). We adopted the Abstract Interpretation theory both to design novel string approximations and to refine the existing ones.

It emerges that the analysis of programs that manipulate strings is a complex problem, far from being solved. One of the main reasons for its difficulty lies in the fact that programming languages represent strings, and their operations, in different ways. This is clearly an obstacle since it enormously increases the number of possible scenarios to be tackled, making the design of effective analyses that would not be language-specific quite difficult.

In this context, the semantic-based approach used in this thesis has led to interesting results detailed throughout this dissertation. We investigated different aspects of the string analysis problem, among which the correct management of character arrays in `C` programs and the enhancement of existing string analyses. This brought us to the conclusion that semantic-based approaches can be effectively implemented and fully incorporated in the software development life-cycle. One of the advantages of the Abstract Interpretation approach is that the degree of freedom in the choice of the abstract domain allows to adequately address the trade-off between efficiency and accuracy of the analysis.

Among the countless scenarios that we could explore in depth, we consider the following logical developments of this thesis, as they are an expansion of the work done so far. In fact, in the short term, we aim to:

- Further enhance the effectiveness of the M-String domain by combining it by reduced product with other either numerical or symbolic domains.
- Design a tailoring method to select what domain combination better suits the string analysis of a given program based on a pre-process of the operations it applies on string and character variables.

- Investigate forward completeness [80], the property that guarantees that no loss of precision occurs during the output abstraction process of a given operation. We aim to integrate the two completeness methodologies within an industrial JavaScript static analyser to deeply evaluate their actual cost and overall impact. Moreover, we aim to define a precision measure between string abstractions.
- Investigate suitable combinations of relational string domains with non-relational ones and/or with complementary numerical domains. We also aim at optimizing the efficiency of the implementation by taking full benefit from the results discussed in [126] for numerical domains by representing relations as multi-valued maps.
- Investigate the possibility to integrate the string domains we have proposed in this thesis, together with the techniques we have used to combine them, into LiSA.<sup>1</sup> LiSA is a library for static analysis, aiming to ease the creation and implementation of static analysers based on the Abstract Interpretation theory.
- Design a module of semantic preservation for numerical and textual data, considering the elements of watermarking techniques and approximate query processing.

---

<sup>1</sup>LiSA is available at <https://github.com/UniVE-SSV/lisa>.

## Appendix A

---

### UNIFICATION ALGORITHM

We recall the unification algorithm of [57].

The first step of the algorithm is checking the compatibility of the two input segmentations to verify that they do have common lower and upper bounds. Then, the unification proceeds recursively from left to right and maintains the invariant that the left part is already unified. Let  $I_l$  ( $I_r$  resp.) denote the left (right resp.) neutral element.

- 1  $\bar{b}[?_1] \bar{p}_1 \bar{b}'_1[?_1] \dots$  and  $\bar{b}[?_2] \bar{p}_2 \bar{b}'_2[?_2] \dots$  have some lower bounds and so keep the first segments as they are and go on with  $\bar{b}'_1[?_1] \dots$  and  $\bar{b}'_2[?_2] \dots$
- 2 In case  $(\bar{b} \cup \bar{b}_1)[?_1] \bar{p}_1 \bar{b}'_1[?_1] \dots$  and  $\bar{b}[?_2] \bar{p}_2 \bar{b}'_2[?_2] \dots$  with  $\bar{b}_1 \neq \emptyset$  and  $\bar{b} \cap \bar{b}_1 = \emptyset$ , let  $\bar{b}_1^*$  be the set of expressions in  $\bar{b}_1$  appearing in the second segmentation blocks  $\bar{b}'_2, \dots$ 
  - 2.1 If  $\bar{b}_1^*$  is empty, then go on with  $\bar{b}[?_1] \bar{p}_1 \bar{b}'_1[?_1] \dots$  and  $\bar{b}[?_2] \bar{p}_2 \bar{b}'_2[?_2] \dots$  following case 1.
  - 2.2 Otherwise, go on with  $\bar{b}[?_1] I_l \bar{b}_1? \bar{p}_1 \bar{b}'_1[?_1] \dots$  and  $\bar{b}[?_2] \bar{p}_2 \bar{b}'_2[?_2] \dots$  as in case 1.
- 3 The symmetrical case is similar.
- 4 In case  $(\bar{b} \cup \bar{b}_1)[?_1] \bar{p}_1 \bar{b}'_1[?_1] \dots$  and  $(\bar{b} \cup \bar{b}_2)[?_2] \bar{p}_2 \bar{b}'_2[?_2] \dots$  with  $\bar{b}_1, \bar{b}_2 \neq \emptyset$  and  $\bar{b} \cap \bar{b}_1 = \bar{b} \cap \bar{b}_2 = \emptyset$ , let  $\bar{b}_1^*$  (resp.  $\bar{b}_2^*$ ) be the set of expressions in  $\bar{b}_1$  (resp.  $\bar{b}_2$ ) appearing in the second (resp. first) segmentation blocks  $\bar{b}'_2, \dots$  ( $\bar{b}'_1, \dots$ ).
  - 4.1 If  $\bar{b}_1^*$  and  $\bar{b}_2^*$  are both empty, go on with  $\bar{b}[?_1] \bar{p}_1 \bar{b}'_1[?_1] \dots$  and  $\bar{b}[?_2] \bar{p}_2 \bar{b}'_2[?_2] \dots$  as in case 1.
  - 4.2 Else if  $\bar{b}_1^*$  is empty (so that  $\bar{b}_2^*$  is not empty) then go on with  $\bar{b}[?_1] \bar{p}_1 \bar{b}'_1[?_1] \dots$  and  $\bar{b}[?_2] I_r \bar{b}_2? \bar{p}_2 \bar{b}'_2[?_2] \dots$  (where  $I_r$  is the right neutral element).
  - 4.3 The symmetrical case is similar.
  - 4.4 Finally if  $\bar{b}_1^*$  and  $\bar{b}_2^*$  are both non-empty, then go on with  $\bar{b}[?_1] I_l \bar{b}_1? \bar{p}_1 \bar{b}'_1[?_1] \dots$  and  $\bar{b}[?_2] I_r \bar{b}_2? \bar{p}_2 \bar{b}'_2[?_2] \dots$  as in case 1.

- 5** In case  $\bar{\mathbf{b}}_1[?_1] \bar{\mathbf{p}}_1 \bar{\mathbf{b}}_1'[?_1'] \dots$  and  $\bar{\mathbf{b}}_2[?_2] \bar{\mathbf{p}}_2 \bar{\mathbf{b}}_2'[?_2'] \dots$  with  $\bar{\mathbf{b}}_1 \cap \bar{\mathbf{b}}_2 = \emptyset$ , we cannot be on the first left segment block so we have on the left  $\bar{\mathbf{b}}_0[?_0] \bar{\mathbf{p}}_0 \bar{\mathbf{b}}_1[?_1] \bar{\mathbf{p}}_1 \bar{\mathbf{b}}_1'[?_1'] \dots$  and  $\bar{\mathbf{b}}_0'[?_0'] \bar{\mathbf{p}}_0' \bar{\mathbf{b}}_2[?_2] \bar{\mathbf{p}}_2 \bar{\mathbf{b}}_2'[?_2'] \dots$  and go on by merging these consecutive blocks  $\bar{\mathbf{b}}_0[?_0] \bar{\mathbf{p}}_0 \sqcup \bar{\mathbf{p}}_1 \bar{\mathbf{b}}_1[?_1 \wedge ?_1'] \dots$  and  $\bar{\mathbf{b}}_0'[?_0'] \bar{\mathbf{p}}_0' \sqcup \bar{\mathbf{p}}_2 \bar{\mathbf{b}}_2[?_2 \wedge ?_2'] \dots$ .
- 6** Finally, at the end either we are left with the right limits that have both been checked to be equal or else we have  $\bar{\mathbf{b}}_1[?_1] \bar{\mathbf{p}}_1 \bar{\mathbf{b}}_1'[?_1']$  and  $\bar{\mathbf{b}}_2[?_2]$  with  $\bar{\mathbf{b}}_1' = \bar{\mathbf{b}}_2$ . Because we have maintained the invariant that  $\bar{\mathbf{b}}_1$  is always equal to  $\bar{\mathbf{b}}_2$  in the concrete (so necessarily  $[?_1] = ?$  since then  $\bar{\mathbf{b}}_1 = \bar{\mathbf{b}}_2 = \bar{\mathbf{b}}_1'$ ), and so we end up with  $(\bar{\mathbf{b}}_1 \cup \bar{\mathbf{b}}_1' \cup \bar{\mathbf{b}}_2)[?_1]$  and  $(\bar{\mathbf{b}}_1 \cup \bar{\mathbf{b}}_1' \cup \bar{\mathbf{b}}_2)[?_2]$ .

## Appendix B

---

### STRING ABSTRACT DOMAINS

We formally define the String Length complete lattice introduced in Section 4.4.1 and its abstract semantics. The soundness of the String Length domain abstract semantics is also proved. Moreover, we recall the Character Inclusion (cf. Section 4.4.2) and the Prefix (cf. Section 4.4.3) complete lattices and their abstract semantics [49].

#### B.1 String Length

The String Length abstract domain  $\overline{\mathbf{SL}}$  (cf. Section 4.4.1) is the complete lattice [11]:

$$(\overline{\mathbf{SL}}, \sqsubseteq_{\overline{\mathbf{SL}}}, \perp_{\overline{\mathbf{SL}}}, \top_{\overline{\mathbf{SL}}}, \sqcap_{\overline{\mathbf{SL}}}, \sqcup_{\overline{\mathbf{SL}}})$$

where:

- $\overline{\mathbf{SL}} \triangleq \{[m, M] : m \in \mathbb{N} \wedge M \in \mathbb{N} \cup \{\infty\} \wedge m \leq M\} \cup \{\perp_{\overline{\mathbf{SL}}}\}$
- Let  $[m_1, M_1]$  and  $[m_2, M_2]$  be two abstract values in the String Length domain then:

$$[m_1, M_1] \sqsubseteq_{\overline{\mathbf{SL}}} [m_2, M_2] \Leftrightarrow [m_1, M_1] = \perp_{\overline{\mathbf{SL}}} \vee (m_2 \leq m_1 \wedge M_1 \leq M_2)$$

- $\perp_{\overline{\mathbf{SL}}}$  is a special element denoting the bottom element of the considered lattice.
- $\top_{\overline{\mathbf{SL}}}$  denotes the top element of the considered lattice, and:

$$\top_{\overline{\mathbf{SL}}} = [0, \infty]$$

- $\sqcap_{\overline{\mathbf{SL}}}$  represents the meet operator, that defines the greatest lower bound between two string approximations in the String Length abstract domain, such that:

$$[m_1, M_1] \sqcap_{\overline{\mathbf{SL}}} [m_2, M_2] = \begin{cases} [\max(\{m_1, m_2\}), \min(\{M_1, M_2\})] & \text{if } \max(\{m_1, m_2\}) \leq \min(\{M_1, M_2\}) \\ \perp_{\overline{\mathbf{SL}}} & \text{otherwise} \end{cases}$$

- $\sqcup_{\overline{\mathbf{SL}}}$  represents the join operator, that defines the least upper bound between two string approximations in the String Length abstract domain, such that:

$$[m_1, M_1] \sqcup_{\overline{\mathbf{SL}}} [m_2, M_2] = [\min(\{m_1, m_2\}), \max(\{M_1, M_2\})]$$

### Abstraction

Let  $X$  be a set of concrete string values. The abstraction function on the String Length abstract domain  $\alpha_{\overline{\mathbf{SL}}}$  maps  $X$  to  $\perp_{\overline{\mathbf{SL}}}$  in the case in which  $X$  is equal to the empty set, otherwise to the interval  $[\min(L_X), \max(L_X)]$ , where  $L_X = \{|\sigma| : \sigma \in X\}$ . Notice that, if  $X$  contains only one string element, i.e., if  $X = \{\sigma\}$ , then  $\alpha_{\overline{\mathbf{SL}}}(X) = [|\sigma|, |\sigma|]$ .

### Concretization

The concretization function on the String Length abstract domain  $\gamma_{\overline{\mathbf{SL}}}$  maps an abstract element to a set of strings as follows:  $\gamma_{\overline{\mathbf{SL}}}(\perp_{\overline{\mathbf{SL}}}) = \emptyset$ , while is  $\{\sigma : \sigma \in \Sigma^* \wedge m \leq |\sigma| \leq M\}$ , i.e., the set of all the possible strings whose length goes from  $m$  to  $M$ .

### Galois connection

We follow Theorem 1.1 of [49] and we prove that  $(\alpha_{\overline{\mathbf{SL}}}, \gamma_{\overline{\mathbf{SL}}})$  is a Galois Connection.

**Theorem B.1.** Let the abstraction function  $\alpha_{\overline{\mathbf{SL}}}$  be defined as:

$$\alpha_{\overline{\mathbf{SL}}} = \lambda X. \sqcap_{\overline{\mathbf{SL}}} \{[m, M] : X \subseteq \gamma_{\overline{\mathbf{SL}}}([m, M])\}$$

Then,  $(\mathcal{P}(\Sigma^*), \subseteq) \xleftrightarrow[\alpha_{\overline{\mathbf{SL}}}]{\gamma_{\overline{\mathbf{SL}}}} (\overline{\mathbf{SL}}, \sqsubseteq_{\overline{\mathbf{SL}}})$ .

*Proof.*

By Theorem 1.1 of [49], we only need to prove that  $\gamma_{\overline{\mathbf{SL}}}$  is a complete meet morphism.

Formally, we have to prove that:

$$\begin{aligned} & \gamma_{\overline{\mathbf{SL}}}\left(\sqcap_{[m, M] \in \overline{X}} [m, M]\right) = \bigcap_{[m, M] \in \overline{X}} \gamma_{\overline{\mathbf{SL}}}([m, M]) \\ & \gamma_{\overline{\mathbf{SL}}}\left(\sqcap_{[m, M] \in \overline{X}} [m, M]\right) \\ & \text{by definition of } \sqcap_{\overline{\mathbf{SL}}} \\ & = \gamma_{\overline{\mathbf{SL}}}([m', M']) \text{ where } m' = \max_{[m, M] \in \overline{X}} m \text{ and } M' = \min_{[m, M] \in \overline{X}} M \end{aligned}$$

$$\begin{aligned}
\mathfrak{S}_{\overline{\mathbf{SL}}}[\mathbf{new\ String}(\sigma)]() &= [|\sigma|, |\sigma|] \\
\mathfrak{S}_{\overline{\mathbf{SL}}}[\mathbf{concat}]([m_1, M_1], [m_2, M_2]) &= [m_1 + m_2, M_1 + M_2] \\
\mathfrak{S}_{\overline{\mathbf{SL}}}[\mathbf{substring}_b^e]([m, M]) &= \begin{cases} [n, n] & \text{if } n \in [m, M] \\ \perp_{\overline{\mathbf{SL}}} & \text{otherwise} \end{cases} \\
&\quad \text{where } n = (e - b) + 1 \\
\mathfrak{B}_{\overline{\mathbf{SL}}}[\mathbf{contains}_c]([m, M]) &= \top_{\mathbf{B}}
\end{aligned}$$

TABLE B.1:  $\overline{\mathbf{SL}}$  abstract semantics.

by definition of  $\gamma_{\overline{\mathbf{SL}}}$

$$= \{\sigma \mid \sigma \in \Sigma^* \wedge m' \leq |\sigma| \leq M'\}$$

by definition of  $\cap$

$$= \bigcap_{[m, M] \in \overline{\mathbf{X}}} \{\sigma : \sigma \in \Sigma^* \wedge m \leq |\sigma| \leq M\}$$

by definition of  $\gamma_{\overline{\mathbf{SL}}}$

$$= \bigcap_{[m, M] \in \overline{\mathbf{X}}} \gamma_{\overline{\mathbf{SL}}}([m, M])$$

□

### Widening operator

The String Length abstract domain has infinite height and does not respect the Ascending Chain Condition. Therefore it needs to be equipped with a widening operator. The widening  $\nabla_{\overline{\mathbf{SL}}}$  for the string length analysis is as follows:

$$\begin{aligned}
[m_1, M_1] \nabla_{\overline{\mathbf{SL}}} [m_2, M_2] &= [\text{if } m_2 < m_1 \text{ then } 0 \text{ else } m_1, \\
&\quad \text{if } M_2 > M_1 \text{ then } \infty \text{ else } M_1]
\end{aligned}$$

### Abstract semantics

We define the abstract semantics of the String Length domain and we prove its soundness. The semantics  $\mathfrak{S}_{\overline{\mathbf{SL}}}$  and  $\mathfrak{B}_{\overline{\mathbf{SL}}}$  are the abstract counterparts of  $\mathfrak{S}$  and  $\mathfrak{B}$  (cf. Section ??).

Table B.1 formalizes the abstract semantics of  $\overline{\mathbf{SL}}$ . Precisely,

- $\mathfrak{S}_{\overline{\mathbf{SL}}}[\mathbf{new\ String}(\sigma)]() = [|\sigma|, |\sigma|]$

Let  $\sigma$  be a sequence of characters.

The semantics  $\mathfrak{S}_{\overline{\mathbf{SL}}}$ , when applied to  $\mathbf{new\ String}(\sigma)$ , returns the interval  $[|\sigma|, |\sigma|]$ , as when a new string constant is evaluated both its minimum and its maximum length coincide and are equal to the length of  $\sigma$ .

- $\mathfrak{S}_{\overline{\mathbf{SL}}}\llbracket \text{concat} \rrbracket([m_1, M_1], [m_2, M_2]) = [m_1 + m_2, M_1 + M_2]$

Let  $[m_1, M_1], [m_2, M_2] \in \overline{\mathbf{SL}}$ .

The semantics  $\mathfrak{S}_{\overline{\mathbf{SL}}}$ , when applied to  $\text{concat}([m_1, M_1], [m_2, M_2])$ , returns the sum of the input intervals, i.e.,  $[m_1, M_1] + [m_2, M_2] = [m_1 + m_2, M_1 + M_2]$ , as the length of a string resulting from the concatenation between two other strings is the sum of the lengths of the latter, i.e.,  $|\sigma_1 \sigma_2| = |\sigma_1| + |\sigma_2|$ .

- $\mathfrak{S}_{\overline{\mathbf{SL}}}\llbracket \text{substring}_b^e \rrbracket([m, M]) = \begin{cases} [n, n] & \text{if } n \in [m, M] \\ \perp_{\overline{\mathbf{SL}}} & \text{otherwise} \end{cases}$

Let  $[m, M] \in \overline{\mathbf{SL}}$  and let  $n$  be equal to  $(e - b) + 1$ .

The semantics  $\mathfrak{S}_{\overline{\mathbf{SL}}}$ , when applied to  $\text{substring}_b^e([m, M])$ , returns the interval  $[n, n]$  if  $n$  belongs to the input interval  $[m, M]$ , otherwise it returns  $\perp_{\overline{\mathbf{SL}}}$ , as the substring of a string from the  $b$ -th to the  $e$ -th index has length equal to  $n$ .

- $\mathfrak{B}_{\overline{\mathbf{SL}}}\llbracket \text{contains}_c \rrbracket([m, M]) = \top_{\mathbf{B}}$

Let  $[m, M] \in \overline{\mathbf{SL}}$ .

The semantics  $\mathfrak{B}_{\overline{\mathbf{SL}}}$  applied to  $\text{contains}_c([m, M])$  returns  $\top_{\mathbf{B}}$ , as we can not infer any information about the content of the strings approximated by  $[m, M]$ .

**Theorem B.2.**  $\mathfrak{S}_{\overline{\mathbf{SL}}}$  and  $\mathfrak{B}_{\overline{\mathbf{SL}}}$  are sound over-approximations of  $\mathfrak{S}$  and  $\mathfrak{B}$  respectively. Formally,

$$\gamma_{\overline{\mathbf{SL}}}(\mathfrak{S}_{\overline{\mathbf{SL}}}\llbracket \text{Stm} \rrbracket([m, M])) \supseteq \{\mathfrak{S}\llbracket \text{Stm} \rrbracket(\sigma) \mid \sigma \in \gamma_{\overline{\mathbf{SL}}}([m, M])\}$$

$$\gamma_{\overline{\mathbf{SL}}}(\mathfrak{B}_{\overline{\mathbf{SL}}}\llbracket \text{Stm} \rrbracket([m, M])) \supseteq \{\mathfrak{B}\llbracket \text{Stm} \rrbracket(\sigma) \mid \sigma \in \gamma_{\overline{\mathbf{SL}}}([m, M])\}$$

*Proof.*

We prove the soundness separately for each operator.

- Consider the **new String** operator and let  $\sigma$  be a sequence of characters. We have to prove that,

$$\gamma_{\overline{\mathbf{SL}}}(\mathfrak{S}_{\overline{\mathbf{SL}}}\llbracket \text{new String}(\sigma) \rrbracket()) \supseteq \{\mathfrak{S}\llbracket \text{new String}(\sigma) \rrbracket()\}$$

This holds from the definitions of  $\mathfrak{S}_{\overline{\mathbf{SL}}}$  and of  $\gamma_{\overline{\mathbf{SL}}}$ .

- Consider the **concat** operation and let  $[m_1, M_1], [m_2, M_2] \in \overline{\mathbf{SL}}$ . We have to prove that,

$$\begin{aligned} \gamma_{\overline{\mathbf{SL}}}(\mathfrak{S}_{\overline{\mathbf{SL}}}\llbracket \text{concat} \rrbracket([m_1, M_1], [m_2, M_2])) \\ \supseteq \\ \{\mathfrak{S}\llbracket \text{concat} \rrbracket(\sigma_1, \sigma_2) \mid \sigma_1 \in \gamma_{\overline{\mathbf{SL}}}([m_1, M_1]) \wedge \sigma_2 \in \gamma_{\overline{\mathbf{SL}}}([m_2, M_2])\} \end{aligned}$$



Let  $\sigma_1\sigma_2$  be a generic element in  $\{\mathfrak{S}[\text{concat}](\sigma_1, \sigma_2) : \sigma_1 \in \gamma_{\overline{\mathbf{SL}}}([m_1, M_1]) \wedge \sigma_2 \in \gamma_{\overline{\mathbf{SL}}}([m_2, M_2])\}$ . As  $\sigma_1$  is any string whose length goes from  $m_1$  to  $M_1$  and  $\sigma_2$  is any string whose length goes from  $m_2$  to  $M_2$ , the concatenation of  $\sigma_1$  and  $\sigma_2$  is any string whose length goes from  $m_1 + m_2$  to  $M_1 + M_2$ . Therefore  $\sigma_1\sigma_2$  belongs to  $\gamma_{\overline{\mathbf{SL}}}(\mathfrak{S}_{\overline{\mathbf{SL}}}[\text{concat}]([m_1, M_1], [m_2, M_2]))$ , because  $\mathfrak{S}_{\overline{\mathbf{SL}}}[\text{concat}]([m_1, M_1], [m_2, M_2]) = [m_1 + m_2, M_1 + M_2]$ , by definition of  $\mathfrak{S}_{\overline{\mathbf{SL}}}$ . Indeed,  $\gamma_{\overline{\mathbf{SL}}}([m_1 + m_2, M_1 + M_2])$  contains all the strings whose length goes from  $m_1 + m_2$  to  $M_1 + M_2$ , by definition of  $\gamma_{\overline{\mathbf{SL}}}$ .

- Consider the **substring** operation and let  $[m, M] \in \overline{\mathbf{SL}}$ . we have to prove that,

$$\gamma_{\overline{\mathbf{SL}}}(\mathfrak{S}_{\overline{\mathbf{SL}}}[\text{substring}_b^e]([m, M])) \supseteq \{\mathfrak{S}[\text{substring}_b^e](\sigma) \mid \sigma \in \gamma_{\overline{\mathbf{SL}}}([m, M])\}$$

Let  $\sigma_b \dots \sigma_e$  be a generic element in  $\{\mathfrak{S}[\text{substring}_b^e](\sigma) \mid \sigma \in \gamma_{\overline{\mathbf{SL}}}([m, M])\}$ . As  $\sigma$  is any string whose length goes from  $m$  to  $M$ , its substring from the  $b$ -th to the  $e$ -th index, if it exists, i.e., if  $n \leq |\sigma|$ , has length equal to  $n$ . Therefore  $\sigma_b \dots \sigma_e$  belongs to  $\gamma_{\overline{\mathbf{SL}}}(\mathfrak{S}_{\overline{\mathbf{SL}}}[\text{substring}_b^e]([m, M]))$ , because  $\mathfrak{S}_{\overline{\mathbf{SL}}}[\text{substring}_b^e]([m, M]) = [n, n]$  if  $n \in [m, M]$ , by definition of  $\mathfrak{S}_{\overline{\mathbf{SL}}}$ . Indeed,  $\gamma_{\overline{\mathbf{SL}}}([n, n])$  contains all the strings of length  $n$  and  $\gamma_{\overline{\mathbf{SL}}}(\perp_{\overline{\mathbf{SL}}})$  is the emptyset, by definition of  $\gamma_{\overline{\mathbf{SL}}}$ .

- Consider the **contains** operator and let  $[m, M] \in \overline{\mathbf{SL}}$ . We have to prove that,

$$\gamma_{\overline{\mathbf{SL}}}(\mathfrak{B}_{\overline{\mathbf{SL}}}[\text{contains}_c]([m, M])) \supseteq \{\mathfrak{B}[\text{contains}_c](\sigma) \mid \sigma \in \gamma_{\overline{\mathbf{SL}}}([m, M])\}$$

This holds as the abstract semantics returns the top element  $\top_{\mathbf{B}}$  that approximates any possible result of the concrete semantics.

□

## B.2 Character Inclusion

The Character Inclusion abstract domain  $\overline{\mathbf{CI}}$  is the complete lattice [49]:

$$(\overline{\mathbf{CI}}, \sqsubseteq_{\overline{\mathbf{CI}}}, \perp_{\overline{\mathbf{CI}}}, \top_{\overline{\mathbf{CI}}}, \sqcap_{\overline{\mathbf{CI}}}, \sqcup_{\overline{\mathbf{CI}}})$$

where:

- $\overline{\mathbf{CI}} \triangleq \{(C, MC) : C, MC \in \mathcal{P}(\Sigma) \wedge C \subseteq MC\} \cup \{\perp_{\overline{\mathbf{CI}}}\}$

Where  $\Sigma$  denotes the alphabet.

- Let  $(C_1, MC_1)$  and  $(C_2, MC_2)$  be two abstract values in the Character Inclusion domain then:

$$(C_1, MC_1) \sqsubseteq_{\overline{\mathbf{CI}}} (C_2, MC_2) \Leftrightarrow (C_1, MC_1) = \perp_{\overline{\mathbf{CI}}} \vee (C_1 \supseteq C_2 \wedge MC_1 \subseteq MC_2)$$

- $\perp_{\overline{\text{CI}}}$  denotes the bottom element of the considered lattice, and:

$$\perp_{\overline{\text{CI}}} = \{(C, MC) \mid C \not\subseteq MC\}$$

- $\top_{\overline{\text{CI}}}$  denotes the top element of the considered lattice, and:

$$\top_{\overline{\text{CI}}} = (\emptyset, \Sigma)$$

- $\sqcap_{\overline{\text{CI}}}$  represents the meet operator, that defines the greatest lower bound between two string approximations in the Character Inclusion abstract domain, such that:

$$\begin{aligned} \sqcap_{\overline{\text{CI}}}((C_1, MC_1), (C_2, MC_2)) = \\ \begin{cases} (C_1 \cup C_2, MC_1 \cap MC_2) & \text{if } C_1 \subseteq MC_2 \text{ and } C_2 \subseteq MC_1 \\ \perp_{\overline{\text{CI}}} & \text{otherwise} \end{cases} \end{aligned}$$

- $\sqcup_{\overline{\text{CI}}}$  represents the join operator, that defines the least upper bound between two string approximations in the Character Inclusion abstract domain, such that:

$$\sqcup_{\overline{\text{CI}}}((C_1, MC_1), (C_2, MC_2)) = (C_1 \cap C_2, MC_1 \cup MC_2).$$

### Abstraction

Let  $X$  be a set of concrete string values. The abstraction function on the Character Inclusion abstract domain  $\alpha_{\overline{\text{CI}}}$  maps  $X$  to  $\perp_{\overline{\text{CI}}}$  in the case in which  $X$  is equal to the empty set, otherwise to the pair of sets  $(\bigcap H_X, \bigcup H_X)$  where  $H_X = \{\text{char}(\sigma) : \sigma \in X\}$ . Notice that, if  $X$  contains just one string element, then  $\alpha_{\overline{\text{CI}}}(X) = (H_X, H_X)$ .

### Concretization

The concretization function on the Character Inclusion abstract domain  $\gamma_{\overline{\text{CI}}}$  maps an abstract element to a set of strings as follows:  $\gamma_{\overline{\text{CI}}}(\perp_{\overline{\text{CI}}}) = \emptyset$ , while is  $\{\sigma : \sigma \in \Sigma^*, C \subseteq \text{char}(\sigma) \subseteq MC\}$ , i.e., the set of all the possible strings certainly containing the elements in  $C$  and possibly containing the elements in  $MC$ .

### Galois connection

The functions  $\gamma_{\overline{\text{CI}}}$  and  $\alpha_{\overline{\text{CI}}}$  form a Galois Connection (Theorem 4.2 of [49]).

### Widening operator

The Character Inclusion abstract domain has finite height. Therefore its widening operator coincides with its least upper bound operator.

$$\begin{aligned}
\mathfrak{S}_{\overline{\text{CI}}}[\text{new String}(\sigma)]() &= (\text{char}(\sigma), \text{char}(\sigma)) \\
\mathfrak{S}_{\overline{\text{CI}}}[\text{concat}]((C_1, \text{MC}_1), (C_2, \text{MC}_2)) &= (C_1 \cup C_2, \text{MC}_1 \cup \text{MC}_2) \\
\mathfrak{S}_{\overline{\text{CI}}}[\text{substring}_b^e]((C, \text{MC})) &= (\emptyset, \text{MC}) \\
\mathfrak{B}_{\overline{\text{CI}}}[\text{contains}_c]((C, \text{MC})) &= \begin{cases} \text{true} & \text{if } c \in C \\ \text{false} & \text{if } c \notin \text{MC} \\ 0.5\text{ex} \uparrow_{\mathbf{B}} & \text{otherwise} \end{cases}
\end{aligned}$$

TABLE B.2:  $\overline{\text{CI}}$  abstract semantics.

### Abstract semantics

We recall the Character Inclusion domain abstract semantics. Table B.2 formalizes the abstract semantics of  $\overline{\text{CI}}$ . Precisely,

- $\mathfrak{S}_{\overline{\text{CI}}}[\text{new String}(\sigma)]() = (\text{char}(\sigma), \text{char}(\sigma))$

Let  $\sigma$  be a sequence of characters.

The semantics of  $\mathfrak{S}_{\overline{\text{CI}}}$ , when applied to  $\text{new String}(\sigma)$  returns the pair  $(\text{char}(\sigma), \text{char}(\sigma))$ , as when a new constant string is evaluated both the sets of characters it certainly and probably contains are those appearing in  $\sigma$ .

- $\mathfrak{S}_{\overline{\text{CI}}}[\text{concat}]((C_1, \text{MC}_1), (C_2, \text{MC}_2)) = (C_1 \cup C_2, \text{MC}_1 \cup \text{MC}_2)$

Let  $(C_1, \text{MC}_1), (C_2, \text{MC}_2) \in \overline{\text{CI}}$ .

The semantics  $\mathfrak{S}_{\overline{\text{CI}}}$ , when applied to  $\text{concat}((C_1, \text{MC}_1), (C_2, \text{MC}_2))$ , returns the pair  $(C_1 \cup C_2, \text{MC}_1 \cup \text{MC}_2)$ , as the characters which occur in a string resulting from the concatenation between two other strings are those appearing in the first string plus those appearing in the second string, i.e.,  $\text{char}(\sigma_1\sigma_2) = \text{char}(\sigma_1) \cup \text{char}(\sigma_2)$ .

- $\mathfrak{S}_{\overline{\text{CI}}}[\text{substring}_b^e]((C, \text{MC})) = (\emptyset, \text{MC})$

Let  $(C, \text{MC}) \in \overline{\text{CI}}$ .

The semantics  $\mathfrak{S}_{\overline{\text{CI}}}$ , when applied to  $\text{substring}_b^e((C, \text{MC}))$  returns the pair  $(\emptyset, \text{MC})$ , as we do not have any information about the position of the characters in  $C$ .

- $\mathfrak{B}_{\overline{\text{CI}}}[\text{contains}_c]((C, \text{MC})) = \begin{cases} \text{true} & \text{if } c \in C \\ \text{false} & \text{if } c \notin \text{MC} \\ \uparrow_{\mathbf{B}} & \text{otherwise} \end{cases}$

Let  $(C, \text{MC}) \in \overline{\text{CI}}$ .

The semantics  $\mathfrak{B}_{\overline{\text{CI}}}$  applied to  $\text{substring}_c((C, \text{MC}))$  returns **true** if the set  $C$  contains the character  $c$ , **false** if the set  $\text{MC}$  does not contain the character  $c$ , otherwise

$\top_{\mathbf{B}}$ , as the characters which occur in  $C$  are those certainly appearing in all the strings approximated by  $(C, MC)$ .

The Character Inclusion domain abstract semantics is a sound over-approximation of the string domain concrete semantics (Theorem 4.3 of [49]).

### B.3 Prefix

A prefix is a sequence of characters followed by an asterisk \*. Since the asterisk at the end of the literal part is always present, it is not included in the domain. The Prefix abstract domain  $\overline{\mathbf{PR}}$  is the complete lattice [49]:

$$(\overline{\mathbf{PR}}, \sqsubseteq_{\overline{\mathbf{PR}}}, \perp_{\overline{\mathbf{PR}}}, \top_{\overline{\mathbf{PR}}}, \sqcap_{\overline{\mathbf{PR}}}, \sqcup_{\overline{\mathbf{PR}}})$$

where:

- $\overline{\mathbf{PR}} \triangleq \Sigma^* \cup \{\perp_{\overline{\mathbf{PR}}}\}$

- Let  $\bar{p}_1$  and  $\bar{p}_2$  be two abstract values in the Prefix domain then:

$$\bar{p}_1 \sqsubseteq_{\overline{\mathbf{PR}}} \bar{p}_2 \Leftrightarrow \bar{p}_1 = \perp_{\overline{\mathbf{PR}}} \vee (\text{len}(\bar{p}_2) \leq \text{len}(\bar{p}_1) \wedge (\forall i \in [0, \text{len}(\bar{p}_2) - 2]: \bar{p}_2[i] = \bar{p}_1[i]))$$

Notice that  $\text{len}(\bar{p})$  corresponds to the number of characters spelled out in the prefix.

- $\perp_{\overline{\mathbf{PR}}}$  is a special element denoting the bottom element of the considered lattice.
- $\top_{\overline{\mathbf{PR}}}$  denotes the top element of the considered lattice and it corresponds to the empty prefix value, such that:

$$\top_{\overline{\mathbf{PR}}} = *$$

- $\sqcap_{\overline{\mathbf{PR}}}$  represents the meet operator, that defines the greatest lower bound between two string approximations in the Prefix abstract domain, such that:

$$\sqcap_{\overline{\mathbf{PR}}}(\bar{p}_1, \bar{p}_2) = \begin{cases} \bar{p}_1 & \text{if } \bar{p}_1 \sqsubseteq_{\overline{\mathbf{PR}}} \bar{p}_2 \\ \bar{p}_2 & \text{if } \bar{p}_2 \sqsubseteq_{\overline{\mathbf{PR}}} \bar{p}_1 \\ \perp_{\overline{\mathbf{PR}}} & \text{otherwise} \end{cases}$$

- $\sqcup_{\overline{\mathbf{PR}}}$  represents the join operator, that defines the least upper bound between two string approximations in the Prefix abstract domain, such that:

$$\sqcup_{\overline{\mathbf{PR}}}(\bar{p}_1, \bar{p}_2) = \begin{cases} \top_{\overline{\mathbf{PR}}} & \text{if } \bar{p}_1[0] \neq \bar{p}_2[0] \\ \text{lcp}(\bar{p}_1, \bar{p}_2) & \text{otherwise} \end{cases}$$

where  $\text{lcp}(\bar{p}_1, \bar{p}_2)$  denotes the longest common prefix of  $\bar{p}_1$  and  $\bar{p}_2$ .

$$\begin{aligned}
\mathfrak{S}_{\overline{\mathbf{PR}}}\llbracket \text{new String}(\sigma) \rrbracket() &= \sigma \\
\mathfrak{S}_{\overline{\mathbf{PR}}}\llbracket \text{concat} \rrbracket(\bar{p}_1, \bar{p}_2) &= \bar{p}_1 \\
\mathfrak{S}_{\overline{\mathbf{PR}}}\llbracket \text{substring}_b^e \rrbracket(\bar{p}) &= \begin{cases} \bar{p}[b\dots e] & \text{if } e < \text{len}(\bar{p}) \\ \bar{p}[b\dots(\text{len}(\bar{p})-1)] & \text{if } e \geq \text{len}(\bar{p}) \wedge b < \text{len}(\bar{p}) \\ \top_{\overline{\mathbf{PR}}} & \text{otherwise} \end{cases} \\
\mathfrak{B}_{\overline{\mathbf{PR}}}\llbracket \text{contains}_c \rrbracket(\bar{p}) &= \begin{cases} \text{true} & \text{if } c \in \text{char}(\bar{p}) \\ \top_{\mathbf{B}} & \text{otherwise} \end{cases}
\end{aligned}$$

TABLE B.3:  $\overline{\mathbf{PR}}$  abstract semantics.

### Abstraction

Let  $X$  be a set of concrete string values. The abstraction function on the Prefix abstract domain  $\alpha_{\overline{\mathbf{PR}}}$  maps  $X$  to  $\perp_{\overline{\mathbf{PR}}}$  in the case in which  $X$  is equal to the empty set, otherwise to the longest common prefix of all the strings in  $X$  (which is followed by  $*$ ). Notice that, if  $X$  contains just one string concrete value, its abstract value will be equal to any possible sequence of characters having that string as prefix.

### Concretization

The concretization function on the Prefix abstract domain  $\gamma_{\overline{\mathbf{PR}}}$  maps an abstract element to a set of strings as follows:  $\gamma_{\overline{\mathbf{PR}}}(\perp_{\overline{\mathbf{PR}}}) = \emptyset$ , while is  $\{\sigma : \sigma \in \Sigma^*, \text{len}(\sigma) \geq \text{len}(\bar{p}) \wedge \forall i \in [0, \text{len}(\bar{p}) - 1] : \sigma[i] = \bar{p}[i]\}$ , i.e., the set of all the possible sequences of characters sharing the same prefix.

### Galois connection

The functions  $\gamma_{\overline{\mathbf{PR}}}$  and  $\alpha_{\overline{\mathbf{PR}}}$  form a Galois Connection (Theorem 4.7 of [49]).

### Widening operator

The Prefix abstract domain has infinite height, but it respects the Ascending Chain Condition. Therefore there is no need to define a widening operator to ensure the convergence of the analysis.

### Abstract semantics

We recall the Prefix domain abstract semantics.

Table B.3 formalizes the abstract semantics of  $\overline{\mathbf{PR}}$ . Precisely,

- $\mathfrak{S}_{\overline{\mathbf{PR}}}\llbracket \text{new String}(\sigma) \rrbracket() = \sigma$

Let  $\sigma$  be a sequence of characters.

The semantics  $\mathfrak{S}_{\overline{\mathbf{PR}}}$ , when applied to `new String( $\sigma$ )`, returns the  $\sigma$ , as when a new string constant is evaluated its most precise prefix corresponds to the string itself.

- $\mathfrak{S}_{\overline{\mathbf{PR}}}[\text{concat}](\bar{p}_1, \bar{p}_2) = \bar{p}_1$

Let  $\bar{p}_1, \bar{p}_2 \in \overline{\mathbf{PR}}$ .

The semantics  $\mathfrak{S}_{\overline{\mathbf{PR}}}$ , when applied to `concat( $\bar{p}_1, \bar{p}_2$ )`, returns  $\bar{p}_1$ , as a string resulting from the concatenation of two other strings will always begin with the first string involved in the concatenation.

- $\mathfrak{S}_{\overline{\mathbf{PR}}}[\text{substring}_b^e](\bar{p}) = \begin{cases} \bar{p}[b\dots e] & \text{if } e < \text{len}(\bar{p}) \\ \bar{p}[b\dots(\text{len}(\bar{p}) - 1)] & \text{if } e \geq \text{len}(\bar{p}) \wedge b < \text{len}(\bar{p}) \\ \top_{\overline{\mathbf{PR}}} & \text{otherwise} \end{cases}$

Let  $\bar{p} \in \overline{\mathbf{PR}}$ .

The semantics  $\mathfrak{S}_{\overline{\mathbf{PR}}}$ , when applied to `substring $_b^e$ ( $\bar{p}$ )`, returns the subprefix from the index  $b$  to the index  $e$  if  $e < \text{len}(\bar{p})$ , as the substring is completely contained in  $\bar{p}$ ; the subprefix from the index  $b$  to the index  $\text{len}(\bar{p})$  minus 1 if  $e \geq \text{len}(\bar{p})$  and  $b < \text{len}(\bar{p})$ , as the substring is not entirely contained in  $\bar{p}$ ; otherwise it returns  $\top_{\overline{\mathbf{PR}}}$ .

- $\mathfrak{B}_{\overline{\mathbf{PR}}}[\text{contains}_c](\bar{p}) = \begin{cases} \text{true} & \text{if } c \in \text{char}(\bar{p}) \\ \top_{\mathbf{B}} & \text{otherwise} \end{cases}$

Let  $\bar{p} \in \overline{\mathbf{PR}}$ .

The semantics  $\mathfrak{B}_{\overline{\mathbf{PR}}}$  applied to `contains $_c$ ( $\bar{p}$ )` returns `true` if the prefix contain the character  $c$ , otherwise  $\top_{\mathbf{B}}$ , as we can not infer any information about the content of the strings approximated by  $\bar{p}$  after the prefix itself.

The Prefix domain abstract semantics is a sound over-approximation of the string domain concrete semantics (Theorem 4.5 of [49]).

## Appendix C

---

# RELATIONAL STRING ABSTRACT DOMAINS

### C.1 Abstract Semantics of Len

Let  $\llbracket \text{st} \rrbracket^{\text{len}} : \text{Len} \rightarrow \text{Len}$  be the function that given an input abstract memory returns an abstract memory containing the new string length relations introduced by  $\text{st}$ .

- **Assignment abstract semantics**

Consider the assignment  $x = \mathbf{s}$  such that  $x \in X_{\text{str}}$  and  $\mathbf{s} \in \text{SE}$ . Formally, its abstract semantics, i.e.,  $\llbracket x = \mathbf{s} \rrbracket^{\text{len}} \mathcal{L}$ , is given by the following steps:

- **[remove]**

$$\mathcal{L}_r = \begin{cases} \mathcal{L} \setminus \{w \preceq_{\text{len}} z \mid w = x\} & \text{if } x \text{ appears at the top-level of } \mathbf{s} \\ \mathcal{L} \setminus \{w \preceq_{\text{len}} z \mid w = x \vee z = x\} & \text{otherwise} \end{cases}$$

- **[add]**

$$\mathcal{L}_a = \mathcal{L}_r \cup \{y \preceq_{\text{len}} x \mid y \in \text{extr}(\mathbf{s})\}$$

where  $\text{extr}$  is the extraction function defined in Section 6.3.2.

- **[closure]**

$$\llbracket x = \mathbf{s} \rrbracket^{\text{len}} \mathcal{L} = \text{Clos}(\mathcal{L}_a)$$

- **Boolean expressions abstract semantics**

- *contains abstract semantics*

$$\llbracket \text{contains}(x, \mathbf{s}) \rrbracket^{\text{len}} \mathcal{L} = \text{Clos}(\mathcal{L} \cup \{y \preceq_{\text{len}} x \mid y \in \text{extr}(\mathbf{s})\})$$

- *equality abstract semantics*

$$\llbracket x == \mathbf{s} \rrbracket^{\text{len}} \mathcal{L} = \begin{cases} \text{Clos}(\mathcal{L} \cup \{y \preceq_{\text{len}} x, x \preceq_{\text{len}} y\}) & \text{if } \mathbf{s} = y \in X_{\text{str}} \\ \text{Clos}(\mathcal{L} \cup \{y \preceq_{\text{len}} x \mid y \in \text{extr}(\mathbf{s})\}) & \text{otherwise} \end{cases}$$

The same applies to  $\llbracket \mathbf{s} == x \rrbracket^{\text{len}} \mathcal{L}$ .

- *conjunction abstract semantics*

$$\llbracket e_1 \ \&\& \ e_2 \rrbracket^{\text{len}} \mathcal{L} = \llbracket e_1 \rrbracket^{\text{len}} \mathcal{L} \sqcup_{|\text{len}} \llbracket e_2 \rrbracket^{\text{len}} \mathcal{L}$$

- *disjunction abstract semantics*

$$\llbracket e_1 \ || \ e_2 \rrbracket^{\text{len}} \mathcal{L} = \llbracket e_1 \rrbracket^{\text{len}} \mathcal{L} \cap_{|\text{len}} \llbracket e_2 \rrbracket^{\text{len}} \mathcal{L}$$

**Theorem C.1.** The abstract semantics of `Len` is sound. Indeed, it holds that:

$$\forall M \in \mathcal{P}(M). \alpha_{|\text{len}}(\llbracket x = \mathbf{s} \rrbracket M) \sqsubseteq_{|\text{len}} (\llbracket x = \mathbf{s} \rrbracket^{\text{len}} \alpha_{|\text{len}}(M))$$

and

$$\forall M \in \mathcal{P}(M). \alpha_{|\text{len}}(\llbracket \mathbf{b} \rrbracket M) \sqsubseteq_{|\text{len}} (\llbracket \mathbf{b} \rrbracket^{\text{len}} \alpha_{|\text{len}}(M))$$

*Proof.*

The proof is analogous to the one provided for Theorem C.3. □

## C.2 Abstract Semantics of Char

Let  $\llbracket \mathbf{st} \rrbracket^{\text{char}} : \text{Char} \rightarrow \text{Char}$  be the function that given an input abstract memory returns an abstract memory containing the new character inclusion relations introduced by `st`.

- **Assignment abstract semantics**

Consider the assignment  $x = \mathbf{s}$  such that  $x \in X_{\text{str}}$  and  $\mathbf{s} \in \text{SE}$ . Formally, its abstract semantics, i.e.,  $\llbracket x = \mathbf{s} \rrbracket^{\text{char}} \mathcal{C}$  is given by the following steps:

- **[remove]**

$$\mathcal{C}_r = \begin{cases} \mathcal{C} \setminus \{w \preceq_{\text{char}} z \mid w = x\} & \text{if } x \text{ appears at the top-level of } \mathbf{s} \\ \mathcal{C} \setminus \{w \preceq_{\text{char}} z \mid w = x \vee z = x\} & \text{otherwise} \end{cases}$$

- **[add]**

$$\mathcal{C}_a = \mathcal{C}_r \cup \{y \preceq_{\text{char}} x \mid y \in \text{extr}(\mathbf{s})\}$$

where `extr` is the extraction function defined in Section 6.3.2.

- **[closure]**

$$\llbracket x = \mathbf{s} \rrbracket^{\text{char}} \mathcal{C} = \text{Clos}(\mathcal{C}_a)$$



- **Boolean expressions abstract semantics**

- *contains abstract semantics*

$$\llbracket \text{contains}(x, s) \rrbracket^{\text{char}} \mathcal{C} = \text{Clos}(\mathcal{C} \cup \{y \preceq_{\text{char}} x \mid y \in \text{extr}(s)\})$$

- *equality abstract semantics*

$$\llbracket x == s \rrbracket^{\text{char}} \mathcal{C} = \begin{cases} \text{Clos}(\mathcal{C} \cup \{y \preceq_{\text{char}} x, x \preceq_{\text{char}} y\}) & \text{if } s = y \in X_{\text{str}} \\ \text{Clos}(\mathcal{C} \cup \{y \preceq_{\text{char}} x \mid y \in \text{extr}(s)\}) & \text{otherwise} \end{cases}$$

The same applies to  $\llbracket s == x \rrbracket^{\text{char}} \mathcal{C}$ .

- *conjunction abstract semantics*

$$\llbracket e_1 \ \&\& \ e_2 \rrbracket^{\text{char}} \mathcal{C} = \llbracket e_1 \rrbracket^{\text{char}} \mathcal{C} \sqcup_{\text{char}} \llbracket e_2 \rrbracket^{\text{char}} \mathcal{C}$$

- *disjunction abstract semantics*

$$\llbracket e_1 \ || \ e_2 \rrbracket^{\text{char}} \mathcal{C} = \llbracket e_1 \rrbracket^{\text{char}} \mathcal{C} \sqcap_{\text{char}} \llbracket e_2 \rrbracket^{\text{char}} \mathcal{C}$$

**Theorem C.2.** The abstract semantics of Char is sound. Indeed, it holds that:

$$\forall M \in \mathcal{P}(\mathbb{M}). \alpha_{\text{char}}(\llbracket x = s \rrbracket M) \sqsubseteq_{\text{char}} (\llbracket x = s \rrbracket^{\text{char}} \alpha_{\text{char}}(M))$$

and

$$\forall M \in \mathcal{P}(\mathbb{M}). \alpha_{\text{char}}(\llbracket b \rrbracket M) \sqsubseteq_{\text{sub}} (\llbracket b \rrbracket^{\text{char}} \alpha_{\text{char}}(M))$$

*Proof.*

The proof is analogous to the one provided for Theorem C.3.

□

### C.3 Abstract Semantics of Sub

Let  $\llbracket \text{st} \rrbracket^{\text{sub}} : \text{Sub} \rightarrow \text{Sub}$  be the function that given an input abstract memory returns an abstract memory containing the new substring relations introduced by **st**.

- **Assignment abstract semantics**

Consider the assignment  $x = s$  such that  $x \in X_{\text{str}}$  and  $s \in \text{SE}$ . Formally, its abstract semantics, i.e.,  $\llbracket x = s \rrbracket^{\text{sub}} \mathcal{S}$  is given by the following steps:

- **[remove]**

$$\mathcal{S}_r = \begin{cases} \mathcal{S} \setminus \{w \preceq_{\text{sub}} z \mid w = x\} & \text{if } x \text{ appears at the top-level of } s \\ \mathcal{S} \setminus \{w \preceq_{\text{sub}} z \mid w = x \vee z = x\} & \text{otherwise} \end{cases}$$

– [add]

$$\mathcal{S}_a = \mathcal{S}_r \cup \{y \preceq_{\text{sub}} x \mid y \in \text{extr}(\mathbf{s})\}$$

where  $\text{extr}$  is the extraction function defined in Section 6.3.2.

– [closure]

$$\llbracket x = \mathbf{s} \rrbracket^{\text{sub}} \mathcal{S} = \text{Clos}(\mathcal{S}_a)$$

• **Boolean expressions abstract semantics**

- *contains abstract semantics*

$$\llbracket \text{contains}(x, \mathbf{s}) \rrbracket^{\text{sub}} \mathcal{S} = \text{Clos}(\mathcal{S} \cup \{y \preceq_{\text{sub}} x \mid y \in \text{extr}(\mathbf{s})\})$$

- *equality abstract semantics*

$$\llbracket x == \mathbf{s} \rrbracket^{\text{sub}} \mathcal{S} = \begin{cases} \text{Clos}(\mathcal{S} \cup \{y \preceq_{\text{sub}} x, x \preceq_{\text{sub}} y\}) & \text{if } \mathbf{s} = y \in X_{\text{str}} \\ \text{Clos}(\mathcal{S} \cup \{y \preceq_{\text{sub}} x \mid y \in \text{extr}(\mathbf{s})\}) & \text{otherwise} \end{cases}$$

The same applies to  $\llbracket \mathbf{s} == x \rrbracket^{\text{sub}} \mathcal{S}$ .

- *conjunction abstract semantics*

$$\llbracket e_1 \ \&\& \ e_2 \rrbracket^{\text{sub}} \mathcal{S} = \llbracket e_1 \rrbracket^{\text{sub}} \mathcal{S} \sqcup_{\text{sub}} \llbracket e_2 \rrbracket^{\text{sub}} \mathcal{S}$$

- *disjunction abstract semantics*

$$\llbracket e_1 \ \|\ e_2 \rrbracket^{\text{sub}} \mathcal{S} = \llbracket e_1 \rrbracket^{\text{sub}} \mathcal{S} \cap_{\text{sub}} \llbracket e_2 \rrbracket^{\text{sub}} \mathcal{S}$$

**Theorem C.3.** The abstract semantics of **Sub** is sound, namely  $\forall M \in \mathcal{P}(\mathbb{M})$

$$\alpha_{\text{sub}}(\llbracket x = \mathbf{s} \rrbracket M) \sqsubseteq_{\text{sub}} \llbracket x = \mathbf{s} \rrbracket^{\text{sub}} \alpha_{\text{sub}}(M)$$

and

$$\alpha_{\text{sub}}(\llbracket \mathbf{b} \rrbracket M) \sqsubseteq_{\text{sub}} \llbracket \mathbf{b} \rrbracket^{\text{sub}} \alpha_{\text{sub}}(M)$$

*Proof.*

The soundness proofs for the Boolean expressions (`contains`, `==`, disjunction and conjunction) are trivial and straightforward. Hence, we focus only on the soundness proof of the assignment.

First, we observe that given a string expression  $\mathbf{s}$ , any sub-expression  $\mathbf{s}'$  of  $\mathbf{s}$  is a substring of  $\mathbf{s}$ , for any  $\mathbf{m} \in \mathbb{M}$ , namely

$$\forall \mathbf{m} \in \mathbb{M} : \llbracket \mathbf{s}' \rrbracket \mathbf{m} \curvearrowright \llbracket \mathbf{s} \rrbracket \mathbf{m} \tag{C.1}$$

Hence, we prove that:

$$\alpha_{\text{sub}}(\llbracket x = s \rrbracket M) \sqsubseteq_{\text{sub}} (\llbracket x = s \rrbracket^{\text{sub}} \alpha_{\text{sub}}(M))$$

Without loss of generality, let us suppose that  $x$  does not appear at the top-level of the expression  $s$ . Note that the case in which  $x$  appears at the top-level of the expression  $s$  is analogous.

$$\begin{aligned}
& \alpha_{\text{sub}}(\llbracket x = s \rrbracket M) = \\
& = \{y \preceq_{\text{sub}} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](y) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z)\} && \text{[Def. } \alpha_{\text{sub}}\text{]} \\
& = \{y \preceq_{\text{sub}} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](y) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), z \neq x\} \\
& \cup \{y \preceq_{\text{sub}} x \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](y) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](x)\} && \text{[Set prop.]} \\
& \sqsubseteq_{\text{sub}} (\{y \preceq_{\text{sub}} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](y) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), z \neq x\} \\
& \quad \setminus \{y \preceq_{\text{sub}} z \mid y = x \vee z = x\}) \\
& \cup \{y \preceq_{\text{sub}} x \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](y) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](x)\} \\
& = (\{y \preceq_{\text{sub}} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](y) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), z \neq x\} \\
& \quad \setminus \{y \preceq_{\text{sub}} z \mid y = x \vee z = x\}) \\
& \cup \{y \preceq_{\text{sub}} x \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](y) \curvearrowright \llbracket s \rrbracket m\} && \text{[Def. } m(x)\text{]} \\
& \sqsubseteq_{\text{sub}} (\{y \preceq_{\text{sub}} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](y) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), z \neq x\} \\
& \quad \setminus \{y \preceq_{\text{sub}} z \mid y = x \vee z = x\}) \cup \{y \preceq_{\text{sub}} x \mid y \in \text{extr}(s)\} && \text{[Eq. C.1]} \\
& = \text{Clos}(\mathcal{S}_a) = \llbracket x = s \rrbracket^{\text{sub}} \alpha_{\text{sub}}(M) && \text{[Def. } \llbracket \cdot \rrbracket^{\text{sub}}\text{]}
\end{aligned}$$

□

## C.4 Abstract Semantics of $\text{Sub}^*$

Let  $\llbracket \text{st} \rrbracket^{\text{sub}^*} : \text{Sub}^* \rightarrow \text{Sub}^*$  be the function that given an input abstract memory returns an abstract memory containing the new substring relations introduced by  $\text{st}$ .

- **Assignment abstract semantics**

Consider the assignment  $x = s$  such that  $x \in X_{\text{str}}$  and  $s \in \text{SE}$ . Formally, its abstract semantics, i.e.,  $\llbracket x = s \rrbracket^{\text{sub}} \mathcal{S}$  is given by the following steps:

– [remove]

$$\mathcal{S}_r^* = \begin{cases} \mathcal{S}^* \setminus \{s' \preceq_{\text{sub}^*} z \mid x \text{ appears in } s'\} & \text{if } x \text{ appears} \\ & \text{at the top-level of } s \\ \mathcal{S}^* \setminus \{s' \preceq_{\text{sub}^*} z \mid z = x \vee x \text{ appears in } s'\} & \text{otherwise} \end{cases}$$

– [add]

$$\mathcal{S}_a^* = \mathcal{S}_r^* \cup \{s' \preceq_{\text{sub}^*} x \mid s' \in \text{extr}^*(s)\}$$

– [inter-asg]

$$\mathcal{S}_i^* = \mathcal{S}_a^* \cup \{x \preceq_{\text{sub}^*} y \mid \mathcal{S}_a^*(x) \subseteq \mathcal{S}_a^*(y)\}$$

– [closure]

$$\llbracket x = s \rrbracket^{\text{sub}^*} \mathcal{S}^* = \text{Clos}(\mathcal{S}_i^*)$$

• Boolean expressions abstract semantics

- *contains abstract semantics*

$$\llbracket \text{contains}(x, s) \rrbracket^{\text{sub}^*} \mathcal{S}^* = \text{Clos}(\mathcal{S}^* \cup \{s' \preceq_{\text{sub}^*} x \mid s' \in \text{extr}^*(s)\})$$

- *equality abstract semantics*

$$\llbracket x == s \rrbracket^{\text{sub}^*} \mathcal{S}^* = \text{Clos}(\mathcal{S}^* \cup \{s' \preceq_{\text{sub}^*} x \mid s' \in \text{extr}^*(s)\})$$

The same applies to  $\llbracket s == x \rrbracket^{\text{sub}^*} \mathcal{S}^*$ .

- *conjunction abstract semantics*

$$\llbracket e_1 \ \&\& \ e_2 \rrbracket^{\text{sub}^*} \mathcal{S}^* = \llbracket e_1 \rrbracket^{\text{sub}^*} \mathcal{S}^* \sqcup_{\text{sub}^*} \llbracket e_2 \rrbracket^{\text{sub}^*} \mathcal{S}^*$$

- *disjunction abstract semantics*

$$\llbracket e_1 \ \|\| \ e_2 \rrbracket^{\text{sub}^*} \mathcal{S}^* = \llbracket e_1 \rrbracket^{\text{sub}^*} \mathcal{S}^* \sqcap_{\text{sub}^*} \llbracket e_2 \rrbracket^{\text{sub}^*} \mathcal{S}^*$$

**Theorem C.4.** The abstract semantics of  $\text{Sub}^*$  is sound, namely  $\forall M \in \mathcal{P}(\mathbb{M})$

$$\alpha_{\text{sub}^*}(\llbracket x = s \rrbracket M) \sqsubseteq_{\text{sub}^*} \llbracket x = s \rrbracket^{\text{sub}^*} \alpha_{\text{sub}^*}(M)$$

and

$$\alpha_{\text{sub}^*}(\llbracket b \rrbracket M) \sqsubseteq_{\text{sub}^*} \llbracket b \rrbracket^{\text{sub}^*} \alpha_{\text{sub}^*}(M)$$

*Proof.* In the following we prove that:

$$\alpha_{\text{sub}^*}(\llbracket x = s \rrbracket M) \sqsubseteq_{\text{sub}^*} (\llbracket x = s \rrbracket^{\text{sub}^*} \alpha_{\text{sub}^*}(M))$$

The proof is similar to the one reported above, except for the phase **[inter-asg]**. As above, we assume that  $x$  does not appear at the top-level of the expression  $s$  and the contrary in analogous.

$$\begin{aligned}
& \alpha_{\text{sub}^*}(\llbracket x = s \rrbracket M) = \\
& = \{s' \preceq_{\text{sub}^*} z \mid \forall m \in M. \llbracket s' \rrbracket(m[x \leftarrow \llbracket s \rrbracket m]) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z)\} && [\text{Def. } \alpha_{\text{sub}^*}] \\
& = \{s' \preceq_{\text{sub}^*} z \mid \forall m \in M. s'(m[x \leftarrow \llbracket s \rrbracket m]) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), s' \neq x\} \\
& \cup \{x \preceq_{\text{sub}^*} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](x) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z)\} && [\text{Set prop.}] \\
& = \{s' \preceq_{\text{sub}^*} z \mid \forall m \in M. s'(m[x \leftarrow \llbracket s \rrbracket m]) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), s' \neq x, z \neq x\} \\
& \cup \{s' \preceq_{\text{sub}^*} x \mid \forall m \in M. s'(m[x \leftarrow \llbracket s \rrbracket m]) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), s' \neq x\} \\
& \cup \{x \preceq_{\text{sub}^*} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](x) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z)\} && [\text{Set prop.}] \\
& \sqsubseteq_{\text{sub}^*} \{s' \preceq_{\text{sub}^*} z \mid \forall m \in M. s'(m[x \leftarrow \llbracket s \rrbracket m]) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), s' \neq x, z \neq x\} \\
& \quad \setminus \{s' \preceq_{\text{sub}^*} x \mid z = x \vee x \text{ appears in } s\} \\
& \cup \{s' \preceq_{\text{sub}^*} x \mid \forall m \in M. s'(m[x \leftarrow \llbracket s \rrbracket m]) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), s' \neq x\} \\
& \cup \{x \preceq_{\text{sub}^*} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](x) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z)\} \\
& \sqsubseteq_{\text{sub}^*} (\{s' \preceq_{\text{sub}^*} z \mid \forall m \in M. s'(m[x \leftarrow \llbracket s \rrbracket m]) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z), s' \neq x, z \neq x\} \\
& \quad \setminus \{s' \preceq_{\text{sub}^*} x \mid z = x \vee x \text{ appears in } s\}) \\
& \cup \{s' \preceq_{\text{sub}^*} x \mid s' \in \text{extr}^*(s)\} \\
& \cup \{x \preceq_{\text{sub}^*} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](x) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z)\} && [\text{Eq. C.1}] \\
& = \mathcal{S}_a^* \cup \{x \preceq_{\text{sub}^*} z \mid \forall m \in M. m[x \leftarrow \llbracket s \rrbracket m](x) \curvearrowright m[x \leftarrow \llbracket s \rrbracket m](z)\} && [\text{Def. } \mathcal{S}_a^*] \\
& \sqsubseteq_{\text{sub}^*} \text{Clos}(\mathcal{S}_a^* \cup \{x \preceq_{\text{sub}^*} z \mid \mathcal{S}_a^*(x) \subseteq \mathcal{S}_a^*(z)\}) = \llbracket x = s \rrbracket^{\text{sub}^*} \alpha_{\text{sub}^*}(M) && [\text{Def. } \llbracket \cdot \rrbracket^{\text{sub}^*}]
\end{aligned}$$

□

## C.5 Len<sup>\*</sup> Relational Abstract Domain

The Len relational string abstract domain can be extended in order to track also relations between string expressions and variables, e.g., if the length of the concatenation between two string variables is smaller than or equal to the length of another string variable. Thus, similarly to Section 6.3.5, we introduce the binary relation  $\preceq_{\text{len}^*} \subseteq \text{SE} \times X_{\text{str}}$  and the extended abstract domain  $\text{Len}^* = \mathcal{P}(\{s \preceq_{\text{len}^*} x \mid x \in X_{\text{str}}, s \in \text{SE}\}) \cup \{\perp_{\text{len}^*}\}$ , which ranges over the meta-variable  $\mathcal{L}^*$ , where the top element, denoted by  $\top_{\text{len}^*}$ , is  $\emptyset$ , and  $\perp_{\text{len}^*}$  is the bottom element.  $\text{Len}^*$  inherits the properties of  $\text{Len}$ .  $\text{Len}^*$  lattice operations are  $\sqcup_{\text{len}^*}$  (lub),  $\sqcap_{\text{len}^*}$  (glb) and  $\preceq_{\text{len}^*}$  (partial order) and are similar to the ones of  $\text{Len}$  reported in Figure 6.4. Its concretization  $\gamma_{\text{len}^*} : \text{Len}^* \rightarrow \mathcal{P}(M)$  and abstraction  $\alpha_{\text{len}^*} : \mathcal{P}(M) \rightarrow \text{Len}^*$

functions, which form a Galois Connection, are defined as follows:

$$\gamma_{\text{len}^*}(\mathcal{L}^*) = \bigcap_{\mathbf{s} \preceq_{\text{len}^*} x \in \mathcal{L}^*} \{ \mathbf{m} \mid \mathbf{m}(x), \llbracket \mathbf{s} \rrbracket \mathbf{m} \in \Sigma^*, \|\llbracket \mathbf{s} \rrbracket \mathbf{m}\| \leq |\mathbf{m}(x)| \}$$

$$\alpha_{\text{len}^*}(M) = \{ \mathbf{s} \preceq_{\text{len}^*} x \mid \forall \mathbf{m} \in M. \|\llbracket \mathbf{s} \rrbracket \mathbf{m}\| \leq |\mathbf{m}(x)|, x \in X_{\text{str}}, \mathbf{s} \in SE \}$$

The  $\text{Len}^*$  abstract semantics extends the one defined for  $\text{Len}$  and it is defined as follows.

Let  $\llbracket \text{st} \rrbracket^{\text{len}^*} \mathcal{L}^* : \text{Len}^* \rightarrow \text{Len}^*$ .

- **Assignment abstract semantics**

Consider the assignment  $x = \mathbf{s}$  such that  $x \in X_{\text{str}}$  and  $\mathbf{s} \in SE$ . Formally, its abstract semantics, i.e.,  $\llbracket x = \mathbf{s} \rrbracket^{\text{len}^*} \mathcal{L}^*$ , is given by the following steps:

– [remove]

$$\mathcal{L}_r^* = \begin{cases} \mathcal{L}^* \setminus \{ \mathbf{s}' \preceq_{\text{len}^*} z \mid x \text{ appears in } \mathbf{s}' \} & \text{if } x \text{ appears} \\ & \text{at the top-level of } \mathbf{s} \\ \mathcal{L}^* \setminus \{ \mathbf{s}' \preceq_{\text{len}^*} z \mid z = x \vee x \text{ appears in } \mathbf{s} \} & \text{otherwise} \end{cases}$$

– [add]

$$\mathcal{L}_a^* = \mathcal{L}_r^* \cup \{ \mathbf{s}' \preceq_{\text{len}^*} x \mid \mathbf{s}' \in \text{extr}^*(\mathbf{s}) \}$$

– [inter-asg]

$$\mathcal{L}_i^* = \mathcal{L}_a^* \cup \{ x \preceq_{\text{len}^*} y \mid \mathcal{L}_a^*(x) \subseteq \mathcal{L}_a^*(y) \}$$

– [closure]

$$\llbracket x = \mathbf{s} \rrbracket^{\text{len}^*} \mathcal{L}^* = \text{Clos}(\mathcal{L}_i^*)$$

- **Boolean expressions abstract semantics**

- *contains abstract semantics*

$$\llbracket \text{contains}(x, \mathbf{s}) \rrbracket^{\text{len}^*} \mathcal{L}^* = \text{Clos}(\mathcal{L}^* \cup \{ \mathbf{s}' \preceq_{\text{len}^*} x \mid \mathbf{s}' \in \text{extr}^*(\mathbf{s}) \})$$

- *equality abstract semantics*

$$\llbracket x == \mathbf{s} \rrbracket^{\text{len}^*} \mathcal{L}^* = \text{Clos}(\mathcal{L}^* \cup \{ \mathbf{s}' \preceq_{\text{len}^*} x \mid \mathbf{s}' \in \text{extr}^*(\mathbf{s}) \})$$

The same applies to  $\llbracket \mathbf{s} == x \rrbracket^{\text{len}^*} \mathcal{L}^*$ .

- *conjunction abstract semantics*

$$\llbracket e_1 \ \&\& \ e_2 \rrbracket^{\text{len}^*} \mathcal{L}^* = \llbracket e_1 \rrbracket^{\text{len}^*} \mathcal{L}^* \sqcup_{\text{len}^*} \llbracket e_2 \rrbracket^{\text{len}^*} \mathcal{L}^*$$

- *disjunction abstract semantics*

$$\llbracket e_1 \parallel e_2 \rrbracket^{\text{len}^*} \mathcal{L}^* = \llbracket e_1 \rrbracket^{\text{len}^*} \mathcal{L}^* \sqcap_{\text{len}^*} \llbracket e_2 \rrbracket^{\text{len}^*} \mathcal{L}^*$$

**Theorem C.5.** The abstract semantics of  $\text{Len}^*$  is sound. Indeed, it holds that:

$$\forall M \in \mathcal{P}(\mathbb{M}). \alpha_{\text{len}^*}(\llbracket x = s \rrbracket M) \sqsubseteq_{\text{len}^*} (\llbracket x = s \rrbracket^{\text{len}^*} \alpha_{\text{len}^*}(M))$$

and

$$\forall M \in \mathcal{P}(\mathbb{M}). \alpha_{\text{len}^*}(\llbracket b \rrbracket M) \sqsubseteq_{\text{len}^*} (\llbracket b \rrbracket^{\text{len}^*} \alpha_{\text{len}^*}(M))$$

*Proof.*

The proof is analogous to the one provided for Theorem C.4. □

## C.6 Char<sup>\*</sup> Relational Abstract Domain

We define  $\preceq_{\text{char}^*}$  as the extension of the binary relation  $\preceq_{\text{char}}$  to string expressions such that  $\preceq_{\text{char}^*} \subseteq \text{SE} \times X_{\text{str}}$ . Moreover,  $\text{Char}$  becomes  $\text{Char}^* = \mathcal{P}(\{\mathbf{s} \preceq_{\text{char}^*} x \mid x \in X_{\text{str}}, \mathbf{s} \in \text{SE}\}) \cup \{\perp_{\text{char}^*}\}$ , where the top element, denoted by  $\top_{\text{char}^*}$ , is  $\emptyset$ , and  $\perp_{\text{char}^*}$  is the bottom element. The abstract domain  $\text{Char}^*$  ranges over the meta-variable  $\mathcal{C}^*$ . The extended relational character inclusion string abstract domain is a complete lattice and its lattice operations ( $\sqcup_{\text{char}^*}$  (lub),  $\sqcap_{\text{char}^*}$  (glb),  $\sqsubseteq_{\text{char}^*}$  (partial order)) are similar to the ones of  $\text{Char}$ . The concretization  $\gamma_{\text{char}^*} : \text{Char}^* \rightarrow \mathcal{P}(\mathbb{M})$  and abstraction  $\alpha_{\text{char}^*} : \mathcal{P}(\mathbb{M}) \rightarrow \text{Char}^*$  functions, defined below, form a Galois Connection.

$$\gamma_{\text{char}^*}(\mathcal{C}^*) = \bigcap_{\mathbf{s} \preceq_{\text{char}^*} x \in \mathcal{C}^*} \{\mathbf{m} \mid \mathbf{m}(x), \llbracket \mathbf{s} \rrbracket \mathbf{m} \in \Sigma^*, \text{char}(\llbracket \mathbf{s} \rrbracket \mathbf{m}) \subseteq \text{char}(\mathbf{m}(x))\}$$

$$\alpha_{\text{char}^*}(M) = \{\mathbf{s} \preceq_{\text{char}^*} x \mid \forall \mathbf{m} \in M. \text{char}(\llbracket \mathbf{s} \rrbracket \mathbf{m}) \subseteq \text{char}(\mathbf{m}(x)), x \in X_{\text{str}}, \mathbf{s} \in \text{SE}\}$$

As usual, the abstract semantics is captured by the function  $\llbracket \text{st} \rrbracket^{\text{char}^*} \mathcal{C}^* : \text{Char}^* \rightarrow \text{Char}^*$ . and it is defined as follows.

- **Assignment abstract semantics**

Consider the assignment  $x = s$  such that  $x \in X_{\text{str}}$  and  $\mathbf{s} \in \text{SE}$ . Formally, its abstract semantics, i.e.,  $\llbracket x = s \rrbracket^{\text{char}^*} \mathcal{C}^*$  is given by the following steps:

- [remove]

$$\mathcal{C}_r^* = \begin{cases} \mathcal{C}^* \setminus \{\mathbf{s}' \preceq_{\text{char}^*} z \mid x \text{ appears in } \mathbf{s}'\} & \text{if } x \text{ appears} \\ & \text{at the top-level of } \mathbf{s} \\ \mathcal{C}^* \setminus \{\mathbf{s}' \preceq_{\text{char}^*} z \mid z = x \vee x \text{ appears in } \mathbf{s}'\} & \text{otherwise} \end{cases}$$

– [add]

$$\mathcal{C}_a^* = \mathcal{C}_r^* \cup \{s' \preceq_{\text{char}^*} x \mid s' \in \text{extr}^*(s)\}$$

– [inter-asg]

$$\mathcal{C}_i^* = \mathcal{C}_a^* \cup \{x \preceq_{\text{char}^*} y \mid \mathcal{C}_a^*(x) \subseteq \mathcal{C}_a^*(y)\}$$

– [closure]

$$\llbracket x = s \rrbracket^{\text{char}^*} \mathcal{C}^* = \text{Clos}(\mathcal{C}_i^*)$$

• Boolean expressions abstract semantics

- *contains abstract semantics*

$$\llbracket \text{contains}(x, s) \rrbracket^{\text{char}^*} \mathcal{C}^* = \text{Clos}(\mathcal{C}^* \cup \{s' \preceq_{\text{char}^*} x \mid s' \in \text{extr}^*(s)\})$$

- *equality abstract semantics*

$$\llbracket x == s \rrbracket^{\text{char}^*} \mathcal{C}^* = \text{Clos}(\mathcal{C}^* \cup \{s' \preceq_{\text{char}^*} x \mid s' \in \text{extr}^*(s)\})$$

The same applies to  $\llbracket s == x \rrbracket^{\text{char}^*} \mathcal{C}^*$ .

- *conjunction abstract semantics*

$$\llbracket e_1 \ \&\& \ e_2 \rrbracket^{\text{char}^*} \mathcal{C}^* = \llbracket e_1 \rrbracket^{\text{char}^*} \mathcal{C}^* \sqcup_{\text{char}^*} \llbracket e_2 \rrbracket^{\text{char}^*} \mathcal{C}^*$$

- *disjunction abstract semantics*

$$\llbracket e_1 \ \|\ \ e_2 \rrbracket^{\text{char}^*} \mathcal{C}^* = \llbracket e_1 \rrbracket^{\text{char}^*} \mathcal{C}^* \sqcap_{\text{char}^*} \llbracket e_2 \rrbracket^{\text{char}^*} \mathcal{C}^*$$

**Theorem C.6.** The abstract semantics of  $\text{Char}^*$  is sound. Indeed, it holds that:

$$\forall M \in \mathcal{P}(\mathbb{M}). \alpha_{\text{char}^*}(\llbracket x = s \rrbracket M) \sqsubseteq_{\text{char}^*} (\llbracket x = s \rrbracket^{\text{char}^*} \alpha_{\text{char}^*}(M))$$

and

$$\forall M \in \mathcal{P}(\mathbb{M}). \alpha_{\text{char}^*}(\llbracket b \rrbracket M) \sqsubseteq_{\text{sub}^*} (\llbracket b \rrbracket^{\text{char}^*} \alpha_{\text{char}^*}(M))$$

*Proof.*

The proof is analogous to the one provided for Theorem C.4.

□



## BIBLIOGRAPHY

- [1] P. A. Abdulla et al. “Norn: An SMT Solver for String Constraints”. In: *Computer Aided Verification - 27th International Conference, CAV 2015, San Francisco, CA, USA, July 18-24, 2015, Proceedings, Part I*. 2015, pp. 462–469. DOI: 10.1007/978-3-319-21690-4\\_29.
- [2] E. Agirre et al. “A Study on Similarity and Relatedness Using Distributional and WordNet-based Approaches”. In: *Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics, Proceedings, May 31 - June 5, 2009, Boulder, Colorado, USA*. The Association for Computational Linguistics, 2009, pp. 19–27.
- [3] R. Agrawal, P. J. Haas, and J. Kiernan. “Watermarking Relational Data: Framework, Algorithms and Analysis”. In: *VLDB J.* 12.2 (2003), pp. 157–169. DOI: 10.1007/s00778-003-0097-x.
- [4] R. Agrawal and J. Kiernan. “Watermarking Relational Databases”. In: *Proceedings of 28th International Conference on Very Large Data Bases, VLDB 2002, Hong Kong, August 20-23, 2002*. Morgan Kaufmann, 2002, pp. 155–166. DOI: 10.1016/B978-155860869-6/50022-6.
- [5] M. T. Ahvanooy et al. “A Comparative Analysis of Information Hiding Techniques for Copyright Protection of Text Documents”. In: *Security and Communication Networks 2018* (2018), 5325040:1–5325040:22. DOI: 10.1155/2018/5325040.
- [6] A. Al-Haj and A. Odeh. “Robust and Blind Watermarking of Relational Database Systems”. In: *Journal of Computer Science* 4.12 (2008), pp. 1024–1029.
- [7] F. E. Allen. “Control flow analysis”. In: *SIGPLAN Notices*. Vol. 5. 1970, pp. 1–19. DOI: 10.1145/390013.808479.
- [8] R. Amadini, G. Gange, and P. J. Stuckey. “Dashed Strings for String Constraint Solving”. In: *Artificial Intelligence*. Vol. 289. 2020. DOI: <https://doi.org/10.1016/j.artint.2020.103368>.

- [9] R. Amadini et al. “A Novel Approach to String Constraint Solving”. In: *Principles and Practice of Constraint Programming - 23rd International Conference, CP 2017, Melbourne, VIC, Australia, August 28 - September 1, 2017, Proceedings*. 2017, pp. 3–20. URL: [10.1007/978-3-319-66158-2\\_1](https://doi.org/10.1007/978-3-319-66158-2_1).
- [10] R. Amadini et al. “Combining String Abstract Domains for JavaScript Analysis: An Evaluation”. In: *Tools and Algorithms for the Construction and Analysis of Systems - 23rd International Conference, TACAS 2017, held as part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2017, Uppsala, Sweden, April 22-29, 2017, Proceedings, Part I*. 2017, pp. 41–57. DOI: [10.1007/978-3-662-54577-5\\_3](https://doi.org/10.1007/978-3-662-54577-5_3).
- [11] R. Amadini et al. “Reference Abstract Domains and Applications to String Analysis”. In: *Fundam. Inform.* 158.4 (2018), pp. 297–326. DOI: [10.3233/FI-2018-1650](https://doi.org/10.3233/FI-2018-1650).
- [12] V. Arceri. “Taming Strings in Dynamic Languages - An Abstract Interpretation-based Static Analysis Approach”. PhD thesis. May 2020. DOI: [10.13140/RG.2.2.27093.45286](https://doi.org/10.13140/RG.2.2.27093.45286).
- [13] V. Arceri and S. Maffei. “Abstract Domains for Type Juggling”. In: *Electron. Notes Theor. Comput. Sci.* 331 (2017), pp. 41–55. DOI: [10.1016/j.entcs.2017.02.003](https://doi.org/10.1016/j.entcs.2017.02.003).
- [14] V. Arceri and I. Mastroeni. “A Sound Abstract Interpreter for Dynamic Code”. In: *SAC '20: The 35th ACM/SIGAPP Symposium on Applied Computing*, online event, Brno, Czech Republic, March 30 - April 3, 2020. Ed. by C. Hung et al. ACM, 2020, pp. 1979–1988. DOI: [10.1145/3341105.3373964](https://doi.org/10.1145/3341105.3373964).
- [15] V. Arceri and I. Mastroeni. “An Automata-based Abstract Semantics for String Manipulation Languages”. In: *Proceedings Seventh International Workshop on Verification and Program Transformation, VPT@Programming 2019, Genova, Italy, 2nd April 2019*. Ed. by A. Lisitsa and A. P. Nemytykh. Vol. 299. EPTCS. 2019, pp. 19–33. DOI: [10.4204/EPTCS.299.5](https://doi.org/10.4204/EPTCS.299.5).
- [16] V. Arceri and I. Mastroeni. “Static Program Analysis for String Manipulation Languages”. In: *Proceedings Seventh International Workshop on Verification and Program Transformation, Genova, Italy, 2nd April 2019*. Ed. by A. Lisitsa and A. Nemytykh. Vol. 299. *Electronic Proceedings in Theoretical Computer Science*. Open Publishing Association, 2019, pp. 19–33. DOI: [10.4204/EPTCS.299.5](https://doi.org/10.4204/EPTCS.299.5).
- [17] V. Arceri, I. Mastroeni, and S. Xu. “Static Analysis for ECMAScript String Manipulation Programs”. In: *Appl. Sci.* 10 (2020), p. 3525. DOI: [10.3390/app10103525](https://doi.org/10.3390/app10103525).
- [18] V. Arceri et al. “Completeness of Abstract Domains for String Analysis of JavaScript Programs”. In: *Theoretical Aspects of Computing - ICTAC 2019 - 16th International Colloquium, Hammamet, Tunisia, October 31 - November 4, 2019, Proceedings*. Ed.

- by R. M. Hierons and M. Mosbah. Vol. 11884. Lecture Notes in Computer Science. Springer, 2019, pp. 255–272. DOI: 10.1007/978-3-030-32505-3\\_15.
- [19] N. Ayewah et al. “Using Static Analysis to Find Bugs”. In: *IEEE Software* 25.5 (2008), pp. 22–29. DOI: 10.1109/MS.2008.130.
- [20] R. Bagnara, P. M. Hill, and E. Zaffanella. “Exact Join Detection for Convex Polyhedra and Other Numerical Abstractions”. In: *Comput. Geom.* 43.5 (2010), pp. 453–473. DOI: 10.1016/j.comgeo.2009.09.002.
- [21] R. Bagnara, P. M. Hill, and E. Zaffanella. “The Parma Polyhedra Library: Toward a Complete Set of Numerical Abstractions for the Analysis and Verification of Hardware and Software Systems”. In: *Sci. Comput. Program.* 72.1-2 (2008), pp. 3–21. DOI: 10.1016/j.scico.2007.08.001.
- [22] C. Baier and J. P. Katoen. “Principles of Model Checking”. MIT Press, 2008.
- [23] Z. Baranová et al. “Model Checking of C and C++ with DIVINE 4”. In: *Automated Technology for Verification and Analysis - 15th International Symposium, ATVA 2017, Pune, India, October 3-6, 2017, Proceedings.* 2017, pp. 201–207. DOI: 10.1007/978-3-319-68167-2\\_14.
- [24] M. Batet and D. Sánchez. “A Review on Semantic Similarity”. In: *Encyclopedia of Information Science and Technology, Third Edition.* IGI Global, 2015, pp. 7575–7583.
- [25] S. Bautista, T. Jensen, and B. Montagu. “Numeric Domains Meet Algebraic Data Types”. In: *Proceedings of the 9th ACM SIGPLAN International Workshop on Numerical and Symbolic Abstract Domains. NSAD 2020.* Virtual, USA: Association for Computing Machinery, 2020, 12–16. DOI: 10.1145/3427762.3430178.
- [26] E. Bertino et al. “Privacy and Ownership Preserving of Outsourced Medical Data”. In: *Proceedings of the 21st International Conference on Data Engineering, ICDE 2005, 5-8 April 2005, Tokyo, Japan.* Ed. by K. Aberer, M. J. Franklin, and S. Nishio. IEEE Computer Society, 2005, pp. 521–532. DOI: 10.1109/ICDE.2005.111.
- [27] A. Bessey et al. “A Few Billion Lines of Code Later: Using Static Analysis to Find Bugs in the Real World”. In: *Commun. ACM* 53.2 (2010), pp. 66–75. DOI: 10.1145/1646353.1646374.
- [28] S. Bhattacharya and A. Cortesi. “A Distortion Free Watermark Framework for Relational Databases”. In: *ICSOFIT 2009 - Proceedings of the 4th International Conference on Software and Data Technologies, Volume 2, Sofia, Bulgaria, July 26-29, 2009.* Ed. by B. Shishkov, J. Cordeiro, and A. Ranchordas. INSTICC Press, 2009, pp. 229–234.

- [29] S. Bhattacharya and A. Cortesi. “A Generic Distortion Free Watermarking Technique for Relational Databases”. In: Information Systems Security, 5th International Conference, ICISS 2009, Kolkata, India, December 14-18, 2009, Proceedings. Ed. by A. Prakash and I. Gupta. Vol. 5905. Lecture Notes in Computer Science. Springer, 2009, pp. 252–264. DOI: 10.1007/978-3-642-10772-6\\_19.
- [30] T. Bultan et al. “String Analysis for Software Verification and Security”. Springer, 2017. DOI: 10.1007/978-3-319-68670-7.
- [31] S. Cass. “IEEE Spectrum Ranking 2019”. <https://spectrum.ieee.org/computing/software/the-top-programming-languages-2019>. Last check: November 23rd, 2020.
- [32] I. Casso et al. “Computing Abstract Distances in Logic Programs”. In: Logic-Based Program Synthesis and Transformation - 29th International Symposium, LOPSTR 2019, Porto, Portugal, October 8-10, 2019, Revised Selected Papers. Ed. by Maurizio Gabbriellini. Vol. 12042. Lecture Notes in Computer Science. Springer, 2019, pp. 57–72. DOI: 10.1007/978-3-030-45260-5\\_4.
- [33] C. Chang, T. S. Nguyen, and C. Lin. “A Blind Robust Reversible Watermark Scheme for Textual Relational Databases with Virtual Primary Key”. In: Digital-Forensics and Watermarking - 13th International Workshop, IWDW 2014, Taipei, Taiwan, October 1-4, 2014. Revised Selected Papers. Ed. by Y. Shi et al. Vol. 9023. Lecture Notes in Computer Science. Springer, 2014, pp. 75–89. DOI: 10.1007/978-3-319-19321-2\\_6.
- [34] A. S. Christensen, A. Møller, and M. I. Schwartzbach. “Precise Analysis of String Expressions”. In: Static Analysis, 10th International Symposium, SAS 2003, San Diego, CA, USA, June 11-13, 2003, Proceedings. 2003, pp. 1–18. DOI: 10.1007/3-540-44898-5\\_1.
- [35] E. M. Clarke, O. Grumberg, and D. E. Long. “Model Checking and Abstraction”. In: ACM Trans. Program. Lang. Syst. 16.5 (1994), pp. 1512–1542. DOI: 10.1145/186025.186051.
- [36] E. M. Clarke et al. “Symbolic Model Checking”. In: Computer Aided Verification, 8th International Conference, CAV ’96, New Brunswick, NJ, USA, July 31 - August 3, 1996, Proceedings. 1996, pp. 419–427. DOI: 10.1007/3-540-61474-5\\_93.
- [37] E. F. Codd. “A Relational Model of Data for Large Shared Data Banks”. In: Commun. ACM 13.6 (1970), pp. 377–387. DOI: 10.1145/362384.362685.
- [38] M. Codish et al. “Improving Abstract Interpretations by Combining Domains”. In: ACM Trans. Program. Lang. Syst. 17.1 (1995), pp. 28–44. DOI: 10.1145/200994.200998.

- [39] Colorado-State-University. “Forest CoverType, The UCI KDD Archive”. <http://kdd.ics.uci.edu/databases/covertypetype/covertypetype.html>. Information and Computer Science. University of California, Irvine. Remote Sensing and GIS Program Department of Forest Sciences. College of Natural Resources. Colorado State University. Fort Collins, CO 80523. June 1999.
- [40] A. Cortesi, B. Le Charlier, and P. Van Hentenryck. “Combinations of Abstract Domains for Logic Programming”. In: Conference Record of POPL’94: 21st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Portland, Oregon, USA, January 17-21, 1994. 1994, pp. 227–239. DOI: 10.1145/174675.177880.
- [41] A. Cortesi, B. Le Charlier, and P. Van Hentenryck. “Combinations of Abstract Domains for Logic Programming: Open Product and Generic Pattern Construction”. In: *Sci. Comput. Program.* 38.1-3 (2000), pp. 27–71. DOI: 10.1016/S0167-6423(99)00045-3.
- [42] A. Cortesi, G. Costantini, and P. Ferrara. “A Survey on Product Operators in Abstract Interpretation”. In: *Semantics, Abstract Interpretation, and Reasoning about Programs: Essays Dedicated to David A. Schmidt on the Occasion of his Sixtieth Birthday*, Manhattan, Kansas, USA, 19-20th September 2013. 2013, pp. 325–336. DOI: 10.4204/EPTCS.129.19.
- [43] A. Cortesi, G. Costantini, and P. Ferrara. “The Abstract Domain of Trapezoid Step Functions”. In: *Computer Languages, Systems & Structures* 43 (2015), pp. 41–68. DOI: 10.1016/j.cl.2015.04.002.
- [44] A. Cortesi, G. Filé, and W. H. Winsborough. “Comparison of Abstract Interpretations”. In: *Automata, Languages and Programming, 19th International Colloquium, ICALP92*, Vienna, Austria, July 13-17, 1992, Proceedings. 1992, pp. 521–532. DOI: 10.1007/3-540-55719-9\_101.
- [45] A. Cortesi, G. Filé, and W. H. Winsborough. “The Quotient of an Abstract Interpretation”. In: *Theor. Comput. Sci.* 202.1-2 (1998), pp. 163–192. DOI: 10.1016/S0304-3975(97)00137-0.
- [46] A. Cortesi and M. Zanioli. “Widening and Narrowing Operators for Abstract Interpretation”. In: *Comput. Lang. Syst. Struct.* 37.1 (2011), pp. 24–42. DOI: 10.1016/j.cl.2010.09.001.
- [47] A. Cortesi and M. Olliaro. “M-String Segmentation: A Refined Abstract Domain for String Analysis in C Programs”. In: *2018 International Symposium on Theoretical Aspects of Software Engineering, TASE 2018*, Guangzhou, China, August 29-31, 2018. Ed. by J. Pang et al. IEEE Computer Society, 2018, pp. 1–8. DOI: 10.1109/TASE.2018.00009.

- [48] A. Cortesi et al. “String Abstraction for Model Checking of C Programs”. In: *Model Checking Software - 26th International Symposium, SPIN 2019, Beijing, China, July 15-16, 2019, Proceedings*. Ed. by F. Biondi, T. G. Wilson, and A. Legay. Vol. 11636. *Lecture Notes in Computer Science*. Springer, 2019, pp. 74–93. DOI: 10.1007/978-3-030-30923-7\5.
- [49] G. Costantini, P. Ferrara, and A. Cortesi. “A Suite of Abstract Domains for Static Analysis of String Values”. In: *Softw., Pract. Exper.* 45.2 (2015), pp. 245–287. DOI: 10.1002/spe.2218.
- [50] P. Cousot. “Abstract Interpretation Based Formal Methods and Future Challenges”. In: *Informatics - 10 Years Back. 10 Years Ahead.* 2001, pp. 138–156. DOI: 10.1007/3-540-44577-3\10.
- [51] P. Cousot and R. Cousot. “Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints”. In: *Conference Record of the Fourth ACM Symposium on Principles of Programming Languages, Los Angeles, California, USA, January 1977.* 1977, pp. 238–252. DOI: 10.1145/512950.512973.
- [52] P. Cousot and R. Cousot. “Abstract Interpretation Frameworks”. In: *J. Log. Comput.* 2.4 (1992), pp. 511–547. DOI: 10.1093/logcom/2.4.511.
- [53] P. Cousot and R. Cousot. “An Abstract Interpretation-based Framework for Software Watermarking”. In: *Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2004, Venice, Italy, January 14-16, 2004.* 2004, pp. 173–185. DOI: 10.1145/964001.964016.
- [54] P. Cousot and R. Cousot. “Comparing the Galois Connection and Widening/Narrowing Approaches to Abstract Interpretation”. In: *Programming Language Implementation and Logic Programming, 4th International Symposium, PLILP’92, Leuven, Belgium, August 26-28, 1992, Proceedings.* 1992, pp. 269–295. DOI: 10.1007/3-540-55844-6\142.
- [55] P. Cousot and R. Cousot. “Constructive Versions of Tarski’s Fixed Point Theorems”. In: *Pacific Journal of Mathematics* 81.1 (1979), pp. 43–57.
- [56] P. Cousot and R. Cousot. “Systematic Design of Program Analysis Frameworks”. In: *Conference Record of the Sixth Annual ACM Symposium on Principles of Programming Languages, San Antonio, Texas, USA, January 1979.* 1979, pp. 269–282. DOI: 10.1145/567752.567778.
- [57] P. Cousot, R. Cousot, and F. Logozzo. “A Parametric Segmentation Functor for Fully Automatic and Scalable Array Content Analysis”. In: *Proceedings of the 38th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages,*

- POPL 2011, Austin, TX, USA, January 26-28, 2011. 2011, pp. 105–118. DOI: 10.1145/1926385.1926399.
- [58] P. Cousot and N. Halbwachs. “Automatic Discovery of Linear Restraints Among Variables of a Program”. In: Conference Record of the Fifth Annual ACM Symposium on Principles of Programming Languages, Tucson, Arizona, USA, January 1978. 1978, pp. 84–96. DOI: 10.1145/512760.512770.
- [59] P. Cousot et al. “The ASTREÉ Analyzer”. In: Programming Languages and Systems, 14th European Symposium on Programming, ESOP 2005, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2005, Edinburgh, UK, April 4-8, 2005, Proceedings. Ed. by S. Sagiv. Vol. 3444. Lecture Notes in Computer Science. Springer, 2005, pp. 21–30. DOI: 10.1007/978-3-540-31987-0\3.
- [60] I. Cox et al. “Digital Watermarking and Steganography”. M. Kaufmann, 2007.
- [61] A. Dey, S. Bhattacharya, and N. Chaki. “Software Watermarking: Progress and Challenges”. In: (2018). Exported from <https://app.dimensions.ai> on 2019/03/16, pp. 1–11. DOI: 10.1007/s41403-018-0058-8.
- [62] N. Dor, M. Rodeh, and S. Sagiv. “CSSV: towards a realistic tool for statically detecting all buffer overflows in C”. In: Proceedings of the ACM SIGPLAN 2003 Conference on Programming Language Design and Implementation 2003, San Diego, California, USA, June 9-11, 2003. 2003, pp. 155–167. DOI: 10.1145/781131.781149.
- [63] D. Evans and D. Larochelle. “Improving Security Using Extensible Lightweight Static Analysis”. In: IEEE Software 19.1 (2002), pp. 42–51. DOI: 10.1109/52.976940.
- [64] M. Fähndrich and F. Logozzo. “Static Contract Checking with Abstract Interpretation”. In: Formal Verification of Object-Oriented Software - International Conference, FoVeOOS 2010, Paris, France, June 28-30, 2010, Revised Selected Papers. Ed. by B. Beckert and C. Marché. Vol. 6528. Lecture Notes in Computer Science. Springer, 2010, pp. 10–30. DOI: 10.1007/978-3-642-18070-5\2.
- [65] M. E. Farfoura et al. “A Blind Reversible Method for Watermarking Relational Databases Based on a Time-stamping Protocol”. In: Expert Syst. Appl. 39.3 (2012), pp. 3185–3196. DOI: 10.1016/j.eswa.2011.09.005.
- [66] P. Ferrara, F. Logozzo, and M. Fähndrich. “Safer unsafe code for .NET”. In: Proceedings of the 23rd Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2008, October 19-23, 2008, Nashville, TN, USA. Ed. by G. E. Harris. ACM, 2008, pp. 329–346. DOI: 10.1145/1449764.1449791.

- [67] D. Filaretti and S. Maffei. “An Executable Formal Semantics of PHP”. In: ECOOP 2014 - Object-Oriented Programming - 28th European Conference, Uppsala, Sweden, July 28 - August 1, 2014. Proceedings. 2014, pp. 567–592. DOI: 10.1007/978-3-662-44202-9\_23.
- [68] G. Filé, R. Giacobazzi, and F. Ranzato. “A Unifying View of Abstract Domain Design”. In: ACM Comput. Surv. 28.2 (1996), pp. 333–336. DOI: 10.1145/234528.234742.
- [69] J. Franco-Contreras and G. Coatrieux. “Robust Watermarking of Relational Databases With Ontology-Guided Distortion Control”. In: IEEE Trans. Inf. Forensics Secur. 10.9 (2015), pp. 1939–1952. DOI: 10.1109/TIFS.2015.2439962.
- [70] J. Franco-Contreras et al. “Ontology-guided Distortion Control for Robust-lossless Database Watermarking: Application to inpatient hospital stay records”. In: 36th Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBC 2014, Chicago, IL, USA, August 26-30, 2014. IEEE, 2014, pp. 4491–4494. DOI: 10.1109/EMBC.2014.6944621.
- [71] X. Fu et al. “A Static Analysis Framework For Detecting SQL Injection Vulnerabilities”. In: 31st Annual International Computer Software and Applications Conference, COMPSAC 2007, Beijing, China, July 24-27, 2007. Volume 1. 2007, pp. 87–96. DOI: 10.1109/COMPSAC.2007.43.
- [72] J. Fulara et al. “Relational Abstract Domain of Weighted Hexagons”. In: Electron. Notes Theor. Comput. Sci. 267.1 (2010), pp. 59–72. DOI: 10.1016/j.entcs.2010.09.006.
- [73] G. Gange et al. “Abstract Interpretation over Non-lattice Abstract Domains”. In: Static Analysis - 20th International Symposium, SAS 2013, Seattle, WA, USA, June 20-22, 2013. Proceedings. Ed. by Francesco Logozzo and Manuel Fähndrich. Vol. 7935. Lecture Notes in Computer Science. Springer, 2013, pp. 6–24. DOI: 10.1007/978-3-642-38856-9\_3.
- [74] H. Gericke. “Alfred Tarski. A Lattice-theoretical Fixpoint Theorem and its Applications”. In: Journal of Symbolic Logic. Pacific journal of mathematics, Bd. 5 (1955) 22.4 (1957), 285–309. DOI: 10.2307/2963937.
- [75] R. Giacobazzi and I. Mastroeni. “Making Abstract Models Complete”. In: Math. Struct. Comput. Sci. 26.4 (2016), pp. 658–701. DOI: 10.1017/S0960129514000358.
- [76] R. Giacobazzi and I. Mastroeni. “Transforming Abstract Interpretations by Abstract Interpretation”. In: Static Analysis, 15th International Symposium, SAS 2008, Valencia, Spain, July 16-18, 2008. Proceedings. Ed. by M. Alpuente and G. Vidal. Vol. 5079. Lecture Notes in Computer Science. Springer, 2008, pp. 1–17. DOI: 10.1007/978-3-540-69166-2\_1.



- [77] R. Giacobazzi and E. Quintarelli. “Incompleteness, Counterexamples, and Refinements in Abstract Model-Checking”. In: *Static Analysis, 8th International Symposium, SAS 2001, Paris, France, July 16-18, 2001, Proceedings*. Ed. by P. Cousot. Vol. 2126. *Lecture Notes in Computer Science*. Springer, 2001, pp. 356–373. DOI: 10.1007/3-540-47764-0\\_20.
- [78] R. Giacobazzi and F. Ranzato. “Functional Dependencies and Moore-Set Completions of Abstract Interpretations and Semantics”. In: *Logic Programming, Proceedings of the 1995 International Symposium, Portland, Oregon, USA, December 4-7, 1995*. Ed. by J. W. Lloyd. MIT Press, 1995, pp. 321–335.
- [79] R. Giacobazzi and F. Ranzato. “Refining and Compressing Abstract Domains”. In: *Automata, Languages and Programming, 24th International Colloquium, ICALP’97, Bologna, Italy, 7-11 July 1997, Proceedings*. 1997, pp. 771–781. DOI: 10.1007/3-540-63165-8\\_230.
- [80] R. Giacobazzi, F. Ranzato, and F. Scozzari. “Making Abstract Interpretations Complete”. In: *J. ACM* 47.2 (2000), pp. 361–416. DOI: 10.1145/333979.333989.
- [81] M. L. P. Gort, E. A. Díaz, and C. F. Uribe. “A Highly-Reliable Virtual Primary Key Scheme for Relational Database Watermarking Techniques”. In: *2017 International Conference on Computational Science and Computational Intelligence (CSCI)*. IEEE. 2017, pp. 55–60. DOI: 10.1109/CSCI.2017.10..
- [82] M. L. P. Gort, C. F. Uribe, and I. Nummenmaa. “A Minimum Distortion: High Capacity Watermarking Technique for Relational Data”. In: *Proceedings of the 5th ACM Workshop on Information Hiding and Multimedia Security*. ACM. 2017, pp. 111–121. DOI: 10.1145/3082031.3083241.
- [83] M. L. P. Gort et al. “HQR-Scheme: A High Quality and Resilient Virtual Primary Key Generation Approach for Watermarking Relational Data”. In: *Expert Systems with Applications* 138 (2019), p. 112770. ISSN: 0957-4174. DOI: 10.1016/j.eswa.2019.06.058.
- [84] M. L. P. Gort et al. “Preventing Additive Attacks to Relational Database Watermarking”. In: *Research and Practical Issues of Enterprise Information Systems - 13th IFIP WG 8.9 International Conference, CONFENIS 2019, Prague, Czech Republic, December 16-17, 2019, Proceedings*. Ed. by P. Doucek et al. Vol. 375. *Lecture Notes in Business Information Processing*. Springer, 2019, pp. 131–140. DOI: 10.1007/978-3-030-37632-1\\_12.
- [85] M. L. P. Gort et al. “Semantic-driven Watermarking of Relational Textual Databases”. In: *Expert Systems with Applications* (2020), p. 114013. DOI: 10.1016/j.eswa.2020.114013.

- [86] C. Gould, Z. Su, and P. Devanbu. “JDBC Checker: a Static Analysis Tool for SQL/JDBC Applications”. In: Proceedings. 26th International Conference on Software Engineering. IEEE, 2004. DOI: 10.1109/ICSE.2004.1317494.
- [87] P. Granger. “Improving the Results of Static Analyses of Programs by Local Decreasing Iterations”. In: Shyamasundar R. (eds) Foundations of Software Technology and Theoretical Computer Science. FSTTCS 1992. Lecture Notes in Computer Science, vol 652. Springer, Berlin, Heidelberg. 1992. DOI: 0.1007/3-540-56287-7\_95.
- [88] S. Gulwani and A. Tiwari. “Combining Abstract Interpreters”. In: Proceedings of the 27th ACM SIGPLAN Conference on Programming Language Design and Implementation. New York, NY, USA: Association for Computing Machinery, 2006, 376–386. DOI: 10.1145/1133981.1134026.
- [89] J. Guo. “Fragile Watermarking Scheme for Tamper Detection of Relational Database”. In: 2011 International Conference on Computer and Management (CAMAN). IEEE. 2011, pp. 1–4. DOI: 10.1109/CAMAN.2011.5778907.
- [90] M. K. Gupta, M. C. Govil, and G. Singh. “Static Analysis Approaches to Detect SQL Injection and Cross Site Scripting Vulnerabilities in Web Applications: a Survey”. In: International Conference on Recent Advances and Innovations in Engineering (ICRAIE-2014). 2014, pp. 1–5. DOI: 10.1109/ICRAIE.2014.6909173.
- [91] S. Gupta and B. B. Gupta. “Cross-Site Scripting (XSS) Attacks and Defense Mechanisms: Classification and State-of-the-Art”. In: Int. J. Systems Assurance Engineering and Management 8.1 (2017), pp. 512–530. DOI: 10.1007/s13198-015-0376-0.
- [92] R. Halder, S. Pal, and A. Cortesi. “Watermarking Techniques for Relational Databases: Survey, Classification and Comparison”. In: J. Univers. Comput. Sci. 16.21 (2010), pp. 3164–3190. DOI: 10.3217/jucs-016-21-3164.
- [93] R. Halder, S. Pal, and A. Cortesi. “Watermarking Techniques for Relational Databases: Survey, Classification and Comparison”. In: J. UCS 16.21 (2010), pp. 3164–3190. DOI: 10.3217/jucs-016-21-3164.
- [94] W. G. J. Halfond and A. Orso. “AMNESIA: Analysis and Monitoring for NEutralizing SQL-injection Attacks”. In: 20th IEEE/ACM International Conference on Automated Software Engineering (ASE 2005), November 7-11, 2005, Long Beach, CA, USA. 2005, pp. 174–183. DOI: 10.1145/1101908.1101935.
- [95] W. G. J. Halfond, J. Viegas, and A. Orso. “A Classification of SQL Injection Attacks and Countermeasures”. 2006.
- [96] A. Hliaoutakis et al. “Information Retrieval by Semantic Similarity”. In: Int. J. Semantic Web Inf. Syst. 2.3 (2006), pp. 55–73. DOI: 10.4018/jswis.2006070104.
- [97] G. J. Holzmann. “UNO: Static Source Code Checking for UserDefined Properties”. In: 6th World Conf. on Integrated Design and Process Technology, IDPT '02. 2002.

- [98] J. M. Howe and A. King. “Logahedra: A New Weakly Relational Domain”. In: *Automated Technology for Verification and Analysis, 7th International Symposium, ATVA 2009, Macao, China, October 14-16, 2009. Proceedings.* 2009, pp. 306–320. DOI: 10.1007/978-3-642-04761-9\_23.
- [99] H. Seidl, R. Wilhelm, and S. Hack. “Compiler Design - Analysis and Transformation”. Springer, 2012. DOI: 10.1007/978-3-642-17548-0.
- [100] Y. W. Huang et al. “Securing Web Application Code by Static Analysis and Runtime Protection”. In: *Proceedings of the 13th international conference on World Wide Web, WWW 2004, New York, NY, USA, May 17-20, 2004.* 2004, pp. 40–52. DOI: 10.1145/988672.988679.
- [101] IEEE. “IEEE Standard for Software Verification and Validation Plans”. IEEE Std 1012-1986. 1986.
- [102] H. Illous, M. Lemerre, and X. Rival. “A Relational Shape Abstract Domain”. In: *NASA Formal Methods - 9th International Symposium, NFM 2017, Moffett Field, CA, USA, May 16-18, 2017, Proceedings.* 2017, pp. 212–229. DOI: 10.1007/978-3-319-57288-8\_15.
- [103] J. R. Ingram. “Digital Piracy”. In: *The Encyclopedia of Criminology and Criminal Justice* (2014). DOI: 10.1002/9781118517383.wbeccj116.
- [104] Z. Jalil and A. M. Mirza. “A Review of Digital Watermarking Techniques for Text Documents”. In: *2009 International Conference on Information and Multimedia Technology.* Dec. 2009, pp. 230–234. DOI: 10.1109/ICIMT.2009.11.
- [105] S. H. Jensen, A. Møller, and P. Thiemann. “Type Analysis for JavaScript”. In: *Static Analysis, 16th International Symposium, SAS 2009, Los Angeles, CA, USA, August 9-11, 2009. Proceedings.* 2009, pp. 238–255. DOI: 10.1007/978-3-642-03237-0\_17.
- [106] C. Jiang, X. Chen, and Z. Li. “Watermarking Relational Databases for Ownership Protection Based on DWT”. In: *Proceedings of the Fifth International Conference on Information Assurance and Security, IAS 2009, Xi’An, China, 18-20 August 2009.* IEEE Computer Society, 2009, pp. 305–308. DOI: 10.1109/IAS.2009.220.
- [107] R. W. M. Jones and P. H. J. Kelly. “Backwards-Compatible Bounds Checking for Arrays and Pointers in C Programs”. In: *Proceedings of the Third International Workshop on Automated Debugging, AADEBUG 1997, Linköping, Sweden, May 26-27, 1997. Vol. 2. Linköping Electronic Articles in Computer and Information Science 009.* Linköping University Electronic Press, 1997, pp. 13–26.

- [108] M. Journault, A. Miné, and A. Ouadjaout. “Modular Static Analysis of String Manipulations in C Programs”. In: *Static Analysis - 25th International Symposium, SAS 2018, Freiburg, Germany, August 29-31, 2018, Proceedings*. 2018, pp. 243–262. DOI: 10.1007/978-3-319-99725-4\\_16.
- [109] N. Jovanovic, C. Kruegel, and E. Kirda. “Pixy: a Static Analysis Tool for Detecting Web Application Vulnerabilities”. In: *2006 IEEE Symposium on Security and Privacy (S P’06)*. 2006, 6 pp.–263. DOI: 10.1109/SP.2006.29.
- [110] N. S. Kamaruddin et al. “A Review of Text Watermarking: Theory, Methods, and Applications”. In: *IEEE Access* 6 (2018), pp. 8011–8028. URL: <https://doi.org/10.1109/ACCESS.2018.2796585>.
- [111] M. Kamran and M. Farooq. “A Formal Usability Constraints Model for Watermarking of Outsourced Datasets”. In: *IEEE Trans. Inf. Forensics Secur.* 8.6 (2013), pp. 1061–1072. DOI: 10.1109/TIFS.2013.2259234.
- [112] V. Kashyap et al. “JSAI: a Static Analysis Platform for JavaScript”. In: *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, (FSE-22), Hong Kong, China, November 16 - 22, 2014*. 2014, pp. 121–132. DOI: 10.1145/2635868.2635904.
- [113] G. A. Kildall. “A Unified Approach to Global Program Optimization”. In: *Conference Record of the ACM Symposium on Principles of Programming Languages, Boston, Massachusetts, USA, 1973*, pp. 194–206. DOI: 10.1145/512927.512945.
- [114] H. Kim, K. Doh, and D. A. Schmidt. “Static Validation of Dynamically Generated HTML Documents Based on Abstract Parsing and Semantic Processing”. In: *Static Analysis - 20th International Symposium, SAS 2013, Seattle, WA, USA, June 20-22, 2013. Proceedings*. Ed. by F. Logozzo and M. Fähndrich. Vol. 7935. *Lecture Notes in Computer Science*. Springer, 2013, pp. 194–214. DOI: 10.1007/978-3-642-38856-9\\_12.
- [115] S. Kim et al. “Inferring Grammatical Summaries of String Values”. In: *Programming Languages and Systems - 12th Asian Symposium, APLAS 2014, Singapore, November 17-19, 2014, Proceedings*. 2014, pp. 372–391. DOI: 10.1007/978-3-319-12736-1\\_20.
- [116] D. A. Kindy and A. K. Pathan. “A Survey on SQL Injection: Vulnerabilities, Attacks, and Prevention Techniques”. In: *2011 IEEE 15th International Symposium on Consumer Electronics (ISCE)*. 2011, pp. 468–471. DOI: 10.1109/ISCE.2011.5973873.
- [117] J. C. King. “Symbolic Execution and Program Testing”. In: *Commun. ACM* 19.7 (1976), pp. 385–394. DOI: 10.1145/360248.360252.

- [118] E. Kneuss, P. Suter, and V. Kuncak. “Phantom: PHP analyzer for type mismatch”. In: Proceedings of the 18th ACM SIGSOFT International Symposium on Foundations of Software Engineering, 2010, Santa Fe, NM, USA, November 7-11, 2010. 2010, pp. 373–374. DOI: 10.1145/1882291.1882355.
- [119] C. Lattner and V. Adve. “LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation”. In: International Symposium on Code Generation and Optimization (CGO’04). Palo Alto, California, Mar. 2004. DOI: 10.1109/CGO.2004.1281665.
- [120] H. Lauko, P. Rockai, and J. Barnat. “Symbolic Computation via Program Transformation”. In: Theoretical Aspects of Computing - ICTAC 2018 - 15th International Colloquium, Stellenbosch, South Africa, October 16-19, 2018, Proceedings. 2018, pp. 313–332. DOI: 10.1007/978-3-030-02508-3\_17.
- [121] H. Lauko et al. “Abstracting Strings for Model Checking of C Programs”. In: Appl. Sci., 10(21), 7853; (2020). DOI: 10.3390/app10217853.
- [122] H. Lee et al. “SAFE: Formal Specification and Implementation of a Scalable Analysis Framework for ECMAScript”. In: Proceedings of the 19th International Workshop on Foundations of Object-Oriented Languages (FOOL’12). 2012.
- [123] T. Liang et al. “An Efficient SMT Solver for String Constraints”. In: Formal Methods in System Design 48.3 (2016), pp. 206–234. DOI: 10.1007/s10703-016-0247-6.
- [124] V. B. Livshits and M. S. Lam. “Finding Security Vulnerabilities in Java Applications with Static Analysis”. In: Proceedings of the 14th USENIX Security Symposium, Baltimore, MD, USA, July 31 - August 5, 2005. Ed. by P. D. McDaniel. USENIX Association, 2005.
- [125] F. Logozzo. “Towards a Quantitative Estimation of Abstract Interpretations”. In: Workshop on Quantitative Analysis of Software. Microsoft, June 2009.
- [126] F. Logozzo and M. Fähndrich. “Pentagons: a Weakly Relational Abstract Domain for the Efficient Validation of Array Accesses”. In: Sci. Comput. Program. 75.9 (2010), pp. 796–807. DOI: 10.1016/j.scico.2009.04.004.
- [127] M. Madsen and E. Andreasen. “String Analysis for Dynamic Field Access”. In: Compiler Construction - 23rd International Conference, CC 2014, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2014, Grenoble, France, April 5-13, 2014. Proceedings. 2014, pp. 197–217. DOI: 10.1007/978-3-642-54807-9\_12.
- [128] S. Maffeis, J. C. Mitchell, and A. Taly. “An Operational Semantics for JavaScript”. In: Programming Languages and Systems, 6th Asian Symposium, AP-LAS 2008, Bangalore, India, December 9-11, 2008. Proceedings. 2008, pp. 307–325. DOI: 10.1007/978-3-540-89330-1\_22.

- [129] MathWorks. “Polyspace”. <http://www.mathworks.com/products/polyspace.html>. Last check: November 23rd, 2020.
- [130] B. B. Mehta and H. D. Aswar. “Watermarking for Security in Database: a Review”. In: 2014 Conference on IT in Business, Industry and Government (CSIBIG) (2014), pp. 1–6. DOI: 10.1109/CSIBIG.2014.7056938.
- [131] S. Melkundi and C. Chandankhede. “A Robust Technique for Relational Database Watermarking and Verification”. In: Communication, Information & Computing Technology (ICCICT), 2015 International Conference on. IEEE. 2015, pp. 1–7. DOI: 10.1109/ICCICT.2015.7045676.
- [132] G. A. Miller. “WordNet: An Electronic Lexical Database”. MIT press, 1998.
- [133] Y. Minamide. “Static Approximation of Dynamically Generated Web Pages”. In: Proceedings of the 14th international conference on World Wide Web, WWW 2005, Chiba, Japan, May 10-14, 2005. Ed. by A. Ellis and T. Hagino. ACM, 2005, pp. 432–441. DOI: 10.1145/1060745.1060809.
- [134] A. Miné. “Field-sensitive Value Analysis of Embedded C Programs with Union Types and Pointer Arithmetics”. In: Proceedings of the 2006 ACM SIGPLAN/SIGBED Conference on Languages, Compilers, and Tools for Embedded Systems (LCTES’06), Ottawa, Ontario, Canada, June 14-16, 2006. 2006, pp. 54–63. DOI: 10.1145/1134650.1134659.
- [135] A. Miné. “The Octagon Abstract Domain”. In: High. Order Symb. Comput. 19.1 (2006), pp. 31–100. DOI: 10.1007/s10990-006-8609-1.
- [136] A. Møller and M. I. Schwartzbach. “Static Program Analysis”. Department of Computer Science, Aarhus University, 2015.
- [137] F. Nielson. “Tensor Products Generalize the Relational Data Flow Analysis Method”. In: Fourth Hungarian Computer Science Conference. 1985, pp. 211–225.
- [138] F. Nielson and H. R. Nielson. “Type and Effect Systems”. In: Correct System Design, Recent Insight and Advances, (to Hans Langmaack on the occasion of his retirement from his professorship at the University of Kiel). Ed. by E. R. Olderog and B. Steffen. Vol. 1710. Lecture Notes in Computer Science. Springer, 1999, pp. 114–136. DOI: 10.1007/3-540-48092-7\_6.
- [139] F. Nielson, H. R. Nielson, and C. Hankin. “Principles of program analysis”. Springer, 1999. DOI: 10.1007/978-3-662-03811-6.
- [140] Aleph One. “Smashing The Stack For Fun And Profit”. Phrack Magazine. 1996.
- [141] OWASP. “OWASP Top Ten”. <https://owasp.org/www-project-top-ten/>. Last check: November 25th, 2020.

- [142] C. Park, H. Im, and S. Ryu. “Precise and Scalable Static Analysis of jQuery Using a Regular Expression Domain”. In: Proceedings of the 12th Symposium on Dynamic Languages, DLS 2016, Amsterdam, The Netherlands, November 1, 2016. 2016, pp. 25–36. DOI: 10.1145/2989225.2989228.
- [143] E. G. M. Petrakis et al. “X-Similarity: Computing Semantic Similarity between Concepts from Different Ontologies”. In: J. Digit. Inf. Manag. 4.4 (2006), pp. 233–237.
- [144] A. Di Pierro and H. Wiklicky. “Measuring the Precision of Abstract Interpretations”. In: Logic Based Program Synthesis and Transformation, 10th International Workshop, LOPSTR 2000 London, UK, July 24-28, 2000, Selected Papers. 2000, pp. 147–164. DOI: 10.1007/3-540-45142-0\9.
- [145] M. Pradel and K. Sen. “The Good, the Bad, and the Ugly: An Empirical Study of Implicit Type Conversions in JavaScript”. In: 29th European Conference on Object-Oriented Programming, ECOOP 2015, July 5-10, 2015, Prague, Czech Republic. Ed. by J. T. Boyland. Vol. 37. LIPIcs. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2015, pp. 519–541. DOI: 10.4230/LIPIcs.ECOOP.2015.519.
- [146] M. Dalla Preda and M. Pasqua. “Software Watermarking: A Semantics-based Approach”. In: Electr. Notes Theor. Comput. Sci. 331 (2017), pp. 71–85. DOI: 10.1016/j.entcs.2017.02.005.
- [147] Princeton-University. “About WordNet. WordNet. Princeton University”. 2010. URL: <http://wordnet.princeton.edu>.
- [148] H. G. Rice. “Classes of Recursively Enumerable Sets and Their Decision Problems”. In: *Transactions of the American Mathematical Society* 74.2 (1953), pp. 358–366.
- [149] X. Rival and L. Mauborgne. “The Trace Partitioning Abstract Domain”. In: ACM Trans. Program. Lang. Syst. 29.5 (2007), p. 26. DOI: 10.1145/1275497.1275501.
- [150] H. Samimi et al. “Automated Repair of HTML Generation Errors in PHP Applications Using String Constraint Solving”. In: 34th International Conference on Software Engineering, ICSE 2012, June 2-9, 2012, Zurich, Switzerland. 2012, pp. 277–287. DOI: 10.1109/ICSE.2012.6227186.
- [151] H. M. Sardroudi and S. Ibrahim. “A New Approach for Relational Database Watermarking Using Image”. In: 5th International Conference on Computer Sciences and Convergence Information Technology. IEEE. 2010, pp. 606–610. DOI: 10.1109/ICCIT.2010.5711126.
- [152] P. Saxena et al. “A Symbolic Execution Framework for JavaScript”. In: 31st IEEE Symposium on Security and Privacy, S&P 2010, 16-19 May 2010, Berkeley/Oakland, California, USA. 2010, pp. 513–528. URL: <https://doi.org/10.1109/SP.2010.38>.

- [153] D. A. Schmidt. “Denotational Semantics: A Methodology for Language Development”. USA: William C. Brown Publishers, 1986. ISBN: 0697068492.
- [154] R. C. Seacord. “Secure Coding in C and C++”. 2nd. Addison-Wesley Professional, 2013. ISBN: 0321822137.
- [155] N. Seco, T. Veale, and J. Hayes. “An Intrinsic Information Content Metric for Semantic Similarity in WordNet”. In: Proceedings of the 16th European Conference on Artificial Intelligence, ECAI’2004, including Prestigious Applicants of Intelligent Systems, PAIS 2004, Valencia, Spain, August 22-27, 2004. Ed. by R. L. De Mántaras and L. Saitta. IOS Press, 2004, pp. 1089–1090.
- [156] H. Shahriar and M. Zulkernine. “Classification of Static Analysis-Based Buffer Overflow Detectors”. In: Fourth International Conference on Secure Software Integration and Reliability Improvement, SSIRI 2010, Singapore, June 9-11, 2010 - Companion Volume. 2010, pp. 94–101. DOI: 10.1109/SSIRI-C.2010.28.
- [157] J. Shanmugam and M. Ponnaivaikko. “Cross Site Scripting-Latest Developments and Solutions: a Survey”. In: Int. J. Open Problems Compt. Math 1 (2008).
- [158] P. Singh and R. S. Chadha. “A Survey of Digital Watermarking Techniques, Applications and Attacks”. In: 2013.
- [159] T. Slimani. “Description and Evaluation of Semantic Similarity Measures Approaches”. In: International Journal of Computer Applications 80.10 (Oct. 2013), 25–33.
- [160] P. Sotin. “Quantifying the Precision of Numerical Abstract Domains”. 2010.
- [161] A. Sotirov. “Automatic Vulnerability Detection Using Static Source Code Analysis”. 2005.
- [162] F. Spoto. “The Julia Static Analyzer for Java”. In: Static Analysis - 23rd International Symposium, SAS 2016, Edinburgh, UK, September 8-10, 2016, Proceedings. Ed. by X. Rival. Vol. 9837. Lecture Notes in Computer Science. Springer, 2016, pp. 39–57. DOI: 10.1007/978-3-662-53413-7\\_3.
- [163] SQLi. “SQL Injection”. [https://www.owasp.org/index.php/SQL\\_Injection](https://www.owasp.org/index.php/SQL_Injection). Last check: December 18th, 2018.
- [164] Z. Su and G. Wassermann. “The Essence of Command Injection Attacks in Web Applications”. In: Proceedings of the 33rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2006, Charleston, South Carolina, USA, January 11-13, 2006. 2006, pp. 372–382. DOI: 10.1145/1111037.1111070.
- [165] J. Sun, Z. Cao, and Z. Hu. “Multiple Watermarking Relational Databases Using Image”. In: MultiMedia and Information Technology, 2008. MMIT’08. International Conference on. IEEE. 2008, pp. 373–376. DOI: 10.1109/MMIT.2008.211.



- [166] M. A. H. Taieb, M. B. Aouicha, and A. B. Hamadou. “A New Semantic Relatedness Measurement Using WordNet Features”. In: *Knowl. Inf. Syst.* 41.2 (2014), pp. 467–497. DOI: 10.1007/s10115-013-0672-4.
- [167] T. Tateishi, M. Pistoia, and O. Tripp. “Path- and Index-sensitive String Analysis Based on Monadic Second-order Logic”. In: *ACM Trans. Softw. Eng. Methodol.* 22.4 (2013), 33:1–33:33. DOI: 10.1145/2522920.2522926.
- [168] TIOBE. “TIOBE Index for November 2020”. <https://www.tiobe.com/tiobe-index/>. Last check: November 25th, 2020.
- [169] U. Topkara, M. Topkara, and M. J. Atallah. “The Hiding Virtues of Ambiguity: Quantifiably Resilient Watermarking of Natural Language Text Through Synonym Substitutions”. In: *Proceedings of the 8th workshop on Multimedia & Security, MM&Sec 2006, Geneva, Switzerland, September 26-27, 2006*. Ed. by S. Voloshynovskiy, J. Dittmann, and J. J. Fridrich. ACM, 2006, pp. 164–174. DOI: 10.1145/1161366.1161397.
- [170] O. Tripp, P. Ferrara, and M. Pistoia. “Hybrid Security Analysis of Web JavaScript Code Via Dynamic Partial Evaluation”. In: *International Symposium on Software Testing and Analysis, ISSTA '14, San Jose, CA, USA - July 21 - 26, 2014*. 2014, pp. 49–59. DOI: 10.1145/2610384.2610385.
- [171] F. Vasilescu, P. Langlais, and G. Lapalme. “Evaluating Variants of the Lesk Approach for Disambiguating Words”. In: *Proceedings of the Fourth International Conference on Language Resources and Evaluation, LREC 2004, May 26-28, 2004, Lisbon, Portugal*. European Language Resources Association, 2004.
- [172] M. Veanes, P. De Halleux, and N. Tillmann. “Rex: Symbolic Regular Expression Explorer”. In: *Third International Conference on Software Testing, Verification and Validation, ICST 2010, Paris, France, April 7-9, 2010*. 2010, pp. 498–507. DOI: 10.1109/ICST.2010.15.
- [173] J. Viega et al. “Token-based Scanning of Source Code for Security Problems”. In: *ACM Trans. Inf. Syst. Secur.* 5.3 (2002), pp. 238–261. DOI: 10.1145/545186.545188.
- [174] P. Vogt et al. “Cross Site Scripting Prevention with Dynamic Data Tainting and Static Analysis”. In: *Proceedings of the Network and Distributed System Security Symposium, NDSS 2007, San Diego, California, USA, 28th February - 2nd March 2007*. The Internet Society, 2007.
- [175] D. A. Wagner et al. “A First Step Towards Automated Detection of Buffer Overrun Vulnerabilities”. In: *Proceedings of the Network and Distributed System Security Symposium, NDSS 2000, San Diego, California, USA*. The Internet Society, 2000.

- [176] G. Wassermann and Z. Su. “Sound and Precise Analysis of Web Applications for Injection Vulnerabilities”. In: Proceedings of the ACM SIGPLAN 2007 Conference on Programming Language Design and Implementation, San Diego, California, USA, June 10-13, 2007. 2007, pp. 32–41. DOI: 10.1145/1250734.1250739.
- [177] G. Wassermann and Z. Su. “Static Detection of Cross-Site Scripting Vulnerabilities”. In: 30th International Conference on Software Engineering (ICSE 2008), Leipzig, Germany, May 10-18, 2008. 2008, pp. 171–180. DOI: 10.1145/1368088.1368112.
- [178] G. Wassermann et al. “Static Checking of Dynamically Generated Queries in Database Applications”. In: ACM Trans. Softw. Eng. Methodol. 16.4 (2007). DOI: 10.1145/1276933.1276935.
- [179] R. Wilhelm, S. Sagiv, and T. W. Reps. “Shape Analysis”. In: Compiler Construction, 9th International Conference, CC 2000, Held as Part of the European Joint Conferences on the Theory and Practice of Software, ETAPS 2000, Berlin, Germany, March 25 - April 2, 2000, Proceedings. 2000, pp. 1–17. DOI: 10.1007/3-540-46423-9\_1.
- [180] K. Winstein. “Lexical Steganography through Adaptive Modulation of the Word Choice Hash”. Secondary education at the Illinois Mathematics and Science Academy. 1999.
- [181] Y. Xie, A. Chou, and D. R. Engler. “ARCHER: Using Symbolic, Path-sensitive Analysis to Detect Memory Access Errors”. In: Proceedings of the 11th ACM SIGSOFT Symposium on Foundations of Software Engineering 2003 held jointly with 9th European Software Engineering Conference, ESEC/FSE 2003, Helsinki, Finland, September 1-5, 2003. 2003, pp. 327–336. DOI: 10.1145/940071.940115.
- [182] XSS. “Cross Site Scripting”. [https://www.owasp.org/index.php/Cross-site\\_Scripting\\_\(XSS\)](https://www.owasp.org/index.php/Cross-site_Scripting_(XSS)). Last check: November 23rd, 2020.
- [183] R. Xu, P. Godefroid, and R. Majumdar. “Testing for Buffer Overflows with Length Abstraction”. In: Proceedings of the ACM/SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2008, Seattle, WA, USA, July 20-24, 2008. 2008, pp. 27–38. DOI: 10.1145/1390630.1390636.
- [184] F. Yu, M. Alkhalaf, and T. Bultan. “Patching Vulnerabilities with Sanitization Synthesis”. In: Proceedings of the 33rd International Conference on Software Engineering, ICSE 2011, Waikiki, Honolulu , HI, USA, May 21-28, 2011. 2011, pp. 251–260. DOI: 10.1145/1985793.1985828.
- [185] F. Yu, T. Bultan, and B. Hardekopf. “String Abstractions for String Verification”. In: Model Checking Software - 18th International SPIN Workshop, Snowbird, UT, USA, July 14-15, 2011. Proceedings. 2011, pp. 20–37. DOI: 10.1007/978-3-642-22306-8\_3.

- [186] F. Yu, T. Bultan, and O. H. Ibarra. “Symbolic String Verification: Combining String Analysis and Size Analysis”. In: Tools and Algorithms for the Construction and Analysis of Systems, 15th International Conference, TACAS 2009, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2009, York, UK, March 22-29, 2009. Proceedings. Vol. 5505. Lecture Notes in Computer Science. Springer, 2009, pp. 322–336. DOI: 10.1007/978-3-642-00768-2\\_28.
- [187] F. Yu et al. “Automata-based Symbolic String Analysis for Vulnerability Detection”. In: Formal Methods Syst. Des. 44.1 (2014), pp. 44–70. DOI: 10.1007/s10703-013-0189-1.
- [188] F. Yu et al. “Symbolic String Verification: An Automata-Based Approach”. In: Model Checking Software, 15th International SPIN Workshop, Los Angeles, CA, USA, August 10-12, 2008, Proceedings. 2008, pp. 306–324. DOI: 10.1007/978-3-540-85114-1\\_21.
- [189] L. Zhang et al. “Relational databases watermarking for textual and numerical data”. In: Aug. 2011, pp. 1633–1636. DOI: 10.1109/MEC.2011.6025791.
- [190] Y. Zhang et al. “A Method of Verifying Relational Databases Ownership with Image Watermark”. In: The 6th International Symposium on Test and Measurement, Dalian, PR China. 2005, pp. 6316–6319.
- [191] Z. Zhang et al. “Watermarking Relational Database Using Image”. In: Proceedings of 2004 International Conference on Machine Learning and Cybernetics (IEEE Cat. No. 04EX826). Vol. 3. IEEE. 2004, pp. 1739–1744. DOI: 10.1109/ICMLC.2004.1382056.
- [192] X. Zhou et al. “Security Theory and Attack Analysis for Text Watermarking”. In: 2009 International Conference on E-Business and Information System Security. May 2009, pp. 1–6. DOI: 10.1109/EBISS.2009.5138072.