



Università  
Ca' Foscari  
Venezia

Master's Degree programme – Second Cycle  
(*D.M. 270/2004*)  
in Informatica - Computer Science

Final Thesis

—  
Ca' Foscari  
Dorsoduro 3246  
30123 Venezia

# Comparing metabolic networks at a global level

**Supervisor**  
Prof. Nicoletta Cocco

**Assistant Supervisor**  
Prof. Sabrina Manente

**Graduand**  
Gianluca Erbosio  
834510

**Academic Year**  
2015 / 2016



# Acknowledgements

I would like to express my special appreciation to my supervisor Prof. Cocco Nicoletta and to Prof. Simeoni Marta. Their support, useful comments and remarks have been essential in writing this thesis. I would also like to thank Prof. Poli Irene, Manente Sabrina and Bocci Martina for their advices in the experiments and for helping me to understand biological concepts.

I am also thankful to my classmates Antonio, Enrico, Alberto and Stefano. I have enjoyed my time with all of them and together, we have passed the most difficult moments during the university studies supporting each other.

Finally, a special thanks to my family, my mother, my father and my brothers who have supported me in these five years both economically and psychologically.



# Abstract

The metabolic networks analysis is an important process to discover similarities and differences between metabolisms of different organisms. When comparing metabolic networks we need to take into account the computational complexity due to the large size of such nets. This thesis has been developed together with [1], with the aim to propose a method for the comparison of metabolic networks which overcomes the problem of complexity. We propose to abstractly represent the metabolism as a graph connecting the metabolic functions, considering uniquely KEGG database information. To compare such representations we define two indices on structural similarity and two global indices which consider both structure and functionalities of such networks. The proposed methods have been implemented in Java and integrated in a tool for metabolic networks comparison. Finally, we discuss some experiments performed with our tool in order to evaluate the results.



# Index

<b>Acknowledgements</b>	<b>i</b>
<b>Abstract</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Metabolism &amp; KEGG Data Base</b>	<b>5</b>
2.1 Metabolism in different organisms . . . . .	5
2.2 Databases for metabolisms . . . . .	9
2.3 KEGG Database . . . . .	10
2.3.1 Informative contents and data organization . . . . .	11
2.3.2 KEGG Metabolism . . . . .	13
2.3.3 KEGG Pathways . . . . .	14
2.3.3.1 KGML . . . . .	19
2.3.3.2 Pathway structure in the XML file . . . . .	21
2.3.3.3 KEGG API . . . . .	23
2.3.3.4 Problems with data . . . . .	24
<b>3 Comparison between metabolism in different organisms</b>	<b>27</b>
3.1 Metabolism representation: State of the Art . . . . .	27
3.2 Comparison between metabolic Networks . . . . .	30
<b>4 Metabolic Networks comparison</b>	<b>37</b>
4.1 Metabolic network construction . . . . .	37

4.1.1	Network construction . . . . .	37
4.1.2	Implementation . . . . .	39
4.1.3	Data structures . . . . .	43
4.2	Comparison of metabolic networks . . . . .	45
4.2.1	Similarity between metabolic networks . . . . .	45
4.2.2	Global similarity indexes . . . . .	47
<b>5</b>	<b>Tool</b>	<b>53</b>
5.1	Requirements analysis . . . . .	53
5.1.1	Functional requirements . . . . .	54
5.1.2	Non-functional requirements . . . . .	55
5.2	Project architecture . . . . .	56
5.3	Libraries and technologies . . . . .	58
5.4	Documentation . . . . .	59
<b>6</b>	<b>Experimenting with the tool</b>	<b>65</b>
6.1	Cluster analysis . . . . .	65
6.2	Experiments . . . . .	68
6.2.1	Experiment 1: Metabolic evolution in a group of species . . .	68
6.2.2	Experiment 2: Yeasts and Molds metabolism . . . . .	70
6.2.3	Experiment 3: Sulfur metabolism in different Kingdoms . . . .	73
6.2.4	Experiment 4: Carbon fixation in photosynthetic organisms .	78
6.2.5	Experiment 5: Glycolysis metabolism . . . . .	82
6.2.6	Conclusion . . . . .	84
<b>7</b>	<b>Conclusion</b>	<b>85</b>
	<b>Bibliography</b>	<b>87</b>

# Chapter 1

## Introduction

In biology, metabolic networks comparison is relevant in studying the evolutionary process, finding similarities or dissimilarities between species, discovering drug targets and more in general in supporting medical science activities.

Generally, in the literature metabolic networks are expressed through graphs, hypergraphs or set structures. The first two representations allow to handle information with a great level of details but during their comparison some simplifications must be performed in order to reduce the complexity. Comparison of metabolic networks modelled as graphs is challenging from a computational point of view. The resolution of the subgraphs isomorphism problem, that represents a NP-complete problem, is infeasible for the huge dimensions of the metabolic nets. The set representation, instead, allows to compare metabolic networks in a easier way but it processes information with much less detail.

This thesis has been developed together with [1].

The aim of both theses, is to propose a method for the comparison of metabolic networks represented as graphs while avoiding the computational problem described above. In our approach this is achieved representing the metabolic network in two distinct levels. At the lower level, we represent each metabolic function of the net as a set of reactions and at the higher level we represent an abstraction of the metabolic networks as a graph in which each node represents a specific metabolic function and the edges represent the relations among them. The independence of these two levels

allows us to define similarity indices on both levels.

In this thesis, we focus on the higher level of our representation, the entire metabolic network, and we propose two indices: the first one characterises the structural similarity of the same metabolic function (node) in two organisms and the second one represents the structural similarity of their entire metabolic networks. Such indexes are then combined with similarity indexes capturing on the metabolic pathways similarities defined in [1] in order to compute a global similarity value on the entire metabolisms based both on structures and on functionalities.

Our method has been implemented in a Java tool that allows one to compare the metabolism of two different organisms. The tool provides some user-friendly interfaces that guide the user during the comparison process. During its development, software engineering principles and parallelization techniques has been adopted. In particular the modularisation of the tool allows for extending of the tool itself with new comparison methods and new functionalities.

The tool has been used to perform some experiments in order to check the quality of the indices and to validate them.

The thesis is organized as follows.

In Chapter 2 we give a general overview of metabolism from a biological point of view. We introduce the KEGG knowledge base, we discuss its databases structures and the technical aspects used for data retrieval and we give some basic notions used in the next chapters.

In Chapter 3 we describe the state of art in metabolism representation and then we focus on the comparison techniques for metabolic networks found in literature.

In Chapter 4 we present the metabolic network construction describing the idea, its implementation and data structure. Then the similarity indices used for comparison are explained.

Chapter 5 describes the tool we developed to implement the proposed approach. The functional and non-functional requirements, the project architecture and the used technologies are discussed. Moreover a documentation is given.

In Chapter 6 we report some experiments performed with our similarity indices.

The results have been discussed using a hierarchical clustering algorithm.

Chapter 7 contains some remarks and possible future developments about the comparison methods and the tool.



# Chapter 2

## Metabolism & KEGG Data Base

In this chapter we want to introduce, briefly, the metabolism to understand what are the elements that interact in this complex process. Moreover, we give an overview of the databases for metabolisms and in particular we examine in depth KEGG database. We describe its structure, the information which is contained, the API used to retrieve them and the methods used to represent them. We focus our attention on the description of the KGML structure to understand how to extract fundamental information used for further analysis.

### 2.1 Metabolism in different organisms

*The **metabolism**[2, 3] is the network of all the chemical and physical reactions that take place within the cells of the organisms.*

The complex set of chemical transformations is responsible of the growth and survival of the cells and the organisms themselves. In general, the metabolism is composed of two different and fundamental phases:

- **catabolism**[3, 4]: it is composed by all the metabolic tasks that produce simpler substances, producing energy (ATP);
- **anabolism**[3, 4]: it is the opposite of the catabolism. It is composed by all the synthesis tasks that produce more complex organic molecules from simpler

ones, consuming the energy released from catabolism;

The metabolism is composed by many different interacting functions called metabolic pathways.

A **metabolic pathway**<sup>[4]</sup> is a sequence of reactions such that the product of a single reaction can be used as reagent for another one.

In this document we refer to metabolic pathway using equivalently the term pathway. A metabolic pathway has an associated function. As an example you may consider the Glicolysis pathway. An interesting aspect is that each function occur in a specific location in the cells. The synthesis of particular substances in distinct compartments requires mechanisms to transport these substances between compartments. For example, to move ATP generated in the mitochondria, to the cytosol where most of it is consumed, a transport of protein is necessary. Figure 2-1 shows some metabolic functions locations in a eukaryotic cell.

In a pathway we may find biochemical reactions through which, using particular enzymes, there is a catalization. A *catalization* is a chemical phenomenon where the speed of reactions undergoes changes for the intervention of one or more substances said catalysts. A *substratum*, a particular molecule where the enzyme operates, is transformed into a product used as substratum in the next step. The *reactants*, the *intermediates* and the *products*, respectively the substance consumed in the chemical reaction, the molecule formed from the reactants and that reacts further to produce the product, the final element of the reaction, are called *metabolites*. The quantitative connection between these elements is specified by the **stoichiometry**. Through *stoichiometric coefficient* it is possible to mathematically represent the reagents and products quantities involved in a reaction<sup>[5]</sup>. According to the classification of the metabolic phases, sequences of catabolic and anabolic reactions are called *catabolic pathways* and *anabolic pathways* respectively. Their continuous overlap forms a complex exchange system which is the basis of growth and survival of cells.

---

<sup>1</sup>Image from: Judith G. Voet, Donald Voet, Charlotte W. Pratt. In Fundamentals of Biochemistry: Life at the Molecular Level, page 442. John Wiley and Sons, 4 edition, 2012

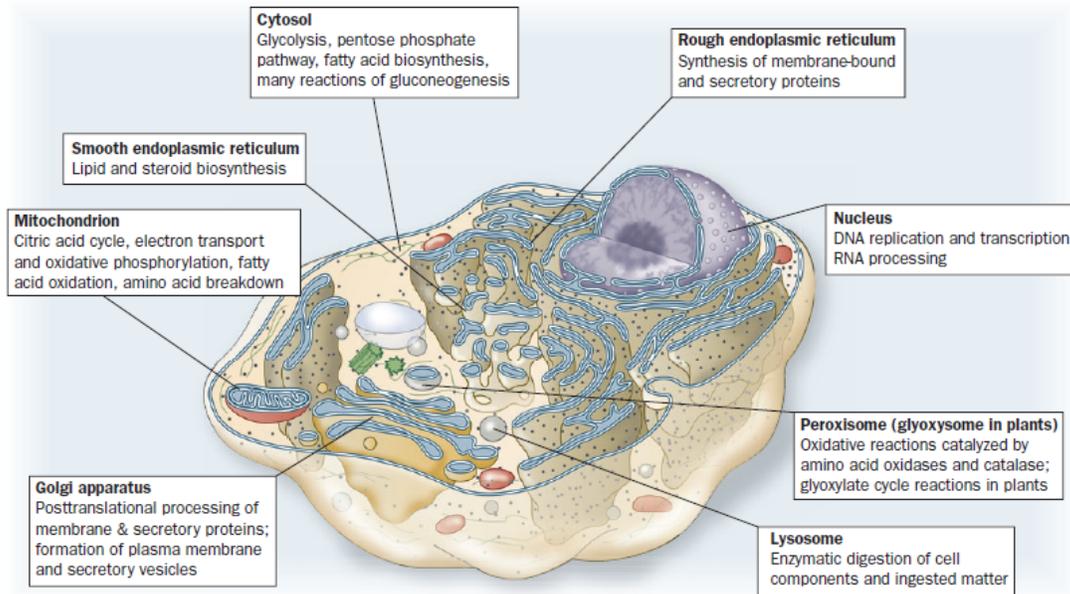


Figure 2-1: Metabolic pathways<sup>1</sup>

The **metabolic network** represents the complete set of metabolic functions and processes that determine the structure and properties of the cells. These functions are not independent but they interact each one with the others creating a more complex structure, representing the network.

In Figure 2-2 the entire metabolic network is represented from a general point of view. In the different living species the metabolism is similar in the components and in the organizations into metabolic pathways. The main metabolic pathways, such as Glycolysis or Citric acid cycle, for instance, probably appeared in one universal ancestor and they have been conserved during evolution because of their efficiency (ability to get to the final products with small number of steps). The analysis and comparison of metabolic networks between different organisms may yield important information on their evolution. Moreover, useful applications of such comparison are related to human disease analysis, drug design and metabolic engineering. However, some problems arise in metabolic networks comparison. Using a graph based modelling system, the resulting graph that represent a metabolic network may be composed by hundreds of nodes and thousands of edges. In this case, the graph matching may represent an infeasible computational problem. From graph theory in fact, the

exact graph matching problems, like isomorphism and sub-graphs isomorphism, are known to be NP-Complete problems.

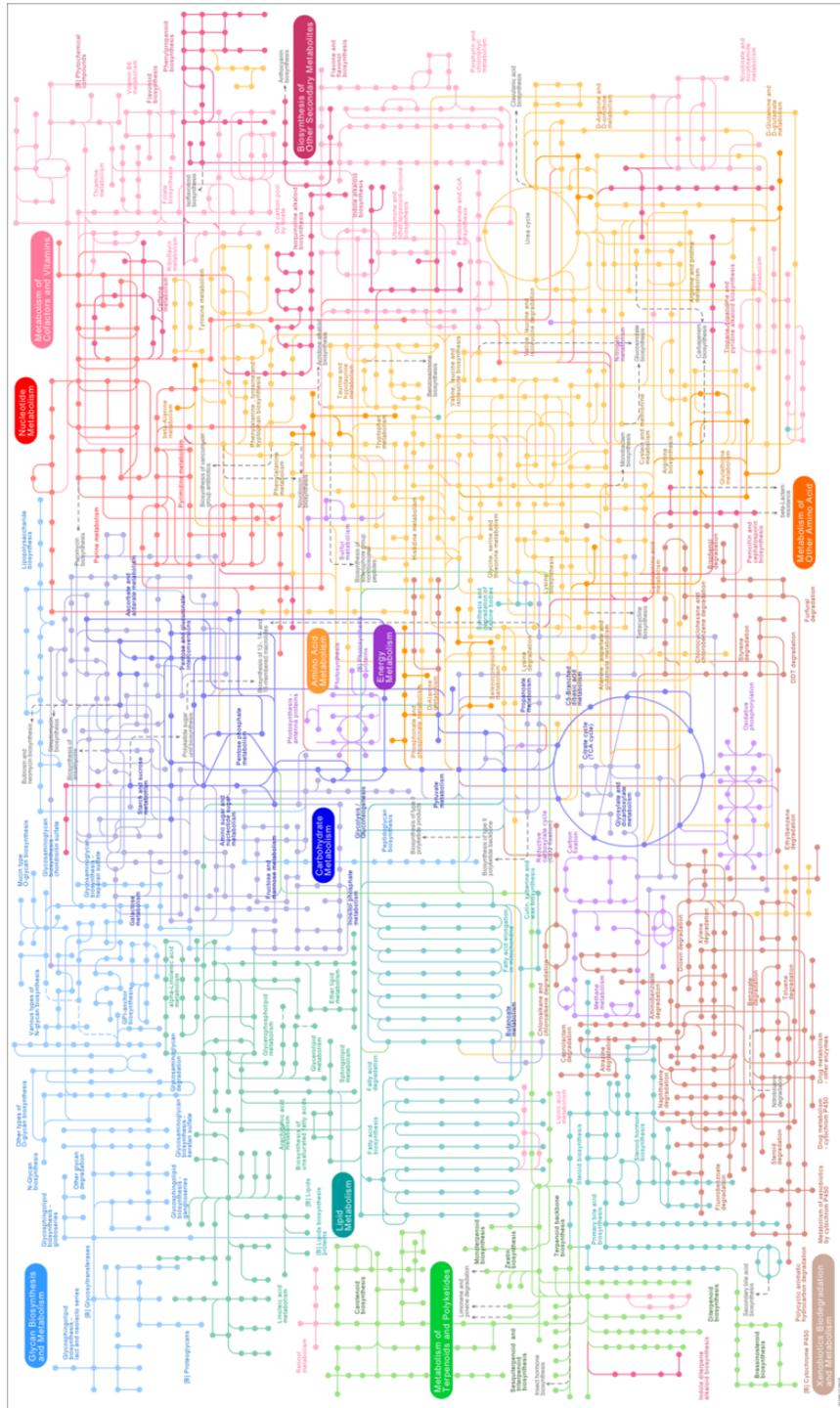


Figure 2-2: Metabolic Network<sup>2</sup>

<sup>2</sup>Image taken from: [http://www.kegg.jp/kegg-bin/show\\_pathway?map01100](http://www.kegg.jp/kegg-bin/show_pathway?map01100)

## 2.2 Databases for metabolisms

Since the early 80's to today, thanks to the development of the technology both in computer science and in biology, the number of databases for biological data is growing.

*A **database** is a collection of related data concerning a same topic that are stored and that can be used from applications.*[6]

In particular, a metabolism database contains data involving metabolities, enzymes, genes and reactions and information about their relations. In general, metabolic databases are incomplete due to the complexity in digitalization of such kind of data.

We can cite as the most used:

- KEGG[7]: a good description of this project is given in the next section of this chapter;
- BioCyc[8]: it contains a collection of databases about metabolic pathways and the genome of thousand of organisms integrating information from other databases, such as protein features and Gene Ontology information. BioCyc databases are divided into three categories related to the manual update frequency. Tier 1 is the most frequent updated database, Tier 2 receives a moderate quantity of manual update and Tier 3 receives only computational updates. In particular Tier 1 contains the following databases: EcoCyc, HumanCyc, MetaCyc, AraCyc, LeishCyc, YeastCyc. The first two contains the entire genome of the Escherichia coli and Human organism, MetaCyc contains representative metabolism sample of more than 2600 organisms and the lasts three contains information about Arabidopsis thaliana, Leishmania major and Saccharomyces cerevisiae organism.
- SEED[9]: it is a project developed by the Fellowship for Interpretation of Genomes (FIG) with the aim to develop a comparative genomics enviroment. The curation of genomic data is performed through subsystems by an expert annotator across many genomes. A subsystem in this case is defined by a set

of functional related roles. Then, a metabolic pathway is represented by subsystems, as the collection of functional roles, creating complex class of proteins. The result given by the subsystems extracts a set of protein families (FIGfams). The latter in turn, create the core component of the RAST subsystem. RAST, or Rapid Annotation using Subsystem Technology, is a rapid and very accurate annotation technology that makes use of data and procedures provided from SEED framework. Therefore, this technology provides a good level of automation in high quality gene calling and functional annotation.

- BioModels[10]: it contains computational models of biological processes. These are collected from literature and they are integrated with other references and stored in a set of MySQL tables. It is divided into three categories: the curated one that contains the curated models, the non-curate one that contain models that cannot be curated or still not curated and the automatic generated one that contains models generated automatically from other databases.

## 2.3 KEGG Database

The KEGG[7, 11, 12] (Kyoto Encyclopedia of Genes and Genomes) Database was started in 1995 by Minoru Kanehisa, Professor at the Institute for Chemical Research at Kyoto University with the purpose to collect all the information on sequenced organisms. It is presently one of the most important collection of biological data, containing information on metabolic pathways, Genomic information, Chemical information and Health information of different organisms.

One of the main efforts of the KEGG project is to standardize gene annotations, providing functional information of cellular processes to genomics. In order to do that, all the available knowledge about systemic functions, biochemical pathways and other kinds of molecular interactions have been taken by hands and then reorganized in a computable way, creating a big digital knowledge base. As a result, KEGG becomes one of the reference knowledge bases for data integration and systematic interpretation of sequence data. KEGG project aims to provide and maintain a reliable

knowledge base, supporting basic research activities in biology. Moreover, thinking about the benefits that information technology has given and can give through digitalization, large-scale data organization and development of tools for data analysis, a natural evolution of the knowledge bases is expected. KEGG in fact, is being expanded exploiting data extraction coming from the use of applications and tools based on its database. Thus, a new kind of information is collected. The latter, is typically related to health information including human diseases, drugs and other health-related substances.

### 2.3.1 Informative contents and data organization

As we said before, the KEGG aim is to automatize the interpretation of biological information encoded in sequence data. The prediction of gene function is also a considered problem. In particular, prediction of gene function is treated like a reconstruction process of biological system functioning, starting from genes and their products. Understanding how genes and molecules interoperate defining a biological system is typically a critical task. Therefore, a good organization of data is necessary.

The data are stored in four different macro categories as we can see in table 2.1. Each of them contains some specific databases and in particular:

- *System information category* contains functional information on how molecules and genes interact (KEGG PATHWAY), functional hierarchies of biological entities and functional units for biological interpretation of genomes. Moreover some other kind of information are stored, for instance: cell cycle, membrane transport, and more in general, information about regulatory aspects of cells function.
- *Genomic information category* contains structural information about genomics for all different organisms, gene catalogs, complete genomes and ortholog groups.
- *Chemical information category* contains information about chemical compounds, enzymes, molecules and reactions.

Category	Database	Content	Color
Systems information	KEGG PATHWAY	KEGG pathway maps	
	KEGG BRITE	BRITE functional hierarchies	
	KEGG MODULE	KEGG modules of functional units	
Genomic information	KEGG ORTHOLOGY	KEGG Orthology (KO) groups	 
	KEGG GENOME	KEGG organisms with complete genomes	
	KEGG GENES	Gene catalogs of complete genomes	
	KEGG SSDB	Sequence similarity database for GENES	
Chemical information KEGG LIGAND	KEGG COMPOUND	Metabolites and other small molecules	
	KEGG GLYCAN	Glycans	
	KEGG REACTION	Biochemical reactions	
	KEGG RPAIR	Reactant pair chemical transformations	
	KEGG RCLASS	Reaction class defined by RPAIR	
	KEGG ENZYME	Enzyme nomenclature	
Health information KEGG MEDICUS	KEGG DISEASE	Human diseases	
	KEGG DRUG	Drugs	
	KEGG DGROUP	Drug groups	
	KEGG ENVIRON	Crude drugs and health-related substances	

Table 2.1: The KEGG database. Table taken from <http://www.kegg.jp/kegg/kegg1a.html>

- *Health information* contains health information including human diseases, drugs and other health-related substances.

All these information are organized and represented in a big wiring diagrams called *Reference Pathway* that constitutes the core of the resource. In Figure 2-2 is represented the entire metabolic network generated by KEGG.

As we can see, the map identifies specific areas using different colours. In turn, each area corresponds to a specific metabolic pathway/function and integrate all the information available in the knowledge base, like interactions and reactions. Later we will see more details about maps, pathways etc.

In biology organisms are divided into categories that compose a taxonomic hierarchy built on different levels. From the top to the bottom the levels are the following: domain, kingdom, phylum, class, order, family, genus and species. Each level inherits the features from the upper one and it adds others to classify more precisely the organism. KEGG database recalls this organization, in fact it is split into macro-categories divided into more smaller ones. Each organism belongs to a specific category. The

mains are:

- *Eukaryotes* divided into animals, plants, fungi and protists;
- *Prokaryotes* divided into bacteria and archaea.

At present, in the database we find 313 Eukaryotes and 3562 Bacteria plus 215 Archea for the Prokaryotes. All these organisms are collected as complete genomes.

It is important to underline that KEGG is a freely accessible resource, that is constantly updated by the staff. So, we may see over the time, changes of data on the basis of new scientific discoveries and integrations.

### 2.3.2 KEGG Metabolism

As depicted in the figure 2-2, there are different metabolisms, each of them includes several distinct metabolic functions. In particular, we find information about: Carbohydrate metabolism, Energy metabolism, Lipid metabolism, Nucleotide metabolism, Amino acids metabolism, Glycan biosynthesis and metabolism, Metabolism of cofactors and vitamins, Metabolism of terpenoids and polyketides, Biosynthesis of other secondary metabolites, Xenobiotics biodegradation and metabolism, Chemical structure transformation maps.

KEGG subdivides the metabolic network into modules that represent the union of reference pathways. This structure does not constitute a partition over the network since each pathway can share parts of it with another one. It is known that the metabolic pathway are quite preserved among organisms and so, KEGG associates to each function, a unique reference pathway which corresponds to the union of the corresponding pathway in different organisms. It is possible to obtain a specific organism pathway from the corresponding reference one. The same concept is applied to the entire metabolic network: from the *reference network* that represent the union of all the reference pathways, it is possible to obtain a specific metabolic network for a choosen organism. In KEGG, the visual representation of an organism specific network, is given by highlighting the interested parts for the organism and by shading

all the rest. Graphically this approach gives to the users a rapid view of the overall functions present in a chosen organism and also the relations between them.

The metabolic network of homo sapiens is represented in Figure 2-3.

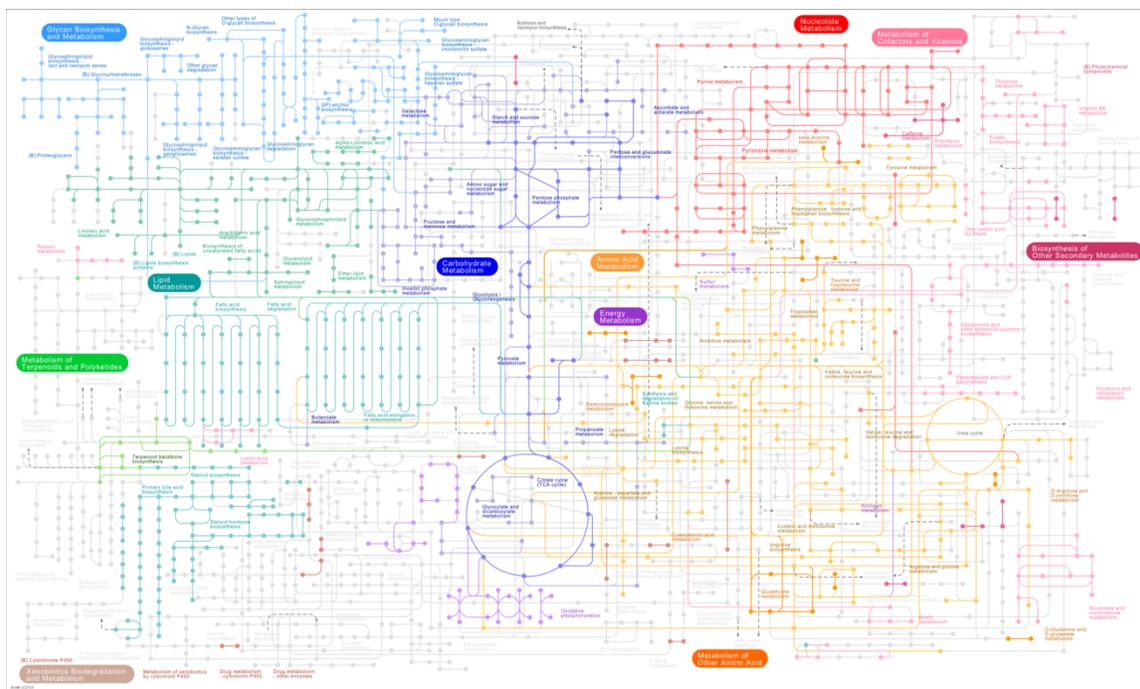


Figure 2-3: Metabolic Pathways of homo sapiens<sup>3</sup>

### 2.3.3 KEGG Pathways

KEGG PATHWAY is a collection of manually drawn diagrams and related textual informations. For each pathway we find a graphical representation typically called pathway map and a textual one written in XML format called KGML file. The latter representation is described in detail in Section 2.3.3.1. Every map graphically represents all the KEGG knowledge about molecular pathways for metabolism, adding important information about molecular interactions and reaction networks. Reaction networks are obtained through an integration of genetic information processing, environmental information processing, cellular processes, organismal systems, human diseases and drug development.

<sup>3</sup>The image is available at the following link: [http://www.kegg.jp/kegg-bin/show\\_pathway?org\\_name=hsa&mapno=01100&mapscale=0.35&show\\_description=hide&show\\_module\\_list=](http://www.kegg.jp/kegg-bin/show_pathway?org_name=hsa&mapno=01100&mapscale=0.35&show_description=hide&show_module_list=)

In the graphical representation, each map is composed by four different objects:

- **boxes** that identify gene products (enzymes);
- **circles** that represent other molecules, typically chemical compounds;
- **rectangles** for other maps representations;
- **lines** for molecular interactions.

Figures 2-4 and 2-5 summarize the entire notation for pathway map representation. They are taken from KEGG documentation<sup>4</sup>. Two kind of maps are provided for

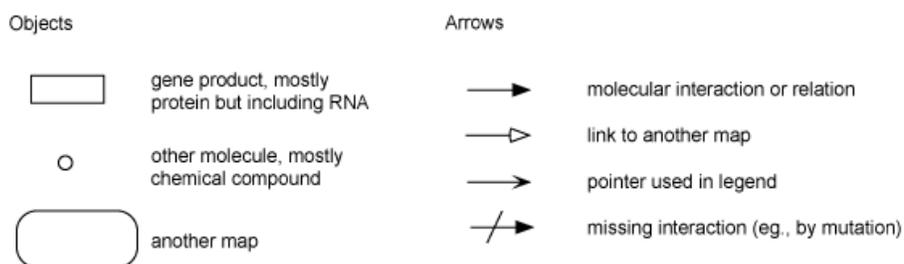


Figure 2-4: Symbols for map notation

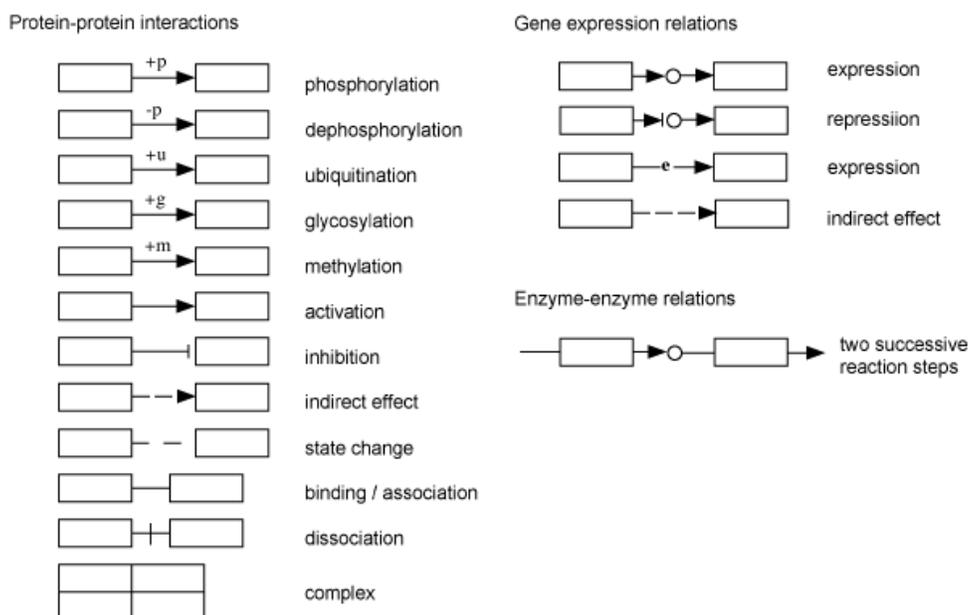


Figure 2-5: Different kind of relations in map Notation

<sup>4</sup>The symbology is available at the address [http://www.genome.jp/kegg/document/help\\_pathway.html](http://www.genome.jp/kegg/document/help_pathway.html)

each metabolic pathway: the reference pathway and organism-specific pathways. At present, there are 176 metabolic pathways, each one with its reference metabolic map, and for each of them we find more specific maps and KMGL files, one for each organism that includes the metabolic function under examination. A pathway has a unique identifier for reference map composed from the keyword *map* followed by a five digit number to distinguish the metabolic function. The organism-specific pathways instead, are identified by using *org* prefix, expressed with three or four letters that specify an organism. For example, the code *map00010* identifies the glycolysis reference pathway, while the code *hsa00010* identifies a more detailed map, glycolysis in the homo sapiens organism. From a technical point of view, reference pathways are the original manually drawn maps and they do not make use of colours. All the specific-organism maps are instead computationally generated, following some rules, from the corresponding reference pathway. Each component of an organism-specific pathway, in turn, has a unique identifier that satisfies one the following patterns:

- *ko* (*KEGG Orthology*) identifier, is expressed by *ko:* followed by *K* and a number of five digits that specifies the Ortholog<sup>5</sup> group for a specific organism. For example *ko:K01568* is the identifier for the pyruvate decarboxylase;
- *rn* identifier, expressed by *rn:* followed by *R* and a five digit number that identifies the reaction in the relative database. For instance: *rn:R00014* corresponds to the thiamin diphosphate acetaldehydetransferase (decarboxylating) in the pyruvate;
- *ec* identifier, expressed by *ec:* followed by an EC number defined by IUBMB-IUPAC commission that identifies a specific enzyme. As an example we can consider: *ec:4.1.1.1* that corresponds to the pyruvate decarboxylase;
- *cpd* identifier, expressed by *cpd:* followed by *C* and a five digit number that specifies a chemical compound. For example: *cpd:C00161* corresponds to the formula: *C2HO3R*.

---

<sup>5</sup>*Orthologs* are collection of genes belonging to different species evolved from a common ancestor. They preserve the same ancestor's function during the evolution process. For that reason, the identification of orthologs is quite critical for reliable prediction of gene function in newly sequenced genomes.

At the beginning, the KEGG project was based on an automatic matching between gene catalogs and enzymes in the reference map using EC<sup>6</sup> numbers. Now, the EC numbers are no longer used as identifiers in KEGG and the system was updated in order to use a different mapping criteria. Taking into account the computation of organism-specific pathway, it uses KEGG Orthology (KO)<sup>7</sup> system as the basis for genome annotation and mapping. Green boxes are linked to genes through a conversion of KO identifiers to gene identifiers in the reference pathway. Figures 2-6 and 2-7 show the reference pathway for the Citrate Cycle and the specific Citrate Cycle for Homo Sapiens organism.

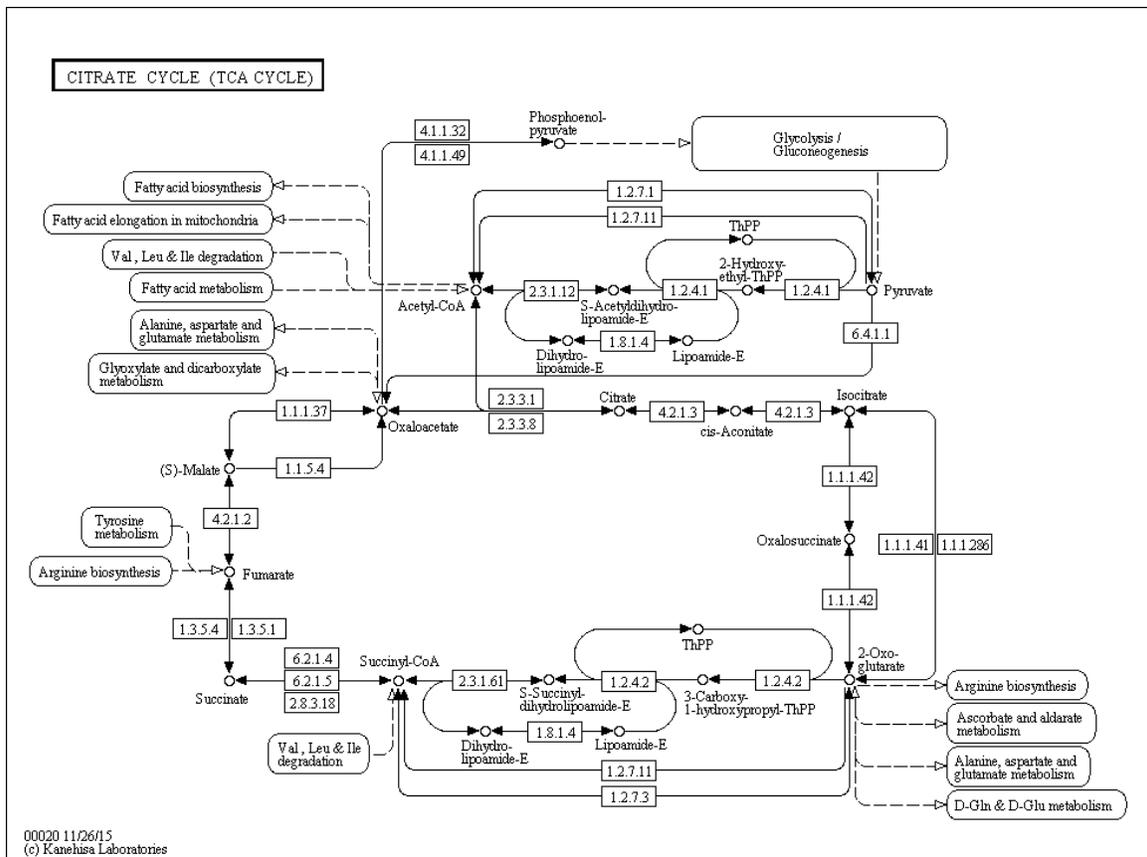


Figure 2-6: Citrate Cycle reference pathway<sup>8</sup>

<sup>6</sup>EC (Enzyme Classification) number: it provides a classification schema for enzymes. The classification is based on their catalizations and on chemical reactions. An EC number is expressed by a set of numbers separated by periods that progressively give a finer classification of the enzyme.

<sup>7</sup>The KEGG Orthology (KO) system is a collection of ortholog groups defined by hand, that capture experimental knowledge from literature and experimental observations. For more info about genome annotation and KO identifiers visit: <http://www.genome.jp/kegg/ko.html>

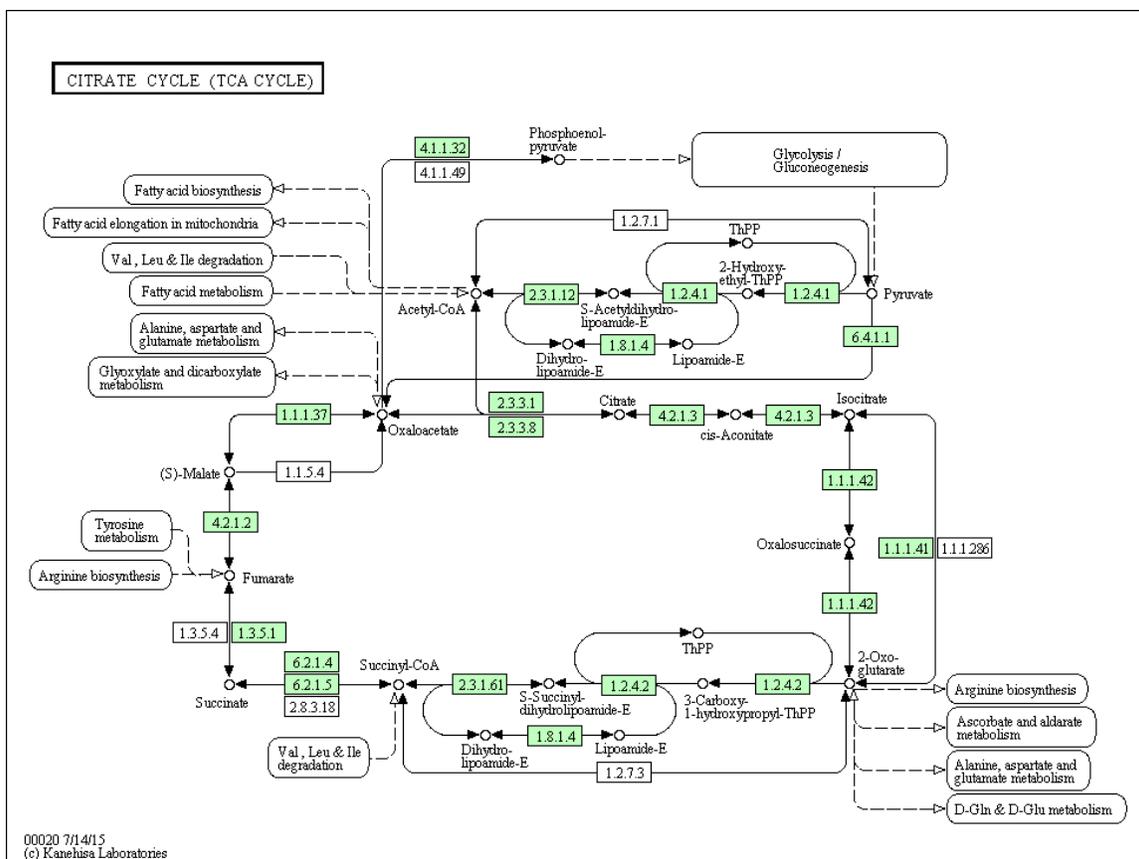


Figure 2-7: Citrate Cycle in the Homo Sapiens organism<sup>9</sup>

Drawing operations are performed using KegSketch software that produces a semi-static image where the map structure is fixed but each element can be coloured using user's preferences. Enzymes, maps and compounds inside a pathway map are all clickable objects and permit one to get more details on molecular structure.

The textual representation and the KGML file structure is described in the next section.

<sup>8</sup>Image taken from: [http://www.kegg.jp/kegg-bin/show\\_pathway?map=map00020&show\\_description=show](http://www.kegg.jp/kegg-bin/show_pathway?map=map00020&show_description=show)

<sup>9</sup>Image is available at the address: [http://www.kegg.jp/kegg-bin/show\\_pathway?org\\_name=hsa&mapno=00020&mapscale=&show\\_description=show](http://www.kegg.jp/kegg-bin/show_pathway?org_name=hsa&mapno=00020&mapscale=&show_description=show)

### 2.3.3.1 KGML

The textual representation contains partial information represented in the respective map. In general, the information are written in the KGML format (KEGG Markup Language) that is based on XML language. By its markup language nature, it allows to define and control the meaning of the elements using customized tags. As we have seen before for the maps, also KGML files have a unique identifier. A specific code corresponds to the relative map. As an example, we can recognize the homo sapiens glycolysis file by the name: *hsa00010.xml*. The first three letters identify the homo sapiens organism, and the code 00010 refers to the glycolysis function.

```
1 <?xml version="1.0"?>
2 <!DOCTYPE pathway SYSTEM "http://www.kegg.jp/kegg/xml/KGML_v0.7.1_.dtd">
3 <pathway name="path:hsa00010" org="hsa" number="00010"
4     title="Glycolysis / Gluconeogenesis"
5     image="http://www.kegg.jp/kegg/pathway/hsa/hsa00010.png"
6     link="http://www.kegg.jp/kegg-bin/show_pathway?hsa00010">
7   <entry id="13" name="hsa:226 hsa:229 hsa:230" type="gene" reaction="rn:R01070" link="http://www.
8     kegg.jp/dbget-bin/www_bget?hsa:226+hsa:229+hsa:230">
9     <graphics name="ALDOA, ALDA, GSD12, HEL-S-87p..." fgcolor="#000000" bgcolor="#BFFFBF" type="
10       rectangle" x="483" y="407" width="46" height="17"/>
11   </entry>
12   ...
13   <entry id="40" name="cpd:C00033" type="compound" link="http://www.kegg.jp/dbget-bin/www_bget?
14     C00033">
15     <graphics name="C00033" fgcolor="#000000" bgcolor="#FFFFFF" type="circle" x="146" y="958"
16       width="8" height="8"/>
17   </entry>
18   ...
19   <entry id="41" name="path:hsa00030" type="map" link="http://www.kegg.jp/dbget-bin/www_bget?
20     hsa00030">
21     <graphics name="Pentose phosphate pathway" fgcolor="#000000" bgcolor="#FFFFFF" type="
22       roundrectangle" x="656" y="339" width="62" height="237"/>
23   </entry>
24   ...
25   <entry id="46" name="ko:K01568" type="ortholog" reaction="rn:R00014" link="http://www.kegg.jp/
26     dbget-bin/www_bget?K01568">
27     <graphics name="K01568" fgcolor="#000000" bgcolor="#FFFFFF" type="rectangle" x="431" y="879"
28       width="46" height="17"/>
29   </entry>
30   ...
```

```

23 <entry id="140" name="hsa:9562" type="gene" reaction="rn:R09532" link="http://www.kegg.jp/dbget-
    bin/www_bget?hsa:9562">
24 <graphics name="MINPP1, HIPER1, MINPP2, MIPP" fgcolor="#000000" bgcolor="#BFFFFB"
25 type="rectangle" x="571" y="630" width="46" height="17"/>
26 </entry>
27 ...
28 <relation entry1="62" entry2="42" type="maplink">
29 <subtype name="compound" value="84"/>
30 </relation>
31 <relation entry1="133" entry2="61" type="ECrel">
32 <subtype name="compound" value="90"/>
33 </relation>
34 ...
35 <reaction id="48" name="rn:R03270" type="irreversible">
36 <substrate id="99" name="cpd:C05125"/>
37 <substrate id="96" name="cpd:C15972"/>
38 <product id="102" name="cpd:C16255"/>
39 <product id="136" name="cpd:C00068"/>
40 </reaction>
41 ...
42 <reaction id="13" name="rn:R01070" type="reversible">
43 <substrate id="104" name="cpd:C05378"/>
44 <product id="130" name="cpd:C00118"/>
45 <product id="88" name="cpd:C00111"/>
46 </reaction>
47 </pathway>

```

Listing 2.1: KGML example from hsa00010.xml file

The textual representation is the most useful for applications that manage data, because allows for data extraction. On the contrary extracting data from the visual representation is infeasible even if it is the most readable and understandable representation for the final users. The KGML files are available only for organism-specific pathways and not for the reference pathways. Each file can be downloaded from the KEGG site. However, in order to select and download KGML files, there is an API suite published by the authors. This is the case of our application that makes use of the KEPP API<sup>10</sup> in order to get all the useful files through an automatic process. This aspect will be describe later. In the next section we present the structure of the KGML files.

<sup>10</sup>APIs are available at the address: <http://www.kegg.jp/kegg/rest/keggapi.html>

### 2.3.3.2 Pathway structure in the XML file

In this section we describe the KGML content and its relation with the corresponding graphic representation. The Figure 2-8 shows an overview of the file structure.

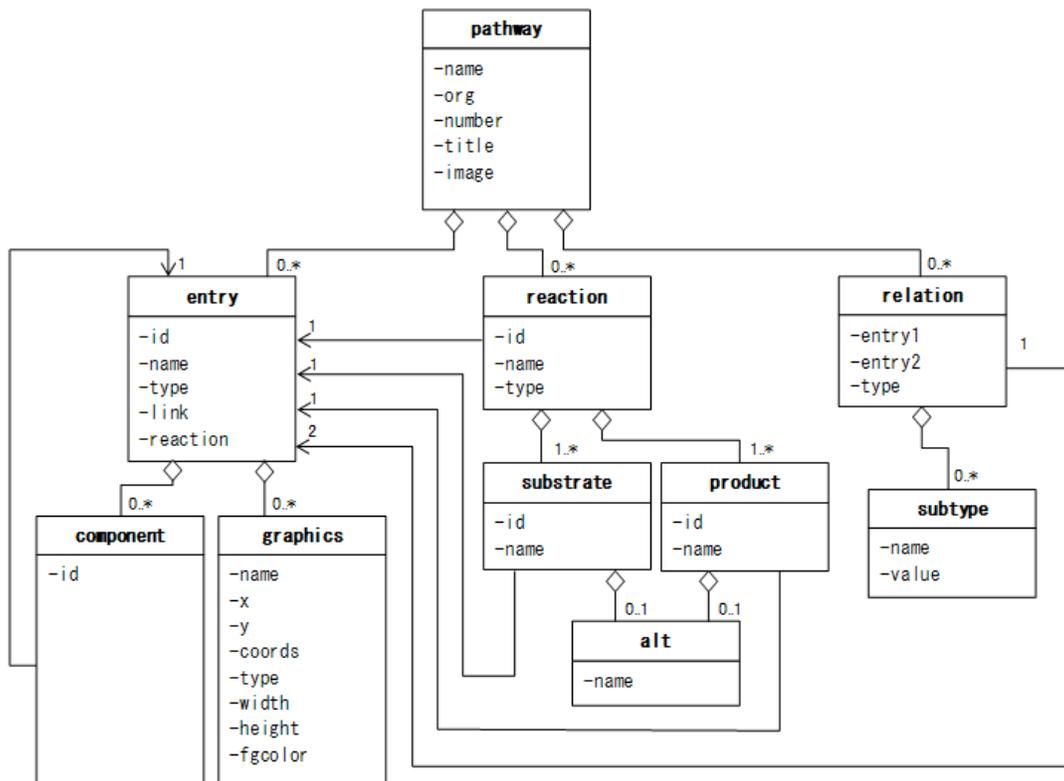


Figure 2-8: KGML Structure<sup>11</sup>

Each node corresponds to an element and the arrows indicate relations between nodes. The numbers associated to the arrows are cardinalities: the minimum and maximum numbers of relation instances. The tags used are:

**pathway:** This is the root element and it is unique in any file. It has many attributes: the name specifies the id of the pathway, the number specifies the map number and the org represents the classification of the map. The last one can assume a value among the following ones: ko, rn, ec or a three/four letter string representing the organism. All these attributes are required.

<sup>11</sup>Image is available at the address: <http://www.kegg.jp/kegg/xml/docs/>

**entry:** This element represents a node in the pathway and it contains all the information about it. The id attribute is the identification number of the entry that is unique in the file in which it is located. The same element in different maps can assume different values. The name attribute is the KEGG identifier of the entry and it is expressed in the form *db:accession*, where *accession* is the specific number of that element in the database *db*, and it can be used to perform a request using API to obtain all the information about it. Finally, the type attribute explains the element type (enzyme, reaction, gene, compound, map, group etc). The entry element has two sub-elements:

**graphics:** It contains all the information needed to draw the object. The name attribute is the label associated to the object, x and y attributes explain the position and the type specifies the object shape. In particular a rectangle represents a gene product, a circle represents others molecules such as compounds, roundrectangle represents linked pathway and lines represent reactions or relations.

**component:** It is used when the entry element is a complex node. For each component that constitutes the complex node a component element is specified with its own ID.

**relation:** It defines a relationship between two proteins or between a protein and a compound. Graphically it is represented as a line that connects two nodes. The direction is specified by entry1 (from) and entry2 (to). It contains also a type attribute that specifies the nature of the relation. In particular the type ECrel specifies an enzyme-enzyme relation and the type maplink specifies a relation between a protein and another one belonging to a different map.

**subtype:** it provides additional information about the nature of the relation such as state transition or molecular events.

**reaction:** A reaction substrate-product is described by this element and it is represented as an arrow between two circles.

**substrate:** The substrate of the reaction.

**product:** The product of the reaction.

**alt:** The alternative name of the parent.

### 2.3.3.3 KEGG API

The KEGG API (Application Programming Interface) allows users and applications to perform operation like searching, analyzing and retrieving biological information from KEGG database. The general structure of the request is the following one:

`http://rest.kegg.jp/< operation >/< argument >[/< argument >][/< option >]`

where:

`< operation > = info | list | find | get | conv | link`

`< argument > = < database > | < dbentries >`

and the *option* parameter can assume different values wrt the operation:

Data search: `< option > = formula | exact_mass | mol_weight`

Data retrieval: `< option > = aaseq | ntseq | mol | kcf | image | kgml`

In particular the `< database >` refers to the specific database name we want to use (e.g. KEGG PATHWAY, KEGG GENES ecc). The *dbentries* are in turn defined as:

`< dbentries > = < dbentry > [+ < dbentry > ...]`

where:

`< dbentry > = < db : entry > | < kid > | < org : gene >`

For our application the KEGG APIs were used to retrieve KGML files for specific organisms. In particular we have used the following type of requests:

`http://rest.kegg.jp/get/org: pathway/kgml`

The *get* operation is used to retrieve data from database, where the *org:pathway* specifies the pathway of the organism and finally *kgml* specifies the format.

### 2.3.3.4 Problems with data

During the preliminary phases of the projet we bump into some data inconsistency problems. These concern, in particular, the graphical representations and the corresponding KGML files.

We can summarize the problem in the following ways:

- Graphical representations are not always complementary: if in a specific pathway there is a graphic connection with another map, in the second map the corresponding link (the complementary information) can be missing;
- Orientation: graphical representations do not always match the orientation described by the corresponding maplink relation entries in the KGML file;
- *Maplink* relations between different pathway maps: their use sometimes is not so clear and it seems to be incomplete or mismatching.

More in general the graphical representation and the corresponding KGML files are not always equivalent.

In order to clarify the previous issues, we signal the problems to KEGG authors. The reply we got shows that there are reasonable motivations to justify the inconsistencies found. We summarize these below: the relations between different pathway maps may not be clear since there are maps that include other maps or they overlap. Therefore, relations between pathway maps are represented so that the users can refer to other pathway maps in order to get detailed information about the connection with the analyzed pathway. For these reasons connections between different pathway maps represented by dashed lines are not always represented. Furthermore a metabolic pathway can be connected to another one via compounds or through reactions.

Hence the relations of type maplink are added to each metabolic pathway map only for visual comprehensibility and they are not complete. Moreover, we cannot expect to find out all the connections and so the complementary between pathway maps both from a visual depiction as well as in the KGML files.

Everything is due to the fact that data are intrinsically incomplete and that the dig-

italization of the hand drawn schemas and the translation of biological data requires huge efforts. Accordingly with these reasons, the development of our project is based on the knowledge currently available in KEGG, taking into account such limitations.



# Chapter 3

## Comparison between metabolism in different organisms

In this chapter, we briefly review the existing methods for the reconstruction and the comparison of metabolic networks and metabolic pathways, which have been developed in the last decade for biochemical applications. These techniques are inter-related each other since metabolic networks comparison implies the reconstruction of the nets themselves from existing data repositories.

### 3.1 Metabolism representation: State of the Art

In recent years the interest of the scientific community in the development of new methods for metabolic network analysis has grown considerably. This is due to the technological evolution that allows us to deal with complex representations of big data. The methods developed for the metabolic network representation use mainly two mathematical structures:

- **Sets:** with this technique a metabolic network (or a metabolic pathway) can be represented as a set of components that can be enzymes, reactions or compounds. The comparison of such structures is generally based on set operations and it is the simplest one. A variant of this approach may be based on multi-set structures.

- **Graphs:** this representation considers both chemical compounds and their relations. There exist different methods based on simple graphs [13][14], hypergraphs [15][16] and bipartite graphs, such as Petri-nets. Graphs are a more representative structure than the sets but they have a drawback related to the complexity of the comparison between such representations. The computation of graphs or subgraphs isomorphism are known as NP-complete problems.

In this brief review of the literature we take into consideration only proposals based on KEGG information since they are freely available and well structured. We are interested in these particular proposals because also our software is based uniquely on the KEGG database. Moreover we consider only proposal based on graph since our method represents metabolic network as graph.

In [13] the authors propose a method based on graphs in order to provide a rational representation of the metabolic network structure. In order to give the metabolic network representation, they use a directed graph where nodes correspond to compounds, the oriented connections, called *arcs*, represent the irreversible reactions and the non oriented connections, called *edges*, represent reversible reactions. The method uses information taken from KEGG and the network reconstruction is performed using enzymes and reactions information.

Another approach to the representation of the entire metabolic network, has been developed by Markus Rohrschneider [15]. His work intends to provide a visualization of the metabolic network and it is based on the use of hierarchical directed hypergraph with two levels. It makes use of KEGG information. At the first level each node of the hypergraph represents a metabolic pathway of an organism and the hyperedges represent the relations between the pathways themselves (maplink information). Each hypernode is then linked to other nodes that constitutes the second level of the data structure. These ones represent the chemical compounds that specifies enzymes. The compound nodes are in turn connected to each other exploiting the enzymes relations. The structure contains also virtual edges that connect identical compounds in different

pathways to allow the user to do interactive operations, like collapse and expansion over the hypergraph.

[16] the authors propose another representation based on direct hypergraphs too, where hypernodes are metabolites and hyperedges are the enzyme-catalized reactions. The aim of this proposal is to demonstrate that metabolic networks contain phylogenetic information analysing the phylogenies obtained from network comparison. An equivalent representation of the hypergraphs is given by bipartite graph where metabolites and reactions represent two different type of vertices.

Zevedei and Schuster [17] propose a solution based on Petri Net where two kinds of nodes are considered: places and transition. Places represent metabolites and transitions represent reactions and enzymes. The static structure of the net is represented by weighted arcs that connect places to transitions and transitions to places. Weights represent the stoichiometric coefficients<sup>1</sup>. From a structural point of view the Petri Net is equivalent to an hypergraph. Beside the structure, PN allows to describe also the dynamic behaviour of a metabolic pathway. Each place (metabolite) can be equipped with tokens, describing the number of molecules of that metabolite, that is, the state of that metabolite. State changing is achieved by the firing of transitions. A transition fires if the number of tokens in the input places is greater or equal to the weights of the corresponding edge. The firing transition produces a new state in the system.

---

<sup>1</sup>The stoichiometric coefficient indicates how many molecules are needed for the reactions to happen.

## 3.2 Comparison between metabolic Networks

In the latest years some approaches have been proposed and implemented to compare the metabolic network of different organisms. These analyses are very useful to understand organisms evolution and the identification of the connections between the topology of the net and the metabolic functions is important in the biological field. Each proposal uses a different database for data retrieving and a different approach to perform the comparison, depending on the information chosen for the analysis.

When metabolic networks are compared, the computational complexity must be taken into account. In fact, some methods are infeasible because of the size of the net and alternative solutions or simplifications must be used.

In [18], Toshato proposes a method to compare organisms belonging to different species for studying interspecies phylogenies. The metabolic networks are characterized by set of enzymatic reactions and the nets are represented by a string of binary values, where the value 1 indicates the presence of the specific reaction inside the network and 0 indicates the absence. All the information are taken from MetaCyc database.

To compare two organisms,  $O$  and  $O'$ , a set  $R$  is created and it contains all the reactions contained in the networks of the two organisms. From the set  $R$ , an organism's reactions profile,  $P_x = b_{x1}b_{x2}...b_{xn}$ , is built, where  $b_{xi} = 1$ , if the reaction  $r_i$  is contained in the metabolism of organism  $X$  and  $b_{xi} = 0$ , if it is not present.

A similarity measure between each pair of organisms  $X$  and  $Y$  is performed by using Tanimoto coefficient  $T(X, Y) = \frac{N_z}{N_z + N_y - N_z}$ , where  $N_x$  and  $N_y$  are the number of bits equal to 1 in the reactions profiles  $P_x$  and  $P_y$  and  $N_z$  is the number of common reactions in the two profiles. The closer to 1 this coefficient, the higher the similarity between the organisms. From this similarity index, a dissimilarity measure,  $D(X, Y) = 1 - T(X, Y)$ , is derived and used for creating a distance matrix useful to run a Clustering algorithm based on Complete-linkage clustering method<sup>2</sup>.

---

<sup>2</sup>The *complete-linkage cluster technique* belongs to the family of hierarchical clustering procedures. This differs from the others in the computation of the distance measure (similarities). The distance between two clusters is equal to the longest distance between all elements of the first cluster and all elements of the second one.

Unfortunately no tool is available for this method.

In [19] the authors proposed a comparison method based on the construction of an evolutionary tree using pathway reaction contents and the relationship between the organisms in the cluster of orthologous group database (COG<sup>3</sup>). In particular the metabolic network of each considered organism, is built by using information retrieved from KEGG and MetaCyc database.

The metabolic net of each organisms is divided into 64 subpathways (or sub-networks), following COG division, which are compared each others to construct the metabolic content matrix. The comparison take into account the subpathways content that is based on the number of reactions included in a pathway and it is performed using the formula  $p_{ij} = 100 * r_{ij} / R_j$ , where  $r_{ij}$  is the number of reactions in the  $j$ -th sub-network of organism  $i$ , while  $R_j$  is the cardinality of the set of non-duplicated reactions found in the same pathway of all organisms.

The next step of the procedure is the Pearson correlation evaluation<sup>4</sup>. This index is performed using the formula  $PC = \frac{1}{N} \sum_{i=1, N} (\frac{X_i - \bar{X}}{\sigma_X})(\frac{Y_i - \bar{Y}}{\sigma_Y})$  considering  $\bar{X}, \bar{Y}$  the average values in  $X = \{X_1, \dots, X_N\}$  and in  $Y = \{Y_1, \dots, Y_N\}$  and  $\sigma_X, \sigma_Y$  their standard deviations. It is considered as a measure of similarity. Finally using these measures, a hierarchical clustering algorithm with the complete-linkage method is run.

Unfortunately no tool is available.

Another approach was developed in [16] in order to find if metabolic networks contain phylogenetic information. The authors propose to compare the compounds and the enzymes of the metabolic reactions in organisms. They use hypergraph as a structure representation and the algebra set operations for defining the comparison methods. This technique is based on the information taken from KEGG databases.

---

<sup>3</sup>The Cluster of Orthologous groups is generated starting from complete genomes. In particular, the protein sequences of these genomes were compared to construct clusters. Each of them contains proteins or orthologous sets of paralogs from at least three lineages [20].

<sup>4</sup>The *Pearson correlation* is an index that indicates the linearity relation between two variables. It is defined as the ratio between the covariance of the two variables and the product of their standard deviations.

The metabolic network is described by an hypergraph  $M(X, \epsilon)$ , where  $X$  is the set of metabolites and  $\epsilon$  is the set of reactions. A reaction  $E$  is expressed by a pair  $(E-, E+)$  where  $E-, E+ \subseteq X$ . For the comparison of two metabolic nets,  $M$  and  $M'$ , the set operations are used:

- **Union:** it is defined as  $M = M' \cup M''$ , where  $M$  is a new network  $(X' \cup X'', \epsilon' \cup \epsilon'')$ ;
- **Intersection:** the intersection of two networks  $M' \cap M''$  is defined as  $[(X' \cap X''), \epsilon' \cap \epsilon'']$ ;
- **Difference:** it is defined as  $M = M' \setminus M''$ , where  $M = [(supp(\epsilon' \setminus \epsilon''), \epsilon' \setminus \epsilon'')]$ . Practically the difference constitutes a new network that contains the reactions that belong to  $M'$  but that don't belong to  $M''$  and all the metabolites involved in the other reactions.
- **Symmetric difference:** it is defined as  $M = M' \Delta M''$ , where  $M = [(M' \cup M'') \setminus (M' \cap M'')]$

The phylogenetic relation among the organisms is computed by using the distance measure  $d(M, M') = \frac{\|M \Delta M'\|_5}{\|M \cup M'\|}$ .

Then a phylogenetic tree is built by using Fitch algorithm from the distance matrix of the organisms.

At present there is no tool available related to this work but the *Vienna Reaction Network Library* is reachable on the following link <http://www.biomedcentral.com/content/supplementary/1471-2105-7-67-S1.gz>. It is an open source library that implements the set operations on reactions networks.

In the paper [21], the authors developed a method for the alignment of molecular networks considering node similarity and the architecture similarity using integer quadratic programming<sup>6</sup>.

<sup>5</sup> $\|M \Delta M'\|$  represents the symmetric difference and  $\|M \cup M'\|$  represents the union

<sup>6</sup>The *integer quadratic programming* is an optimization problem with the aim of minimizing or maximizing an objective function subject respect to some constraint. In particular the purpose is find a vector  $X$  such that minimizes  $\frac{1}{2}X'QX + C'X$  subject to  $AX \leq B$  where  $Q$  is a symmetric matrix,  $A$  is a real matrix,  $B$  and  $C$  are two vectors and  $'$  is the transposition operator.

The information are taken from Metacyc database and SGD database. A molecular network is represented by a graph  $G(V, E)$ , where  $V$  is the set of nodes (proteins or genes) and  $E$  is the set of edges (interactions between nodes).

This graph structure is expressed by an adjacency matrix  $A = (a_{ij})_{m \times n}$ , where  $a_{ij} = 1$  if there is an interaction between two proteins and  $a_{ij} = 0$  otherwise. This value can also be modified to be a real number in the range  $[0,1]$  to express not only the presence or absence of an interaction but to give a confidence ratio for that interaction. The similarity,  $s_{ij} = S(v_i^1, v_j^2)$ , between two nodes in two networks  $G_1(V_1, E_1)$  and  $G_2(V_2, E_2)$ , where  $v_i^1 \in V_1$  and  $v_j^2 \in V_2$  is based on the alignment of their sequences and it is expressed by a function:  $S : V_1 \times V_2 \rightarrow [0, 1]$ .

The matching between nodes belonging to two different nets is represented as:

$$x_{ij} = \begin{cases} 1 & \text{if } v_i^1 \text{ matches } v_j^2 \\ 0 & \text{otherwise} \end{cases}$$

The optimal alignment  $X = \{x_{ij}\}$  is obtained by the integer quadratic programming model, whose aim is to maximize the similarity score between  $G_1$  and  $G_2$  wrt. all  $X$  combinations.

The model is the following:

$$\max_X f(G_1, G_2) = \lambda \sum_{i=1}^m \sum_{j=1}^n s_{ij} x_{ij} + (1 - \lambda) \sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^m \sum_{l=1}^n a_{ik} b_{jl} x_{ij} x_{kl}$$

$$\text{such that } \begin{cases} \sum_{j=1}^n x_{ij} \leq 1 & i = 1, 2, \dots, m \\ \sum_{i=1}^m x_{ij} \leq 1 & j = 1, 2, \dots, n \\ x_{ij} = 0, 1 & i = 1, 2, \dots, m, j = 1, 2, \dots, n \end{cases} .$$

$\lambda$  is simply a coefficient that allow to set a weight to nodes or edges.

The authors have built a command line tool, written in Matlab language, that allows one to align networks formulating the problem as a quadratic programming problem and solving it. The program implements the interior points method to solve the quadratic problem model through the *spphase1* subroutine and then it finds a

local optimal solution through *spphase2* subroutine.

Ay, Dang and Kahveci [22] proposed an approach for aligning metabolic networks solving the complexity due to the large size of the considered nets. Their main aim is to reduce the scale of the networks in order to use existing alignment algorithms in an efficient way. The procedure is divided into three phases:

1. *compression phase*: in this step the size of the network is reduced, moving from the *original domain* to a *compressed domain*. During this phase a set of nodes and the corresponding relations in the first domain are substituted by a *supernode* in the second domain. Thus, a supernode represents a summary of the previous nodes.
2. *alignment phase*: in this step subMAP alignment algorithm is applied to compress domain data. As a result, a mapping between supernodes in the compressed domain is obtained.
3. *refinement phase*: in this final step each mapping of supernodes found in the previous steps, is considered, one by one. In fact, they correspond to an instance of the alignment problem. The subMAP algorithm is applied again at each of these mappings to obtain the final result.

In [23] many metabolic networks were compared in order to show that the *scope* is robust with respect to the reactions network. The authors define the scope as the synthesizing network capacity with respect to specific substrates. It allows to characterize a biological function of the network. The final aim of the work is to evaluate similarity or differences of functional and structural proprieties in different organisms. The compared metabolic networks are constructed starting from an initial set of compounds, *seed*, and a set of reactions, *base*, following these steps:

- Identify reactions from base set that use as substrates compounds included in the seed set;
- Add these reactions and their products to the net;

- Repeat these steps until no new reaction to add can be founded.

At the end all the compounds constitute the scope of the seed. All the information are taken from KEGG databases.

The two networks  $M$  and  $M'$  of two organisms  $O$  and  $O'$  are compared using a distance measure that measures the dissimilarity of their reactions:  $d(O, O') = |M \cup M'| - |M \cap M'|$ .  $M \cup M'$  indicates all the reactions present in either of the two networks and  $M \cap M'$  indicates reactions included in both organisms. This formula can be considered as an evolutionary distance measure for two organisms, where reaction acquisition and reaction loss are the events.

Furthermore another measure to compare biological functions have been defined:  $d(O, O') = |S \cup S'| - |S \cap S'|$ , where  $S$  and  $S'$  are the scopes of the seed.

In the Table 3.1 we summarize the main characteristic of the works just described.

Reference	Representation	Database	Tool
[18]	Reaction profile	MetaCyc	/
[19]	Graph	KEGG, Metacyc	/
[16]	Hypergraph	KEGG	/
[21]	Graph	Metacyc, SGD	MNAligner
[22]	Graph	KEGG	/
[23]	Set	KEGG	/

Table 3.1: Summary of proposed method

All the comparison methods found in the literature in order to deal with the huge size of metabolic networks use techniques that can be applied to large data sets, like set operations, or use simplifications to the detriment of some accuracy.



# Chapter 4

## Metabolic Networks comparison

### 4.1 Metabolic network construction

In this chapter we describe how we perform the metabolic network reconstruction starting from the KEGG database information. In particular we explain our algorithm and the data structures used. We propose a comparison method for such metabolic networks which consider both the structure of the network and the similarity between corresponding pathways. The similarity indexes which are computed are also described. Moreover, we discuss some troubles found during the development of the method and the solutions adopted.

#### 4.1.1 Network construction

In the following section we introduce the methodology that allows for reconstructing metabolic networks. Our work relies uniquely on KEGG database information for many reasons. One of these is that the KEGG project has proved to be a reliable knowledge base during the time and it is growing steadily. Another reason is that KEGG provides a digitization of information that are particularly complex. KEGG gives a global metabolic network representation which resumes the metabolisms of all the catalogued organisms. We refer to this data structure as the *reference metabolism*. The net is composed by the union of all the reference pathways thus giving an implicit

partition of the the whole metabolism into metabolic pathways which is standardized wrt. all organisms. This choice allows us to analyse and compare specific metabolic functions or entire organisms. However, KEGG data can be incomplete and this leads to issues related to data completeness.

In the literature the majority of the approaches represent the metabolic pathways as graphs of reactions, in order to keep a good level of details, then the metabolic network is obtained by the union of the involved pathway's graphs. The consequence of this approach is that the resulting graph of the metabolism is very huge. The comparison of such large graphs requires to compute some kind of graph isomorphism and it becomes infeasible.

The aim of our thesis is to propose a new method to compare the entire metabolism of different organisms by considering both topology and functionality. We model the metabolic networks using graphs with a certain level of abstraction. This choice simplifies the problem of comparing graphs of big dimensions. We propose the following graph representation.

*Let  $O$  be a specific organism, then  $G_O = (V_O, E_O)$  is the **metabolic graph** of  $O$ , where  $V_O = \{P_1, \dots, P_n\}$  is the set of nodes which represent the metabolic pathways of  $O$ , namely each  $P_i$ , with  $i \in [1, n]$ , is the  $i$ -th pathway represented as a set (or multiset) of reactions,  $P_i = \{r_1, \dots, r_m\}$ , and  $E_O$  is the set of edges that represents the relations between the pathways of  $O$ .*

Our representation of metabolism is organized into two levels:

- **Lower level:** it represents a metabolic pathway  $P_i$  in terms of set/multiset of chemical reactions;
- **Higher level:** it represents the entire metabolism by a graph  $G_O$ , considering the pathways and the relations among them.

This representation fits perfectly the KEGG database organization since each specific pathway  $P_i$ , is represented in all the organisms in a standardized way (reference pathway) and the metabolism considers each metabolic function and the interactions

defined between them. Our representation guarantees the independence between the two levels, global network level and pathway level. This gives the possibility to use more detailed representations both at pathway level and at network level. At present, we decided to start from the simplest representations leaving more complex representations for future developments.

By adopting a two levels representation, we are able to reduce the size of the graph representing the metabolic network, since nodes represent the pathways rather than the reactions. Hence, comparison between graphs becomes feasible.

### 4.1.2 Implementation

In order to build the metabolic network of a specific organism, the first step is the data retrieval. We have considered 159 pathways belonging to the following categories in KEGG:

- Carbohydrate metabolism;
- Energy metabolism;
- Lipid metabolism;
- Nucleotide metabolism;
- Amino-acid metabolism;
- Metabolism of other amino-acids;
- Glycan biosynthesis and metabolism;
- Metabolism of cofactors and vitamins;
- Metabolism of Terpenoids and polyketides;
- Biosynthesis of other secondary metabolites;
- Xenobiotics biodegradation and metabolism.

The pathways that are included in the listed categories and that belong to the considered organism  $O$  constitute the set of nodes,  $V_O$ . KEGG provides KGML files for pathways that include gene/protein network or chemical network, the other pathways

are depicted using images. Pathways that don't have a KGML file and use images are not considered in our comparison method. Moreover, we check specific cases in which the KGML files don't contain declaration of chemical reactions. Finally, we consider the set of connections between the pathways themselves and collect such relationship into  $E_O$ .

We download through the public KEGG's API, the kgml files for each pathway. The requests are done by using the following URL: <http://rest.kegg.jp/get/org:pathway/kgml> as described in section 2.3.3.3. Then we perform a sequential parsing procedure, reading each KGML file iteratively. In this step the essential information (reactions and maplinks) are extracted for the construction of  $G_O$ . In listing 4.1 we can see a fragment of KGML file with some of such data. For convenience we shown only a snippet of code containing the necessary information treated during the parsing phase.

```

1 <?xml version="1.0"?>
2 <!DOCTYPE pathway SYSTEM "http://www.kegg.jp/kegg/xml/KGML_v0.7.1_.dtd">
3 <!-- Creation date: Apr 22, 2016 16:49:05 +0900 (GMT+9) -->
4 <pathway name="path:hsa00010" org="hsa" number="00010"
5     title="Glycolysis / Gluconeogenesis"
6     image="http://www.kegg.jp/kegg/pathway/hsa/hsa00010.png"
7     link="http://www.kegg.jp/kegg-bin/show_pathway?hsa00010">
8     ...
9     <entry id="41" name="path:hsa00030" type="map"
10        link="http://www.kegg.jp/dbget-bin/www_bget?hsa00030">
11        <graphics name="Pentose phosphate pathway" fgcolor="#000000" bgcolor="#FFFFFF"
12            type="roundrectangle" x="656" y="339" width="62" height="237"/>
13    </entry>
14    <entry id="56" name="hsa:2597 hsa:26330" type="gene" reaction="rn:R01061"
15        link="http://www.kegg.jp/dbget-bin/www_bget?hsa:2597+hsa:26330">
16        <graphics name="GAPDH, G3PD, GAPD, HEL-S-162eP..." fgcolor="#000000" bgcolor="#BFFFBF"
17            type="rectangle" x="458" y="484" width="46" height="17"/>
18    </entry>
19    <entry id="61" name="hsa:2821" type="gene" reaction="rn:R02740"
20        link="http://www.kegg.jp/dbget-bin/www_bget?hsa:2821">
21        <graphics name="GPI, AMF, GNPI, NLK, PGI, PHI, SA-36, SA36" fgcolor="#000000" bgcolor="#BFFFBF"
22            "
23            type="rectangle" x="483" y="265" width="46" height="17"/>
24    </entry>
25    ...
26    <relation entry1="61" entry2="41" type="maplink">

```

```

26     <subtype name="compound" value="90"/>
27 </relation>
28 <relation entry1="41" entry2="56" type="maplink">
29     <subtype name="compound" value="130"/>
30 </relation>
31 ...
32 </pathway>

```

Listing 4.1: Example of maplink information extracted from KGML files

To represent a pathway we extract the set of its reactions. We collect the set of reaction attributes related to the entries of type *gene* (see the lines 14 and 19 in Listing 4.1). In this way we consider only genes belonging to the specific organism without considering orthologs. For the network construction, we extract the tags relation of type *maplink* (rows 25 and 28). The maplink relations have two relevant attributes, *entry1* and *entry2*, that contain the IDs of specific entries (see the lines 25 and 28 in Listing 4.1). Such attributes specify also the orientation of the connection:

- **entry1**: it is the start element of the relation;
- **entry2**: it is the end element of the relation.

By analysing the type attribute of the two entries, we can understand which are the pathways involved. If the entry is of type *gene*, the corresponding pathway is the pathway tag of the KGML file in analysis, else if the entry is of type *map*, the corresponding pathway is specified by the name attribute of the entry itself.

The information listed in 4.1 is visualized (see the part highlighted in red) in Fig. 4-1.

As an example let us consider the entries with ID 56 and 61 that represent two distinct enzymes which correspond to EC numbers 1.2.1.12 and 5.3.1.9 respectively. The third entry with ID 41, instead, represents the Pentose phosphate pathway. The maplink relation at row 25 in Listing 4.1, connects the enzyme 5.3.1.9 with the Pentose map. The orientation of the connection is given following the order of the entries, as defined before. Since the enzyme 5.3.1.9 constitutes an element of the Glycolysis pathway, an edge from the Glycolysis node to the Pentose phosphate node is created.

Our tool allows for representing metabolic network either as directed graphs (i.e. maplink are translated into oriented edges) or as undirected graphs (i.e. maplink are represented as undirected arcs). Concerning pathways, the tool offers the possibility to represent the either as set of reactions, or as multiset of reactions (i.e. multiple occurrences of the same reaction are considered).

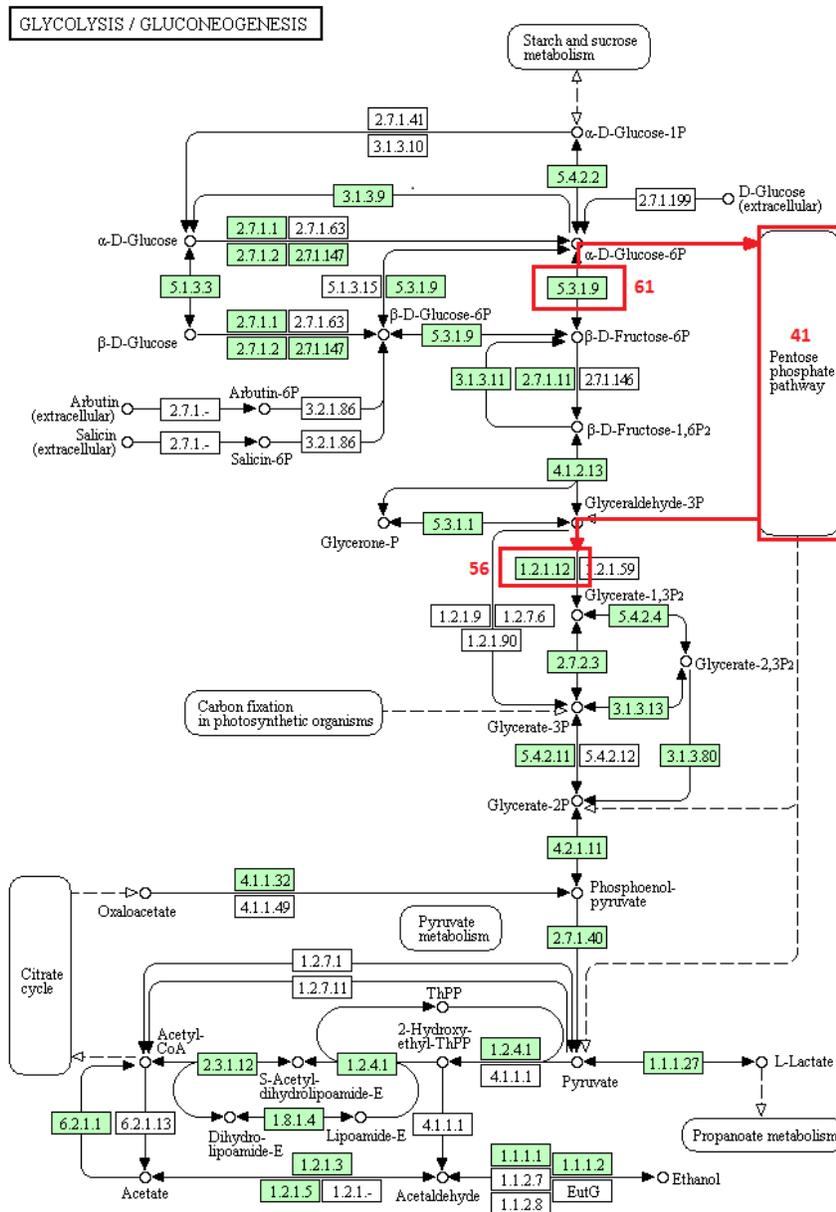


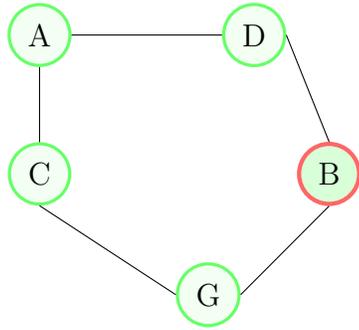
Figure 4-1: Example of maplink detection on Homo Sapiens Glycolysis

### 4.1.3 Data structures

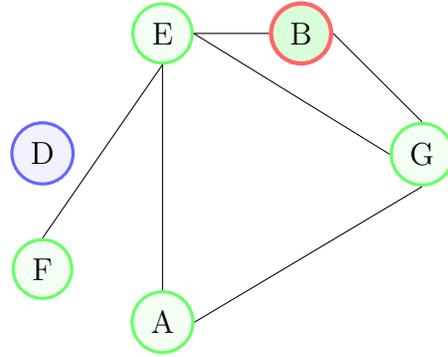
The implementation of the graph of a metabolic network of an organism is done by using a modified adjacency matrix. We use a square matrix of size  $n$ , representing all the pathways listed in Section 4.1.2. This standardized data structure correspond to a mapping between the same metabolic functions (nodes) in different organisms (matrices). A value 1 on the diagonal of an adjacency matrix indicates the presence of a loop in the corresponding graph, but loops are actually never present in our graphs of metabolic networks because edges correspond to KEGG maplinks. Hence a diagonal would be composed only by 0 values. For that reason we exploit the diagonals of the adjacency matrices to represent further information on the nodes of the graphs with the following conventions:

- 1 represents a connected node (pathway);
- 0 represents an isolated node (pathway);
- -1 represents a pathway which is not present in the metabolism of the organism.

This choice allows us to check quickly whether the nodes are connected. 1 values, represent metabolic functions that are connected with at least one other function. 0 values represent metabolic pathways with no connections. The -1 values, indicate that a specific pathway is not present in the metabolism of the organism. The values outside the diagonal represent the edges of the graphs. 0 values represent missing edges and 1 values represent the existing ones. Such matrices represent an abstraction of the reference metabolism given in KEGG. Let us consider a simplified example of metabolic networks of two organisms,  $O$  and  $O'$ . The set of metabolic pathways, in this artificial example is represented by  $\{A, B, C, D, E, F, G, H\}$ . We show the two graphs and the matrices used in our approach.



(a) Graph of the metabolic network of  $O$



(b) Graph of the metabolic network of  $O'$

	A	B	C	D	E	F	G	H
A	1	0	1	1	0	0	0	0
B	0	1	0	1	0	0	1	0
C	1	0	1	0	0	0	1	0
D	1	1	0	1	0	0	0	0
E	0	0	0	0	-1	0	0	0
F	0	0	0	0	0	-1	0	0
G	0	1	1	0	0	0	1	0
H	0	0	0	0	0	0	0	-1

(a) Matrix of  $G_O$

	A	B	C	D	E	F	G	H
A	1	0	0	0	1	0	1	0
B	0	1	0	0	1	0	1	0
C	0	0	-1	0	0	0	0	0
D	0	0	0	0	0	0	0	0
E	1	1	0	0	1	1	1	0
F	0	0	0	0	1	1	0	0
G	1	1	0	0	1	0	1	0
H	0	0	0	0	0	0	0	-1

(b) Matrix of  $G_{O'}$

We consider the sets of metabolic pathways present in the metabolism of each organism:

$$M(O) = \{A, B, C, D, G\} \quad \text{and} \quad M(O') = \{A, B, D, E, F, G\}.$$

We use three different colours in the picture: green represents the connected pathways, red represents the metabolic pathways not present in the organism and blue represents the isolated pathways. The correspondence between nodes that represent the same metabolic pathway is given for free since these nodes have the same indexes in the two matrices. This implicit matching is based on KEGG's reference pathways and it allows us to simplify the matrices comparison.

## 4.2 Comparison of metabolic networks

We present now the technique we propose for comparing metabolic networks of two different organisms. Our comparison method follows our metabolic network representation: it uses a bottom-up approach and it is developed in two distinct levels.

- **Low level:** we perform a comparison between pairs of corresponding pathways in the two organisms. Each metabolic function can be represented either as a set or as a multiset of reactions, depending on the user's choice. We execute a comparison on these sets (multisets) providing a similarity value.
- **High level:** we compare the topologies of the metabolic networks (i.e. their modified adjacency matrices) taking into account also the similarity values computed at the first level. Networks can be modeled as directed or undirected graphs and their comparison produces a similarity value for the entire metabolic networks.

We discuss now the similarity indexes computation that allow us to compare two metabolic networks. In particular we describe in detail the similarity computation implemented in order to compare the topologies of nets. After that we define two different similarities about the overall metabolic networks. As explained before in order to compute such indexes we evaluate both the topology and the similarities between metabolic pathways. Concerning the evaluation between each pair of pathways we refer to the definition of the similarity index described in [1].

### 4.2.1 Similarity between metabolic networks

We propose an index to represent the topological similarity between two metabolic networks. Before defining such index, we need some further concepts. Given two graphs, representing metabolic networks,  $G = (V, E)$  and  $G' = (V', E')$  we have that:

- Two nodes  $v \in V$  and  $u \in V'$  are *matching nodes* if they represent the same metabolic function. In particular if there exists a matching function that maps

$v$  into  $u$ . In our case, we want to compare matching nodes, where the matching function is defined implicitly by KEGG through the reference pathways.

- The *degree* of a vertex  $v \in V$  ( $v' \in V'$ ) is defined as the number of edges incident to it and it is represented by  $\deg(v)$  ( $\deg(v')$ ). If the graph is oriented, then it represents the sum of the ingoing and outgoing edges. This information is used in the similarity definition to compare a connected node to an isolated node.
- $G$  and  $G'$  are *isomorphic*, if there exists a bijective function  $f$  such that for each  $v, w \in V$  then  $(v, w) \in E \Leftrightarrow (f(v), f(w)) \in E'$ . In other words, matching nodes in the two graphs must be connected to corresponding nodes. Isomorphism would be the maximal similarity between  $G$  and  $G'$  and this is the idea underlying our similarity measure which considers how many edges are in commons between  $G$  and  $G'$ .

Let us consider two organisms  $O$  and  $O'$  to be compared and their corresponding graphs of metabolic network,  $G = (V, E)$  and  $G' = (V', E')$ , where  $V$  is the set of metabolic pathways of the first organism and  $V'$  is the set of pathways of the second one, while  $E$  and  $E'$  are the corresponding sets of edges that represent connections between such pathways. Let us consider the  $i$ -th pathway in the two graphs where the order is determined by KEGG reference pathways,  $P_i \in V$  and  $P'_i \in V'$ . Let  $E_i, E_i \subseteq E$  be the set of edges that connect  $P_i$  with other nodes. Similarly we define  $E'_i, E'_i \subseteq E'$ , for  $P'_i$ . The **structural similarity index** wrt. the  $i$ -th pathway,  $SimS_i$ , is defined as:

$$SimS_i = \begin{cases} 0 & \text{if } P_i \text{ is not present in } O \text{ or } P'_i \text{ is not present in } O' \\ 1 & \text{if } P_i \text{ and } P'_i \text{ are both isolated nodes} \\ \frac{1}{1+\deg(P_i)} & \text{if only } P'_i \text{ is an isolated node} \\ \frac{1}{1+\deg(P'_i)} & \text{if only } P_i \text{ is an isolated node} \\ \frac{|E_i \cap E'_i|}{|E_i \cup E'_i|} & \text{if } P_i \text{ and } P'_i \text{ are both connected nodes} \end{cases}$$

This definition is based on the Jaccard index  $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$ , but it considers also the case of isolated or missing nodes in one of the two graphs.

Note that,  $SimS_i \in [0, 1]$ , if its value is close to 1, it means that the similarity is high, conversely if it is close to 0, it means that the similarity is low. Value 1 is associated to matching nodes that have all matching edges or are both isolated.

This similarity measure wrt. a pair of matching nodes, can be extended to the entire network to obtain a global value representing the structural net similarity. The **structural network similarity index** is defined as:

$$SimS = \frac{\sum_{i=1}^n SimS_i}{n}$$

where  $n = |V \cup V'|$ . This is an arithmetic mean whose value belongs to the range  $[0, 1]$ .

Computing the global similarity index  $SimS$  corresponds to a specific comparison of the modified adjacency matrices we introduced in Section 4.1.3. It corresponds to a sequence alignment with the scoring system given by the definition of  $SimS_i$ .

## 4.2.2 Global similarity indexes

We define two different indexes for comparing metabolic networks in different organisms which are based both on the similarity of their structure and on the similarity of the corresponding functions (pathways). This means to compare the metabolic networks using both  $SimS_i$  defined before and  $SimP_i$  defined in [1]. We briefly recall the definition of  $SimP_i$ .

Given two organisms  $O$  and  $O'$ , we consider  $P_i$ , the  $i$ -th pathway in the KEGG order of reference pathway. Let us assume that a pathway defined as a set of reactions, then  $SimP_i$  is defined as:

$$SimP_i = \begin{cases} 0 & \text{if } P_i \text{ is missing in } O \text{ or in } O' \\ 1 & \text{if } P_i \text{ is present in both } O, O' \text{ but there are no reactions to compare} \\ \frac{|R_i \cap R'_i|}{|R_i \cup R'_i|} & \text{otherwise, where } R_i, R'_i \text{ are the sets of reactions of } P_i \text{ in } O, \text{ and of } P'_i \text{ in } O' \text{ respec} \end{cases}$$

The first global similarity index we define, is called the *Combined Similarity Index*

since  $SimS_i$  and  $SimP_i$  are related to each other. Given two organisms  $O$  and  $O'$ , we define the **combined similarity index** as follows:

$$CI = \frac{\sum_{i=1}^n SimS_i * SimP_i}{n}$$

where  $n = |M|$  and  $M$  is the union of their metabolic pathways. In order to normalize the index we divide the summation by  $n$ , thus the value of  $CI$  is in  $[0, 1]$ .

The second global similarity index is called *Separated Similarity Index* since we introduce an  $\alpha$  parameter that allows us to weight the structural similarity index,  $SimS$ , and the weighted functional similarity index,  $SimPW$ .  $SimPW$  is defined in [1] and it is the weighted similarity index wrt. the metabolic functions of an organism:

$$SimPW = \frac{\sum_{i=1}^n SimP_i * |R_i \cup R'_i|}{\sum_{i=1}^n |R_i \cup R'_i|}$$

where  $R_i$  is the set of reactions of  $P_i$ ,  $R'_i$  is the set of reactions of  $P'_i$  and  $n$  is the union of the pathways of the organisms  $O$  and  $O'$ .

We define the **separated similarity index** as follow:

$$SI = \alpha * SimS + (1 - \alpha) * SimPW$$

where  $\alpha \in [0, 1]$ . The values assumed by the index  $SI$  are in  $[0, 1]$ . Choosing different values of  $\alpha$  allows us to give more relevance either to structural similarity or to pathway similarity. Particular cases are for  $\alpha = 0$  or  $\alpha = 1$ . When  $\alpha = 0$  we consider uniquely  $SimPW$  and exclude  $SimS$ , on the contrary, when  $\alpha = 1$ , we consider uniquely  $SimS$  and exclude  $SimPW$ .

The choice of the global index, either  $CI$  or  $SI$ , is determined by the context in which the metabolic network comparison is done. If we compare two organisms belonging to the same phylum<sup>1</sup>, the topology of their metabolic networks should be almost the same. In this case the use of  $SI$  is more suitable, since with  $\alpha < 0,5$  we

---

<sup>1</sup>In biology, the phylum is the primary subdivision of a taxonomic kingdom, grouping together all classes of organisms that have the same body plan [24]

can give more relevance to the comparison of metabolic functions. The use of *CI* could be more useful when comparing two distant organisms, by considering both the relative topologies and pathways. In Table 4.1 and in Table 4.2 we summarize the local/global similarity indices respectively.

	Index	Description
$SimP_i =$	$\begin{cases} 0 & \text{if } P_i \text{ is missing in } O \text{ or in } O' \\ 1 & \text{if } P_i \text{ is in } O, O' \text{ but there} \\ & \text{are no reactions to compare} \\ \frac{ R_i \cap R'_i }{ R_i \cup R'_i } & \text{otherwise} \end{cases}$	The <b>pathway similarity index</b> considers the union of the metabolic pathways of the organisms, the similarity value of the corresponding pathways is defined in term of reactions.
$SimS_i =$	$\begin{cases} 0 & \text{if } P_i \text{ or } P'_i \text{ is not present} \\ 1 & \text{if } P_i \text{ and } P'_i \text{ are both} \\ & \text{isolated} \\ \frac{1}{1+deg(P_i)} & \text{if only } P'_i \text{ is isolated} \\ \frac{1}{1+deg(P'_i)} & \text{if only } P_i \text{ is isolated} \\ \frac{ E_i \cap E'_i }{ E_i \cup E'_i } & \text{if } P_i \text{ and } P'_i \text{ are both} \\ & \text{connected} \end{cases}$	The <b>structural similarity index</b> defines the similarity between two matching nodes in terms of connections

Table 4.1: Summary of the local similarity indexes

Let us briefly discuss the complexity of the main functions implemented in our tool. The functions used in the comparison procedures are:

- **SetCompare**: this function allows one to compare the reactions of the same metabolic pathway in two different organisms. We store the reactions into HashMap data structures. Usually, basic operations like insertion, deletion and search in such data structure have a constant complexity  $O(1)$ . In the worst cases they have  $O(n)$  complexity. In the simplest cases, the SetCompare returns 0 value if the pathways is missing in one of the two organisms or 1 value if the pathway is present in both organisms but there are no reactions to compare. The more complex case is verified when both pathways contain reactions and thus, when the function returns the ratio between the intersection and the union of the reactions involved in the comparison. The computation of the union is performed using HashSet that correspond to set data structure

Index	Description
$SimPA = \frac{\sum_{i=1}^n SimP_i}{n}$	The <b>functional similarity index</b> is the mean similarity over the union of the pathways of the organisms $O$ and $O'$
$SimPW = \frac{\sum_{i=1}^n SimP_i *  R_i \cup R'_i }{\sum_{i=1}^n  R_i \cup R'_i }$	The <b>weighted functional similarity index</b> is the weighted mean similarity over the union of the pathways of $O$ and $O'$ wrt. the number of reactions
$SimS = \frac{\sum_{i=1}^n SimS_i}{n}$	The <b>structural network similarity index</b> represents the topological similarity of the entire metabolic networks
$CI = \frac{\sum_{i=1}^n SimS_i * SimP_i}{n}$	The <b>combined similarity index</b> provides a global measure comparing the similarities of both topology and functionalities of the metabolic networks
$SI = \alpha * SimS + (1 - \alpha) * SimPW$	The <b>separated similarity index</b> provides a global measure combining with a weight the similarities of both topology and functionalities of the metabolic networks

Table 4.2: Summary of the global similarity indexes

developed in Java. The complexity of *union* function is  $O(m + n)$  where  $m$  and  $n$  are the number of reactions in the two pathways, since a scan of both HashMaps is required. Each element is added to a set with constant time  $O(1)$ . The complexity of the *intersection* function is  $O(m \cdot n)$  where  $m$  and  $n$  are the number of reactions in the two pathways respectively.  $O(m)$  is the complexity to scan all the elements in the first HashMap and  $O(n)$  is the time, in the worst case, for searching the corresponding element in the second HashMap. Considering that the number of reactions for each pathway is on average less than one hundred, the computational complexity becomes reasonable. These operations are repeated for all the pathways of the two organisms. The same considerations can be done for the complexity in using multiset data structures.

- **NetworkCompare**: this function allows one to compare the topology of two metabolic networks. The networks are represented by using square matrices

$N * N$ , where  $N$  is 159, namely the cardinality of the pathways taken from KEGG. The complexity of this function is  $O(N^2)$  since all the elements in the matrices need to be inspected. If the metabolisms are represented as undirected graphs, the complexity is  $O(\frac{N^2}{2})$  since only half of the matrices need to be considered.

The complexity of the metabolic network comparison, is the sum of the procedures' complexity described above. Thanks to the two-level representation and comparison performed by our approach, the order of all the elements treated is reasonable. This permit us to perform the comparison in a reasonable time respect to the existing method that model the entire metabolism as graph of reactions.

In Chapter 6 we discuss some experiments with the two indexes done in order to validate their application.



# Chapter 5

## Tool

The goal of our project is to create an application that allows the user to compare whole metabolisms or specific set of metabolic functions of different species. This kind of comparisons are important to find out differences between the metabolic functions of different organisms. The analysis is useful to identify important information that can be used in some branches like drug engineering and medical science. Our application permits the choice of two organisms, performs a fast comparison for which it is possible to select the comparison method on two distinct levels (pathway level and network level) and provides as a result some similarity measures. In this chapter we describe the requirements of the project, the software architecture, the technologies and libraries used and finally we present a brief documentation.

### 5.1 Requirements analysis

The first step in the development of a software project is the requirement analysis. During this phase we consider the software system requirements as functional ones, which describe the services and the features of the application, and the non functional ones that describe the constraints on the product and the process development.

### 5.1.1 Functional requirements

Functional requirements permit us to identify functionalities of the software system in terms of services, system reactions under specific inputs and general behaviour of the system. Below we list the functional requirements of our application:

- **Download of the KEGG organisms information:** this functionality should permit the update of the local database with all the information about the organisms;
- **Selection of the comparison of a specific pathway or metabolic network:** the software gives the possibility to compare either the entire metabolism of the organisms or only a subset of the metabolic functions;
- **Selection of the two organisms:** the user should select two organisms from the list of all organisms present in the KEGG database;
- **Download of KGML files:** the application should download automatically the KGML files when they are not already present in the local folders. If the files are already locally present, the user should choose if to update them or to use the existing ones for the comparison;
- **Choice of the comparison methods:** the user must have the possibility to select different methods of comparison, either for the metabolic functions or for the metabolic network;
- **Choice of the  $\alpha$  value:** this functionality should allow the user to set a value for the alpha parameter in order to tune the separated similarity index;
- **Automatic exportation of the results as .xls file:** the application must save the computation results in a .xls file for an instant retrieval in a second moment. An .xls file should be saved for each comparison executed by the application;
- **Visualization of data results:** the tool must provide clear and readable results about pathways and networks comparison;

- **Navigation between views:** the user should be able to move back and forward between the windows of the application by using specific buttons;
- **Consecutive comparisons on the same selected organisms:** at the end of an execution, the software must give the possibility to the user to select different comparison methods and execute another run on the same selected organisms.

### 5.1.2 Non-functional requirements

Non-functional requirements are not directly concerned with the specific services and functionalities defined by functional requirements, but they define constraints on the system or on the development process of the software. They are classified in three main classes that are:

- **Product Requirements:** they allow us to define constraints on the services offered by the system specifying the usability, efficiency, reliability and portability of the software;
- **Organizational requirements:** they specify process standards, platforms, delivery requirements, etc, to be used;
- **External requirements:** they stem from factors external to the system and its development process (such as the interoperability requirements, legislative, ethical, etc.).

We define a list of non-functional requirements for our application as follows:

- **Fast comparison:** the computation of the similarity indexes must be done in a reasonable time;
- **Parallelized computation:** the software must be developed using threads in order to parallelize the computation as much as possible;
- **Portability:** the application must run on different heterogeneous environments.

## 5.2 Project architecture

Our project is developed using the MVC (Model-View-Controller) pattern [25]. It is the most used programming pattern to manage software that makes use of GUI (Graphical User Interface). The three main components of the MVC are:

- **Model:** the code collected under this module of the pattern handles data and business logic of the application. In particular, here we find the set of classes that define the context of the application and all the methods that allow the interactions with the databases;
- **View:** the view module collect the set of the GUIs and it is the main responsible for the logic of data presentation. Each view represents the way through which the users interact with the system;
- **Controller:** it reacts to the interactions of the users on the views and it executes the corresponding actions in the model that allow the update of the views.

The three modules interact with each other starting from the main controller that represents the entry point of the application. Then, it initializes the view and it interacts with the model in order to update the view with the data. Every time that a user executes an action through the view, the controller checks the correctness of the inputs blocking bad requests or calling the related methods in the model. In this last case, the execution of the procedure updates the view. The architecture of our application is given in Fig.5-1.

The advantages in using this programming pattern is to achieve a good modularization of the code which gives ease of maintenance, clear separation of tasks during the development process and possibility to work with a certain level of independence on the components to develop. Moreover, the development of additional features and functionalities are possible and made easier thanks to this kind of software architecture.

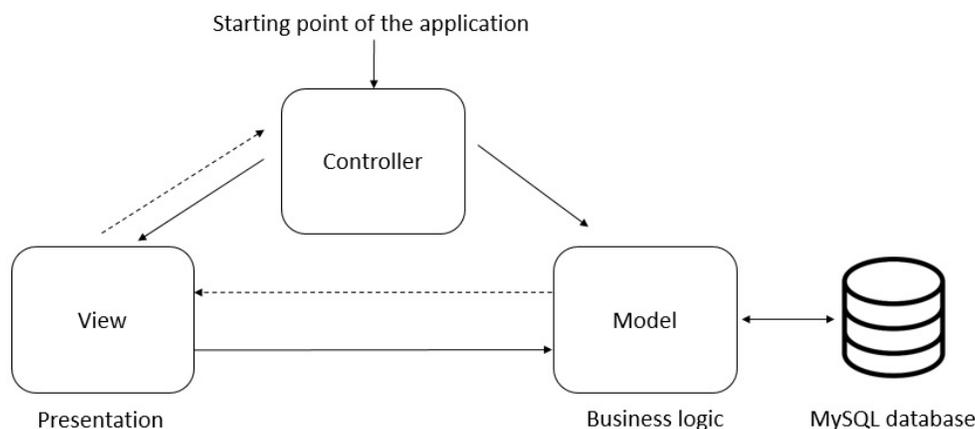


Figure 5-1: Application architecture

The implementation of our tool makes also use of multithreading, a programming technique that allows the execution of distinct threads in a concurrent way in the context of a process<sup>1</sup>. A thread constitutes a part of a process that executes a small part of the code and makes use of shared resources. Thus, processes can be divided in more threads whose execution can be parallelized using the resources of the process itself. Such approach is relevant in our application, since we implement some procedures by dividing the workload in different tasks that can be executed in parallel. Processes like the download of the KGML files and the parsing procedures are implemented as threads. In this way, we increase the performance in the execution of such tasks. Technically, starting from a request made by the user through the GUI of our tool, we create an instance of a thread for each selected organisms whose task is to retrieve all KGML files related to the organism and we store them in an organized structure of folders. In the next step, when the user requires to start the comparison, we create other two threads in order to parallelize the sequential parsing operations for each single file in the corresponding folders.

The tool relies on a MySQL database [26] in which we store all the KEGG organisms information. This choice gives us some advantages. The retrieval of such

---

<sup>1</sup>A process is an instance of a program that is executed by the CPU. It consists of resources like an image of the code that should be executed, security attributes and the context of a process. Since the CPU handles the processes concurrently, when a process is pre-empted from the CPU, some information must be stored in order to allow a correct resume of the process itself when it comes running again. Such information define the context of the process.

data from our local database allows us to populate dynamically the views in order to support the user during the organisms selection. This task is performed by using dynamic queries instead of multiple connections through the KEGG service. Moreover, assuming the local presence of the KMGL files of the selected organisms, we can use the tool offline. More details about the use of the tool are given in Section 5.4.

### 5.3 Libraries and technologies

The software was developed using NetBeans IDE [27], a Java-based integrated development environment. It offers an interface to assist developers during coding. The choice to develop the tool in Java is related to the non-functional requirement concerning portability. In this way the application can run in every environment in which the JRE (Java Runtime Environment) is installed, thus ensuring its portability.

For the tool development, we have used external libraries written in Java language:

- **MySQL JDBC Driver** [28]: it allows to create an object for the connection to a MySQL database. In particular, it contains all the methods to perform operations on a specific database like insertions, deletions and data fetching.
- **Guava** [29]: it is an open source library developed by Google company. It contains methods to manage concurrency, I/O operations, string processing and so on. In our case, we use it because it allows to define multi-set structures with all the standard multi-set operations.
- **Poi** [30]: it is a library that belongs to the Apache POI Project and that is developed by the Apache software foundation. It represents the master project for the creation of Microsoft Office documents. This library is used in our application to manage the creation and modification of .xls files.
- **SaxParser** [31]: it is a Java library that allows one to perform XML data processing. It is more efficient wrt. a standard DOM parser since it doesn't load the document into memory and it doesn't create a representation of its file. In

fact it uses some callback functions to process the XML structure. The main functions are `startDocument()`, `endDocument()`, `startElement()` and `endElement()`. In particular, the last two methods are used to inform the client when a specific element tag is open or closed and to fetch all its attributes. The files are scanned sequentially.

- **Seaglass look and feel** [32]: this library is used as an alternative GUI style given by default from Swing Framework. The use of this package gives a better look and feel to the program.

## 5.4 Documentation

In this section we describe a typical example of use of our tool with the aim to provide a guideline for the users. The tool was thought to guide the user starting from the choice of the comparison to perform, passing through the selection of the organisms to be compared and ending with the selection of the comparison methods.

When we start the application, we see the main view shown in Fig. 5-2.

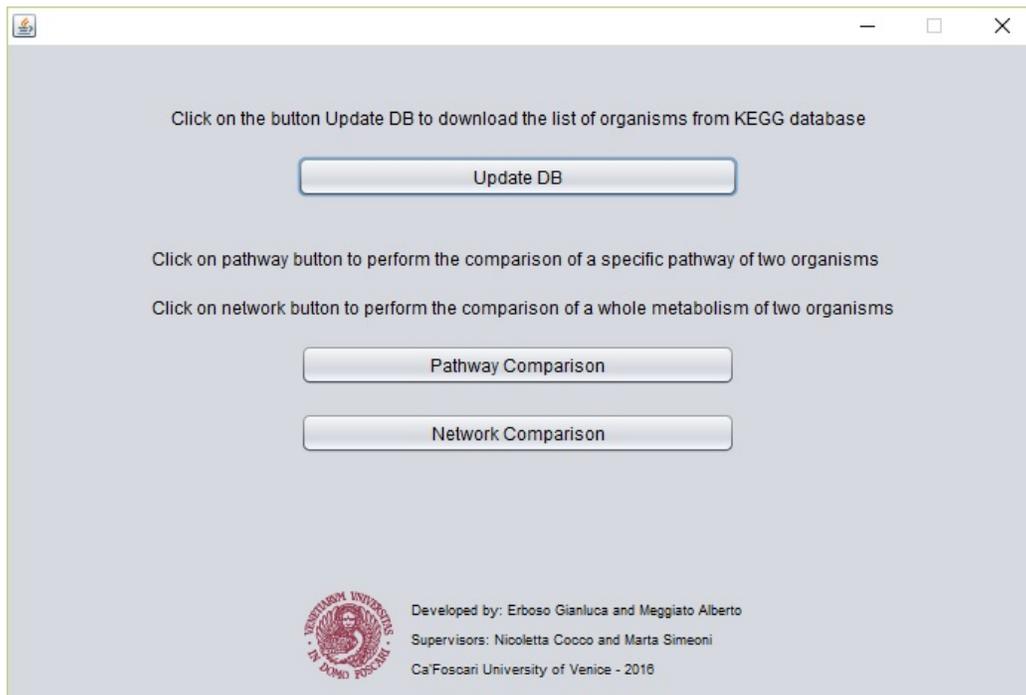


Figure 5-2: Main window of the tool

In this first step the user has three different choices:

- **Update Database:** the main advantage in doing such operation is that it maintains the application up to date, synchronizing the information on catalogued organisms in KEGG. In this way the application is not bound to a specific set of organisms, and it can be used always with the updated information given by KEGG. The updating of the information is a free choice of the user and it is not an automatic procedure;
- **Pathway:** this action allows the comparison of one or more metabolic pathways instead of comparing the entire metabolic networks. The user can select the metabolic function(s) from a predefined list of pathways and then it can select the organisms to compare;
- **Network:** this choice allows the comparison between entire metabolisms of different species.

Below we describe in detail the next steps when the user chooses to compare entire metabolisms.

Select the two organisms to analyze. If the KGML files are not present locally they will be automatically downloaded from the KEGG server, otherwise you can choose to download them again or use the local files.

Dominion	Eukaryotes	Eukaryotes
kingdom	Animals	Animals
subphylum	Vertebrates	Vertebrates
class	Mammals	Mammals
organism	hsa: Homo sapiens (human)	hsa: Homo sapiens (human)

Previous Next

Developed by: Gianluca Erbozo and Alberto Meggiato  
Supervisors: Nicoletta Cocco and Marta Simeoni  
Ca' Foscari University of Venice - 2018

Figure 5-3: Organisms selection

After the choice of the type of comparison, the view shown to the user is given in Fig. 5-3. In this step we give the possibility to select the organisms to compare. Due to the high number of the species catalogued in KEGG, we decided to support the user with a fast selection providing a hierarchical classification of the organisms. The classification is divided by dominion, kingdom, subphylum, class and organism. The population of each pull-down menu is dynamically executed according to the choices performed by the user. Once the organisms selection is done, the user is driven to the next step. Clicking on the next button, the tool checks if the KGML files are already present in the local folders. If the files are present, a new window is shown in order to provide the possibility to download the files again or not. In this way the user can decide either to execute a comparison by keeping the information up to date, wrt. the frequency of the KEGG updates, or to execute the comparison with the existing files. If the selected organisms have never been used in a comparison, the downloading procedure of the KGML files starts automatically before passing to the next window.

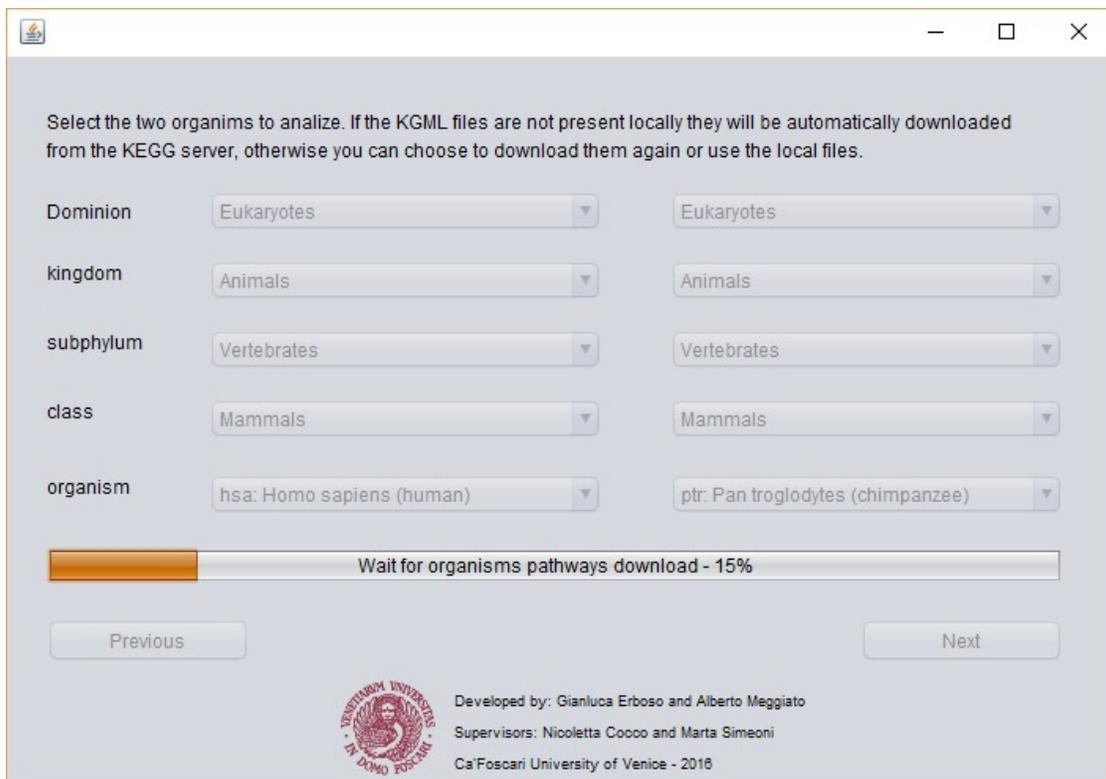


Figure 5-4: Downloading file for selected organisms

The last step is visible in Fig. 5-5. This window represents the core of the application since it permits the choice of the comparison measures both at pathway level and at network level. At pathway level it offers the possibility to compare pathways represented either as sets or as multisets of reactions. At the network level, the comparison methods are based either on directed or on undirected graphs. Moreover the user can set the  $\alpha$  parameter in order to associate a weight to the measures involved in the separated similarity index. Setting  $\alpha = 0,5$  the same relevance is given to *SimS* and *SimPW* indexes. Analogously, setting  $\alpha < 0,5$  more significance is given to *SimPW*, while with  $\alpha > 0,5$  more significance is given to *SimS*. After the setup of these parameters, the comparison can be launched by clicking on the start button.

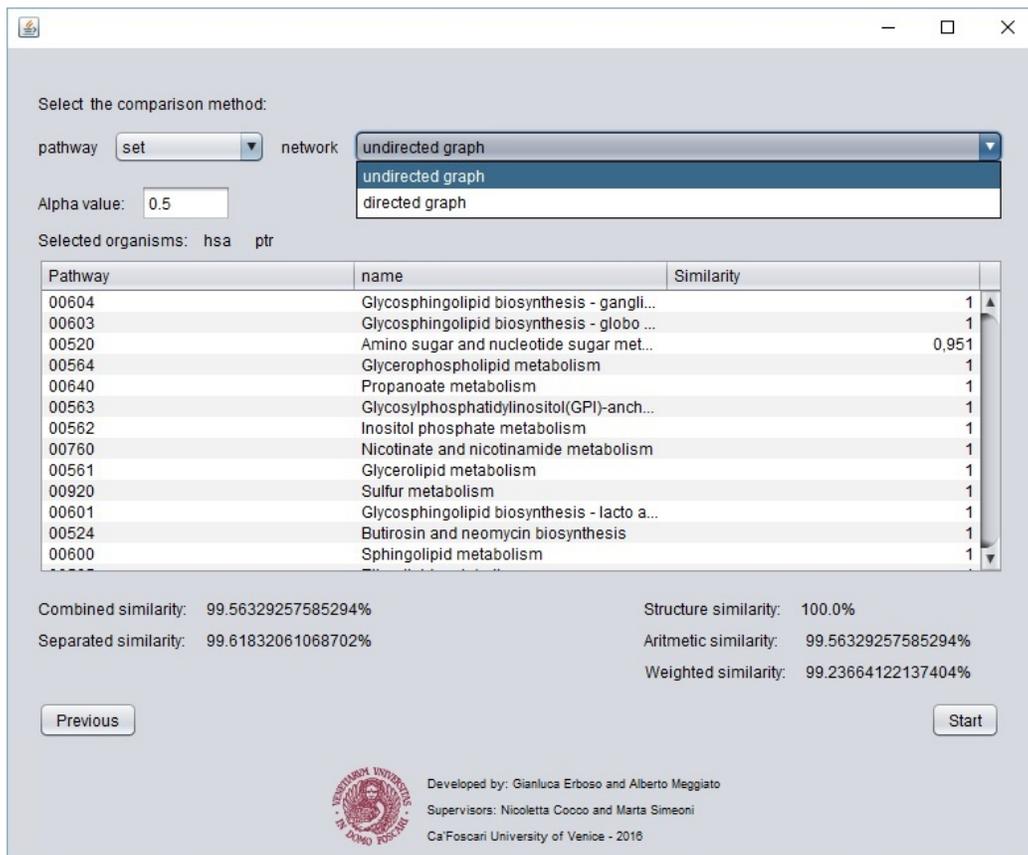


Figure 5-5: Results of networks comparison of hsa and ptr organisms with set representation of pathways and the separated index with  $\alpha = 0.5$

The computation takes a time that is related to the complexity of the networks and, from our tests, the average execution time is included between 20 and 90 seconds.

Then, the results are displayed on the same window through a summary table of similarities. The table has three columns showing the KEGG pathway number, the relative name and the similarity values computed for the same metabolic function in the two selected organisms. The results are automatically exported as .xsl file in the main tool's folder. After the comparison, further runs can be performed on the same organisms, by selecting a new method and clicking on start again.



# Chapter 6

## Experimenting with the tool

The comparison of different metabolisms as well as the comparison of metabolic pathways can be useful to discover similarities among organisms. In this chapter we describe the experiments performed with our tool to validate it. This validation is necessary on one hand since it is not possible to compare our results with other proposals in the literature either because their tools are not available or because their network are not based on KEGG's data on the other hand because there is no benchmark on which to perform a data comparison. We use a hierarchical clustering technique in order to provide a results classification and representation.

### 6.1 Cluster analysis

Clustering analysis is the process of organizing data into groups of observations related to each other. A cluster represents a collection of elements that are similar between them and dissimilar wrt. the elements contained in other clusters. This technique exploits a similarity measure in order to define the concepts of intracluster and intercluster distances. The intracluster measure represents a distance between inner elements of a cluster while the intercluster measure gives the distance wrt. the elements of the other groups. In our case we use a clustering algorithm in order to minimize the intracluster distance (high similarity between cluster's elements) and maximize the intercluster one (low similarity wrt. other cluster's elements). We use a

hierarchical clustering, a technique that gives a hierarchical organization of clusters. It produces a set of nested clusters that can be represented by dendrograms. A simple example is given in Fig.6-1.

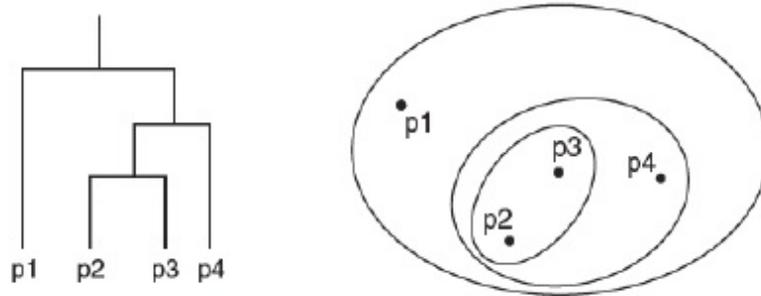


Figure 6-1: Example of nested cluster diagram and the corresponding dendrogram.[33]

There are some advantages in using this kind of algorithms: no assumption on clusters number are needed and moreover they are well suited in taxonomies representation. The hierarchical algorithms may use two different approaches:

- **Agglomerative:** it uses a bottom-up approach in which in the initial phase each element represents a singleton cluster. Then, at each iteration it merges pairs of nearest clusters until a unique cluster is obtained. This implies the use of proximity notion in order to define when two clusters can be merged or not;
- **Divisive:** it uses a top-down approach in which in the initial phase there is a unique cluster containing all the elements. Then, at each iteration it splits a cluster until it reaches the singletons. In this case, the algorithm chooses which cluster to split and how to perform the split [33].

Hierarchical algorithms require a similarity (or distance) matrix as input. Moreover, the key of these procedures is the computation of the proximity measure. Different definitions of proximity provide variants to the algorithms which can be used in specific cases. Examples of measures for proximity are the minimum, maximum or average distances between clusters. Below we give a pseudo-code of a basic agglomerative hierarchical clustering algorithm.

---

**Algorithm 1** Basic agglomerative hierarchical clustering algorithm.

---

```
1 Compute the proximity matrix, if necessary.  
2 Repeat:  
3   Merge the closest two clusters.  
4   Update the proximity matrix to reflect the proximity between the new cluster and the original  
   cluster.  
5 Until: Only one cluster remains.
```

---

To measure the distance between clusters we use the *complete linkage* method where the distance between the observations of the two clusters is the maximum one. In order to perform this analysis we exploit an existing implementation of linkage method given in MATLAB software. The *linkage* function and *dendrogram* function are used together to plot the phylogenetic tree. The similarity matrix given as input to the linkage method is created from our tool. Our experiments compare groups of organisms hence clustering techniques are a good way to represent the results. The similarity matrix for applying the clustering, represents the comparison between all possible pairs of such organisms. Our tool is fit to produce such similarity matrix.

## 6.2 Experiments

We discuss the experiments performed with our tool in order to evaluate the results. We conducted different kinds of experiments considering the entire metabolism of specific sets of organisms, selected by using different criteria. Generally, we use the default configuration of the tool both for pathways and networks representation. Namely we use sets as data structures for metabolic pathways and undirected graphs for metabolic networks and the CI as the similarity index. In the experiments in which we use the SI index, the default value for  $\alpha$  is 0.5. Moreover we consider SimPW instead of SimPA since it takes into consideration the number of reactions in the pathways providing a more refined measure. Different configurations are used to perform the experiment 2.

### 6.2.1 Experiment 1: Metabolic evolution in a group of species

The aim of the first experiment is to verify if the similarities in the metabolism's of a group of organisms find a correspondence in the phylogenesis due to evolution found in the literature [34] [35]. The experiment is executed considering organisms belonging to different taxonomic groups described in the Table 6.1, using the default configuration.

Code	Organism	Kingdom	Taxonomic group
<i>hsa</i>	<i>Homo sapiens</i> (human)	Animals	Mammals
<i>ptr</i>	<i>Pan troglodytes</i> (chimpanzee)	Animals	Mammals
<i>nle</i>	<i>Nomascus leucogenys</i> (gibbon)	Animals	Mammals
<i>mcf</i>	<i>Macaca fascicularis</i> (crab-eating macaque)	Animals	Mammals
<i>rno</i>	<i>Rattus norvegicus</i> (rat)	Animals	Mammals
<i>fca</i>	<i>Felis catus</i> (domestic cat)	Animals	Mammals
<i>gga</i>	<i>Gallus gallus</i> (chicken)	Animals	Birds
<i>cmg</i>	<i>Chelonia mydas</i> (green sea turtle)	Animals	Reptiles
<i>xla</i>	<i>Xenopus laevis</i> (African clawed frog)	Animals	Amphibians
<i>ola</i>	<i>Oryzias latipes</i> (Japanese medaka)	Animals	Fishes
<i>crg</i>	<i>Crassostrea gigas</i> (Pacific oyster)	Animals	Mollusks
<i>fve</i>	<i>Fragaria vesca</i> (woodland strawberry)	Plants	Rose family
<i>pti</i>	<i>Phaeodactylum tricornutum</i>	Chromista	Chromalveolata
<i>eco</i>	<i>Escherichia coli K-12 MG1655</i>	Bacteria	Proteobacteria

Table 6.1: Group of selected organisms.

What we expect is that our similarity indices produces a classification close to the phylogenetic one. The results of our tool with CI index are shown in Figure 6-2. As we can see, the main groups are clearly separated. There is a clear discrimination between animal's Kingdom and the other ones. Furthermore, in the Animals all the Mammals are grouped together and they are separated from Birds, Reptiles, Fishes and Mollusc. In more details in the Mammals, the distinction between primates and non-primates (*rno*, *fca*) is highlighted. The organisms more distant wrt. Animals (Plants, Protists and Bacteria) are split into another group. Moreover, we note that organisms that perform photosynthesis function are grouped together (*fve* and *pti*).

This experiment shows also some unexpected relations. The *nle* organism should be more similar to the *hsa* wrt. the *mcf* [36]. The same consideration is valid for the *xla* wrt. *ola*. In the last case, from a behavioural point of view, the two organisms have developed the ability to resist at the environmental changes.

We can conclude that our tool with the default setting allows organisms to be grouped together according to main taxonomy.

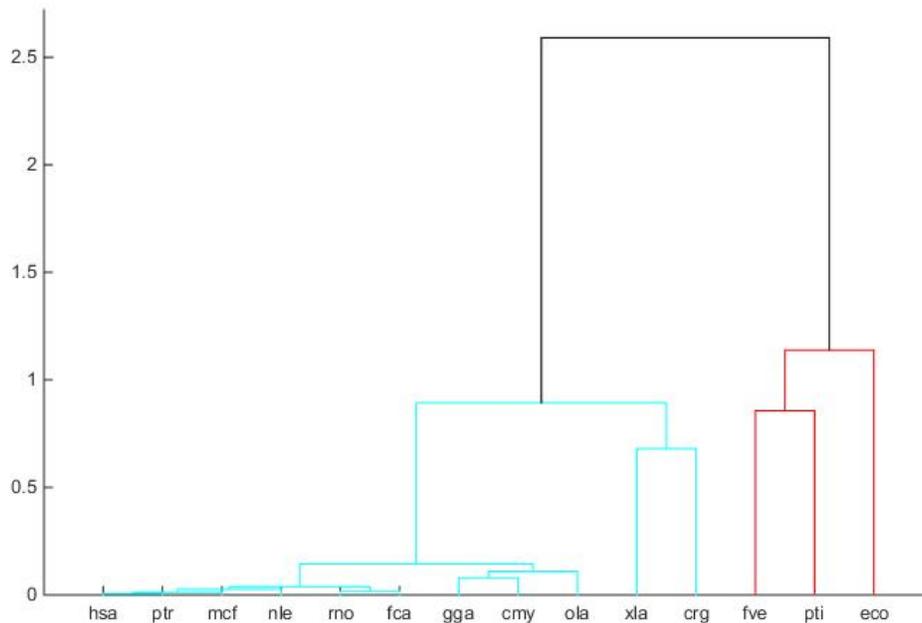


Figure 6-2: Phylogenetic tree produce by clustering with index CI.

## 6.2.2 Experiment 2: Yeasts and Molds metabolism

The second experiment is more refined since it is meant to test our tool wrt. the classification of a specific group of organisms belonging to the same Kingdom. Differently from the previous one, we select a specific group of organisms of the same Kingdom whose metabolism presents some differences. We select eight organisms among Fungi. The organisms used in the experiment are listed in table 6.2. In particular, we choose four yeasts (*sce*, *zro*, *tpf*, *cal*) and four molds (*fgr*, *tre*, *afm*, *abp*).

Code	Organism	Kingdom	Taxonomic group
<i>sce</i>	<i>Saccharomyces cerevisiae</i> (budding yeast)	Fungi	Saccharomycetes
<i>zro</i>	<i>Zygosaccharomyces rouxii</i>	Fungi	Saccharomycetes
<i>tpf</i>	<i>Tetrapisispora phaffii</i>	Fungi	Saccharomycetes
<i>cal</i>	<i>Candida albicans</i>	Fungi	Saccharomycetes
<i>fgr</i>	<i>Fusarium graminearum</i>	Fungi	Sordariomycetes
<i>tre</i>	<i>Trichoderma reesei</i>	Fungi	Sordariomycetes
<i>afm</i>	<i>Aspergillus fumigatus</i>	Fungi	Eurotiomycetes
<i>abp</i>	<i>Agaricus bisporus</i> var. <i>burnettii</i> JB137-S8	Fungi	Basidiomycetes

Table 6.2: Molds and Yeasts considered in the second experiment.

In this experiment we perform three tests using both the *CI* index and *SI* index with different  $\alpha$  value, in order to check if differences are detected. In the first experiment we use CI index, in the second one we use SI index with  $\alpha = 0.5$  and in the last one SI index with  $\alpha = 0.2$ . The following images Figure 6-3 and Figure 6-4 show the results achieved in the first two cases.

We note that the classification due to the clustering, produces two identical phylogenetic trees. Both the indices produce good results since we have an optimal separation at the top level between Yeasts and Molds, as expected from a phylogenetic point of view. The results with the two indices are different in the distance values as shown by the *y* axis.

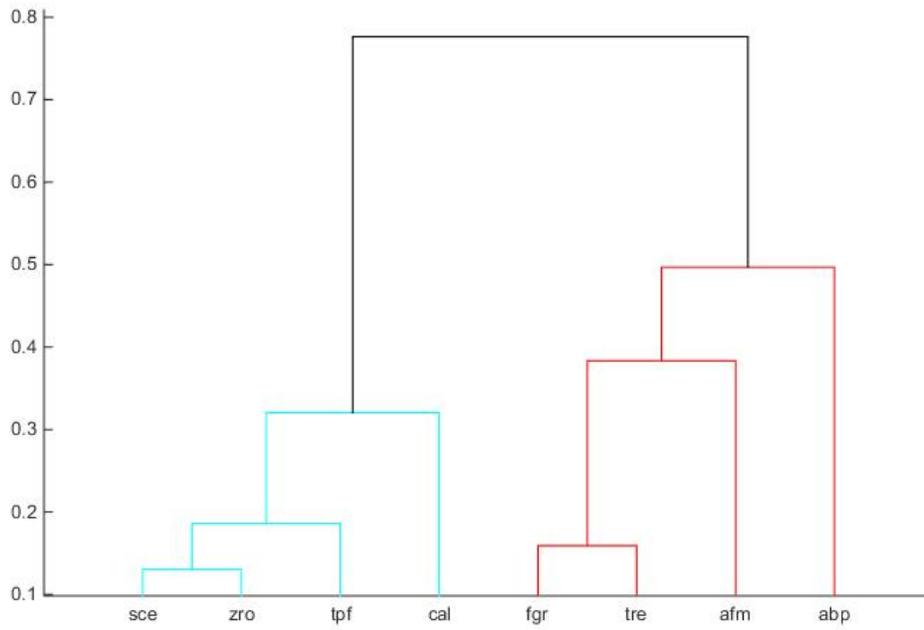


Figure 6-3: Clustering obtained using combined similarity index CI.

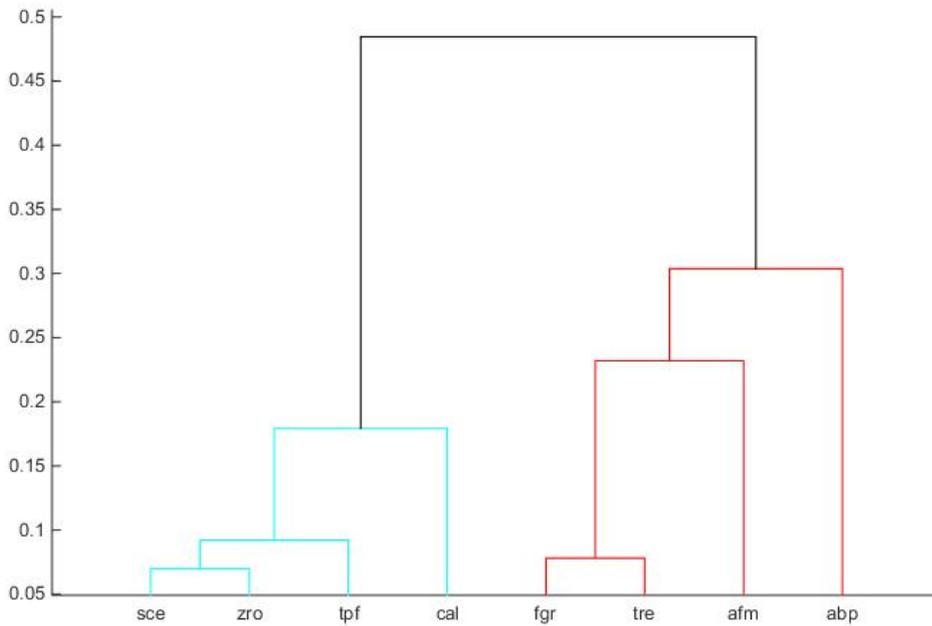


Figure 6-4: Clustering obtained using separated similarity index SI.

On the same group of organisms we consider also the classification due to structural similarity alone (SI with  $\alpha = 1$ ). The results are summarized in Table 6.3. Generally we note that the similarity values are high since they belong to the same Kingdom. For these reasons we perform a further test in which we give different weights to structure and pathways similarities. In particular, we use SI index with  $\alpha = 0.2$  in order to give lower weight to the structure (20%). The resulting dendrogram obtained from clustering is shown in Figure 6-5.

	sce	zro	tpf	cal	fgr	tre	afm	abp
sce	1	0,9583	0,9712	0,9365	0,8774	0,8761	0,7925	0,9185
zro	0,9583	1	0,9710	0,9774	0,8777	0,8764	0,8243	0,9183
tpf	0,9712	0,9710	1	0,9491	0,8526	0,8513	0,8007	0,9042
cal	0,9365	0,9774	0,9491	1	0,8991	0,8978	0,8438	0,9326
fgr	0,8774	0,8777	0,8526	0,8991	1	0,9953	0,8832	0,9346
tre	0,8761	0,8764	0,8513	0,8978	0,9953	1	0,8788	0,9335
afm	0,7925	0,8243	0,8007	0,8438	0,8832	0,8788	1	0,8295
abp	0,9185	0,9183	0,9042	0,9326	0,9346	0,9335	0,8295	1

Table 6.3: Structural similarities matrix

From the dendrogram we see that a separation between Yeasts and Molds is performed. However a distortion is introduced in the group of Yeasts. In particular, the *zro* organism is placed distant from *sce*. This is due to the fact that metabolisms structures have lower weights in the comparison. Therefore, we conclude that the structural similarity plays an important role in order to classify correctly the organisms [37].

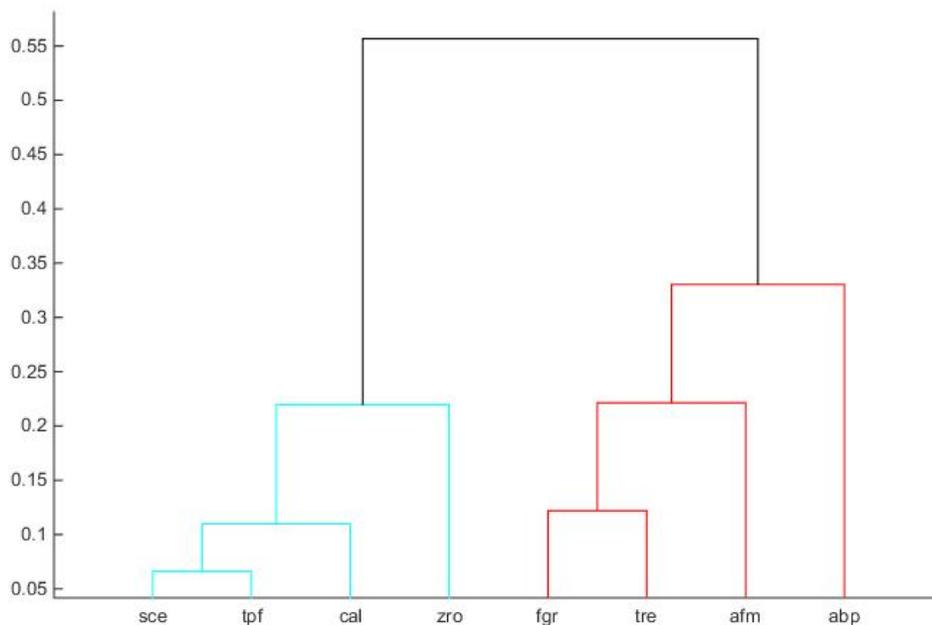


Figure 6-5: Clustering analysis using SI and  $\alpha = 0.2$ .

### 6.2.3 Experiment 3: Sulfur metabolism in different Kingdoms

In this experiment we consider specifically the *Sulfur metabolism* pathway (map: 00920 in KEGG). The Sulfur metabolism plays an important role on the amino acid construction, like the cysteine and methionine and other important molecules for the metabolism. Organisms take sulfur in different ways. Plants, Fungi and Bacteria take it and reduce it to sulfide, that is the simplest form of sulfur that can be used for the construction of the amino acids. The other Animals instead, take it indirectly from proteins that they assume through their diet [38].

For this experiment we choose organisms belonging to different Kingdoms considering their behaviour in sulfur reduction. We list the selected organisms in Table 6.4. For this experiment we use  $SimP_i$  index in order to compute the similarity of the Sulfur pathway in the two organisms and CI index to analyse their entire metabolism.

The expectations from this test are to obtain high similarity values for organisms belonging to the same Kingdom and low similarity values for organisms of different

Code	Organism	Kingdom	Taxonomic group
<i>hsa</i>	<i>Homo sapiens</i> (human)	Animals	Mammals
<i>ecb</i>	<i>Equus caballus</i> (horse)	Animals	Mammals
<i>gga</i>	<i>Gallus gallus</i> (chicken)	Animals	Birds
<i>tgu</i>	<i>Taeniopygia guttata</i> (zebra finch)	Animals	Birds
<i>ath</i>	<i>Arabidopsis thaliana</i> (thale cress)	Plants	Mustard family
<i>osa</i>	<i>Oryza sativa japonica</i> (Japanese rice)	Plants	Grass family
<i>bdi</i>	<i>Brachypodium distachyon</i>	Plants	Grass family
<i>nfi</i>	<i>Aspergillus fischeri</i>	Fungi	Eurotiomycetes
<i>ang</i>	<i>Aspergillus niger</i>	Fungi	Eurotiomycetes
<i>cpw</i>	<i>Coccidioides posadasii</i>	Fungi	Eurotiomycetes
<i>cow</i>	<i>Caldicellulosiruptor owensensis</i>	Bacteria	Caldicellulosiruptor
<i>toc</i>	<i>Thermosediminibacter oceani</i>	Bacteria	Thermosediminibacter
<i>hsl</i>	<i>Halobacterium salinarum</i>	Archaea	Halobacterium
<i>hvo</i>	<i>Haloferax volcanii</i>	Archaea	Haloferax
<i>pto</i>	<i>Picrophilus torridus</i>	Archaea	Picrophilus

Table 6.4: Set of considered organisms on *Sulfur metabolism*.

taxonomic groups. The results of the computation are shown in Table 6.5. We obtain expected results, coherent with our previous considerations: higher similarities are reached by the organisms belonging to the same Kingdom while lower similarities are found between organisms of different Kingdoms.

We represent the groups with different colours in the table. As we can see, the Archea group is not well distinguished since the comparison between the *Picrophilus torridus* organism and the other two Archea, produces low similarities. This is due to the fact that Archea considered in our experiment, constitute extreme ecological niches<sup>1</sup>. In particular, *hsl* and *hvo* are associated thanks to the ability to manage/resist to environments with high level of salinity.

For these reasons the metabolism of these Archea can be rather different and the comparison between them can produce low similarities.

We have performed the clustering using the similarity matrix in Table 6.5. The resulting dendrogram in Figure 6-6 shows that the tool provides a good classification

<sup>1</sup>Ecological niches [39] indicate the role, the chemical and the biological properties that permit the existence of an organism within an ecosystem. Extreme niches are organisms that live in extreme environments in which the biological life is constrained by particular conditions. Their survival is given by their adaptability.

	hsa	ecb	gga	tgu	ath	osa	bdi	nfi	ang	cpw	cow	toc	hsl	hvo	pto
hsa	1	1	1	0,7500	0,5385	0,5385	0,5385	0,5833	0,5833	0,5000	0,1000	0,2000	0,0833	0,1667	0,1667
ecb	1	1	1	0,7500	0,5385	0,5385	0,5385	0,5833	0,5833	0,5000	0,1000	0,2000	0,0833	0,1667	0,1667
gga	1	1	1	0,7500	0,5385	0,5385	0,5385	0,5833	0,5833	0,5000	0,1000	0,2000	0,0833	0,1667	0,1667
tgu	0,7500	0,7500	0,7500	1	0,3846	0,3846	0,3846	0,5455	0,5455	0,4545	0,1250	0,2500	0,1000	0,2000	0,2000
ath	0,5385	0,5385	0,5385	0,3846	1	1	1	0,5333	0,5333	0,5714	0,2500	0,3333	0,2143	0,2000	0,2000
osa	0,5385	0,5385	0,5385	0,3846	1	1	1	0,5333	0,5333	0,5714	0,2500	0,3333	0,2143	0,2000	0,2000
bdi	0,5385	0,5385	0,5385	0,3846	1	1	1	0,5333	0,5333	0,5714	0,2500	0,3333	0,2143	0,2000	0,2000
nfi	0,5833	0,5833	0,5833	0,5455	0,5333	0,5333	0,5333	1	1	0,9091	0,1667	0,2500	0,2308	0,3077	0,3077
ang	0,5833	0,5833	0,5833	0,5455	0,5333	0,5333	0,5333	1	1	0,9091	0,1667	0,2500	0,2308	0,3077	0,3077
cpw	0,5000	0,5000	0,5000	0,4545	0,5714	0,5714	0,5714	0,9091	0,9091	1	0,1818	0,2727	0,2500	0,2308	0,2308
cow	0,1000	0,1000	0,1000	0,1250	0,2500	0,2500	0,2500	0,1667	0,1667	0,1818	1	0,7500	0,3333	0,2857	0,1250
toc	0,2000	0,2000	0,2000	0,2500	0,3333	0,3333	0,3333	0,2500	0,2500	0,2727	0,7500	1	0,5000	0,4286	0,2500
hsl	0,0833	0,0833	0,0833	0,1000	0,2143	0,2143	0,2143	0,2308	0,2308	0,3333	0,3333	0,5000	1	0,8333	0,3750
hvo	0,1667	0,1667	0,1667	0,2000	0,2000	0,2000	0,2000	0,3077	0,3077	0,2308	0,2857	0,4286	0,8333	1	0,5000
pto	0,1667	0,1667	0,1667	0,2000	0,2000	0,2000	0,2000	0,3077	0,3077	0,2308	0,1250	0,2500	0,3750	0,5000	1

Table 6.5: Similarity matrix of the sulfur metabolism experiment

of the organisms. As we can see there is a clear distinction between the Kingdoms. Organisms belonging to same Kingdom are grouped together and they are discriminated wrt. the others. At the top level of the tree we find a discrimination between the Bacteria and all the other organisms. At the lower levels instead, Plants and Fungi are separated from Animals.

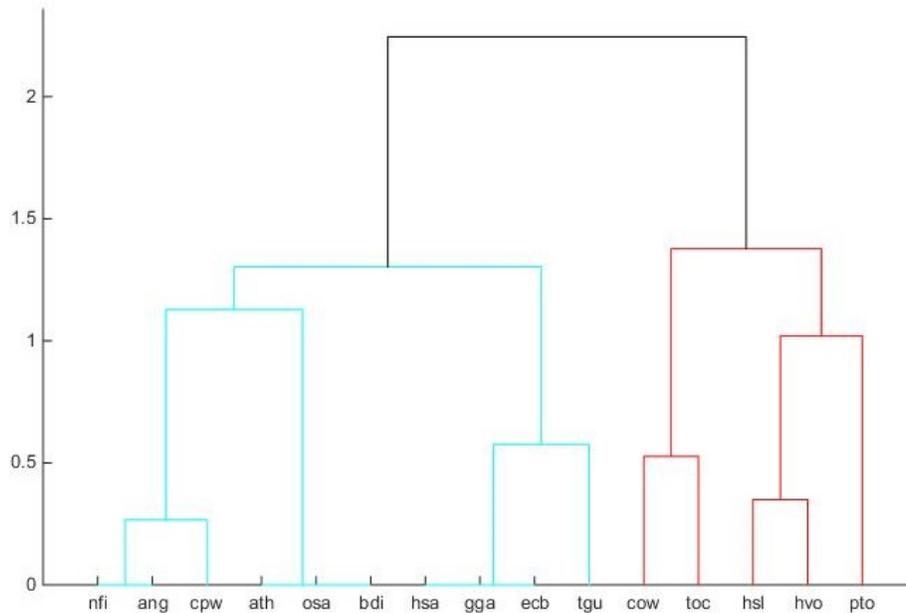


Figure 6-6: Clustering results on Sulfur metabolism.

We perform a further experiment using the same group of organisms and considering the entire metabolisms. The result of this experiment is shown in Figure 6-7. The tree underlines a significant difference in the classification of the organisms. In particular, at the highest level the algorithm provides a discrimination between Animals and all the other organisms. Moreover, Plants and Fungi are separated from Bacteria and Archea.

Considering the *hsa* and *gga* organisms, some differences are present. The analysis of these two organisms, tell us that they are more similar considering only Sulfur metabolism rather than the entire metabolism. Thus, the dimension of the considered dataset is relevant in the comparison.

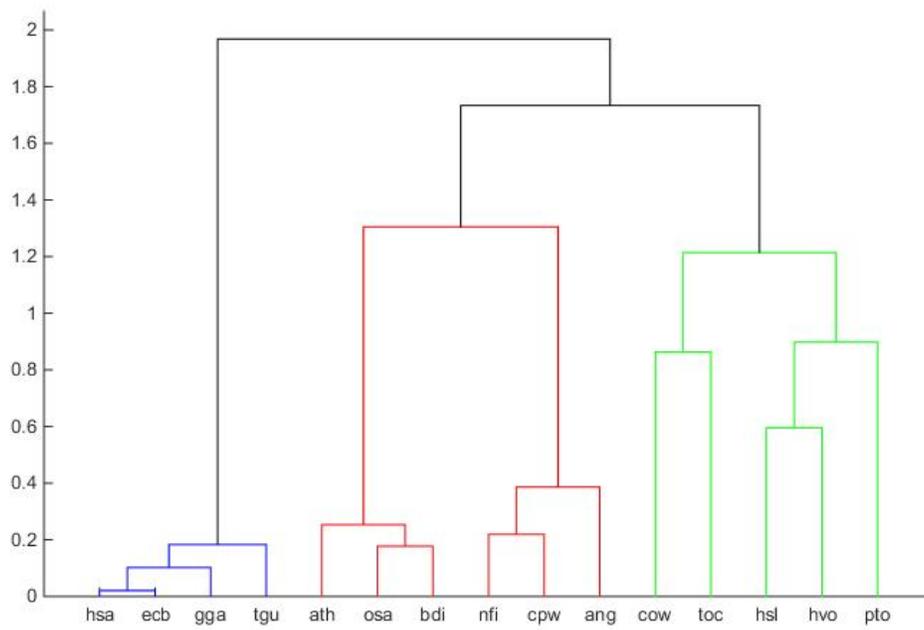


Figure 6-7: Organism classification obtained considering the entire metabolisms in experiments.

## 6.2.4 Experiment 4: Carbon fixation in photosynthetic organisms

This experiment considers the pathway *Carbon fixation in photosynthetic organisms* (map: 00710 in KEGG). This metabolic function refers to the conversion process of carbon dioxide to organic compound in photosynthetic organisms<sup>2</sup>. Since variants of this metabolic pathway exists due to environmental adaptations, we select a list of organisms that live in different environments. In Table 6.6 we give the organisms selected for the experiments.

Code	Organism	Kingdom	Taxonomic group
<i>gmx</i>	<i>Glycine max</i> (soybean)	Plants	Pea family
<i>pop</i>	<i>Populus trichocarpa</i> (black cottonwood)	Plants	Willow family
<i>vvi</i>	<i>Vitis vinifera</i> (wine grape)	Plants	Grape family
<i>osa</i>	<i>Oryza sativa japonica</i> (Japanese rice)	Plants	Grass family
<i>zma</i>	<i>Zea mays</i> (maize)	Plants	Grass family
<i>bdi</i>	<i>Brachypodium distachyon</i>	Plants	Grass family
<i>cre</i>	<i>Chlamydomonas reinhardtii</i>	Plants	Green algae
<i>vcn</i>	<i>Volvox carteri f. nagariensis</i>	Plants	Green algae
<i>npu</i>	<i>Nostoc punctiforme</i>	Bacteria	Nostoc
<i>acy</i>	<i>Anabaena cylindrica</i>	Bacteria	Anabaena
<i>oni</i>	<i>Oscillatoria nigro-viridis</i>	Bacteria	Oscillatoria
<i>mar</i>	<i>Microcystis aeruginosa</i>	Bacteria	Microcystis

Table 6.6: Selected organisms for *Carbon fixation* experiment.

The resulting similarity matrix is given in Table 6.7. As we can see, there is a clear separation between organisms that belong to the same Kingdom. Moreover, we note that the *Volvox carteri f. nagariensis* has a low similarity wrt. the other Plants which are coloured in green. The result can be reasonable since we are considering a particular organism, namely a *Green algae*. In general, Green algae should not be considered as Plants due to the fact that they don't have neither roots nor leaves. Furthermore, considering *cre* and *vcn* organisms, differences are related to the multicellular specie (*vcn*) that assume a simplified carbon fixation cycle wrt. to the others.

<sup>2</sup>Organisms that are able to synthesize organic compounds using the sunlight energy.[40]

	gmx	pop	vvi	osa	zma	bdi	cre	vcn	npu	acy	oni	mar
gmx	1	1	1	1	1	1	1	0,9524	0,6250	0,5833	0,6250	0,6250
pop	1	1	1	1	1	1	1	0,9524	0,6250	0,5833	0,6250	0,6250
vvi	1	1	1	1	1	1	1	0,9524	0,6250	0,5833	0,6250	0,6250
osa	1	1	1	1	1	1	1	0,9524	0,6250	0,5833	0,6250	0,6250
zma	1	1	1	1	1	1	1	0,9524	0,6250	0,5833	0,6250	0,6250
bdi	1	1	1	1	1	1	1	0,9524	0,6250	0,5833	0,6250	0,6250
cre	1	1	1	1	1	1	1	0,9524	0,6250	0,5833	0,6250	0,6250
vcn	0,9524	0,9524	0,9524	0,9524	0,9524	0,9524	0,9524	1	0,6522	0,6087	0,6522	0,6522
npu	0,6250	0,6250	0,6250	0,6250	0,6250	0,6250	0,6250	0,6522	1	0,9444	1	1
acy	0,5833	0,5833	0,5833	0,5833	0,5833	0,5833	0,5833	0,6087	0,9444	1	0,9444	0,9444
oni	0,6250	0,6250	0,6250	0,6250	0,6250	0,6250	0,6250	0,6522	1	0,9444	1	1
mar	0,6250	0,6250	0,6250	0,6250	0,6250	0,6250	0,6250	0,6522	1	0,9444	1	1

Table 6.7: Resulting similarity matrix in comparing organisms wrt. Carbon fixation.

Clustering produces the classification shown in Figure 6-8, and it is good since it groups correctly Plants and Bacteria at the top level and separates at lower level organisms with differences wrt. the metabolic function in analysis. In Plants we can see that the *Volvox carteri f. nagariensis* presents the problem described before: it is separated from the other Plants according with its simplified carbon fixation cycle.

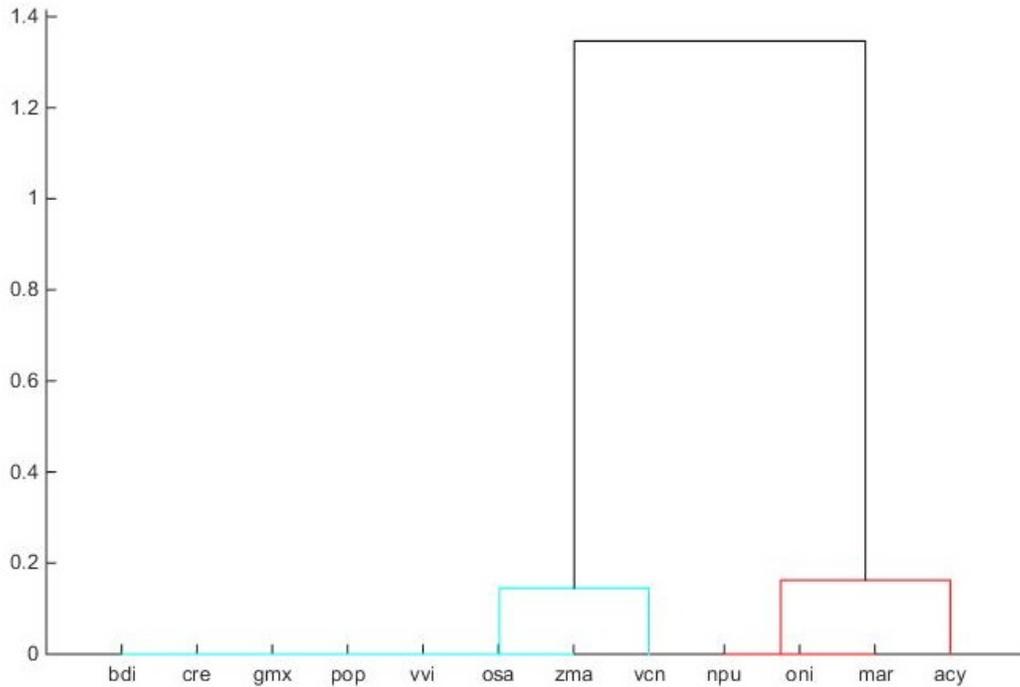


Figure 6-8: Clustering based on Carbon fixation in photosynthetic organisms.

The same group of organisms was compared considering the entire metabolisms. The result given in Figure 6-9 shows a phylogenetic tree in which we have a good separation between Kingdoms. Plants are discriminated from Bacteria at the top level. Then, a separation of Green algae from the other Plants is performed according to the initial considerations. In general, we can conclude that the CI index provides a good classification of the organisms in their Kingdoms. We also note that the *osa* organism is less similar than other plants. This classification can be reasonable because it is the unique plant that lives in highly hydrated environments (paddy field). Other consideration can be done considering Bacteria. *Npu* and *acy* are nitrogen-fixing cyanobacteria, *oni* and *mar*, instead, are cyanobacteria that produce toxins.

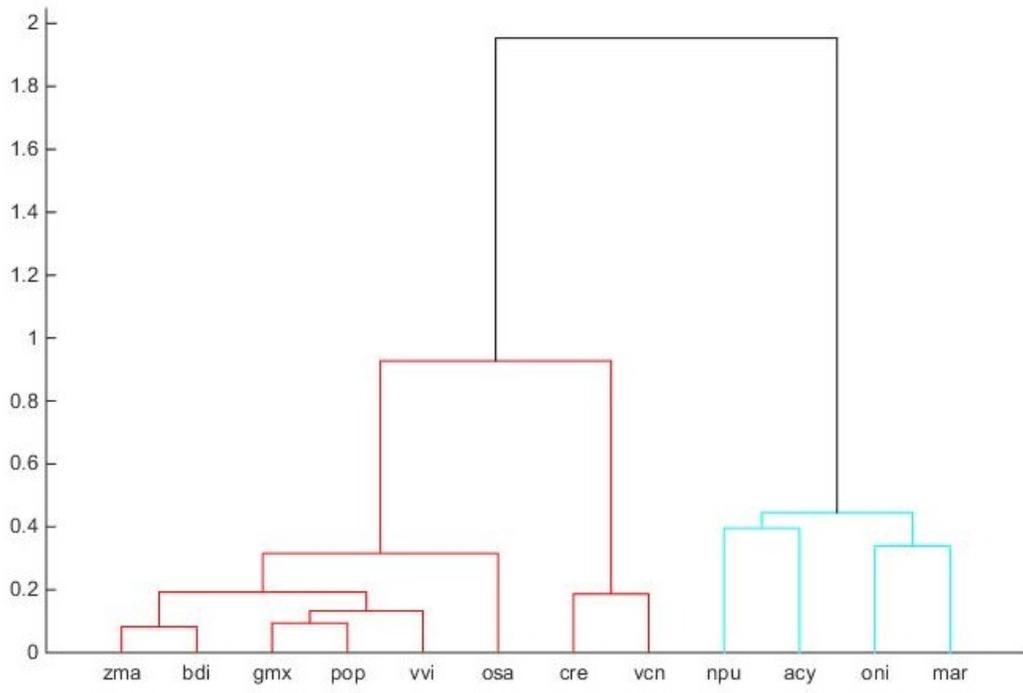


Figure 6-9: Phylogenetic tree from cluster analysis of the entire metabolisms in experiment 6.2.4.

## 6.2.5 Experiment 5: Glycolysis metabolism

The aim of this experiment is to give a classification of some organisms wrt. the Glycolysis pathway (map:00010 in KEGG). This pathway is responsible to convert the glucose into the pyruvate and during this process it generates energy in the form of ATP. For this experiment we choose a set of organisms which differ wrt. sugar metabolism. We use different configurations in order to perform the experiment. In particular, for the specific pathway analysis we use both set and multiset data structure, and undirected graph for networks. For the global similarity indices we consider both CI and SI with different values of  $\alpha$  (0.25, 0.5, 0.75). Below we list the organisms considered for the experiment. They can be divided in four different groups: nitrogen-fixing Bacteria, methanogen Archaea, sulfate-reducing Bacteria, sulfate-reducing Archaea.

Code	Organism	Kingdom	Taxonomic group
<i>dvu</i>	<i>Desulfovibrio vulgaris Hildenboroug</i>	Bacteria	Desulfovibrio family
<i>sfu</i>	<i>Syntrophobacter fumaroxidans</i>	Bacteria	Syntrophobacter
<i>rsp</i>	<i>Rhodobacter sphaeroides 2.4.1</i>	Bacteria	Rhodobacter
<i>cdf</i>	<i>Peptoclostridium difficile 630</i>	Bacteria	Peptoclostridium
<i>drm</i>	<i>Desulfotomaculum reducens</i>	Bacteria	Desulfotomaculum
<i>ana</i>	<i>Nostoc sp. PCC 7120</i>	Bacteria	Nostoc
<i>npu</i>	<i>Nostoc punctiforme</i>	Bacteria	Nostoc
<i>tye</i>	<i>Thermodesulfovibrio yellowstonii</i>	Bacteria	Thermodesulfovibrio
<i>msi</i>	<i>Methanobrevibacter smithii</i>	Archaea	Methanobrevibacter
<i>mel</i>	<i>Methanobacterium lacus</i>	Archaea	Methanobacterium
<i>afu</i>	<i>Archaeoglobus fulgidus DSM 4304</i>	Archaea	Archaeoglobus
<i>thg</i>	<i>Thermogladius cellulolyticus</i>	Archaea	Thermogladius
<i>cma</i>	<i>Caldivirga maquilingensis</i>	Archaea	Caldivirga

Table 6.8: Organisms considered for experiment 6.2.5

In particular, *dvu*, *sfu*, *drm*, *tye* are sulfate-reducing Eubacteria, *afu*, *thg*, *cma*, are sulfate-reducing Archaeabacteria, *ana*, *npu*, *cdf*, *rsp* are nitrogen-fixing Bacteria and *msi*, *mel* are methanogen Archaeabacteria. From the test we expect to obtain a good distinction of the above groups.

As we can see, the results in Figure 6-10 give a classification of the organisms with some distortions. In facts, *sfu* and *thg* are placed inside the wrong group. However

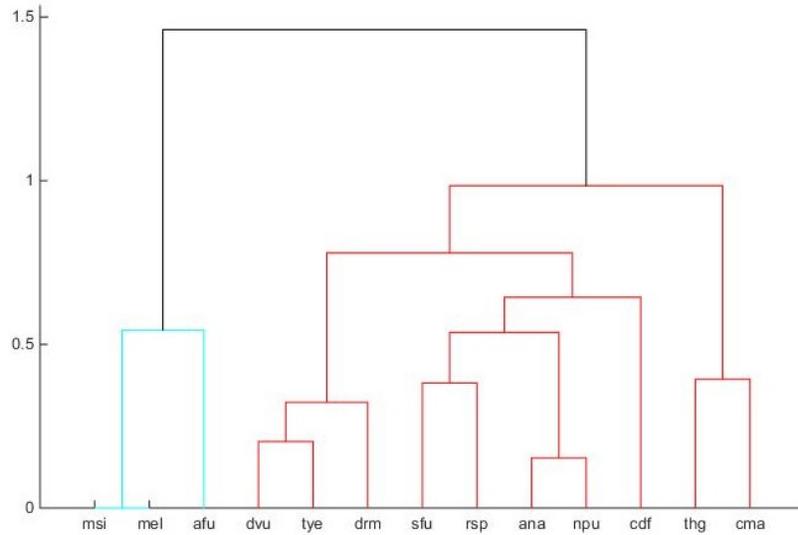


Figure 6-10: Clustering on Glycolysis pathway.

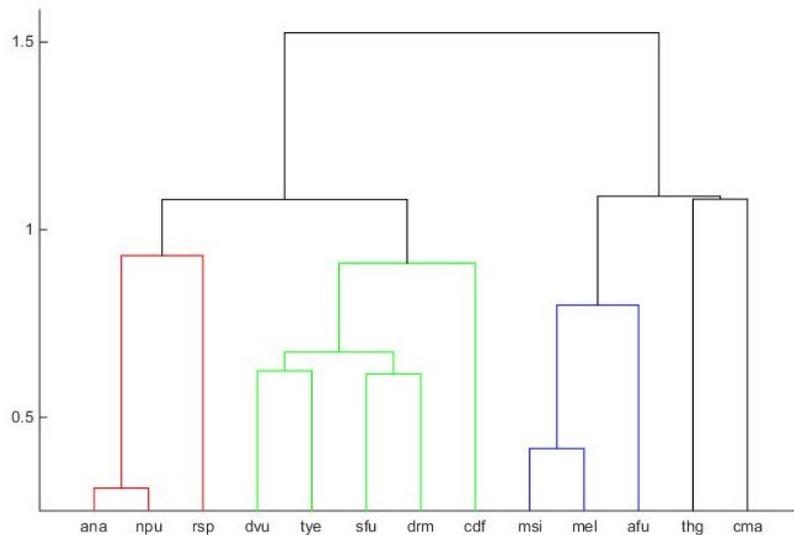


Figure 6-11: Classification based on Bacteria and Archea on the entire metabolism.

these two organisms are grouped correctly at lower level with organisms of the same Kingdom. In particular, in biology, *thg* and *cma* are inclined to degrade carbohydrate-based compounds. Considering the results we decided to perform another test on the same group of organisms taking into account the entire metabolism. In this case the result given in Figure 6-11 provides a clear discrimination between Kingdoms.

Analysing the results obtained with the SI index and different  $\alpha$  values, the Kingdoms' discrimination is maintained.

### 6.2.6 Conclusion

Analysing the behaviour of the tool in the experiments, we can make some considerations. First of all, we note that the structure of the metabolic networks is relevant in order to obtain a good classification. In particular, considering only the metabolic networks functionalities, the results present some distortions. This may be due since the measure in metabolic pathway comparison assumes a certain level of abstraction. Considering the two global similarity indices, the CI index is in general better than the SI index. However, the SI index permits us to tune the  $\alpha$  value in order to weight structure and functionality of the network. Finally, we can conclude that the CI index, in all experiments, provides a good classification between Kingdoms.

# Chapter 7

## Conclusion

The aim of this thesis is to propose a new approach to compare the entire metabolism between different species considering both the topology of the metabolic network and its functionalities. This comparison is useful to discover similarities or dissimilarities between organisms providing information about the evolutionary process.

In the literature, the proposed techniques build and compare the metabolic networks in detail. This fact produces many problems from a complexity point of view.

Our approach, is based on KEGG database information and on the implicit mapping between metabolic pathways defined by the reference pathways. In the approach, two independent levels are defined to reduce the complexity of the network and for exploiting the standardized modularization of the data adopted by KEGG.

The metabolism comparison, we propose, is performed combining two independent measures. The first one, developed in [1], evaluates the similarity between metabolic pathways, the second one, described in this thesis, evaluates the similarity of the topology of the metabolic networks represented as graphs.

Four similarity indexes have been defined:  $SimS_i$ , that computes the similarity between matching nodes (i.e. nodes representing same metabolic function in the two metabolisms) in terms of connections;  $SimS$  that represents the topological similarity between the two entire nets;  $CI$  and  $SI$  that provide a global similarity measure combining the indexes defined in [1], based on corresponding pathway similarities, with the above ones.

Our method has been implemented in a Java tool that relies on KEGG database information. The program allows to compare the metabolism between pairs of organisms, selected by the user, providing different similarity measures. Some experiments have been executed considering both the entire metabolisms and specific metabolic functions on selected sets of organisms. The results are represented in a tree using a hierarchical clustering algorithm. Our analysis of the experiments permits us to conclude that our algorithm is able to classify correctly wrt. evolution organisms belonging to the same Kingdom. In specific cases some distortions are verified in comparing organism to the same taxonomic group. This is probably due to the level of abstraction of the measure on the metabolic pathways.

Further developments of our proposal can be considered. More experiments could be executed to determine a threshold value on the similarity measure for each group of organisms belonging to various Kingdom. Furthermore the tool can be extended, thanks to its modular structure, implementing new comparison methods both for networks and pathways and implementing new functionalities in order to allow comparison of specific pathways on specific groups of organisms. Again, a clustering algorithm can be integrated in order to provide the phylogenetic tree of the considered organisms.

# Bibliography

- [1] Alberto Meggiato. Comparing metabolic networks at pathway level. Master's thesis, Ca'Foscari University of Venice, 6 2016.
- [2] Christophe H. Schilling, Stefan Schuster, Bernhard O. Palsson, and Reinhart Heinrich. *Metabolic Pathway Analysis: Basic Concepts and Scientific Applications in the Post-genomic Era*.
- [3] Michael Palmer. In *Human Metabolism*, chapter Introduction, pages 1–2. Department of Chemistry, University of Waterloo, 2015.
- [4] Donald Voet, Charlotte W. Pratt, and Judith G. Voet. In *Fundamentals of Biochemistry: Life at the Molecular Level*, pages 436–439, 442. John Wiley and Sons, 4 edition, 2012.
- [5] Paolo Baldan, Nicoletta Cocco, Andrea Marin, and Marta Simeoni. *Petri nets for modelling metabolic pathways: a survey*. Natural Computing, 9(4):955–989, 2010.
- [6] Antonio Albano, Giorgio Ghelli, and Renzo Orsini. Fondamenti di basi di dati. page 8. Zanichelli, 2005.
- [7] Kanehisa Laboratories. *Kyoto Encyclopedia of Genes and Genomes*. <http://www.genome.jp/kegg/>, 2015.
- [8] SRI International. Biocyc database collection. <http://biocyc.org/>, 2015.
- [9] Fellowship for Interpretation of Genomes. Seed. <http://www.theseed.org/>, 2003.

- [10] European Molecular Biology Laboratory. The european bioinformatics institute. <https://www.ebi.ac.uk>, 2016.
- [11] Hiroyuki Ogata, Susumu Goto, Kazushige Sato, Wataru Fujibuchi, Hidemasa Bono, and Minoru Kanehisa. *KEGG: Kyoto Encyclopedia of Genes and Genomes*. Oxford University Press, 28, 2000.
- [12] Minoru Kanehisa, Susumu Goto, Yoko Sato, Masayuki Kawashima, Miho Furumichi, and Mao Tanabe. *KEGG: Data, information, knowledge and principle: back to metabolism in KEGG*. Nucleic Acids Research, 42, 2014.
- [13] Hongwu Ma and An-Ping Zeng. *Reconstruction of metabolic networks from genome data and analysis of their global structure for various organisms*. Bioinformatics, 19(2):270–277, 2003.
- [14] H. Jeong, B. Tombor, R. Albert, Z. N. Oltvai, and A. L. Barabasi. *The large scale organization of metabolic networks*. Nature, 407:651–654, 2000.
- [15] Markus Rohrschneider. *Visualization of Metabolic Networks*. Master’s thesis, Universität Leipzig, 2015.
- [16] C. V. Forst, C. Flamm, I. L. Hofacker, and P. F. Stadler. *Algebraic comparison of metabolic networks, phylogenetic inference, and metabolic innovation*. BMC Bioinformatics, 7(1):1–11, 2006.
- [17] I. Zevedei-Oancea and S. Schuster. *Topological analysis of metabolic networks based on Petri Net theory*. In Silico Biology, 3(3):323–345, 2003.
- [18] Y. Toshato. *A Method for Species Comparison of Metabolic Networks Using Reaction Profile*. IPSJ Digital Courier, 2(0):685–690, 2006.
- [19] S. H. Hong, T. Y. Kim, and S. Y. Lee. *Phylogenetic analysis based on genome-scale metabolic pathway reaction content*. Applied Microbiology and Biotechnology, 65(5):203–210, 2004.

- [20] Roman L. Tatusov, Eugene V. Koonin, and David J. Lipman. *A Genomic Perspective on Protein Families*. Science, 278(5338):631–637, 1997.
- [21] Z. Li, S. Zhang, Y. Wabg, X. Zhang, and L. Chen. *Alignment of molecular networks by integer quadratic programming*. Bioinformatics, 23(13):1631–1639, 2007.
- [22] Ay et al. *Metabolic network alignment in large scale by network compression*. BMC Bioinformatics, 13(suppl 3):1–19, 2012.
- [23] O. Ebenhoh, T. Handorf, and R. Heinrich. *A Cross Species Comparison of Metabolic Network Functions*. Genome Informatics, 16(1):203–213, 2005.
- [24] Random House Unabridged Dictionary. Dictionary.com. <http://www.dictionary.com/browse/phylum>, 2016.
- [25] Robert Eckstein. Java se application design with mvc. <https://www.oracle.com/technetwork/articles/javase/index-142890.html>, 2007.
- [26] Oracle Corporation. Mysql database. <https://www.mysql.it>, 2016.
- [27] Sun Microsystems. Netbeans ide. <https://netbeans.org>, 2016.
- [28] Oracle Corporation. Mysql jdbc. <https://www.mysql.it/products/connector>, 2016.
- [29] Google. Mysql jdbc. <https://github.com/google/guava/wiki>, 2015.
- [30] The Apache Software Foundation. Apache poi. <https://poi.apache.org>, 2016.
- [31] Oracle Corporation. Sax. <http://www.saxproject.org/quickstart.html>, 2016.
- [32] Kathryn Huxtable. Seaglass look and feel. <https://github.com/khuxtable/seaglass/wiki>, 2015.
- [33] Pang-Ning Tan, Michael Steinbach, and Vipin Kumar. *Introduction to Data Mining*. 2005.

- [34] Mitchell L. G. and Reece J. B. *Biology: concepts and connections*. chapter 14, 19. Benjamin and Cummings Pub. Comp., 2003.
- [35] Arnold M. L. *Evolution through genetic exchange*. chapter 2, 9. Oxford University Press Inc., 2008.
- [36] R. Martin R. Classification of primates. In *The Cambridge Encyclopedia of Human Evolution*. Cambridge University Press, 1994.
- [37] Kurtzman C. P. and Robnett C. J. *Phylogenetic relationships among yeasts of the Saccharomyces complex determined from multigene sequence analyses*. FEMS Yeast Research, 4(3):417–432, 2006.
- [38] Thomas Leustek, Melinda N. Martin, Julie-Ann Bick, and John P. Davies. Pathways and regulation of sulfur metabolism revealed through molecular and genetic studies. *Annual Review of Plant Physiology and Plant Molecular Biology*, 51(1):141–165, 2000.
- [39] Eric R. Pianka. *Evolutionary ecology*. Benjamin Cummings, 6 edition, 1999.
- [40] Stefania Azzolini. Treccani.it - enciclopedia della scienza e della tecnica. [http://www.treccani.it/enciclopedia/autotrofo\\_\(Enciclopedia-della-Scienza-e-della-Tecnica\)/](http://www.treccani.it/enciclopedia/autotrofo_(Enciclopedia-della-Scienza-e-della-Tecnica)/), 2016.