# A network-based approach to cell metabolism: from structure to flux balances

Oriol Güell Riera

Doctoral program: Ciència i Tecnologia de Materials

---

# A network-based approach to cell metabolism: from structure to flux balances

---

**Doctoral thesis**

*Author:*

Oriol Güell Riera

*Advisors:*

Dr. M. Ángeles Serrano Moral

Departament de Física Fonamental, Universitat de Barcelona

Prof. Francesc Sagués Mestre

Departament de Química Física, Universitat de Barcelona

Universitat
de Barcelona

**Aquesta tesi doctoral està dedicada a les persones que estimo**

*"One of the principal objects of theoretical research is to find the point of view from which the subject appears in the greatest simplicity."*

Josiah Willard Gibbs

# *Acknowledgements*

Vull donar les gràcies als meus pares, Ramon i Teresa, també a l'Anna i l'Albert, sense oblidar a la resta de la família, els avis i les àvies, els tetes Pere i Ferran i els seus fills! Maite Zaitut! Una especial menció també per la Lis i la Nina, que heu estat allà fent companyia durant tota l'escriptura d'aquesta tesi sense cap queixa. Sou unes gossetes increïbles!

Vull agrair a la Irene, per estar sempre allà quan fas falta i ajudar sempre que hi ha problemes.

També vull agrair al Dani, l'Édgar, l'Sveto, l'Edo, l'Àngel i el Rubén. Sois grandes amigos y los viajes a Madrid para ir a veros siempre han sido grandes. Eusebio y Gañán! Una altra menció al Fran i al Francesco per els cops que hem anat a sopar per parlar, d'entre molts temes, de la feina que hem fet junts. També vull agrair al Pol, l'Ula i al Mario pels bons moments que hem passat.

Una altra especial i important menció als col·legues d'institut, el Gerard, el Guilla, l'Armand, l'Uri, el Fido, el Duñó, el Víctor i l'Àlex. Sempre vau estar, esteu i estareu allà quan feu falta. Ja sabeu que no falta tall.

Una altra menció important pel Trapote i pel Dani Igual. Sempre ajudeu quan podeu i això no s'oblida mai. Sou molt grans! També vull donar les gràcies al

Yerai per totes les rutes amb moto que hem fet, siguin ON o OFF. Gas! També agrair a l'Eric i al Jordi per altres rutes amb moto!

També vull donar les gràcies a la gent de la uni. Al Manu per la seva constant i gran ajuda, a la Cristina per respondre tantes preguntes, a la Isabel per haver estat sempre allà per parlar i explicar el que sigui, i al Francesc Mas, per la seva constant ajuda i pels seus consells. Faig una especial menció pel Vela, la Maria, l'Anna, el Sergi, el Fran i la Laura, pels bons moments i per les nombroses partides, tant al "Bang" com "La Resistencia" o altres jocs.

# Contents

# Abbreviations

**ORGANISMS**

*E. coli*    *Escherichia coli*

*M. pneumoniae*  *Mycoplasma pneumoniae*

*S. aureus*    *Staphylococcus aureus*


**METHODS**

**DF**      **D**isparity **F**ilter

**DP**      **D**egree **P**reserving randomization

**FBA**      **F**lux **B**alance **A**nalysis

**FBA-MBR**   **FBA - M**aximum **B**iomass **R**ate

**FFP**      **F**easible **F**lux **P**henotypes

**FVA**      **F**lux **V**ariability **A**nalysis

**HR**      **H**it-And-**R**un

**K-S**      **K**olmogorov-**S**mirnov test

**MB**      **M**ass-**B**alanced randomization

**PCA**      **P**rincipal **C**omponent **A**nalysis


**CONCEPTS**

**GCC**      **G**iant **C**onnected **C**omponent

| | |
|---|---|
| **GENRE** | **GEN**ome-scale **RE**construction |
| **OMP** | **O**ne-Mode **P**rojection |
| **PSL** | **P**lasticity **S**ynthetic **L**ethality |
| **RSL** | **R**edundancy **S**ynthetic **L**ethality |
| **SCC** | **S**trongly **C**onnected **C**omponent |
| **SL** | **S**ynthetic **L**ethality |

**PATHWAYS**

| | |
|---|---|
| **ACM** | **A**lternate **C**arbon **M**etabolism |
| **AR** | **A**naplerotic **R**eactions |
| **B** | **B**iomass |
| **CAC** | **C**itric **A**cid **C**ycle |
| **CEB** | **C**ell **E**nvelope **B**iosynthesis |
| **GM** | **G**lutamate Metabolism |
| **GG** | **G**lycolysis/**G**luconeogenesis |
| **IITM** | **I**norganic **I**on **T**ransport and **M**etabolism |
| **MLM** | Membrane **L**ipid **M**etabolism |
| **NSP** | Nucleotide **S**alvage **P**athway |
| **OP** | **O**xidative **P**hosphorylation |
| **PPB** | **P**urine and **P**yrimidine **B**iosynthesis |
| **PPP** | Pentose **P**hosphate **P**athway |
| **PM** | **P**yruvate **M**etabolism |
| **TE** | **T**ransport, **E**xtracellular |
| **TIM** | **T**ransport, **I**nner Membrane |

**UNITS**

| | |
|---|---|
| **gDW** | **g**rams **D**ry **W**eight |

# Chapter 1

# Cellular metabolism at the systems level

This chapter reviews basic concepts of cellular metabolism. First, an overall view of the architecture of cellular metabolism is given, from the large-scale of Catabolism and Anabolism to biochemical pathways, reactions, and metabolites. Fundamental concepts of chemical kinetics and thermodynamics are mentioned, followed by a brief consideration of key ideas about regulation, control, and evolution of metabolism. Finally, the need for a systems-level approach is discussed. Aims and objectives, together with an outline of this thesis, are included at the end of the chapter.

Cellular metabolism is composed of enzyme-controlled biochemical reactions. They form a densely-connected metabolic network which is responsible of maintaining cells alive by generating chemical energy and by synthesizing important metabolic intermediates from nutrients taken from the environment. Over the

years, cellular metabolism has attracted the attention of many researchers. At the end of the 19th century, the view of metabolism was dominated by studies of specific biochemical reactions or processes. It is worth mentioning in this respect the work of Eduard Buchner who, based on previous work by Louis Pasteur, demonstrated that cell-free biochemical extracts of yeast -known today as enzymes- could catalyze alcoholic fermentation. This put an end to vitalism-based ideas and boosted the then emerging field of biochemistry [1]. Later on, with the help of experimental techniques such as NMR sprectroscopy add X-ray diffraction, the idea of the organization of reactions into sequences of consecutive transformations or pathways arose, creating the basis of modern biochemistry [2]. In principle, pathways were treated as entities with a definite function which operated independently of each other. Despite the enormous success achieved by biochemistry, studies focusing on single reactions, enzymes, or even single pathways are not sufficient to explain most experimental results on metabolism at the functional level, which require a high knowledge of the entire map of metabolic interactions and their interplay with other cellular components. Examples of these results are the identification of redundant metabolic pathways [3], or the observation of the effect known as synthetic lethality [4], which arises when a combination of mutations leads to cell death, whereas the individual mutations are not lethal.

Since metabolic phenotypes[1] and behavior emerge from the interactions of many metabolic reactions and other cell components, understanding them at the systems level is crucial for our understanding of living cells. Metabolism is not isolated from the rest of the cell machinery. Therefore, a key challenge in biology is to integrate all the knowledge about the constituents of cells, from genes, to proteins, to metabolites, and reactions, in order to understand how they

---

[1]A phenotype is the composite of the observable characteristics of an organism, such as its morphology, development, biochemical or physiological properties.

interact and how these interactions determine the behavior of cells [5]. This implies a wide knowledge on how reactions are interconnected with metabolites to integrate a whole metabolic network. One can use this metabolic map to study, for example, how different pathways interact [6, 7]. A clear understanding of all these metabolic interactions, and their linkages and interdependencies with other biological scales like genetic networks, will allow us to decipher crucial questions, such as how cells are able to adapt to their environment, or in which way evolutionary processes led to the properties of metabolism as we currently observe them.

The study of integrated metabolic maps is difficult due to the inherent complexity of these intricate systems composed of thousands of interacting reactions. To ease the understanding of cellular metabolism as a complex system, the classical reductionist approach has given way to the so-called *systems-level approach*, which studies metabolism as a whole, taking into account the largest number of experimentally known constituents of the metabolic network, their interactions, and the linkages to other cell constituents such as enzymes, proteins, and genes. This emerging paradigm for the study of cell metabolism is at the core of an emerging interdisciplinary field called Systems Biology [8–11], which uses a holistic approach to understand the relationships between structure and function in biological systems, an impossible endeavor for studies that focus on specific reactions, enzymes, or metabolic processes. The use of this approach has provided a large amount of new validated hypotheses, like the heterogeneity of physiological metabolic fluxes in cells [12], or the phylogenetic analysis of metabolic environments that determine which components must be exogenously acquired [13]. Along with the development of *Complex Network Science* [14, 15], the systems-level approach has led to a huge increase in our understanding of how metabolic networks operate.

## 1.1    A brief introduction to cellular metabolism

Cellular metabolism comprises the complete set of chemical reactions at the cell level needed for life. While chemical syntheses in laboratory focus on specific sequences of chemical reactions in order to optimize processes, thousands of reactions, tightly interconnected through common metabolites, take place simultaneously in cells, forming a network that is precisely controlled by the combined action of enzymes, genes, etc., in order to secure functions. This network takes part in the growth of cells, in the maintenance and construction of their structures, and in the response and adaptation of the cell to different environmental conditions or internal changes [16].

Cellular metabolism is divided in two big blocks. The first is called *Catabolism*, whose processes are related with the degradation of nutrients and intermediate substrates to provide energy and basic building blocks coming from the rupture of chemical bonds of nutrients. The second is referred to as *Anabolism*, whose processes are related fundamentally to the synthesis of complex organic molecules. Notice that Catabolism supplies Anabolism with the necessary energy and basic compounds or elements to synthesize new molecules. At a different scale, biochemical reactions have been classically classified into different biochemical pathways, which are sequences of consecutive reactions that transform certain metabolites into specific products. Pathways are traditionally associated with definite functions, like Glycolysis which breaks down glucose into other small compounds to extract chemical energy and basic building blocks for anabolic reactions in the synthesis of fatty acids or amino acids. Currently, we know that pathways are not isolated entities and, instead, they constantly interact with each other [6, 7].

Focusing on individual reactions, one must notice that they require the action of catalysts -called enzymes- to take place. Enzymes are a special class of proteins. Proteins are macromolecules composed of amino acids, which perform a large number of functions in living cells, participating for example in the responses to stimuli, the replication of DNA, and transportation of molecules. It is worth stressing that, even though biochemical reactions may be thermodynamically spontaneous, they would not take place without enzymes because the activation energy required inside cells is very large. To ensure that all reactions occur, enzymes decrease the necessary activation energy by generating feasible chemical mechanisms that allow these reactions to take place in a controlled way and in reasonable amounts of time [17]. The action of enzymes helps also to control reaction fluxes, *i.e.*, the rates of biochemical reactions. Not all reactions in metabolism proceed with the same speed or are always on. Biochemical fluxes present a broad distribution of values [18] that reconfigure in response to internal or external changes and signals.

### 1.1.1   Key compounds

Biochemical reactions are connected by their participating chemical species, the products of one reaction are the substrates of subsequent reactions, and so on. These compounds -metabolites- participate in many different cell functions, including catalytic activity of their own. Five different general categories of metabolites are described in the following paragraphs (see Figure 1.1).

- Amino acids [19, 20] are compounds composed by amines (-$NR_2$), carboxylic groups (-COOH), and a different side chain for each amino acid. The polymerization of different amino acids generates short chains called

FIGURE 1.1:  Examples of classes of compounds that can be found in the metabolism of cells.

peptides, or long chains called polypeptides that can be arranged in one or more biological functional way to form proteins.

- Lipids are amphyphilic molecules, like fats or sterols, that contain both a polar and an apolar part. This implies that they can be in contact with water (polar part), whereas at the same time are soluble in substances like oil through its hydrophobic part. The main uses of lipids are to store energy [21], signaling [22], and being constituents of membranes [23].

- Carbohydrates are large biological molecules consisting of carbon (C), hydrogen (H), and oxygen (O) arranged on a carbon backbone possibly containing in addition aldehydes (-CHO), ketones (-CO-) and hydroxyl (-OH) groups. They fulfill many roles like energy source [24], storage of energy in the form of glycogen [25], or structural functions.

- Nucleotides are organic molecules containing a nitrogenated base (an aromatic compound containing a basic[2] nitrogen), a ribose or deoxyribose sugar, and a phosphate group ($-PO_4^{3-}$). They are building blocks of the two nucleic acids DNA and RNA [26]. Genes are fragments of DNA that contain the hereditary information in order to code for polypeptides or for RNA chains. At the same time, RNA performs multiple vital roles in the coding, decoding, regulation, and expression of genes. Nucleotides are obtained from the phosphorylation of nucleosides, and in the form of nucleoside triphosphates, nucleotides play central roles in metabolism [27]. One of these roles is to act as coenzymes, which are important metabolic intermediates that bond loosely to enzymes so as they can perform their catalytic activity. For instance, coenzymes serve to carry energy within the cell. An important coenzyme is adenosine triphosphate (ATP). It is one of the energy currencies of the cell [28]. Many reactions depend on ATP to become thermodynamically spontaneous, taking advantage of the large content of free energy that is released when the high-energy oxygen-phosphate bond of ATP is broken. Another example of the importance of nucleotides as coenzymes is nicotine adenine dinucleotide $NAD^+$, a derivative of vitamine $B_3$, along with its reduced form nicotine adenine dinucleotide - hydrogen (NADH), which are in charge of balancing the quantity of reduced / oxidized species inside the cell [29, 30].

- Inorganic compounds like water ($H_2O$), or ionic species like potassium ($K^+$), sodium ($K^+$), chlorine ($Cl^-$), calcium ($Ca^{2+}$), etc., are simple but not less important components of metabolism. Some of them are abundant, like sodium or potassium, whereas others are present at very low concentrations (traces) [31]. They appear in the form of electrolytes,

---

[2]In this context, basic refers to acid-base behavior.

and thus their concentrations play a key role for example in fixing the osmotic pressure, pH, or the cell membrane potential [32, 33]. Some transition heavy metals like iron ($Fe^{2+}$/$Fe^{3+}$) or zinc ($Zn^{2+}$) are cofactors, compounds which are essential for the activity of proteins like hemoglobin [34].

### 1.1.2  Biochemical reactions

Metabolites are the substrates or products of biochemical reactions in the cell. These can be classified in different categories. An important kind of metabolic reactions is a redox process, which involve the transfer of electrons from reduced species, like ammonia or hydrogen sulfide, to oxidized ones, like oxygen or nitrates. Redox reactions play fundamental roles in respiration, where glucose reacts with oxygen, the final products being carbon dioxide coming from the oxidation of glucose, and water, obtained by reduction of oxygen, along with a large quantity of free energy, which is mainly used for non-spontaneous anabolic processes.

Another type of reactions in metabolism involves the transference of entire chemical groups, like a phosphate group in a phosphorylation reaction. Other reactions involve the direct breakage of chemical bonds, like the rupture of carbon-carbon bonds in the decarboxylation of pyruvate. This is a principal process in fermentation which, in order to obtain energy and avoid pyruvate accumulation, transforms a carboxylic group in the form of carbon dioxide, generating acetaldehyde that finally gets reduced into ethanol by a redox reaction [35]. Decarboxylations are also important, for example, in the intermediate step between Glycolysis and the Citric Acid Cycle[3] to obtain acetyl-CoA, or in

---

[3]Also called Krebs Cycle or Tricarboxylic Acid Cycle (TCA Cycle).

subsequent steps of this last pathway to generate new intermediates. Transport reactions deserve special attention, since they are responsible for the entrance of nutrients and the excretion of waste products.

An important feature of biochemical reactions is reversibility. Depending on the value of $\Delta G^{o4}$, reactions can be considered as reversible or irreversible [36]. More precisely, for $\Delta G^o \approx 0$ reactions can be considered reversible, meaning that both directions of the reaction are thermodynamically favored; generically one would write $aA + bB \rightleftharpoons cC + dD$. On the contrary, if $\Delta G^o < 0$ and significantly negative, reactions are considered irreversible, and one direction is favored $aA + bB \rightarrow cC + dD$. For $\Delta G^o > 0$, the reaction takes place mostly in the opposite direction $aA + bB \leftarrow cC + dD$.

### 1.1.3   Biochemical pathways

Traditionally, sequences of consecutive biochemical reactions that transform a principal chemical into specific products are called *pathways*. In cell metabolism, there are several universal pathways that when interconnected form a complex metabolic network. Next, the central pathways of metabolism are briefly reviewed.

Glycolysis is the pathway that degrades carbohydrates. It takes place in the cytosol, and its main fuel is glucose. Basically, Glycolysis contains enzyme-catalyzed chemical reactions which transform glucose into pyruvate. In its more common form, this process generates the necessary free energy in order to form two molecules of ATP along with NADH. Glycolysis contains two phases, the first one where energy must be invested, which costs two ATP molecules but

---

[4]Thermodynamically speaking one should refer to $\Delta G$, the change in Gibbs free energy (SI units J mol$^{-1}$). An approximate but convenient way is however to refer to $\Delta G^o$, which denotes the free energy change in standard conditions of a reaction.

that generates important intermediate compounds. On the contrary, the second phase produces energy, since four ATP molecules are generated, along with two pyruvate molecules and two NADH molecules. Therefore, Glycolysis is important not only to obtain energy but to generate important biosynthetic precursors [16]. Notice that the inverse process, which generates glucose from pyruvate is called Gluconeogenesis and corresponds to Anabolism.

Pyruvate obtained from Glycolysis can be metabolized in two different ways. The first way corresponds to anaerobic processes, when no oxygen is available. This is called fermentation, and consists in reducing pyruvate into several components like ethanol, lactate or acetate by oxidizing NADH into $NAD^+$. Fermentation generates two ATP molecules [16].

In case that oxygen is present, the main fate of pyruvate is to become acetyl-CoA, a chemically activated compound formed by a cofactor, called coenzyme A, and an acetyl group. Acetyl-CoA enters the Citric Acid Cycle, a route that takes simple carbon compounds and transforms them into $CO_2$ in order to obtain energy. The Citric Acid Cycle not only accepts acetyl-CoA from Glycolysis, but also from other routes like lipid or protein metabolism, which emphasizes the importance and centrality of this pathway (see Figure 1.2) [16]. In eukaryotic cells, the Citric Acid Cycle occurs in the matrix of mitochondria, whereas in prokaryotic cells it takes place in the cytosol, like Glycolysis. The Citric Acid Cycle generates $CO_2$, guanosine-5'-triphosphate (GTP), NADH, and flavin adenine dinucleotide in hydroquinone form ($FADH_2$). GTP is transformed directly into ATP. NADH and $FADH_2$ are two reduced species that, by being oxidized, generate also ATP. This oxidation takes place in the process called Oxidative Phosphorylation.

FIGURE 1.2: Schematic representation of biochemical routes in central metabolism. Notice that the size of the arrows is only determined by aesthetics and does not contain any information about the magnitude of the fluxes through these routes. Thin arrows are not explained in the text but are added in this figure for completeness and to remark the interconnectivity present in cell metabolism. Blue circles denote major families of compounds. The Citric Acid Cycle (due to its cyclic form) is represented also with a circle. The other processes are represented by squares, orange color denoting pathways and green color denoting specific metabolites.

Organisms take advantage of the processes in the electron respiratory chain called Oxidative Phosphorylation in order to oxidize the reduced species coming from the Citric Acid Cycle to generate energy. In eukaryotic cells, Oxidative Phosphorylation takes place inside mitochondria. In prokaryotic organisms, where no mitochondria are present, it takes place across the prokaryotic cell membrane. To summarize, by coupling Glycolysis to the Citric Acid Cycle and Oxidative Phosphorylation, organisms are be able to generate up to 38 ATP molecules [16], which compared to two ATP molecules generated by fermentation, represents a great advantage in order to obtain ATP, whenever oxygen is present.

Other important pathways in cell metabolism comprise the degradation of fatty acids inside mitochondria, a process called $\beta$-oxidation, which is another source of acetyl-CoA apart from Glycolysis. Another source of acetyl-CoA comes from the degradation of amino acids, which can be synthesized by transamination [16]. Basically, transamination transforms $\alpha$-ketoacids coming from the Citric Acid Cycle to generate amino acids, which emphasizes again the centrality of the Citric Acid Cycle.

There are two main routes for the synthesis of purines and pyrimidines, the building blocks of nucleic acids or coenzymes like $NAD^+$: the *de novo*, which refers to the synthesis from simple molecules, and the salvage pathways, where purines and pyrimidines are recycled from intermediates coming from the routes that degrade nucleotides. The *de novo* route of nucleotide synthesis has a high energetic requirement as compared to the salvage pathway. The enzymes that synthesize purines and pyrimidines perform basic, cellular activities and it is thought that are present in low, constitutive levels in all cells [37].

### 1.1.4   Classical studies of metabolism

Traditionally, metabolism has been studied using a biochemical reductionist approach focused mainly on the study of the role of biomolecules and the kinetics and on the thermodynamics of particular metabolic reactions. As an example, processes like the non-spontaneous transport across the membrane -which takes advantage of the free energy coming from a proton gradient [38] or from ATP hydrolysis- have been studied using irreversible thermodynamics. Classical questions in biochemistry that prompt new systems-level studies refer to regulation and control of metabolism, the interplay and adaptation to the environment, and the effects of evolutionary pressure.

#### 1.1.4.1   Kinetics and thermodynamics

Classic metabolic studies have usually focused on the kinetics of reactions. The traditional approach was to discover the chemical mechanism by which reactions take place. In this way, kinetic constants were measured for specific reactions using *in vitro* experimental techniques in order to obtain a velocity law.

As mentioned before, the action of enzymes decreases the necessary activation energy of a reaction, so that the reaction rate increases (otherwise it would take place more slowly or even it would take place so slowly that any progression of the reaction would be unnoticeable). A scheme of this decrease in the energy barrier is shown in Figure 1.3. The best-known kinetic enzymatic mechanism in biochemistry is the famous Michaelis-Menten kinetics [39]. In fact, biochemical reactions involving a single substrate are often assumed to follow Michaelis–Menten kinetics. This model assumes that the minimal equation to describe a simple reaction with one reactant $S$ and one product $P$ catalyzed by

FIGURE 1.3: Energy diagram showing the dependence of the energy required for the reactants in order to be transformed into products as a function of the reaction coordinate. This is an abstract coordinate that represents the progress of a reaction along the complete path.

one enzyme $E$ is

$$S + E \underset{k_{-1}}{\overset{k_1}{\rightleftharpoons}} ES \overset{k_2}{\to} P + E \tag{1.1}$$

where $k_1$, $k_{-1}$, and $k_2$ are rate constants. The model relates the overall reaction rate $v$ to the concentration of substrate $[S]$ and the concentration of enzyme $[E]$ under assumptions like steady-state conditions and low enzyme concentration. The rate $v$ is given by the expression $v = v_{max} \frac{[S]}{K_m + [S]}$, where $v_{max} = k_2[E]$ and $K_m$ is the substrate concentration at which the reaction rate is at half-maximum. Michaelis–Menten kinetics reaches a saturation of the velocity as a function of the substrate concentration due to the limited availability of enzyme that can bind to the substrate.

Apart from Michaelis-Menten kinetics, other mechanisms were described for reactions involving more than one substrate or even for reactions with one substrate that do not follow the Michaelis-Menten mechanism. One of these examples is cooperation, which happens when the binding of one substrate

molecule to the enzyme affects the binding of subsequent substrate molecules. This effect is modeled by the Hill equation [40], which has the form $\theta = \frac{[L]^n}{K_a^n + [L]^n}$, where $\theta$ is the fraction of occupied sites and the Hill coefficient $n$ measures how much the binding of substrate to one active site affects the binding of substrate to the other active sites. The case $n < 1$ indicates that once one substrate molecule is bound to the enzyme, its affinity for other substrate molecules decreases, whereas $n > 1$ indicates that once one substrate molecule is bound to the enzyme, its affinity for other substrate molecules increases. The case $n = 1$ indicates that the binding of one substrate does not affect the binding of other ligands. The other parameters $[L]$ and $K_a$ are, respectively, the free unbound substrate concentration and the apparent dissociation constant derived from the law of mass action.

Other kinetic mechanisms, involving multi-substrate reactions, are the so-called ternary-complex mechanisms and ping–pong mechanisms [16]. These mechanisms describe the kinetics of an enzyme that takes two substrates, namely A and B, and turns them into two products, namely P and Q. Ternary-complex mechanisms imply that the substrates bind to the enzyme forming a ternary complex, where the reaction takes place. After this transformation, the complex dissociates, giving products P and Q. Ping–pong mechanisms consist on sequences of enzyme transformations due to interactions with the substrates. First, the enzyme binds to one substrate and one product is formed. After this process, the second substrate binds to the enzyme giving the second product.

Specific applications of thermodynamics to cell metabolism can be found for example in the description of transport of molecules across the cell membrane. On the one side, passive transport implies a movement of compounds which involves no energy supply, happening spontaneously. On the other side, active transport accounts for the movement of compounds across the cell membrane

in the direction against a concentration gradient. Active transport is usually associated to the accumulation of high concentrations of molecules that the cell needs, such as ions, glucose and amino acids. If this process uses chemical energy in the form of ATP, it is termed as primary active transport. Secondary active transport involves the use of an electrochemical gradient. Examples of active transport include uptake of glucose in the human intestines [41].

Kinetics describes the rates of reactions and how fast equilibrium is reached, but it gives no information about conditions once the reaction reaches equilibrium. At the systems level, several aspects must be taken into consideration in relation to its second law. In simple terms, the second law of thermodynamics states that in a closed system entropy tends to increase. An increase in the entropy of a system implies an increase of the number of its possible reachable states. However, organisms seem to contradict this law, since biological systems are complex but ordered structures. To obey the second law and, at the same time, to generate these structures, organisms must exchange matter and energy with their surroundings (see Figure 1.4). In this way, organisms are not in thermodynamic equilibrium, but they are dissipative systems which, to maintain their high degree of complexity and order, increase the entropy of their surroundings whereas their internal entropy is decreased. Thus, the necessary free energy required by Anabolism to generate complex molecules is obtained by coupling it to Catabolism. For example, nutrients are metabolized and small molecules like $CO_2$, whose entropy is much larger than that of nutrients [42, 43], are expelled as waste.

Another thermodynamic discussion concerns energy balance. The intake of energy is equal to the sum of the energy expended in the form of heat or work, and the stored energy. Energy balance states that no energy can be created or destroyed, but it can be transformed. This is indeed the first law of

FIGURE 1.4: Schematic example of an open system, with exchange of matter and energy, and a closed system, where there are no exchanges of any type.

thermodynamics. For example, when a cell consumes nutrients, a part of the energy content of the nutrients will be diverted towards the storage as fat, or transferred inside the cell as chemical energy in the form of ATP, or immediately dissipated as heat.

### 1.1.4.2    Regulation and control

The environment of organisms is constantly changing. In fact, organisms themselves modify their own surroundings by consuming nutrients and expelling waste. Therefore, organisms must be regulated in order to avoid large imbalances within themselves. Furthermore, possible internal perturbations can also lead to imbalances inside an organism. Hence, organisms have developed different regulation strategies to be able to maintain *homeostatic* states in which internal conditions remain stable [44]. Regulation requires that a system operates near steady-state conditions, which means that the temporal variation of the properties through time is practically null, except for adjustments to internal or external perturbations. This implies that concentrations of internal metabolites

are maintained steady in front of variations in metabolic fluxes. This entails the regulation of enzymes by increasing or decreasing their response to signals.

A real example of homeostatic readjustment is the regulation of glucose concentration by insulin [45, 46]. When large levels of glucose are present in the blood, insulin binds to its receptors, which generates a cascade of protein kinases[5] that cause the consumption of glucose into fatty acids or glycogen. Therefore, the increase in the concentration of glucose is regulated by the control of fluxes of catabolic biochemical reactions, so as to decrease the concentration of glucose until a stable steady-state is reached.

Control has been differentiated from regulation. Metabolic control refers to the ability to change a metabolic state as a response to an external signal [47]. In this way, control can be assessed in terms of the intensity of the response to the external factor without the need of knowing how the organism is able to achieve internally this state. This implies that control is simpler than regulation, because no judgment about the function of the system is needed. For example, an enzyme may show large changes in activity due to some external signal, but these changes may have little effect on the overall flux of a certain set of reactions or pathway. Therefore, this enzyme is not involved in its control.

### 1.1.4.3    Evolution

Through the process of descend with modifications, organisms evolve and change in time under the driving force of survival. In cell metabolism, there are central pathways that have been conserved through evolution and that are present

---

[5]A protein kinase is a kind of enzyme which transfers phosphate groups from high-energy phosphate donor molecules to specific substrates. This process is called phosphorylation, not to be confused with the Oxidative Phosphorylation pathway described in Section 1.1.3.

in practically all kinds of organisms. In fact, these pathways were in the so-called last universal ancestor, which is the most recent organism from which all organisms that now live on Earth descend [48]. Pathways like Glycolysis and the Citric Acid Cycle have been retained probably due to their optimality when producing their products and intermediates in a relatively small number of steps, which then can act as precursors for other biochemical routes. Many studies support the theory that organisms have evolved towards the maximization of the growth rate, *i.e.*, organisms tend to reproduce as much as possible [49, 50].

There have been proposals in recent years in order to understand how metabolism might have evolved including the retention of ancestral pathways. Different mechanisms have been proposed for the evolution of metabolic pathways, for instance (1) sequential addition of old or new enzymes within short ancestral pathways, (2) duplication and then divergence of pathways, and (3) recruitment of enzymes that are already present to be assembled into a novel pathway [51]. Horizontal gene transfer is another way that organisms use to evolve, consisting on the transfer of genes between organisms. In fact, bacteria acquire resistance to antibiotics due to horizontal gene transfer [52]. This process implies modifications in the metabolic network, in the form of alterations of pathways, to generate by-passes in order to avoid the effect of the antibiotic.

Heritable epigenetic effects have also impact on evolution. Epigenetics studies the changes in gene expression that cannot be explained by changes in DNA sequences. There are two ways in which epigenetic inheritance may be different from traditional genetic inheritance. The first way corresponds to the situation where the rates of epimutation are much faster than the rates of mutation [53]. Alternatively, epimutations are more easily reversible [54]. The existence of these possibilities implies that epigenetics, and thus metabolic effects, can increase the evolvability of species.

Evolution can cause not only the gain of new metabolic functions but also the loss of functions which are not useful anymore for cells. *Mycoplasmas*, a kind of bacteria without cellular wall that act as parasites, have lost those processes and pathways that are essential for survival as independent entities, since these microorganisms obtain compounds from their hosts [55, 56].

## 1.2    Genome-scale models

A systems-level approach to the study of cell metabolism takes into account the entire set of biochemical reactions and their interactions at different levels of organization. At the core of this approach, *genome-scale metabolic networks* [10] provide high quality representations of cell metabolism which integrate biochemical information with genome annotations, physiological requirements, and constraint-based modeling refinements. These genome-scale models, after experimental validation, have predictive capacity and can be used for detailed analysis of metabolic capabilities, with applications in a range of fields like biomedicine or biotechnology [57, 58].

### 1.2.1    Reconstructing metabolism

Nowadays, genome-scale metabolic models have been reconstructed and experimentally validated for different organisms like *Escherichia coli*, *Saccharomyces cerevisiae*, *Mycoplasma pneumoniae*, and *Homo sapiens*, among others. These reconstructions are called *GENome-scale metabolic REconstructions* (GENREs) (see Figure 1.5) [10]. In GENREs, reactions are typically stoichiometrically

FIGURE 1.5: Simplified representation of a genome-scale model. Reactions are catalyzed by enzymes, whereas enzymes are codified by genes. Reactions are represented by blue squares, metabolites by green circles, enzymes by red rhombus and genes by yellow triangles. Note that enzymes 8 and 9 form a complex and the latter is the catalyst of reaction *j*.

balanced and categorized into their corresponding pathway, for example, reactions belonging to Glycolysis, Oxidative Phosphorylation, or Citric Acid Cycle. Reactions are also associated to their corresponding enzyme and metabolic gene.

Generating these representations is a difficult task and several steps are needed in the protocol [59]. First, an initial reconstruction is proposed from gene-annotation data coupled with biochemical information from databases like the Kyoto Encyclopedia of Genes and Genomes (KEGG) [60], or BioCyc [61], among others. In these databases, reactions are linked with metabolic genes, enzymes, and also to functional categories like pathways. Second, the obtained reconstruction is curated by checking it against experimental evidence in the

existing literature, including for instance physiological requirements. This revised reconstruction is further translated into a computational mathematical model using constraint-based approaches. Third, the reconstruction is validated by comparing the results obtained by the model with experimental evidence. After curation of inconsistencies, models are , see for instance the BiGG database [62]. Finally, one has to remember that GENREs are constantly improved in new versions as new experimental results become available.

Among all metabolic network simulation techniques for model refinement, Flux Balance Analysis (FBA) [63] is probably the most widespread. Very briefly, FBA uses constraint-based analysis to compute a metabolic phenotype, in the form of the set of fluxes of reactions, which maximizes biomass production given a set of external bounds typically referring to nutrient amounts.

Since the first GENRE reconstructed a decade ago [64], there has been a huge expansion on the construction and use of GENREs [65–69]. Their applications can be divided into four categories [57, 70, 71].

- Many advances in Biology are the result of *hypothesis-driven discoveries*. Metabolic GENREs enable the identification and confirmation of new or existing hypotheses, representing an important framework for the incorporation of cell biological data. The key to unlock the potential of GENREs for the discovery of unknown metabolic mechanisms is to ask feasible questions and to know the limitations of the used methodology, since one must always have in mind that in real living cells, many biological levels act together (metabolism, regulation, signaling, gene regulation, etc.) creating a complex system, and GENREs are after all simplified models [57].

- Many characteristic phenotypes of several organisms arise when they *interact with other species* [72, 73]. GENREs enable to analyze interactions between organisms, like for example mutualism, comensalism, parasitism, etc [73]. It is worth mentioning in this respect the work of Bordbar *et al.* [74], where the authors developed a model of parasitism between a human cell and the bacterium that causes tuberculosis, *Mycobacterium tuberculosis.*

- Metabolic reconstructions serve as a framework for the contextualization of data obtained using *high-throughput techniques* [75]. A functional way to apply GENREs for contextualization of experimental data, like gene expression or C13 flux data, is by imposing constraints on the fluxes of GENRE based on experimental values. If experiments suggest for instance that reactions of a particular pathway carry large fluxes, one can force the GENRE to have a minimal bound for these fluxes so as to fit the experimental observations. Then, changes in the global flux structure are studied and evaluated.

- *Metabolic engineering* involves the use of recombinant DNA technology[6] to selectively alter metabolism and improve a targeted cellular function [57, 76]. The use of GENREs for metabolic engineering has led to what has been termed as *Systems Metabolic Engineering* [77]. An example of the new advances in metabolic engineering achieved using GENREs is the modification of *Saccharomyces cerevisiae* to increase the production of industrially important intermediates of the Citric Acid Cycle [78]. Another possibility is to study gene knockouts. More precisely, in Reference [79],

---

[6]Recombinant DNA molecules are DNA molecules engineered to assemble genetic material from multiple sources, creating sequences that would not otherwise be found in biological organisms.

the authors performed gene knockouts in *Geobacter sulfurreducens* to maximally increase its respiration rate.

### 1.2.2   The systems-level approach

GENRE reconstructions and the systems-level approach have led to the development of the field called Systems Biology. It is an emerging interdisciplinary field applied to biological systems that focuses on complex interactions using a holistic approach [80]. It is not easy to have a precise and unique definition encompassing all the concepts underlying Systems Biology.

A possible definition was stated by Ideker *et al.* [81]:

> "Systems biology studies biological systems by systematically perturbing them (biologically, genetically, or chemically); monitoring the gene, protein, and informational pathway responses; integrating these data; and ultimately, formulating mathematical models that describe the structure of the system and its response to individual perturbations."

An alternative was given by Kitano *et al.* [9]:

> "To understand complex biological systems requires the integration of experimental and computational research — in other words a systems biology approach."

These definitions share common features. On the one side, a systems-level approach considers all the components and linkages constituting the system. On the other side, the properties of the components and interactions must be

integrated in a computational mathematical model. It is worth stressing the importance of the assembly of these components, *i.e.*, how components interact between them. This can be understood with the analogy of a road-map as given in Reference [8]. In order to understand traffic patterns, it is necessary to know not only the static road-map but also how cars interact to generate the observed final traffic patterns. Thus, to fully understand a system in a systems-level approach, one needs the diagram with all the connections of all components but also the knowledge of why, how, and to which extent components interact.

Systems Biology can therefore be defined as an approach whose aim is to study biological systems focusing on all the constituents and interactions. In this way, emergent properties which are not present at the level of the single components of the system can be discovered, and phenotype and behavior can be related to the underlying systems architecture. Central to Systems Biology is the holistic approach. *Holism* is based on the idea that natural systems and their properties should be viewed as a whole instead on focusing on the parts that constitute the system (see Figure 1.6). Contrarily, the focusing on single parts is called *reductionism.* Examples of traditional reductionist approaches are the study of a single protein or a single chemical mechanism, and they have dominated Biochemistry [82] and Molecular Biology [83] for decades.

Systems Biology represents a paradigm shift that requires the interplay between different disciplines, *e.g.*, Biology, Physics, Mathematics, Chemistry, Computer Science, etc. [84, 85]. Systems Biology foments interactions from traditional computational scientists, modeling experts, and experimental researchers. Research developed to date typically requires powerful computational tools, and this particular emphasis in Systems Biology has given rise to the subfield known as Computational Systems Biology or Computational Biology [84].

# REDUCTIONISM

| Molecular Biology Biochemistry | Metabolites Enzymes / reactions |
|---|---|
| Metabolic Control Analysis | Pathways |
| Systems Biology | Complete network |

# HOLISM

FIGURE 1.6: Scheme of the different levels in the holist vs the reductionist approaches.

Systems Biology has grown in parallel to the development of the *omics* fields. Omics are different disciplines integrating and analyzing different kinds of data. Systems Biology combines the datasets obtained in these disciplines in order to achieve the maximum knowledge to model an organism (see 1.7). Examples of omics related to genes are *Genomics*, which involves sequencing an organism genome, and *Transcriptomics*, which evaluates gene transcription. In relation with proteins, the field called *Proteomics* measures protein abundance. Regarding metabolism, *Metabolomics* deals with the study of the concentration of all the compounds present in a organism. There is another important omic field, in line with this thesis, which studies the chemical fluxes in metabolism, *Fluxomics* [86, 87]. Fluxomics provides a measure of a metabolic phenotype as the set of fluxes going through all reactions in a metabolic network.

FIGURE 1.7: Interplay of the different layers involved in the final phenotype of an organism.

The holistic view of metabolism including reactions, metabolites, enzymes, genes, and fluxes, represents a new paradigm that requires new tools. Complex Network Science [14, 15] has become a new promising domain for the study of biological systems. Metabolism is formed by a large amount of components and interactions and can be categorized as a complex network and, thus, many applicable techniques that belong to the complex network field are appropriate for the the study of metabolism.

In fact, some of the applications of Complex Network Science ideas to cell metabolism have led to the discovery of many unenvisaged properties such as the existence of loops [88], optimal pathway usage [89], and metabolite connectivity [90]. Other possible discoveries are the exploration of evolutionary relationships [91]. In addition, complex networks applied to metabolism serve as a tool for the identification of how evolutionary pressure has shaped the

topological features of metabolic networks, such as the degree distribution [7, 92–94]. Therefore, the joint use of complex network methodologies and Systems Biology provides an excellent arena to study metabolic capabilities and the evolutionary forces that shape metabolic networks.

## 1.3    Aims and objectives

This thesis aims at studying cell metabolism from a systems-level perspective, *i.e.*, taking metabolism as a whole.

In particular, one of the questions is how metabolism responds as a whole when some of its constituents fail, *i.e.*, when reactions or genes are non-operative by removal or mutation. It is important to mention that the aim is not to focus on the study of how to perform biochemically the perturbation or the analysis of biochemical failures at a molecular level. Instead, the investigation focus on the impact on the whole system of harmful situations and how metabolism is able to overcome them as a whole entity. In this way, one can study how different pathways reorganize to adapt to perturbations, something impossible to understand by typical molecular biology studies centered on single constituents.

Another question addressed in this thesis is focused on filtering metabolic networks in order to extract metabolic backbones providing valuable biological information. To do this, FBA and the disparity filter [95] are used. The disparity filter allows to extract backbones of the metabolic network containing the significant links. The analysis of metabolic backbones allows the identification of pathways with important roles in survival. The first role corresponds to pathways that have been present in organisms since the first stages of life, *i.e.*, pathways central in long-term evolution. The second role corresponds to

pathways more sensitive to external stimuli, *i.e.*, pathways displaying short-term adaptation.

The last question addressed in this thesis is the assessment of FBA solutions in relation to all the feasible flux space, so as to identify whether solutions obtained with this technique describe reliably the set of possible metabolic states or, on the contrary, the FBA solution is uninformative of the entire set of metabolic phenotypes. The space of metabolic flux states can be exploited with different strategies. It can be used as a benchmark to calibrate the distance of FBA fluxes as compared to experimental measures, or to identify metabolic phenotypes unreachable by constraint-based techniques.

The main objectives of this thesis are summarized in the following bullet list:

- To study whether the structure of metabolic networks has evolved towards robustness resisting external perturbations, like gene or reaction removals or mutations.

    - To study the spreading of a cascade when a reaction or a pair of reactions fail, unveiling the interplay between multiple cascades.

    - To study the propagation of the damage to metabolism when genes fail.

    - To discuss the findings in terms of an evolutionary perspective.

- To study the effects on fluxes of individual and pairs of reactions knockouts using FBA.

    - To study the activity and essentiality of reactions.

    - To understand the mechanisms of synthetic lethality, unveiling the plasticity and the redundancy capabilities displayed by the metabolic networks of bacteria.

- – To study the dependence of plasticity and redundancy on the environment.

- To identify those pathways that perform important roles for the survival of an organism.

  - – To check the efficiency of the disparity filter on metabolic networks.

  - – To analyze backbones in terms of the long-term evolution of organisms

  - – To extract information about the short-term adaptation of metabolism to the external environment.

- To assess the FBA solution in the entire space of metabolic solutions.

  - – To demonstrate that solutions obtained using FBA as a constraint-based technique may be uninformative of typical behaviors.

  - – To provide a benchmark to calibrate FBA.

  - – To recover phenotypes not attainable by constraint-based techniques by using the full metabolic solution map.

## 1.4   Outline

After this introduction to cellular metabolism and its genome-scale models, Chapter 2 presents the general tools, methodologies, and GENREs used in this thesis.

Chapter 3 starts by considering a structural study of how metabolic networks of the bacteria *Mycoplasma pneumoniae*, *Escherichia coli*, and *Staphylococcus aureus* respond to internal perturbations, like removals of reactions or genes individually or in pairs, *i.e.*, how the structure of the metabolic network is damaged following an internal failure which propagates as a cascade, by which

the metabolic capabilities of an organism are weakened. Further, these results are linked to evolutionary explanations, *i.e.*, how evolution has shaped and dictated the form of metabolic networks so as to respond to perturbations. This discussion is related with the robustness of organisms, in order to unveil whether the structure of the metabolic network is prepared to suffocate the advance of a damage cascade.

Chapter 4 extends the structural study of perturbations to flux distributions obtained using FBA. This study allows, on the one side, to know whether there are important reactions that must be always active in order to guarantee the survival of an organism and, on the other side, to check whether cell metabolism has developed protection mechanisms when some of its parts are unable to work. In this respect, synthetic lethal reaction pairs are analyzed. These are pairs of reactions whose removal from is lethal, but metabolism is still able to survive when each reaction forming the pair is removed individually. This allows to identify two different mechanisms, plasticity and redundancy, which have helped to protect metabolism against possible reaction failures.

Chapter 5 analyzes metabolic fluxes so as to extract more biological information on how organisms adapt to external environments and evolve. To perform this analysis, the disparity filter is used in order to obtain backbones as reduced versions of metabolism without losing its properties as a complex network. The structure of these backbones unveils pathways with a prominent role in the long-term evolution of the organisms and in their short-term adaptation to the environment.

Chapter 6 revises the FBA technique in relation to the whole set of feasible flux states in a metabolic network. FBA uses a strong assumption -organisms try to grow as much as possible- allowing to solve the mass action equations

at steady state describing metabolism without the need of kinetic parameters. This assumption is commonly applied due to the lack of availability of kinetic constants of reactions. It is worth exploring the distribution of possible fluxes without making use of the assumption of maximal growth. This allows to perform a mapping of all the feasible flux solutions in metabolism and thus to assess the relevance of the solution obtained by FBA compared to all the other possible solutions.

General conclusions are given in Chapter 7. At the end of the thesis, there are four appendixes reviewing the basics of some specific tools used in Chapters 3, 4, 5, and 6. Finally, the list of references is included. A CD containing supplementary tables and data is also provided with the book. In addition to the CD, the supplementary files corresponding to each chapter can be found in the following links:

- Chapter 3 file (Supplementary Tables C3): http://tinyurl.com/nc2vmhf

- Chapter 4 file (Supplementary Tables C4): http://tinyurl.com/ktswl4k

- Chapter 5 file (Supplementary Tables C5): http://tinyurl.com/m7nh4b6

- Chapter 6 file (Supplementary Tables C6): http://tinyurl.com/k2bqtug

# Chapter 2

# Methods and data

This chapter describes the basics of the fundamental techniques used in this thesis. It is divided in three parts: (1) complex network tools applied to metabolism, (2) description of Flux Balance Analysis (FBA) -used to compute metabolic fluxes at steady state- and of Flux Variability Analysis -a variant of FBA to bound minimum and maximum fluxes for each reaction- and (3) a description of all the genome-scale metabolic reconstructions analyzed in this thesis.

Nowadays, the explosion in computational power has allowed us to deal with systems of thousands or even millions of constituents and interactions, boosting the degree of our understanding on how these systems are structured and behave. Complex Network Science comprises a large amount of techniques and models which help us to study these intricate systems as a whole [15, 96]. These methodologies can be applied to any system which can be modeled as a network. Networks can be briefly defined as a set of items that interact, like for example the World Wide Web and the Internet and, in a biological

context, metabolic networks [97] or protein-protein interaction networks [98]. Complex Network Science has led to an important advance in the understanding of metabolic networks [5, 97, 99–102] which is in line with the systems-level view of metabolism in fields like Systems Biology [10, 11].

When dealing with metabolic networks, the complex network approach has to be combined with other techniques coming from Systems Biology in order to understand functional or behavioral features, for example why the inability to operate of some reactions leads to cell death, or why some reactions carry a determinate flux given a set of external nutrients. The most widespread mathematical approach used for the systems-level analysis of metabolic networks is Flux Balance Analysis (FBA) [63]. This technique is based on constraint-based analysis and optimization of an objective function, usually the biomass formation function of the cell. In this way, the fluxes through all the biochemical reactions of cell metabolism that maximize the biomass formation rate or, equivalently, the specific growth rate, can be computed. Apart from the mentioned reaction fluxes and growth rate, this technique allows to compute, for instance, the maximum yield of important compounds such as ATP or NADH [63, 103], and the effects of knockouts of genes or reactions [104, 105]. Related to FBA, other related techniques like Flux Variability Analysis (FVA) [106, 107] allow to identify possible alternate solutions and, in conjunction with FBA, allow to perform a deep study of the flux capabilities of metabolic networks.

Complex network methodologies and constraint-based techniques applied to metabolic reconstructions represent a powerful tool for the analysis and development of new insights into metabolic functions and mechanisms that cells have developed from the earliest stages of life to the current days.

# 2.1 Structural properties of metabolic networks as complex networks

Networks are discrete systems of elements that interact. These systems are represented by graphs of *nodes* (or vertices) -which represent elements- connected by *links* (or edges) -which represent interactions. The presence of a large number of nodes interacting in non-trivial connectivity patterns between order and disorder is what gives to networks their intrinsic complexity.

It is important to distinguish between complex and complicated. The main difference between these two words is better explained by a single example: solving a whole metabolic network composed of thousands of reactions is a complex problem in the sense that the large amount of interactions leads to emerging unexpected behaviors, like the effect of the removal of a biochemical reaction on other reactions, which can increase or decrease their fluxes depending on their biological activity. On the contrary, the study of a typical chemical engineering process to obtain a precise output may be a complicated problem, since one needs to draw a flowchart of all the chemical reactions and involved intermediate species that participate in the chemical synthesis. This may require a wide knowledge of the system, implying a large degree of control on all the processes, but the final behavior of the system will be what is expected in a well-designed process.

## 2.1.1 Basic representation frameworks

Links in networks can have either a defined direction or may lack it. Therefore, when links are directed, they are depicted by arrows, specifying a source and a target. A directed link can represent, for example, a transformation between

two metabolites, typically a reactant and a product with the link pointing to the product. When no specification source / target is prescribed, the interaction is mutual, like in a protein-protein interaction[1], and links without direction are used. Associated to this, networks are classified as *directed*, *undirected*, or *semidirected*. It is worth mentioning that links can also be *bidirectional*, meaning that the interaction allows either the forward or backward direction at the same time and this interconnection is thus reciprocal. This is specially important in the context of metabolic networks, where reactions can be either reversible (bidirected links, meaning that both directions of the reaction are possible) or irreversible (directed links, meaning that only one direction is thermodynamically favored). Moreover, a link can carry a weight, representing the intensity of the interaction. Therefore, networks can be *weighted* or *unweighted*. In the case of metabolic networks, weights usually correspond to fluxes of the biochemical reactions. Metabolic networks typically display a probability distribution of fluxes (or weights) that follows a power law, meaning that fluxes spanning different orders of magnitude coexist in the same metabolic state [12].

Mathematically, unweighted undirected networks are described by the adjacency matrix, a square symmetric matrix $\{a_{ij}\}$ of binary values with an entry of 1 whenever there is a link between nodes $i$ an $j$ and 0 otherwise. In directed networks, the matrix is instead non-symmetric. Weighted networks are encoded by the weighted adjacency matrix $\{\omega_{ij}\}$, in which the values correspond to the weight of the edge between nodes $i$ and $j$.

Furthermore, networks can have different classes of nodes, leading to the so-called multipartite graphs. In multipartite graphs, links happen only between nodes in different categories. Networks with one kind of node are called *unipartite*, whereas

---

[1]Protein–protein interactions refer to physical contacts established between two or more proteins as a result of biochemical events and/or electrostatic forces.

FIGURE 2.1: Examples of different types of networks. a) Undirected unipartite. b) Directed unipartite. c) Undirected bipartite. d) Semidirected bipartite network. Notice that connections involving node *e* are bidirectional. e) Semidirected weighted bipartite. The thickness of the links is proportional to their weight. f) Example of the transformation into a one-mode projected network of metabolites from a semidirected bipartite metabolic network containing metabolites and reactions. Metabolites are represented by circles and reactions by squares.

networks with two kinds of nodes are called *bipartite* [108]. An important thing to notice is that bipartite networks can be projected into unipartite networks by performing a *one-mode projection*. To do this, one chooses a particular type of node and, in the projected reduction, places a link between two such nodes if there is at least one node of the complementary type connected to both of them.

In the real world, one can find networks combining all the mentioned properties (see Figure 2.1). Metabolic networks are usually represented as bipartite semidirected networks, with metabolites and reactions belonging to different node

categories with no direct connections between any two metabolites or any two reactions [109, 110] (see Figure 2.1d). Although a bipartite representation is more accurate, it is sometimes preferable and always simpler to work with one-mode projections based on metabolites, which can be either directed or undirected (see Figure 2.1a and b) depending on the reversibility of reactions, and weighted or unweighted depending on whether fluxes are taken into account. In such a projection, two metabolites get directly connected if there is at last one reaction in which they both participate (see Figure 2.1f).

### 2.1.2    Degree distribution

Nodes in networks are locally characterized by the number of their surrounding neighbors. This magnitude is called the *degree* of a node $k$ (see Figure 2.2). The probability of nodes having a certain degree $k$ is written $P(k)$ and named *degree distribution*, and can be computed from the fraction of nodes in the network that has degree $k$.

Usually, real world networks show degree distributions $P(k)$ that are highly skewed with long tails that reach values far above the mean [111]. In most cases, degree distributions follow a power-law, $P(k) \propto k^{-\gamma}$, where $\gamma$ is the characteristic exponent and it has values in the range $2 < \gamma < 3$. Networks with a degree distribution described by a power-law are called *scale-free*[2]. Networks with power-law degree distributions have attracted much attention and have been studied intensively [112–114]. Notice that, usually, it is useful to work with the complementary cumulative probability distribution function $P(k' \geq k)$ in order to avoid noise effects present for large values of $k$.

---

[2]This name is refers to the scale-invariance that power-laws display: if $f(x) = a(x)^{\gamma}$, then $f(cx) = a(cx)^{\gamma} = c^{\gamma} f(x)$.

FIGURE 2.2: Schematic example of a degree of a node (left) and a path between two nodes (right). Left: example of the degrees in a undirected, semidirected, and directed networks. Right: path between node $a$ and $b$, highlighted in green. In this case, the shortest path length between nodes $a$ and $b$ is $\ell_{ab} = 5$.

In semidirected networks, the degrees of nodes are defined in relation to incoming $(k_{in})$, outgoing $(k_{out})$ and bidirectional $(k_b)$ links. Correspondingly, nodes have a total degree expressed as a sum of contributions $k = k_{in} + k_{out} + k_b$. These degrees can present local correlations and so the degrees of nodes are described by the joint probability $P(k_{in}, k_{out}, k_b)$. In addition, for bipartite networks, nodes of each kind have also their own degree distribution.

Regarding specifically metabolic networks, the total degree of metabolites $k_M$ in bipartite representations follows a power-law degree distribution $P(k_M) \propto k_M^{-\gamma}$ [5, 96]. In Reference [97] it is found that in the organism *Escherichia coli*, the probability $P(k_{in})$ that a metabolite participates as a product and the probability $P(k_{out})$ that a metabolite participates as a reactant have both a value of $\gamma$ of 2.2. Similarly, Reference [115] shows that for the organism *Helicobacter pylori*, the exponent has a value of 2.32. The fact that metabolites display a scale-free degree distribution means that there is a high diversity in the number of reactions in which metabolites participate. The largest part of metabolites have a few connections, whereas a few metabolites, generically called hubs, have many of them. Examples of these highly-connected metabolites are ATP, $H_2O$,

FIGURE 2.3: Features of the networks of *Escherichia coli i*JO1366 (see Section 2.3.1), *Mycoplasma pneumoniae i*JW145 (see Section 2.3.2), and *Staphylococcus aureus i*SB619 (see Section 2.3.3). a) Complementary cumulative probability distribution function of metabolites. b) Degree distribution of reactions.

or $H^+$, which can participate in up to 50% of the total number of reactions for the case of $H^+$ in the organism *Escherichia coli* [116]. On the contrary, reactions show a peaked distribution of total degree, the peak being located at an average degree $< k_R > \sim 4$. The bounded form of the distribution arises from the fact that reactions have a limited number of participants, typically from 2 to 12.

In Figure 2.3a, the bipartite cumulative probability distribution function $P(k'_M \geq k_M)$ of metabolites and the bipartite probability distribution function $P(k_R)$ for reactions of the three organisms analyzed in this thesis, *Escherichia coli* [117–119], *Mycoplasma pneumoniae* [56, 120], and *Staphylococcus aureus* [66] are shown. Clearly, metabolites show a power-law degree distribution and reactions a peaked distribution, as mentioned above. In fact, all networks studied here have similar tendencies for both distribution functions, showing that metabolic networks, in spite of corresponding to quite different microorganisms, display often universal properties [97].

### 2.1.3    Average path length

Another common feature of complex networks, and in particular of metabolic networks, is the fact that any two nodes are connected by *paths* of links that are typically very short in the number of intermediate steps [5]. This is called the *small-world* property. In technical terms, the distance $\ell$ between two nodes is defined as the number of jumps or hops along the shortest path that connects them (see Figure 2.2). Hence, it is possible to define the average shortest path length $< \ell >$, which is the average of all the shortest distances between pairs of nodes. The small-world property is stated in the fact that $< \ell >$ increases as the logarithm of the network size $N$ (number of nodes) [111, 113].

Small average path lengths indicate that the network contains highly-connected nodes that act as shortcuts, reducing the average number of steps needed to go from one node to another. This is crucial in many real contexts, and in particular for cell metabolism. In Reference [97], the authors measured the average path length for 43 organisms and found a similar value for all of them, $< \ell > \sim 3.2$. This value was explained by the role of hubs, which decrease dramatically the number of steps needed to travel from one node to another. When hubs are not taken into account, longer and variable path lengths are obtained [84, 121], depending on the biological domain where organisms belong to. Typical values are 9.57, 8.50, and 7.22 for eukaryotes, archaea, and bacteria, respectively, with the differences due to evolutionary processes. Nevertheless, there remains some controversy about the small-world property in metabolic networks. In Reference [122], it is stated that usually paths are computed by directly linking metabolites through reactions and that this is not adequate, since pathways computed in this way do not conserve their structural moieties[3]

---

[3]According to the IUPAC, a moiety is a part of a molecule that may include either whole functional groups or parts of functional groups as substructures.

and thus they do not correspond to pathways on a traditional metabolic map. Therefore, in Reference [122] metabolites are linked depending on the conserved structural moieties in the adjacent reactions and, as a result, it is stated that the average path length of *Escherichia coli* metabolism is longer than it was previously thought and, consequently, the *Escherichia coli* metabolic network is not small in terms of biosynthesis and degradation of metabolites. However, it is generally accepted that metabolic networks show indeed the small-world property at the structural level. In this thesis, path lengths will be computed in Chapter 4.

### 2.1.4    Communities at the mesoscale

It is thought that biological networks are composed by subsets of nodes that are functionally separable called *modules* [113, 123]. In general, this idea corresponds to the concept of communities in networks. The organization of a network into communities does not imply fragmentation. Instead, communities are subsets of a network which contain a dense interconnection pattern between nodes inside the community and lower interconnection levels with nodes outside. This can be related with the presence of a large clustering (see Section 2.1.6) between nodes inside the community.

Community detection [124] represents an active field in Complex Network Science motivated by the potential identification of communities with functional or operational units. Several methods, based on different exploratory techniques, have been proposed. Among the most successful community detection methods one finds, for instance, algorithms that use random walkers to partition the network into communities, like Infomap [125]. Other methods are based on the optimization of modularity. Modularity is a measure of the quality of a

community structure [126]. It measures the internal connectivity of identified communities with reference to a randomized null model with the same degree distribution. Algorithms based on modularity optimization try to find the best community structure in terms of the modularity measure. Examples of successful algorithms based on this measure are SpinGlass [127] or Louvain [128]. On what follows, the three methods used in Chapter 3 of this thesis to detect communities are explained.

- *Distance hierarchical clustering*: this method starts by defining a distance between pairs of nodes in the network. Then, once the pairs of nodes have a defined distance, one groups similar nodes into communities according to this distance. There are different schemes based on distances to group nodes intro communities. The two simplest methods are single-linkage clustering, in which two sets of nodes are considered separate communities if and only if all pairs of nodes in the different sets have distance larger than a given threshold, and complete linkage clustering, in which all nodes of a community have a distance smaller than a threshold [129] (see Figure 2.4).

- *Infomap algorithm* [125]: the main idea of this algorithm is that a random walker will tend to flow at different paces within a network, spending more time inside communities and less time to pass between them (see Figure 2.4). The way in which the random walker moves around communities can be compared to the flow of messages between individuals. In this way, there is a strong current of messages between individuals inside a community, and a weaker current of messages between individuals of different communities.

- *Recursive percolation*: this method has been developed in a work related to this thesis [94]. Recursive percolation identifies components in which the network is fragmented just below the percolation threshold (see

**a   Distance hierarchical clustering**



**b   Infomap algorithm**



**c   Recursive percolation**



FIGURE 2.4: Examples of the clustering methods. a) Example of the distance hierarchical clustering method. Modules are formed by nodes that are nearer. Notice that with this method it is necessary to apply a threshold depending on the distances. In this example, the threshold is represented by the green rectangle. At this level, three communities are detected. b) Example of the Infomap algorithm. Clusters are found with a random walker. Communities are found depending on the frequency of times that each random walker visits a set of nodes. c) Example of the application of Recursive percolation. The first step leads to 10 clusters. Among these 10 clusters, the largest are fragmented, leading to more clusters. This partition is iterated until the distribution of sizes is similar to that in other methods.

Section 2.1.5), where the connected network disaggregates into smaller components. To find them, links are removed sequentially from lower to higher weights until the percolation transition is detected. Then, clusters are identified using a burning algorithm [130]. This procedure is applied to each component until the distribution of sizes of the obtained communities reaches some thresholds, for instance, to be similar to those given by the distance hierarchical clustering technique and Infomap. A schematic example of this process is shown in Figure 2.4.

### 2.1.5   Large-scale connected components

Global connectedness is one of the most fundamental properties of complex systems. The theory that describes the behavior of network connected components is *percolation theory.* Briefly, percolation theory states that there exists a critical point, called *percolation threshold* denoted as $p_c$, where a transition in the global connectedness of the network occurs, from a state where the network is formed by small isolated components to the emergence of a *giant connected component* (GCC) spanning a macroscopic fraction of the network. This means that it is always possible to find a path connecting every pairs of nodes inside the GCC.

This concept can be extended to networks with directed links. The connectivity of directed networks presents special features since the path between two nodes $i$ and $j$ can be different when going from $i$ to $j$ or vice versa. This fact leads to the existence of a bow-tie structure inside the GCC [110, 131, 132]. The main feature of the bow-tie structure of a GCC in a directed network is that one can detect the presence of a *strongly connected component* (SCC), which is a region of the network where any node is reachable from any other by a directed path. It can happen that directed networks contain more than one SCC.

Apart from the SCC, one of the other significant regions that can be found in the bow-tie structure of directed networks is called *IN component*, with nodes that can reach the SCC but that cannot be reached from the SCC. Analogously, the *OUT component* contains nodes that can be reached from the SCC but that cannot return to it. *Tubes* are sequences of nodes that connect the IN with the OUT component without going through the SCC. Finally, *tendrils* are composed by nodes that have no access to the SCC and that are not reachable from it, similarly to tubes. They go out from the IN component and come in from the OUT component. A visual scheme of the bow-tie structure of directed networks is shown in Figure 2.5a. The bow-tie structure of *Escherichia coli* and *Mycoplasma pneumoniae* will be explicitly considered in Chapter 5.

Metabolic networks show a bow-tie structure typically with a large SCC connected to non-structured IN and OUT components (see Figure 2.5b) [131, 133]. The SCC contains the largest part of metabolites and reactions composing the network, representing thus the entire metabolic machinery of cells. IN and OUT components are formed of, respectively, nutrients and waste products directly connected to the SCC (see Figure 2.5b).

### 2.1.6   Other structural properties of complex networks

Real networks exhibit also the presence of non trivial correlations in their connectivity. At the level of two nodes, it is convenient to characterize degree correlations with the average nearest neighbor degree $\bar{k}_{nn}(k) = \sum_{k'} k' P(k'|k)$, where $P(k'|k)$ is the probability of having a node with degree $k'$ given that it is connected to a node with degree $k$. It basically considers the mean degree of the neighbors of a node as a function of its degree $k$. If $\bar{k}_{nn}(k)$ increases with $k$, it is said that the network is *assortative*, with nodes that connect preferentially to

FIGURE 2.5: Examples of connected components. (a) Schematic example of a bow-tie structure. (b) Example of the bow-tie structure of *Mycoplasma pneumoniae* [120], an organism studied in this thesis. Blue nodes compose the SCC, red nodes compose the IN component, and green nodes compose the OUT component.

other nodes of similar degree. If $\bar{k}_{nn}(k)$ decreases with $k$, the network is named *disassortative*, with high-degree nodes attached preferentially to nodes with low degrees. Biological networks, and in particular metabolic networks, usually show a disassortative pattern [5].

Correlations among three nodes can be measured by means of the concept of *clustering*, which refers to the tendency to form triangles between the neighbors

of a vertex. Watts and Strogatz [134] proposed a measure known as *clustering coefficient*, $c_i = \frac{2E_i}{k_i(k_i-1)}$, where $E_i$ is the number of edges that exist between neighbors of the node $i$ and $k_i$ denotes the degree of the node $i$. Although this measure is helpful as a first indication for clustering, it is more informative to work with quantities which depend explicitly on the degree $k$. Therefore, a degree-dependent clustering coefficient $\bar{c}(k)$ is calculated as the clustering coefficient of nodes averaged for each degree class $k$. Metabolic networks tend to display high levels of clustering [99, 113] with $\bar{c}(k)$ having a decreasing dependence on $k$ [100].

A final mention is deserved to structures called *motifs* [11]. Motifs are small subsets of connected nodes that are found in networks more often than expected at random. They are considered as elementary functional units, and each real network has its own set of distinct motifs. Their identification provides useful insights into the typical local connectivity patterns in the network.

### 2.1.7   Null model networks and randomization methods

Null models in Complex Network Science serve to study fundamental properties of complex networks and to asses the statistical significance of a property, first measuring it in the real network and then comparing the original results to the ones obtained in the randomized versions. These models can be used to prove the existence of graphs satisfying various properties, or to provide a rigorous definition of what it means for a property to hold for almost all graphs or, finally, to act as a benchmark for specific features of real networks.

One of the most known models was the graph structure proposed by Paul Erdös and Alfréd Rényi. The *Erdös-Rényi* model [135, 136] consists on generating realizations of random networks given the total number of nodes $N$ and a total

number of links $L$, and connecting every pair of them with probability $p$. This leads to a binomial degree distribution, that can be approximated by a Poisson distribution for realizations with a large number of nodes.

Another important method to construct random networks is the *Configuration model*, an algorithm to construct random networks with a degree sequence or degree distribution $P(k)$ settled *a priori* [137, 138]. The total number of nodes $N$ remains constant. For each node, a random number $k$ is drawn from the probability distribution $P(k)$ and it is assigned to the node in the form of half-edges. The network is then constructed by connecting pairs of these link ends chosen uniformly at random. These realizations, like the Erdös-Rényi networks described above, are uncorrelated and have no clusters in the thermodynamic limit $N \to \infty$.

Instead of comparing real networks with null models as those described above, it is sometimes preferable to randomize a network obtained from real data by rewiring, *i.e.*, by picking two links at random and swapping their end [139]. While randomizing, one can preserve different properties, for instance the degrees of all nodes. Two rewiring randomization methods have been used in this thesis, one that preserves the degrees of all nodes -similar to comparing with the Configuration model- called *degree-preserving* randomization, and another that generates randomized versions taking into account that new reactions must be stoichiometrically balanced, called *mass-balanced* randomization.

### 2.1.7.1    Degree-preserving randomization

In metabolic networks, the degree-preserving randomization method is similar to the Configuration model in bipartite networks. Degree-preserving randomization works by choosing two pairs of connected nodes (metabolites and reactions) of

the bipartite network at random and swapping their ends, unless this would lead to a repeated metabolite in a reaction (see Figure 2.6 left). The steps of the algorithm are:

1. Pick two links at random: $m_1 \rightarrow r_1$ and $m_2 \rightarrow r_2$ or $r_1 \rightarrow m_1$ and $r_2 \rightarrow m_2$, where $m$ are metabolites and $r$ reactions.

2. Swap the end of the links avoiding repeated links and self-production: $(m_1 \rightarrow r_2$ and $m_2 \rightarrow r_1$ or $r_1 \rightarrow m_2$ and $r_2 \rightarrow m_1)$.

3. Repeat until $L^2$ swappings are performed, where $L$ is the total number of links in the network.

4. Make several realizations of the randomized metabolic network following the three previous steps.

Reversible reactions are rewired independently of the irreversible ones in order to preserve the degrees of metabolites which correspond to reversible and irreversible reactions. This method gives networks which preserve the degrees of metabolites and reactions and it is useful, for instance, to determine the role of the degree distribution in large failure cascades in bacterial organisms, which may have evolved towards reducing the probability of having large cascades that produce metabolic damage, increasing thus robustness [140]. This method will be used in Chapter 3.

### 2.1.7.2   Mass-balanced randomization

Mass-balanced randomization generates randomized networks by rewiring the links corresponding to substrate-reaction or product-reaction relationships, while preserving atomic mass balance of the reactions [141]. Given a reaction $r$, its

FIGURE 2.6: Left: Scheme of the degree-preserving randomization algorithm. IN and OUT degrees are conserved, but mass balance is not satisfied. Right: Scheme of the mass-balanced randomization. In this case metabolites are switched only if the new reaction is mass balanced; while reaction degrees are kept constant, the degrees of metabolites are not preserved.

atomic mass balance is given by:

$$\sum_{e \in E_r} s_{e,r} \cdot m_e = \sum_{p \in P_r} s_{p,r} \cdot m_p \qquad (2.1)$$

where $E_r$ denotes the set of substrates and $P_r$ the set of products in $r$, and $m_e, m_p$ are the mass vectors $(m_{H_2O} = (0, 2, 0, 1, 0, 0) \cdot (C, H, N, O, P, S)^T$ as an example) of $e$ and $p$, respectively. Finally, $s_{e,r}$, $s_{p,r}$ are their stoichiometric coefficients. For instance, consider the reaction $A \rightarrow B$, with $m_A = m_B = C_6H_{12}O_6$. Then, $A$ may be substituted by a compound $C$ with $m_C = C_3H_6O_3$ from within the network, resulting in the randomized reaction $2\ C \rightarrow B$, which satisfies Equation 2.1 since $2\ C_3H_6O_3 = C_6H_{12}O_6$ (see Figure 2.6 right). In addition to substituting individual substrates or products, the method also allows more complex substitutions involving pairs of substrates or products, yielding a large number of possible substitutions.

The motivation for preserving atomic mass balance of reactions, a fundamental physico-chemical constraint, is that the resulting null model allows estimating the importance of network properties with respect to evolutionary pressure. As biological systems and their properties evolve under physical constraints and evolutionary pressure, a null model which satisfies physical principles but does not account for evolutionary pressure differs from a metabolic network only in the properties which are affected by evolutionary pressure. Thus, a property deemed statistically significant following mass-balanced randomization is beyond basic physical constraints and likely to be a result of evolutionary pressure [142]. The method preserves mass balance and reaction degrees but not the degrees of metabolites, since the stoichiometric coefficients and metabolite degrees are changed. This method will be used in Chapter 3.

## 2.2   Flux Balance Analysis

A general aim of the study of a metabolic network is to characterize and understand the configuration of fluxes of the reactions constituting the network in connection to phenotype and behavior. The study of fluxes in metabolic networks deserves a special treatment more biochemically focused than in usual chemical kinetics schemes. With the knowledge of the kinetic constants of the reactions, it would be possible to solve the equations associated to the fluxes of reaction and the concentrations of metabolites in the metabolic network using proper mathematical methods. However, there is a lack in the availability of kinetic parameters [143] due to the difficulty in measuring them experimentally. As an alternative, computational techniques have been proposed in order to estimate fluxes through reactions of metabolic networks at steady-state.

*Flux Balance Analysis* is maybe the most successful and widely used approach to compute the fluxes through metabolic reactions of an organism. In addition, FBA also estimates its growth rate by maximizing the flux through the biomass reaction of the network. This technique will be used in Chapters 4, 5, and 6.

To be more specific, metabolic reactions can be represented in terms of a *stoichiometric matrix*, this being the fundamental basis in FBA and other modeling approaches [63, 107, 144, 145]. To construct a stoichiometric matrix [57, 71, 146], one must first write the typical kinetic equations which describe the temporal variation of the concentration of metabolites, which are derived from the mass conservation principle,

$$\frac{dc_i}{dt} = \sum_{j=1}^{N_R} S_{i,j} \nu_j \qquad (2.2)$$

$$\frac{dc_A}{dt} = -\nu_1 - \nu_2$$

$$\frac{dc_B}{dt} = +\nu_1 - \nu_4$$

$$\frac{dc_C}{dt} = +\nu_2 - \nu_3$$

$$\frac{dc_D}{dt} = +\nu_3 + \nu_4$$

$$\frac{d\vec{c}}{dt} = \mathbf{S} \cdot \vec{\nu}$$

$$\vec{c} = \begin{pmatrix} c_A \\ c_B \\ c_C \\ c_D \end{pmatrix} \qquad \mathbf{S} = \begin{pmatrix} -1 & -1 & 0 & 0 \\ 1 & 0 & 0 & -1 \\ 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix} \qquad \vec{\nu} = \begin{pmatrix} \nu_1 \\ \nu_2 \\ \nu_3 \\ \nu_4 \end{pmatrix}$$

FIGURE 2.7: Equations derived from mass-balance associated to a simple metabolic network. Matrix $\mathbf{S}$ is the so-called stoichiometric matrix, $\vec{\nu}$ is a vector containing all the fluxes of the metabolic network, and $\vec{c}$ denotes the vector with concentrations of metabolites.

The concentration of metabolite $i$ is denoted by $c_i$, $N_R$ is the total number of reactions, $S_{i,j}$ is the stoichiometric coefficient of metabolite $i$ in reaction $j$, and $\nu_j$ stands for the flux of reaction $j$. Note that, typically, reaction fluxes have units of mmol gDW$^{-1}$h$^{-1}$, where gDW means grams Dry Weight. Notice that the values of the $\mathbf{S}$ matrix correspond to the stoichiometric coefficients of each metabolite in each reaction. Thus, each row represents a metabolite, whereas each column represents a reaction. Therefore, if a metabolite $i$ does not participate in a reaction $j$, its stoichiometric coefficient will be 0, $S_{ij} = 0$. Otherwise, if the metabolite is a reactant, the stoichiometric coefficient will be negative, $S_{ij} < 0$, and if it is a product, it will be positive, $S_{ij} > 0$ (see Figure 2.7).

Metabolic networks are open-systems, which implies that some metabolites can leave or enter the organism. Therefore, it is not possible to arrive to a thermodynamic equilibrium state. However, it is possible to attain a non-equilibrium steady state, where the concentrations of metabolites do not change with time, forcing the system to exchange metabolites with the environment. This steady-state condition simplifies the system of coupled differential Equations 2.2 derived from mass balance into an ordinary linear system of equations, which can be written as a product of the stoichiometric matrix $\mathbf{S}$ by the vector of fluxes $\vec{\nu}$,

$$\mathbf{S} \cdot \vec{\nu} = \vec{0} \tag{2.3}$$

This is the typical form of the equation to be solved by the FBA technique. As mentioned before, it is important to notice that no kinetic parameters [147, 148] appear explicitly in Equation 2.3 and, thus, they are not needed in relation to FBA applications.

It is important to precise that, apart from the intrinsic constraints imposed by the steady-state condition, other bounds of the form $\alpha_i \leq \nu_i \leq \beta_i$ may be imposed on the values of the fluxes to render the whole scheme both chemically and biologically realistic. These upper and lower bounds may depend on the thermodynamics of reactions, more precisely on their reversibility. If reactions are reversible, fluxes can have positive or negative fluxes, whereas for the case of irreversible reactions, reactions must have only positive fluxes. Further, since the steady-state condition forces the system to exchange metabolites with the environment, constraints on *exchange* fluxes are imposed for metabolites that can either enter or leave the organism. These exchange fluxes are taken positive from the system to the environment. Notice that fluxes obtained using FBA will depend on the particular chosen external medium.

FIGURE 2.8: Example of the optimization of an objective function on a system of two variables.

In metabolic networks, there are usually more reactions than metabolites. The system of Equations 2.3 is thus underdetermined, *i.e.*, there are multiple solutions even after imposing the mentioned constraints. Therefore, a biological objective function is introduced to restrict the solution space to a single biologically meaningful solution. Technically, this means that FBA selects the state in the solution space that maximizes the value of the objective function (see Figure 2.8). This objective function depends on the biological information that one wants to extract, but usually one chooses to optimize biomass formation adjusted to be equivalent to maximize the specific growth rate of the organism. To do this, a *biomass reaction* is added to the network which simulates the biomass production. Other possible objective functions are ATP or NADH production or yield.

Often, other auxiliary reactions are needed apart from exchange and the biomass formation reactions. The first category includes physiological requirements, like the *ATP maintenance* reaction, which is a reaction which consumes ATP in order to simulate biological energetic costs for the organism which are not associated to

FIGURE 2.9: Example of a FBA calculation in a metabolic network. Reactions are denoted by squares and metabolites by circles. The biomass production reaction (red square) is labeled as $\nu_g$. Exchanges fluxes for interactions with the environment (orange arrows) are denoted with $b$ labels. A sink reaction (cyan square) is shown with a $s$ label. The ATP maintenance reaction (green square) is also shown denoted with a $M$ label.

growth. A second category are the so-called *sink* reactions, which are reactions that have not been identified yet and that consume some metabolites to avoid accumulation. A generic sink reaction has the simple form $A \to \emptyset$.

In this way, a consistent system of equations representing the whole cell metabolism is obtained and one tries to find a solution that optimizes the value of an objective function (see Figure 2.9 for a schematic picture of a FBA computation). If no solution exists for optimization of biomass production in a particular medium condition, one can assume that the system is not able to grow and therefore one can conclude that the organism is not able to survive in this medium.

The mathematical notation to denote a standard FBA problem choosing to optimize the specific growth rate is

$$
\begin{aligned}
\textbf{maximize} \quad & \nu_g \\
\textbf{subject to} \quad & \mathbf{S} \cdot \vec{\nu} = \vec{0} \\
\textbf{and} \quad & \vec{\alpha} \leq \vec{\nu} \leq \vec{\beta}
\end{aligned}
$$

where $\vec{\alpha}$ and $\vec{\beta}$ represent the vectors determining the lower and upper bounds of the reactions, and $\nu_g$ denotes the specific growth rate.

The software used to perform these calculations in this thesis is GNU Linear Programming Kit (GLPK) [149–152], through its associated solver GLPSol. This solver uses a dual *simplex* algorithm to compute the solutions. It is a variant of the normal simplex algorithm [153]. The latter is an iterative algorithm which is based on finding first feasible solutions and then finding the most optimal solution based on these feasible solutions. On the contrary, dual simplex works by first finding optimal solutions and then finding a feasible solution, again, if it exists.

### 2.2.1    Formulation of the biomass reaction

FBA problems are usually solved by maximizing the flux through the biomass reaction [56, 117, 119]. This typically gives a particular flux state of the metabolic network compatible with the constraints. However, the solution obtained by FBA is often not unique. In some cases, the metabolic network is able to achieve the same specific growth rate by using alternate reactions and pathways. Therefore, phenotypically different solutions that optimize the specific growth rate are possible, implying that FBA solutions can be degenerate [63].

Technically, the biomass reactions is modeled as a reaction, $aA + bB + cC + dD... \longrightarrow xX + yY + zZ$, which produces and consumes some specific metabolites (see Figure 2.9). These metabolites are known biosynthetic precursors present in the metabolic network under consideration. The key point is given by their stoichiometric coefficients in the biomass reaction, which are experimentally measured proportions in the biomass of the organism measured in dry weight conditions. The stoichiometric coefficients of the metabolites participating in the biomass reaction have units of mmol gDW$^{-1}$, and the biomass reaction has units of $h^{-1}$. It is worth stressing that this reaction simulates the growth of an

organism given a set of external nutrients and that its coefficients are adjusted so that its flux is equivalent to the specific growth rate of the organism.

FBA can also maximize the biomass yield, which is the equivalent to maximize the specific growth rate but taking into account that the maximum uptake of the carbon source, for example glucose, must be set to 1 mmol gDW$^{-1}h^{-1}$ to set the maximum amount of biomass that can be produced per 1 mole of nutrient.

### 2.2.2   Simulation of different environments

It is important to make explicit the way to simulate changes in the environment using FBA. To do this, one must tune the upper and lower bounds of the values of the exchange reactions of the metabolites that are present in the environment. As an example, suppose that one wants to model that glucose is present in the environment and that, therefore, the organism consumes it in order to obtain energy. The explicit form of the constraint of the exchange flux of glucose will be $-10 \leq \nu_{glucose}^{exchange} \leq \infty$, which means that the organism can expel as much as glucose as it wants but that it can eat glucose with a maximum uptake of 10 mmol gDW$^{-1}$h$^{-1}$.

Notice that nutrients have a negative lower bound and an unlimited upper bound, whereas waste products have a value of the lower bound of 0 and unlimited upper bound, which means that the organism cannot uptake it but, if the compound is generated inside the organism, it can be expelled to the exterior as waste. As an example, this would be the case for $CO_2$ in *Escherichia coli*, which is not eaten by the organism but that is expelled due to respiration.

To summarize, an environment is simulated by choosing a set of nutrients and assigning a lower bound $-\alpha_i$ to each nutrient, which is the maximum uptake

of each nutrient, and assigning a lower bound of 0 to components not present in the environment. For all external metabolites, the upper bound is set to $\infty$. Therefore, for nutrients one has $-\alpha \leq \nu_{nutrient}^{exchange} \leq \infty$, whereas for waste products one has $0 \leq \nu_{waste}^{exchange} \leq \infty$. The rest of reactions are modeled as told in the previous section.

#### 2.2.2.1    Construction of minimal media

A minimal medium is the minimal set of metabolites which ensure the viability of an organism. The modelization of these media can be made as in Reference [119]. Minimal media consist of a set of mineral salts, and one source of carbon, of nitrogen, of sulfur and of phosphorus, from four families representing carbon, nitrogen, phosphorus, and sulfur compounds, respectively. To construct different minimal media, the set of mineral salts is always the same -which contains, for example, magnesium sulfate, iron chloride, and calcium chloride-, but each source family is browsed while the other three sources are fixed to the standard metabolites of each kind (C*: glucose, N*: ammonia, P*: phosphate, S*: sulfate) (see Figure 2.10).

#### 2.2.2.2    Construction of rich media

Sometimes it can be useful to perform FBA computations in a medium with more components that the ones present in a minimal medium. These media containing more nutrients than a minimal medium are called rich media. One of this rich media is an amino acid-enriched medium. This medium can be constructed from a minimal medium with the standard metabolites explained in Section 2.2.2.1 (glucose, ammonia, phosphate, and sulfate), by adding the following set of amino acids: D-Alanine, L-Alanine, L-Arginine, L-Asparagine, L-Aspartate,

| Variation of carbon sources | | | | | Variation of phosphorus sources | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Medium 1 | $C_1$ | N* | P* | S* | Medium 1 | C* | N* | $P_1$ | S* |
| Medium 2 | $C_2$ | N* | P* | S* | Medium 2 | C* | N* | $P_2$ | S* |
| Medium 3 | $C_3$ | N* | P* | S* | Medium 3 | C* | N* | $P_3$ | S* |

| Variation of nitrogen sources | | | | | Variation of sulfur sources | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Medium 1 | C* | $N_1$ | P* | S* | Medium 1 | C* | N* | P* | $S_1$ |
| Medium 2 | C* | $N_2$ | P* | S* | Medium 2 | C* | N* | P* | $S_2$ |
| Medium 3 | C* | $N_3$ | P* | S* | Medium 3 | C* | N* | P* | $S_3$ |

FIGURE 2.10: Examples of the construction of minimal media. Asterisks denote the standard metabolite of each kind. To construct carbon media, the sources of nitrogen, phosphorus and sulfur are set to the standard components of each kind whereas the carbon sources are varied. The same procedure applies to construct nitrogen, phosphorus and sulfur media.

D-Cysteine, L-Cysteine, L-Glutamine, L-Glutamate, Glycine, L-Histidine, L-Homoserine, L-Isoleucine, L-Leucine, L-Lysine, L-Methionine, L-Phenylalanine, L-Proline, D-Serine, L-Serine, L-Threonine, L-Tryptophan, L-Tyrosine, L-Valine. This set of amino acids enriches the minimal medium allowing the organism to take them as nutrients. Otherwise the organism would have to synthesize them, resulting in a more stringent environment for the organism. To simulate the presence of this set of amino acids in the medium, the exchange constraints bounds of these amino acids are set to -10 mmol/(gDW·h).

Another famous rich medium is called Luria-Bertani Broth [154]. The Luria-Bertani Broth used in this thesis contains all the nutrients present in the amino acid-enriched medium, but it contains as additional compounds purines and pyrimidines, vitamins (namely biotin, pyridoxine, and thiamin), and also the nucleotide nicotinamide monocleotide [155]. The exchange constraints bounds of these compounds are usually set to -10 mmol/(gDW·h) ($\nu_{compound}^{exchange} \geq -10$) for *Escherichia coli*.

### 2.2.3   Activity and essentialify of genes and reactions

An important application of FBA is to compute the *activity* and *essentiality* of reactions in a network. These concepts can be applied either to genes or reactions, since a reaction is catalyzed by an enzyme which at the same time is codified by a gene or a set of genes. Both concepts will be analyzed in Chapter 4.

The concept of activity is quite simple. A reaction is said to be active when, given an external environment, the chosen reaction carries a non-zero flux. The concept of essentiality is more subtle. It refers to how a network, and thus the growth rate, is affected when one reaction is forced to be non-operative through the knockout of a reaction or of the corresponding gene.

To calculate the effect of the knockout of a reaction, the selected reaction is removed from the network, which is equivalent to force the chosen reaction to have a null flux. The new system is usually called a mutant. In terms of the notation used before, this is modeled as $\nu_i = 0$ with $i$ the removed reaction and $\nu_i = 0$ its flux. Thus, a FBA problem with a reaction $i$ constrained to be non-active is

$$
\begin{aligned}
\textbf{maximize} \quad & \nu'_g \\
\textbf{subject to} \quad & \mathbf{S} \cdot \vec{\nu} = \vec{0} \\
\textbf{and} \quad & \vec{\alpha} \leq \vec{\nu} \leq \vec{\beta} \\
& \nu_i = 0
\end{aligned}
$$

where $\nu'_g$ denotes the growth rate of the mutant. As a consequence, the system can respond in three different ways as compared to the non-perturbed case $\nu_g$:

1. The growth rate is unaltered $\nu'_g = \nu_g$.

2. The growth rate is decreased $0 < \nu'_g < \nu_g$, which means that the biomass formation of the organism is reduced but the organism is still alive at the

expense of losing some performance.

3. The growth rate takes a null value $\nu'_g = 0$, meaning that the performed knockout is lethal for the organism. This is the signature of essentiality.

It has been shown that FBA predicts gene essentiality with an accuracy of 90% [117] in *Escherichia coli* under glucose aerobic conditions, which means that FBA is a reliable tool to predict whether a knockout will be lethal or not in this particular condition.

### 2.2.4   Flux Variability Analysis

Sometimes it is useful to identify which are the minimum and maximum bounds that each reaction can take independently of the growth optimality condition. In this way, one can have an idea of the flux space for a particular environmental condition, and in particular which reactions can have a non-zero flux in a given environment, since some reactions may be active for low values of the growth rate but the same reactions must have a zero flux in order to ensure growth optimality. This may happen due to the fact that some reactions can compete with the growth reaction by consuming metabolites needed to grow and therefore this would reduce the flux through the biomass reaction. As a consequence, when one optimizes the flux through the biomass reaction, all reactions whose activation competes with the flux of the biomass reaction will have a null value in order to assure maximum growth conditions.

In addition to identifying those reactions that can compete with the growth rate, reactions whose minimum and maximum values are close indicate that they may be important for the organism since those reactions are allowed only to have a low variability in their fluxes. To know the minimum and maximum flux values

of a reaction, one applies the technique called Flux Variability Analysis (FVA) [106, 107, 156].

In most applications of FVA, the biomass reaction is imposed to have a minimum value $\nu_g \geq \nu_g^{min}$ to ensure viability. Hence, one can consider that the limiting fluxes correspond to states where the organism is alive, even if the growth rate is not the maximum value that the organism can achieve. Using the mathematical notation used in Linear Programming computations, FVA for each flux of a metabolic reaction can be written as follows

$$
\begin{array}{ll}
\textbf{minimize} & \nu_i \\
\textbf{subject to} & \mathbf{S} \cdot \vec{\nu} = 0 \\
& \vec{\alpha} \leq \vec{\nu} \leq \vec{\beta} \\
& \nu_g \geq \nu_g^{min}
\end{array}
\qquad
\begin{array}{ll}
\textbf{maximize} & \nu_i \\
\textbf{subject to} & \mathbf{S} \cdot \vec{\nu} = 0 \\
& \vec{\alpha} \leq \vec{\nu} \leq \vec{\beta} \\
& \nu_g \geq \nu_g^{min}
\end{array}
$$

However, it may happen that one is interested in capturing all the possible scenarios independently of the value of the flux of the biomass reaction, since in this way non-optimal/low-growth scenarios can be taken also into account. Therefore, FVA can be modified to compute the minimum and maximum possible values of the flux of each reaction regardless of the value of the biomass formation rate. To this end, the value of the flux of the biomass reaction is not constrained and any positive value is allowed, $\nu_g \geq 0$. Under this condition, one will obtain the maximal set of reactions that can be active in the considered medium independently of the rate of biomass formation. This variation of FVA [157, 158] will be used in Chapter 4 and 5. Using the previous notation, this version of FVA, which we call *Biomass unconstrained Flux Variability Analysis*, can be written as

$$
\begin{array}{ll}
\textbf{minimize} & \nu_i \\
\textbf{subject to} & \mathbf{S} \cdot \vec{\nu} = 0 \\
& \vec{\alpha} \leq \vec{\nu} \leq \vec{\beta} \\
& \nu_g \geq 0
\end{array}
\qquad
\begin{array}{ll}
\textbf{maximize} & \nu_i \\
\textbf{subject to} & \mathbf{S} \cdot \vec{\nu} = 0 \\
& \vec{\alpha} \leq \vec{\nu} \leq \vec{\beta} \\
& \nu_g \geq 0
\end{array}
$$

## 2.3 Model organisms

Information about metabolism of specific organisms [56, 65, 66, 68, 69, 119, 159–163] -most single cell- are gathered in databases, like the BiGG database [62], Kyoto Encyclopedia of Genes and Genomes (KEGG) [60], BioCyc/EcoCyc/MetaCyc [61], BRENDA [164], etc. The BiGG database deserves special attention in this thesis, since it has been extensively used as it contains full reconstructions of metabolic networks for specific organisms including all the biochemical reactions and the biomass formation function in order to compute FBA solutions for different organisms.

The BiGG database provides high-quality curated information. Network reconstructions coming from this database are structured in compartments like cytosol inside cells or periplasm -the space bordered by the inner and the outer membranes in Gram-negative bacteria. Therefore, metabolites present in different compartments of the organisms are treated as different nodes. Using different compartments allows the inclusion of transport systems in both the inner and outer membrane and thus the metabolic machinery of organisms is more accurately represented. As an example, water in the periplasm will be a different metabolite than water in the cytosol. In addition, a directed bipartite representation of the metabolic network can be constructed since in the databases reactants, products, reversible, and irreversible reactions are distinguished. Further, the BiGG database specifies which enzyme catalyzes each reaction and also which gene or set of genes codifies each enzyme. However, reactions are also listed which have neither associated enzymes nor genes. It may be that, for these particular reactions, enzymes have not been identified yet or that some reactions are spontaneous and they can take place without the need of an enzyme.

It is important to notice that there exist different versions for each metabolic network of each organism. This happens due to the fact that the reconstructions of metabolic networks are constantly improved and, therefore, versions are constantly updated. As an example, the first version of *Escherichia coli* [165] contained 660 genes, 627 reactions, and 438 metabolites, while the last version of *Escherichia coli* [119] contains 1366 genes, 2250 reactions, and 1805 metabolites.

### 2.3.1    *Escherichia coli*

*Escherichia coli*, abbreviated as *E. coli*, is the most studied prokaryotic organism and it is the bacterial model that is most frequently used in experiments due to the ease of its manipulation. More precisely, the strain studied in this thesis is K-12 MG1655. This strain colonizes the lower gut of animals. Moreover, it has been maintained as a laboratory strain with minimal genetic manipulation.

Three versions of this strain have been used in this thesis. The first one is *i*AF1260, which can be obtained either from Reference [117] or directly from the BiGG database. This version is based on an earlier reconstruction called *i*JR904 [166], on the annotation of the genome of *E. coli* from Reference [167], on contents from the EcoCyc (an *E. coli* version of BioCyc) database [168] and on specific biochemical characterization studies from Reference [117]. The *i*AF1260 version contains 2077 reactions, 1669 metabolites, and 1260 genes [117] (see Table 2.1). Metabolites are located in three compartments: exterior, periplasm and cytosol. Notice that although the exterior is not a real compartment, it is treated in this way in order to be able to use the exchange reactions explained before.

The most recent version of *E. coli* is *i*JO1366 [119]. It is an update of the *i*AF1260 version. EcoCyc [169] and the KEGG database [170] were used in order

to improve the *i*AF1260 version, in addition to experimental techniques [119]. It contains 2250 biochemical reactions, 1805 metabolites, and 1366 genes [119] (see Table 2.1). Like in *i*AF1260, metabolites are located in cytosol, periplasm and exterior.

A simplified version called core *E. coli* metabolic model is also used, which can be obtained either from References [63, 118] or the BiGG database. It is a condensed version of the genome-scale metabolic reconstruction *i*AF1260 that contains 73 metabolic reactions in central metabolism, 72 metabolites, and 136 genes (see Table 2.1). This network is complemented with a biomass formation reaction and an ATP maintenance reaction.

### 2.3.2   *Mycoplasma pneumoniae*

*Mycoplasma pneumoniae*, abbreviated as *M. pneumoniae*, is a human pathogen of primary atypical pneumonia that has recently been proposed as a genome-reduced model organism for bacterial and archaeal systems biology [56, 120, 171, 172]. Interest in this organism has grown recently since it lacks many anabolic processes and rescue pathways compared to more complex organisms. This in turn translates into a highly linear metabolism singularly suited to study basic metabolic functions [56]. This property will be again mentioned in Chapters 3 and 4.

The first version of the metabolic network of *M. pneumoniae* used in this thesis was published in Reference [120], where the authors integrated biochemical and computational studies, complementing the information using the KeGG database. Its metabolic reconstruction contains 187 reactions taking place in cytosol and in exterior, the number of metabolites is 228, and the number of genes is 140 (see Table 2.1).

| Organism | $N_R$ | $N_M$ | $N_G$ | Source |
|---|---|---|---|---|
| *E. coli i*AF1260 | 2077 | 1669 | 1260 | Ref. [117], BiGG |
| *E. coli i*JO1366 | 2250 | 1805 | 1366 | Ref. [119] |
| *E. coli* core model | 73 | 72 | 136 | Refs. [63, 118], BiGG |
| *M. pneumoniae* | 187 | 228 | 140 | Ref. [120] |
| *M. pneumoniae i*JW145 | 240 | 266 | 145 | Ref. [56] |
| *S. aureus i*SB619 | 642 | 644 | 619 | Ref. [66], BiGG |

TABLE 2.1: Summary of the properties of all metabolic reconstructions used in this thesis. $N_R$, $N_M$, and $N_G$ stand for the number of reactions, metabolites, and metabolic genes respectively. Metabolites in different compartments are treated as different metabolites.

The *i*JW145 version of *M. pneumoniae* is the last update [56]. This network was constructed by determining the behavior of the organism under different nutrition conditions, using literature information and experimental data. It contains 240 biochemical reactions, 266 metabolites, and 145 genes (see Table 2.1). Metabolites can be located in cytosol and exterior.

### 2.3.3    *Staphylococcus aureus*

*Staphylococcus aureus*, abbreviated as *S. aureus*, is found in the human respiratory tract and on the skin. It is an anaerobic bacterium which is present world-wide, and it is a common cause of skin infections, respiratory disease, and food poisoning. The strain used in this thesis is N315, a major pathogen which is able to acquire antibiotic-resistance [173].

The *i*SB619 version of *S. aureus* can be obtained either from the BiGG database or from Reference [66]. To construct this model, the authors used the KeGG database and the Comprehensive Microbial Resource (CMR) at The Institute for Genomic Research (TIGR) website [174]. Missing functions were annotated based on reported evidence from this organism, as well as for *Bacillus subtilis* and *E. coli*. The number of reactions is 642 and the number of metabolites is

644 (see Table 2.1). Like in *M. pneumoniae*, there are only cytosol and exterior compartments.

# Chapter 3

# Structural knockout cascades in metabolic networks

This chapter presents the analysis of the response of metabolic networks of model organisms to different forms of structural stress, including removals of individual and pairs of reactions and knockouts of single or co-expressed genes. Local metabolite motifs can be used as predictors of failure cascade sizes caused by individual failures, and for amplification effects in cascades caused by multiple failures. Correlation between gene essentiality and damages produced by single gene knockouts is detected, which points out that genes controlling high-damage reactions tend to be expressed independently of each other. This study is carried out for three characteristic organisms: *Mycoplasma pneumoniae*, *Escherichia coli*, and *Staphylococcus aureus*.

The architecture of complex networks is imprinted with universal features that affect their resilience and condition their behavior [14, 96, 175]. Most relevant, the scale-free connectivity (see Chapter 2, Section 2.1.2) of many natural and man-made networks explains their fragility in front of attacks to the most connected nodes, while they are able to deal with accidental failures of single components [176, 177]. A manifestation of this fragile yet robust nature of complex networks is that the failure cascade triggered by a local shock rarely propagates to the whole system [178–181]. At the same time, it is worth to remember that network studies have mainly focused on single node failures, and that systemic responses to more globalized forms of structural and functional stress still remain to be explored.

In a more biological context, metabolic networks are among the best probed in terms of robustness in front of a variety of *in silico* perturbation experiments. They have been found to comply with the design principles of error-tolerant scale-free networks [5], and recent progress in network dynamics is also starting to portray the concept of stress-induced network rearrangements [102, 182]. The exploration of single biochemical reaction inactivations has shown that when a reaction is forced to be non-operative, a cascade of consequent failures propagates to a variable extent trough the whole network, and that the structural organization of metabolic networks reduce the likelihood of large damaging cascades [140]. At the same time, many individual mutations affecting enzyme-coding genes seem to have very little effect on cell growth [183, 184]. By contrast, the impact of multiple failures could go beyond the mere accumulation of individual effects, producing amplified damage due to peculiar biochemical interweaving or gene epistatic interactions [185].

The analysis presented in this chapter considers the removal of single and pairs

of biochemical reactions and the knockout of individual genes and clusters of co-expressed genes in three bacteria, *Mycoplasma pneumoniae*, *Escherichia coli*, and *Staphylococcus aureus*. To simulate the effect of reaction knockouts, a cascading failure algorithm [140] is used and the significance of the obtained results is assessed using two different null models called degree-preserving randomization (DP) (see Chapter 2, Section 2.1.7.1) and mass-balanced randomization (MB) (see Chapter 2, Section 2.1.7.2). One finds that, for the three organisms, the sizes of cascade distributions span a broad range of values, with many short propagations but a few that spread at the systems level. *M. pneumoniae* exhibits similar network responses to *E. coli* and *S. aureus*, although its increased linearity and reduced redundancy [120] threaten its robustness against individual reaction removals. For all three organisms, the impact of failure cascades can be predicted in terms of local network motifs. In this way, targets prone to introduce structural vulnerability can be readily detected prior to experimental testing without expensive computations, even for large and complex organisms. This chapter also reports the effects of single and multiple gene knockouts in *M. pneumoniae* by coupling, through enzyme activity, its metabolic network to the experimentally measured gene co-expression network. One observes that genes related to high-damage reactions are essential for the organism and that their expression tends to be isolated from that of other genes. This hints at the interplay between metabolism and genome, apparently evolved to favor the robustness of this organism by avoiding the potentially catastrophic effect of coupling the co-expression of structurally vulnerable metabolic genes. At the same time, one finds that this enables the organisms the ability to perform more efficient metabolic regulation at the expense of losing some of the maximum attainable robustness determined by physico-chemical constraints.

The contents of this chapter correspond to References [94, 186].

## 3.1    Cascading failure algorithm

It is important to start by explaining how the cascading failure algorithm works after a reaction or a set of them are inactivated. First of all, the metabolic networks of *M. pneumoniae*, *E. coli* (*i*AF1260), and *S. aureus* (*i*SB619) (see Chapter 2, Section 2.3, their networks can be seen in Supplementary Tables C3-1, C3-2, and C3-3) are modeled as a bipartite semidirected network (see Chapter 2, Section 2.1.1), with two specific criteria:

1. All biochemical reactions in the genome-scale metabolic reconstructions (GENREs) are considered except exchange, sink, biomass formation, and ATP maintenance reactions.

2. All metabolites involved in the reactions included in the network representation are considered. In particular, hubs participating in a huge number of reactions are not excluded. Hubs stay neutral with respect to structural cascades and do not contribute to propagate them. Due to their large number of connections, they are highly unlikely to become nonviable as a consequence of single or double cascades reaching them.

The cascading failure algorithm [140] is based on the states of the nodes on the network, *i.e.*, nodes can be viable or non-viable. Non-viable nodes spread the perturbation, whereas viable do not. To define viability, two aspects are considered since a bipartite representation of the metabolic network is used. The first one refers to metabolites, and consists on the fact that each viable metabolite must have at least one outgoing and one incoming connections so as to prevent accumulation or depletion of the metabolite. For reactions, the criterion is that all metabolites participating in a reaction must be viable.

FIGURE 3.1: Example of how the cascading failure algorithm is applied to a metabolic network. 1) For clarity, metabolites 4 and 5 are labeled with R and 7 and 8 with P depending on whether they are reactants or products of the reversible reaction denoted $d$ (for simplicity, only a reversible reaction is considered in this illustration, the rest being assumed to be irreversible). The cascade starts when reaction $c$ fails. 2) Therefore, metabolites 3 and 6 become non-viable. Because metabolite 6 is connected to reaction $g$, the later becomes non-viable, turning also metabolite 12 non-viable. Notice that metabolite 11 loses one IN connection, but it is still viable, meaning that one of the waves of the cascade stops here. However the other wave keeps spreading. 3) Metabolite 4 causes the reversible reaction $d$ to remain viable only towards the production of metabolites 7 and 8. 4) Consequently, metabolite 4 becomes non-viable, and so its associated reactions also become non-viable. 5) The cascade spreads until all metabolites and reactions affected by the cascade remain viable. Finally, note that metabolites 1, 2, 3, 13, and 14, which initially have no incoming or outgoing connections, are not considered nonviable by the algorithm.

The algorithm works by removing one or more reactions, and then checking the viability of its surrounding metabolites. If they are viable, the cascade stops, otherwise the cascade is spread into other reactions and metabolites until all remaining nodes satisfy the mentioned criteria (see Figure 3.1). When the cascade stops, the corresponding damage is quantified as the number of reactions turned non-operational.

Reversible reactions (see Chapter 2, Section 1.1.2) deserve a special treatment in this algorithm. They are decoupled in two half-nodes, the forward and the reverse direction. A cascade propagating to a metabolite of a reversible reaction fixes it in the forward or reverse direction depending on whether the single incoming or outgoing link left to the affected metabolite is connected to the forward or reverse half of the reaction (see Figure 3.1, step 3). In all cases, when any metabolite of a reversible reaction has this reaction as the single one producing and consuming it, the reaction must be removed to satisfy the viability criterion.

## 3.2    Impact of reaction failures

This first result is the distribution of damages for cascades triggered by individual and by pairs of reactions in the metabolic networks of *M. pneumoniae*, *E. coli*, and *S. aureus*. Later, local network motifs responsible for the propagation of cascades are identified, and a local predictor for damage is proposed.

### 3.2.1    Impact of individual reactions failures

Although close to 50% of all individual reaction failures in the three organisms considered do propagate cascades, most cascades are indeed small (59% of the

cascades in *M. pneumoniae*, 38% in *S. aureus*, and 55% in *E. coli* propagate to only one or two reactions). However, the removal of some particular reactions may trigger relatively far reaching damages. This is shown in Figures 3.2a-c, that display the cumulative probability distributions $P(d'_r \geq d_r)$ that the failure of a reaction $r$ attains at least $d_r - 1$ other reactions in each metabolic network. All species show similar broad distributions, although the crossover in the tail of the distribution from power-law-like to exponential-like is not evident in *M. pneumoniae* probably due to its limited redundancy. In order to assess the significance of cascades, the computed distributions are compared with those corresponding to DP randomized variants of the metabolic networks taken as null models (see Chapter 2, Section 2.1.7.1).

To check consistency, Kolmogorov-Smirnov (K-S) tests [187] (see Appendix A) are performed measuring the maximum absolute difference between the null model and the empirical distributions (see caption of Figure 3.2 for specific values). This difference is transformed into a significance level directly compared to a chosen threshold, typically $\alpha = 0.05$. If the significance associated to the K-S test statistic is equal or smaller than $\alpha$, the compared distributions cannot be considered consistent. Both *E. coli* and *S. aureus* display values much below the threshold, meaning that the empirical distributions are not determined just by the connectivity imposed by the degrees of metabolites. Comparing both distributions, the metabolic organization of the organisms appears to have evolved towards reducing the likelihood of large failure cascades (probably lethal for the organisms) or, equivalently, towards increased structural robustness, as previously seen for *S. aureus* and for an older version of the metabolic network of *E. coli* in [140]. In contrast, the value of the associated significance level for *M. pneumoniae* is very similar to the threshold. As a consequence, one cannot say that the difference between cascade size distributions in the original network

FIGURE 3.2: Damage in cascades triggered by individual reactions. a-c)
Cumulative probability distribution functions of damages in *M. pneumoniae*,
*E. coli*, and *S. aureus*. Results are compared with damages produced in DP
randomized versions of the metabolic networks in order to discount structural
effects. In each case, the solid black curve is the average over 100 realizations.
Results for *S. aureus* and an older version of *E. coli* were already presented
in Reference [140]. The results of the Kolmogorov-Smirnov tests are given in
terms of the K-S statistic and its associated significance level (K-S statistic/as-
sociated significance level) (see Appendix A): 0.095/0.07, 0.086/0.0002, and
$0.079/1.4 \cdot 10^{-11}$ for *M. pneumoniae*, *S. aureus*, and *E. coli* respectively. With
a significance value of $\alpha = 0.05$, distributions of damages can be considered
not consistent with those for randomized variants, except for *M. pneumoniae*.
d) Spearman's rank correlation coefficient $\rho_S$ between predictors and damages,
plotted against metabolic network size (number of reactions R). Results are
compared to random reshuffling of the predictor value associated to reactions
(100 realizations for each organism). Average Spearman's rank correlation
coefficients for the randomizations appear in black, and error bars delimit the
maximum and the minimum values obtained.

and in the randomized counterparts is statistically significant, even though the probability for large cascades is still smaller in the original metabolic network. This can be explained by the increased linearity and limited redundancy of *M. pneumoniae* metabolic network structure, according to available data [120].

Along with structure, biochemical insight contributes to explain why some reactions trigger larger cascades. For *M. pneumoniae*, the most vulnerable reactions can be classified into four groups related to vital functions. One group is associated to metabolites phosphoenolpyruvate and protein L-histidine, each solely produced by one generating reaction and both of them directly related to phosphorylation processes, vital for instance in the synthesis of ATP. The second group relates to formate, which has a prominent role in the energy metabolism on many bacteria. The third group involves reactions where the important metabolite is thioredoxin, an antioxidant protein essential to reduce oxidized metabolites, along $NADP^+$. Finally, the failure of reactions in the fourth group trigger large cascades that affect the synthesis of fatty acids by turning acyl carrier proteins inviable.

Prediction of the size of the cascades is possible by looking to the local information corresponding to the triggering reaction. An expression for the predictor $P_r$ for the damage spreading from the triggering reaction $r$ which is surrounded by $m$ metabolites is:

$$
\begin{aligned}
P_r = \sum_{m \in r} \Big[ & (k_i + k_b)\delta_{k_o}^0(\delta_{k_b}^1 + \delta_{k_b}^0)(\delta_{k'_o - k_o}^1 + \delta_{k'_b - k_b}^1) \\
& + (k_o + k_b)\delta_{k_i}^0(\delta_{k_b}^1 + \delta_{k_b}^0)(\delta_{k'_i - k_i}^1 + \delta_{k'_b - k_b}^1) \\
& - \delta_{k_i}^0 \delta_{k_o}^0 \delta_{k_b}^1 \delta_{k'_b - k_b}^1 \Big] .
\end{aligned}
\tag{3.1}
$$

Degrees $k_i$, $k_o$ and $k_b$ refer respectively to the number of incoming, outgoing

and bidirectional links of metabolite $m$ (reactant or product) associated to the triggering reaction $r$ after discounting the links used to propagate the cascade, and $k_i'$, $k_o'$ and $k_b'$ denote the original values before the cascade is triggered. $\delta_a^b$ are used for the Kronecker's delta function. Basically, this predictor identifies metabolites susceptible to propagate the cascade, which are those having originally just one IN or just one OUT link, which is the one connecting them to the triggering reaction, or those connected to the triggering reaction by a *bidirectional* link and lacking in or out connections. The contribution to the predictor of one of those metabolites counts the number of connections of this metabolite with the rest of reactions, which can then be considered susceptible to become non-viable and propagate the cascade (see Figure 3.3 for illustrations showing how the measure works for some particular cases).

Propagator motifs are represented by branched metabolites with just one in or out connection that happens to be attached to the triggering reaction. The higher the branching ratio of these metabolites, the higher the likelihood that the reaction propagates a large cascade, and thus to become a target for structural vulnerability in the network. To give an example, the two most vulnerable reactions in *M. pneumoniae* produce phosphoenolpyruvate, a compound involved in Glycolysis and Gluconeogenesis that acts as a source of energy. It happens to be a highly-branched cascade propagator motif connected to two reversible reactions and, as a product, to eight irreversible reactions (see Figure 3.4 for a categorization of cascade propagator motifs in bipartite networks).

To check the predictive power of our predictor $P_r$, Spearman's rank correlation coefficient $\rho_S$ between predictors and damages are measured for each organism (see Appendix B). Basically, Spearman's correlation [188] is the Pearson correlation coefficient between two ranks, here given by the positions in ordered lists of reactions according to predictor values $P_r$ and damages $d_r$. A high ranking

FIGURE 3.3: Examples of application of Equation 3.2 to several configurations of metabolites and reactions. Triggering reactions are colored yellow, whereas metabolites which spread the cascade are colored red. For clarity, the contribution of each metabolite to the value of $P_r$ is also given.

position by predictor value is expected to correlate with vulnerable reactions at the top of the damage ranking. For all three organisms, very high values of the correlation coefficient are found, which are statistically significant (see Figure 3.2d). This evidences the ability of this predictor, calculated on the basis of local information, to rank reactions by damage without directly computing the effect of the failure.

### 3.2.2 Non-linear effects triggered by pairs of reactions cascades

As expected, the simultaneous failure of two reactions leads to higher damages compared to single reaction failures as shown in Figure 3.5. The graphs display

## Local motifs triggering individual cascades



FIGURE 3.4: Motifs of cascade propagation after failure of individual reactions. Cases a-j result into cascades with $d_r$ larger than 1, while cases k-p correspond to potential transmitters in the sense that they may or may not spread the cascade.

the cumulative probability distributions $P(d'_{rr'} \geq d_{rr'})$ calculated from all possible pairs of reactions initiating the cascades. It is worth stressing that the order of initiation is irrelevant. Notice that the exponential cut-off is still present, and becomes prominent even for *M. pneumoniae*. Again, metabolic robustness is assessed by comparing cascades in the original networks with those in DP randomized counterparts using K-S tests (see caption of Figure 3.5 for specific values). One finds that, for all three organisms including *M. pneumoniae*, the probability for large cascades triggered by pairs of reactions is significantly smaller in the original metabolic networks as compared to those in the randomized variants, suggesting that the organization of metabolic networks has evolved towards protecting metabolism against multiple reaction failures.

FIGURE 3.5: Damage in cascades triggered by pairs of reactions. a-c) Cumulative probability distribution functions of damages in *M. pneumoniae*, *E. coli*, and *S. aureus*. Results are compared with damages produced in DP randomized versions of the metabolic networks in order to discount structural effects. In each case, the solid black curve is the average over 100 realizations. Results of the Kolmogorov-Smirnov tests (K-S statistic/associated significance level) (see Appendix A): 0.15/0, 0.14/0, and 0.13/0 for *M. pneumoniae*, *S. aureus*, and *E. coli* respectively. Taking $\alpha = 0.05$, distributions of damages can be considered not consistent with those for randomized variants. d) Most frequent double cascades output. Solid line: interference without amplification. It is related with cases b and c in Figure 3.6. Dashed line: no interference, which is related with case a in Figure 3.6. e) Non-linear effects in double cascades. Solid line: overlap. It is related with cases c and e in Figure 3.6. Dashed line: amplification. Amplification is related with cases d and e in Figure 3.6.

It can also be observed that cascades caused by individual reactions combine in different ways when two reactions fail simultaneously (see Figure 3.6). The crucial concept here is that of the pattern of interference of the respective areas of influence of the two individually considered cascades. By that, one refers to all metabolites and reactions altered[1], removed or not, by each single cascade. If there is no interference, the total damage $d_{rr'}$ is additive and equal to the sum of the two single damages $d_r$ and $d_{r'}$. Otherwise, different situations are possible leading to a combined damage that can be equal, larger or smaller than the single added values. The latter case is a univocal signature of cascade overlapping $o_{rr'}$, pointing to the existence of a common subset of reactions that fail in both cascades (the most extreme realization is when one cascade is totally contained in the other). More interesting is the situation when, irrespectively of the presence or absence of overlap, a non-linearly amplified damage is detected, involving a number $a_{rr'}$ of new reactions that break down under simultaneous black outs. For all cascades,

$$d_{rr'} = d_r + d_{r'} - o_{rr'} + a_{rr'} \tag{3.2}$$

Interference without amplification is the most common situation, followed by the absence of interference (see Figure 3.5d). In contrast, overlap and amplification happen for a very small fraction of all double cascades, and their occurrence decreases with the size of the organism (see Figure 3.5e). In particular, the reduced incidence of amplification represents a new signature that organizational principles at play ensure the robustness of the organisms, despite increasing complexity and interweaving.

---

[1]Reactions altered but not removed are reversible reactions that become directed by effect of the cascade.

## Patterns of interference



## Interference metabolic network motifs



FIGURE 3.6: Cascade propagator network motifs and typology of double cascades. a-e) Illustration of possible interference patterns between individual cascades: additive, interference without overlap or amplification, interference with overlap and without amplification, interference without overlap and with amplification, interference with overlap and amplification, respectively. Blue and yellow stand for single cascades, green for interference, and red for overlap and amplification, depending on whether the red zone is in the interference zone (green) or not. f-k) Metabolic network motifs in the interference of two individual cascades that induce amplification. Cases f-g) Motif caused by a metabolite which loses its only generating reaction and at the same time it is the reactant of several reactions. These reactions are going to be become non-viable. Case g is equivalent to f but inverting the sense of the links. Case h) Metabolite which has been left with one connection to a reversible reaction. This reversible reaction has zero net flux and becomes inviable. Cases i-j) This motif appears when a modified metabolite is lead with only one incoming connection coming from a reversible direction. This fixes the reversible reaction towards the production of this metabolite. If this step turns a metabolite of the reversible reaction inviable, the reversible reaction becomes inviable (*follows to next page*).

FIGURE 3.6: (*Follows from previous page*) Therefore, this motif is a potential trigger of amplification. Case j is equivalent to case i when the senses of the reactions are inverted. Case k) The individual cascades fix the sense of a reversible reaction oppositely, one cascade forwards (k top) and the other backwards (k bottom) (note that the pictures illustrate the effects of both cascades individually). After superimposing the effects of the two cascades, one can see that this reversible reaction becomes inviable. Thus, metabolites surrounding the reaction may become inviable as well, depending on their degrees. It is also a potential trigger, as in cases i-j.

However, amplification may have a very large impact when it occurs. For instance, pyruvate (a product of glucose metabolism and a key intersection in several metabolic pathways) provides energy by fermentation. This process reduces pyruvate into lactate, a reaction that does not trigger any black out cascade when it fails, so $d_r = 1$. At the same time, pyruvate can also be decarboxylated to produce acetyl groups, the building blocks of a large number of molecules that are synthesized in cells. The failure of the first reaction in such pathway triggers a cascade of length $d_{r'} = 3$. In contrast, the simultaneous failure of both the fermentation and the reduction of pyruvate induces a large cascade of size $d_{rr'} = 36$, most likely lethal. As a biological explanation, one could argue that both processes are strongly interdependent to maintain the oxidation-reduction balance when fermentation is in action.

Collateral effects offer the clue to understand this amplification phenomenon. In parallel to rendering non-operational some reactions and their corresponding metabolites, a cascade can reduce the connectivity and increase the branching ratio of other viable metabolites in its influence area. When stricken by the propagation front of a second cascade, these metabolites are susceptible of becoming inviable, further spreading the failure wave. In this way, interference is a necessary but not a sufficient condition for amplification, and a large amplification can be possible even when there is no overlap and the interference

between the individual cascades is small. To predict which pairs will trigger amplification, one must focus on metabolites in the interference of the influence areas of the two individual cascades. Those metabolites that remain viable after each individual cascade but become inviable when the two effects are superposed will produce amplification, propagating the double cascade to new reactions. In Figures 3.6f-k, the connectivity structure of all interference cascade propagator motifs responsible for amplification is provided.

## 3.3    Impact of gene knockouts in metabolic structure

Reaction failures are usually associated to the disruption of an enzyme due to knockout, inhibition, or deleterious mutation of the corresponding gene. iIn *M. pneumoniae*, enzyme multi-functionality and gene essentiality are higher as compared to other prokaryotic bacteria, so gene malfunctioning can potentially produce an acuter stress response at the level of metabolism. To address this issue, the metabolic network of *M. pneumoniae* is coupled to its gene co-expression network through the activity of enzymes, and knockouts of individual genes and clusters of co-expressed genes are performed. Inherent to this analysis is the potential occurrence of individual, double, or multiple cascades simultaneously. Multiple knockouts are algorithmically handled as an obvious extension of the previously considered situation of pair cascades.

The genome of *M. pneumoniae* [172] comprises 688 genes, 140 of which have a metabolic function. Except for one spontaneous reaction and 20 reactions with unknown regulation, these metabolic genes codify 142 enzymes that catalyze reactions in the metabolic network of this organism.

### 3.3.1    Metabolic effects of individual mutations

Individual metabolic gene knockouts or mutations inhibit the production of catalytic enzymes and induce black outs of reactions propagating in the metabolic network as a failure cascade (see Figure 3.7). From existing data, 71% of the 140 metabolic genes in *M. pneumoniae* have a one-to-one relation with reactions, and 21% of the genes regulate multiple reactions. Seldom the same reaction may be individually regulated by different enzymes produced by different genes, which happens for only four non-damaging reactions. More often, several genes are necessary to regulate the activity of a single reaction through an enzymatic complex. Twelve complexes codified by 26% of genes regulate the activity of 10% of metabolic reactions in *M. pneumoniae*. The removal of any of the genes involved in a complex is expected to cause the inactivation of the reaction controlled by the complex, which in principle may increase vulnerability. However, it can be observed that almost all complexes are associated to low damage reactions, which indicates a certain degree of structural robustness.

To study the metabolic effects of individual gene mutations, the knockout of all reactions associated to the gene under consideration are simulated. As explained, most often this corresponds to one single reaction but sometimes multiple reactions are removed simultaneously. The first observation is that metabolic genes affecting vulnerable reactions trigger large failure cascades. More interestingly, genes with large associated damages in metabolism turn out to be essential or conditionally essential for *M. pneumoniae* (see Table 3.1), with a unique exception discussed below. The classification given in Reference [120] is used, where essentiality is defined according to the measured metabolic map and the definition of a minimal medium which allows *M. pneumoniae* to grow. Essential genes are those that are required for the survival of the organism,

FIGURE 3.7: Left: Scheme of genes connected to reactions. Direct connections between genes and reactions are shown, but notice that connections between genes and reactions can be done only due to the existence of enzymes. Right: effect of how performing a knockout of a gene, labeled as *g8*, spreads a cascade in the metabolic network. Squares denote reactions, circles denote metabolites and triangles genes. Black nodes denote nodes that have become non-operational, whereas gray nodes are viable nodes that have reduced their connectivity due to the effect of the cascade.

meaning that the products of the reactions that they control are essential for life and cannot be produced by alternative pathways, while conditional means that essentiality depends on the media composition available.

In fact, all conditionally essential genes with the potential of producing high damage in the metabolism of *M. pneumoniae* have been found to have an essential orthologue (essentiality determined by loss-of-function experiments) in *Mycoplasma genitalium* [189], a comparable genome-reduced bacterium. The only exception to essentiality in Table 3.1 is gene MPN062, considered as non-essential in Reference [120], while in this study it triggers a large failure cascade and so it can be classified as a vulnerable target for metabolic structure. Its damaging potential can be explained by the fact that each of the four reactions controlled by the gene has a contribution that, although not extremely high individually, adds to the total damage and interferes to produce amplification. Therefore, MPN062 can be proposed as an important gene for metabolic function in *M.*

| Gene | Essentiality | Damage | Reactions |
|---|---|---|---|
| **MPN429** | yes | 49 | 4 (1,1,1,1) |
| MPN606 | yes | 32 | 1 (32) |
| MPN628 | yes | 32 | 1 (32) |
| MPN017 | yes | 25 | 3 (14,1,9) |
| MPN303 | yes | 18 | 8 (1,1,1,1,8,1,2,3) |
| *MPN062* | *no* | 17 | 4 (6,3,2,3) |
| MPN576 | cond | 16 | 2 (13,2) |
| **MPN005** | yes | 13 | 1 (13) |
| **MPN336** | yes | 13 | 3 (4,3,6) |
| **MPN354** | yes | 13 | 1 (13) |
| **MPN627** | yes | 11 | 1(11) |
| MPN066 | yes | 9 | 4 (1,1,2,5) |
| **MPN240** | cond | 9 | 1 (9) |
| MPN299 | cond | 9 | 1 (9) |
| MPN322 ⎫ | cond | 9 | 4(1,1,2,1) |
| MPN323 ⎬ | cond | 9 | 4 (1,1,2,1) |
| MPN324 ⎭ | cond | 9 | 4 (1,1,2,1) |
| **MPN034** ⎫ | yes | 7 | 4 (1,1,2,3) |
| **MPN378** ⎭ | yes | 7 | 4 (1,1,2,3) |

TABLE 3.1: Largest structural damages produced in metabolism by gene knockouts and correspondence with gene essentiality as given in Reference [172]. Damage in metabolic structure caused by gene knockout (third column) is measured in number of deleted reactions. In the fourth column, the number of reactions regulated by the corresponding gene is given, and in parentheses the damage associated to each of these reactions is also given. Genes in monocomponent clusters are highlighted in boldface, and braces are used to denote genes that form complexes. Note that the complex at the end of the list is not detected by any of the three clustering procedures. Finally, gene MPN062 is the only one in the table annotated as non-essential although it is associated to a large failure cascade.

*pneumoniae*, a conjecture that is supported by the essentiality of its orthologue in *M. genitalium* [189].

Another interesting case is essential gene MPN429, whose knockout triggers the largest cascade in *M. pneumoniae*. Each of the four affected reactions in the Glycolysis pathway is not able to propagate a cascade individually. However, when they all are removed simultaneously, the strongest amplification effect

is observed. The biochemical explanation is that the non-linear interaction of the cascades stops the production of phosphoenolpyruvate, which disrupts the synthesis of ATP, a circumstance particularly harming to the organism.

The obtained results of the study of structural cascades to predict gene essentiality in *M. pneumoniae* is in agreement with the gene essentiality computed in Reference [56].

### 3.3.2   Metabolic effects of knocking out gene co-expression clusters

Groups of co-expressed genes in *M. pneumoniae* can be identified from gene expression data under different conditions [190–192], which reveals a complex gene regulatory machinery [172]. The functional deactivation of these clusters might be produced by the failure of common regulatory elements and important damage could be transmitted to metabolism.

In this subsection, results on the effects on the metabolic structure of *M. pneumoniae* by suppressing gene co-expression clusters are shown. Information about gene expression is provided in Reference [172]. Correlations in the expression of genes were measured from tilling arrays under 62 different environmental conditions. This matrix of correlations between the expression levels of pairs of genes gives a fully connected network where the link between two genes carries a weight ranging from -1 to 1. This gene correlation matrix can be coupled to the metabolic network of *M. pneumoniae* through the activity of enzymes to produce a multilevel network representation.

To detect gene co-expression clusters, three different strategies -distance hierarchical clustering, Infomap, and recursive percolation (see Chapter 2, Section 2.1.4)-

FIGURE 3.8: Pictures of co-expressed genes and distribution of sizes of gene clusters. Left: Groups of co-expressed clusters regulating a metabolic network. Right: Distribution of sizes of the clusters obtained using distance hierarchical clustering (blue), Infomap (red) and recursive percolation (green).

are applied to the gene expression correlation matrix in order to discount biases introduced by the specifications of the community detection method (in Supplementary Table C3-4 it is possible to see a table showing which cluster each gene belongs to depending on the method used to detect the clusters). The distributions of sizes of the obtained clusters with the three strategies are shown in the right panel of Figure 3.8, where it is indeed possible to see a certain degree of similarity between the distribution of sizes of the clusters obtained using the three different methods.

The comparative analysis of the detected clusters of genes showed that, although the partitions found by each algorithm may differ in their composition and in the maximum size of the clusters, there are preserved commonalities independently of the method. One of them is that all methods are able to detect seven of the twelve complexes, since the related genes always appear classified in the same cluster. Another remark is that, as explained in the previous paragraph, the three detection methods result in qualitatively similar power-law-like cluster size distributions (see Figure 3.8 right), with most clusters having small size while

some are relatively big. Interestingly, genes related to high damage spreading reactions are secluded into mono-component clusters. To be more precise, eight of the nineteen genes in Table 3.1 are recognized by all three methods as having an expression profile that is not correlated to other gene activity levels. This is surprising since, in principle, high-damage genes might be expected to be co-regulated with other genes, as influencing big parts of metabolism usually requires coordinated gene activity. The fact that these genes appear isolated pinpoints them as potentially important metabolic regulator targets, since the alteration of only one gene may affect a large number of metabolic reactions. In any case, the lack of co-regulation of genes related to high damage spreading reactions is again an indication that the structural organization of the organism has evolved towards protecting the system against multiple failures.

Taking averages for equally sized clusters, it can be found that knockouts of co-expression clusters produce a damage on metabolic structure that increases with the number of affected metabolic genes, except when most metabolic genes in a cluster codify an enzymatic complex regulating one reaction (see Figure 3.9, left panels). The damage produced by the failure of the cluster also increases with the number of associated reactions (right panels in Figure 3.9). In order to discount structural effects, these results are compared with those measured on DP randomized versions of the metabolic network of the genome-reduced bacterium. As evidenced in Figure 3.9, all cluster detection methods identify clusters that produce lower damages in the real metabolic network of *M. pneumoniae* as compared to the randomized network. This supports the idea that the regulatory machinery that controls the coupled-to-metabolism co-expression of genes has evolved towards robustness.

Finally, since the three cluster detection methods propose different forms of aggregating metabolic genes, it is relevant to consider whether cluster composition

FIGURE 3.9: Damages as a function of the number of metabolic genes and reaction failures in gene co-expression cluster knockouts. Clusters are defined according to three different methods: Hierarchical Clustering (HC), Infomap (I), and Recursive Percolation (RP). Results are compared with damages produced in randomized versions of the metabolic networks in order to discount structural effects. In each case, the solid black curve is the average over 100 realizations.

is relevant for failure propagation. As a null model, one can consider randomization restricted not to the network itself but to the specific gene metabolic composition, while maintaining the total number of metabolic genes in each cluster. It can be observed that such a reshuffling of metabolic genes in clusters has no relevant effect on the damages measured on the metabolic network (see Figure 3.10). This means that, surprisingly, the composition of the clusters is not as statistically relevant for metabolic vulnerability as the distribution of the cluster sizes itself. This feature, together with the large detected amount of mono-component clusters, point out to the existence of multiple levels of regulation, depending on experimental conditions and, at the same time, explains why genes controlling high damage spreading reactions operate preferentially under functional isolation as a metabolism protection mechanism.

## 3.4   Robustness vs regulation in metabolic networks

The null model used in the first part of this chapter, called degree-preserving randomization, does not account for the most basic physico-chemical constraints and may lead, in the case of metabolic networks, to consideration of reactions which are not mass (*i.e.*, stoichiometrically) balanced (which do not preserve the same type and number of atoms on the substrate and product sides). As a result, the randomized networks may not be chemically feasible. As an alternative, the null model called mass-balanced randomization [141] accounts for this issue (see Chapter 2, Section 2.1.7.2). It is worth stressing that this method preserves the degrees of reactions but not the degrees of metabolites (see Figure 3.11).

In this section, cascades originated by single reaction and pair of reaction failures in the original networks of the three bacteria are compared with those obtained from two null models: DP, already used in the previous section, and

FIGURE 3.10: Damage distributions as a function of the number of genes and reaction failures, similar to Figure 3.9, but now randomizing the specific genetic contents of each cluster while maintaining the total number of metabolic genes in each cluster. The size of the clusters are defined according to three different methods: Hierarchical Clustering (HC), Infomap (I), and Recursive Percolation (RP).

FIGURE 3.11: Comparison of the degrees of reactions and metabolites obtained by the two null models applied to *E. coli* network. In this representation, each point is a reaction or a metabolite with coordinates $(k_{real}, k_{randomized})$, where $k_{real}$ is the metabolite/reaction degree in the original network, and $k_{randomized}$ the corresponding degree in a randomized network. Points fall in the diagonal if degrees are preserved in the randomized networks. a) and b) MB randomization. This method gives networks in which the degrees of the reactions are preserved. However, degrees of metabolites are not conserved. c) and d) DP randomization. This method gives networks with preserved degrees of reactions. Degrees of metabolites are also preserved with DP randomization, however at the expense of violating mass balances of reactions.

MB randomization. As in the first part, K-S tests are used to statistically assess whether the null models are relevant to explain the resulting damage distributions in the original networks. The analysis reinforces the importance of choosing an appropriate null model according to the question at hand, since the null model ultimately affects the interpretation of the findings [186].

First, cascades triggered by individual removal of reactions are studied, each cascade having its associated damage $d_r$. When comparing the cumulative distributions $P(d_r' \geq d_r)$ of the damage $d_r$ produced by individual removal of reactions between the original and randomized networks (see Figure 3.12 left panels), it can be observed that the distributions of the original networks lie in between the distributions of the two null models. To check whether or not the cumulative probability distributions are significantly different in the original networks and in their randomized variants, K-S tests are performed (see Table 3.2), taking as the standard significance level $\alpha = 0.05$. The compared distributions are considered significantly different from the null models because their associated significance is smaller than 0.05, except for *M. pneumoniae*, whose distribution can be considered consistent with the DP model as seen in Section 3.2, probably due to its linearity. Both for *E. coli* and *S. aureus*, damages are smaller compared to their DP randomizations but larger when compared to their MB randomizations. Thus, the robustness of the analyzed networks cannot be explained by the distribution of degrees or by basic physical constraints. For the DP null model, this finding indicates that robustness is positively influenced by factors other than the degrees. The results from the MB null model suggest that, for all three organisms, evolutionary pressure leads to larger cascades of non-viable reactions as compared to those imposed by physico-chemical constraints, and thus lower robustness.

After performing single reaction removals, the same analysis for the removal of each possible pair of reactions is done. Similar to the single reaction case, the cumulative probability distributions $P(d_{rr'}' \geq d_{rr'})$ of the damage $d_{rr'}$ resulting from the knockout of two reactions is determined (see Figure 3.12). K-S tests with a standard significance level $\alpha = 0.05$ are again applied (see Table 3.2), finding that the distributions of the original networks are significantly

FIGURE 3.12: Distributions of damage caused by removal of reactions. a, c, e) Cumulative probability distributions for *M. pneumoniae* (blue), *S. aureus* (green), and *E. coli* (red). Averaged distributions over 100 randomizations of the original networks are shown for DP (dashed line) and MB randomization (continuous line). b, d, f) Damages caused by pairs of removal of reactions.

| Organism | SR | | PR | |
|---|---|---|---|---|
| | MB | DP | MB | DP |
| *M. pneumoniae* | 0.10/0.03 | 0.095/0.07 | 0.15/0 | 0.15/0 |
| *S. aureus* | 0.19/0 | 0.086/0.0002 | 0.27/0 | 0.14/0 |
| *E. coli* | 0.19/0 | $0.079/1.4 \cdot 10^{-11}$ | 0.21/0 | 0.13/0 |

TABLE 3.2: Kolmogorov-Smirnov tests for comparing single reaction (SR) and pairs of reactions (PR) failure cascades in the three metabolic networks with both randomization methods, MB and DP. The values of the K-S statistic / associated significance level are given.

different from those of both randomization methods. All organisms display in this case similar results, the distributions of the original networks lie again between the distributions of the two null models, and all of them can be considered inconsistent with both null models. Consequently, the observations for individual failures also hold for the failure of reaction pairs: robustness is positively influenced by factors other than degrees, but negatively influenced by evolutionary pressure.

The cascade algorithm produces larger damages in the original networks as compared to those in MB randomized networks, but smaller cascades as compared to those in DP randomized counterparts. A possible explanation is offered by the difference in global properties of the networks obtained from the two randomization methods [142]. DP randomization decreases the average path length and increases the clustering coefficient of the randomized network, increasing its small-world property. Consequently, such networks are more interconnected and, thus, a cascade may in principle propagate further in the network. The opposite holds for MB randomization, which increases the average path length while decreasing the clustering coefficient of the randomized network so that the spread of the damage is less likely. Although the average path length does not resemble the length of metabolic inter-conversion, the small-world property may still affect the impact of removal of reactions due to its functional importance.

It can also be pointed out that the principle of cascade propagation relies on violation of a structural precondition for a steady-state, namely that all metabolites can be produced and consumed in order to avoid their depletion or accumulation. However, the steady-state assumption is only meaningful for networks which satisfy fundamental physical principles. Therefore, the use of MB randomization, which guarantees preservation of mass balance, allows to discern whether the measured property is a result of basic physical principles, or, instead, whether it is affected by evolutionary pressure. Since the size of cascades in MB randomized networks is significantly lower than those in real networks, evolutionary pressure may indeed lead to larger cascades.

Consequently, this finding indicates that evolutionary pressure may favor lower robustness of metabolic networks with respect to the failure of reactions, seemingly contradicting the general requirement of robustness in biological systems. On the one hand, this finding may be a result of the evolutionary versatility of metabolic networks, which favors organisms that are able to evolve quickly, *i.e.*, by few modifications to their metabolic networks. On the other hand, it is worth stressing that a cascade may not only be interpreted as the harmful spreading of failure, but also as the ability to regulate metabolism by activating/deactivating reactions, *e.g.*, by transcriptional regulation [193]. Thus, large cascades, favored by evolutionary pressure, may point at the evolutionary requirement of regulating large parts of metabolism through the regulation of small sets of enzyme-coding genes. The ability to regulate the activity of metabolic reactions by deactivating competing reactions is a well-known principle of metabolism. These results thus indicate that evolutionary pressure may favor the ability of efficient metabolic regulation at the expense of robustness to reaction or gene knockouts, pointing at the necessary integration of trade-offs from various cellular functions.

## 3.5    Conclusions

Results obtained in this chapter demonstrate that when *E. coli* and *S. aureus* are subjected to reaction failures, their metabolic networks have a structure that minimizes the number of large cascades. In this way, the largest part of reaction failures lead to small cascades, resulting in a small damage for the metabolic network. Hence, one can conclude that these organisms have a robust metabolic network against reaction failures. *M. pneumoniae* exhibits network responses that are qualitatively comparable to *E. coli* or *S. aureus*, although it is found that it less robust against individual reaction removals with reactions more prone to trigger large metabolic failure cascades identified as key participants in the regulation of energy and fatty acid synthesis.

The concept of cascade amplification has been for the first time formulated and interpreted as a signature of the subtle non-linearities underlying the structure of complex networks. Specific scenarios in *M. pneumoniae* have been discussed. In addition, there is a motivation to assess the predicting power of the used formalism. In this sense, a predictor of damage propagation for single cascades, and structural motifs underlying amplified failure patterns in situations of concurrent spreading have been proposed.

On what respects to the analysis of single gene knockouts, it reveals its potentiality in capturing most of the scenarios of experimentally determined lethality for *M. pneumoniae*. Moreover, when clustered and knocked together new trends of the complex genomic regulation of the metabolism emerge. First, the distribution of cluster sizes seems to matter more than the actual composition of the clusters. This is connected to the fact that the regulation of high-damage genes tends to appear isolated from that of other genes, a kind of functional switch in metabolic networks that at the same time acts as a kind of genetic firewall.

The introduction of a randomization model that generates new realizations of the network which are mass balanced indicates that evolutionary pressure favors the ability of efficient metabolic regulation at the expense of robustness to gene knockouts. This is explained because it favors organisms to evolve quickly by little modifying their metabolic networks, and because a failure cascade can be interpreted as an ability to regulate metabolism by activating/deactivating reactions, apart from being interpreted as a harmful spreading of a failure.

## 3.6   Summary

- The metabolic networks of three bacteria, *M. pneumoniae*, *E. coli*, and *S. aureus*, have been found to be robust against reaction failures, although *M. pneumoniae* is less robust against individual reaction removals due to its simplicity [94].

- A predictor of damage propagation for cascades produced by single reaction failures and the structural motifs underlying amplified failure patterns have been proposed. It has been checked that the predictor successfully predicts damage without the need of computing cascades [94].

- The concept of cascade amplification has been formulated and interpreted as a signature of the subtle non-linearities underlying the structure of complex networks [94].

- The study of structural stress at the level of metabolic genes reveals its potentiality in capturing most of the scenarios of experimentally determined lethality for *M. pneumoniae* [94].

- The distribution of gene cluster sizes seems to matter more than the actual composition of the clusters in relation to failure propagation in the metabolic network [94].

- The studied organisms show a trade-off between robustness and efficient regulation of their metabolic networks [186].

# Chapter 4

# Effects of reaction knockouts on steady states of metabolism

The activity and essentiality of metabolic reactions of two model organisms, *Escherichia coli* and *Mycoplasma pneumoniae*, are studied using Flux Balance Analysis in different environments. In particular, synthetic lethal pairs correspond to combinations of active and active or inactive non-essential reactions whose simultaneous deletion causes cell death. Lethal knockouts of pairs of reactions separate in two different groups depending on whether the pair of reactions works as a backup or as a parallel use mechanism, the first corresponding to essential plasticity and the second to essential redundancy. Within this perspective, functional plasticity and redundancy are essential mechanisms underlying the ability to survive of metabolic networks.

The previous chapter reported the study of structural perturbations modeled by the removal of a reaction or a set of them and the application of ta viability

criterion at the structural level. This chapter goes from structure to function by using the technique called Flux Balance Analysis (FBA) (see Chapter 2, Section 2.2) to implement reaction knockouts. A FBA analysis goes beyond the structural characterization of a cascade triggered by a reaction knockout in the sense that FBA intrinsically assigns zero fluxes to all the reactions in the network that turn out to be non-viable, *i.e.*, that are not able to maintain a balanced steady state in a certain environmental condition. In addition, using FBA one can compute how the environment affects the fluxes of reactions in metabolic networks. In particular, FBA allows to compute the activities and essentialities of reactions at steady state (see Chapter 2, Section 2.2.3), and to study the concept of synthetic lethality and how it is related to concepts such as *plasticity* and *redundancy*.

The computation of the activity of reactions using FBA has permitted a better understanding of how metabolism adapts to environmental changes by means of modifications in the biochemical fluxes [12, 194]. Beyond the concept of activity, the study of essentiality can help to understand how metabolism adapts to an internal failure, analyzing the adaptation of the fluxes when one reaction is forced to be non-operative. In fact, the concept of essentiality has been studied extensively, from single reaction failures [105, 195, 196] to multiple failures [104, 197].

Plasticity and redundancy are large-scale strategies that offer the organism the ability to exhibit no or only mild phenotypic variation in front of environmental changes or upon malfunction of some of its parts. In particular, these mechanisms protect metabolism against the effects of single enzyme-coding gene mutations or reaction failures, the final outcome being that most metabolic genes result to be not essential for cell viability. However, some mutants fail when an additional gene is knocked out, so that specific pair combinations of non-essential metabolic

genes or reactions become essential for biomass formation. As an example, double mutants defective in the two different phosphoribosylglycinamide transformylases present in *Escherichia coli* -with catalytic action in purine biosynthesis and thus important as crucial components of DNA, RNA or ATP- require exogenously added purine for growth, while single knockout mutants do not result in purine auxotrophy [198].

These synthetic lethal (SL) combinations [4, 199–201] have recently attracted attention because of their utility for identifying functional associations between gene functions and, in the context of human genome, for the prospects of new targets in drug development. However, inviable synthetic lethal mutants are difficult to characterize experimentally despite the high-throughput techniques developed recently [202]. We are still far from a comprehensive empirical identification of all SL metabolic gene or reaction pairs in a particular organism [104], even more when considering different growth conditions. Metabolic screening based on computational methods becomes then a powerful complementary technique particularly suited for an exhaustive *in silico* prediction of SL pairs in high-quality genome-scale metabolic reconstructions.

This chapter unveils how functional plasticity and redundancy are essential systems-level mechanisms underlying the viability of metabolic networks. In previous works on cellular metabolism [194, 203], plasticity was some times associated to changes in the fluxes of reactions when an organism is shifted from one growth condition to another. Instead, here functional plasticity is discussed as the ability of reorganizing metabolic fluxes to maintain viability in response to reaction failures when the environment remains unchanged. On the other hand, functional redundancy applies to the simultaneous use of alternative fluxes in a given medium, even if some can completely or partially compensate for the other [101]. An exhaustive computational screening of SL reaction pairs is

performed in *E. coli* in glucose minimal medium and it is found that SL reaction pairs divide in two different groups depending on whether the SL interaction works as a backup or as a parallel use mechanism, the first corresponding to essential plasticity and the second to essential redundancy. When comparing the metabolisms of *E. coli* and *Mycoplasma pneumoniae*, one can find that the two organisms exhibit a large difference in the relative importance of plasticity and redundancy. In *E. coli*, the analysis of how pathways are entangled through SL pairs supports the view that redundancy SL pairs preferentially affect a single function or pathway [199]. In contrast -and in agreement with reported SL genetic interactions in yeast [204]- essential plasticity, which is the dominant class in *E. coli*, tends to be inter-pathway but concentrated and unveils Cell Envelope Biosynthesis as an essential backup for Membrane Lipid Metabolism. Finally, different environmental conditions are tested to explore the interplay between these two mechanisms in coessential reaction pairs. Knockouts of genes are not considered because approaching directly pairs of reactions without the scaffold of enzymes and genes allows to determine in a clean and systematic way the minimal combinations of reactions that turn out to be essential for an organism.

The contents of this chapter correspond to References [158, 205].

## 4.1    Activity and essentiality of single reactions of *E. coli* across media

This section summarizes the results of the study of how the activity and essentiality of reactions in the *i*JO1366 version of *E. coli* (see Section 2.3.1, its network can be seen in Supplementary Table C4-1) depend on the nutrient composition

of the environment [205]. To this end, the activity and essentiality of all active reactions in a set of minimal media are computed, and then, depending on their behavior on each environment, each reaction is classified according to four general categories. This study also allows to identify reactions as eventual candidates to form part of SL pairs.

A total number of 555 minimal media can be constructed as proposed in Chapter 2, Section 2.2.2.1, with a final number of 333 which allow growth for the *i*JO1366 version of *E. coli*. In addition to these minimal media, 10000 random media are also analyzed, of which 3707 give a non-zero growth. To construct these random media, one considers all metabolites present in the extracellular environment of *E. coli*. Then, one chooses the number of these metabolites that can act as nutrients. In this case, 90% of the total number of external metabolites are allowed to act as nutrients. Once the number of nutrients is selected, one chooses at random metabolites until one reaches the selected number of nutrients, and the lower bound of the exchange reactions of each metabolite is changed to a value of -10 mmol gDW$^{-1}$h$^{-1}$.

### 4.1.1   Quantifying activity and essentiality

The activity and essentiality of each reaction are computed in every medium, with the obvious constraint that essentiality is computed only if the reaction is active. On what follows, an explanation to compute the accumulated values of the activity and essentiality is given.

The activity $a_{i,j}$ of a reaction $i$ in a medium $j$ is defined as

$$a_{i,j} \equiv \begin{cases} 1 & \text{if } \nu_{i,j} > 0 \\ 0 & \text{if } \nu_{i,j} = 0 \end{cases} \tag{4.1}$$

where $\nu_{i,j}$ denotes the flux of reaction $i$ in medium $j$. To obtain a representative value of the activity, FBA calculations are performed in both minimal and random media. In addition, the activity is normalized by the number of media in which the calculations have been performed. Therefore, the activity of a reaction $i$ for a given set of media $n_{media}$ will be obtained according to

$$a_i \equiv \frac{1}{n_{media}} \sum_{j=1}^{n_{media}} a_{i,j} \tag{4.2}$$

with $0 \leq a_i \leq 1$.

Essentiality is defined on the subset of active reactions. To compute the essentiality of a particular reaction, the FBA growth rate is examined after removing the corresponding reaction. An expression of the essentiality of a reaction $i$ in a medium $j$ is given as

$$e_{i,j} \equiv \begin{cases} 0 & \text{if } \nu'_{g,j} > 0 \\ 1 & \text{if } \nu'_{g,j} = 0 \end{cases} \tag{4.3}$$

where $\nu_{g,j'}$ denotes the flux of the reaction of production of biomass in medium $j$ when reaction $i$ is constrained to have zero flux. Again, the results are averaged on several media and normalized by dividing by the number of media. In this way, the bounds of essentiality of a reaction lay between 0 and the corresponding activity of the reaction, $0 \leq e_i \leq a_i$,

$$e_i \equiv \frac{1}{n_{media}} \sum_{j=1}^{n_{media}} e_{i,j} \tag{4.4}$$

Another useful magnitude to be used later on is the ratio of media where reaction $i$ is essential with respect to the number of media where it is active. This measure is trivially computed according to $p_i = \frac{e_i}{a_i}$.

## 4.1.2   Characterization of the reactions

After computing FBA on all environments and for all mutants, essentiality *vs* activity is plotted for all reactions. All points must fall on the diagonal or under it. This plot is shown in Figure 4.1 for both minimal and random media.

Reactions can be classified into four categories:

1. **Essential whenever active reactions:** $0 < a_i = e_i$. They are essential in all media where they are active. These reactions lay on the diagonal of the aforementioned plot.

2. **Always active reactions:** $a_i = 1$, $0 < e_i < a_i$. They are always active and sometimes essential. These reactions are located in the opposite $y$ axis.

3. **Never essential reactions:** $0 < a_i < 1$, $e_i = 0$. They are never essential but sometimes active. These reaction are located in the $x$ axis.

4. **Partially essential reactions:** $0 < a_i < 1$, $0 < e_i < a_i$. They are essential only a fraction of times when they are active. These reactions are located inside the triangle formed by the diagonal, and the $y$ and $x$ axes.

To understand the obtained results, the study focuses on the different subnetworks obtained by filtering the complete original network according to the four basic explained categories. More precisely, these subnetworks are obtained by maintaining in the network only those reactions within the respective mentioned categories. Once these subnetworks are obtained, the number of connected components in the subnetwork are computed in order to know whether the selected subnetwork is fragmented or not. In particular, the giant connected component (GCC) and the strongly connected component (SCC) (see Chapter 2,

FIGURE 4.1: Representation of essentiality *vs* activity. a) Minimal media. b)
Random media. In both pictures the four different categories can be clearly
differentiated. Diagonal: *essential whenever active reactions*. Opposite $y$ axis:
*always active reactions*.  $x$ axis: *never essential reactions*.  Inside triangle:
*partially essential reactions*.

Section 2.1.5) of the subnetworks are computed. This study is done in order to
detect whether reactions within a specific type are responsible for the percolation
state of the network.

In Table 4.1, the statistics of active and essential reactions are summarized
together with values of the sizes of the connected components of the subnetworks.
Results correspond to the set of minimal media. A precise discussion of such
statistics is provided on what follows. Notice first that there are several reactions
which are strictly never active (902). This may be explained by the fact that
these computations have been done in minimal media, which may only activate
a few number of reactions needed to survive. In addition, it can be seen that
the complete network, which corresponds to values of activity $0 \leq a_i \leq 1$ and
essentiality $0 \leq e_i \leq a_i$, is constituted by a single GCC and that, in addition, it
has a large SCC, a typical situation in metabolic networks.

| Category | CC | $N_R$ |
|---|---|---|
| **Essential whenever active reactions** | **Total** | **665** |
| $a_i > 0$ | **GCC** | **611(91.9)** |
| $e_i = a_i$ | **SCC** | **409(66.4)** |
| Essential and active in some media | Total | 458 |
| $0 < a_i < 1$ | GCC | 409(89.3) |
| $e_i = a_i$ | SCC | 200(48.8) |
| Essential and active in all media | Total | 207 |
| $a_i = 1$ | GCC | 198(95.7) |
| $e_i = a_i$ | SCC | 174(86.1) |
| **Always active reactions** | **Total** | **37** |
| $a_i = 1$ | **GCC** | **34(91.9)** |
| $0 < e_i < a_i$ | **SCC** | **29(85.3)** |
| **Never essential reactions** | **Total** | **494** |
| $0 < a_i < 1$ | **GCC** | **494** |
| $e_i = 0$ | **SCC** | **476(96.4)** |
| **Partially essential reactions** | **Total** | **152** |
| $0 < a_i < 1$ | **GCC** | **145(95.4)** |
| $0 < e_i < a_i$ | **SCC** | **129(90.0)** |
| All reactions | Total | 2250 |
| $0 \leq a_i \leq 1$ | GCC | 2250 |
| $0 \leq e_i \leq a_i$ | SCC | 2076(92.0) |
| Never active | | |
| $a_i = 0$ | Total | 902 |
| $e_i = 0$ | | |

TABLE 4.1: Connected components and number of reactions $N_R$ in each subnetwork. Values in parentheses correspond to percentages. GCC percentages are computed by dividing the number of reactions in GCC relative to the total number of reactions in each category, whereas SCC percentages are computed by dividing the number of reactions in SCC relative to the number of reactions in the GCC subnetwork. Categories in bold correspond to the four basic categories mentioned in the text.

### 4.1.2.1    Essential whenever active reactions

A histogram of the values of the essentiality of the set of reactions *essential whenever active reactions* is shown in Figures 4.2a and b. A bimodal distribution is clearly displayed, with peaks at extreme values, $a_i = e_i \simeq 0$ and the other

at $a_i = e_i \simeq 1$. This means that there is a core of reactions that are always active and essential, as pointed out in Reference [194], and there is another set of reactions that are active very few times. This histogram coincides with the classification of the dependence of essentiality on the environment given in Reference [105]. The peak at values of activity $\sim 0$ corresponds to *environment-specific* essential reactions, whereas the peak at values of activity $\sim 1$ corresponds to *environment-general* essential reactions. The first region includes reactions whose deletion abolishes growth in specific environments, whereas the second one corresponds to reactions whose deletion suppresses growth in all environments.

A deeper characterization of this set of reactions is made in Table 4.1, which shows that this subnetwork has a large GCC with nearly 90% of the subnetwork. If reactions with activity-essentiality index of 1 are excluded from this subnetwork, another subset is obtained which has also a large GCC (89.3% of the total 458 reactions). This means that reactions with $a_i = e_i = 1$ are not responsible for the percolation state of the subnetwork of *essential whenever active reactions*, which points out to a large degree of redundancy.

An illustrative example of a particular reaction in this subcategory is *Potassium transport* (Ktex). This is a reaction which supplies the organism with potassium. This mineral salt is an important metabolite which influences the osmotic pressure through the cell membrane and also secures the propagation of electric impulses. Since these are important processes for organisms, this reaction is always active in order to secure that these processes are done properly and that the organism has a non-zero growth.

For random media (see Figure 4.3a and b), a similar behavior is obtained, with larger probabilities at the extrema, but an extra peak is obtained for low values of activity and essentiality. This means that there are some reactions which

are not as specific as *environment-specific* reactions because they are active and essential in more than one medium, loosing in this way their specificity. This makes sense for random media, since they contain many metabolites that may activate many reactions and, in this way, they lose the specificity of a minimal medium, which triggers only the reactions that allow an organism to grow on it.

### 4.1.2.2   Always active reactions

The set of reactions called *always active reactions* contains reactions with $a_i = 1$ and $0 < e_i < a_i$. In Figure 4.2c and d one can see that, in this case, there is a large peak at values of $e_i = 0$, meaning that the largest part of reactions with $a_i = 1$ have a value of $e_i = 0$. This means that, although these reactions are always active, they are not essential. One may be tempted to think that reactions with very low values of essentiality are useless and hence they could be removed from the network. Nevertheless, there are two reasons that justify their consideration.

- The first one is that these reactions may improve the life conditions of the organism. These reactions, in spite of being non-essential, might be active in order to increase the growth of the organism. As a matter of fact, to survive to hard conditions, an organism which is able to reproduce fast and efficiently will, with large probability, survive to unfriendly life conditions.

- The second one is more subtle. These reactions could form SL pairs. As an example, two reactions regulated by the genes *tktA* and *tktB*, which are in the peak at $e_i = 0$ and $a_i = 1$, form a synthetic lethal pair and the removal of these reactions would abolish growth by impeding the synthesis of nucleotides, nucleic acids, and aromatic amino acids. Briefly, the reactions

FIGURE 4.2: Histograms (fraction) and complementary cumulative probability distribution function of activity or essentiality (depending on the category) for minimal media. a, b) Essential whenever active reactions. c, d) Always active reactions. e, f) Never essential reactions. g, h) Partially essential reactions.

FIGURE 4.3: Histograms (fraction) and complementary cumulative probability distribution function of activity / essentiality (depending on the category) for random media. a, b) Essential whenever active reactions. c, d) Always active reactions. e, f) Never essential reactions. g, h) Partially essential reactions.

regulated by these mentioned genes, called TKT1 and TKT2 and both with a complete name of *Transketolase*, are reactions which belong to the Pentose Phosphate Pathway. This pathway generates NADPH and pentoses phosphate, the latter being a precursor used in the synthesis of nucleotides, nucleic acids and aromatic amino acids. Both reactions are always active to ensure a sufficient production of these mentioned products, and when one of these reactions is knocked out, the other reaction is in charge to restore this function.

In Table 4.3 one can see that, as in *essential whenever active reactions*, this category of *always active reactions* form a subnetwork with a GCC which is almost the full subnetwork with also a large SCC.

Note that for random media (see Figure 4.3c and d), a similar trend to minimal media is obtained.

### 4.1.2.3    Never essential reactions

*Never essential reactions* have values of activity and essentiality which satisfy $e_i = 0$ and $0 < a_i < 1$. The histogram of the values of the activities for these reactions is shown in Figures 4.2e and f. A similar histogram to that corresponding to *always active reactions* is recovered again. This means that, not surprisingly, the largest part of *never essential reactions* are not much active. The individual removal of these reactions will leave the growth rate unaltered or only reduced. The existence of these reactions could be explained again in terms of improving the growth of the organism and, again, for the possibility of participating in SL pairs.

An example of a reaction of this kind is *Manganese transport in via permease (no H+)*, MN2tpp, a reaction which pumps manganese into the organism. Its non-essentiality comes from the fact that there exists an alternative reaction called *Manganese (Mn+2) transport in via proton symport (periplasm)*, MNt2pp, which also pumps manganese into the organism, but the latter uses a proton gradient to perform the transport.

In Table 4.1 one can see again the same trend as the other categories of reactions. This subnetwork contains a GCC that is almost the full subnetwork with a large SCC.

For random media, different results are obtained in this case (see Figure 4.3e and f). The largest peak is located at large values of activity, which means that there is a large set of reactions which are mainly active but never essential. The peak located slightly above 0.8 could appear due to the fact that the random media are in fact rich media. Hence, it is possible that a common set of metabolites activate the same reactions in many media. These reactions are responsible for the increase of the value of the flux of the biomass reaction.

### 4.1.2.4 Partially essential reactions

*Partially essential reactions* contain reactions with activity and essentiality values of $0 < a_i < 1$ and $0 < e_i < a_i$. Since these reactions have both values of essentiality and activity different from zero, the histogram is represented in terms of $\frac{e_i}{a_i}$ as shown in Figures 4.3g and h. The distribution is rather homogeneous, meaning that these two quantities may be largely uncorrelated, their ratio spanning the whole range of allowed values.

Table 4.1 shows again a large GCC containing a large SCC. Notice that this trend has been maintained for all categories of reactions.

Again, different results are obtained for random media (see Figure 4.3g and h). In this case, homogeneously distributed values as for minimal media are not obtained. Instead, the behavior resembles that of the *essential whenever active* subset, they are concentrated at low values and at a value of $\frac{e_i}{a_i} = 0.8$. This means that reactions are essential in fewer environments as compared to those in which they are active, showing again that activity does not imply essentiality.

## 4.2  SL pairs and plasticity and redundancy of metabolism

In metabolism, *synthetic lethality* arises when the individual failures of two reactions are not essential for cell growth but, contrarily, their simultaneous removal causes cell death [4, 199–201, 206, 207].

Synthetic lethality has been originally proposed in relation to genes [4, 199–202]. Its definition is that two genes are synthetic lethal when their individual knockout does not lead to the death of the organism but when both genes are removed simultaneously the organism is not able to overcome which leads to the death of organism (see Figure 4.4). Genes code for enzymes, and enzymes determine the kinetics of reactions and thus whether reactions take place in a feasible amount of time. Therefore, as for essentiality of individual genes and reactions, it is possible to extend the concept of synthetic lethality to reactions.

FBA is a powerful technique particularly suited for an exhaustive *in silico* prediction of SL pairs [104, 208]. Using FBA, a reaction pair deletion is annotated

FIGURE 4.4: Synthetic lethality schemes. a) Simplified scheme to illustrate the concept of synthetic lethality of two genes. b) Different possible organization of genes, enzymes and reactions in SL pairs.

as inviable, and so as a synthetic lethal, if the double mutant shows a no-growth phenotype.

This section presents the study of plasticity and redundancy of metabolism by directly computing the effects of double reaction knockouts, excluding those reactions that are individually essential in order to identify SL pairs. On what follows, a detailed analysis of the classification of identified SL reaction pairs into plasticity and redundancy subtypes in the *i*JO1366 version of *E. coli* and in the *i*JW145 version of *M. pneumoniae* (see Section 2.3.2, its network can be seen in Supplementary Table C4-8) is presented.

### 4.2.1   Classification of SL pairs

Some considerations are needed in relation to the space of reactions to be considered in forming potential SL pairs, the set of reactions that can be active but not essential in glucose minimal medium (see Chapter 2, Section 2.2.2.1). Different from the analysis in the previous section, in this section the study is primarily focused in one medium, not in a set of environments. In addition, the space of reactions to be considered is preliminary reduced using a method that we call "Biomass unconstrained Flux Variability Analysis", where Flux Variability Analysis (FVA) is applied irrespective of the level of attainable growth (see Chapter 2, Section 2.2.4). The final ensemble, formed of 1176 reactions in *E. coli* (see Supplementary Table C4-2) and 66 in *M. pneumoniae* (see Supplementary Table C4-9), is a subset of the original reconstruction that includes but that is not limited to the set of FBA active reactions under maximum growth constraint [106, 107].

An important remark is worth mentioning at this point. Some FBA computationally predicted SL pairs can be inconsistent with experimental data since they may contain at least one gene reported as essential *in vivo*. For *E. coli*, results are checked with essentiality information given in Reference [119]. Given the lack of direct evidence, results for *M. pneumoniae* are compared to a genome-wide transposon study in *Mycoplasma genitalium* given in Reference [189]. Since a functional ortholog in *M. genitalium* can be assigned to 128 metabolic genes in *i*JW145 (of a total of 145 genes), the essentiality of that ortholog can be associated to the corresponding gene in *M. pneumoniae*. The other 17 genes are assumed, similarly to Reference [56], to be not essential for growth due to their absence in *M. genitalium* and the high similarity of the metabolic networks of

both organisms [120]. Three cases may occur when FBA *in silico* results are compared to experimental essentiality:

- Both reactions in the *in silico* SL pair involve non-essential genes. In this case, the pair can be considered a potential synthetic lethal (see Figure 4.5a).

- One reaction involves a non-essential gene whereas the other is regulated by an essential one. In this case, if the essential gene regulates more than one reaction, one can consider that the *in silico* prediction is not an inconsistency (see Figure 4.5c), since the essentiality might refer to the rest of regulated reactions. Otherwise, the pair is considered as inconsistent with experimental data (see Figure 4.5b).

- Both reactions are regulated by essential genes. With the same argument as before, for the case that both reactions have associated genes which regulate more than one reaction, one can still consider the pair to be a potential synthetic lethal (see Figure 4.5d). The other possible combinations are considered inconsistent with empirical evidence (see Figure 4.5e).

Detected SL pairs associated to isoenzymes (see Figure 4.4b) and multifunctional enzymes (see Figure 4.4b) are also classified as inconsistencies. Isoenzymes (also known as isozymes) are enzymes that differ in amino acid sequence but that catalyze the same chemical reaction. In this way, a reaction can be catalized by two different enzymes in case that one of them becomes non-operative. Multifunctional enzymes are those that can catalyze more than one reaction at the same time. They are very important for organisms, since they are responsible of the catalysis of more than one reaction and their failure may cause important damage to organisms, since many reactions can become non-operative.

FIGURE 4.5: Schematic representation of the identification synthetic lethal inconsistencies.

### 4.2.2    Classification of SL reactions pairs into plasticity and redundancy

Of all reaction pair deletions in *E. coli*, 0.04% are *in silico* synthetic lethals and can be separated in two different subtypes. In the biggest group, having a relative size of 91%, one of the paired reactions is active in the medium under evaluation while the second reaction has no associated flux. The rest of SL reaction pairs are formed by two active reactions. Moreover, in accordance with results in Reference [104], it is found that inconsistencies correspond to 4% of all identified *in silico* SL pairs in *E. coli*.

Active-inactive coessential reaction pairs are referred to as *plasticity synthetic lethal* (PSL) pairs (see Figure 4.6a). 219 PSL reaction pairs are found in *E. coli* (see Supplementary Table C4-3), 86% of all diagnosed SL pairs in the *i*JO1366 version of *E. coli* (see Figure 4.7). Coessential inactive and active reactions in these pairs have zero and non-zero FBA flux respectively. When the active

FIGURE 4.6: Schematic representation of plasticity and redundancy synthetic lethality subtypes in metabolic networks. Metabolites are represented by circles and reactions by squares. Colored reactions with black arrows represent active reactions, whereas gray discontinuous lines are used for inactive reactions and metabolites and black for knockouts. The biomass production reaction is represented as a larger square with an associated flux $\nu_g$. When it turns to inactive, meaning that it has no associated flux, the organism is not able to grow. For simplicity, SL reaction pairs are illustrated in this figure as having a common metabolite, although this is not necessarily always the case. a) Initial configuration of a plasticity synthetic lethality reaction pair (reaction 2 active and reaction 3 inactive). b) Initial configuration of a redundancy synthetic lethality reaction pair (both reactions 2 and 3 active). c) Final configuration after knockout of reaction 2 in a) or b). d) Final configuration after knockout of reaction 3 in a) or b). e) Final configuration after simultaneous knockout of reactions 2 and 3 in a) or b).

reaction is removed from the metabolic network, fluxes reorganize such that the zero-flux reaction in the pair turns on as a backup of the removed reaction to ensure viability of the organism, even though the growth is generally lowered. In contrast, the level of growth is unperturbed when the inactive reaction is removed. As an example, the SL pair valine-pyruvate aminotransferase and valine transaminase form a PSL pair, the second reaction being the backup of the first, whose simultaneous knockout produces auxotrophic mutants requiring isoleucine to grow [209].

While the single activation of one of the reactions in a PSL pair is enough to ensure viability in front of single reaction disruptions, the parallel use of both coessential reactions may happen in other cases. *Redundancy synthetic lethal* (RSL) pairs are those in which both reactions are active and used in parallel (see Figure 4.6b). Of all SL reaction pairs in the *i*JO1366 version of *E. coli*, one finds that 15 (6%) are RSL (see Figure 4.7) (see Supplementary Table C4-3). Indeed, for 13 of the 15 RSL pairs the simultaneous use of both reactions increases fitness as compared to the situation when only one of the reactions is active (fitness is here understood as the maximal FBA biomass production rate for the organism). For the remaining two pairs growth remains unchanged. As an illustrative example of parallel use, oxygen transport combines with reactions in the ATP forming phase of Glycolysis to form RSL reaction pairs. If Oxydative Phosphorylation is blocked by the absence of oxygen and no alternative anaerobic process like Glycolysis is used, the energy metabolism of *E. coli* collapses and so the whole organism.

It is interesting to compute the shortest path length (see Chapter 2, Section 2.1.3) between reactions in SL pairs. It is found that network distances between reaction counterparts is slightly shorter in RSL pairs than in PSL pairs. Indeed, not all reactions in RSL or PSL pairs are directly connected through common metabolites. Direct connections happen for 60% and 38% of pairs respectively, while the rest can be separated by up to four other intermediate reactions so that the average shortest paths are 3.33 and 3.80, respectively (the average shortest path of the whole metabolic network is 5.02). Both essential plasticity and redundancy display overlap in reactions and associated genes. In the 15 RSL pairs, one can identify 17 different reactions controlled by 15 genes or gene complexes. The 219 PSL pairs involve 108 different reactions controlled by 61 genes or gene complexes.

FIGURE 4.7: Histogram for the four different categories of SL pairs in *E. coli* (left) and *M. pneumoniae* (right). RSL (blue): redundancy synthetic lethal pairs, RSL I (blue points): redundancy synthetic lethal pairs showing inconsistencies, PSL (orange): plasticity synthetic lethal pairs, PSL I (orange points): plasticity synthetic lethal pairs showing inconsistencies.

Although this analysis refers to reactions, specific signatures of enzyme activity may be worth stressing in connection with the analysis of coessential reaction pairs. For some of the identified SL pairs, direct experimental evidence is reported in the literature [209, 210]. Other experimental results support the buffering activity of reactions in some SL pairs, like in the aerobic/anaerobic synthesis of Heme [211, 212] and in the oxidative/non-oxidative working phases of the Pentose Phosphate Pathway [213]. Enzymatic degeneracy can be responsible for explaining two of the *in silico* detected RSL reaction pairs in *E. coli*. One RSL reaction pair, which produces isopentenyl diphosphate and its isomer dimethylallyl diphosphate -biosynthetic precursors of terpenes in *E. coli* that have the potential to serve as a basis for advanced biofuels [214]- is catalyzed by a single enzyme encoded by an essential gene (one-to-many enzyme multifunctionality (see Figure 4.6f)). Conversely, isoenzymes are encoded by different genes but can catalyze the same biochemical reactions. This many-to-one relationship

ensures that single deletion mutants lacking any of the genes encoding one of the isoenzymes can still be viable (see Figure 4.6f). This case happens in one RSL reaction pair catalyzed by isoenzymes encoded by nonessential genes associated to transketolase activity in the Pentose Phosphate Pathway [207].

Finally, a comparative study shows that coessential reaction pairs are 50 times more abundant in a much simpler genome-reduced organisms of increased linearity and reduced complexity such as *M. pneumoniae*. To perform the computations, the medium given in Table S5 of the Supplementary Information of Reference [56] is used. Constraints corresponding to the category called defined medium have been used, adding also D-ribose. 2% of all potential candidate reaction pairs in *M. pneumoniae* are synthetic lethals, vs solely the 0.04% in *E. coli*. Inconsistencies are also much more abundant relatively to *E. coli* and the balance of RSL vs PSL reaction pairs is also different (see Figure 4.7). Parallel use happens as frequently as the backup mechanism in coessential reactions, with 42% of all synthetic lethals being RSL pairs and 58% being PSL pairs (see Supplementary Table C4-10). As compared to results reported in Reference [56] for the synthetic lethality of genes, the used methodology detects the same 29 SL gene pairs and 15 new SL gene pairs. Since the 8 different genes in these pairs form two different complexes of four and three genes and one gene remains isolated, the 15 SL gene pairs reduce to just 2 SL reaction pairs (in the RSL and RSL I categories) sharing one of the reactions. The three reactions involved in the pairs are uptake of G3P (glycerol 3-phosphate), G3P oxidation to dihydroxyacetone phosphate, and uptake of orthophosphate. As reported in Reference [56], two independent routes through third-party pathways connect Glycolysis to Lipid Biosynthesis. The first two reactions above, R1 and R2, are involved in one of the routes, while the last reaction R3 influences the flux through the other route. When R1 and R3 or R2 and R3 are removed from

*i*JW145 model, the organism collapses due to the simultaneous failure of both routes.

### 4.2.3 Pathways entanglement

To investigate further the role of essential plasticity and redundancy in the global organization of metabolic networks, one can study the entanglement of biochemical pathways [7] through synthetic lethality. To do this, it is necessary to annotate all reactions in synthetic lethal pairs in terms of the standard metabolic pathway classification and to count the frequencies of dual pathways combinations both for plasticity and redundancy subtypes. In Figure 4.8, a visual summary of pathways entanglement through essential plasticity and redundancy is given. A graph representation is used, where pathways are linked whenever they participate together in a SL interaction (discontinuous lines represent redundancy SL interactions (see Figure 4.8c) and continuous arrows stand for plasticity SL interactions (see Figure 4.8d)). The frequency of a given pathway combination in RSL or PSL pairs defines the weight of the corresponding link.

In *E. coli* (see Figure 4.8a), one can observe that the synthetic lethality entanglement of pathways is in general very low, with the exception of the entanglement between Cell Envelope Biosysthesis and Membrane Lipid Metabolism. Redundancy SL pairs are basically intra-pathway, with only 3 of 15 being inter-pathway. Of all intra-pathway RSL pairs, 75% concentrate in the Pentose Phosphate pathway. Interestingly, the distribution of PSL reaction pairs avoids that of RSL pairs and, in contrast, tends to be inter-pathway. Of all PSL pairs, 67% include zero-flux reactions in Cell Envelope Biosysthesis and active reactions in the Membrane Lipid Metabolism, which unveils Cell Envelope Biosysthesis as an essential backup for Membrane Lipid Metabolism. Intra-pathway plasticity

FIGURE 4.8: Metabolic pathways entanglement through essential plasticity and redundancy in *E. coli* and *M. pneumoniae*. Nodes represent pathways and two pathways are joined by a link whenever there exists a SL pair containing one reaction in each pathway. Links corresponding to plasticity SL pairs are represented by green continuous arrows pointing from backup to active. Redundancy SL pairs are represented by discontinuous red lines. Labels correspond to the number of pairs which generate this combination of pathways, being thicker those links with more associated pairs. Self-loops correspond to SL pairs with both reactions in the same associated pathway. a) Pathways entanglement in *E. coli*. b) The same for *M. pneumoniae*. c) Scheme of how pathways entanglement is derived from RSL pairs. d) The same for PSL pairs.

coessentiality amounts to 29% of PSL pairs and is concentrated in Cofactor and Prosthetic Group and Cell Envelope Biosynthesis.

In *M. pneumoniae* (see Figure 4.8b) pathways entanglement through coessentiality of reactions is very low as in *E. coli*. Redundancy SL pairs can be intra-pathway (4 of 10) or inter-pathway (6 of 10) and PSL pairs are basically intra-pathway (12 of 14). Redundancy SL pairs denote the parallel use of reactions in Folate Metabolism and reactions in Nucleotide and Cofactor Metabolism. These two pathways, Folate and Nucleotide Metabolism, are also linked by 2 PSL pairs with non-essential reactions in Folate Metabolism and

essential reaction backups in Nucleotide Metabolism. Nucleotide Metabolism is also the pathway that concentrates most PSL pairs. Both RSL and PSL reaction pairs unveil Nucleotide and Folate Metabolism as the most entangled pathways. Taken together, these results indicate that Folate and Nucleotide Metabolic pathways preserve most rescue routes for reaction deletion events, in accordance with results in Reference [56]. The fact that the proportion of plasticity SL pairs is considerably decreased in *M. pneumoniae* as compared to *E. coli* could be indicative that, even if both plasticity and redundancy serve an important function in achieving viability, essential plasticity is a more sophisticated mechanism that requires a higher degree of functional organization, using at the same time less resources for maximum growth. At the same time, this can also be explained by the relative unchanging environmental conditions of *M. pneumoniae* in the lung, that could have induced the elimination of pathways not required in that medium [56]. This suggests that the adaptability of *M. pneumoniae* is very much reduced and its behavior could not be resilient to environmental changes.

## 4.2.4   Sensitivity to differences on environmental conditions

The last part of this section presents the analysis of plasticity and redundancy depending on the growth condition under evaluation. Environmental specificity of genes and reactions has been explored experimentally [119, 215, 216] and *in silico* [105] for different organisms and for random viable metabolic network samples, and it has also been extended to multiple knockouts in yeast [203, 208] and *E. coli* [217].

To investigate the sensitivity of SL reaction pairs in *E. coli* to changes in minimal medium composition, the study focuses on the 234 SL pairs detected in glucose

minimal medium and checks their classification over the 333 minimal media constructed as in the previous Section 4.1. Figure 4.9a shows the SL reaction pairs ranked by the fraction of media in which the pairs are synthetically lethal. For most pairs, coessentiality is not specific of an environment and only a minimal number of pairs shows environmental specificity. In particular, 53% coessential pairs are lethal in all media and 95% are lethal in more than 95% of environments. For each SL pair, one can count the number of media in which the SL pair is classified in the plasticity subtype as compared to the total number of media in which the pair is predicted to be coessential. Results are shown in Figure 4.9b. Nearly all SL pairs, 93%, are in the plasticity subclass for more than 93% of the media, while 12 pairs display a switching behavior between plasticity and redundancy. Noticeably, these pairs are intra-pathway and share common metabolites. Of them, three pairs contribute to biosynthesis of amino acids (Valine, Leucine, and Isoleucine Metabolism and Glycine and Serine Metabolism) and five pairs belong to the Pentose Phosphate Pathway and are related to the production of carbon backbones used in the synthesis of aromatic amino acids. Finally, five reaction pairs maintain in the redundancy subclass across all conditions in which are coessential.

The behavior of *E. coli* can be explored in an amino acid-enriched medium (see Section 2.2.2.2). Comparing with glucose minimal medium, the first observation is that 223 of the 234 SL pairs detected in glucose minimal medium are also found to be lethal in amino acid-enriched medium, which means that 11 pairs are rescued (see Supplementary Table C4-4). Of the 11 RSL pairs in amino acid-enriched medium, eight are conserved and three switch from plasticity in the minimal to redundancy in the amino acid-enriched medium. On the other hand, 208 of the 212 PSL pairs are conserved and four change from redundancy in the minimal to plasticity in the amino acid-enriched medium.

FIGURE 4.9: Synthetic lethal reaction pairs in minimal media. a) Synthetic lethal reaction pairs ranked by the fraction of minimal media for which the pair is synthetically lethal. b) Synthetic lethal reaction pairs ranked by the fraction of minimal media in which the SL pairs are classified as essential plasticity and, complementary, as essential redundancy, provided that the pairs remain synthetically lethal.

Noticeably, only in one of the 208 conserved PSL pairs the pattern of activity changes from the reductase reaction producing dimethylallyl diphosphate to the isomerization of the less reactive isopentenyl pyrophosphate. In addition, a new set of 12195 lethal reaction pairs occurs, all of them involving however one essential reaction in glucose minimal medium that in amino acid-enriched medium becomes nonessential and instead takes part in SL pairs. Apart from those, no other new SL pairs are found.

In addition, this study also considers a rich medium. To construct this rich medium, a Luria-Bertani Broth (see Section 2.2.2.2) has been taken into consideration. In this rich medium, 13 new rescues are found when compared to the minimal medium (two new rescues as compared to the amino acid-enriched medium) and only three SL pairs change their plasticity/redundancy category (see Supplementary Table C4-5).

Plasticity and redundancy are still conserved when the growth maximization requirement is loosen. To implement the relaxation of the growth maximization requirement, again the glucose minimal medium is taken as a reference and the biomass production or the basic nutrients uptake rates are limited. In the first case, a FVA calculation is performed fixing the growth of the biomass to 30% of the maximal growth in glucose minimal medium and the exchange bounds of all nutrient uptakes are obtained. In comparison to the reference values, one can observe that the only metabolites which lower their maximal uptakes are the mineral salts (approximately reduced also to a 30%), while the uptake rates for the rest of compounds remained with the same bounds. Then, it is possible to perform FBA calculations in this overconstrained condition and compute SL pairs and their classification in RSL and PSL. If growth is relaxed in *E. coli* to 30% of its maximum value in glucose minimal medium by doing this, *in silico* essentiality of individual reactions does not change but activation of reactions increases. It is found, however, that the effect of this reorganization is indeed mild for plasticity and redundancy. All SL pairs are conserved and 82% of them maintain their PSL or RSL classification (see Supplementary Table C4-6). The absolute number of RSL pairs increases from 15 to 50 since four RSL pairs in the reference condition given by glucose minimal medium change to plasticity in the overconstrained medium, and at the same time 39 PSL pairs change to RSL. On the other hand, 180 SL pairs of 219 in the reference medium remain as

PSL pairs in the overconstrained condition. However, the pattern of activity in the pair has switched in 14% of the PSL pairs in this case, which indicates that the specific selection of the active reaction in a PSL pair can have an impact in the level of attainable growth.

If instead of limiting the uptake of mineral salts, the uptake rates of basic nutrients providing sources of carbon, nitrogen, phosphorus and sulfur are overconstrained, the effect is even softer and indeed negligible as compared to the reference medium. To do this overcostraining, it is necessary to first apply FVA setting the value of biomass growth to the maximum in glucose minimal medium in order to determine an upper uptake limit. Then, the maximum rate uptake of glucose and of the other three basic compounds is constrained to 30% of the maximum possible values while keeping the reference values for the mineral salts. FBA is then applied in the resulting overconstrained medium and SL pairs and their classification in RSL and PSL are computed. The number of active reactions only increases in three, the essentiality of individual reactions and SL pairs is conserved, and 99% of them maintain their PSL or RSL classification with only three SL pairs that switch class and only one PSL pair that changes the active reaction (see Supplementary Table C4-7).

In both overconstrained modifications of the glucose minimal medium, the number of active reactions changed from 412 to 490 in the mineral salts overconstrained medium and to 415 in the basic nutrients overconstrained medium. It is important to stress that, in both cases, the essentiality of individual reactions and all SL pairs were conserved (except for two new RSL inconsistencies in the basic nutrients overconstrained medium).

## 4.3   Conclusions

The first part of this chapter presents the results of a study of the activity and the essentiality of single reactions of *E. coli* in different environments. Reactions can be divided in four categories depending on their values of essentiality and activity. By doing this, one recovers *environment-specific* and *environment-general* reactions as given in Reference [105]. These correspond to the bimodal behavior in the category called *essential whenever active reactions*. Given their importance, these reactions can be selected as drug targets since they are fundamental constituents of the metabolism of *E. coli*. Another important feature that can be observed is the fact that some reactions, in spite of being never essential, are always active, which may favor an increase of the growth rate of the organism and the robustness of metabolism through redundancy. The categories of reactions which show this behavior are *always active reactions* and *never essential reactions*. The last feature that one can extract from the category *partially essential reactions* is that active reactions are not necessarily essential. Therefore, in general extrapolating activity to essentiality is not correct.

Beyond the essentiality of single reactions, SL pairs are complex functional combinations of reactions (or genes) that denote at the same time both vulnerability in front of double deletions and robustness in front of the failures of any of the two counterparts. Working at the level of reactions, synthetic lethality is meditated by two different mechanisms, essential plasticity and essential redundancy, depending on whether one reaction is active for maximum growth in the medium under consideration and the second inactive, or in contrast both reactions have non-zero flux. Plasticity sets up as a sophisticated backup mechanism (mainly inter-pathway in *E. coli*) that is able to reorganize metabolic fluxes turning on inactive reactions when coessential counterparts are removed

in order to maintain viability in a specific medium. Redundancy corresponds to a simultaneous use of different flux channels (mainly intra-pathway in *E. coli*) that ensures viability and besides increases fitness. Apparently, it could seem extremely improbable that the removal of an inactive reaction together with a non-essential active one, like in PSL pairs, could have any lethal effect on an organism. However, it is found that this situation is indeed overwhelmingly dominant in *E. coli* as compared to redundancy synthetic lethality, and it is still relatively frequent even in a less complex organism like *M. pneumoniae*.

Synthetic lethal mutations have been assumed to affect a single function or pathway [199], which reinforces the idea that pathways act as autonomous self-contained functional subsystems. In contrast, other investigations in yeast [204] report that synthetic-lethal genetic interactions are approximately three and a half times as likely to span pairs of pathways than to occur within pathways. In this chapter, it is found that RSL pairs in *E. coli* are predominantly intra-pathway while PSL pairs, more abundant, tend to be inter-pathway although concentrated in the entanglement of just two pathways, Cell Envelope Biosynthesis and Membrane Lipid Metabolism. The comparative study here shows that although pathways entanglement through coessentiality of reactions is low in both organisms, RSL pairs in *M. pneumoniae* can be intra-pathway or inter-pathway, linking Folate Metabolism and Nucleotide and Cofactor Metabolism, and PSL pairs are basically intra-pathway and located in Nucleotide Metabolism. Taken together, these results indicate that Folate and Nucleotide Metabolic pathways preserve most rescue routes for reaction deletion events, in accordance with results in Reference [56]. The fact that the proportion of PSL pairs is considerably decreased in *M. pneumoniae* as compared to *E. coli* could be indicative that, even if both plasticity and redundancy serve an important function in achieving viability, essential plasticity is a more sophisticated mechanism that

requires a higher degree of functional organization, using at the same time less resources for maximum growth. At the same time, this can also be explained by the relative unchanging environmental conditions of *M. pneumoniae* in the lung, that could have induced the elimination of pathways not required in that medium [56]. This suggests that the adaptability of *M. pneumoniae* is very much reduced and its behavior could not be resilient to environmental changes.

It has also been found that SL reaction pairs and their subdivision in plasticity and redundancy are highly conserved independently of the composition of the minimal medium that acts as environmental condition for growth, and even when this environment is enriched with nonessential compounds or overconstrained to decrease the maximum biomass production. These environment unspecific SL pairs can thus be selected as potential drug targets operative regardless of the chemical environment of the cell.

## 4.4   Summary

- There exists a set of reactions, and thus enzymes and genes, that must be always active in order to ensure the viability of an organism [205].

- Non-essential reactions deserve special attention for two causes: their role as growth enhancers and for their potential participation in synthetic lethal pairs [205].

- Synthetic lethality is meditated by two different mechanisms, essential plasticity and essential redundancy, depending on whether one reaction is active for maximum growth in the medium under consideration and the second inactive, or in contrast both reactions have non-zero flux [158].

- Plasticity sets up as a sophisticated backup mechanism that is able to reorganize metabolic fluxes turning on inactive reactions when coessential counterparts fail in order to maintain viability in a specific medium [158].

- Redundancy corresponds to a simultaneous use of different flux channels that ensures viability and besides increases fitness [158].

- Plasticity and redundancy are highly conserved independently of the composition of the minimal medium that acts as environmental condition for growth, and even when this environment is enriched with nonessential compounds or overconstrained to decrease the maximum biomass production [158].

# Chapter 5

# Detection of evolution and adaptation fingerprints in metabolic networks

Metabolic fluxes present an heterogeneity that can be exploited to construct metabolic backbones as reduced versions of metabolic networks. These backbones can be analyzed to extract important biological information. In this chapter, the disparity filter is applied to two organisms, *Escherichia coli* and *Mycoplasma pneumoniae*. Backbones offer information about long-term evolution since they contain the core of ancestral pathways related with energy obtainment optimized by evolution to maximize growth. At the same time, backbones unveil short-term adaptation capabilities to variable external stimuli.

The analysis of metabolic networks is a difficult task which requires a mixed use of tools that belong to Systems Biology, such as Flux Balance Analysis (FBA) (see Chapter 2, Section 2.2), and tools that belong to complex network science, such as modelization of metabolic networks as bipartite semidirected networks (see Chapter 2, Section 2.1.1). The combination of these approaches has enabled a huge step further towards the elucidation of important biological information hidden in the complexity of genome-scale metabolic reconstructions.

A useful tool in the endeavor of extracting useful biological information is the concept of backbone. Backbones maintain relevant biological information while displaying a substantially decreased number of interconnections and, hence, can provide accurate but reduced versions of the whole system. In particular, the work by Almaas *et al.* [12] introduced a filtering technique that selects the reaction that dominates the production or consumption of each metabolite such that a high-flux backbone can be retrieved. Although this method recovers pathways, the obtained backbones present a linear structure with very little interconnectivity and lack many of the features of real metabolic networks [5, 131].

Filtering approaches have also interested researchers working on networks in a more general context. A filtering method for weighted networks based on the disparity measure [218, 219] was developed in Reference [95]. This approach exploits the heterogeneity present in the intensity of interactions in real networks both at the global and local levels [220] to extract the dominant set of connections for each element. Typically, the obtained disparity backbones preserve almost all nodes in the initial network and a large fraction of the total weight, while reducing considerably the number of links that pass the filter. At the same time, disparity backbones preserve the heterogeneity and cutoff of the degree distribution, the level of clustering, and the bow-tie structure (see Chapter 2,

Section 2.1.5), and other characteristic features of the original networks [95]. Hence, the complex features of the original networks are preserved.

In this chapter, FBA is used to determine reaction fluxes and the disparity filter (see Appendix D) [95] is applied to extract the metabolic backbones of two organisms: *Escherichia coli* and *Mycoplasma pneumoniae.* These backbones are investigated for fingerprints of evolution and adaptation. One finds that the metabolic backbones of both organisms in minimal medium are mainly composed of a core of reactions belonging to ancient pathways. This means that the significant fluxes in these bacterial metabolic backbones are associated to reactions which have been present from the earliest stages of their life and still remain at present significant for biomass production. At the same time, external conditions modify the structure of the backbones, which allows to identify pathways that are more sensitive to changes in the environment and so prone to short-term adaptation.

The contents of this chapter correspond to Reference [221].

## 5.1 Identification of the disparity backbones of metabolic networks

FBA is used to compute the fluxes of the reactions composing the metabolic networks. These fluxes are treated as weights by the disparity filter. In this chapter, the *i*JO1366 version of *E. coli* K-12 MG1655 and the *i*JW145 version of *M. pneumoniae* are used (see Chapter 2, Section 2.3, their networks can be seen in Supplementary Tables C5-1 and C5-3). FBA calculations are performed in glucose minimal medium with a maximum uptake of glucose limited to 10 mmol gDW$^{-1}$ h$^{-1}$ for *E. coli* and 7.37 mmol gDW$^{-1}$ h$^{-1}$ for *M. pneumoniae*

(D-ribose is added to enrich the medium for *M. pneumoniae*). Once the fluxes are computed, the disparity filter is applied to the incoming and outgoing connections of each metabolite, such that only those links to reactions which concentrate a significant amount of flux are selected for the backbone (see Appendix D). The connectivity structure (see Chapter 2, Section 2.1.5) of the obtained backbones is analyzed from an evolutionary perspective, and additional media are considered to analyze environmental sensitivity (see Chapter 2, Section 2.2.2).

An important feature of flux solutions obtained using FBA is the heterogeneity of the flux distributions. In the same state, fluxes of reactions can span several orders of magnitude [12, 222]. To check this statement, the probability distribution functions of the obtained fluxes are shown (disregarding zero-flux reactions) in the insets of Figure 5.1b and c, confirming that, indeed, fluxes show an heterogeneous distribution at the global level. The set of metabolites in non-zero flux reactions is considerably reduced from the original total number, from 1805 to 445 metabolites in *E. coli*, and from 266 to 227 metabolites in *M. pneumoniae*. To characterize the existence of such heterogeneity also at the local level, the disparity measure [12, 95] is calculated for every metabolite $i$, $\Upsilon_i(k) = k \sum_{\forall j \in \Gamma(i)} (\nu_j / \sum_j \nu_j)^2$ (see Appendix D), accounting for the $k$ reactions $j$ in its neighboring set $j \in \Gamma(i)$ with corresponding fluxes $\nu_j$. Figures 5.1b and 5.1c display the disparity values for all metabolites as a function of their incoming and outgoing degree in *E. coli* and *M. pneumoniae*, respectively. The shadowed areas correspond to values compatible with a random distribution of fluxes among the reactions producing or consuming a metabolite and help to discount local heterogeneities produced by random fluctuations (see caption of Figure 5.1). As shown, most metabolites present flux disparity values that cannot be explained by random fluctuations meaning that the local distribution

of the fluxes of reactions associated to metabolites is significantly heteroge-
neous. One concludes then that the disparity filter will be able to efficiently
extract a backbone with the most relevant connections for both organisms, while
preserving the characteristic features of metabolism as a complex network.

Briefly, the disparity filter works by comparing weights of links with a random
assignment. The filter preserves a link in the backbone if the probability that
its normalized weight $\alpha_{ij}$ is compatible with the random assignment ($p$-value)
is smaller than a chosen threshold $\alpha$ which determines the filtering intensity
(see Appendix D). One proceeds to filter the metabolic networks with fluxes
of reactions as weights of the connections between metabolites and reactions.
For each metabolite $i$, the $\alpha_{ij}$ of each connection between metabolite $i$ and
its neighboring reactions $j$ is computed and the obtained $p$-value is compared
with the significance level $\alpha$. The disparity filter can be adjusted by tuning
this threshold to observe how the metabolic networks of both *E. coli* and *M.
pneumonaie* are reduced as $\alpha$ is decreased from 1 to 0, both of them included,
$\alpha = 1$ meaning the complete network. Notice that, after applying the filter,
one recovers a bipartite representation of the metabolic backbone. To avoid
working with stoichiometrically non-balanced reactions, the filtered bipartite
representation is transformed into a one-mode projection of metabolites placing
a directed link between two metabolites if there is a reaction whose flux is
simultaneously relevant for the consumption of one metabolite and for the
production of the other [12]. In this one-mode projected backbone, one computes
how many links $E$, nodes $N$ and total weight $W$ remain. These magnitudes
are normalized by dividing them by the corresponding values in the original
network, $E_T$, $N_T$, and $W_T$.

Figures 5.1d and e show the dependencies $N/N_T$ vs $E/E_T$, and $W/W_T$ vs $E/E_T$
in the associated insets, for the one-mode metabolic projections of the backbones

FIGURE 5.1: Scheme of the application of the disparity filter and measures of the heterogeneity of reaction fluxes in *E. coli* and *M. pneumoniae*. a) Scheme of the filtering method. Blue nodes are metabolites and green squares denote reactions. Incoming connections to metabolites are represented by red arrows, outgoing connections with blue arrows, and bidirectional connectiosns with dark yellow arrows. OMP denotes one-mode projection. b) Disparity measure as a function of incoming and outgoing degrees ($k$) in *E. coli*. The shadowed area corresponds to the average plus 2 standard deviations given by the null model, meaning that points which lie outside this are can be considered heterogeneous [95]. Inset: global distribution of fluxes of *E. coli*. c) Disparity measure as a function of IN and OUT degrees ($k$) for *M. pneumoniae* (*follows to next page*).

FIGURE 5.1: (*Follows from previous page*) Again, the shadowed area corresponds to the average plus 2 standard deviations given by the null model. Inset: global distribution of fluxes of *M. pneumoniae*. d) Fraction of nodes as a function of the fraction of links in *E. coli*. Inset: remaining weight as a function of the fraction of links in the network. e) Fraction of nodes as a function of the fraction of links in *M. pneumoniae*. Inset: remaining weight as a function of the fraction of links in the network.

of both *E. coli* and *M. pneumoniae*. While the filter can reduce considerably the fraction of links, the corresponding fraction of nodes is maintained at almost the original value. In addition, the total weight in the backbone only starts to drop appreciably after more than 50% of the links are removed. One takes the critical value $\alpha_c$ as the point where the fraction of nodes starts to decay (see Figures 5.1d and e). This critical value can be seen as an optimal point which reduces greatly the number of links in the network preserving at the same time most nodes and so as much biochemical and structural information as possible. The values are $\alpha_c = 0.21$ for *E. coli* and $\alpha_c = 0.37$ for *M. pneumoniae*.

## 5.2    Evolutionary signatures in the backbones of metabolites

The metabolic backbones of both *E. coli* and *M. pneumoniae* are constructed using the identified critical values for the significance level. The backbones retain all the 445 and 227 metabolites present in active reactions respectively. Next, one analyzes their structure in terms of connectedness. Metabolic networks have been found to display typical large-scale connectivity patterns of directed complex networks, called the bow-tie structure, with most reactions in a interconnected core, named the strongly connected component (SCC), together with in (IN) and out (OUT) components formed mainly by nodes directly connected to the

SCC component [131, 133] (see Chapter 2, Section 2.1.5). This is the case of the original metabolic networks of both organisms, whose SCCs contain the largest part of the metabolites and reactions of the network, and whose IN and OUT components are formed, respectively, by nutrients and waste metabolites.

Metabolites in the backbone of *E. coli* (it can be seen in Supplementary Tables C5-2) are arranged in a connected component of 178 nodes and several disconnected small components (51). Three different SCCs can be identified in the connected part of the backbone, each with 25%, 10%, and 6% of the nodes in the connected component (see Figure 5.2a). The two smallest SCCs are in the OUT component of the largest SCC. For the three of them, the IN and OUT components and tendrils are recovered. Metabolites corresponding to central compounds of metabolism are identified in these SCCs: protons, water, ATP, glutamate, phosphate, $NAD^+$, diphosphate, ADP and $FAD^+$. These metabolites are highly-connected metabolites even in the metabolic backbone, helping to preserve the same structural features of the complete metabolic network.

Since links in the metabolic backbone denote reactions transforming metabolites, it is interesting to annotate links with the pathway associated to the corresponding reaction. In this way, it is possible to count the composition of the three SCCs in terms of pathways. Starting with the largest SCC (see Figure 5.2a), one finds that the major contributions are Oxidative Phosphorylation (26%), Citric Acid Cycle (16%), Glycolysis/Gluconeogenesis (15%), Pentose Phosphate Pathway (9%), and Glutamate Metabolism (9%) (see Figure 5.2c). It has been demonstrated that these routes are ancient pathways that have been conserved through evolution. More precisely, Glycolysis and Pentose Phosphate Pathway take place without the need of enzymes in a mimetic Archean ocean [223]. Concerning the Citric Acid Cycle, it is also an ancient pathway that has evolved in order to achieve maximum ATP efficiency [224] by being coupled to Oxidative

FIGURE 5.2: SCCs of the backbone of metabolites and corresponding pathways. a) Connected component in the metabolic backbone of *E. coli*. The colors of the nodes depend on the component each node belongs to (yellow: SCC$_1$, green: IN component of SCC$_1$, red: OUT component of SCC$_1$, violet: tendrils of SCC$_1$, cyan: SCC$_2$, blue: SCC$_3$) (*follows to next page*).

FIGURE 5.2: (*Follows from previous page*) The color of the links, and its association given in the legend, depends on the functional categories given in Reference [119], where each category contains pathways that realize similar tasks. b) Connected component of the metabolic backbone of *M. pneumoniae*. The color of the nodes denote again the component each node belongs to (red: $SCC_1$, turquoise: IN component of $SCC_1$, green: OUT component of $SCC_1$, dark yellow: tendrils of $SCC_1$, violet: $SCC_2$). The color of the links, and its association given in the legend, depends on the pathway each reaction belongs to. c) Percentage of links in pathways for the largest SCC in the metabolic backbone of *E. coli*. d) The same for *M. pneumoniae*.

Phosporylation and Glycolysis [225], in addition to help the organism to decrease their quantity of reactive oxygen species by modulation of their participating metabolites [226]. Another pathway significantly present in the largest SCC is Glutamate Metabolism. Glutamate has been reported to be one of the oldest amino acids used in the earliest stages of life [227].

Links in the other two SCCs correspond also to reactions belonging to ancestral pathways. The second largest SCC contains links that belong mainly to Purine and Pyrimidine Biosynthesis (91%). Purines and pyrimidines serve as activated precursors of RNA and DNA, glycogen, etc. [228, 229], and it has been found that the synthesis of purines and pyrimidines was the first pathway involving enzyme-based metabolism [230]. Interestingly, the other contribution to this SCC is Glycine and Serine Metabolism. Glycine is a precursor of purines and pyrimidines. Pathways related to the third SCC are Membrane Lipid Metabolism (97%) and Cofactor and Prosthetic Group Biosynthesis (3%). Membrane Lipid Metabolism supplies the necessary lipids to generate the cell membrane needing the participation of the cofactor $FAD^+/FADH_2$. It has been shown that the pathways involved in lipid metabolism exhibit differences between different lineages in organisms [231], whereas pathways related to central metabolism are more conserved and are transversal [231].

When considering $\alpha$ values smaller than the critical one, implying that the filter is more restrictive and more heterogeneity is needed to overcome it, we observe that the smallest SCCs discussed above disappears. More precisely, it happens for a value of $\alpha = 0.19$. Decreasing even more the significance level to $\alpha = 0.15$ the SCC containing reactions in the Purine and Pyrimidine Biosynthesis pathway retains the 30% of the nodes for $\alpha_c = 0.21$, whereas the largest SCC still contains a 86%, showing the large resistance of this large core to lose nodes. At a value of $\alpha = 0.14$, the second SCC finally disappears and there only remains a single SCC, still preserving 82% of the nodes in it for $\alpha_c = 0.21$. Hence, energy metabolism shows a large resistance to get fragmented even though the filter becomes progressively more and more restrictive.

To contrast the obtained results in *E. coli*, the same analysis in *M. pneumoniae* is performed (its backbone can be seen in Supplementary Tables C5-4). Its critical value $\alpha_c$ is 0.37 (see Figure 5.1). The connected component of its metabolic backbone is shown in Figure 5.2b. It contains two SCCs, one of them being irrelevant with only two nodes (see Figure 5.2b). The relevant SCC contains 21% of the nodes in the connected component, and the largest part of its links are related also with energy metabolism as in *E. coli*. The dominant pathways in this core are Glycolysis and Pyruvate Metabolism (see Figure 5.2d). Along Glycolysis, Pyruvate Metabolism is also an ancestral pathway that was present in the earliest stages of life [35], when no oxygen was present in the early atmosphere.

## 5.3    The metabolic backbones of *E. coli* encode its short-term adaptation capabilities

The previous section analyzes the metabolic backbone of *E. coli* in glucose minimal medium in terms of the long-term evolution of the organism. In this section, the study is focused on how changes in the environment modify this backbone, which exposes short-term adaptation capabilities. First, FBA fluxes that maximize the growth rate of *E. coli* in the rich medium Luria-Bertani (LB) Broth [154, 158] are calculated. Afterwards, the disparity filter is applied to extract the metabolic backbone in this new environment, that is obtained for a significance level threshold $\alpha_c = 0.4$. This value is noticeably larger than $\alpha_c = 0.21$ identified for the glucose minimal medium. Interestingly, this rich medium activates 400 reactions, 11 less than in glucose minimal medium. Of them, 279 are active in both media, of which 247 have a larger flux in LB Broth. An analysis of the connected components in the metabolic backbone of *E. coli* in rich medium is also performed. One finds that it contains a large connected component with 449 metabolites and 60 small disconnected components. The connected component contains also three SCCs. However, two of them are tiny with only two nodes, whereas the largest one encloses 34% of the nodes in the connected component. Interestingly, the pathway contributing more reactions to this large SCC is Membrane Lipid Metabolism (see Figure 5.3a). This fact is in accordance with Reference [232], where the authors found that the expression of the genes which synthesize fatty acids was generally elevated in rich medium. Another important difference is the loss of prominence of Oxidative Phosphorylation and the Pentose Phosphate Pathways.

Next, the set of minimal media given in Reference [119] (see Chapter 2, Section 2.2.2.1) are considered, where different carbon, nitrogen, phosphorus and

FIGURE 5.3: Dependence of the distribution of pathways in the metabolic backbone of *E. coli* with the composition of environment. a) Histogram of the fraction of links belonging to each pathway (x axis) for the 333 minimal media (left) and in the rich medium (right). b) Probability distribution function of $\alpha_c$ for all minimal media. c) Probability distribution function of the fraction of links in the metabolic backbones for all minimal media. d) Histogram of weights of links in the metabolic superbackbone.

sulfur sources are alternated. For each minimal medium, $\alpha_c$ is scanned as in Figure 5.1b and c. In Figure 5.3b and c, one plots, respectively, the probability distribution functions of the tuned $\alpha_c$ values and of the fraction of links remaining in the metabolic backbones for all media. One finds that there is a characteristic value of these magnitudes with no outliers, meaning that the flux structure is very similar across media in spite of the difference in the composition of nutrients. The presence of these characteristic values of $\alpha_c$ and the retained fraction of links in the metabolic backbones motivates to merge all of them into a single merged metabolic backbone. The links in this superbackbone correspond

to reactions that passed the filter in any of the external media considered and are annotated with a weight that corresponds to the number of media in which the corresponding metabolic backbone contains the link. The histogram of the distribution of these weights is shown in Figure 5.3d, characterized by a clear bimodal behavior. One peak corresponds to links being common to all media, and the other corresponds to the most common situation of links specific to a few media.

An analysis of connectedness shows that this superbackbone contains a large connected component and 11 disconnected components. The connected component is composed by a large SCC with 43% of its nodes, in addition to three small SCCs containing only two nodes each. A pathway composition analysis in the large SCC indicates that, again, one obtains significantly different results from the glucose minimal medium (see Figure 5.3a). The most prominent pathway is Alternate Carbon Metabolism, in agreement with Reference [233], where the authors found that Alternate Carbon Metabolism is related to genes whose expression depends on external stimuli, particularly on alteration of carbon sources. It is also in agreement with results in Reference [234], where the authors hypothesize that Alternate Carbon Metabolism can adapt to different nutritional environments, and also with results in Reference [7], where Alternate Carbon Metabolism is found to be an important intermediate pathway in the network of pathways. The second most abundant pathway corresponds to Transport, Inner Membrane, which again is in agreement with Reference [233] and Reference [7]. It is a transversal pathway which is in charge of the transport of metabolites between periplasm and cytosol. Finally, if one retains links present at least in 25% of the minimal media, the network fragments into 40 components with the largest one containing five SCCs, which indicates that links with small weight,

*i.e.* links specific for a few media, have an important role in providing global connectivity to the superbackbone.

## 5.4   Conclusions

Identifying high-flux routes in metabolic networks has been useful in order to, for example, identify principal chains of metabolic transformations [12, 235, 236]. In this chapter, one goes beyond the mere identification of high-flux routes with metabolic pathways. Using a high-flux fluctuation analysis, it is possible to identify ancestral pathways and, on the other hand, pathways with capabilities to adapt to short-term external changes. At the core of the high-flux fluctuation analysis, a filtering tool which needs no *a priori* assumptions for the connectivity of the filtered subnetworks is used, but that produces reduced versions which are globally connected and retain the characteristic complex features of the original network. This procedure allows to extract a metabolic backbone which contains all relevant connections given a set of external nutrients, recovering both intra- and inter- pathway connections which can be understood as the superhighways of metabolism. Further, an evolutionary explanation can also be given for this identification of both intra- and inter- pathway connections since the cooperation between reaction inside and outside pathways implies that the overall performance of a cell will be improved due to a better and more efficient utilization of the available resources. This fact reinforces the idea that pathways are not isolated identities performing their tasks independently of others [7].

As stated in Reference [141], properties that originate from evolutionary pressure should not be observed in random networks. Due to the fact that the disparity filter identifies links that deviate from a random null model, it allows to identify those reactions for which evolutionary pressure has had a large incidence. Since

FBA flux solutions are used, in this chapter the effect of evolutionary pressure is understood to favor the maximization of the growth of the organism [49, 237, 238]. The evolutionary analysis of the metabolic backbones of the two considered organisms in minimal medium shows that their SCCs are composed by reactions that belong to ancient pathways. In *E. coli*, each SCC has different and definite metabolic functions. In both *E. coli* and *M. pneumoniae*, the largest SCC contains pathways related to energy metabolism, meaning that these organisms have evolved towards maximum efficiency in obtaining chemical energy, something very important in case of nutrient scarcity. A smaller SCC is responsible for the synthesis of purines and pyrimidines, vital for DNA / RNA synthesis. The third SCC corresponds to the metabolism of lipids, the most important constituents that compose the cell membrane. Two findings relating the two small SCCs deserve also special attention. Firstly, the two small SCCs are located in the OUT component of the large SCC. Secondly, as the filter becomes more restrictive, the small SCCs fragment, while the large SCC still maintains a large part of links and nodes. These features could be explained in terms of the functional requirements of the small SCCs. On the one side, they need chemical energy to perform their tasks and, on the other side, they need also basic building blocks. These tasks are performed in the large SCC by, for example, Glycolysis/Gluconeogenesis or the Citric Acid Cycle. Therefore, it suggests that those SCCs were added to the OUT component of the large SCC in later steps of evolution. A simpler organism, *M. pneumoniae*, has no other relevant SCCs apart from energy metabolism, as a result of its parasitism, which has led to the loss of many metabolic functions [56]. More precisely, in *M. pneumoniae* the Citric Acid Cycle and Oxidative Phosphorylation do not take place [56, 239], meaning that it must rely on organic acid fermentation to obtain energy. Moreover, changes in the growth rate greatly affect the fluxes

through Glycolysis and Pyruvate Metabolism [56].

The study of the dependence on the environment of the *E. coli* metabolic backbone allows to identify short-time adaptation capabilities. Regarding rich medium, one observes that the critical value of $\alpha$ is substantially different than the one in glucose minimal medium, suggesting that this enriched medium modifies significantly the flux structure compared to the glucose minimal medium. The bacterium in rich medium displays less active reactions than in glucose minimal medium since, in minimal medium, many reactions must be active in order to synthesize biosynthetic precursors that in the rich medium can be obtained from the environment, in agreement with Reference [232]. The pathway called Membrane Lipid Metabolism achieves a high relevancy, being the most abundant pathway in the largest SCC of the rich medium metabolic backbone. This happens because the instantaneous response of *E. coli* to this rich medium, which induces a large increase in the growth rate of the organism due to nutrient abundance, is to synthesize as much as membrane lipids as possible, since fast-growing cells must synthesize membrane components more rapidly to satisfy the high lipid demand to generate new cells [232]. The analysis of the adaptation of *E. coli* to 333 different minimal media shows that the distribution of fluxes is practically independent on the composition of the nutrients present in these environments, allowing to extract characteristic features that describe the backbones of the metabolic network independently of the environment. This permits the construction of a merged backbone that comprises all the links composing the metabolic backbone in each media. This leads to the identification of pathways whose associated reactions are more sensitive to changes in the environment, unveiling Alternate Carbon Metabolism as the pathway with more capabilities to respond to external stimuli, in accordance with previously reported results [233, 234].

The use of filtering methods usually imply a drastic reduction of the complexity of metabolic maps, which weakens the validity of potentially inferred conclusions. The application of the disparity filter based on a high-flux fluctuation analysis to produce metabolic backbones enables to reduce the system while maintaining all relevant interactions and so it becomes a useful tool to unveil sound biological information. For instance, the investigation of *E. coli* and *M. pneumoniae* revealed metabolic backbones in minimal medium mainly composed of a core of reactions belonging to ancient pathways, for which the effects of evolutionary pressure are higher, and unveiled pathways with high capacity to respond to external stimuli.

## 5.5    Summary

- The disparity filter is very efficient in order to compute metabolic backbones as reduced versions of metabolism which retain its complexity [221].

- The study of the bow-tie structure of the backbones in a glucose minimal medium reveals that pathways related with energy obtainment have an important evolutionary role in *E. coli* and *M. pneumoniae* [221].

- The study of the backbone of *E. coli* in rich medium identifies the pathway Membrane Lipid Metabolism as relevant for growth in the nutritionally rich medium, due to the necessity of large amounts of lipids to generate the cell membrane [221].

- The analysis of the superbackbone, constructed by merging all the backbones corresponding to different minimal media, recognizes the pathway Alternate Carbon Metabolism as the most relevant pathway to respond to external stimuli [221].

# Chapter 6

# Assessing FBA optimal states in the feasible flux phenotypic space

Optimal growth solutions can be confronted with the whole set of feasible flux phenotypes (FFP), which provides a reference map that helps to assess the likelihood of optimal and high-growth states and their extent of conformity with experimental results. In addition, FFP maps are able to uncover metabolic behaviors that are unreachable using models based on optimality principles. The information content of the full FFP space of metabolic states provides with an entire map to explore and evaluate metabolic behavior and capabilities, opening new avenues for biotechnological and biomedical applications.

The results presented in previous chapters required an extensive use of Flux Balance Analysis (FBA) (see Section 2.2) in order to extract backbones or to compute the effect of failures of reactions. If the removal of a reaction or a pair of reactions is not lethal for the organism, *i.e.*, the growth rate is not zero, there can exist many flux solutions for the organism to be alive. As it has been already explained, the FBA solution is a possible solution, the one which maximizes the growth rate. One may be tempted to ask where the solutions given by FBA lay in the whole space of possible flux solutions of a metabolic network. In this way, it will be possible to know whether the state given by FBA is indeed representative of the system or, on the contrary, it is not a representative solution of the flux space, this eventually being interpreted for example due to evolutionary effects.

FBA studies, like in the previous Chapter 4, reveal that metabolism is a dynamically regulated system that reorganizes to safeguard survival [49, 237], implying that metabolic phenotypes directly respond to environmental conditions. For instance, unicellular organisms can be stimulated to proliferate by controlling the abundance of nutrients available. In rich media, cells reproduce as quickly as possible by fermenting glucose, a process which produces high specific growth rates as well as large quantities of excess carbon in the form of ethanol and organic acids [240]. To survive the scarcity of nutrients during starvation periods, Glycolysis is hypothesized to switch to oxidative metabolism, which no longer maximizes the specific growth rate, but instead the ATP yield needed for cellular processes. Cells of multicellular organisms show similar metabolic phenotypes, relying primarily on Oxidative Phosphorylation when not stimulated to proliferate and changing to nonoxidative glycolytic metabolism during cell proliferation, even if this process -known in cancer cells as the Warburg effect [241, 242]- is much less efficient at the level of energy yield.

These metabolic phenotypes are captured by FBA. However, the identified solutions are frequently inconsistent with the biological reality since no single objective function describes successfully the variability of flux states under all environmental conditions [243, 244], and in fact the highest accuracy of FBA predictions is achieved whenever the most relevant objective function is tailored to particular environmental conditions according to the empirical evidence for a very specific metabolic phenotype. For instance, FBA requires either a rich medium or a manual limitation of the oxygen uptake to a physiological enzymatic limit to mimic the observed fermentation of glucose to formate, acetate, or ethanol typical of proliferative metabolism, while in minimal medium optimization of growth rate relies primarily on Oxidative Phosphorilation, which increases ATP production converting glucose to carbon dioxide, as in starvation metabolism. However, along optimal metabolic phenotypes, there is a whole space of possible states non-reachable by invoking optimality principles that prevent non-optimal biological states. Optimization of a biological function in the absence of *a priori* biological justification, which happens for instance under conditions for proliferative or starvation metabolism, may weaken *in silico* predictions.

In this chapter, optimal growth rate solutions are confronted to the whole set of feasible flux phenotypes (FFP) of core *Escherichia coli* metabolism in minimal medium, which provides a reference map that helps to assess the likelihood of optimal and high-growth states [245]. The whole set of feasible flux phenotypes is determined by mass-balance conditions and the bounds imposed on metabolites. Mathematically, it constitutes a convex finite polytope, and it is sampled using an algorithm called Hit-And-Run (HR) (see Appendix E) [246]. One can quantitatively and visually show that optimal growth flux phenotypes are eccentric with respect to the bulk of states, statistically represented by

the feasible flux phenotypic mean, which suggests that optimal phenotypes are uninformative about the more probable states, most of them low-growth rate. Feasible flux phenotypic space is proposed as a benchmark to calibrate the deviation of optimal phenotypes from experimental observations. Finally, the analysis of the entire high-biomass production region of the feasible flux phenotypic space unveils metabolic behaviors observed experimentally but unreachable by models based on optimality principles, like FBA, which forbid aerobic fermentation -a typical pathway utilization of proliferative metabolism- in minimal medium with unlimited oxygen uptake.

The contents of this chapter correspond to Reference [245].

## 6.1    Optimal growth is eccentric with respect to the full FFP space

As in FBA, feasible flux states of a metabolic network are those that fulfill stoichiometric mass balance constraints together with imposed upper and lower bounds on the reaction fluxes. These constraints restrict the number of solutions to a compact convex set which contains all possible flux steady states in a particular environmental condition. In glucose minimal medium (see Chapter 2, Section 2.2.2.1), the FFP space of the core *E. coli* model (its network can be seen in Supplementary Table C6-1) is determined by 70 potentially active reactions, including biomass formation and the ATP maintenance reaction, and 68 metabolites. Using the HR algorithm, a raw sample of $10^9$ feasible states is obtained, from which a final uniform representative set of $10^6$ feasible states is extracted.

Notice that the used approach is suitable for genome-scale network sizes beyond the reduced size of the core *E. coli* model. There is not any fundamental or technical bottleneck that prevents its application to complete metabolic descriptions at the cell level. In this chapter, the core *E. coli* model is used due to a matter of computational time and ease of visualization.

From the sampled set of core *E. coli* metabolic states in minimal medium of glucose bounded to 10 mmol gDW$^{-1}$h$^{-1}$, the metabolic flux profiles of each individual reaction is collected as the set of its feasible metabolic fluxes. From such profile, one can compute the probability density function $f(\nu)$ which describes the likelihood for a reaction to take on a particular flux value. In Figure 6.1, the profiles of all reactions are shown. One can observe a variety of shapes, all of them low-variance, most displaying a maximum probability for a certain value of the flux inside the allowed range[1], and many being clearly asymmetric. The allowed range is computed using Biomass unconstrained Flux Variability Analysis (see Chapter 2, Section 2.2.4).

To characterize the dispersion of the possible fluxes for each reaction, one can measure its coefficient of variation $CV(f(\nu))$ calculated as the ratio between the standard deviation of possible fluxes and their average (see Supplementary Table C6-2). For all but three reversible reactions (Malate dehydrogenase, Glucose-6-Phosphate isomerase, and Glutamate dehydrogenase), the only reversible reactions having a low associated flux mean and thus a higher $CV(f(\nu))$, this metric is below one and when ranked for all reactions it steadily decreases to almost zero, Figure 6.2a. Interestingly, it can be found that this coefficient is significantly anticorrelated with the essentiality of reactions as observed experimentally [119] (point-biserial correlation coefficient -0.29 with *p*-value

---

[1]Notice that none of these histograms can have more than one peak due to the convexity of the steady-state flux space.

FIGURE 6.1: Probability density functions of metabolic fluxes values for all reactions in core *E. coli* under glucose minimal conditions (*follows to next page*).

FIGURE 6.1: (*Follows from previous page*) Each graph shows the reaction label, the flux variability range (values inside parentheses), and each associated pathway (acronyms in italics). Notice that the range plotted in the axes does not coincide with the flux variability range, since in the axes an optimal x range for each reaction is chosen to distinguish the shape of each profile. In addition, in each profile the position of the FBA point (blue marker) and the position of the Mean (green marker) are also shown.

0.01, see Appendix C). This means that essential reactions tend to have a highly concentrated profile of feasible fluxes. Besides, and only for the glucose transferase reaction GLCpts, one finds a zero probability of having a zero flux, which indicates that this reaction is essential in glucose minimal medium as expected. The asymmetry of each profile is characterized by the distance between the more probable flux in the FFP space and the lower flux bound of the flux variability range rescaled by the flux variability range of the reaction (see Chapter 2, Section 2.2.4). In Figure 6.2b, a scatterplot of values for all 68 core reactions is shown. Strikingly, the rescaled distances cluster in three regions around 0, 0.5 and 1 forming groups of sizes 38, 15 and 17 respectively. This indicates that the most probable flux is close to either the lower or upper bound or, conversely, the probability distribution function tends to be quite symmetric. Moreover, it can be also observed that an anticorrelation between the length of the flux range and the position of the most probable flux is present, so that the closer is this to its maximum value the shorter is the allowed range of fluxes.

In order to assess the likelihood of flux states corresponding to FBA maximization of the flux through the biomass reaction (FBA-MBR) (or equivalently of the growth rate) in relation to typical[2] points within the whole FFP space, one can calculate the average flux value for each reaction, the mean, and compare it to the FBA optimal biomass production flux. The complementary cumulative

---

[2]In the mathematical/computational context, typical means statistically representative in relation to the whole set of flux states contained in the FFP space.

FIGURE 6.2: Analysis of reaction profiles and visualization of the FFP space. a) Coefficient of variation for all core reactions ranked by value. b) Scatterplot of distances between the more probable flux in the FFP space and the lower flux bound rescaled by flux variability range for each reaction. c) Complementary cumulative distribution function of distances between FBA maximal growth flux and FFP space mean flux rescaled by flux variability range for each reaction, in log-log scale. d) Matrix of Pearson correlation coefficients measuring the degree of linear associations between feasible fluxes of reactions (acronyms of the pathways are shown in Abbreviations). e) Projection of the FFP space onto the two principal component vectors of the correlation matrix in e). All sampled flux phenotypes are normalized and projected along the first ($\rho_1$) and second ($\rho_2$) principal components. The plot is in polar coordinates, with the negative logarithm of the radius. The majority of points lies in a circle close to the origin (the darker area). The FBA solution (green circle) is, conversely, rather eccentric.

distribution function of the distances between these two characteristic fluxes rescaled by the flux variability range of reactions is shown in Figure 6.2c (see Supplementary Table C6-2). A broad distribution of values can be observed over several orders of magnitude with no mean value actually very close to the FBA maximal solution except for a few reactions, typically working at maximum growth. At the other end of the spectrum, deviated reactions include for instance excretion of acetate and phosphate exchange. As a summary, one can conclude that the mean and the FBA biomass optimum are rather distant, which suggests that FBA optimal states are uninformative about phenotypes in the bulk of states in the FFP space.

To visualize neatly the eccentricity of the FBA maximum growth state with respect to the bulk of metabolic flux solutions, Principal Component Analysis [247, 248] is used in order to reduce the high-dimensionality of the full flux solution space projecting it onto a two-dimensional plane from the most informative viewpoint (see Appendix F). Reaction profiles are taken in pairs to calculate the matrix of Pearson correlation coefficients measuring their degree of linear association (see Figure 6.2d, the matrix is provided in Supplementary Table C6-3). Note that an ordering of reactions by pathways allows to have a clear visual feedback of intra- and inter-pathway correlations taking place in the core *E. coli* metabolic network, such that clusters of highly correlated reactions appear as bigger darker squares. The two axes of our projection correspond to the two first principal components of this profile correlation matrix $\rho_1$ and $\rho_2$, which account for most of the variability in profile correlations. Each sampled metabolic flux state has been rescaled as a z-score centered around the mean and projected onto these axes, as shown in the scatterplot Figure 6.2e in polar coordinates, where a negative logarithmic transformation to the radial coordinate for ease of visualization has been applied. The majority of phenotypes have a radius

FIGURE 6.3: a) Probability distribution function of the radii of all solutions before applying the negative logarithmic transformation. The red area denotes the probability of having a smaller radius than the FBA solution. This fraction of area is the 3% of the total area, which means that the 97% of the solutions have a larger radius than the FBA solution. b) Cumulative probability distribution function of the radius. The blue region denotes the range of solutions with a radius smaller than FBA. The probability of having a radius smaller than FBA is the y-value of the curve at the rightmost side of the region.

close to zero. Since points closer to the origin are better described by the two principal components (see Appendix F), this implies that $\rho_1$ and $\rho_2$ capture the largest variability of the sampled FFP. Clearly, the FBA optimal growth solution is rather eccentric with respect to typical solutions, with an associated radius of 0.98 in this representation. In fact, 97% of states have a smaller radius than the optimal growth solution (see Figure 6.3).

## 6.2 The FFP space gives a standard to calibrate the deviation of optimal phenotypes from experimental observations

This section focuses on the relationship between primary carbon source uptake and oxygen need to illustrate the potential of the FFP space as a benchmark

to calibrate the deviation of *in silico* predicted optimal phenotypes from experimental observations. Sampled FFP states of core *E. coli*, in particular FFP mean values, as a function of the upper bound uptake rate of the carbon source are compared with reported experimental data for oxygen uptakes in minimal medium with glucose, pyruvate, or succinate as a primary carbon source (see Figure 6.4). The line of optimality representing FBA optimal growth solutions is also considered. Glucose experimental data points were used from Reference [49], experimental results for pyruvate are reported in Reference [50], and experimental results in Reference [249] report the quantitative relationship between oxygen uptake rate and acetate production rate as a function of succinate uptake rate.

In all cases, FBA-MBR reproduces well experimental data points in the low carbon source uptake region [249], where *E. coli* is indeed optimizing biomass yield. However, oxygen uptake rate saturates after some critical threshold of carbon source uptake rate (which depends on the carbon source) reaching a plateau which, among other possibilities, could be explained by the existence of a physiological enzymatic limit in oxygen uptake that lessens the capacity of the respiratory system [250]. The plateau levels are $18.8 \pm 0.7$ mmol $gDW^{-1}$ $h^{-1}$ for glucose [249], $16.8 \pm 0.4$ mmol $gDW^{-1}$ $h^{-1}$ for pyruvate [50], and $19.49 \pm 0.78$ mmol $gDW^{-1}$ $h^{-1}$ for succinate [249]. In this region of high carbon source uptake, FBA-MBR predicts an oxygen uptake overestimated by around 25% with respect to the values reported from experiments. While this amount is in principle large, the FFP space gives a standard that helps to calibrate it.

The eccentricity of experimental observations is measured as their distance to the FFP mean. For glucose, this value is 19.4 mmol $gDW^{-1}$ $h^{-1}$, which makes the distance of 5.3 mmol $gDW^{-1}$ $h^{-1}$ between the FBA-MBR prediction and experimental data relatively low (see Figure 6.4a). The distance of 8.2

FIGURE 6.4: Comparison of predicted phenotypes and experimental data. Sampled points in the FFP space with maximum carbon source upper bound are plotted in shaded grey, darkness is proportional to the number of points. Experimental data points are red circles. The *in silico*-defined line of optimality, representing FBA optimal growth solutions as a function of the upper bound uptake rate of the carbon source, is shown in orange. Blue squares correspond to FFP mean values for different carbon source upper bound uptake rates. a) Oxygen vs. glucose uptake rates, experimental data from [49]. The FFP space is sampled with glucose bounded to 12 mmol gDW$^{-1}$ h$^{-1}$. b) Oxygen vs. pyruvate uptake rates, experimental data from Reference [50]. The FFP space is sampled with pyruvate bounded to 23 mmol gDW$^{-1}$ h$^{-1}$. c) Oxygen vs. succinate uptake rates, experimental data from Reference [249]. The FFP space is sampled with succinate bounded to 15 mmol gDW$^{-1}$ h$^{-1}$. *Inset* Acetate production rate vs. succinate uptake rate, experimental data from Reference [249].

mmol gDW$^{-1}$ h$^{-1}$ between the FBA-MBR prediction and experimental data is slightly worse for pyruvate (see Figure 6.4b), in which case the eccentricity of experimental observations is of 18.4 mmol gDW$^{-1}$ h$^{-1}$. The disagreement between optimality predictions and experimental data is much more significative in the case of succinate (see Figure 6.4c), for which the eccentricity of experimental observations is only of 4.3 mmol gDW$^{-1}$ h$^{-1}$, while the distance between the FBA-MBR prediction and experimental data is of 5.4 mmol gDW$^{-1}$ h$^{-1}$, meaning that the FFP mean is indeed more adjusted to observations. The case of acetate production for this carbon source is even more conspicuous (see Figure 6.4c *Inset*). While FBA-MBR is still reproducing well the experimental results of no acetate production in the low succinate uptake region, it cannot predict production of acetate at any succinate uptake rate due to the fact that FBA-MBR in minimal medium with unlimited oxygen does not capture the enzymatic oxygen limitation. The FBA-MBR solution diverts resources to the production of ATP entirely through the Oxidative Phosphorylation pathway. Thus, it fails to reproduce experimental observations of acetate production in the region of high succinate uptake rates [249, 251–253]. In contrast, most metabolic states in the FFP space are consistent with acetate production, so that in this case the FFP mean turns out as a good predictor of the experimentally observed metabolic behavior.

In summary, while FBA-MBR predictions seem accurate for low carbon source uptake rate states in minimal medium as seen previously [249], the experimental points diverge from the FBA-MBR prediction state when increased values of carbon source uptake rates are considered. Note that, in general, it is not straightforward to quantify the significance of the divergence. Here, the FFP space is proposed as a benchmark. According to this calibration, one finds that FBA optimal growth predictions of oxygen needs versus glucose, pyruvate, or

succinate uptake are worse the more downstream the position of the carbon source into catalytic metabolism. Using the core *E. coli* model, it has been checked that the ratio of the maximum ATP production rate to the maximum oxygen uptake (both calculated by FBA optimization of ATP production rate) for the three carbon sources glucose, pyruvate, and succinate are respectively 2.9, 2.6, and 2.4, so this ratio decreases as more downstream in the catalytic metabolism.

## 6.3    The high-biomass production region of the FFP space displays aerobic fermentation in minimal medium with unlimited oxygen uptake

The high-growth metabolic region of the core *E. coli* FFP space is resampled in glucose minimal medium with a glucose upper bound of 10 mmol gDW$^{-1}$ h$^{-1}$. This region is defined by setting a minimal threshold for the biomass production of $\geq 0.4$ mmol gDW$^{-1}$ h$^{-1}$ [254], and the new sample has a final size of $10^5$ states. Note that phenotypes in this high-growth sample remain very close to the biomass yield threshold due to the exponential decrease of the number of feasible flux states with increased biomass production, as shown in the biomass flux profile in Figure 6.1.

In this region, one can identify pathway utilization typical of proliferative microbial metabolism, even when considering a minimal medium and unlimited oxygen uptake. This metabolic behavior is consistent with experimental data [49, 249, 255] but it is unreachable by FBA models based on optimality principles (unless optimization is accompanied by auxiliary constraints not assumed in standard FBA implementations, like the solvent capacity constraint [254], or

FIGURE 6.5: Schematic of pathway utilization in high-growth vs low-growth conditions.

by modelization beyond stoichiometric mass balance, for instance, thermodynamically feasible kinetics or enzyme synthesis [256, 257]). These by-products cannot be explained by FBA-MBR in minimal medium with unlimited oxygen supply since, in this optimization framework, metabolic fluxes are basically forced to ATP production through Oxidative Phosphorylation with excretion of $CO_2$ as waste. However, increasing the oxygen limitation in FBA-MBR results in secretion of formate, acetate, and ethanol -in that order-, with corresponding shifts in metabolic behavior [250].

According to the FFP space of core *E. coli*, one can observe that the high-biomass production FFP subsample is characterized by the secretion of small organic molecules, even when the supply of oxygen is unlimited. This fact points to the simultaneous utilization of Glycolysis and Oxidative Phosphorylation to produce biomass and energy, as illustrated in the schematic shown in Figure 6.5. Quantitative relationships between the production of small organic molecules and glucose and oxygen uptake rates are shown in the remaining panels of Figure 6.6. Three-dimensional scatterplots for the production rates of formate, acetate, ethanol, and lactate are shown in Figures 6.6a, 6.6c, 6.6e, and 6.6g respectively, with projections into the three possible two-dimensional planes

FIGURE 6.6: High growth phenotypes of core *E. coli* on glucose minimal medium. a) 3-dimensional scatterplot of formate production rate vs glucose and oxygen uptake rates. b) Density projections of a) on each of the possible 2D planes, formate-glucose, formate-oxygen, and glucose-oxygen. c) 3-dimensional scatterplot of acetate production rate vs glucose and oxygen uptake rates. d) Density projections of c) on each of the possible 2D planes, acetate-glucose, acetate-oxygen, and glucose-oxygen. e) 3-dimensional scatterplot of ethanol production rate vs glucose and oxygen uptake rates. f) Density projections of e) on each of the possible 2D planes, ethanol-glucose, ethanol-oxygen, and glucose-oxygen. f) 3-dimensional scatterplot of lactate production rate vs glucose and oxygen uptake rates. g) Density projections of f) on each of the possible 2D planes, lactate-glucose, lactate-oxygen, and glucose-oxygen.

shown in Figures 6.6b, 6.6d, 6.6f, and 6.6h respectively. As the levels of glucose and oxygen uptakes are raised, metabolic phenotypes can achieve an increased production of formate, acetate, ethanol, and lactate even though the majority of feasible phenotypes remain at low production values. Due to the high-growth requirement, oxygen uptake is always high but its variability increases with glucose uptake increase around a value of approximately 41.2 mmol gDW$^{-1}$ h$^{-1}$, which clusters the majority of high-growth metabolic phenotypes. Interestingly, this oxygen uptake rate value marks a region in the FFP space with maximum potential production rates of formate, acetate, ethanol, and lactate. Above and below that value most states are concentrated in the range [39.0,42.0] mmol gDW$^{-1}$ h$^{-1}$.

Taken together, these results indicate that, contrarily to standard FBA predictions, a high level of glucose uptake combined with enough oxygen can maintain the requirements of proliferative metabolism for biomass formation through aerobic fermentation even if the rest of nutrients are scarce and restricted to the minimum. At the same time, additional oxygen uptake diverts glucose back towards more efficient ATP production through Oxidative Phosphorylation. Hence, oxygen has the potential of regulating the glucose metabolic switch in which glucose uptake rates larger than a critical threshold around 5.0 mmol gDW$^{-1}$ h$^{-1}$ [254] lead to a linearly increasing maximum organic by-products production by a gradual activation of aerobic fermentation and a slight decrease of Oxidative Phosphorylation.

## 6.4   Conclusions

The information content of the full FFP space of metabolic states in a certain environment provides with an entire map to explore and evaluate metabolic

behavior and capabilities. While optimality goals need to be tailored to conditions and produce singular optimal solutions that may not be consistent with experimental observations, we have nowadays sufficient computational and methodological capacity to produce and analyze full FFP maps. The latter offer a reference framework to put into perspective the likelihood of particular phenotypic states that, as shown, enables to uncover metabolic behaviors that are unreachable using standard models based on optimality principles. In fact, the location of metabolic flux distributions into precise optimal states has been challenged recently by the proposal that metabolic flux evolve under the trade-off between two forces, optimality under one given condition and minimal adjustment between conditions [244]. In this way, resilience to changing environments necessarily forces flux states to near-optimal but suboptimal regions of feasible flux states in order to maintain adaptability.

In the FFP map of core *E. coli* in aerobic minimal medium, optimal growth states appear as eccentric and far from the bulk of more probable phenotypes represented by the FFP mean, which offers an ergodic perspective of the FFP space in which all states can be explored at random with equal probability. One of the uses of the method is precisely to evidence the effects of evolutionary pressure on organisms, which may actually result in eccentric flux states. On the other hand, the FFP space gives a standard to calibrate the deviation of optimal phenotypes from experimental observations. Oxygen consumption is a particularly interesting target for analysis since it has been identified as a trigger of metabolic shifts [250, 258]. Interestingly, according to the FFP map as a reference standard, it is found that, in high-growth conditions, FBA-MBR predictions of experimental observations for unlimited oxygen needs versus glucose, pyruvate, or succinate uptakes are worse the more downstream the uptake of the carbon source into the catalytic metabolic stream. This is consistent

with the fact that the FBA-MBR solution diverts resources to the production of ATP entirely through the Oxidative Phosphorylation pathway, so that the more is the effective potential of the carbon source to recombine with oxygen to produce energy the more convergent will be the *in silico* prediction and the observed states.

In order to correct FBA in high-growth conditions, some investigations restricted the solution space beyond mass balance and uptake bounds through additional thermodynamic, kinetic or physiological constraints, like the solvent capacity constraint quantifying the maximum amount of macromolecules that can occupy the intracellular space [254]. Alternatively, the objective function implemented in FBA has been modified to nonlinear maximization of the ATP or biomass yield per flux unit [243], or modelization beyond stoichiometric mass balance, like thermodynamically feasible kinetics or enzyme synthesis, has been considered [256, 257]. While these FBA modifications enhance some predictions, their effectiveness depends on the estimation of kinetic coefficients using empirical or experimental data. In contrast, the FFP map naturally displays all high-growth feasible states which show characteristic metabolic behaviors, like aerobic fermentation with unlimited oxygen uptake even in minimal medium, without the need to determine additional constants. This aerobic fermentation, apparently inefficient in terms of energy yield as compared to Oxidative Phosphorylation, has been demonstrated to be a favorable catabolic state for all rapidly proliferating cells with high glucose uptake capacity [254], and from this analysis it turns out as a probable metabolic phenotype even in minimal medium.

Beyond theoretical implications, FFP maps of microbial organisms can be of particular interest as tools for biotechnological applications, for instance in the engineering of *E. coli* fermentative metabolism as a fundamental cellular capacity for valuable industrial biocatalysis [259]. In biomedicine, the investigation

of FBA optimal phenotypes in the framework of the FFP map can help to contextualize disease phenotypes in comparison to normal states. For instance, FBA proved suitable for modeling complex diseases like cancer as it assumes that cancer cells maximize growth searching for metabolic flux distributions that produce essential biomass precursors at high rates [185, 260]. The analysis of the entire region of high-growth phenotypes will allow to reach and study a variety of suboptimal feasible flux states close to optimality but which cannot be reproduced by optimality principles, and so it opens new avenues for the understanding of general and fundamental mechanisms that characterize this disease across subtypes.

## 6.5   Summary

- FFP maps offer a reference framework to put into perspective the likelihood of particular phenotypic states. It enables to uncover metabolic behaviors that are unreachable using standard models based on optimality principles [245].

- Optimal FBA growth states are eccentric and appear far from the bulk of more probable phenotypes represented by the FFP mean [245].

- The FFP space gives a standard to calibrate the deviation of extreme phenotypes from experimental observations [245].

- The FFP map naturally displays all high-growth feasible states which show characteristic metabolic behaviors like aerobic fermentation with unlimited oxygen uptake even in minimal medium without the need to force additional constants [245].

# Chapter 7

# Conclusions

This thesis presents a study of cell metabolism from a systems-level approach trying to unveil new mechanisms and responses impossible to reach by traditional reductionist procedures. Different methods and analysis techniques have been used, and each one has allowed to extract new insights about the properties of cell metabolism. Tools that belong to the complex network science and Systems Biology have been used. On what follows, the conclusions of this thesis are given, answering to the objectives stated in Chapter 1.

The thesis starts by considering the study of the topology, *i.e.*, the connectivity pattern, of metabolic networks. From this point of view, it is possible to check whether the structure of metabolic networks has evolved towards increasing its robustness against external perturbations. It is important to notice that, at this stage, reaction fluxes are not considered.

From the obtained results of the first chapter of this thesis, one can conclude that the structure of the metabolic networks of *Escherichia coli*, *Staphylococcus aureus*, and *Mycoplasma pneumoniae* has evolved towards robustness against individual

and multiple reaction failures, which produce a reduced damaged compared to failures in degree-preserving randomized counterparts. *M. pneumoniae* is an exception in relation to individual reaction failures. This feature can be explained in terms of its simpler structure. Moreover, it is found that failures provoked by pairs of reactions generate an amplification effect which arises due to the non-linear interactions between the two damaging cascades propagating in the networks. In addition, a predictor of damage propagation for single cascades computed locally accounts for damage spreading. Also at the local level, a series of structural motifs can explain amplified failure patterns in double reaction cascades.

When the study is extended to gene failures, one finds that the method to compute cascades captures most of the scenarios of experimentally determined lethality in *M. pneumoniae*. Furthermore, when referring to multiple failures, the proposed analysis allows to find that (1) for failure cascade spreading, the distribution of cluster sizes is more important than the actual composition of the clusters, and (2) the regulation of high-damage genes tends to appear isolated from that of other genes, a kind of functional switch in metabolic networks that at the same time acts as a kind of genetic firewall. In any case, it is important to notice that a cascade may not only be interpreted as the harmful spreading of failures, but also as the ability to efficiently regulate metabolism. Large cascades may point at the evolutionary requirement of regulating large parts of metabolism through the regulation of small sets of enzyme-coding genes. Therefore, evolutionary pressure seems to favor the ability of efficient metabolic regulation at the expense of robustness to reaction knockouts.

This study can be complemented taking into account the fluxes flowing through the biochemical reactions with the aim to describe more appropriately real features of metabolic operation. These investigations permit to know how

reactions adapt to different situations, extending for instance the previous study of gene knockouts, or additionally, looking at responses to changes of the composition of the external environment. Flux Balance Analysis is used to compute fluxes of biochemical reactions. This method is based on different suppositions, principally that (1) metabolic networks work at steady state and that (2) the biological target of organisms is to grow as much as possible. In this way, FBA can be used to go beyond the mere analysis of the structure of metabolic networks and to identify metabolic fluxes that cannot be resolved using only a topological analysis. When FBA is applied to single reaction knockouts in *E. coli*, the main conclusion is that there exists a set of reactions which must be always active in order to ensure viability. However, non-essential reactions deserve special attention, either considering their role as growth enhancers or their potential participation in synthetic lethal pairs.

The study of synthetic lethal pairs allows to understand new protection mechanisms that metabolism has developed to survive. Synthetic lethal reaction pairs can be classified into two classes, plastic and redundant, depending on whether one reaction is active for maximum growth in the medium under consideration and the second inactive (plasticity) or, conversely, both reactions have simultaneously non-zero fluxes (redundancy). This particular study is made in both *E. coli* and *M. pneumoniae*. On the one hand, plasticity is a sophisticated mechanism that is able to reorganize metabolic fluxes turning on inactive reactions when coessential counterparts are removed so as to maintain viability, working as a backup mechanism. On the other hand, redundancy corresponds to a simultaneous use of different flux channels, ensuring in this way viability and increasing the growth rate of the organism. Furthermore, plasticity requires a higher degree of functional organization, using at the same time less

resources for maximum growth. It takes place more often in *E. coli* than in *M. pneumoniae*.

The previous study is completed by analyzing how plasticity and redundancy depend on the external environment for *E. coli*. One finds that plasticity and redundancy are conserved independently of the composition of the medium which acts as environmental condition for growth. Moreover, this conservation takes place also when this environment is enriched with non-essential compounds or overconstrained to decrease the maximum growth rate.

One can further exploit FBA, assuming conditions of growth optimality, in order to assess evolution or adaptation characteristics of metabolic networks. A filtering method called disparity filter allows to reduce the density of links of metabolic networks while preserving their main features. The metabolic networks of *E. coli* and *M. pneumoniae* are filtered to extract their backbones. First of all, it is checked that the disparity filter is, indeed, very efficient in order to decrease the link density of the studied metabolic networks using FBA fluxes as the weights of the links.

The analysis of the connected components of the metabolic backbones of both *E. coli* and *M. pneumoniae* in a glucose minimal medium allows to identify that these components mainly contain reactions that belong to ancient pathways, *i.e.*, pathways showing long-term evolution. Moreover, for both organisms, the presence of pathways related to energy metabolism -like Glycolysis, Citric Acid Cycle, and Oxidative Phosphorylation for *E. coli*, or Glycolysis and Pyruvate Metabolism for *M. pneumoniae*- could mean that these pathways have an important role in maximizing the growth and have evolved towards maximum efficiency to obtain chemical energy, something very important in case of nutrient scarcity and hence energy deficiency.

In addition, the study of the dependence of *E. coli* backbones on different environments allows to identify environment specific pathways displaying short-term adaptation. First, the analysis of the metabolic backbone obtained in a rich medium allows to demonstrate that the nutritionally-rich medium induces a large increase in the growth rate of *E. coli* due to nutrient abundance. The instantaneous response of *E. coli* to environment is to synthesize as much as membrane lipids as possible, since fast-growing cells must synthesize membrane components more rapidly to satisfy the high lipid demand to generate new cells. Second, with the study of the different backbones obtained from different minimal media, one finds that the distribution of the fluxes is little dependent on the nutrients present in the environment. In addition, it is also possible to extract that the pathway Alternate Carbon Metabolism is, for *E. coli*, the pathway with more capabilities to respond to external stimuli.

It is worth remarking that FBA makes the supposition that the biological target of organisms is to grow as much as possible. This may be plausible in some situations but there exist other in which the biological target of an organism is not to maximize growth. Hence, a study of the entire space of possible flux solutions can help to assess whether the FBA solution is representative of the whole space or not. The whole space encompassing the entire set of flux solutions, referred to as the full feasible flux phenotypes (FFP) space, is computed for *E. coli*. The information contents of the FFP space of metabolic states in a certain environment provides with an entire map to explore and evaluate metabolic behavior and capabilities. In fact, FFP maps can answer the question of whether FBA gives a representative solution of the flux space. The main conclusion is that optimal growth states obtained via FBA computations appear as eccentric and far from the bulk of more probable phenotypes.

In addition to the eccentricity of the FBA solution, the FFP space also gives a standard to calibrate the deviation of phenotypes obtained using FBA from experimental observations. Thus, it serves to compare FBA predictions with experimental results. For instance, the analysis of oxygen needs versus glucose, pyruvate, or succinate uptakes show that FBA results are worse the more downstream the uptake of the carbon source into the catalytic metabolic stream. This is explained due to the fact that the FBA solution diverts resources to the production of ATP entirely through Oxidative Phosphorylation. In this way, the more the effective potential of the carbon source to recombine with oxygen to produce energy using Oxidative Phosphorylation, the more convergent will be the FBA prediction with respect to experimental results.

On the other hand, the FFP space naturally displays all high-growth feasible states which show characteristic metabolic behaviors, like aerobic fermentation with unlimited oxygen uptake even in minimal medium. This is an important feature, since these metabolic behaviors cannot be obtained under FBA maximum growth computations without using additional constraints. This reinforces the idea that the FFP map contains valuable information about metabolic states.

It is important to point out that the used methodology in this thesis is not restricted to bacteria, and that it could also be applied to metabolic networks of other species. In particular, the results of the study of structural stress may have potential implications in areas like metabolic engineering or disease treatment. The study of complex systems under structural stress poses a number of formidable challenges critical to understand their behavior as well as towards proposing successful strategies for prediction and control. In this framework, the study of structural stress in human pathogens may help to develop more sophisticated forms of identifying new and more efficient drug targets.

Plasticity and redundancy are very important concepts for biological complex systems in general. Whether they are adaptive in cell metabolism or, as it has been argued for metabolism in changing environments [197, 203], they are rather a byproduct of the evolution of biological networks toward survival, these regulatory mechanisms are key to understand how complex biological systems protect themselves against malfunction. Among the many different applications of synthetic lethality, one of them is to determine the accuracy of gene essentiality of new genome-scale reconstructions of metabolic networks [261].

Since the application of the disparity filter in metabolic networks can be used to recognize pathways and reactions which (1) are more sensitive to environmental changes, and (2) which are involved in the maximization of the growth rate of an organism due to evolutionary pressure, its use could be appropriate in the field of biotechnology. For example, it could be useful for the targeting of the most important pathways present in cancer cells which are in charge of their high growth rate. Therefore, this could help to understand the biochemical mechanisms that cancer cells use to proliferate. In this way, it will be possible to find a way to decrease the high performance achieved by cancer cells in terms of growth efficiency.

Finally, FFP maps of microbial organisms can be of particular interest as tools for biotechnological applications, for instance in the engineering of *E. coli* fermentative metabolism as a fundamental cellular capacity for valuable industrial biocatalysis [259]. In biomedicine, the investigation of FBA phenotypes in the framework of the FFP map can help to contextualize disease phenotypes in comparison to normal states. For instance, FBA proved suitable for modeling complex diseases like cancer as it assumes that cancer cells maximize growth searching for metabolic flux distributions that produce essential biomass

precursors at high rates [185, 260]. The analysis of the entire region of high-growth phenotypes will allow to reach and study a variety of suboptimal feasible flux states close to optimality but which cannot be reproduced by optimality principles, and so it opens new avenues for the understanding of general and fundamental mechanisms that characterize this disease across subtypes.

# Appendix A

# Kolmogorov-Smirnov test

The Kolmogorov-Smirnov test [187] is a test used in statistics which compares the probability distribution obtained from a sample with a reference probability (one-sample K-S test), or which compares two samples (two sample K-S test). It basically quantifies a distance between the cumulative distribution function of the sample and the cumulative distribution function of the reference distribution, or between the cumulative distribution functions of two samples. The null hypothesis of this test assumes that the samples are obtained from the same distribution (two sample K-S test) or that the sample is drawn from the reference distribution (one sample K-S test). The two-sample KS test is one of the most useful methods for comparing two samples, which is the variant that has been used in this thesis.

To compare two samples, first of all one has to compute the maximum distance $K - S$ between the two cumulative distribution functions

$$K - S = \max |F_{1,n}(x) - F_{2,n'}(x)| \tag{A.1}$$

FIGURE A.1: Visualization of the value $K - S$ used in the K-S test computed using two Log-normal distributions [262] with different means and the same standard deviation. After computing the maximum difference, this value is transformed into a *p*-value.

where $F_{1,n}(x)$ and $F_{2,n'}(x)$ are the cumulative distribution functions of the fist and second sample, and $n$ and $n'$ are the sizes of each sample respectively. To compute the associated significance of the value of $K - S$, one has to calculate the *p*-value applying the following expression:

$$p = 2 \sum_{j=1}^{\infty} (-1)^{j-1} \exp(-2 j l^2) \tag{A.2}$$

where $l = K - S \cdot (\sqrt{N} + 0.12 + \frac{0.11}{\sqrt{N}})$, and $N = \frac{n \, n'}{n+n'}$. Then, one compares this *p*-value with the chosen reference, usually $\alpha = 0.05$. If $p < \alpha$, one can consider that both distributions are drawn from the same distribution, otherwise they are considered significantly different.

# Appendix B

# Spearman's rank correlation coefficient

The Spearman's rank correlation coefficient [188], often denoted by the Greek letter $\rho$, is a nonparametric measure used in statistics which measures statistical dependence between two variables. It assesses how well the relationship between two variables can be described using a monotonic function. Spearman's coefficient can be used both for continuous and discrete variables, including ordinal variables.

The Spearman's coefficient is basically the Pearson correlation coefficient between the ranked variables. The ranks of both samples are compared and the value of $\rho_S$ is computed with the following expression:

$$\rho_S = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 (y_i - \bar{y})^2}} \tag{B.1}$$

where $x_i$ and $y_i$ are the ranks of the values of the sample $X_i$ and $Y_i$.

To assess the significance of the measure, a permutation test is done in this thesis, where the values of $X_i$ and $Y_i$ are reshuffled and then, for each realization, $\rho_S$ is calculated. After doing this for all realizations, one keeps the maximum and minimum value of the obtained $\rho_S$, which gives the interval that belongs to the null model. Thus, if the value of $\rho_S$ of the original sample lies within this interval, it implies that there is no correlation between the ranks of both samples. Otherwise, if the value of $\rho_S$ of the original sample lies outside the range of the null model, one can consider that there exists a correlation between both samples.

# Appendix C

# Point-biserial correlation coefficient

The point biserial correlation coefficient ($r_{pb}$) is a correlation coefficient which is used when one variable is continuous and the other is dichotomous. This dichotomous variable can either be a truly dichotomous variable, like male / female, or an artificially dichotomized variable, obtained by using a threshold on a continuous variable. However, in most situations it is not advisable to dichotomize variables artificially and thus it is more appropriate to use specific statistical tests for continuous variables.

The point-biserial correlation is equivalent to the Pearson correlation. To calculate the point-biserial correlation coefficient, one assumes that the dichotomous variable can have the values 0 and 1. Therefore, one can divide the data between two groups, the first group which corresponds to the value 1 on the dichotomous variable, and the second group which corresponds the value 0 on the dichotomous

variable. Thus, the point-biserial correlation coefficient is calculated as follows:

$$r_{pb} = \frac{M_1 - M_0}{s_n} \sqrt{\frac{n_1 n_0}{n^2}} \tag{C.1}$$

where $M_1$ is the mean value of the continuous variable for all data points in the first group, and $M_0$ is the mean value of the continuous variable for all data points in the second group. Further, $n_1$ is the number of data points in the first group, $n_0$ is the number of data points in the second group, and $n$ is the total sample size. $s_n$ is the standard deviation computed as follows:

$$s_n = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (x_i - \nu)^2} \tag{C.2}$$

where $x_i$ is continuous variable and $\nu$ is its average value. It is possible to compute a $t$-value (associated to a Student's $t$-distribution) from this correlation coefficient:

$$t = r_{pb} \sqrt{\frac{n_1 + n_0 - 2}{1 - r_{pb}^2}} \tag{C.3}$$

where $r_{pb}$ is the the point-biserial correlation coefficient. From this value of $t$, a $p$-value of the significance can be obtained by computing the area of the Student's $t$-test from $-\infty$ to the computed value of $t$ with $(n_1 + n_0 - 2)$ degrees of freedom. If the $p$-value is lower than a chosen critical value of the significance (usually 0.05), one can consider that there is a significant correlation between the continuous and the dichotomous variable. Otherwise, one must conclude that there is not a significant correlation between both variables.

# Appendix D

# Disparity filter

The disparity filter [95] takes advantage of the local fluctuations present in the weights of the links between nodes. It is useful to define the strength $s_i$ of a node $i$ as the sum of the weights ($\nu_{ij}$) of the links associated to this node, $s_i = \sum_j \nu_{ij}$. The filtering method starts by normalizing the weight of the nodes $p_{ij} = \frac{\nu_{ij}}{s_i}$, where $\nu_{ij}$ is the weight of a link $j$ of the node $i$, since one needs a measure of the fluctuations of the weights attached to a node at the local level. The key point is that a few links have a large value of $p_{ij}$ being thus more significant than the others, as computed by the disparity measure defined as $\Upsilon_i(k) \equiv k \sum_j p_{ij}^2$, where $k$ is the degree of the node and $p_{ij}$ is the normalized weight of the link between node $i$ and node $j$.

In the application of this method to metabolic networks, $\Upsilon_i(k)$ characterizes the level of local heterogeneity of a metabolite $i$, and so $p_{ij}$ stands for the normalized weight of the link between metabolite $i$ and reaction $j$, with $\nu_{ij}$ the flux of reaction $j$. Under perfect homogeneity, when all the links share the same amount of the strength of the node, $\Upsilon_i(k)$ equals 1 independently of $k$, whereas

for perfect heterogeneity, when one of the links carries the whole strength of the node, $\Upsilon_i(k)$ equals $k$. Usually, an intermediate behavior is observed in real systems.

To assess the deviations of the weights of the links, a null model is used which provides the expectation of the disparity measure of a node in a random case. The null hypothesis consists on the fact that the normalized weights that correspond to a certain node are produced by a random assignment coming from a uniform distribution. Notice that, since in this chapter directed metabolic networks are used, one has three kinds of links. Bidirectional links are decoupled into incoming and outgoing links, leading to a network where nodes have incoming and outgoing links. Each kind of links are treated independently, each one having its own probability density function. The filter then proceeds by identifying which links must be preserved. To do this, one computes the probability $\alpha_{ij}$ that a weight $p_{ij}$ is non-compatible with the null model. This probability is compared to a significance level $\alpha$, and thus links that carry weights with a probability $\alpha_{ij} < \alpha$ can be considered non-consistent with the null model and they are considered significant for the metabolite. The probability $\alpha_{ij}$ is computed with the expression $\alpha_{ij}^{in/out} = (1 - p_{ij}^{in/out})^{k^{in/out}-1}$. Note that, for nodes with only one incoming or outgoing connection, one uses the prescription to preserve those links.

# Appendix E

# Hit-And-Run algorithm

The feasible flux phenotypes (FFP) space of different metabolic models in specific environments has been explored using different sampling techniques [12, 263–265]. Here, the Hit-And-Run (HR) algorithm is used, tailoring it to enhance its sampling rate and to minimize its mixing time [265]. On what follows, the key points and ideas behind the HR algorithm are stated.

One must start by noticing that all points in the FFP space must simultaneously satisfy mass balance conditions and uptake limits for internal and exchanged metabolites. The former requirement defines a set of homogeneous linear equalities, whose solution space is $K$, while the latter defines a set of linear inequalities, whose solutions lie in a convex compact set $V$. From a geometrical point of view, the FFP space is thus given by the intersection $S = K \cap V$. A key step of the HR approach used here consists on realizing that one can directly work in $S$ by sampling $V$ in terms of a basis spanning $K$. This allows to retrieve all FFPs that satisfy mass balance in the medium conditions under consideration, without rejection. Additionally, sampling in $S$ allows to perform

a drastic dimensional reduction and to decrease considerably the computation time. Indeed, assuming to have $N$ reactions, $I$ internal metabolites, and $E$ exchanged metabolites ($N > I + E$), one has that $S \subset \mathbb{R}^{N-I}$, which is typically a space with greatly reduced dimensionality with respect to $V \subset \mathbb{R}^N$.

Once a basis for $K$ is found, the main idea behind HR is fairly simple. Given a feasible solution $\boldsymbol{\nu}_\mathrm{o} \in S$, a new, different feasible solution $\boldsymbol{\nu}_\mathrm{n} \in S$ can be obtained as follows:

1. Choose a random direction $\boldsymbol{u}$ in $\mathbb{R}^I$

2. Draw a line $\ell$ through $\boldsymbol{\nu_o}$ along direction $\boldsymbol{u}$:

$$\ell : \boldsymbol{\nu_o} + \lambda \boldsymbol{u}, \qquad \lambda \in \mathbb{R}$$

3. Compute the two intersection points of $\ell$ with the boundary of $S$, parametrized by $\lambda = \lambda_-, \lambda_+$:

$$\boldsymbol{\nu}_- = \boldsymbol{\nu_o} + (\lambda_-)\boldsymbol{u}$$
$$\boldsymbol{\nu}_+ = \boldsymbol{\nu_o} + (\lambda_+)\boldsymbol{u}$$

4. Choose a new point $\boldsymbol{\nu}_\mathrm{n}$ from $\ell$, uniformly at random between $\boldsymbol{\nu}_-$ and $\boldsymbol{\nu}_+$. In practice, this implies choosing a value $\lambda_\mathrm{n}$ in the range $(\lambda_-, \lambda_+)$ uniformly at random, and then

$$\boldsymbol{\nu}_\mathrm{n} \equiv \boldsymbol{\nu}_\mathrm{o} + \lambda_\mathrm{n} \boldsymbol{u}$$

This procedure is repeated iteratively so that, given an initial condition, the algorithm can produce an arbitrary number of feasible solutions (see Figure E.1 for an illustrative representation of the algorithm). The initial condition, which must be a feasible metabolic flux state itself (*i.e.*, it must belong to $S$), is

FIGURE E.1: Illustrative representation of the HR fundamental step, which generates a new feasible state $\nu_n$ $\nu_{\mathbf{n}}$ from a given one $\nu_{\mathbf{o}}$.

obtained by other methods. In this work, the algorithm called MinOver is used, see [265, 266], but any other technique is valid. In particular, in cases where small samples of the FFP space have been already obtained by other sampling techniques, such points can be used to feed the HR algorithm and produce a new, larger sample.

It was proven [246] that, by iterating steps (1-4), the samples obtained are asymptotically unbiased, in the sense that the whole FFP space is explored with the same likelihood in the limit of very large samples. In practice, one must always work with a finite sample, and hence the following additional measures are taken so as to ensure that the used samples were truly representative of the whole FFP space. In particular:

1. Only one every $10^3$ points generated by HR is included in the final sample. This effectively decreases the "mixing time" of the algorithm, since the correlation among the points that are actually retained decays fast.

2. Different initial conditions are used. Results show no dependence on the initial condition, as expected for large samples. Even so, the first 30% of

points are discarded, in order to rule out any subtler effect of the initial condition on the final results.

3. Results are recalculated using subsamples of size 10% of the original sample. Qualitative differences between the two sets are not found.

Since the HR algorithm is very efficient itself and due to the dimensionality reduction that this implementation adds, very large samples can be generated in reasonable time. For each model, samples of size $10^9$ are initially created, giving rise to a final set of $10^6$ feasible solutions uniformly distributed along the whole FFP space.

# Appendix F

# Principal Component Analysis

The computation of reaction pairs correlations may be exploited to detect how global flux variability emerges in the system through Principal Component Analysis (PCA) [247, 248] and to quantify, in turn, the closeness of optimal phenotypes to the bulk of the feasible flux phenotypes (FFP) space. On what follows, PCA is briefly described, while an illustrative example is also provided (see Figure F.1).

One starts by writing down the matrix $C_{ij}$ of correlations between all reaction pairs $i, j$. In doing this, one measures how much the variability of a reaction flux $\nu_i$ affects the flux $\nu_j$ (and viceversa). In mathematical terms, for each pair of reactions $i, j$, one has:

$$C_{ij} = \frac{\langle \nu_i \nu_j \rangle - \langle \nu_i \rangle \langle \nu_j \rangle}{\sqrt{\left(\langle \nu_i^2 \rangle - \langle \nu_i \rangle^2\right)\left(\langle \nu_j^2 \rangle - \langle \nu_j \rangle^2\right)}}, \tag{F.1}$$

where $\langle \ldots \rangle$ denotes an average over the sampled set and the denominator of the fraction is simply the product of the standard deviations of $\nu_i$ and $\nu_j$. This matrix is shown in Figure 6.4e in Chapter 5.

Matrix $\boldsymbol{C}$ is real and symmetric by definition and, thus, diagonalizable. This means that, for every eigenvector $\rho_\kappa$, one has $\boldsymbol{C}\rho_\kappa = \lambda_\kappa \rho_\kappa$. Note that matrix $\boldsymbol{C}$ describes paired flux fluctuations in a reference frame centered on the mean flux vector. The eigenvectors $\rho_\kappa$ of $\boldsymbol{C}$ express, in turn, the directions along which such fluctuations are taking place. In particular, the eigenvectors $\rho_1, \rho_2$ associated with the first two largest (in modulo) eigenvalues dictate the two directions in space where the sampled FFP displays the greatest variability (see Figure F.1). This implies that sampled phenotypes lie closer to the plane spanned by $\rho_1$ and $\rho_2$ than the ones produced by any other linear combination of $\boldsymbol{C}$ eigenvectors. Projecting all sampled FFP onto this plane allows thus to perform a drastic dimensional reduction yet retaining much of the original variability and allows to have a direct graphical insight on where phenotypes lie, on where the bulk of the FFP is located, and on how the Flux Balance Analysis (FBA) solution compares to them. In such plot, each phenotype $\jmath$ is described by two coordinates that may be parametrized via a radius $r_\jmath$ and an angle $\theta_\jmath$. Since the projection is normalized, it follows that $r_\jmath \leq 1$. Furthermore, the closer $r_\jmath$ to one, the better the phenotype $\jmath$ is described by only looking at variability along $\rho_1, \rho_2$. As $r_\jmath$ is one at the most and since one has so many phenotypes clustered together, it is possible to choose to plot the PCA projection by using an effective radius $r'_\jmath = -\log r_\jmath$, as in Figure 6.4e. In this way one could better discriminate among different phenotypes and got a 'closest to the origin, closest to the $\rho_1, \rho_2$–plane' setup.

FIGURE F.1: An example to describe PCA analysis. a) FFP sampling produces a cloud of points in a multidimensional space that, when projected along the $(x, z)$, $(y, z)$ and $(x, y)$ planes, is seen to span a wide range of values. Finding the eigenvectors of the correlation matrix, one can see that such points are actually clustered around a plane (plotted as a yellow grid). By diagonalizing the $(3 \times 3)$ correlation matrix, one finds the three vectors (plotted in blue, red and green, respectively) identifying the direction in space where the points show most variation, in a decreasing manner. A black square is also plotted as a reference eccentric point. b) By projecting the sampled FFP along vectors $\rho_1$, $\rho_2$, and $\rho_3$, all points are squeezed in a thin region close to the $(\rho_1, \rho_2)$ plane. This shows that the greatest variability of the sampled points actually occurs in the $\rho_1$, $\rho_2$ directions. In this representation, the eccentric black square point is seen to lie far from the plane with a large $\rho_3$ coordinate. c) Normalizing the projection in b) over the modulus of the vector identifying the point coordinate allows to quantify the closeness to the $(\rho_1, \rho_2)$ plane. In such way all points are projected over the unit radius sphere, with the majority of points scattered near the equator, *i.e.*, the $(\rho_1, \rho_2)$ plane. Therefore, in this representation, eccentric points like the black square are close to the pole. d) Points on the unit sphere may in turn be projected on the $(\rho_1, \rho_2)$ plane only. In this way all points are constrained within the unitary radius circle, with points close to the equator in plot c) now close to the circle and the ones close to the pole in c) near the origin. In this representation, typical points, *i.e.*, those originally closer to the yellow plane in a), have larger radius (close to one, but smaller than that) and eccentric points have a smaller radius, like the black square (*follows to next page*).

FIGURE F.1: (*Follows from previous page*) e) Plotting the distribution of the points radius, as in Figure 6.3, one sees that $P(r)$ has indeed a peak in one, with very low probability of finding a point with a radius close to zero. Similarly to Figure 6.3 the radius of the eccentric point is indicated, highlighting how low $r$, eccentric points are indeed unlikely. f) Similarly to Figure 6.2e in Chapter 5, the points on the $(\rho_1, \rho_2)$ plane are re-projected, but with a negative log radius. Here all points plotted in panel d) appear with the same angular coordinate they have in d) but with a radius $r' = -\log(r)$. In this way, typical points that in d) have almost unitary radius now coalesce towards the origin and atypical points, that in d) lie close to zero, are now pushed away from the origin, like the black square. A similar pattern is observed in Figure 6.2e in Chapter 5, where the majority of points converge towards the origin and FBA is seen to be a rather eccentric outlier.

# Resum

El metabolisme cel·lular està format per un gran nombre de reaccions bioquímiques que formen una xarxa densament connectada. Les xarxes metabòliques són les encarregades de que les cèl·lules puguin mantenir-se vives, generant energia química mitjançant diversos processos bioquímics. Com que els fenotips metabòlics apareixen degut a les interaccions entre els diferents components del metabolisme, estudiar aquestes interaccions des d'una perspectiva global és fonamental per entendre els organismes vius i com aquests han evolucionat al llarg dels anys.

Tot i això, l'estudi de xarxes metabòliques completes és difícil degut a la complexitat que es genera quan milers de reaccions s'acoblen i actuen simultàniament. Per facilitar l'estudi del metabolisme com a sistema complex, s'ha desenvolupat l'anomenat *systems-level approach*, que intenta estudiar els sistemes biològics tenint en compte el màxim nombre de constituents coneguts experimentalment. Aquesta manera d'estudiar el metabolisme és la base d'una disciplina emergent anomenada Biologia de Sistemes [10, 11], una branca de la Biologia que últimament té molt de protagonisme per estudiar les causes dels fenòmens fisicoquímics que ocorren en les cèl·lules. Junt a la Biologia de Sistemes, el

desenvolupament de la Ciència de les Xarxes Complexes [14, 15], un camp interdisciplinari que permet tractar sistemes de molts constituents que interaccionen entre ells, ha ajudat a la comprensió de com funciona el metabolisme.

Aquest resum mostra com l'ús combinat tant d'eines provinents de la Biologia de sistemes com de la ciència de xarxes complexes permet extreure i analitzar noves propietats del metabolisme. El resum comença donant les idees bàsiques sobre el metabolisme i les eines que s'usen per analitzar-lo, conceptes que corresponen als dos primers capítols d'aquesta tesi. Després, es resumeixen els quatre capítols de resultats, seguit de les conclusions generals que es poden extreure.

## Què és el metabolisme?

El metabolisme és el conjunt de reaccions químiques que tenen lloc en un organisme i que el mantenen viu. Aquests processos permeten als organismes créixer i reproduir-se, mantenir estructures bioquímiques i respondre al seu medi exterior. El metabolisme es divideix en dues categories: el catabolisme i l'anabolisme. El catabolisme s'encarrega de descompondre la matèria orgànica, com per exemple per extreure energia en la respiració cel·lular. I l'anabolisme, contràriament, utilitza aquesta energia per construir components de les cèl·lules, com ara proteïnes i àcids nucleics.

Les reaccions químiques que es donen en el metabolisme s'organitzen en rutes metabòliques, on les substàncies químiques es transformen en altres mitjançant una seqüència de reaccions catalitzades per enzims. Els enzims són molt importants pel metabolisme, ja que són els responsables de que les reaccions siguin cinèticament possibles, és a dir, actuen de catalitzadors, permetent que les velocitats de les reaccions tinguin lloc en quantitats de temps raonables.

Algunes de les rutes més importants del metabolisme són la Glucòlisi, el Cicle de Krebs, i la Fosforilació Oxidativa. La primera degrada la glucosa mitjançant una sèrie de reaccions a piruvat. Un dels destins més importants del piruvat és generar acetil CoA, el qual és capaç de seguir altres vies metabòliques. Una de les més importants és el Cicle de Krebs, que acaba de transformar, mitjançant oxigen, els productes que provenen de la Glucòlisi a diòxid de carboni. Aquestes rutes generen productes químicament reduïts. En la Fosforilació Oxidativa, els electrons que han anat a parar a aquestes molècules reduïdes són transferits a l'oxigen i l'energia alliberada és utilitzada per crear ATP.

En general, el metabolisme s'ha estudiat utilitzant un enfocament reduccionista, centrat principalment en l'estudi del paper de les biomolècules, la cinètica i la termodinàmica de les reaccions metabòliques. Com a exemple, la termodinàmica de processos irreversibles estudia processos com el transport no espontani a través de la membrana, que aprofita l'energia lliure procedent d'un gradient de protons [38] o de la hidròlisi d'ATP per poder realitzar el transport de components a través de la membrana. Tot i això, hi ha moltes preguntes fetes en Biologia que no es poden respondre amb un enfoc reduccionista. Degut a això, ara s'està desenvolupant el systems-level approach, que té en compte de manera conjunta el màxim nombre de components possibles del metabolisme.

Un enfocament global del metabolisme té en compte tot el conjunt de reaccions bioquímiques i les seves interaccions. Les *genome-scale metabolic networks* [10] proporcionen representacions d'alta qualitat del metabolisme que integren informació bioquímica amb informació genòmica. Un cop s'han validat experimentalment aquests genome-scale models, es poden utilitzar per realitzar prediccions i per a l'anàlisi detallat de les capacitats metabòliques, amb aplicacions en una varietat de camps com la biomedicina o la biotecnologia [57, 58]

# El metabolisme vist com un sistema complex

Les xarxes metabòliques s'han estudiat en el context de les xarxes complexes. Bàsicament, una xarxa complexa [15, 96] és un sistema discret d'elements, també anomenats nodes, que interaccionen entre ells. Les xarxes metabòliques es modelitzen com a xarxes complexes que contenen dos tipus de nodes, anomenades xarxes *bipartites*, i els nodes són els metabòlits i les reaccions. És important tenir en compte que les connexions poden tenir direccionalitat, depenent de si les reaccions són reversibles o irreversibles (veure Figura F.2). Amb la informació sobre la reversibilitat de les reaccions, la xarxa obtinguda és més precisa i descriu el sistema d'una manera més realista. A més, els fluxos de les reaccions determinen els *pesos* de les connexions, cosa que implica que el pes és més gran quan el flux és més gran.

No obstant això, es poden construir representacions més simples de les xarxes metabòliques, anomenades projeccions *one-mode*. En aquestes versions, les xarxes tenen un sol tipus de nodes. Típicament, els metabòlits es trien com a nodes, i es col·loca un enllaç entre dos metabòlits si hi ha almenys una reacció que els connecta. Per contra, si les reaccions són elegides com a nodes, dues reaccions estan connectades si hi ha almenys un metabòlit comú entre elles.

Les xarxes metabòliques estudiades com a xarxes complexes mostren propietats característiques. Els nodes estan caracteritzats pel nombre de veïns, una magnitud anomenada *grau*, $k$, d'un node. Cal tenir en compte que per les xarxes dirigides, és a dir, les que tenen connexions amb direccionalitat, el grau està format per tres contribucions, *entrada*, *sortida* i *bidireccional*. Una mesura important és la distribució de graus $P(k)$, que dóna la probabilitat que un node seleccionat a l'atzar tingui un grau $k$. En la majoria dels casos, els metabòlits mostren un distribució de grau que segueix una llei de potències, $P(k_M) \propto k_M^{-\gamma}$,

FIGURA F.2: Exemple d'una xarxa metabòlica modelitzada com a xarxa complexa. Les reaccions estan representades per quadrats blaus, mentre que els metabòlits estan representats per cercles verds. Les fletxes direccionals són connexions que provenen de reaccions termodinàmicament irreversibles, mentre que fletxes bidireccionals impliquen reaccions termodinàmicament reversibles.

on $\gamma$ és l'exponent característic de la llei de potències. Les xarxes amb una distribució de grau descrit per una llei de potències s'anomenen *scale-free*. Per a les xarxes metabòliques, $\gamma$ és típicament $\sim 2.2$ [97]. El fet que els metabòlits mostrin una distribució que segueix una llei de potències vol dir que hi ha una alta diversitat en els graus dels nodes de la xarxa. Com a conseqüència, la majoria dels metabòlits tenen poques connexions, mentre que pocs metabòlits, anomenats *hubs*, en tenen moltes. Un exemple d'aquests metabòlits altament connectats és l'ATP, ja que participa en moltes reaccions amb la finalitat de subministrar energia lliure quan es realitza el trencament dels seus enllaços. En canvi, les reaccions tenen una distribució de grau punxeguda. El pic està situat al valor mig del grau de totes les reaccions $< k_R >$. Aquesta propietat sorgeix del fet de que les reaccions tenen un nombre limitat de participants, típicament d'entre dos a dotze. El cas més típic és quan les reaccions tenen quatre metabòlits, el que porta a una distribució de grau amb el pic al voltant de $k_R = 4$.

Una altra característica és la propietat anomenada *small-world* [227]. Aquesta propietat significa que és possible anar d'un node a un altre de la xarxa mitjançant pocs salts entre nodes seguint les connexions de la xarxa. Per xarxes metabòliques, a la Referència [97] es calcula la mitjana de les longituds entre nodes a projeccions one-mode de metabòlits, i es troba un valor mig de $\sim 3.2$. Això implica que a les xarxes metabòliques partint d'un metabòlit es pot aconseguir un altre mitjançant un nombre petit de reaccions.

Una altra propietat és que les xarxes poden estar formades per diferents zones (o components) que no estan connectades entre elles. Quan una d'aquestes zones és tant gran que abasta una fracció macroscòpica de la xarxa, la zona en qüestió s'anomena *giant connected component* (GCC). Cal remarcar també que la connectivitat de les xarxes on les connexions tenen direccionalitat presenta característiques especials, ja que el camí entre dos nodes $i$ i $j$ pot ser diferent en passar de $i$ a $j$ o de $j$ a $i$. Aquest fet dóna lloc a l'existència d'una estructura anomenada *bow-tie* [131]. La característica principal d'aquesta estructura és que es pot detectar la presència d'un component anomenat *strongly connected component* (SCC), que és una regió de la xarxa, i en concret que forma part de la GCC, on qualsevol node és accessible des de qualsevol altre, tenint en compte un altre cop que les connexions de la xarxa tenen direccionalitat i que els camins entre els nodes han de satisfer la direcció que imposen les connexions.

Per últim, es creu que les xarxes biològiques estan formades per subconjunts de nodes anomenat *mòduls* [113]. En general, aquesta idea es correspon amb el concepte de comunitats. L'organització d'una xarxa en comunitats no implica que la xarxa estigui fragmentada en diferents components, ja que les comunitats són subconjunts d'una xarxa que contenen un patró d'interconnexió molt dens entre nodes dins de la comunitat, mentre que els nivells d'interconnexió amb els nodes externs a la comunitat són més baixos.

# Flux Balance Analysis

Molts estudis requereixen un càlcul de les velocitats de les reaccions de la xarxa metabòlica. Per calcular els fluxos a través de les reaccions d'una xarxa metabòlica, s'utilitza una tècnica anomenada *Flux Balance Analysis* (FBA). Aquesta tècnica no necessita constants cinètiques per calcular els fluxos metabòlics. Per fer-ho, només usa *constrained-optimization* [63]. Per aplicar FBA es comença per construir la matriu estequiomètrica $\mathbf{S}$ de la xarxa metabòlica, que conté els coeficients estequiomètrics dels metabòlits a les reaccions de la xarxa. Després, es multiplica aquesta matriu pel vector de fluxos $\vec{\nu}$. Aquest producte és, degut al principi de conservació de massa, igual al vector de la variació en el temps de les concentracions $\frac{d\vec{c}}{dt} = \mathbf{S} \cdot \vec{\nu}$. Suposant estat estacionari, $\mathbf{S} \cdot \vec{\nu} = \vec{0}$. Cal remarcar que la suposició d'estat estacionari implica que, a diferència de la cinètica química tradicional, no tractem amb un sistema d'equacions diferencials, sinó amb un simple sistema d'equacions algebraiques.

Atès que en general, les xarxes metabòliques contenen més reaccions que metabòlits, tenim un sistema d'equacions indeterminat. Per tant, es defineix una funció objectiu amb la finalitat d'escollir, d'entre totes les solucions que compleixen les equacions, una solució biològicament significativa. En general, la funció objectiu escollida és la velocitat de creixement de l'organisme. Això significa que tractem de trobar la solució que optimitza el creixement de l'organisme, que és equivalent a maximitzar el flux de la reacció de formació de biomassa, una reacció virtual que s'afegeix la xarxa per simular el creixement. Cal dir també que cal imposar un valor mínim $\nu_{min}$ i màxim $\nu_{max}$ als fluxos de les reaccions, $\nu_{min} \leq \nu_i \leq \nu_{max}$, on $\nu_i$ és el flux d'una reacció qualsevol $i$. Tècnicament, com que tenim un sistema lineal d'equacions amb restriccions lineals, es pot utilitzar *linear programming* per tal de calcular una solució que optimitzi el creixement

en una petita quantitat de temps (de l'ordre d'1 s), cosa que implica que és un mètode computacionalment barat.

# Objectius

Un cop feta una introducció sobre metabolisme i la metodologia emprada, s'especifiquen els objectius de la tesi:

- Estudiar si l'estructura de les xarxes metabòliques ha evolucionat cap a la robustes contra inactivacions de reaccions o gens.

  - Estudi de la inactivació de reaccions individuals i parelles.

  - Estudi de la propagació del dany quan s'inactiven gens.

  - Discussió dels resultats en termes d'una perspectiva evolutiva.

- Analitzar els efectes dels fluxos de les reaccions individuals i de parelles de reaccions individuals usant FBA.

  - Estudi simultani de l'activitat i essencialitat de reaccions individuals.

  - Caracterització dels mecanismes de plasticitat i redundància.

  - Estudi de la dependència en els nutrients de la plasticitat i la redundància.

- Identificar, usant FBA i un mètode de filtratge, les rutes metabòliques amb un paper important per la supervivència dels organismes.

  - Comprovació de l'eficiència del mètode de filtratge en xarxes metabòliques.

  - Anàlisi de les versions filtrades de les xarxes en termes d'evolució a llarg termini.

- – Obtenció de informació sobre l'adaptació a curt termini del metabolisme als medis exteriors dels organismes.

- Avaluar la representativitat de la solució FBA en l'espai complet de solucions metabòliques.

  - – Avaluació de l'excentricitat de la solució FBA respecte a la resta de solucions.

  - – Obtenció d'un esquema de referència per tal de calibrar solucions FBA.

  - – Recuperació de fenotips que no es poden obtenir amb tècniques simples de constrained-optimization.

## Cascades estructurals en xarxes metabòliques

El primer capítol de resultats de la tesi, basat en les Referències [94, 186], estudia com responen les xarxes metabòliques dels organismes *Mycoplasma pneumoniae* [120], *Escherichia coli* [117] i *Staphylococcus aureus* [66] quan les diferents reaccions que les composen són forçades a ser inactives. Es consideren inactivacions de reaccions individuals i parelles de reaccions. A més, només per *M. pneumoniae*, s'estudia l'efecte de les inactivacions de gens i grups de gens. Una inactivació d'un gen causa que els enzims que són codificats per aquest gen no es produeixin, donant com a resultat que les reaccions controlades per aquests enzims esdevinguin no-operatives.

## Mètodes

L'algorisme per calcular cascades estructurals [140] considera que el fet que una reacció s'inactivi desencadena una cascada que es propaga fins que totes les reaccions i metabòlits de la xarxa són capaços de mantenir un estat estacionari. Pels metabòlits, un estat estacionari s'aconsegueix quan aquests participen com a mínim en una reacció que el consumeix i en una reacció que el produeix. Per a les reaccions, es tradueix en que tots els metabòlits que participen en una reacció han de ser capaços de mantenir un estat estacionari. El dany produït per una cascada es quantifica amb el nombre de reaccions que s'han tornat no-operatives.

Els resultats de les cascades en els organismes es comparen amb dos models nuls. Un model nul és una xarxa complexa que s'aconsegueix partint d'una xarxa complexa inicial i aleatoritzant alguna de les seves propietats. Així, s'aconsegueix una nova xarxa que conserva algunes de les propietats estructurals de la xarxa inicial. Per realitzar l'estudi, primer s'obtenen les versions aleatoritzades de les xarxes metabòliques i s'hi aplica l'algorisme de cascada, comparant després els resultats entre les xarxes aleatoritzades i l'original. S'usen dos models nuls, *degree-preserving* (DP) [94, 140, 186] i *mass-balanced* (MB) [141, 186]. El model nul DP obté xarxes on les connexions s'han anat aleatoritzat preservant el grau dels nodes de la xarxa. Aquest mètode serveix per saber si a la xarxa original, el fet de que els nodes de la xarxa mantinguin els mateixos graus minimitza el dany causat per les cascades. L'altre mètode, anomenat mass-balanced, és capaç d'aleatoritzar les xarxes metabòliques complint l'estequiometria de les reaccions, és a dir, s'alteren aleatòriament els reactius o productes només quan aquest canvi dóna lloc a reaccions que compleixen l'estequiometria. Aquest mètode preserva el nombre de metabòlits que participen en cada reacció, i pot discernir si una

propietat prové de la pressió evolutiva, ja que al mantenir l'estequiometria, la limitació fisicoquímica més bàsica, la única font de variació restant provindrà de l'evolució d'un organisme al medi exterior.

## Impacte de la inactivació de reaccions

Primer de tot, s'aplica l'algorisme de cascada i el model nul DP per estudiar com afecta a l'estructura de les xarxa metabòliques dels tres organismes el fet de que una reacció sigui inactivada. És a dir, s'inactiven totes les reaccions de les xarxes metabòliques, una per una, i es calcula el nombre de reaccions que s'han inactivat a causa de que la primera hagi estat forçada a ser inactiva. Així, per cada reacció tenim el dany $d_r$ que ha causat. Als panells esquerra de la Figura F.3 es mostren les distribucions de probabilitat acumulada $P(d'_r \geq d_r)$ que la inactivació d'una reacció $r$ arribi com a mínim a $d_r - 1$ altres reaccions de la xarxa metabòlica de cada organisme. Es duen a terme tests de Kolmogorov-Smirnov [187] per a la comparació amb els models nuls. Tant per *E. coli* com per *S. aureus* es troba que les probabilitats de les cascades calculades a les xarxes originals són incompatibles amb les cascades a les xarxes obtingudes usant el mètode DP. Això significa que la seva organització metabòlica ha evolucionat cap a una major robustesa estructural. D'altra banda, per *M. pneumoniae* no es pot dir que la diferència entre les distribucions de la xarxa original i el model DP sigui estadísticament significativa, tot i que la probabilitat de cascades grans és menor en la xarxa metabòlica original. Això es pot explicar per l'augment de la linealitat de l'estructura de la xarxa metabòlica de *M. pneumoniae*.

Un cop considerades les inactivacions de reaccions individuals, es considera la inactivació de parelles de reaccions. Els panells dreta de la Figura F.3 mostren les distribucions de probabilitat acumulada $P(d'_{rr'} \geq d_{rr'})$ calculades a partir de

FIGURA F.3: Dany en les cascades produïdes per reaccions individuals i parelles de reaccions. a, c, e) Funcions de distribució de probabilitat acumulada dels danys en els organismes *M. pneumoniae*, *E. coli*, i *S. aureus*. Els resultats es comparen amb els danys produïts en les xarxes aleatoritzades amb el mètode DP (100 realitzacions, línia negra discontínua) i amb el mètode MB (100 realitzacions, línia negra contínua). En cada cas, la distribució de color negre sòlid és la mitjana de 100 realitzacions. b, d, f) Comparacions entre les funcions de distribució de probabilitat acumulada i els danys produïts en les xarxes aleatoritzades per a la inactivació de parelles de reaccions.

la inactivació de totes les parelles de reaccions que generen les cascades. Es troba que pels tres organismes, incloent *M. pneumoniae*, la probabilitat de cascades grans provocades per parelles de reaccions és significativament més petita en les xarxes metabòliques originals que en les xarxes obtingudes amb el mètode DP. Això suggereix que l'organització de les xarxes metabòliques ha evolucionat protegint-se contra les inactivacions de més d'una reacció.

## Robustesa vs regulació en les xarxes metabòliques

A la Figura F.3 també es pot veure la comparació dels resultats de cascades amb el model nul MB. Es comprova que l'algorisme de cascada produeix danys més grans a les xarxes originals que a les xarxes obtingudes amb el mètode MB, però cascades més petites en comparació amb les xarxes obtingudes amb el mètode DP. Atès que la mida de les cascades a les xarxes aleatòries fetes amb MB és significativament menor que els de les xarxes reals, es pot dir que la pressió evolutiva ha conduït a cascades més grans. Per tant, cascades grans, afavorides per la pressió evolutiva, poden apuntar al requisit evolutiu de regular grans parts del metabolisme a través de la regulació de petits conjunts de gens. Aquests resultats indiquen que la pressió evolutiva pot afavorir a la regulació eficient del metabolisme a costa de perdre robustesa en front de inactivacions de reaccions o gens.

## Impacte de la inactivació de gens a *M. pneumoniae*

Per tal d'estudiar els efectes metabòlics de mutacions genètiques individuals a *M. pneumoniae*, es simula la inactivació total de les reaccions associades al gen en qüestió. S'obté que els gens amb danys més grans són essencials o condicionalment essencials per a *M. pneumoniae* [120].

Un cop s'ha estudiat l'essencialitat de gens individuals, s'analitza l'efecte quan s'inhibeixen grups de gens co-expressats. Els grups de gens co-expressats a *M. pneumoniae* poden ser identificats a partir de les dades d'expressió de gens en diferents condicions [172] mitjançant tres estratègies: *distance hierarchical clustering* [129], *Infomap* [125] i *recursive percolation* [94]. Es troba que els gens relacionats amb reaccions que generen cascades grans estan aïllats en grups mono-component. Això és sorprenent, ja que, en principi, caldria esperar que els gens que generen cascades grans s'expressaran junt a altres gens, ja que, en general, controlar parts grans del metabolisme requereix l'expressió de varis gens. El fet que l'expressió d'aquests gens estigui aïllada fa que es puguin identificar com a importants regulador metabòlics, ja que l'alteració d'un sol gen pot afectar a un gran nombre de reaccions metabòliques. En qualsevol cas, l'expressió individual d'aquests gens és de nou una indicació que l'organització estructural de l'organisme ha evolucionat per protegir el sistema contra inactivacions de múltiples reaccions.

# Efecte de les inactivacions de reaccions en estats estacionaris del metabolisme

El segon capítol de resultats, basat en les Referències [158, 205], estén l'estudi estructural de pertorbacions a càlculs dinàmics, més precisament al càlcul de fluxos metabòlics obtinguts mitjançant la tècnica Flux Balance Analysis. Aquest estudi permet, per una banda, predir si hi ha reaccions importants que han de ser sempre actives per tal de garantir la supervivència d'un organisme i, d'altra banda, per comprovar si el metabolisme ha desenvolupat mecanismes de protecció quan algunes dels seus constituents són inactivats. Aquest últim estudi es porta a terme amb un anàlisi de parelles de reaccions anomenades

*synthetic lethal* (SL), que són parelles de reaccions que la seva inactivació és letal, però de manera que la inactivació individual de les reaccions que formen la parella no ho és. Això permet identificar dos mecanismes diferents que el metabolisme ha desenvolupat per protegir-se contra possibles inactivacions de reaccions, anomenats

it plasticitat i *redundància*. L'estudi es realitza en dos bacteris: *E. coli* [119] i *M. pneumoniae* [56]. Per *E. coli* s'analitza en detall l'estudi de inactivacions de reaccions individuals i parelles de reaccions, mentre que per *M. pneumoniae* s'estudien bàsicament parelles de reaccions.

## Mètodes

La tècnica FBA s'usa per saber si una reacció tindrà un flux nul o no nul en funció dels nutrients presents en el medi, és a dir, per saber l'activitat d'una reacció. Bàsicament, si una reacció té un flux no nul l'activitat serà no nul·la, mentre que si el flux és nul l'activitat serà nul·la. FBA també serveix per saber si un organisme podrà sobreviure a la inactivació d'una reacció, és a dir, per determinar l'essencialitat d'una reacció. Per fer això, s'imposa que la reacció a inactivar ha de tenir per força un flux nul, i es computa la velocitat de creixement de l'organisme. Si aquesta és nul·la, l'organisme es considera mort i la reacció és essencial. En canvi, si la velocitat de creixement no és nul·la, es pot considerar que l'organisme ha sobreviscut a la inactivació i, per tant, la reacció no és essencial.

## Activitat i essencialitat de les reaccions de *E. coli*

L'estudi de l'anàlisi de l'activitat i essencialitat de les reaccions de *E. coli* permet dividir les reaccions en quatre categories de reaccions, en funció dels valors

d'activitat i essencialitat que s'obtenen. La primera són les *reaccions sempre essencials quan són actives*. Donada la seva importància, aquestes reaccions es poden seleccionar com a dianes de fàrmacs, ja que són components fonamentals del metabolisme de l'organisme *E. coli*, degut a que algunes han de ser sempre actives per assegurar la vida de l'organisme, i les altres, tot i no ser sempre actives, sempre són essencials quan són actives. Una altra categoria són les *reaccions sempre actives*. Aquestes reaccions són sempre actives a tots els medis, però són essencials en una fracció dels medis. La categoria *reaccions mai essencials* són reaccions que són actives en una fracció de medis, i mai són essencials encara que siguin actives. Les reaccions de les categories *reaccions sempre actives* i *reaccions mai essencials* són actives per tal d'augmentar la velocitat de creixement de l'organisme i per prevenir la formació de parelles SL. L'última categoria de reaccions que s'obté són les *reaccions parcialment essencials*. Aquestes reaccions tenen la propietat de que el fet de que siguin actives, no implica ser també essencials. Per tant, l'extrapolació de l'activitat a l'essencialitat, un fet bastant comú, no és correcta en base a aquest anàlisi.

## Classificació de les parelles SL en plàstiques i redundants

Un cop considerades les reaccions individuals, s'estudien parelles de reaccions SL per tal d'analitzar els mecanismes de plasticitat i redundància. El primer tipus són les parelles SL plàstiques, on una de les reaccions de la parella té un flux FBA no nul mentre que la segona reacció té un flux FBA nul. El segon tipus de parelles SL s'anomenen parelles SL redundants, on les dues reaccions tenen fluxos FBA no nuls.

Per l'organisme *E. coli* es troba que la majoria de parelles SL són plàstiques. Per aquestes parelles, quan s'inactiva la reacció activa de la xarxa metabòlica, els

FIGURA F.4: Representació esquemàtica de les parelles SL plàstiques i redundants. Els metabòlits estan representats per cercles i les reaccions per quadrats. Les reaccions acolorides amb fletxes negres representen reaccions actives, mentre que les línies discontínues grises s'utilitzen per a les reaccions i metabòlits que són inactius. Els negres denoten nodes que han sigut desactivats. La reacció de producció de biomassa es representa com un quadrat més gran amb un flux associat $\nu_g$. Quan aquesta reacció està inactiva, l'organisme es considera mort.

fluxos es reorganitzen per tal que la reacció amb flux nul de la parell s'activi. Això implica que la reacció inicialment no activa actua com una còpia de seguretat de la reacció activa, garantitzant la viabilitat de l'organisme (veure Figura F.4). Pel que fa a les parelles SL redundants, per la gran majoria l'ús simultani d'ambdues reaccions augmenta la velocitat de creixement en comparació amb la situació en que només una de les reaccions és activa (veure Figura F.4). Per *M. pneumoniae*, es troba que la freqüència de les parelles de reaccions SL redundants és més gran que a *E. coli*, degut a l'augment de la linealitat i la reducció de la complexitat de l'organisme, i la seva freqüència s'assembla a la de plasticitat. En conseqüència, l'equilibri de parelles SL redundants vs parelles SL plàstiques és diferent respecte *E. coli*.

**Sensitivitat a canvis en el medi exterior de les parelles SL**

L'última part d'aquesta secció analitza si la plasticitat o la redundància d'una parella SL canvia en funció dels nutrients del medi exterior. D'aquesta manera, s'investiga la sensibilitat de les parelles SL de *E. coli* als canvis en la composició dels nutrients del medi exterior. Es troba que per a la majoria de parelles, l'essencialitat no és específica d'un medi mínim, i només un nombre petit de parelles mostra especificitat en el medi mínim. Per a cada parella de reaccions SL, es troba que la gran majoria de parelles són quasi sempre plàstiques, un petit nombre són sempre redundants, i un altre petit nombre mostra un comportament que varia entre redundància i plasticitat. Per tant, les parelles de reaccions SL i la seva divisió en plasticitat i redundància són aspectes altament conservats independentment de la composició del medi exterior.

# Detecció d'empremtes d'evolució i adaptació en xarxes metabòliques

L'últim capítol de resultats, basat en la Referència [221], analitza la importància dels fluxos metabòlics amb la finalitat d'extreure informació biològica relativa a tendències adaptatives i evolutives. Per realitzar aquest anàlisi, s'utilitza un mètode de filtratge anomenat *disparity filter* [95]. La seva aplicació permet disminuir el nombre d'enllaços d'una xarxa metabòlica, mantenint només aquells que són estadísticament importants, obtenint així els anomenats *backbones* metabòlics. L'estudi d'aquestes backbones metabòlics és important per tal de quantificar la quantitat de vegades que una ruta bioquímica està present en els backbones, i per tal de relacionar la seva presència amb el seu paper en

l'adaptació de l'organisme al medi exterior. L'estudi es realitza pels organismes *E. coli* [119] i *M. pneumoniae* [56].

## Mètodes

El disparity filter funciona mitjançant la comparació dels pesos dels enllaços amb un model nul aleatori [95]. Després, es preserven els enllaços si la probabilitat que els pesos dels enllaços no siguin compatibles amb el model nul és menor que un llindar triat $\alpha$. D'aquesta manera, el valor del llindar $\alpha$, que es pot moure en l'interval [0,1], determina la intensitat de filtrat. Si $\alpha \to 1$, la intensitat de filtratge serà petita. D'altra banda, si $\alpha \to 0$, la intensitat de filtratge serà molt gran, és a dir, el pes de l'enllaç haurà de ser estadísticament molt rellevant per poder passar el filtratge. Una condició imprescindible perquè el disparity filter sigui capaç de filtrar amb eficiència és que la funció de distribució de probabilitat dels pesos dels enllaços sigui heterogènia, és a dir, contingui valors de diversos ordres de magnitud tant a nivell global, que es pot comprovar calculant la funció de distribució de probabilitat, com a nivell local, que implica que els el conjunt de pesos tant d'entrada com de sortida de cada node de la xarxa és heterogeni.

## Identificació dels backbones de les xarxes metabòliques

Primer es comprova el rendiment del disparity filter en xarxes metabòliques. Per fer-ho, es porta a terme una exploració del paràmetre $\alpha$ per tal d'obtenir com les xarxes metabòliques d'ambdós *E. coli* i *M. pneumonaie* es redueixen en funció del valor de $\alpha$. Cal tenir en compte que, un cop filtrada, i per evitar treballar amb reaccions estequiomètricament no equilibrades, la xarxa obtinguda es transforma en una projecció one-mode. Per fer-ho, es connecten dos metabòlits amb un enllaç dirigit si hi ha una reacció que el seu flux és al mateix temps rellevant per
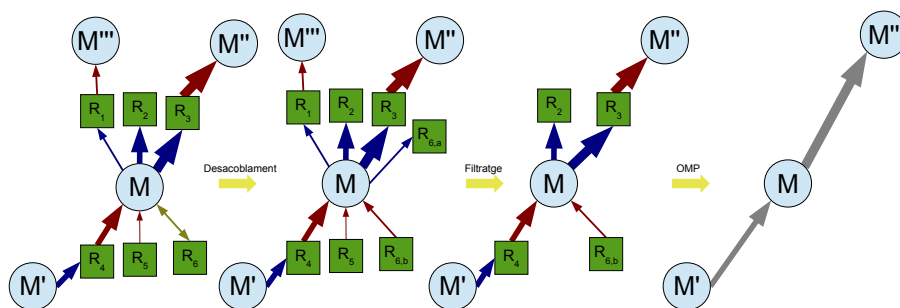
FIGURA F.5: Esquema del mètode de filtratge. Els nodes blaus són metabòlits i els quadrats verds denoten reaccions. Les connexions d'entrada dels metabòlits estan representades per fletxes vermelles, les connexions de sortida per fletxes blaves i les connexions reversibles amb fletxes de color groc fosc.

al consum d'un metabòlit i també rellevant per la producció de l'altre (veure Figura F.5). En aquesta projecció one-mode, es calcula el nombre d'enllaços $E$ i de nodes $N$ que romanen després de filtrar. S'obté que, tot i que la fracció d'enllaços es pot reduir en gran mesura, la fracció de nodes es manté bàsicament inalterada en comparació amb el cas no filtrat. El punt on comença a decaure la fracció de nodes es defineix com un valor crític, i es pot entendre com un punt òptim que redueix les connexions presents a la xarxa però que conté el màxim d'informació bioquímica i estructural possible.

## Signes d'evolució en els backbones de metabòlits

Els backbones metabòlics d'ambdós *E. coli* i *M. pneumoniae* s'analitzen en relació a una perspectiva evolutiva a llarg termini. El backbone de *E. coli* s'analitza en termes de l'estructura bow-tie, i conté una GCC formada per tres SCC. És interessant identificar a quines rutes metabòliques pertanyen les connexions entre metabòlits de les SCCs. A la SCC més gran es pot veure que la major part dels enllaços corresponen al metabolisme energètic. D'aquesta manera, es recuperen les rutes involucrades amb l'obtenció d'energia que han

mostrat processos evolutius per maximitzar l'eficiència en la producció d'ATP [224, 226], i que degut a la seva importància, han estat presents des dels primers moments de la vida dels organismes. Les altres dues SCC també estan formades per enllaços que pertanyen a rutes metabòliques que ja existien en els primers moments de vida dels organismes. Les rutes presents a aquestes SCC són la biosíntesi de purina i pirimidina, el metabolisme de lípids i la biosíntesi de cofactors i grups prostètics. De fet, s'ha trobat que la síntesi de purines i pirimidines va ser la primera ruta metabòlica que implicava l'ús d'enzims [230] [230]. El metabolisme de lípids subministra els lípids necessaris per generar la membrana cel·lular, i s'ha demostrat que presenten diferències entre els diferents llinatges en els organismes [231].

Per *M. pneumoniae*, el seu backbone està format per una GCC que conté dues SCC, tot i que una d'elles és irrellevant. En la SCC més rellevant, la major part dels enllaços estan relacionats també amb el metabolisme de l'energia, igual que *E. coli*. La diferència principal de *M. pneumoniae* és que en el metabolisme de l'energia hi participen només processos de fermentació, presents des dels primers moments de vida dels organismes [35], degut a l'absència del Cicle de Krebs en aquest organisme.

## Els backbones metabòlics de *E. coli* codifiquen la seva capacitat d'adaptació a curt termini

La última part d'aquesta secció estudia com els canvis en la composició dels nutrients modifiquen el backbone metabòlic de *E. coli*, amb la finalitat d'extreure resultats relacionats amb l'adaptació a curt termini a diferents medis exteriors. Es comença analitzant com el backbone metabòlic de *E. coli* es modifica quan s'utilitza un medi nutricional ric anomenat Luria-Bertani Broth [154, 158]. La

GCC del backbone conté tres SCC, tot i que dos d'elles són irrellevants. La majoria de les connexions de la SCC rellevant corresponen a la ruta metabòlica encarregada del metabolisme dels lípids [232].

Finalment, s'estudia com es veu afectat el backbone metabòlic de *E. coli* a causa dels canvis en la composició dels nutrients presents en el medi exterior. Amb aquesta finalitat, es consideren els medis mínims donats a la Referència [119]. Per a cada medi mínim, s'analitza com es modifica el valor crític $\alpha_c$. S'obté que els valors de $\alpha_c$ són semblants, és a dir, s'obté un valor característic en tots els medis, cosa que indica que les solucions de fluxos metabòlics de la major part dels medis mínims és molt similar. Aquest fet suggereix la possibilitat de fusionar tots els backbones metabòlics de cada medi obtenint un *superbackbone* global. La GCC d'aquest superbackbone està composada per una SCC, a més de tres petites SCC irrellevants. Un anàlisi de les rutes metabòliques presents en la SCC indica que la ruta més abundant és en aquest cas el Metabolisme del Carboni Alternatiu [234], una ruta transversal [7] que conté gens que la seva expressió depèn dels estímuls externs, en particular en l'alteració de les fonts de carboni [233].

## Estudi dels estats òptims en l'espai FFP

Aquesta part, basada en la Referència [245], revisa la tècnica FBA en relació amb el conjunt de tots els estats de fluxos possibles en una xarxa metabòlica mínima de *E. coli* (*E. coli core*) [63, 118]. FBA utilitza la suposició que els organismes tracten de créixer tant com sigui possible. Degut a que aquesta suposició pot no ser sempre correcta, és important explorar la resta de possibles solucions sense fer ús de la hipòtesi de creixement màxim. Això permet avaluar

la rellevància de la solució obtinguda per FBA en comparació amb totes les altres possibles solucion.

## Mètodes

Per calcular totes les possibles solucions en una xarxa metabòlica, s'usa el mètode Hit-and-Run (HR). Molt resumidament, aquesta tècnica és capaç d'obtenir totes les solucions que compleixen les condicions FBA, sense optimitzar la velocitat de creixement. Cal dir també que, tot i que el mètode HR dóna una mostra de l'espai i no es calcula l'espai complet, s'ha demostrat que aquesta mostra és representativa de l'espai complet de solucions [246].

## La solució FBA és excèntrica respecte de l'espai FFP

L'espai de solucions que conté totes les solucions de fluxos metabòlics s'anomena *feasible flux phenotypes space* (espai FFP). En aquest cas, el medi exterior és un medi mínim on la font de carboni és glucosa.

Per tal d'avaluar la representativitat de la solució FBA en relació amb les altres solucions de l'espai FFP, es calculen primer els valors mitjos dels fluxos de cada reacció a l'espai FFP i es compara amb els valors que prediu FBA optimitzant el creixement. Fent això, s'obté que la majoria dels valors mitjos de les reaccions queden lluny del que prediu FBA. Per visualitzar l'excentricitat de la solució FBA s'utilitza un anàlisi de components principals [247, 248], amb la finalitat de reduir l'alta dimensionalitat de l'espai FFP ($n_{reaccions} = 70$) i així poder representar l'excentricitat en un gràfic de dues dimensions. Usant aquesta tècnica, es computa la matriu de covariàncies dels fluxos de les reaccions de l'espai FFP, i d'allà es calculen els dos primers vectors principals, que són els que

FIGURA F.6: Projecció de l'espai FFP en els dos components principals de la matriu de correlació. La majoria dels punts es troben en un cercle prop de l'origen (l'àrea més fosca). El cercle verd representa la solució FBA.

contenen la major variabilitat de l'espai. Cada estat de flux metabòlic possible es reescala com un *z-value* centrat al voltant del valor mig i es projecta sobre els dos eixos. Les projeccions es representen en coordenades polars, amb una transformació logarítmica negativa a la coordenada radial original per facilitar la visualització. A la Figura F.6 es pot veure que, clarament, la solució prevista optimitzant el creixement amb FBA és excèntrica comparant-la a la resta de solucions.

## L'espai FFP es pot usar per calibrar la desviació de fenotips òptims respecte resultats experimentals

Aquesta part es centra en la relació entre el flux d'entrada de diferents fonts de carboni (glucosa, piruvat i succinat) i el flux d'entrada d'oxigen, per mostrar que l'espai FFP es pot usar com a punt de referència per calibrar la desviació de les prediccions fetes amb FBA optimitzant el creixement amb els resultats experimentals. Per fer-ho, primer es calculen els fluxos d'entrada de cada font de carboni i d'oxigen de dues maneres, mitjançant FBA optimitzant el creixement,

i calculant els valors migs dels fluxos d'entrada dels metabòlits a l'espai FFP. Un cop es tenen els valors, es comparen amb les dades experimentals reportades pel consum d'oxigen en un medi mínim amb glucosa, piruvat o succinat com a fonts de carboni primàries.

En tots els casos, les prediccions optimitzant el creixement amb FBA reprodueixen bé els punts de les dades experimentals en la regió de valors petits de fluxos d'entrada de totes les fonts de carboni [249]. No obstant això, el flux d'entrada d'oxigen es satura quan s'arriba a fluxos d'entrada grans de les fonts de carboni. En aquesta regió de fluxos d'entrada grans de les fonts de carboni, la solució FBA maximitzant el creixement prediu un consum d'oxigen excessiu, al voltant d'un 25% respecte els valors reportats experimentalment.

## Una regió de l'espai FFP mostra fermentació aeròbica

La regió metabòlica d'alt creixement de l'espai FFP de l'organisme *E. coli* es pot calcular en medi mínim de glucosa posant un valor mínim pel flux de la reacció de biomassa. En aquesta regió, es pot identificar la utilització de rutes metabòliques típiques del metabolisme microbià proliferatiu[1], tot i que es considera un flux màxim del flux d'entrada d'oxigen il·limitat i, per tant, permet arribar a fluxos d'entrada d'oxigen molts grans. Aquest comportament metabòlic és consistent amb les dades experimentals [49, 249, 255] però és inabastable per càlculs FBA basats en els típics principis de optimalitat.

D'acord amb l'espai FFP de *E. coli*, es pot observar que la regió d'alt creixement de l'espai FFP es caracteritza per la secreció de molècules orgàniques típiques de la fermentació fins i tot quan el subministrament d'oxigen és il·limitat. Aquest

---

[1]El metabolisme microbià proliferatiu produeix productes de fermentació com l'acetat o l'etanol fins i tot en condicions aeròbiques.

fet apunta a la utilització simultània de la Glucòlisi i la Fosforilació Oxidativa per produir biomassa i energia. Algunes de les molècules orgàniques que s'expulsen són el formiat, l'acetat, l'etanol i el lactat. Aquests resultats indiquen que, contràriament a les prediccions FBA, que no dóna la producció d'aquestes molècules si no hi ha limitacions extres [254], un alt flux d'entrada de glucosa combinat amb prou oxigen pot mantenir els requisits del metabolisme proliferatiu per formar biomassa a través de la fermentació aeròbica. Aquesta fermentació aeròbica, que hom pot pensar que és ineficient en termes de rendiment energètic en comparació amb la Fosforilació Oxidativa, s'ha demostrat que és un estat catabòlic favorable per totes les cèl·lules que proliferen ràpidament amb un flux alt d'entrada de glucosa [254]. Per això, es pot dir que és un fenotip metabòlic probable.

## Conclusions

Les conclusions que es poden extreure d'aquesta tesi són:

- L'estructura de les xarxes metabòliques dels organismes estudiats ha evolucionat cap a la robustesa contra inactivacions de reaccions.

- L'essencialitat dels gens de *M. pneumoniae* amb l'algorisme de cascada coincideix amb resultats experimentals. A més, la regulació de gens associats a danys grans tendeix a donar-se de forma isolada. Això es pot entendre com a mecanisme de protecció per evitar danys metabòlics grans.

- La pressió evolutiva afavoreix la regulació metabòlica eficient a canvi de perdre un cert grau de robustesa.

- L'estudi amb FBA permet descobrir que hi ha un conjunt de reaccions que han de ser sempre actives per tal de garantir la viabilitat d'un organisme. Les reaccions no-essencials poden formar part de combinacions SL.

- Les parelles SL es poden classificar en funció de dos mecanismes diferents, la plasticitat, on una reacció és activa l'altra inactiva, i la redundància, on les dues reaccions tenen fluxos diferents de zero. La relació entre plasticitat i redundància es fortament organisme depenent.

- La plasticitat i la redundància són altament conservades encara que es canviïn els nutrients del medi exterior.

- El disparity filter és eficient per filtrar xarxes metabòliques.

- L'estudi de les SCC dels backbones metabòlics en medi mínim permet identificar quines rutes metabòliques han jugat un paper important en l'evolució a llarg termini.

- L'estudi dels backbones en diferents medis exteriors permet identificar rutes metabòliques que mostren adaptació a curt termini.

- Els estats de creixement òptims són excèntrics respecte els altres fenotips que formen l'espai complet.

- L'espai FFP es pot usar per calibrar la desviació entre càlculs FBA i observacions experimentals.

- Sense la necessitat d'afegir restriccions extres, l'espai FFP recupera estats que mostren fermentació aeròbica.

# Bibliography

[1] U. Lagerkvist. *The enigma of ferment: From the philosopher's stone to the first biochemical Nobel Prize.* 2005.

[2] K. E. Van Holde. *Physical Biochemistry, Foundations of Modern Biochemistry Series.* 1971.

[3] D. Segura, R. Mahadevan, K. Juárez, and D. R. Lovley. Computational and experimental analysis of redundancy in the central metabolism of *Geobacter sulfurreducens. PLoS Comput. Bio.*, 4(2):e36, 2008.

[4] C. L. Tucker and S. Fields. Lethal combinations. *Nat. Genet.*, 35:204–205, 2003.

[5] A. L. Barabási and Z. N. Oltvai. Network biology: understanding the cells functional organization. *Nat. Rev. Genet.*, 5:101–113, 2004.

[6] Y. Deville, D. Gilbert, J. Van Helden, and S. J. Wodak. An overview of data models for the analysis of biochemical pathways. *Brief. Bioinform.*, 4(3):246–259, 2003.

[7] M. Á. Serrano and M. Boguñá and F. Sagués. Uncovering the hidden geometry behind metabolic networks. *Mol. BioSyst.*, 8:843–850, 2012.

[8] H. Kitano. Systems biology: a brief overview. *Science*, 295(5560):1662–1664, 2002.

[9] H. Kitano. Computational Systems Biology. *Nature*, 420(6912):206–210, 2002.

[10] B. Ø. Palsson. *Systems Biology: Properties of Reconstructed Networks*. Cambridge University Press, 2006.

[11] U. Alon. *An Introduction to Systems Biology: Design Principles of Biological Circuits*. CRC press, 2006.

[12] E. Almaas, B. Kovacs, T. Vicsek, Z. N. Oltvai, and A. L. Barabási. Global organization of metabolic fluxes in the bacterium *Escherichia coli*. *Nature*, 427(6977):839–843, 2004.

[13] E. Borenstein, M. Kupiec, , M. W. Feldman, and E. Ruppin. Large-scale reconstruction and phylogenetic analysis of metabolic environments. *Proc. Natl. Acad. Sci. USA*, 105(38):14482–14487, 2008.

[14] A. Barrat, D. Barthélemy, and A. Vespignani. *Dynamical Processes on Complex Networks*. Cambridge University Press, 2008.

[15] M. Newman. *Networks: an introduction*. Oxford University Press, 2010.

[16] C. K. Mathews, K. E. Van Holde, and K. G. Ahern. *Bioquímica*. Pearson Education, 2002.

[17] H. Lineweaver and D. Burk. The determination of enzyme dissociation constants. *J. Am. Chem. Soc.*, 56(3):658–666, 1934.

[18] M. Emmerling et al. Metabolic flux responses to pyruvate kinase knockout in *Escherichia coli*. *J. Bacteriol.*, 184(1):152–164, 2002.

[19] T. L. Steck. The organization of proteins in the human red blood cell membrane. A Review. *J. Cell Biol.*, 62(1):1–19, 1974.

[20] G. Wu. Amino acids: metabolism, functions, and nutrition. *Amino Acids*, 37(1):1–17, 2009.

[21] L. C. Fitzpatrick. Life history patterns of storage and utilization of lipids for energy in amphibians. *Amer. Zool.*, 16(4):725–732, 1976.

[22] Y. Nishizuka. Protein kinase C and lipid signaling for sustained cellular responses. *FASEB J.*, 9(7):484–496, 1995.

[23] S. J. Marrink and H. J. Berendsen. Simulation of water transport through a lipid membrane. *J. Phys. Chem.*, 98(15):4155–4168, 1994.

[24] T. Bauchop and S. R. Elsden. The growth of micro-organisms in relation to their energy supply. *J. Gen. Microiol.*, 23(3):457–469, 1960.

[25] C. A. Good, H. Kramer, and M. Somogyi. The determination of glycogen. *J. Biol. Chem.*, 100(2):485–491, 1933.

[26] S. Neidle. *Principles of nucleic acid structure*. Academic Press, 2010.

[27] J. D. Carver and W. Allan Walker. The role of nucleotides in human nutrition. *J. Nutr. Biochem.*, 6(2):58–72, 1995.

[28] J. Bergman. ATP: The perfect energy currency for the cell. *Creation Research Society Quarterly*, 36(1):2–9, 1999.

[29] P. Belenky, K. L. Bogan, and C. Brenner. $NAD^+$ metabolism in health and disease. *Trends Biochem. Sci.*, 32(1):12–19, 2007.

[30] N. Pollak, C. Dolle, and M. Ziegler. The power to reduce: pyridine nucleotides-small molecules with a multitude of functions. *Biochem. J.*, 402:205–218, 2007.

[31] G. M. Gadd. Heavy metal accumulation by bacteria and other microorganisms. *Experientia*, 46(8):834–840, 1990.

[32] J. Ariño, J. Ramos, and H. Sychrová. Alkali metal cation transport and homeostasis in yeasts. *Microbiol. Mol. Biol. Rev.*, 74(1):95–120, 2010.

[33] H. Sychrova. Yeast as a model organism to study transport and homeostasis of alkali metal cations. *Physiol. Res.*, 53:S91–98, 2004.

[34] H. Chen, M. Ikeda-Saito, and S. Shaik. Nature of the Fe-$O_2$ bonding in oxy-myoglobin: effect of the protein. *J. Am. Chem. Soc.*, 130(44): 14778–14790, 2008.

[35] M. Tadege, I. Dupuis, and C. Kuhlemeier. Ethanolic fermentation: new functions for an old pathway. *Trends Plant Sci.*, 4(8):320–325, 1999.

[36] M. Criado-Sancho and J. Casas-Vázquez. *Termodinámica química y de los procesos irreversibles*. 1998.

[37] B. A. Moffatt and H. Ashihara. Purine and pyrimidine nucleotide synthesis and metabolism. In *The Arabidopsis book*, pages 1–20. Amercican Society of Plant Physiologists, 2002.

[38] W. R. Harvey, M. Cioffi, J. A. Dow, and M. G. Wolfersberger. Potassium ion transport ATPase in insect epithelia. *J. Exp. Biol.*, 106(1):91–117, 1983.

[39] L. Michaelis and M. L. Menten. Die Kinetik der Invertinwirkung. *Biochem. Z*, 49(333-369):352, 1913.

[40] A. V. Hill. The possible effects of the aggregation of the molecules of haemoglobin on its dissociation curves. *J. Physiol. (Lond)*, 40:4–7, 1910.

[41] R. K. Crane, D. Miller, and I. Bihler. The restrictions on possible mechanisms of intestinal active transport of sugars. *Membrane transport and metabolism*, 439, 1961.

[42] I. Prigogine. *Thermodynamics of irreversible processes*. Thomas, 1955.

[43] Y. Demirel and S. I. Sandler. Thermodynamics and bioenergetics. *Biophys. chem.*, 97(2):87–111, 2002.

[44] B. Desvergne, L. Michalik, and W. Wahli. Transcriptional regulation of metabolism. *Physiol. Rev.*, 86:465–514, 2006.

[45] D. R. Matthews, J. P. Hosker, A. S. Rudenski, B. A. Naylor, D. F. Treacher, and R. C. Turner. Homeostasis model assessment: insulin resistance and $\beta$-cell function from fasting plasma glucose and insulin concentrations in man. *Diabetologia*, 28(7):412–419, 1985.

[46] G. E. Lienhard, J. W. Slot, D. E. James, and M. M. Mueckler. How cells absorb glucose. *Sci. Am.*, 266(1):86–91, 1992.

[47] D. Fell. *Understanding the control of metabolism*. Portland Press, 1997.

[48] D. L. Theobald. A formal test of the theory of universal common ancestry. *Nature*, 465(7295):219–222, 2010.

[49] R. U. Ibarra, J. S. Edwards, and B. Ø. Palsson. *Escherichia coli* K-12 undergoes adaptive evolution to achieve *in silico* predicted optimal growth. *Nature*, 420:186–189, 2002.

[50] S. Fong, J. Y. Marciniak, and B. Ø. Palsson. Description and interpretation of adaptive evolution of *Escherichia coli* K-12 MG1655 by using a genome-scale in silico metabolic model. *J. Bacteriol.*, 185(21):6400–6408, 2003.

[51] S. Schmidt, S. Sunyaev, P. Bork, and T. Dandekar. Metabolites: a helping hand for pathway evolution? *Trends Biochem. Sci.*, 28(6):336–341, 2003.

[52] E. Kay, T. M. Vogel, F. Bertolla, R. Nalin, and P. Simonet. *In situ* transfer of antibiotic resistance genes from transgenic (transplastomic) tobacco plants to bacteria. *Appl. Environ. Microb.*, 68(7):3345–3351, 2002.

[53] O. J. Rando and K. J. Verstrepen. Timescales of genetic and epigenetic inheritance. *Cell*, 128(4):655–668, 2007.

[54] A. K. Lancaster and J. Masel. The evolution of reversible switches in the presence of irreversible mimics. *Evolution*, 63(9):2350–2362, 2009.

[55] J. G. Lawrence. Common themes in the genome strategies of pathogens. *Curr. Opin. Genet. Dev.*, 15(6):584–588, 2005.

[56] J. A. H. Wodke et al. Dissecting the energy metabolism in *Mycoplasma pneumoniae* through genome-scale metabolic modeling. *Mol. Syst. Biol.*, 9:653, 2013.

[57] M. A. Oberhardt, B. Ø. Palsson, and J. A. Papin. Applications of genome-scale metabolic reconstructions. *Mol. Syst. Biol.*, 5(1), 2009.

[58] T. Y. Kim, S. B. Sohn, Y. B. Kim, W. J. Kim, and S. Y. Lee. Recent advances in reconstruction and applications of genome-scale metabolic models. *Curr. Opin. Biotech.*, 23(4):617–623, 2012.

[59] I. Thiele and B. Ø. Palsson. A protocol for generating a high-quality genome-scale metabolic reconstruction. *Nat. Protoc.*, 5:93–121, 2010.

[60] M. Kanehisa and S. Goto. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.*, 28(1):27–30, 2000.

[61] R. Caspi et al. The metacyc database of metabolic pathways and enzymes and the biocyc collection of pathway/genome databases. *Nucl. Acids Res.*, 40(D1):D742–D753, 2012.

[62] J. Schellenberger, J. O. Park, T. C. Conrad, and B. Ø. Palsson. BiGG: a biochemical genetic and genomic knowledgebase of large scale metabolic reconstructions. *BMC Bioinformatics*, 11:213, 2010.

[63] J. D. Orth, I. Thiele, and B. Ø. Palsson. What is flux balance analysis? *Nat. Biotechnol.*, 28:245–248, 2010.

[64] J. S. Edwards and B. Ø. Palsson. Systems properties of the haemophilus influenzaerd metabolic genotype. *J. Biol. Chem.*, 274(25):17410–17416, 1999.

[65] N. C. Duarte, M. J. Herrgard, and B. Ø. Palsson. Reconstruction and validation of *Saccharomyces cerevisiae i*ND750, a fully compartmentalized genome-scale metabolic model. *Genome Res.*, 14:1298–1309, 2004.

[66] S. A. Becker and B. Ø. Palsson. Genome-scale reconstruction of the metabolic network in *Staphylococcus aureus* N315: an initial draft to the two-dimensional annotation. *BMC Microbiol.*, 5:8, 2005.

[67] I. Thiele, T. D. Vo, N. D. Price, and B. Ø. Palsson. Expanded metabolic reconstruction of *Helicobacter pylori* (*i*IT341 GSM/GPR): an *in silico* genome-scale characterization of single- and double-deletion mutants. *J. Bacteriol.*, 187:5818–5830, 2005.

[68] A. M. Feist, J. C. M. Scholten, B. Ø. Palsson, F. J. Brockman, and T. Ideker. Modeling methanogenesis with a genome-scale metabolic reconstruction of *Methanosarcina barkeri*. *Mol. Syst. Biol.*, 2:2006.0004, 2006.

[69] N. Jamshidi and B. Ø. Palsson. Investigating the metabolic capabilities of *Mycobacterium tuberculosis* H37Rv using the *in silico* strain *i*NJ661 and proposing alternative drug targets. *BMC Syst Biol.*, 1:26, 2007.

[70] A. M. Feist and B. Ø. Palsson. The growing scope of applications of genome-scale metabolic reconstructions using *Escherichia coli. Nat. Biotechnol.*, 26(6):659–667, 2008.

[71] D. McCloskey, B. Ø. Palsson, and A. M. Feist. Basic and applied uses of genome-scale metabolic network reconstructions of *Escherichia coli. Mol. Syst. Biol.*, 9(1), 2013.

[72] N. Fernández, E. E. Díaz, R. Amils, and J. L. Sanz. Analysis of microbial community during biofilm development in an anaerobic wastewater treatment reactor. *Microb. Ecol.*, 56(1):121–132, 2008.

[73] A. R. Zomorrodi and C. D. Maranas. OptCom: a multi-level optimization framework for the metabolic modeling and analysis of microbial communities. *PLoS Comp. Biol.*, 8(2):e1002363, 2012.

[74] A. Bordbar, N. E. Lewis, J. Schellenberger, B. Ø. Palsson, and N. Jamshidi. Insight into human alveolar macrophage and *M. tuberculosis* interactions via metabolic reconstructions. *Mol. Syst. Biol.*, 6(1), 2010.

[75] A. Persidis. High-throughput screening. *Nat. Biotechnol.*, 16(5):488, 1998.

[76] J. E. Bailey, S. Birnbaum, J. L. Galazzo, C. Khosla, and J. V. Shanks. Strategies and challenges in metabolic engineering. *Ann. NY Acad. Sci.*, 589(1):1–15, 1990.

[77] J. H. Park and S. Y. Lee. Towards systems metabolic engineering of microorganisms for amino acid production. *Curr. Opin. Biotechnol.*, 19 (5):454–460, 2008.

[78] R. M. Zelle et al. Malic acid production by *Saccharomyces cerevisiae*: engineering of pyruvate carboxylation, oxaloacetate reduction, and malate export. *Appl. Environ. Microbiol.*, 74(9):2766–2777, 2008.

[79] M. Izallalen, R. Mahadevan, A. Burgard, B. Postier, Jr R. Didonato, J. Sun, and C. H. Schilling. *Geobacter sulfurredonces* strain engineered for increased rates of respiration. *Metab. Eng.*, 10:267–275, 2008.

[80] M. D. Mesarović. *Systems theory and biology—view of a theoretician.* Springer, 1968.

[81] T. Ideker, T. Galitski, and L. Hood. A new approach to decoding life: systems biology. *Annu. Rev. Genom. Hum. G.*, 2(1):343–372, 2001.

[82] L. M. Gierasch and A. Gershenson. Post-reductionist protein science, or putting Humpty Dumpty back together again. *Nat. Chem. Biol.*, 5(11): 774–777, 2009.

[83] H. V. Westerhoff and B. Ø. Palsson. The evolution of molecular biology into systems biology. *Nat. Biotechnol.*, 22(10):1249–1252, 2004.

[84] A. Kriete and R. Eils. *Computational Systems Biology.* Academic press, 2005.

[85] H. Kitano. International alliances for quantitative modeling in systems biology. *Mol. Syst. Biol.*, 1(1), 2005.

[86] J. Krömer, L. E. Quek, and L. Nielsen. $^{13}$C-Fluxomics: a tool for measuring metabolic phenotypes. *Australian Biochemist*, 40(3):17–20, 2009.

[87] G. Winter and J. O. Krömer. Fluxomics–connecting 'omics analysis and phenotypes. *Environ. Microbiol.*, 15(7):1901–1916, 2013.

[88] A. Kun, B. Papp, and E. Szathmáry. Computational identification of obligatorily autocatalytic replicators embedded in metabolic networks. *Genome Biol.*, 9(3):R51, 2008.

[89] T. Nishikawa, N. Gulbahce, and A. E. Motter. Spontaneous reaction silencing in metabolic optimization. *PLoS Comput. Bio.*, 4(12):e1000236, 2008.

[90] R. Guimerà, M. Sales-Pardo, and L. A. N. Amaral. A network-based method for target selection in metabolic networks. *Bioinformatics*, 23(13): 1616–1622, 2007.

[91] B. Papp, B. Teusink, and R. A. Notebaart. A critical view of metabolic network adaptations. *HFSP J.*, 3(1):24–35, 2009.

[92] R. Guimerà and L. A. N. Amaral. Functional cartography of complex metabolic networks. *Nature*, 433:895–900, 2005.

[93] M. Á. Serrano and M. Boguñá and F. Sagués. Network-based confidence scoring system for genome-scale metabolic reconstructions. *BMC Syst. Biol.*, 5:76, 2011.

[94] O. Güell, F. Sagués, and M. Á. Serrano. Predicting effects of structural stress in a genome-reduced model bacterial metabolism. *Sci. Rep.*, 2:621, 2012.

[95] M. Á. Serrano, M. Boguñá, and A. Vespignani. Extracting the mutiscale backbone of complex weighted networks. *Proc. Natl. Acad. Sci. USA*, 106: 6483–6488, 2009.

[96] R. Albert and A. L. Barabási. Statistical mechanics of complex networks. *Rev. Mod. Phys*, 74:47–97, 2002.

[97] H. Jeong, B. Tombor, R. Albert, Z. N. Oltvai, and A. L. Barabási. The large-scale organization of metabolic networks. *Nature*, 407:651–654, 2000.

[98] H. Jeong, S. P. Mason, A. L. Barabási, and Z. N. Oltvai. Lethality and centrality in protein networks. *Nature*, 411(6833):41–42, 2001.

[99] A. Wagner and D. A. Fell. The small world inside large metabolic networks. *Proc. R. Soc. Lond., B, Biol. Sci.*, 268:1803–1810, 2001.

[100] E. Ravasz, A. L. Somera, D. A. Mongru, Z. N. Oltvai, and A. L. Barabási. Hierarchical organization of modularity in metabolic networks. *Science*, 297:1551–1555, 2002.

[101] A. Wagner. Distributed robustness versus redundancy as causes of mutational robustness. *BioEssays*, 27:176–188, 2005.

[102] A. E. Motter, N. Gulbahce, E. Almaas, and A. L. Barabási. Predicting synthetic rescues in metabolic networks. *Mol. Syst. Biol.*, 4:168, 2008.

[103] A. Varma and B. Ø. Palsson. Metabolic capabilities of *Escherichia coli*: I. Synthesis of biosynthetic precursors and cofactors. *J. Theor. Biol.*, 165 (4):477–502, 1993.

[104] P. F. Suthers, A. Zomorrodi, and C. D. Maranas. Genome-scale gene/reaction essentiality and synthetic lethality analysis. *Mol. Syst. Biol.*, 5:301, 2009.

[105] A. Barve, J. F. M. Rodrigues, and A. Wagner. Superessential reactions in metabolic networks. *Proc. Natl. Acad. Sci. USA*, 1091:E1121–E1130, 2012.

[106] R. Mahadevan and C. H. Schilling. The effects of alternate optimal solutions in constraint-based genome-scale metabolic models. *Metab. Eng.*, 5:264–276, 2003.

[107] S. Gudmundsson and I. Thiele. Computationally efficient flux variability analysis. *BMC Bioinformatics*, 11:489, 2010.

[108] J. L. Guillaume and M. Latapy. Bipartite graphs as models of complex networks. *Physica A*, 371:795–813, 2006.

[109] P. Holme, F. Liljeros, C. R. Edling, and B. J. Kim. Network bipartivity. *Phys. Rev. E*, 68(5):056107, 2003.

[110] H. W. Ma and A. P. Zeng. Reconstruction of metabolic networks from genome data and analysis of their global structure for various organisms. *Bioinformatics*, 19:270–277, 2003.

[111] M. E. J. Newman. The structure and function of complex networks. *SIAM Rev.*, 45:167–256, 2003.

[112] A. L. Barabási and E. Bonabeau. Scale-free networks. *Sci. Am.*, 288(5): 50–59, 2003.

[113] R. Albert. Scale-free networks in cell biology. *J. Cell Sci.*, 118(21): 4947–4957, 2005.

[114] E. F. Keller. Revisiting "scale-free" networks. *Bioessays*, 27(10):1060–1068, 2005.

[115] R. Tanaka. Scale-rich metabolic networks. *Phys. Rev. Lett.*, 94(16):168101, 2005.

[116] P. Kim et al. Metabolite essentiality elucidates robustness of *Escherichia coli* metabolism. *Proc. Natl. Acad. Sci. USA*, 104(34):13638–13642, 2007.

[117] A. M. Feist et al. A genome-scale metabolic reconstruction for *Escherichia coli* K-12 MG1655 that accounts for 1260 ORFs and thermodynamic information. *Mol. Syst. Biol.*, 3:121, 2007.

[118] J. D. Orth, R. M. Fleming, and B. Ø. Palsson. *EcoSal - Escherichia coli and Salmonella: Cellular and Molecular Biology*. ASM Press, Washington D.C., 2009.

[119] J. D. Orth et al. A comprehensive genome-scale reconstruction of *Escherichia coli* metabolism - 2011. *Mol. Syst. Biol.*, 7:535, 2011.

[120] E. Yus et al. Impact of genome reduction on bacterial metabolism and its regulation. *Science*, 326:1263–1268, 2009.

[121] H. W. Ma and A. P. Zeng. Reconstruction of metabolic networks from genome information and its structural. In *Computational Systems Biology*, page 169. Academic Press, 2005.

[122] M. Arita. The metabolic world of *Escherichia coli* is not small. *Proc. Natl. Acad. Sci. USA*, 101(6):1543–1547, 2004.

[123] J. T. Gao, R. Guimerà, H. Li, I. M. Pinto, M. Sales-Pardo, S. C. Wai, B. Rubinstein, and R. Li. Modular coherence of protein dynamics in yeast cell polarity system. *Proc. Natl. Acad. Sci. USA*, 108(18):7647–7652, 2011.

[124] S. Fortunato. Community detection in graphs. *Phys. Rep.*, 486(3):75–174, 2010.

[125] M. Rosvall and C. T. Bergstrom. Maps of random walks on complex networks reveal community structure. *Proc. Natl. Acad. Sci. USA*, 105: 1118–1123, 2008.

[126] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Phys. Rev. E*, 69(2):026113, 2004.

[127] J. Reichardt and S. Bornholdt. Statistical mechanics of community detection. *Phys. Rev, E*, 74(1):016110, 2006.

[128] V. D. Blondel, J. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *J. Stat. Mech.*, 2008(10):P10008, 2008.

[129] T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning.* Springer, 2009.

[130] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin. *Network flows: theory, algorithms, and applications.* Prentice hall, 1993.

[131] H. W. Ma and A. P. Zeng. The connectivity stucture, giant strong component and centrality of metabolic networks. *Bioinformatics*, 19: 1423–1430, 2003.

[132] M. Boguñá and M. Ángeles Serrano. Generalized percolation in random directed networks. *Phys. Rev. E*, 72:016106, 2005.

[133] M. Á. Serrano and P. De Los Rios. Structural efficiency of percolated landscapes in flow networks. *PLoS ONE*, 3:e3654, 2008.

[134] D. J. Watts and S. H. Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 393:440–442, 1998.

[135] P. Erdös and A. Rényi. On random graphs I. *Publ. Math. Debrecen*, 6: 290–297, 1959.

[136] P. Erdös and A. Rényi. On the evolution of random graphs. *Publ. Math. Inst. Hungar. Acad. Sci*, 5:17–61, 1960.

[137] M. Molloy and B. Reed. A critical point for random graphs with a given degree sequence. *Random Struct. Algor.*, 6:161–179, 1995.

[138] M. E. J. Newman, S. H. Strogatz, and D. J. Watts. Random graphs with arbitrary degree distributions and their applications. *Phys. Rev. E*, 64, 2001.

[139] R. Milo et al. Network motifs: simple building blocks of complex networks. *Science*, 298(5594):824–827, 2002.

[140] A. G. Smart, L. A. N. Amaral, and J. Ottino. Cascading failure and robustness in metabolic networks. *Proc. Natl. Acad. Sci. USA*, 105:13223–13228, 2008.

[141] G. Basler, O. Ebenhöh, J. Selbig, and Z. Nikoloski. Mass-balanced randomization of metabolic networks. *Bioinformatics*, 27:1397–1403, 2011.

[142] G. Basler, S. Grimbs, O. Ebenhöh, J. Selbig, and Z. Nikoloski. Evolutionary significance of metabolic network properties. *J. R. Soc. Interface*, 9:1168–1176, 2012.

[143] R. S. Costa, D. Machado, I. Rocha, and E. C. Ferreira. Critical perspective on the consequences of the limited availability of kinetic data in metabolic dynamic modelling. *IET Syst. Biol.*, 5(3):157–163, 2011.

[144] S. Schuster, T. Dandekar, and D. A. Fell. Detection of elementary flux modes in biochemical networks: a promising tool for pathway analysis and metabolic engineering. *Trends Biotechnol.*, 17(2):53–60, 1999.

[145] N. Price, J. Reed, J. Papin, S. Wiback, and B. Ø. Palsson. Network-based analysis of metabolic regulation in the human red blood cell. *J. Theor. Biol.*, 225(2):185–194, 2003.

[146] M. Terzer, N. D. Maynard, M. W. Covert, and J. Stelling. Genome-scale metabolic networks. *Wiley Interdiscip. Rev. Syst. Biol. Med.*, 1(3):285–297, 2009.

[147] C. H. Schilling and B. Ø. Palsson. The underlying pathway structure of biochemical reaction networks. *Proc. Natl. Acad. Sci. USA*, 95:4193–4198, 1998.

[148] C. H. Schilling, J. S. Edwards, D. Letscher, and B. Ø. Palsson. Combining pathway analysis with flux balance analysis for the comprehensive study of metabolic systems. *Biotechnol. Bioeng.*, 71:286–306, 2000.

[149] A. Makhorin. GNU linear programming kit. *Moscow Aviation Institute, Moscow, Russia*, 38, 2001.

[150] R. Ceron. *The GNU Linear Programming Kit, Part 1: Introduction to linear optimization.* IBM, 2006.

[151] R. Ceron. *The GNU Linear Programming Kit, Part 2: Intermediate problems in linear programming.* IBM, 2006.

[152] R. Ceron. *The GNU Linear Programming Kit, Part 3: Advanced problems and elegant solutions.* IBM, 2006.

[153] K. G. Murty. *Linear programming*, volume 57. Wiley New York, 1983.

[154] G. Sezonov, D. Joseleau-Petit, and R. D'Ari. *Escherichia coli* physiology in Luria-Bertani Broth. *J. Bacteriol.*, 189:8746–8749, 2007.

[155] Z. Wunderlich and L. A. Mirny. Using the topology of metabolic networks to predict viability of mutant straints. *Biophys. J.*, 91:2304–2311, 2006.

[156] A. C. Müller and A. Bockmayr. Fast thermodynamically constrained flux variability analysis. *Bioinformatics*, 29:903–909, 2013.

[157] J. L. Reed and B. Ø. Palsson. Genome-scale in silico models of *E. coli* have multiple equivalent phenotypic states: assessment of correlated reaction subsets that comprise network states. *Genome Res.*, 14(9):1797–1805, 2004.

[158] O. Güell, F. Sagués, and M. Á. Serrano. Essential plasticity and redundancy of metabolism unveiled by synthetic lethality analysis. *PLoS Comput. Bio.*, 10(5):e1003637, 2014.

[159] N. C. Duarte, S. S. Becker, N. Jamshidi, I. Thiele, M. L. Mo, T. D. Vo, R. Srivas, and B. Ø. Palsson. Global reconstruction of the human metabolic network based on genomic and bibliomic data. *Proc. Natl. Acad. Sci. USA*, 104(6):1777–1782, 2007.

[160] A. M. Feist, M. J. Herrgård, I. Thiele, J. L. Reed, and B. Ø. Palsson. Reconstruction of biochemical networks in microorganisms. *Nat. Rev. Microbiol.*, 7(2):129–143, 2008.

[161] R. S. Senger and E. T. Papoutsakis. Genome-scale model for *Clostridium acetobutylicum*: Part I. Metabolic network resolution and analysis. *Biotechnol. Bioeng.*, 101(5):1036–1052, 2008.

[162] A. Raghunathan, J. Reed, S. Shin, B. Ø. Palsson, and S. Daefler. Constraint-based analysis of metabolic capacity of salmonella typhimurium during host-pathogen interaction. *BMC Syst. Biol.*, 3(1):38, 2009.

[163] I. Thiele et al. A community-driven global reconstruction of human metabolism. *Nature Biotechnol.*, 31(5):419–425, 2013.

[164] I. Schomburg et al. BRENDA in 2013: integrated reactions, kinetic data, enzyme function data, improved disease classification: new options and contents in BRENDA. *Nucl. Acids Res.*, 41:D764–D772, 2012.

[165] J. S. Edwards, R. U. Ibarra, and B. Ø. Palsson. The *Escherichia coli* MG1655 *in silico* metabolic genotype: its definition, characteristics, and capabilities. *Proc. Natl. Acad. Sci. USA*, 97:5528–5533, 2000.

[166] J. L. Reed, T. D. Vo, C. H. Schilling, and B. Ø. Palsson. An expanded genome-scale model of *Escherichia coli* K-12 (*i*JR904 GSM/GPR). *Genome Biol.*, 4(9):R54, 2003.

[167] M. Riley et al. *Escherichia coli* K-12: a cooperatively developed annotation snapshot—2005. *Nucl. Acids Res.*, 34(1):1–9, 2006.

[168] I. M. Keseler et al. EcoCyc: a comprehensive database resource for *Escherichia coli*. *Nucl. Acids Res.*, 33(suppl 1):D334–D337, 2005.

[169] I. M. Keseler et al. EcoCyc: a comprehensive view of *Escherichia* coli biology. *Nucl. Acids Res.*, 37(suppl 1):D464–D470, 2009.

[170] M. Kanehisa, S. Goto, M. Furumichi, M. Tanabe, and M. Hirakawa. KEGG for representation and analysis of molecular networks involving diseases and drugs. *Nucl. Acids Res.*, 38(suppl 1):D355–D360, 2010.

[171] S. Kühner et al. Proteome organization in a genome-reduced bacterium. *Science*, 326:1235–1240, 2009.

[172] M. Güell et al. Transcriptome complexity in a genome-reduced bacterium. *Science*, 326:1268–1271, 2009.

[173] M. Kuroda et al. Whole genome sequencing of meticillin-resistant *Staphylococcus aureus*. *Lancet*, 357(9264):1225–1240, 2001.

[174] J. D. Peterson, L. A. Umayam, T. Dickinson, E. K. Hickey, and O. White. The comprehensive microbial resource. *Nucl. Acids Res.*, 29(1):123–125, 2001.

[175] S. N. Dorogovtsev, A. V. Goltsev, and J. F. F. Mendes. Critical phenomena in complex networks. *Rev. Mod. Phys.*, 80:1275–1335, 2008.

[176] R. Cohen, K. Erez, D. ben Avraham, and S. Havlin. Resilience of the internet to random breakdown. *Phys. Rev. Lett.*, 85:4626, 2000.

[177] R. Albert, H. Jeong, and A. L. Barabási. Error and attack tolerance of complex networks. *Nature*, 406:378–382, 2000.

[178] D. J. Watts. A simple model of global cascades on random networks. *Proc. Natl. Acad. Sci. U.S.A.*, 99:5766–5771, 2002.

[179] Y. Moreno, J. B. Gómez, and A. F. Pacheco. Instability of scale-free networks under nodebreaking avalanches. *Europhys. Lett.*, 58:630–636, 2002.

[180] A. E. Motter. Cascade-based attacks on complex networks. *Phys. Rev. E*, 66:065102(R), 2002.

[181] S. V. Buldyrev, R. Parshani, G. Paul, H. E. Stanley, and S. Havlin. Catastrophic cascade of failures in interdependent networks. *Nature*, 464: 1025–1028, 2010.

[182] M. S. Szalay, I. A. Kovacs, T. Korcsmaros, C. Bode, and P. Csermely. Stress-induced rearrangements of cellular networks: consequences for protection and drug design. *FEBS Lett.*, 581:3675–3680, 2007.

[183] J. S. Edwards and B. Ø. Palsson. Robustness analysis of the *Escherichia coli* metabolic network. *Biotechnol. Prog.*, 16:927–939, 2000.

[184] D. Segrè, D., and G. M. Church. Analysis of optimality in natural and perturbed metabolic networks. *Proc. Natl. Acad. Sci. USA*, 99:15112–15117, 2002.

[185] O. Folger, L. Jerby, C. Frezza, E. Gottlieb, E. Ruppin, and T. Shlomi. Predicting selective drug targets in cancer through metabolic networks. *Mol. Syst. Biol.*, 7:501, 2011.

[186] O. Güell, F. Sagués, G. Basler, Z. Nikoloski, and M. Á. Serrano. Assessing the significance of knockout cascades in metabolic networks. *J. Comp. Int. Sci.*, 3(1-2):45–53, 2012.

[187] N. V. Smirnov. Tables for estimating the goodness of fit of empirical distributions. *Ann. Math. Stat.*, 19:279, 1948.

[188] C. Spearman. The proof and measurement of association between two things. *Amer. J. Psychol.*, 15:72–101, 1904.

[189] J. I. Glass et al. Essential genes of a minimal bacterium. *Proc. Natl. Acad. Sci. USA*, 103:425–430, 2006.

[190] B. R. Borate et al. Comparison of threshold selection matrices for microarray gene co-expression matrices. *BMC Res. Notes*, 2:240, 2009.

[191] B. Zhang and S. Horvath. A general framework for weighted gene co-expression network analysis. *Stat. Appl. Genet. Mol. Biol.*, 4:1544–6115, 2005.

[192] R. Khanin and E. Wit. Construction of malaria gene expression network using partial correlations. *Methods of Microarray Data Analysis V*, pages 1544–6115, 2005.

[193] J. L. DeRisi, V. Iyer, and P. O. Brown. Exploring the metabolic and genetic control of gene expression on a genomic scale. *Science*, 278(5338): 680–686, 1997.

[194] E. Almaas, Z. N. Oltvai, , and A. L. Barabási. The activity reaction core and plasticity of metabolic networks. *PLoS Comput. Biol.*, 1:0557–0563, 2005.

[195] T. Baba et al. Construction of *Escherichia coli* K-12 in-frame, single-gene knockout mutants: the Keio collection. *Mol. Syst. Biol*, 2(1), 2006.

[196] A. R. Joyce et al. Experimental and computational assessment of conditionally essential genes in *Escherichia coli. J. Bacteriol.*, 188(23):8259–8271, 2006.

[197] Z. Wang and J. Zhang. Abundant indispensable redundancies in cellular metabolic networks. *Genome Biol. Evol.*, 1:23–33, 2009.

[198] P. Nygaard and J.M. Smith. Evidence for a novel glycinamide ribonucleotide transformylase in *Escherichia coli. J. Bacteriol.*, 175:3591–3597, 1993.

[199] J. L. Hartman, B. Garvik, and L. Hartwell. Principles for the buffering of genetic variation. *Science*, 291:1001–1004, 2001.

[200] J. Masel and M. L. Siegal. Robustness: mechanisms and consequences. *Trends Genet.*, 25:395–403, 2009.

[201] S. M. B. Nijman. Synthetic lethality: general principles, utility and detection using genetic screens in human cells. *FEBS Lett.*, 585:1–6, 2011.

[202] W. G. Kaelin. The concept of synthetic lethality in the context of anticancer therapy. *Nat. Rev. Cancer*, 5:689–698, 2005.

[203] R. Harrison, B. Papp, C. Pál, S. G. Oliver, and D. Delneri. Plasticity of genetic interactions in metabolic networks of yeast. *Proc. Natl. Acad, Sci. USA*, 104:2307–2312, 2007.

[204] R. Kelley and T. Ideker. Systematic interpretation of genetic interactions using protein networks. *Nat. Biotechnol.*, 23:561–566, 2005.

[205] O. Güell, M. Á. Serrano, and F. Sagués. Environmental dependence of the activity and essentiality of reactions in the metabolism of *Escherichia coli*. In *Engineering of Chemical Complexity II*, pages 39–56. World Scientific, 2014. ISBN 978-981-4616-14-0.

[206] P. Novick, B. C. Osmond, and D. Botstein. Suppressors of yeast actin mutations. *Genetics*, 121:659–674, 1989.

[207] G. Zhao and M. E. Winkler. An *Escherichia coli* K-12 *tktA tktB* mutant deficient in transketolase activity requires pyridoxine (vitamin B$_6$) as well as the aromatic amino acids and vitamins for growth. *J. Bacteriol.*, 176: 883–891, 1995.

[208] D. Deutscher, I. Meilijson, M. Kupiec, and E. Ruppin. Multiple knockout analysis of genetic robustness in the yeast metabolic network. *Nat. Genet.*, 38:993–998, 2006.

[209] W. A. Whalen and C. M. Berg. Analysis of *avtA*::Mu *d*1(ap *lac*) mutant: metabolic role of transaminase C. *J. Bacteriol.*, 150:739–746, 1982.

[210] P. Nygaard and J.M. Smith. Evidence for a novel glycinamide ribonucleotide transformylase in *Escherichia coli*. *J. Bacteriol.*, 175:3591–3597, 1993.

[211] B. Troup, C. Hungerer, and D. Jahn. Cloning and characterization of the *Escherichia coli hemN* gene encoding the oxygen-independent coproporphyrinogen III oxidase. *J. Bacteriol.*, 177:3326–3331, 1995.

[212] A. Rompf et al. Regulation of *Pseudomonas aeruginosa hemF* and *hemN* by the dual action of the redox response regulators Anr and Dnr. *Mol. Microbiol.*, 29:985–997, 1998.

[213] Z. Jiao, T. Baba, H. Mori, and K. Shimizu. Analysis of metabolic and physiological responses to *gnd* knockout in *Escherichia coli* by using C-13 tracer experiment and enzyme activity measurement. *FEMS Microbiol. Lett.*, 220:295–301, 2003.

[214] M. A. Rude and A. Schirmer. New microbial fuels: a biotech perspective. *Curr. Opin. Microbiol.*, 12:274–281, 2009.

[215] G. Giaever et al. Functional profiling of the *Saccharomyces cerevisiae* genome. *Nature*, 418:387–391, 2002.

[216] L. M. Steinmetz et al. Systematic screen for human disease genes in yeast. *Nat. Genet.*, 31:400–404, 2002.

[217] K. Nakahigashi et al. Systematic phenome analysis of *Escherichia coli* multiple-knockout mutants reveals hidden reactions in central carbon metabolism. *Mol. Syst. Biol.*, 5:306, 2009.

[218] O. C. Herfindahl. *Copper Costs and Prices: 1870-1957.* John Hopkins University Press, Baltimore, MD, USA, 1959.

[219] A. O. Hirschman. The paternity of an index. *Am. Econ. Rev.*, 54:761–762, 1964.

[220] A. Barrat, M. Barthélemy, R. Pastor-Satorras, and A. Vespignani. The architecture of complex weighted networks. *Proc. Natl. Acad. Sci. USA*, 101(11):3747–3752, 2004.

[221] O. Güell, F. Sagués, and M. Á. Serrano. Detecting evolution and adaptation fingerprints in bacterial metabolic backbones. *arXiv:1412.3353*, 2014.

[222] G. Bianconi. Flux distribution of metabolic networks close to optimal biomass production. *Phys. Rev. E*, 78(3):035101, 2008.

[223] M. A. Keller, A. V. Turchyn, and M. Ralser. Non-enzymatic glycolysis and pentose phosphate pathway-like reactions in a plausible Archean ocean. *Mol. Syst. Biol.*, 10(4), 2014.

[224] E. Meléndez-Hevia, T. G. Waddell, and M. Cascante. The puzzle of the Krebs citric acid cycle: assembling the pieces of chemically feasible reactions, and opportunism in the design of metabolic pathways during evolution. *J. Mol. Evol.*, 43(3):293–303, 1996.

[225] O. Ebenhöh and R. Heinrich. Evolutionary optimization of metabolic pathways. Theoretical reconstruction of the stoichiometry of ATP and NADH producing systems. *B. Math. Biol.*, 63(1):21–55, 2001.

[226] R. J. Mailloux, R. Bériaut, J. Lemire, R. Singh, D. R. Chénier, R. D. Hamel, and V. D. Appanna. The tricarboxylic acid cycle, an ancient metabolic network with a novel twist. *PLoS ONE*, 2(8):e690, 2007.

[227] D. A. Fell and A. Wagner. The small world of metabolism. *Nat. Biotechnol.*, 18:1221–1122, 2000.

[228] D. R. Evans and H. I. Guy. Mammalian pyrimidine biosynthesis: fresh insights into an ancient pathway. *J. Biol. Chem.*, 279(32):33035–33038, 2004.

[229] M. W. Powner, B. Gerland, and J. D. Sutherland. Synthesis of activated pyrimidine ribonucleotides in prebiotically plausible conditions. *Nature*, 459(7244):239–242, 2009.

[230] G. Caetano-Anolles, H. S. Kim, and J. E. Mittenthal. The origin of modern metabolic networks inferred from phylogenomic analysis of protein architecture. *Proc. Natl. Acad. Sci. USA*, 104(22):9358–9363, 2007.

[231] S. Suen, H. H. S. Lu, and C. H. Yeang. Evolution of domain architectures and catalytic functions of enzymes in metabolic systems. *Genome Biol. Evol.*, 4(9):976–993, 2012.

[232] H. Tao, C. Bausch, C. Richmond, F. R. Blattner, and T. Conway. Functional genomics: expression analysis of *Escherichia coli* growing on minimal and rich media. *J. Bacteriol.*, 181(20):6425–6440, 1999.

[233] A. Lourenço, S. Carneiro, J. P. Pinto, M. Rocha, E. C. Ferreira, and I. Rocha. A study of the short and long-term regulation of *E. coli* metabolic pathways. *J. Integr. Bioinform.*, 8(3):183, 2011.

[234] J. M. Monk et al. Genome-scale metabolic reconstructions of multiple *Escherichia coli* strains highlight strain-specific adaptations to nutritional environments. *Proc. Natl. Acad. Sci. USA*, 110(50):20338–20343, 2013.

[235] R. Bourqui et al. Metabolic network visualization eliminating node redundancy and preserving metabolic pathways. *BMC Syst. Biol.*, 1:29, 2009.

[236] K. Faust, P. Dupont, J. Callut, and J. van Helden. Pathway discovery in metabolic networks by subgraph extraction. *Bioinformatics*, 26:1211–1218, 2010.

[237] L. M. Blank, L. Kuepfer, and U. Sauer. Large-scale $^{13}$C-flux analysis reveals mechanistic principles of metabolic network robustness to null mutations in yeast. *Genome Biol.*, 6(6):R49, 2005.

[238] F. Llaneras and J. Picó. Stoichiometric modelling of cell metabolism. *J. Biosci. Bioeng.*, 105(1):1–11, 2008.

[239] J. T. Manolukas, M. F. Barile, D. K. Chandler, and J. D. Pollack. Presence of anaplerotic reactions and transamination, and the absence of the tricarboxylic acid cycle in mollicutes. *J. Gen. Microbiol.*, 134(3):791–800, 1988.

[240] O. Frick and C. Whittmann. Characterization of the metabolic shift between oxidative and fermentative growth in Saccharomyces cerevisiae by comparative 13C flux analysis. 4:30, 2005.

[241] M. G. Vander Heiden, L. C. Cantley, and C. B. Thompson. Understanding the warburg effect: The metabolic requirements of cell proliferation. *Science*, 324:1029–1033, 2009.

[242] J. Menendez, J. Joven, S. Cufí, B. Corominas-Faja, C. Oliveras-Ferraros, E. Cuyàs, B. Martin-Castillo, E. Lopez-Bonet, T. Alarcón, and A. Vazquez-Martin. The Warburg effect version 2.0. *Cell Cycle*, 12(8):1166–1179, 2013.

[243] R. Schuetz, L. Kuepfer, and U. Sauer. Systematic evaluation of objective functions for predicting intracellular fluxes in *Escherichia coli. Mol. Syst. Biol.*, 3:119, 2007.

[244] R. Schuetz, N. Zamboni, M. Zampieri, M. Heinemann, and U. Sauer. Multidimensional optimality of microbial metabolism. *Science*, 336:601–604, 2012.

[245] O. Güell, F. A. Massucci, F. Font-Clos, F. Sagués, and M. Á. Serrano. Mapping high-growth phenotypes in the flux space of microbial metabolism. *arXiv:1409.4595*, 2014.

[246] L. Lovász. Hit-and-Run mixes fast. *Math. Program.*, 86(3):443–461, 1999.

[247] K. Pearson. LIII. On lines and planes of closest fit to systems of points in space. *Philos. Mag. Series 6*, 2(11):559–572, 1901.

[248] I. T. Jolliffe. *Principal Component Analysis*. Springer-Verlag, second edition, New York, 2002.

[249] J. S. Edwards and B. Ø. Palsson. *In silico* predictions of *Escherichia coli* metabolic capabilities are consistent with experimental data. *Nat. Biotechnol.*, 19:125–130, 2001.

[250] A. Varma, B. W. Boesch, and B. Ø. Palsson. Stoichiometric interpretation of *Escherichia coli* glucose catabolism under various oxygenation rates. *Appl. Environ. Microb.*, 59:2465–2473, 1993.

[251] H. E. Reiling, H. Laurila, and A. Fiechter. Mass culture of *Escherichia coli*: medium development for low and high density cultivation of *Escherichia coli* B/r in minimal and complex media. *J. Biotechnol.*, 2:191–206, 1985.

[252] E. M. El-Mansi and W. H. Holms. Control of carbon flux to acetate excretion during growth of *Escherichia coli* in batch and continuous cultures. *J. Gen. Microbiol.*, 135:2875–2883, 1989.

[253] A. J. Wolfe. The acetate switch. *Microbiol. Mol. Biol. R.*, 69:12–50, 2005.

[254] A. Vázquez, Q. K. Beg, M. A. deMenezes, J. Ernst, Z. Bar-Joseph, A. L. Barabási, L. G. Boros, and Z. N. Oltvai. Impact of the solvent capacity constraint on *E. coli* metabolism. *BMC Syst. Biol.*, 2:7, 2008.

[255] A. Varma and B. Ø. Palsson. Stoichiometric flux balance models quantitatively predict growth and metabolic by-product secretion in wild-type *Escherichia coli* W3110. *Appl. Environ. Microb.*, 60(10):3724–3731, 1994.

[256] D. Molenaar, R. van Berlo, D. de Ridder, and B. Teusink. Shifts in growth strategies reflect tradeoffs in cellular economics. *Mol. Syst. Biol.*, 5(1), 2009.

[257] M. T. Wortel, H. Peters, J. Hulshof, B. Teusink, and F. J. Bruggeman. Metabolic states with maximal specific rate carry flux through an elementary flux mode. *FEBS J.*, 281(6):1547–1555, 2014.

[258] M. Losen, B. Frolich, M. Pohl, and J. Buchs. Effect of oxygen limitation and medium composition on *Escherichia coli* fermentation in shake-flask cultures. *Biotechnol. Progr.*, 20:1062–1068, 2004.

[259] M. Orencio-Trejo, J. Utrilla, M.T. Fernández-Sandoval, G. Huerta-Beristain, G. Gosset, and A. Martinez. Engineering the *Escherichia coli* fermentative metabolism. *Adv. Biochem. Eng. Biot.*, 121:71–107, 2010.

[260] N. D. Price and I. Shmulevich. Biochemical and statistical network models for systems biology. *Curr. Opin. Biotech.*, 18(4):365–370, 2007.

[261] M. Larocque, T. Chénard, and R. Najmanovich. A curated *C. difficile* strain 630 metabolic network: prediction of essential targets and inhibitors. *BMC Syst. Biol.*, 8(1):117, 2014.

[262] J. Aitchison and J. A. C. Brown. The lognormal distribution with special reference to its uses in economics. 1957.

[263] N. D. Price, J. L. Reed, and B. Ø. Palsson. Genome-scale models of microbial cells: evaluating the consequences of constraints. *Nat. Rev. Microbiol.*, 2:886–897, 2004.

[264] S. J. Wiback, I. Famili, H. J. Greenberg, and B. Ø. Palsson. Monte carlo sampling can be used to determine the size and shape of the steady-state flux space. *J. Theor. Biol.*, 228:437–447, 2004.

[265] F. A. Massucci, F. Font-Clos, A. De Martino, and I. Pérez Castillo. A novel methodology to estimate metabolic flux distributions in constraint-based models. *Metabolites*, 3:838–852, 2013.

[266] W. Krauth and M. Mezard. Learning algorithms with optimal stability in neural networks. *J. Phys. A*, 20(11):L745–L752, 1987.

# List of Publications

O. Güell, F. Sagués, and M. Á. Serrano. Predicting effects of structural stress in a genome-reduced model bacterial metabolism. *Sci. Rep.*, 2:621, 2012.

O. Güell, F. Sagués, G. Basler, Z. Nikoloski, and M. Á. Serrano. Assessing the significance of knockout cascades in metabolic networks. *J. Comp. Int. Sci.*, 3 (1-2):45–53, 2012.

O. Güell, F. Sagués, and M. Á. Serrano. Essential plasticity and redundancy of metabolism unveiled by synthetic lethality analysis. *PLoS Comput. Bio.*, 10(5): e1003637, 2014.

O. Güell, M. Á. Serrano, and F. Sagués. Environmental dependence of the activity and essentiality of reactions in the metabolism of *Escherichia coli*. In *Engineering of Chemical Complexity II*, pages 39-56. World Scientific, 2014. ISBN 978-981-4616-14-0.

O. Güell, F. A. Massucci, F. Font-Clos, F. Sagués, and M. Á. Serrano. Mapping high-growth phenotypes in the flux space of microbial metabolism. *arXiv*:1409. 4595 [q-bio.MN].

O. Güell, F. Sagués, and M. Á. Serrano. Detecting evolution and adaptation fingerprints in bacterial metabolic backbones. *arXiv*:1142.3353 [q-bio.MN].