



## GLOBAL OPTIMIZATION APPLIED TO KINETIC MODELS OF METABOLIC NETWORKS

**Carlos Pozo Fernández**

Dipòsit Legal: T.1469-2012

**ADVERTIMENT.** L'accés als continguts d'aquesta tesi doctoral i la seva utilització ha de respectar els drets de la persona autora. Pot ser utilitzada per a consulta o estudi personal, així com en activitats o materials d'investigació i docència en els termes establerts a l'art. 32 del Text Refós de la Llei de Propietat Intel·lectual (RDL 1/1996). Per altres utilitzacions es requereix l'autorització prèvia i expressa de la persona autora. En qualsevol cas, en la utilització dels seus continguts caldrà indicar de forma clara el nom i cognoms de la persona autora i el títol de la tesi doctoral. No s'autoritza la seva reproducció o altres formes d'explotació efectuades amb finalitats de lucre ni la seva comunicació pública des d'un lloc aliè al servei TDX. Tampoc s'autoritza la presentació del seu contingut en una finestra o marc aliè a TDX (framing). Aquesta reserva de drets afecta tant als continguts de la tesi com als seus resums i índexs.

**ADVERTENCIA.** El acceso a los contenidos de esta tesis doctoral y su utilización debe respetar los derechos de la persona autora. Puede ser utilizada para consulta o estudio personal, así como en actividades o materiales de investigación y docencia en los términos establecidos en el art. 32 del Texto Refundido de la Ley de Propiedad Intelectual (RDL 1/1996). Para otros usos se requiere la autorización previa y expresa de la persona autora. En cualquier caso, en la utilización de sus contenidos se deberá indicar de forma clara el nombre y apellidos de la persona autora y el título de la tesis doctoral. No se autoriza su reproducción u otras formas de explotación efectuadas con fines lucrativos ni su comunicación pública desde un sitio ajeno al servicio TDR. Tampoco se autoriza la presentación de su contenido en una ventana o marco ajeno a TDR (framing). Esta reserva de derechos afecta tanto al contenido de la tesis como a sus resúmenes e índices.

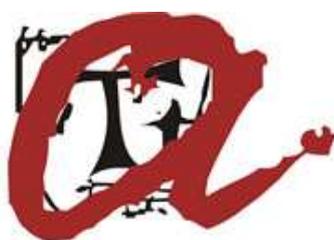
**WARNING.** Access to the contents of this doctoral thesis and its use must respect the rights of the author. It can be used for reference or private study, as well as research and learning activities or materials in the terms established by the 32nd article of the Spanish Consolidated Copyright Act (RDL 1/1996). Express and previous authorization of the author is required for any other uses. In any case, when using its content, full name of the author and title of the thesis must be clearly indicated. Reproduction or other forms of for profit use or public communication from outside TDX service is not allowed. Presentation of its content in a window or frame external to TDX (framing) is not authorized either. These rights affect both the content of the thesis and its abstracts and indexes.

# DOCTORAL THESIS

Carlos Pozo Fernández

# **Global optimization applied to kinetic models of metabolic networks**

Department of Chemical Engineering



**Rovira i Virgili University**



Carlos Pozo Fernández

# Global optimization applied to kinetic models of metabolic networks

DOCTORAL THESIS

Supervised by: Dr. Gonzalo Guillén-Gosálbez  
Dr. Albert Sorribas Tello

Department of Chemical Engineering  
SUSCAPE



**Rovira i Virgili University**

Tarragona  
2012



**Rovira i Virgili  
University**

**Department of Chemical Engineering**

Av. Països Catalans, 26

43007 Tarragona

Phone +34 977 55 86 18

I, Gonzalo Guillén Gosálbez, associate professor in the Department of Chemical Engineering of the Rovira i Virgili University,

CERTIFY:

That the present study, entitled “Global optimization applied to kinetic models of metabolic networks”, presented by Carlos Pozo Fernández for the award of the degree of Doctor, has been carried out under my supervision at the Department of Chemical Engineering of this university.

Tarragona, 27th July 2012,

Dr. Gonzalo Guillén Gosálbez

## **Agradecimientos**

Esta tesis no habría sido posible sin la valiosísima ayuda que me han brindado mis supervisores Dr. Gonzalo Guillén y Dr. Albert Sorbías, y el director del grupo SUSCAPE Dr. Laureano Jiménez. Gracias por guiar mi formación como investigador y como persona. Gracias por estimular el desarrollo de mi creatividad y mi espíritu crítico. Y sobre todo, gracias por vuestra dedicación y apoyo. Estaré eternamente en deuda con vosotros.

No menos agradecido estoy hacia mis compañeros de SUSCAPE. Gracias por compartir conmigo vuestro conocimiento sobre la vida. Gracias por haberme ayudado a desarrollarme en el plano profesional y en el personal. Habéis conseguido que ir a la facultad cada día fuera más un placer que una responsabilidad.

A mis padres les agradezco su educación y su cariño. Gracias por haberme ayudado a convertirme en quién quería ser. Gracias por vuestro amor. Espero que algún día podáis estar tan orgullosos de mí como yo lo estoy de vosotros. Os quiero mucho.

Gracias también a mis hermanos por haber estado a mi lado siempre. Gracias por compartir mis momentos dulces y los amargos. Gracias David por todas las partidas que hemos jugado. Sin esas risas nunca habría podido terminar esta tesis. Gracias MariEli y David por todas las cenas y las comidas. Gracias por haberme ayudado a relajarme cuando lo necesitaba. Gracias también al resto de mi familia por haberme proporcionado un entorno amable en el que todo ha sido más fácil.

Le agradezco a mi amada pareja, Damaris, su paciencia en los momentos de más tensión de estos cuatro años. Gracias por tu apoyo y tus ideas. Gracias por transmitirme tu energía y tu determinación. Gracias también a vosotros Silas, M<sup>a</sup>José y Priscila por haberme acogido en vuestra familia.

Me gustaría también mostrar mi agradecimiento a Núria Juanpere, Núria Golobardes i Laura Cortés por haberme ayudado con todos los trámites administrativos. De sobras sabéis lo torpe que soy con esas cosas. Sin vosotras algunos desastres hubieran sido inevitables.

Por último me gustaría disculparme antes todos aquellos que también me ayudaron y a quiénes esta introducción no ha hecho justicia: gracias a vosotros también.

Tarragona, 27 de Julio de 2012

Carlos Pozo

## **Summary**

In recent years, the use of genetic manipulation techniques has opened the door for obtaining microorganisms with enhanced phenotypes, which has in turn led to significant improvements in the synthesis of certain biochemical products. However, in most cases mutation and selection of these new processes has been performed in a trial-and-error basis. Hence, it is expected that these processes could be further improved if quantitative design principles were used to guide the search towards the ideal enzymatic profiles.

Mathematical programming and particularly optimization have the potential to provide such guidelines yet standard optimization techniques fail at solving the problems arising when kinetic models are used. The reason for this is that the associated problems may contain multiple local optima in which standard optimizers may get trapped during the search. Hence, one must resort to global optimization techniques in order to identify the so-called global optimum of the problem. Even though some algorithms have already been developed for this, they are general purpose and usually fail at solving realistic problems. This thesis is devoted to overcoming such limitations by developing a set of advanced global optimization tools to assess metabolic engineering problems and other questions arising in systems biology.

Specifically, an outer approximation-based algorithm was developed in [1] with the aim of addressing metabolic engineering problems. The framework proposed relies on representing the metabolic network via the Generalized Mass Action (GMA) model and then performing a tailored global optimization of the system in order to obtain the enzymatic profile leading to an optimal product yield. Two case studies consisting of the fermentation of *Saccharomyces cerevisiae* and the citric acid production in *Aspergillus niger* were addressed.

Later, in [2], a similar strategy was devised for the same purpose, but this time the outer approximation was replaced by a customized spatial branch-and-bound. This method exploits the mathematical structure of GMA models to expedite calculations in the more complex cases. This is illustrated by means of comparison between this strategy and both, our previously developed outer approximation and the state-of-art global optimization solver (BARON) for a case study involving the citric acid production in *Aspergillus niger*.

Metabolic engineering studies are not the only ones that benefit from optimization techniques. For instance, these methods can also help understanding the evolution of the strategies that organisms employ to adapt to different environmental situations. Mathematically, the conditions that ensure survival of the microorganism can be modeled as a set of constraints limiting different variables (i.e., enzyme activities) of its metabolic network. Hence, predicting an adaptive response involves determining the set of feasible changes in the enzyme activities that fulfill this set of constraints, that is, characterizing the feasible space of the mathematical problem associated. In [3] we developed a method based on our outer approximation algorithm that enables a systematic characterization of the feasible space of a problem, and thus, of the physiological requirements that may underlie the evolution of adaptive strategies. Specifically, the methodology was used to explain the adaptation of yeast to a heat shock by means of comparing the model predictions with the experimental observations. Additionally, different objective

functions were considered and studied as potential drivers of the actual adaptive process of the microorganism.

Methods presented so far can be applied when the metabolic network is described via the GMA representation. In order to extend this framework to other kinetic models, we resort to a symbolic reformulation strategy under the name of recasting. This technique permits the exact transformation of a model of arbitrary form into a canonical GMA model, at the expense of increasing the number of variables of the original model. Nonetheless, once a model has been transformed into its GMA equivalent, we can effectively apply the optimization and feasibility analysis originally devised for GMA models. In [4], we showed how to perform the recasting for a particular formalism, the Saturable and Cooperative (SC), which is an even more accurate representation than GMA. The usefulness of the approach was illustrated by means of solving problems embedding SC models that standard global optimizers failed to solve efficiently.

Another strategy to further extend these frameworks consists of incorporating more objectives into the associated problems formulations. For instance, biotechnology studies typically seek optimizing a single flux in the metabolic network as unique criterion. In practice, however, there are other criteria of interest for experimentalists, such as minimizing the number of enzymatic changes, metabolic concentration of intermediates or transient times. The incorporation of these functional criteria as constraints ensuring cell viability does not allow for the identification of solutions in which cell viability is further improved at the expense of marginal reductions in other objectives, which is something that can only be achieved via multi-objective optimization (MOO). For this, we developed in [5] a global optimization framework capable of efficiently dealing with several biological criteria simultaneously. The proposed strategy makes use of a heuristic approach based on the epsilon constraint method that reduces the computational burden of generating a set of Pareto optimal alternatives. Furthermore, with the aim of facilitating the post-optimal analysis of these solutions and narrow down their number prior to being tested in the laboratory, we implemented two Pareto filters that systematically identify the preferred subset of enzymatic profiles. The usefulness of our approach was demonstrated by means of a case study that optimizes the ethanol production in the fermentation of *Saccharomyces cerevisiae* considering 14 different objectives.

One of the main advantages of optimization techniques is that they are versatile tools that can be applied to problems arising in different areas. In this thesis [6], we also discuss an application of mathematical programming to the life cycle assessment (LCA) methodology. For this, we considered a supply chain management (SCM) problem in which the goal is to determine the set of Pareto optimal supply chain configurations that maximize the net present value (NPV) and minimize a set of environmental metrics. The multi-objective framework we proposed incorporates a heuristic approach based on Principal Component Analysis (PCA) that allows for uncovering hidden relationships between different LCA metrics. This knowledge is then used to obtain, for a similar computational burden, a better representation of the Pareto set through a reduction of the problem dimensionality.

## **Table of contents**

<b>1. Introduction</b>	<b>1</b>
1.1. Objectives of the thesis	1
<b>2. Systems biology</b>	<b>1</b>
2.1. Metabolic engineering (papers [1,2], also papers [3,4,5])	2
2.2. Evolution studies (paper [3])	2
<b>3. Mathematical programming: optimization</b>	<b>3</b>
3.1. Global optimization	4
3.2. Relaxations in global optimization	4
3.3. Degeneracy	5
3.4. Multi-objective optimization	5
3.4.1. Epsilon constraint	6
3.4.2. Challenges in MOO	6
<b>4. Application of optimization to systems biology problems</b>	<b>7</b>
4.1. Metabolic engineering with GMA models	8
4.2. Metabolic engineering with other kinetic models	9
4.3. Multi-objective global optimization in metabolic engineering	10
4.4. Feasibility analysis in evolution studies	11
4.4.1. Extension of feasibility analysis for metabolic engineering problems	12
<b>5. Other applications of mathematical programming</b>	<b>12</b>
<b>6. Conclusions</b>	<b>12</b>
<b>7. Future work</b>	<b>13</b>
<b>8. Nomenclature</b>	<b>14</b>
8.1. Abbreviations	14
8.2. Indices	15
8.3. Sets/Subsets	15
8.4. Parameters	15
8.5. Variables	15
<b>9. References</b>	<b>16</b>
<b>10. Paper 1: Outer approximation-based algorithm for biotechnology studies in systems biology</b>	<b>19</b>
10.1. Introduction	19

10.2.	Problem statement.....	20
10.3.	Mathematical formulation.....	20
10.3.1.	GMA representation.....	20
10.3.2.	NLP formulation.....	20
10.4.	Solution strategy.....	21
10.4.1.	Upper level master problem.....	21
10.4.2.	Lower level slave problem.....	22
10.4.3.	Algorithm steps.....	22
10.4.4.	Remarks.....	22
10.5.	Case studies.....	23
10.5.1.	Ethanol production in <i>S. cerevisiae</i> .....	25
10.5.2.	Citric acid production in <i>A. niger</i> .....	26
10.6.	Conclusions.....	27
10.7.	Acknowledgements.....	27
10.8.	Appendic A.....	27
10.9.	References.....	30
<b>11.</b>	<b>Paper 2: A spatial branch-and-bound framework for the global optimization of kinetic models of metabolic networks.....</b>	<b>31</b>
11.1.	Introduction.....	31
11.2.	Problem statement.....	33
11.3.	Mathematical formulation.....	33
11.3.1.	GMA representation.....	33
11.3.2.	MINLP formulation.....	33
11.4.	Solution strategy.....	34
11.4.1.	Relaxed subproblem.....	34
11.4.2.	Customized spatial branch-and-bound.....	35
11.5.	Computational results.....	39
11.6.	Conclusions.....	41
11.7.	Author information.....	42
11.8.	Acknowledgement.....	42
11.9.	Nomenclature.....	42
11.10.	References.....	42

<b>12. Paper 3: Optimization and evolution in metabolic pathways: Global optimization techniques in Generalized Mass Action models.....</b>	<b>45</b>
12.1. Introduction.....	45
12.2. Methods.....	46
12.2.1. Generalized Mass Action models .....	46
12.2.2. Characterization of the effect of changes in enzyme activities.....	46
12.2.3. Criteria for functional effectiveness in cellular metabolism.....	47
12.3. Feasibility regions in biochemical pathways: definition and their practical significance	48
12.3.1. Definitions.....	48
12.3.2. Characterization of feasibility regions in GMA models .....	48
12.3.3. Utility of feasibility regions characterization .....	49
12.3.4. Example .....	50
12.4. Discussion .....	55
12.5. Acknowledgements.....	56
12.6. References.....	56
<b>13. Paper 4: Steady-state global optimization of metabolic non-linear dynamic models through recasting into power-law canonical models. ....</b>	<b>58</b>
13.1. Background .....	58
13.2. Results.....	60
13.2.1. Global optimization of non-linear models through recasting .....	60
13.2.2. Optimization goals .....	61
13.2.3. Global optimization of SC models using BARON .....	61
13.2.4. Recasting SC models into GMA models .....	61
13.2.5. Steady-state optimization of SC models through recasting .....	62
13.2.6. Difficult optimization tasks can be solved via recasting .....	63
13.3. Discussion .....	64
13.4. Conclusions.....	65
13.5. Methods.....	65
13.5.1. Modeling strategies.....	65
13.5.2. Recasting non-linear models into power-law canonical models by increasing the number of variables .....	67
13.6. Acknowledgements.....	68
13.7. Author details.....	68

13.7.1.	Author's contributions .....	68
13.7.2.	Competing interests .....	68
13.7.3.	References.....	68
<b>14.</b>	<b>Paper 5: Identifying the preferred subset of enzymatic profiles in nonlinear kinetic metabolic models via multiobjective global optimization and Pareto filters .....</b>	<b>70</b>
14.1.	Introduction.....	71
14.2.	Results.....	75
14.2.1.	Obtention of the Pareto set.....	75
14.2.2.	Selection of preferred subset of solutions.....	76
14.3.	Discussion.....	78
14.4.	Methods.....	79
14.4.1.	Mathematical model: GMA representation .....	79
14.4.2.	Mutiobjective global optimization of metabolic networks described by a GMA model .....	82
14.4.3.	Normalization of the Pareto optimal solutions .....	84
14.4.4.	Pareto filters .....	84
14.5.	Acknowledgements.....	87
14.6.	References.....	87
14.7.	Figure legends.....	92
14.8.	Tables.....	99
<b>15.</b>	<b>Paper 6: On the use of Principal Component Analysis for reducing the number of environmental objectives in multi-objective optimization: Application to the design of chemical supply chains.....</b>	<b>100</b>
15.1.	Introduction.....	100
15.2.	Materials and methods .....	101
15.2.1.	Dimensionality reduction.....	101
15.2.2.	Principal Component Analysis .....	101
15.2.3.	Combined use of PCA and MOO .....	102
15.3.	Results and discussion .....	103
15.3.1.	Numerical examples.....	103
15.4.	Conclusions.....	108
15.5.	Nomenclature .....	108
15.6.	Acknowledgements.....	110

15.7.	Appendix A.....	110
15.8.	Mathematical formulation.....	110
15.9.	References.....	111
	<b>Appendices.....</b>	<b>113</b>
	<b>1. List of publications.....</b>	<b>114</b>
1.1.	Research articles.....	114
1.2.	Book chapters.....	114
	<b>2. Congress contributions.....</b>	<b>115</b>
2.1.	Keynotes.....	115
2.2.	Oral presentations.....	115
2.3.	Posters.....	116

## **1. Introduction**

This thesis introduces a set of advanced mathematical programming tools for systems biology, where problems of interest range from designing microorganisms with enhanced phenotypes to understanding the evolution of cellular metabolism. Most of the tools presented are based in global optimization. The document is organized as follows. The challenging biological problems addressed in this thesis are briefly presented in section 2. The following section (section 3) provides a general background on mathematical programming and some of the techniques used in this thesis. Then, in section 4, we discuss how these techniques can be used to tackle the problems introduced in section 2, whereas in section 5, we illustrate the capabilities of these approaches as applied to similar problems. Finally, the conclusions of the work are drawn and future research lines that could extend the framework proposed herein are outlined.

### **1.1. Objectives of the thesis**

The objectives of this thesis are:

- Devise a systematic framework for the global optimization of metabolic networks described by the Generalized Mass Action (GMA) formalism.
- Extend this framework so as to solve to global optimality other types of kinetic models (i.e., Saturable and Cooperative systems).
- Develop a systematic framework for characterizing the feasible space of a biological optimization problem taking as a basis the strategy proposed by Sorribas and Guillén-Gosálbez [7].
- Use this framework to identify the best combination of constraints and objective function that shapes a prescribed adaptive process.
- Propose a multi-objective optimization (MOO) framework for metabolic engineering.

## **2. Systems biology**

The study of complex biological systems requires the integration of experimental and computational research by adopting a systems biology approach. Systems biology studies the interactions between the individual components of a biological system and how they determine the function and behavior of this system. Here, computational biology plays a major role by developing mathematical tools that aim to provide a powerful foundation from which to address critical scientific questions. In particular, the optimization of metabolic networks has emerged as a very important goal in biotechnology [8-12]. In addition, these techniques can also help in understanding the evolution of cellular metabolism under different conditions [7]. We review next the main topics studied in this thesis

## **2.1. Metabolic engineering (papers [1,2], also papers [3,4,5])**

In recent years, the use of genetic manipulation techniques has led to significant improvements in the production of certain biochemical products. However, in most cases mutation and selection of new processes has been made in a trial-and-error basis [13], which leads to sub-optimal solutions. As actual biological processes are operating far from their (mathematical) global optimum, one expects that they could be further improved if quantitative design principles for gene modification were provided by a more rational approach like optimization. This optimization is known as *engineering design* and consists of, given a model, finding the appropriate changes in the enzyme activities that optimize (maximize) a certain objective function (typically, the synthesis rate of the desired product). The enzyme activities obtained in the optimization can be reproduced in the real system by genetically modifying the expression of the associated genes.

## **2.2. Evolution studies (paper [3])**

While in optimization scientists establish the objective and search the way to accomplish it, in natural systems, the emergence of new designs results from the evolution driven by natural selection [14-17]. Physiological constraints force cellular mechanisms to modify the expression of their genes and their enzyme activities. The observed response should be an optimal (according to some criterion) assuring survival in a range of conditions [18-19]. Thus, while one can argue that natural systems are optimized by natural selection, it is not so clear which is the objective function in which this optimization is based or how close/far is this hypothetical optimum from the theoretical one that would result from a mathematical optimization. Determining these design principles can be posed as a reverse optimization problem, that is, we know the actual solution (the actual system) but we do not know which is the criterion (if any) optimized by this solution. Additionally, since the response has to fulfill some requirements, evolution may be even more closely related to feasibility than to optimization itself. Thus, a more complete statement would be that the adaptation of cellular metabolism consists of an optimal (in some sense) response that accomplishes a set of physiological constraints. These particular features of the evolution studies prevent standard global optimization tools from being directly applied to them. Hence, there is a need to customize such tools so they can deal with the associated problem complexity of the problem.

### 3. Mathematical programming: optimization

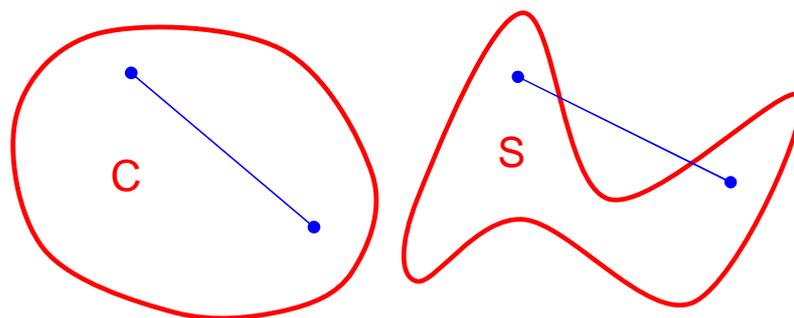
Although optimization started as a methodology of academic interest, it has become a useful technology that makes significant impact in various areas, including industrial applications [20]. In mathematical programming, optimization problems are generally posed as minimizations:

$$\begin{aligned} \text{SOO} \quad & \min \quad f_1 \\ & \text{s.t.} \quad h(x, y) = 0 \\ & \quad \quad g(x, y) \leq 0 \\ & \quad \quad x \in \mathfrak{X}, y \in Z \end{aligned}$$

Optimization problems as SOO are composed of different parts. On the one hand, the objective function  $f_1$  can be understood as the performance index of a given solution. Feasible alternatives (the set of which is sometimes referred to as search space or solution space) are defined by the constraints in the problem. In particular,  $h(x, y)$  represents equality constraints whereas  $g(x, y)$  refers to inequality constraint. Regarding the decision variables, these can either be continuous (denoted by  $x$ ) or integer (represented by  $y$ ). Note that widely-used binary variables are a particular case of the more general integer ones.

The nature of an optimization problem is given by the particular combination of variables and equations it embeds. As a result, one may face linear programming problems (LP, continuous variables and linear equations), non-linear programming problems (NLP, continuous variables and one or more non-linear equations), mixed-integer linear programming problems (MILP, continuous and integer variables, and linear equations) and mixed-integer non-linear programming problems (MINLP, continuous and integer variables, and at least one non-linear equation) among others. Special distinction needs to be made regarding whether the NLP is convex or not, as this second case may give rise to multiple local optimal solutions (i.e., multimodality). The existence of multiple sub-optima is a handicap when addressing these problems as standard algorithms may get trapped in them during the search, reporting a solution far from the global one.

An optimization problem is said to be convex when its objective function and its feasible space are both convex. A feasible space is convex if and only if the inequality constraints are convex and the equality constraints are affine (i.e., linear). In a convex search space, any linear combination of two points of the feasible space leads to a point belonging to the same space, whereas in a non-convex one, it does not (Figure 1). Note that according to this definition any problem involving integer variables is non-convex, since its solution space is defined by disjoint regions. In practice, however, MINLP formulations are in general referred to as non-convex only when the NLP resulting from fixing the values of their integer variables is non-convex (see section 3.2 for further information). Similarly, MILPs are non-convex because of the presence of binary variables.



*Figure 1. Example of a convex space  $C$  and a non-convex space  $S$ .*

### **3.1. Global optimization**

Deterministic global optimization strategies are the only ones that can ensure convergence to the global optimum of a non-convex problem within a desired tolerance in a finite number of iterations. Some of these methods have been implemented in software applications (for instance, a spatial branch and reduce algorithm is implemented in BARON, the state-of-art global optimization solver).

Here, we should distinguish between stochastic and deterministic approaches. Stochastic methods rely on meta-heuristics in order to guide the search for “good” solutions from a series of pseudorandom generated points. These methods are often based on physical and biological analogies and are capable of obtaining near optimal solutions in low CPU times, yet they offer no guarantee of global optimality for these solutions. On the other hand, as already mentioned, deterministic methods are rigorous and, thus, can guarantee global optimality within a desired optimality gap. These methods are based on calculating valid lower and upper bounds on the global optimum of the problem that are gradually tightened until a desired optimality criterion is satisfied. The main drawback of such strategies is that they require a large number of iterations to converge, and sometimes, even after large CPU times, they cannot close the optimality gap (defined as the absolute value of the relative difference between the upper and the lower bounds) below certain limits [21]. The search for the global optimum can be expedited by exploiting the mathematical properties of the specific problem. Hence, there is still room for improvement in this area by devising customized algorithms for specific applications. In this thesis, we have developed efficient deterministic global optimization techniques for non-convex NLPs and MINLPs arising in metabolic engineering studies. From now on, we will refer to deterministic global optimization simply as global optimization.

### **3.2. Relaxations in global optimization**

One key feature of any global optimization algorithm is its capability of predicting valid lower bounds on the global optimum. This is usually accomplished by solving a so-called relaxation. A relaxation is an auxiliary problem obtained with an objective function that underestimates the original one and a search space that contains that of the original problem. Because of these two properties, the relaxed problem provides a rigorous lower bound on the global optimum.

The objective of a global optimization algorithm is to approach the lower and upper bounds it produces to the globally optimal solution. In the case of the lower bound, this can only be accomplished by means of tight relaxations. Hence, in this thesis we studied how to obtain tight relaxations for the problems of interest.

### 3.3. Degeneracy

Different solutions in multimodal problems do not necessarily map into different objective function values. When this happens, that is, when different designs lead to the same objective function value, we say that the problem is degenerated. Accounting for degeneracy is of paramount importance for many applications, since different design may be associated with different practical implications that may make one solution particularly appealing among the others.

In this thesis, we worked on an iterative process to describe the feasible space of a model that can be used for studying degenerated problems. First, a grid is defined for each of the variables of interest and a binary variable is associated to each hyper-rectangle resulting from the intersection of the different grids. Next, the problem is optimized to obtain a given solution, which is allocated in a specific hyper-rectangle (note that the value of the binary variables defined before allows to identify the region in which that solution is located). Then, an integer cut, which is a special type of constraint, is added to the problem formulation to exclude from the search the hyper-rectangles found in previous iterations. The procedure is repeated again, and subsequent integer cuts are added until a predefined stop criterion is satisfied. The explanation of this algorithm is given in more detail in section *Feasibility approach* in [7] and in section 3.2 in [3].

### 3.4. Multi-objective optimization

Sometimes it might be interesting to evaluate alternatives considering more than one criterion. This can be accomplished by appending additional objectives to the problem formulation and by solving the resulting multi-objective optimization (MOO) problem:

$$\begin{aligned}
 MOO \quad \min \quad & F = \{f_1, \dots, f_B\} \\
 s.t. \quad & h(x, y) = 0 \\
 & g(x, y) \leq 0 \\
 & x \in \mathfrak{X}, y \in \{0, 1\}
 \end{aligned}$$

Recall that the difference between problem *MOO* and problem *SOO* relies on the objective function. In particular, in problem *SOO*,  $f_i$  can be regarded as a single objective function whereas in problem *MOO*,  $F$  is a vector containing a set of  $B$  objectives ranging from  $f_1$  to  $f_B$ .

The vector containing the individual minimum of all the objectives is regarded as the utopia point. This point is in general unattainable due to the trade-off existing between the different objectives. As a result, the solution to this kind of problems is usually composed by a set of points instead of a single one. These points are known as Pareto optimal solutions and form the so-called Pareto frontier. A solution is said to be Pareto optimal when it is not possible to improve one of the objectives without worsening any of

the others. For this reason, points in a given Pareto set are all considered to be equally optimal (see [22] for further information).

Many different methods, stochastic and deterministic, have been devised for obtaining Pareto optimal solutions in MOO problems. Here, we are only interested in the deterministic methods, which are the only ones that offer a theoretical guarantee of optimality in their solutions. In particular, the most popular deterministic MOO strategies are the weighted sum and the epsilon constraint method. The former method suffers from a well-reported inability to obtain non-convex parts of the Pareto frontier, which has motivated abandoning any further analysis on this alternative in this thesis (see [22] for a description of this method). The epsilon constraint does not show this limitation, and for this reason has been adopted in this work.

### 3.4.1. Epsilon constraint

In this method, one objective (main objective) is left as the only objective function of the problem, whereas the rest of the objectives (secondary objectives) are transferred to auxiliary constraints that impose bounds  $\varepsilon$  on them:

$$\begin{aligned}
 EC \quad & \min && f_1 \\
 & s.t. && f_b \leq \varepsilon_b \quad b = 2, \dots, B \quad n = 1, \dots, N \\
 & && h(x, y) = 0 \\
 & && g(x, y) \leq 0 \\
 & && x \in \mathfrak{X}, y \in \{0, 1\}
 \end{aligned}$$

The values of the epsilon parameters are obtained by first optimizing each objective individually and then splitting the interval defined by the best ( $\underline{f}_b$ ) and worst ( $\overline{f}_b$ ) values obtained for each objective in this optimization, into a set of subintervals (see section *Multiobjective global optimization of metabolic networks described by a GMA Model* in [5] for further details of the procedure).

One important feature of the epsilon constraint method is that it transforms a MOO problem into a set of single-objective problems, which can be solved by means of any global optimization method, if required. The number of instances of problem *EC* that must be solved is given by all possible combinations of epsilon parameters (i.e.,  $N^B$ ). With the aim of alleviating the computational burden, some authors [5] have proposed an alternative procedure in which only all bi-objective combinations are solved, which leads to  $\binom{B}{2}N$  instances.

### 3.4.2. Challenges in MOO

There are some well-reported difficulties which are intrinsic to MOO problems. For instance, even if the Pareto frontier can be computed, it would in general contain an infinite number of points for continuous problems, which inherently implies a big computational burden to deal with. Furthermore, visualization of the Pareto set becomes a

difficult task especially in problems with more than 3 objectives, which hampers the subsequent decision-making process [23].

Research efforts have been made for alleviating these difficulties. On the one hand, some authors have attempted to reduce the dimensionality of the problem by identifying redundant objectives that can be left out of the analysis. Some of these methods have been reviewed elsewhere [6]. On the other hand, other authors have developed strategies to produce only certain parts of the Pareto frontier, thereby reducing the number of solutions and facilitating the decision making procedure. This issue is further discussed in section *Pareto filters* in [5]. In this thesis work, we have integrated these methods into a complete MOO framework to tackle multi-objective metabolic engineering problems (see [5,6]).

#### 4. Application of optimization to systems biology problems

The use of mathematical optimization to improve biotechnological processes has the potential to produce significant economical savings. This is due to the reduction in the number of experiments required to improve the performance of the microorganisms and obtain higher yields. Additionally, the solutions of the optimization procedure can provide valuable insight on the behavior of the biological systems, thereby enhancing our understanding of cellular metabolism.

One of the key steps in this approach is the selection of the appropriate mathematical model among the different representations available. Although we can distinguish between different types of models (see for instance section 1 in [1] for further discussion), kinetic models based on the so-called power-law formalism show a good compromise between accuracy and simplicity. Among them, we find the S-System and General Mass Action (GMA) representations, which seem a promising alternative in the area. The main advantage of these models is that they can capture the non-linearities required to describe the regulatory processes of the networks while still showing some linear properties in the logarithmic space (i.e., when logarithms are taken in the power-law equations).

GMA models only differ from S-System models in the way in which the branching points are handled. In S-System models all the input flows in the branching point are collected and modeled together as if they were a single flow. The same procedure is followed for the outputs so that, finally, the concentration of the metabolite being balanced is the result of just two contributions. On the other hand, in GMA models each process is approximated separately so that there are as many contributions as actual flows in the real system. In particular, the GMA mathematical representation of a metabolic network containing  $n$  internal metabolites whose concentration varies due to the action of  $p$  flows can be expressed as follows:

$$\frac{dX_i}{dt} = \sum_{r=1}^p \left( \mu_{ir} K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \right) = 0 \quad i = 1, \dots, n \quad (1)$$

where  $\mu_{ir}$  is the stoichiometric coefficient of metabolite  $i$  in process  $r$ ,  $K_r$  is the fold-change produced over the basal-state enzyme activity  $\gamma_r$ ,  $X_j$  corresponds to the

concentration of metabolite  $j$  and  $f_{ij}$  is the kinetic order of metabolite  $X_j$  in process  $r$ , and quantifies its effect on the considered rate. Note that contributions of the  $m$  (independent) external metabolites are also accounted for in this representation. The reader is referred to section 3.1 in [1] and section 3 in [2] for a more detailed development of this expression.

#### 4.1. Metabolic engineering with GMA models

The identification of the enzymatic profile leading to an enhanced phenotype in a given microorganism can be obtained by solving an optimization problem where the maximum flux or yield is sought, subject to the equations describing the microorganisms' metabolism.

$$\begin{aligned}
 \text{GMAO} \quad & \min \quad f_1 = - \sum_{r \in FP_i} \mu_{ir} \nu_r \quad i \in FP \\
 \text{s.t.} \quad & \frac{dX_i}{dt} = \sum_{r=1}^p \left( \mu_{ir} K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{ij}} \right) = 0 \quad i = 1, \dots, n \\
 & h(K_r, X_j^{f_{ij}}, y) = 0 \\
 & g(K_r, X_j^{f_{ij}}, y) \leq 0 \\
 & K_r, X_j^{f_{ij}} \in \mathfrak{R}, y \in \{0,1\}
 \end{aligned}$$

Here  $FP$  is the set of metabolites  $i$  that are final products and  $FP_i$  is the set of processes  $r$  contributing to the production of metabolite  $i$ . Equality and inequality constraints ( $h(K_r, X_j^{f_{ij}}, y)$  and  $g(K_r, X_j^{f_{ij}}, y)$ , respectively) may be used to impose bounds on metabolites concentrations, enzymatic changes and the number of enzymes that can be modified simultaneously, among other things. Continuous variables denote metabolite concentrations and fold-changes in enzyme activities, whereas binary variables represent the number of enzymatic modulations simultaneously allowed. A detailed description of different variations of this model can be found in section 3.2 in [1] and section 3 in [2]. Recall that if all the enzymes can be modified at will, then problem *GMAO* leads to a non-convex NLP since binary variables are dropped from the formulation.

Regardless of the particular instance addressed, the complexity of this MINLP (or NLP) formulation that embeds a GMA model stems from the non-convexities introduced by the sigmoidal terms of the power-law formalism. The resulting (non-convex) MINLP problem may contain multiple local optima where standard optimization packages may converge during the search. In contrast, deterministic global optimization methods ensure convergence to the global optimum within a desired tolerance in a finite number of iterations. In particular, we have devised two global optimization algorithms for these problems: an outer approximation (OA, see section 4 in [1]) and a spatial branch-and-bound (sBB, refer to section 4 in [2]). The first one, which is based on the works by Polisetty [14] and Bergamini [24], has been used to maximize the ethanol synthesis rate in *Saccharomyces cerevisiae* and the citric acid production in *Aspergillus niger* (section 5 in [1]) whereas the later has also been used for the citric acid optimization in *Aspergillus niger* (section 5 in [2]). Furthermore, the performance of both algorithms has also been compared to that of BARON in [2] (see sections 5 and 6).

## 4.2. Metabolic engineering with other kinetic models

In some metabolic networks, it may be necessary to adopt representations other than the GMA formalism in order to reproduce the network's behavior precisely. These new models are likely to be more complex than the GMA formalism. Hence, in such cases, it may be convenient to resort to recasting, a technique which transforms the non-linear model with an arbitrary form into a canonical GMA model [25,26]. Through this technique, which requires definition of additional variables, arbitrary non-linear models can be represented using canonical forms such as GMA or S-system which are mathematically less demanding and hence can be used for simulation and optimization purposes. This opens the door for effectively extending any optimization tool originally devised for GMA models to other more detailed kinetic models. In particular, recasting is immediate for a specific type of model known as Saturable and Cooperative (SC). This formalism, which was proposed by Sorribas et al. [27], extends the power-law representation used in GMA models to account for cooperativity and saturation, which leads to more accurate predictions over a wider range of conditions than both the S-System and GMA representations.

The SC representation of a metabolic network with  $n$  internal metabolites and  $p$  flows is as follows:

$$\frac{dX_i}{dt} = \sum_{r=1}^p \left( \frac{K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}}}{\prod_{j=1}^{n+m} (\delta_{rj} + X_j^{f_{rj}})} \right) = 0 \quad i = 1, \dots, n \quad (2)$$

where  $\delta_{rj}$  is a parameter that depends on the saturation fractions (refer to section 5.1 in [4] for further details on this issue). Recall that variable  $K_r$ , which models changes in enzymes activities, is also included in this expression. The main difference between this equation and eqn 1, that is, between the GMA model and the SC model, is given by the denominator in the right-hand side of eqn 2. This denominator prevents the direct application of the strategies devised for the global optimization of GMA models for optimizing SC models. Nevertheless, a variable change is enough to by-pass this difficulty. In particular, a new variable  $Z_{rj}$  defined as in eqn 3 must be introduced in the model.

$$Z_{rj} = \delta_{rj} + X_j^{f_{rj}} \quad r = 1, \dots, p \quad j = 1, \dots, n + m \quad (3)$$

With this variable change, eqn 4 becomes analogous to that of the GMA model (eqn 1), which permits the application of the global optimization algorithm devised for GMA models to the SC formulation.

$$\frac{dX_i}{dt} = \sum_{r=1}^p \left( \mu_{ir} K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} Z_{rj}^{-1} \right) = 0 \quad i = 1, \dots, n \quad (4)$$

Note that eqn 3 must also be included in the final optimization problem, which can be finally posed as follows:

$$\begin{aligned}
 rGMA \min \quad & f_1 = - \sum_{r \in FP_i} \mu_{ir} \nu_r \quad i \in FP \\
 s.t. \quad & \frac{dX_i}{dt} = \sum_{r=1}^p \left( \mu_{ir} K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} Z_{rj}^{-1} \right) = 0 \quad i = 1, \dots, n \\
 & Z_{rj} = \delta_{rj} + X_j^{f_{rj}} \quad r = 1, \dots, p \quad j = 1, \dots, n+m \\
 & h(x, y) = 0 \\
 & g(x, y) \leq 0 \\
 & x \in \mathfrak{R}, y \in \{0, 1\}
 \end{aligned}$$

Model *rGMA* is mathematically identical to the model embedding the SC equations, but its structure is much more suitable for optimization purposes, as it allows the application of the global optimization strategies devised for GMA models. This is explained in more detail in [4], where some benchmark case studies considering different optimization goals were solved. To our knowledge, this is the first work that introduces a deterministic global optimization method for kinetic models other than GMA or S-Systems.

#### 4.3. Multi-objective global optimization in metabolic engineering

The design of strains with enhanced performance must be accomplished bearing in mind several conflicting criteria. For instance, some authors have argued that the minimization of metabolite concentrations should be regarded as an optimality principle in metabolic networks [28]. If this criterion were to be included in a traditional metabolic engineering optimization problem, then we would face a MOO problem. In previous attempts, this difficulty had been by-passed by sticking to the single-objective problem and simply imposing constraints to ensure cell viability. Note however that this strategy provides no information regarding the trade-off between objectives, and hence, may miss solutions in which cell viability is significantly strengthened to attain a marginal payoff in the final product synthesis rate.

Although MOO contributions are abundant in the chemical engineering field, they are much more uncommon among the metabolic engineering community. One significant contribution was made by Sendín et al. [29], who posed and solved by means of different MOO methods a 6-objectives MOO problem with two different models (S-System and an ad-hoc) of the *Saccharomyces cerevisiae* metabolic network. The objectives considered there were the ethanol synthesis rate and the concentration of 5 dependent metabolites. Most of the strategies compared therein are not suitable for problems based on GMA models as they either rely on local solvers or on stochastic methods which cannot offer any proof of global optimality (and hence, of obtaining the true Pareto front). In this thesis work, we have filled this gap by proposing a novel systematic framework for the deterministic multi-objective global optimization of metabolic networks described by the GMA model (work accepted in PLoS ONE [5]). In particular, a 14-objectives optimization problem considering the ethanol synthesis rate, the concentration of 5 dependent metabolites and the fold-change of the activity of 8 enzymes in the

*Saccharomyces cerevisiae* metabolic pathway has been solved using the epsilon constraint method coupled with an outer approximation algorithm. We have also explored the use of Pareto filters as a manner to reduce the size of the final set of candidates to be tested in the laboratory. The integration of these filters within a unique deterministic MOO framework constitutes one of the contributions of this thesis.

#### **4.4. Feasibility analysis in evolution studies**

The advantages of using systematic tools in evolution studies have already been acknowledged in the literature [13]. Nevertheless, there is still a lack of global optimization strategies in the area. For instance, Vilaprinyo et al. [16] studied the adaptation of yeast to heat shock by combining experimental data and computational work. Experimental observations were used to identify the physiological constraints regulating the problem, while brute force calculations were employed to determine feasible combinations of enzyme activities fulfilling these constraints. The procedure followed consists of choosing a set of discrete values for each of the different enzyme activities and then solving a system of non-linear equations (note that once the enzyme activities' fold-change  $K_r$  have been fixed, the degrees of freedom in the GMA model are 0) to determine the concentration of intermediate metabolites for each of the combinations of enzyme activities. If the system of equations rendered infeasible, the combination was regarded as infeasible. The main limitation of this approach is that the choice of discrete values prevents the proper exploration of the feasible space of enzyme activities as feasibility is not checked between two consecutive values. Thus, the analysis is not exhaustive and hence, the conclusions drawn are not general.

This thesis works provides a new generation of global optimization tools for evolution studies in systems biology. These strategies are based on those used in mathematical programming to deal with degenerated problems. In particular, the algorithm described in section 3.3 is used to characterize the feasible region defined by a set of biological conditions that model the adaptive process of a microorganism to a given stress. If a set of constraints leads to a feasible region that does not contain the experimental observations, then it can be discarded. The biological implications of this methodology are further described in section 3.4.4 in [3], whereas section 3.2 of the same publication is devoted to the technical details of the algorithm.

Once a valid set of constraints has been identified, different single-objective optimization problems considering alternative objective functions can be solved. The proximity between these theoretical optima and real observations may be an indicator of the driving force of these adaptive processes. To illustrate this methodology, we solved a complete example based on the adaptation of yeast to heat shock that can be found in [3]. Note that in this particular example the set of candidate constraints was retrieved from [16].

Note that the use of global optimization tools is especially important in systems biology studies, in which one aims to draw general conclusions from the specific properties of the solution found. In this context, local solutions must be avoided since they might hamper the whole biological analysis.

#### 4.4.1. Extension of feasibility analysis for metabolic engineering problems

The methodology described in the previous section does not require the identification of the global optimum of a model, since only feasibility must be checked. However, one could use the same methodology to characterize the region of the feasible search space that shows a given objective function value. This may be particularly interesting for metabolic engineering studies, by identifying different enzymatic profiles leading to similar yields. One could then choose enzymatic modifications that do not compromise cell survival. This strategy was presented [3], where a case study involving the ethanol production rate in *Saccharomyces cerevisiae* was included (section 3.4.2).

### 5. Other applications of mathematical programming

Optimization provides set of powerful tools intended to give answer to problems that otherwise would be very difficult (if not impossible) to solve. As a result, these techniques are applicable to many different and diverse areas. So far, we have shown how global optimization and related strategies can be effectively applied to address problems that are relevant in the context of systems biology. Next, we discuss how mathematical programming can be applied to other unrelated areas such as the design and planning of suitable processes using life cycle assessment (LCA) principles.

Consider a supply chain management (SCM) problem in which the goal is to determine the set of Pareto optimal supply chain configurations that maximize the net present value (NPV) and minimize a set of environmental metrics (see section 3 in [6] for a detailed problem statement). Intuitively, one may think that there exists a trade-off between the economic and the environmental objectives, but it is not so straight-forward to predict whether the environmental metrics are introducing independent or redundant information to the problem. To detect redundant objectives in a MOO, Deb and Saxena [30] proposed a heuristic method which is based on Principal Component Analysis (PCA). In their work, they apply stochastic MOO methods. As previously discussed, these methods may lead to spurious Pareto frontiers. Hence, in this thesis we overcome this limitation by developing a deterministic MOO algorithm that is coupled with a dimensionality reduction procedure (see CES [6]). This framework is complementary to that presented for the MOO of metabolic network (section 4.3), as it attempts to reduce the computational burden (but this time by reducing the number of objectives rather than by filtering Pareto solutions).

### 6. Conclusions

The results obtained after developing and applying the techniques presented in this work have provided a set of conclusions which are listed below.

- Two systematic frameworks for the global optimization of metabolic networks described by the GMA representation have been developed: one based on an OA and another on a customized sBB.
- Our customized algorithms show better numerical performance than other approaches proposed so far in the literature for the global optimization of metabolic networks [8].

- Numerical results show that the two proposed algorithms outperform the state-of-the-art commercial global optimization solver BARON. This is due to the high quality relaxations obtained by exploiting the mathematical structure of GMA models.
- Although none of the customized methods (sBB and OA) proved to be superior in all of the cases, we observed that sBB shows better performance in the most complicated instances.
- It has been shown that it is possible to solve highly non-linear kinetic models (like the SC formalism) efficiently by recasting them into GMA-like canonical forms.
- A systematic framework for mapping the objective function surface and simultaneously determining the problem's feasible region has been devised. This methodology can be applied to evolution studies as well as to metabolic engineering problems, providing valuable insights into biological systems.
- After comparing different sets of constraints in the adaptive response of yeast to heat shock, we conclude that constraints proposed by Vilaprinyo et al. [10] are the ones which better explain the experimental observations. Besides, results indicate that the actual process may be driven by the minimization of the total enzymatic cost. In other words, it seems that cells try to minimize the number of metabolic changes required to survive in the new conditions.
- Two systematic frameworks have been proposed for MOO problems. On the one hand, the use of PCA allows overcoming the numerical difficulties that arise when dealing with a large number of environmental objectives in a SCM problem. Results indicate that, in general, few environmental metrics suffice to characterize the environmental performance of a given solution. According to this finding, one could exclude many LCA metrics from the pool without compromising the quality of the Pareto structure.
- On the other hand, the application of two Pareto filters in metabolic engineering can reduce the number of solutions to a large extent, thereby leading to a reasonable number of alternatives to be tested in the laboratory. This illustrates the usefulness of the proposed approach.

## **7. Future work**

We next introduce a set of potential research lines related to the material presented in this thesis.

- The systematic framework proposed for the multi-objective global optimization of GMA-described metabolic networks could be extended to other kinetic models by combining it with the recasting strategy. In the particular case of SC models, for which the recasting has already been reported, only linking of both methods is missing.

- Evolution studies could be tackled by means of a MOO approach since there might be no single objective driving the adaptive processes of living-cells, but rather a set of biological criteria.
- The two systematic frameworks devised for MOO problems could be coupled into a single holistic framework. That is, after solving the MOO problem by means of the epsilon constraint method, the dimensionality of the problem could be reduced by removing redundant objectives. This would reduce the computational effort wasted in calculating duplicated information, and more diverse alternatives will be generated. Then, different Pareto filters could be applied in order to select a final set of candidate solutions. Furthermore, this strategy could be integrated with a recasting to tackle any kind of kinetic model.
- We would like to apply the tools developed herein to other biological systems. The final goal should be to develop genome-wide models.
- The technique employed in this thesis to relax the logarithmic function by the combined use of supporting hyper-planes and piecewise linear functions is equally valid for other monotonic decreasing/increasing functions. Hence, it may be interesting to compare the quality of the relaxation provided by this strategy with that obtained using other known relaxation techniques (for instance, with the McCormick's envelopes for bilinear terms).
- The customized algorithms presented in this thesis include a set of parameters which can be configured at will. For instance, in the OA and the sBB one can select the number of piecewise sections that are initially employed for the calculations. It would be interesting to seek for systematic tuning strategies for these parameters.
- All the strategies presented in this thesis are valid for steady-state conditions. It may hence be interesting to extend the proposed methodologies to deal with dynamic systems.
- The development of a software application implementing the algorithms developed herein is crucial for extending their use.

## 8. Nomenclature

### 8.1. Abbreviations

GMA	Generalized Mass Action
LCA	Life Cycle Assessment
LP	Linear Programming
MILP	Mixed-integer Linear Programming
MINLP	Mixed-integer Non-linear Programming
MOO	Multi-objective Optimization

NLP	Non-linear Programming
NPV	Net Present Value
OA	Outer approximation
PCA	Principal Component Analysis
sBB	Spatial branch-and-bound
SC	Saturable and Cooperative
SCM	Supply Chain Management

## 8.2. Indices

$b$	Objective function
$i$	Internal (dependent) metabolite
$j$	Metabolite (can either be internal or external)
$n$	Epsilon constraint subinterval
$r$	Process (flux) in a metabolic network

## 8.3. Sets/Subsets

$FP$	Set of metabolites $i$ that are final products
$FP_i$	Set of processes $r$ contributing to the production of metabolite $i$

## 8.4. Parameters

$\delta_{rj}$	Characteristic SC parameter that depends on the saturation fraction of metabolite $j$ in process $r$
$\varepsilon_b^n$	Epsilon parameter for subinterval $n$ on objective $b$
$\gamma_r$	Basal-state activity of enzyme governing process $r$
$\mu_{ir}$	Stoichiometric coefficient of the metabolite $i$ in process $r$
$B$	Total number of objectives
$\underline{f}_b$	Lower bound on objective $f_b$
$\overline{f}_b$	Upper bound on objective $f_b$
$f_{rj}$	Kinetic order of metabolite $X_j$ in process $r$
$m$	Number of external (independent) metabolites
$n$	Number of internal (dependent) metabolites
$N$	Total number of epsilon subintervals

## 8.5. Variables

$f_b$	Individual objective function
-------	-------------------------------

$F$	Vector of objectives functions
$K_r$	Fold-change produced over the basal-state enzyme activity $\gamma_r$
$x$	Generic continuous variable
$X_j$	Concentration of metabolite $j$
$y$	Generic binary variable
$Z_{rj}$	Continuous variable introduced to recast SC models into GMA models

## 9. References

- [1] Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: Outer approximation-based algorithm for biotechnology studies in systems biology. *Computers and Chemical Engineering* 2010, 34(10), 1719-1730.
- [2] Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: A spatial branch-and-bound framework for the global optimization of kinetic models of metabolic networks. *Industrial and Engineering Chemistry Research* 2011, 50(9), 5225-5238
- [3] Sorribas, A., Pozo, C., Vilaprinyo, E., Guillén-Gosálbez, G., Jiménez, L., Alves, R.: Optimization and evolution in metabolic pathways: Global optimization techniques in Generalized Mass Action models. *Journal of Biotechnology* 2010, 149(3), 141-153.
- [4] Pozo, C., Marín-Sanguino, A., Alves, R., Guillén-Gosálbez, G., Jiménez, L., Sorribas, A.: Steady-state global optimization of metabolic non-linear dynamic models through recasting into power-law canonical models. *BMC Systems Biology* 2011, 5, art. no. 137.
- [5] Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: Identifying the preferred subset of enzymatic profiles in nonlinear kinetic metabolic models via multiobjective global optimization and Pareto filters. Accepted for publication in *PLoS ONE*.
- [6] Pozo, C., Ruíz-Femenia, R., Caballero, J., Guillén-Gosálbez, G., Jiménez, L.: On the use of Principal Component Analysis for reducing the number of environmental objectives in multi-objective optimization: Application to the design of chemical supply chains. *Chemical Engineering Science* 2012, 69(1), 146-158.
- [7] Guillén-Gosálbez, G., Sorribas, A., 2009. Identifying quantitative operation principles in metabolic pathways: a systematic method for searching feasible enzyme activity patterns leading to cellular adaptive responses. *BMC Bioinformatics* 10, 386.
- [8] Grossmann, I.E., Biegler, L.T.: Part II. Future perspective on optimization. *Computers and chemical engineering* 2004, 28, 1193-1218.
- [9] Bailey, J., Birnbaum, S., Galazzo, J., Khosla, C., Shanks, J.: Strategies and challenges in metabolic engineering. *Ann. NY Acad. Sci.* 1990, 589:1-15.
- [10] Cameron, D., Chaplen, F.: Developments in metabolic engineering. *Curr. Opin. Biotechnol.* 1997, 8, 175-180.

- [11] Cameron, D., Tong, J.: Cellular and metabolic engineering: an overview. *Appl. Biochem. Biotechnol.* 1993, 38, 105–140.
- [12] Savageau, M.A.: Biochemical Systems Analysis: A Study of Function and Design in *Molecular Biology*. Reading, Mass.: Addison-Wesley 1976.
- [13] Vera, J., de Atauri, P., Cascante, M., Torres, N.V.: Multicriteria optimization of biochemical systems by linear programming: application to production of ethanol by *Saccharomyces cerevisiae*. *Biotechnol Bioeng* 2003, 83(3), 335–43.
- [14] Polisetty, P.K., Gatzke, E.P., Voit, E.O.: Yield optimization of regulated metabolic systems using deterministic branch-and-reduce methods. *Biot Bioeng* 2008, 99(5), 1154–69.
- [15] Voit, E.O.: Optimization in integrated biochemical systems. *Biotechnol Bioeng* 1992, 40(5), 572–82.
- [16] Vilaprinyo, E., Alves, R., Sorribas, A.: Use of physiological constraints to identify quantitative design principles for gene expression in yeast adaptation to heat shock. *BMC Bioinformatics* 2006, 7, 184.
- [17] Causton, H.C., Ren, B., Koh, S.S., Harbison, C.T., Kanin, E., Jennings, E.G., Lee, T.I., True, H.L., Lander, E.S., Young, R.A.: Remodeling of yeast genome expression in response to environmental changes. *Molecular biology of the cell* 2001, 12(2), 323–337.
- [18] Gasch, A.P., Spellman, P.T., Kao, C.M., Carmel-Harel, O., Eisen, M.B., Storz, G., Botstein, D., Brown, P.O.: Genomic expression programs in the response of yeast cells to environmental changes. *Molecular biology of the cell* 2000, 11(12), 4241–4257.
- [19] Molina-Navarro, M.M., Castells-Roca, L., Belli, G., Garcia-Martinez, J., Marin-Navarro, J., Moreno, J., Perez-Ortin, J.E., Herrero, E.: Comprehensive transcriptional analysis of the oxidative response in yeast. *The Journal of biolog. chem.* 2008, 283(26), 17908–17918.
- [20] Vecchiotti, A., Sangbum, L., Grossmann, I.: Modeling of discrete/continuous optimization problems: characterization and formulation of disjunctions and their relaxations. *Computers and chemical engineering* 2003, 27, 433–448.
- [21] Marin-Sanguino, A., Voit, E.O., Gonzalez-Alcon, C., Torres, N.V.: Optimization of biotechnological systems through geometric programming. *Theor Biol Med Model* 2007, 4, 38.
- [22] Ehrgott, M.: *Multicriteria Optimization*. Berlin: Springer, 1998..
- [23] Das, I.: A preference ordering among various pareto optimal alternatives. *Structural Optimization*, 1999, 18, 30–35.
- [24] Bergamini, M., Aguirre, P., Grossmann, I.: Logic-based outer approximation for globally optimal synthesis of process networks. *Computers and Chemical Engineering* 2005, 29, 1914–1933.
- [25] Savageau, M.A., Voit, E.O.: Recasting nonlinear differential equations as Ssystems: a canonical nonlinear form. *Mathematical Biosciences* 1987, 87, 83–115.

[26] Voit, E.O.: Recasting nonlinear models as S-systems. *Mathematical and Computer Modelling* 1988, 11(C), 140-145.

[27] Sorribas, A., Hernandez-Bermejo, B., Vilaprinyo, E., Alves, R.: Cooperativity and saturation in biochemical networks: a saturable formalism using Taylor series approximations. *Biotechnology and bioengineering* 2007, 97(5), 1259-1277.

[28] Schuster, S., Schuster, R., Heinrich, R.: Minimization of intermediate concentrations as a suggested optimality principle for biochemical networks. II. Time hierarchy, enzymatic rate laws, and erythrocyte metabolism, *J. Math. Biol.* 1991, 29, (5), 443-455.

[29] Sendín, J., Vera, J., Torres, N., Banga, J.: Model-based optimization of biochemical systems using multiple objectives: a comparison of several solution strategies. *Math Comput Model Dyn Syst* 2006, 12(5), 469-487.

[30] Deb, K., Saxena, D.: On finding pareto-optimal solutions through dimensionality reduction for certain large-dimensional multi-objective optimization problems. *Technical Report. Kanpur Genetic Algorithms Laboratory (KanGAL)*, Indian Institute of Technology, Kanpur, 2005.



## Outer approximation-based algorithm for biotechnology studies in systems biology

Carlos Pozo<sup>a</sup>, Gonzalo Guillén-Gosálbez<sup>a,\*</sup>, Albert Sorribas<sup>b</sup>, Laureano Jiménez<sup>a</sup>

<sup>a</sup> Departament d'Enginyeria Química, Universitat Rovira i Virgili, Avinguda Paisos Catalans 26, 43007 Tarragona, Spain

<sup>b</sup> Departament de Ciències Mèdiques Bàsiques, Institut de Recerca Biomèdica de Lleida (IRBLLEIDA), Universitat de Lleida, Montserrat Roig 2, 25008 Lleida, Spain

### ARTICLE INFO

#### Article history:

Received 31 July 2009

Received in revised form 15 February 2010

Accepted 1 March 2010

Available online 7 March 2010

#### Keywords:

Global optimization

Generalized Mass Action (GMA)

Metabolic engineering

### ABSTRACT

Optimization methods play a central role in systems biology studies as they can help in identifying key processes that can be experimentally changed so that specific biological goals can be attained. Standard optimization methods used in this field rely on simplified linear models that may fail in capturing the underlying complexity of the target metabolic network. Within this general context, we present a novel approach to globally optimize metabolic networks. The approach presented relies on (1) adopting a general modeling framework for metabolic networks: the Generalized Mass Action (GMA) representation; (2) posing the optimization task as a non-convex nonlinear programming (NLP) problem; and (3) devising an efficient solution method for globally optimizing the resulting NLP that embeds a GMA model of the metabolic network. The capabilities of our method are illustrated through two case studies: the anaerobic fermentation pathway in *Saccharomyces cerevisiae* and the citric acid production using *Aspergillus niger*. Numerical results show that the method presented provides near optimal solutions in low CPU times even in cases where the commercial global optimization package BARON fails to close the optimality gap.

© 2010 Elsevier Ltd. All rights reserved.

### 1. Introduction

The study of complex biological systems requires the integration of experimental and computational research by adopting a systems biology approach. Systems biology addresses the study of the interactions between the individual components of a biological system through the integration of data and mathematical models. Here, computational biology plays a major role by developing mathematical tools that aim to provide a powerful foundation from which to address critical scientific questions. In particular, the optimization of metabolic networks has emerged as a very important goal in biotechnology (Bailey, Birnbaum, Galazzo, Khosla, & Shanks, 1990; Banga, 2008; Cameron & Chaplen, 1997; Cameron & Tong, 1993; Mendes & Kell, 1996; Torres & Voit, 2002).

In recent years, the use of genetic manipulation techniques has led to significant improvements in the production of certain biochemical products. However, in most cases mutation and selection of new processes have been made in a trial-and-error basis, which has led to local optimal solutions. Hence, one expects that actual biological processes could be further improved if quantitative design principles for the modification of the genes were provided by a more rational approach like optimization (Banga, 2008; Chang &

Sahinidis, 2005; Hatzimanikatis, Floudas, & Bailey, 1996; Polisetty, Gatzke, & Voit, 2008; Vera, de Atauri, Cascante, & Torres, 2003; Voit, 1992). This optimization is known as *metabolic engineering* (Bailey et al., 1990; Bailey, 1991, 1999) and consists of, given a model, finding the appropriate changes in the enzyme activities that optimize (maximize) a certain objective function (typically, the synthesis rate of the desired product). The enzyme activities obtained in the optimization solution can be implemented in the real system by tuning the expressions of the corresponding genes.

The use of mathematical optimization to improve biotechnological processes is nowadays gaining wider acceptance given their potential to produce significant economical savings. These may be achieved by reducing the number of experiments required to find those microorganisms that lead to higher yields. Furthermore, as manipulation of many enzymes at once may be prohibitive, a theoretical analysis on the more promising alternative combinations of limited changes is of great practical interest. Additionally, the solutions of the optimization procedure can provide valuable insights on the behavior of the biological systems, making these techniques useful in other applications such as evolution studies (Guillén-Gosálbez & Sorribas, 2009).

One of the key steps in this approach is the selection of the appropriate mathematical model among the different representations available. Here, we can distinguish between three main groups of models. The first group corresponds to stoichiometric models. These models constitute simple linear representations of

\* Corresponding author.

E-mail address: [Gonzalo.Guillen@urv.cat](mailto:Gonzalo.Guillen@urv.cat) (G. Guillén-Gosálbez).

the stoichiometry of the network (i.e. network structure). However, their simplicity becomes at the same time their main limitation as they fail to capture the non-linear behavior of some key processes of the networks such as regulation (Gavalas, 1968; Heinrich & Schuster, 1996). On the other extreme of accuracy, we would find *ad hoc* models. These models rely on the formulation of detailed kinetics equations, such as Michaelis–Menten, that allow accounting for modulating effects. Unfortunately, optimizing these systems is not a straightforward task as it usually leads to complex mathematical formulations (Polisetty et al., 2008). A third group of models includes representations that result from the combination of linear stoichiometric descriptions and non-linear approximate representations to express the velocities of the metabolic reactions (Alves, Vilaprinyo, Hernández-Bermejo, & Sorribas, 2009; Sorribas, Hernández-Bermejo, Vilaprinyo, & Alves, 2007). Among them, models using the so called power-law formalism show a good compromise between accuracy and simplicity (Marin-Sanguino, Voit, Gonzalez-Alzon, & Torres, 2007). This group includes the S-System and the General Mass Action (GMA) models, which seem a promising alternative in the area (Voit, 1992, 2003). The main advantage of these models is that they can capture the non-linearities required to describe the regulatory processes of the networks while showing linear properties in the logarithmic space. Additionally, these models constitute a very general framework since any kind of metabolic network can be represented through their formulations (Alves et al., 2009).

GMA models only differ from S-System models in the way in which the branching points are handled (Curto, Sorribas, & Cascante, 1995). In S-System models, all the input flows in the branching point are collected and modeled together as if they were a single flow. The same procedure is followed for the outputs so that, finally, the concentration of the metabolite being balanced is the result of just two contributions. On the other hand, in GMA models each process is approximated separately so that there are as many contributions as actual flows in the real system (Voit, 2000 and references therein). If the metabolic network only contains nodes that result from the contribution of an input flow and an output flow, the S-System and GMA representation coincide.

Models based on the power-law formalism were first used in metabolic optimization problems by Voit (1992). The choice of an S-Systems representation allowed him to obtain a linear representation by a simple logarithmic transformation performed on some variables of the model (Alvarez-Vasquez, Canovas, Iborra, & Torres, 2002; Marin-Sanguino & Torres, 2003; Marin-Sanguino et al., 2007; Vecchiotti, Sangbum, & Grossmann, 2003). However, this is not possible in GMA models, since some equations cannot be directly reformulated using the logarithmic transformation. The optimization task then gives rise to a non-convex NLP that may show multiple local optima in which standard commercial packages can get trapped during the search.

In the context of performing a systems biology study, global optimality is particularly important, as one aims to draw general conclusions from the specific properties of the solution found. Hence, local solutions should be avoided, since they might hamper the entire biological analysis by providing insights that are not meaningful at all. A literature review in the area of global optimization of metabolic networks (Banga, 2008) reveals that this is indeed a ripe field for research. In a recent and pioneering work Polisetty et al. (2008) addressed the global optimization of GMA models (see also Marin-Sanguino et al., 2007; Marin-Sanguino & Torres, 2003). The main drawback of the strategy presented by Polisetty et al. (2008) is that it provides solutions with large optimality gaps (i.e., large differences between the best solution that could be found and the one calculated during the execution of the algorithm). More recently Guillén-Gosálbez and Sorribas (2009) presented a novel algorithm that makes use of global optimization techniques for per-

forming feasibility analysis in evolution studies. The tool developed by these authors allowed characterizing the feasible space of optimization problems with embedded GMA models (Sorribas et al., 2010).

The aim of this work is to provide a systematic modeling framework and solution strategy for metabolic optimization problems arising in systems biology studies. The approach presented relies on posing the optimization task as a NLP with an embedded GMA model of the metabolic network under study. An outer-approximation algorithm is presented to solve this type of models to global optimality. We provide a theoretical analysis on some details of the algorithm and illustrate its capabilities through two examples, comparing our results with those produced by BARON, nowadays regarded as the “state of the art” global optimization package.

## 2. Problem statement

Given a metabolic network described by a GMA model, the optimization aims to determine the appropriate changes in enzyme activities and in the internal metabolite concentrations so that the synthesis rate of the desired product is maximized in steady state. Given data for the problem are: (1) the stoichiometry of the reactions involved in the production/consumption of each internal metabolite in the metabolic network; and (2) the value of the parameters of the power-law formalism representing the kinetics of each of these particular reactions at the basal state.

## 3. Mathematical formulation

### 3.1. GMA representation

The GMA representation of a metabolic network containing  $n$  internal metabolites whose concentration  $X_i$  can vary with the time  $t$  due to the action of  $p$  flows can be expressed as follows:

$$\frac{dX_i}{dt} = \sum_{r=1}^p \mu_{ir} v_r \quad i = 1, \dots, n \quad (1)$$

where  $\mu_{ir}$  is the stoichiometric coefficient of the metabolite  $i$  in the process  $r$  and indicates the number of molecules of metabolite  $i$  involved in such a process. Hence, it is always an integer value that is positive when process  $r$  contributes to the production of metabolite  $i$ , negative when process  $r$  consumes metabolite  $i$  and 0 otherwise (i.e., if process  $r$  does not participate in the production/consumption of metabolite  $i$ ). The velocity  $v_r$  can be described using different representations, but, as stated previously, the so-called power-law formalism (Savageau, 1969a,b; Voit, 2000) is an appropriate one:

$$v_r = \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \quad r = 1, \dots, p \quad (2)$$

In this representation,  $\gamma_r$  is an apparent rate constant for flow  $r$ .  $f_{rj}$  is the kinetic order of metabolite  $j$  in process  $r$  and quantifies its effect on the considered rate. Note that contributions of the  $m$  (independent) external metabolites are also accounted for in this representation.

By introducing Eq. (2) into Eq. (1) and assuming steady state conditions for the network, one obtains a GMA model as follows:

$$\frac{dX_i}{dt} = \sum_{r=1}^p \left( \mu_{ir} \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \right) = 0 \quad i = 1, \dots, n \quad (3)$$

### 3.2. NLP formulation

In order to compute the changes in the enzyme activities, we shall rewrite the apparent rate constant  $\gamma_r$  in Eq. (3) as a product of the basal state enzyme activity  $\gamma_r$  (constant parameter) and its fold-change  $K_r$  (continuous variable):

$$\sum_{r=1}^p \left( \mu_{ir} K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \right) = 0 \quad i = 1, \dots, n \quad (4)$$

The final goal of the optimization task is to find the appropriate changes to be performed in the enzyme activities in order to optimize a given biological criteria (typically a flow) described through algebraic equations. This requires the determination of the optimal values of  $K_r$ ,  $v_r$  and  $X_j$  that maximize/minimize the given objective function while fulfilling the GMA model equations in steady state. In general, it will be possible to express the desired criterion in mathematical terms using a specific mathematical function  $U(K_r, v_r, X_j)$ , so that the optimization task can be posed as a non-linear programming problem (NLP) of the following form:

$$\begin{aligned} (\text{ONLP}) \quad & \min U(K_r, v_r, X_j) \\ \text{s.t.} \quad & \sum_{r=1}^p \mu_{ir} v_r = 0 \quad i = 1, \dots, n \\ & v_r = K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \quad r = 1, \dots, p \\ & K_r, v_r, X_j \in \mathbb{R}_+ \end{aligned}$$

Note that maximization problems can be easily reformulated into minimization ones by changing the sign of the objective function. The nonlinear equality constraints that define the velocity terms in **ONLP** give rise to a non-convex search space. Hence, to solve **ONLP** to global optimality, it is necessary to resort to global optimization techniques (see Grossmann & Bigler, 2004; Floudas & Gounaris, 2009) that can provide solutions to the problem with a desired optimality tolerance. These methods can handle a wide variety of non-convex formulations arising in many types of applications. Unfortunately, in practice, their numerical performance may vary drastically depending on the specific problem being solved, leading in some cases to prohibitive CPU times (Grossmann & Bigler, 2004). A possible way to overcome this limitation consists of devising customized algorithms that exploit the mathematical properties of the specific problem under study. This is indeed the underlying idea of our approach.

### 4. Solution strategy

The method we propose to globally optimize **ONLP** is an outer-approximation algorithm based on the works of Bergamini, Aguirre, and Grossmann (2005) and Polisetty et al. (2008). Our method relies on decomposing the original problem **ONLP** into two problems at different hierarchical levels: an upper level master problem **CMILP** and a lower level slave problem **RNLP**. The master level entails the solution of a mixed-integer linear programming (MILP) problem that is a relaxation of **ONLP**. This implies that **CMILP** will predict valid lower bounds on the solution of **ONLP** (the solution of the relaxation will be, at least, as good as that of the original problem). In the lower level, the original problem is locally optimized in a reduced search space (**RNLP**) providing a valid upper bound on its global optimum. These two problems are solved iteratively until the optimality gap is reduced below a given tolerance. A detailed description of the algorithm is given in the following sections.

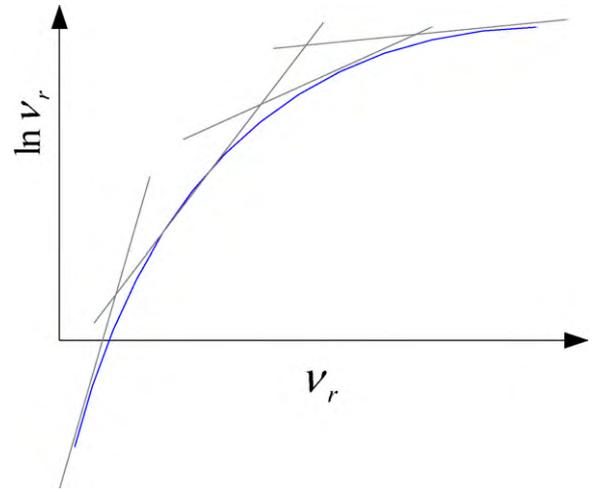


Fig. 1. Natural logarithm overestimation by a 1st degree Taylor series.

#### 4.1. Upper level master problem

To construct a valid relaxation of **ONLP** (i.e., **CMILP**), we first reformulate the equations arising from the power-law formalism via a logarithmic transformation:

$$\ln v_r = \ln K_r + \ln \gamma_r + \sum_{j=1}^{n+m} f_{rj} \ln X_j \quad r = 1, \dots, p \quad (5)$$

We then introduce two new auxiliary variables,  $k_r$  and  $x_j$ , which are defined as follows:

$$\begin{aligned} k_r &= \ln K_r \\ x_j &= \ln X_j \end{aligned}$$

By replacing the original variables in Eq. (5) by the reformulated ones, the following equality can be obtained.

$$\ln v_r = k_r + \ln \gamma_r + \sum_{j=1}^{n+m} f_{rj} x_j \quad r = 1, \dots, p \quad (6)$$

Eq. (6) can then be expressed via the following inequalities:

$$\ln v_r \geq k_r + \ln \gamma_r + \sum_{j=1}^{n+m} f_{rj} x_j \quad r = 1, \dots, p \quad (7)$$

$$\ln v_r \leq k_r + \ln \gamma_r + \sum_{j=1}^{n+m} f_{rj} x_j \quad r = 1, \dots, p \quad (8)$$

The logarithmic terms appearing in the left-hand side of these equations can be replaced by valid upper and lower estimators (note that  $\gamma_r$  is a known model parameter). Specifically, in Eq. (7), the logarithmic function can be approximated by  $L$  supporting hyper-planes (see Fig. 1), which are first order Taylor expansions of the natural logarithm at different points  $l$  of the domain of  $v_r$ :

$$\begin{aligned} \ln v_r^l + \frac{1}{v_r^l} (v_r - v_r^l) &\geq k_r + \ln \gamma_r + \sum_{j=1}^{n+m} f_{rj} x_j \quad r = 1, \dots, p \\ l &= 1, \dots, L \end{aligned} \quad (9)$$

Since the logarithmic function is concave, these hyper-planes constitute valid overestimators that do not chop off any feasible solution of **ONLP**.

Furthermore, the left-hand side of Eq. (8) can be underestimated by a piecewise linear approximation. For that, we consider a partition of the original domain  $[v_r, \bar{v}_r]$  defined by a set of grid

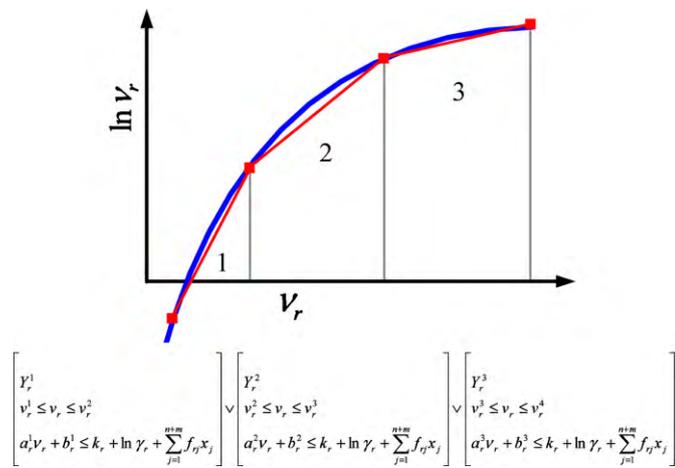


Fig. 2. Example of natural logarithm underestimation by piecewise linear functions.

points  $v_r^1, v_r^2, \dots, v_r^{H+1}$ , being  $v_r^1 = v_r$ ,  $v_r^{H+1} = \bar{v}_r$  and  $v_r^{h+1} \geq v_r^h$  for  $h = 1, \dots, H$ . The piecewise linear approximation can then be modeled via a disjunction with  $H$  terms as follows:

$$\bigvee_{h=1, \dots, H} \left[ \begin{array}{l} Y_r^h \\ v_r^h \leq v_r \leq v_r^{h+1} \\ a_r^h v_r + b_r^h \leq k_r + \ln \gamma_r + \sum_{j=1}^{n+m} f_{rj} x_j \end{array} \right] \quad r = 1, \dots, p$$

$$Y_r^h \in \{True, False\} \quad r = 1, \dots, p \quad h = 1, \dots, H$$

where  $a_r^h$  and  $b_r^h$  are the coefficients of the straight line equation in the  $h^{th}$  interval and  $Y_r^h$  indicates whether the  $h^{th}$  term in the disjunction of the  $r^{th}$  velocity is active or not. Fig. 2 shows an illustrative example of a piecewise function with three terms.

The disjunction can be reformulated using either the big-M or convex hull reformulations (see Vecchiotti et al., 2003). The latter technique allows translating the disjunction into a set of equalities and inequalities using auxiliary (disaggregated) variables as follows:

$$\sum_{h=1}^H z_r^h = v_r \quad r = 1, \dots, p \quad (10)$$

$$v_r^h y_r^h \leq z_r^h \leq v_r^{h+1} y_r^h \quad r = 1, \dots, p \quad h = 1, \dots, H \quad (11)$$

$$\sum_{h=1}^H y_r^h = 1 \quad r = 1, \dots, p \quad (12)$$

$$\sum_{h=1}^H (a_r^h z_r^h + b_r^h y_r^h) \leq k_r + \ln \gamma_r + \sum_{j=1}^{n+m} f_{rj} x_j \quad r = 1, \dots, p \quad (13)$$

where  $z_r^h$  is the new disaggregated variable and  $y_r^h$  is a new binary variable that takes a value of 1 if the  $h^{th}$  interval of the  $r^{th}$  velocity is active and 0 otherwise. Thus, the overall master problem can be finally expressed as follows:

$$\begin{array}{ll} \text{(CMILP)} & \min U(k_r, x_j, v_r, z_r^h, y_r^h) \\ & \text{s.t. constraints 1, 9, 10 to 13} \\ & k_r, x_j \in \mathbb{R} \\ & v_r, z_r^h \in \mathbb{R}_+ \\ & y_r^h \in \{0, 1\} \end{array}$$

Model CMILP takes the form of a mixed-integer linear programming (MILP) problem. These problems can be solved efficiently via standard branch & bound (B&B) techniques.

## 4.2. Lower level slave problem

The slave problem in the lower level of the algorithm, **RNLP**, is obtained by tightening the search space of **ONLP**. This is accomplished by adding lower and upper bounds on the velocity terms  $v_r$ . The associated mathematical formulation is as follows:

$$\begin{array}{ll} \text{(RNLP)} & \min U(K_r, v_r, X_j) \\ & \text{s.t.} \quad \sum_{r=1}^p \mu_{ir} v_r = 0 \quad i = 1, \dots, n \\ & v_r = K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \quad r = 1, \dots, p \\ & \underline{v}_r \leq v_r \leq \bar{v}_r \quad r = 1, \dots, p \\ & K_r, v_r, X_j \in \mathbb{R}_+ \end{array}$$

Hence, the search space of **RNLP** is tighter than that of **ONLP**. For this reason, **RNLP** provides an upper bound on the solution of **ONLP**. Note that in model **RNLP**, bounds on  $v_r$  (third group of constraints) can be obtained from the active intervals of the disjunctions of **CMILP**. For instance, let  $v_r^*$  be the solution of the master problem. We know that  $v_r^*$  must fall within the active interval of the term of the disjunction defined by  $[v_r^h, v_r^{h+1}]$ . Hence, we can set  $\underline{v}_r = v_r^h$  and  $\bar{v}_r = v_r^{h+1}$ .

## 4.3. Algorithm steps

The detailed algorithmic steps of the proposed strategy are as follows:

- (1) Set iteration count  $it = 0$ ,  $UB = \infty$ ,  $LB = -\infty$  and tolerance error =  $tol$ .
- (2) Set  $it = it + 1$ . Solve master problem **CMILP**.
  - (a) If **CMILP** is infeasible, stop. **ONLP** is infeasible.
  - (b) Otherwise, update the current  $LB$  as  $LB = \max_{it}(LB_{it})$ , where  $LB_{it}$  is the value of the objective function of **CMILP** in the  $it^{th}$  iteration. Set bounds on  $v_r$  for the slave problem according to the solution of the master problem ( $\underline{v}_r = v_r^h$  and  $\bar{v}_r = v_r^{h+1}$ ).
- (3) Solve the slave problem **RNLP**.
  - (a) If **RNLP** is infeasible update the grid (see remark 5) and go to step 2 of the algorithm.
  - (b) Otherwise, update the current  $UB$  as  $UB = \min_{it}(UB_{it})$ , where  $UB_{it}$  is the value of the objective function of **RNLP** in the  $it^{th}$  iteration.
- (4) Calculate the optimality gap  $OG$  as  $OG = (|UB - LB|)/UB$ .
  - (a) If  $OG \leq tol$ , then stop. The current  $UB$  can be regarded as the global optimal solution of **ONLP** within the predefined tolerance.
  - (b) Otherwise, update the grid and go to step 2 of the algorithm.

## 4.4. Remarks

- The reformulation of Eq. (6) into two inequalities is only required for those velocities that are involved in balances at branching points, that is, where Eq. (3) includes more than two terms. In equations with only two terms, the logarithmic transformation is enough to obtain a linear constraint (note that the stoichiometric coefficients  $\mu_{ir}$  are known). Hence, in mathematical terms, we have:

$$\mu_{ir} v_r = -\mu_{ir'} v_{r'} \quad i \in XT \quad r, r' \in VT_i \quad (14)$$

$$\ln \mu_{ir} + \ln v_r = \ln (-\mu_{ir'}) + \ln v_{r'} \quad i \in XT \quad r, r' \in VT_i \quad (15)$$

$$\ln \mu_{ir} + k_r + \ln \gamma_r + \sum_{j=1}^{n+m} f_{rj} x_j$$

$$= \ln (-\mu_{ir'}) + k_{r'} + \ln \gamma_{r'} + \sum_{j=1}^{n+m} f_{r'j} x_j \quad i \in XT \quad r, r' \in VT_i \quad (16)$$

where  $XT$  is the set of equations involving only two terms and  $VT_i$  is the set of velocities that appear in those equations in  $XT$ . Note that in S-System models, all the balances include only two terms. This allows reformulating the model into a linear equivalent form, which greatly helps computations (Voit, 1992). Another major advantage of the logarithmic transformation is that it gives rise to concave univariate terms (i.e., logarithmic functions) for which tight under and over estimators can be defined.

- Supporting hyper-planes can be located following different patterns. It can be shown that the one that minimizes the rectilinear distance (i.e.,  $L_1$  norm) between the hyper-planes and the actual logarithmic function is that in which this distance is the same at every interjection of two adjacent hyper-planes (see proof in Appendix A). This allocation can be obtained by solving an optimization problem.
- Similarly, the grid points of the piecewise approximation can be selected according to different criteria. One possible strategy consists of splitting the range  $[v_r, \bar{v}_r]$  into  $H$  intervals with the same width. It can be shown that in order to minimize the rectilinear distance (i.e.,  $L_1$  norm) between the piecewise approximation and the logarithmic function, one needs to define intervals of equal width in the logarithmic space (see proof in Appendix A). Hence, we would have:

$$\ln v_r^{h+1} - \ln v_r^h = \ln v_r^{h+2} - \ln v_r^{h+1} = \dots = \ln v_r^{H+1} - \ln v_r^H$$

$$r = 1, \dots, p \quad h = 1, \dots, H \quad (17)$$

- Increasing the number of terms of the piecewise function and supporting hyper-planes leads to tighter bounds and hence to less iterations. Unfortunately, this is accomplished at the expense of adding more variables to the original problem. This is specially critical in the case of the piecewise approximation, which requires the definition of binary variables that increase considerably the computational burden of the master problem and consequently the time required by each iteration. Hence, a compromise should be found between the number of iterations and the time spent in each of them.
- There are different ways to update the piecewise grid of **CMILP** (steps 3a and 4b in the algorithm). One possible strategy is the division of the active interval into 2 sub-intervals with the same width either in the Cartesian space,  $(v_r^h + v_r^{h+1})/2$ , (see Fig. 3) or in the logarithmic space,  $(\ln v_r^h + \ln v_r^{h+1})/2$ . Another possibility is to split the active interval by adding the point corresponding to the solution of **RNLP** in the last iteration.
- Additional supporting hyper-planes can be iteratively added to **CMILP** in order to improve the overestimation of the logarithmic function. Again, the points where the new supporting hyper-planes will be allocated can be selected following different criteria.

### 5. Case studies

As benchmark problems to test the capabilities of the approach presented, we propose to use the ethanol production in the fermentation of *Saccharomyces cerevisiae* (case study 1) and the citric acid production by *Aspergillus niger* (case study 2) (see Figs. 4 and 5<sup>1</sup>).

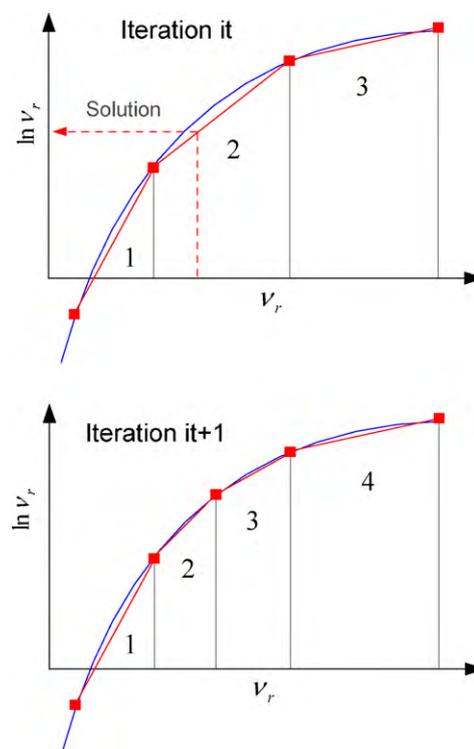


Fig. 3. Piecewise grid update example. As the solution of the first iteration is found in the second interval, it is split into two sub-intervals for the next iteration.

These two problems are convenient since their optimal solutions have already been published in the literature (Polisetty et al., 2008), the GMA models for the two systems can also be found in the same reference).

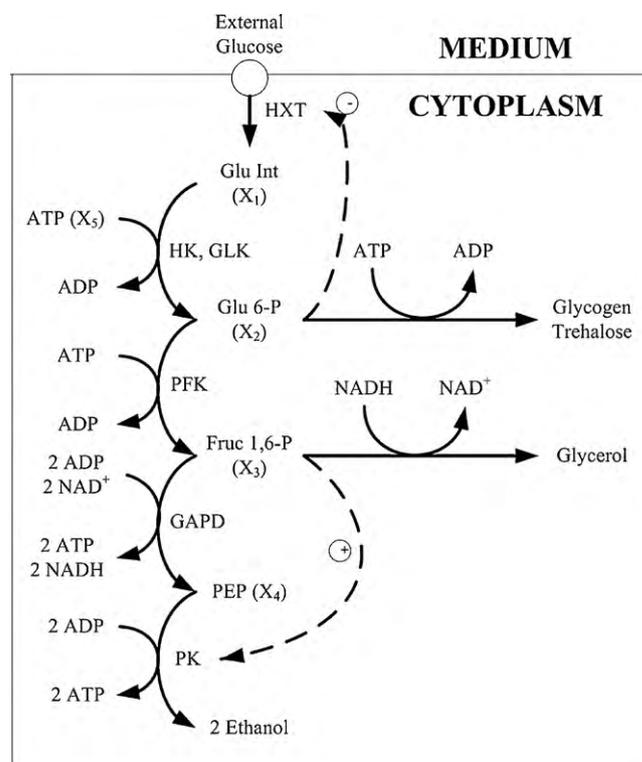


Fig. 4. Metabolic pathway of the fermentation of *Saccharomyces cerevisiae*.

<sup>1</sup> Figures adapted from the original work by Polisetty et al. (2008).



**Table 2**

Results of the global optimization of the ethanol production in *Saccharomyces cerevisiae* (GMA models from Polisetty et al., 2008). Gap: optimality gap.

	Polisetty et al. <sup>a</sup>	BARON	Proposed algorithm
Synthesis rate of ethanol (mM min <sup>-1</sup> )	157.59	157.59	157.59
UB	Not available	–	157.88
LB	157.59	–	157.59
Gap (%)	Not available	0.20	0.18
Iterations	–	–	3
Time (CPU s)	Not available	0.17	0.37

<sup>a</sup> Data termed as “Not available” is not shown in original work by Polisetty et al. (2008).

**Table 3**

Enzyme activities and metabolite concentrations (mM) in the global optimum for the ethanol synthesis rate in *Saccharomyces cerevisiae*.

$i, r$	1	2	3	4	5	6	7	8
$K_r$	5.00	0.89	5.00	0.20	1.25	0.20	5.00	5.00
$X_i$ (mM)	0.35	1.06	91.44	0.01	1.25	–	–	–

### 5.1. Ethanol production in *S. cerevisiae*

This case study was solved using a tolerance (*tol*) of 0.20% and considering that all the enzymes of the network are subject to modification. For comparison purposes, we solved the same problem with the standard global optimization package BARON. In order to provide the solver with a feasible starting point, the basal state solution was used. Note that this point can be easily computed before the optimization takes place by fixing all the  $K_r$  to 1 in the original model and solving the resulting system of nonlinear equations.

As it can be observed in Table 2, the results produced by our algorithm and BARON are in consonance with those reported in the literature by Polisetty et al. (2008): 157.59 mM/min (see Table 3 for the enzyme activities and metabolite concentrations in the solution). This is indeed a problem of small size (see Table 1 for details) for which both algorithms are able to identify the global optimal solution in few seconds of CPU time.

In order to further illustrate the capabilities of our algorithm, we have reproduced (Table 4) some of the results reported in Polisetty et al. (2008) where only a set of reactions are allowed to be modified, whereas the remaining enzyme activities are constrained to their

basal state ( $K_r = 1$ ). These calculations provide valuable information as the implementation of solutions requiring a large number of genetic manipulations might be impractical due to their elevated cost and complexity. Again, we have chosen a tolerance of 0.20% for both, our algorithm and BARON.

As observed in Table 4, the three methods were capable of determining the global optimal solution in a similar CPU time for the 8 combinations of free reactions. BARON showed to be slightly faster than the other two algorithms. With regard to the quality of the solutions found, it is interesting to notice that the method proposed by Polisetty et al. (2008) provides very loose optimality gaps. Particularly, although the method finds the global optimum in all the cases, the reported optimality gaps are very large (i.e., more than 40%). This constitutes a major limitation of this strategy. Interestingly, we identified two cases (5 and 7) where the same values of the objective function were obtained through three different enzyme activities combinations. These results suggest that the problem possess a certain degree of degeneracy. This issue should be carefully studied before attempting to reproduce any of these solutions in the laboratory, as there might be some particular features not considered in the analysis that would make the imple-

**Table 4**

Results of the global optimization of the ethanol production in *Saccharomyces cerevisiae* when fixing all the enzyme activities, but two, to their basal state. LB: lower bound in mM min<sup>-1</sup> (solution of ONLP). Gap: optimality gap.

Case	Modified reactions $r$	Polisetty et al.				BARON				Proposed method			
		$[K_r]$	LB	Gap (%)	Time <sup>a</sup> (CPU s)	$[K_r]$	LB	Gap (%)	Time (CPU s)	$[K_r]$	LB	Gap (%)	Time (CPU s)
1	[1, 3]	[5.00, 2.85]	103.66	21.66	0.81	[5.00, 2.85]	103.66	0.20	0.11	[5.00, 2.85]	103.66	0.09	0.94
2	[1, 4]	[5.00, 5.00]	73.18	48.96	0.26	[5.00, 5.00]	73.18	0.20	0.22	[5.00, 5.00]	73.18	0.16	2.03
3	[1, 7]	[5.00, 5.00]	73.41	47.46	0.20	[5.00, 5.00]	73.41	0.20	0.16	[5.00, 5.00]	73.41	0.11	2.17
4	[1, 6]	[5.00, 0.20]	73.41	47.15	0.24	[5.00, 0.20]	73.41	0.20	0.12	[5.00, 0.20]	73.41	0.11	2.72
5	[1, 5]	[5.00, 1.65]	72.68	48.47	0.22	[5.00, 0.63]	72.68	0.20	0.14	[5.00, 1.00]	72.68	0.11	2.63
6	[1, 8]	[5.00, 5.00]	87.77	22.13	0.19	[5.00, 5.00]	87.77	0.20	0.12	[5.00, 5.00]	87.77	0.14	2.59
7	[1, 2]	[5.00, 1.97]	72.68	47.48	0.24	[5.00, 5.00]	72.68	0.20	0.16	[5.00, 1.00]	72.68	0.11	2.49
8	[3, 7]	[5.00, 5.00]	44.67	76.18	0.09	[5.00, 5.00]	44.67	0.20	0.2	[5.00, 5.00]	44.67	0.08	1.82

<sup>a</sup> The author only reported the CPU time of the master MILP.

**Table 5**

Results of the global optimization of the citric acid production in *Aspergillus niger* (GMA models from Polisetty et al., 2008). Gap: optimality gap.

	Polisetty et al.	BARON	Proposed algorithm
Synthesis rate of citric acid (mM min <sup>-1</sup> )	384.23	Failed <sup>a</sup>	384.22
UB	384.23	–	390.66
LB	384.23	–	384.22
Gap (%)	0.00	–	1.68
Iterations	–	–	4
Time (CPU s)	5.68 <sup>b</sup>	24,000	33.37

<sup>a</sup> Failed means that the optimality gap was higher than 100%.

<sup>b</sup> The author only reported the CPU time of the master MILP.

**Table 6**  
 Metabolite concentrations (mM) in the global optimum for the citric acid synthesis rate in *Aspergillus Niger*.

<i>i</i>	$X_i$	<i>i</i>	$X_i$	<i>i</i>	$X_i$	<i>i</i>	$X_i$	<i>i</i>	$X_i$
1	0.02	7	0.20	13	5.60	19	0.04	25	1.41
2	1.00	8	$1.00 \times 10^{-4}$	14	0.01	20	0.85	26	0.98
3	0.12	9	21.21	15	31.26	21	0.71	27	0.30
4	0.02	10	130.00	16	0.52	22	0.35	28	0.26
5	0.71	11	0.01	17	1.70	23	0.01	29	$6.00 \times 10^{-8}$
6	0.01	12	0.01	18	22.55	24	0.01	30	0.21

**Table 7**  
 Enzyme activities in the global optimum for the citric acid synthesis rate in *Aspergillus Niger*.

<i>r</i>	$K_r$								
1	4.16	13	0.20	25	2.60	37	0.20	49	5.00
2	1.00	14	0.20	26	0.20	38	2.57	50	0.20
3	0.20	15	5.00	27	2.46	39	3.23	51	2.07
4	5.00	16	0.20	28	4.38	40	5.00	52	5.00
5	0.26	17	3.91	29	5.00	41	2.69	53	0.20
6	1.47	18	2.60	30	0.20	42	5.00	54	2.20
7	0.44	19	5.00	31	5.00	43	0.20	55	2.55
8	4.99	20	2.40	32	3.96	44	0.20	56	0.46
9	5.00	21	1.03	33	0.20	45	5.00	57	0.20
10	5.00	22	1.00	34	5.00	46	1.00	58	0.20
11	5.00	23	1.72	35	5.00	47	0.20	59	5.00
12	2.60	24	2.72	36	0.20	48	0.20	60	5.00

mentation of one of them advantageous when compared to the others.

## 5.2. Citric acid production in *A. niger*

The procedure explained for the first case study has been applied to solve the second case study with the only change of using a tolerance of 2.00%. The results obtained in the optimization are reported in Table 5 (the optimal values of the metabolite concentrations and the enzyme activities can be found in Tables 6 and 7, respectively).

This second case study considers a more complex network (4235 equations and 471 integer variables in the master problem vs 518 equations and 53 variables in the ethanol case). In this case, BARON failed at reducing the optimality gap below the specified tolerance (i.e., 2.00%) after 24,000 s of CPU time, whereas our algorithm was able to identify a solution in a relatively low computational time (i.e., less than 35 CPU seconds). In fact, after the aforementioned CPU time, BARON could only attain an optimality gap above 100%, which is very far away from the desired tolerance.

**Table 8**  
 Results of the global optimization of the citric acid production in *Aspergillus niger* when fixing some enzyme activities to their basal state. The number of reactions allowed to be modified depends on the case: Case B: one reaction; Case C: two reactions; Case D: three reactions; Case E: five reactions. LB: Lower bound in  $\text{mM min}^{-1}$  (solution of ONLP). Gap: optimality gap.

Case	Modified reactions <i>r</i>	Polisetty et al.			BARON Results	Proposed method				
		$[K_r]$	LB	Gap (%)		$[K_r]$	LB	Gap (%)	Time (CPU s)	
B	[40]	[5.00]	25.82	1234.12	9.11	Failed <sup>b</sup>	[5.00]	25.82	1.97	16.69
B	[59]	[1.00]	12.35	871.17	30.13	Failed	[1.00]	12.35	1.33	45.15
C	[40, 59]	[5.00, 1.00]	25.78	1254.54	13.27	Failed	[5.00, 1.00]	25.78	1.41	52.53
C	[1, 40]	[1.00, 5.00]	25.82	1241.75	26.4	Failed	[1.00, 5.00]	25.82	1.97	17.93
D	[1, 40, 60]	[1.27, 5.00, 1.12]	40.88	765.46	30.49	Failed	[1.27, 5.00, 1.12]	40.88	1.33	167.35
D	[1, 40, 59]	[1.16, 5.00, 5.00]	176.8	98.63	9.75	Failed	[1.16, 5.00, 5.00]	176.79	1.82	18.07
E	[1, 39, 40, 59, 60]	[1.24, 0.88, 5.00, 5.00, 1.01]	347.32	3.23	231.97	Failed	[1.40, 0.92, 5.00, 5.00, 1.07]	347.93	1.56	6.48
E	[1, 28, 40, 59, 60]	[1.46, 1.01, 5.00, 5.00, 1.11]	256.59	38.81	46.22	Failed	[1.46, 1.01, 5.00, 5.00, 1.11]	256.59	1.84	1093.09

<sup>a</sup> The author only reported the CPU time of the master MILP.

<sup>b</sup> Failed means that the optimality gap was higher than 100%.

**Table 9**  
 Local optimal solutions obtained by solving RNLP from different starting points.

Case	1	2	3	4	5
RNLP solution ( $\text{mM min}^{-1}$ )	354.87	379.75	384.22	372.86	384.21

of variables and constraints and also by the quality of the relaxation (i.e., the difference between the lower bound obtained in the slave problem and the upper bound predicted by the master problem), there are other facts that can affect it. For instance, the way in which the branch and bound is implemented to solve the MILPs (i.e., branching rules, order in which the nodes are explored, derivation of cutting planes, etc.) can have a major influence on the total CPU time.

Finally, it is interesting to notice that during the calculations we confirmed the existence of multiple local optimal solutions in the *A. niger* model. For that, we solved the original non-convex **ONLP** with a local optimizer (i.e., CONOPT) using five different starting points that were calculated by solving different master problems **CMILP**, each of them with a different initial number of piecewise terms (from 1 to 5). The results obtained, which are given in Table 9, show that different local optima can be obtained depending on the starting point used in the initialization of the algorithm. This observation justifies the use of global optimization tools to avoid falling in local optima during the search (see Table 5 for the global optimum obtained).

## 6. Conclusions

This paper has addressed the development of a systematic framework for the global optimization of metabolic networks that can be described by the Generalized Mass Action model. The strategy proposed is based on reformulating the original GMA model via a logarithmic transformation, which gives rise to a non-convex NLP. This model is globally optimized by an outer approximation algorithm that exploits its specific structure.

The capabilities of the proposed method have been illustrated by globally optimizing the fermentation pathway of *S. cerevisiae* and the metabolic network associated with the citric acid production in *A. niger*. For both cases, we have obtained the appropriate changes that need to be performed in the corresponding enzyme activities in order to maximize the production of ethanol and citric acid, respectively. Our algorithm has been able to reproduce the results previously reported by Polisetty et al. (2008), but achieving significant improvements in the optimality gaps of the final solutions. Besides, the method proposed has shown promising results even when applied to a moderately complex network (case study 2), absolutely surpassing the performance of BARON, which failed to solve that particular example within the predefined tolerance.

The generality of the optimization framework introduced in this paper makes it very interesting for biotechnological applications. At this point, the major drawbacks for getting practical results are: (1) the ability of obtaining appropriate models; and (2) the possibility of effectively manipulating the required enzymes. The main limitation for obtaining good mathematical models is the lack of experimental data that can be used for parameter estimation (Chou & Voit, 2009). Unfortunately, most of the Systems Biology effort has focused on gene sequences, protein structures, and so on, with relatively few results on actual data on intact systems. The kind of data required for this task would involve measuring metabolites and fluxes in vivo, a problem that is not totally solved yet. The optimization method presented here can yield valuable insights if and only if the underlying model is a good approximation to reality. Besides this problem, optimization results require experimental confirmation; that is manipulation of enzymes for obtaining the desired optimal increment on the objective function. However, although gene expression changes can be easily introduced in living cells, there is no guaranty that an appropriate change in enzyme activity is also obtained.

In conclusion, our results show that it is possible to appropriately analyze a highly non-linear mathematical model and obtain optimal solution for a given objective function. This should

stimulate experimentalists for developing appropriate tools for measuring living cells and for manipulating them so that practical results can be obtained.

## Acknowledgements

The authors wish to acknowledge support of this research work from the Spanish Ministry of Education and Science (projects BFU2005-00234/BMC, BFU2008-00196/BMC, DPI2008-04099, PHB2008-0090-PC, BFU2008-00196 and CTQ2009-14420-C02), the Spanish Ministry of External Affairs (projects HS2007-0006, A/020104/08 and A/023551/09), and the Generalitat de Catalunya (FI programs).

## Appendix A.

Lemma 1 shows that the maximum error between the linear outer approximation and the logarithmic function lies in a vertex. Proposition 1 uses the results of Lemma 1 to show that the allocation of hyper-planes that minimizes the rectilinear distance (i.e.,  $L_1$  norm) between the hyper-planes and the actual logarithmic function is that in which this distance is the same at every intersection of two adjacent hyper-planes. Lemma 2 and Proposition 2 are similar to Lemma 1 and Proposition 1 but apply to the piece-wise approximation. Finally, Proposition 3 complements Proposition 2 and provides a direct way of defining a piece-wise approximation with minimum error.

**Lemma 1.** Consider an outer approximation of the function  $\ln v_r$  with  $L$  supporting hyper-planes (see Fig. 6). The maximum error,  $error_{max}$ , (defined as the linear distance,  $L_1$  norm), between the hyper-planes and the logarithmic function is attained in a vertex.

**Proof.** We first show that the point with the maximum error lies in a hyperplane, and then prove that it must correspond to one of its intersections with adjacent hyperplanes. Consider problem **PA**, which seeks the maximum difference between a set of hyper-planes and the logarithmic function:

$$\begin{aligned}
 \text{(PA)} \quad & \min \quad \ln v_r - y \\
 \text{s.t.} \quad & y - \left( \ln v_r^l + \frac{1}{v_r^l} (v_r - v_r^l) \right) \leq 0 \quad l = 1, \dots, L \\
 & v_r - \bar{v}_r \leq 0 \\
 & \underline{v}_r - v_r \leq 0 \\
 & y \in \mathbb{R} \\
 & v_r \in \mathbb{R}_+
 \end{aligned}$$

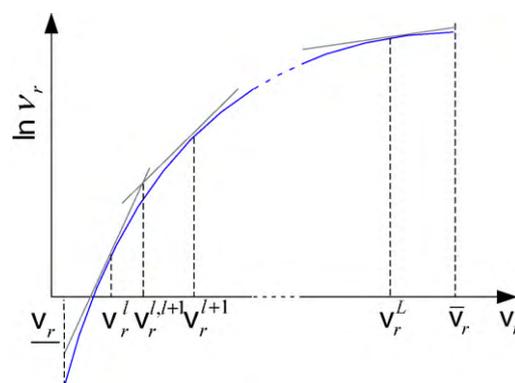


Fig. 6. Approximation of the  $\ln v_r$  function by  $L$  supporting hyper-planes.

where  $\underline{v}_r \leq v_r^l \leq \bar{v}_r$ . The Karush-Kuhn-Tucker (KKT) conditions of **PA** are:

$$-1 + \sum_{l=1}^L u_1^l = 0 \quad (18)$$

$$\frac{1}{v_r} - \sum_{l=1}^L \frac{u_1^l}{v_r^l} + u_2 - u_3 = 0 \quad (19)$$

$$u_1^l \left( y - \left( \ln v_r^l + \frac{1}{v_r^l} (v_r - v_r^l) \right) \right) = 0 \quad (20)$$

$$u_2(v_r - \bar{v}_r) = 0 \quad (21)$$

$$u_3(\underline{v}_r - v_r) = 0 \quad (22)$$

$$u_1^l \geq 0 \quad (23)$$

$$u_2 \geq 0 \quad (24)$$

$$u_3 \geq 0 \quad (25)$$

From Eq. (18), it follows that at least one supporting hyper-plane (*SH*) must be active in the optimal solution. Now, consider problem **PB** that seeks the maximum error along the active *SH<sub>l</sub>* between its extremes  $v_r^{l,o}$  and  $v_r^{l,p}$ , which are given by the intersection of the hyper-plane with either an adjacent *SH<sub>j</sub>* or a limit of the interval  $[\underline{v}_r, \bar{v}_r]$ .

$$\begin{aligned} \text{(PB)} \quad \min \quad & \ln v_r - \left( \ln v_r^l + \frac{1}{v_r^l} (v_r - v_r^l) \right) \\ \text{s.t.} \quad & v_r - v_r^{l,p} \leq 0 \\ & v_r^{l,o} - v_r \leq 0 \\ & v_r \in \mathbb{R}_+ \end{aligned}$$

The KKT conditions of **PB** are:

$$\frac{1}{v_r} - \frac{1}{v_r^l} + u_1 - u_2 = 0 \quad (26)$$

$$u_1(v_r - v_r^{l,p}) = 0 \quad (27)$$

$$u_2(v_r^{l,o} - v_r) = 0 \quad (28)$$

$$u_1 \geq 0 \quad (29)$$

$$u_2 \geq 0 \quad (30)$$

There are 3 possible solutions to this problem.

**Case 1:**  $u_1 = u_2 = 0$ . From Eq. (26), we have:

$$v_r^* = v_r^l$$

And the resulting value of the objective function is:

$$OF = 0$$

It is easy to see that this point is a maximum of **PB** in which the error is minimum. Note that this is the point in which the hyper-plane touches the logarithmic function.

**Case 2:**  $u_1 = 0, u_2 \neq 0$ . From Eqs. (26) and (28), we get:

$$v_r^* = v_r^{l,o}; u_2 = \frac{1}{v_r^{l,o}} - \frac{1}{v_r^l} \geq 0 \quad OF = \ln \left( \frac{v_r^{l,o}}{v_r^l} \right) - \frac{v_r^{l,o}}{v_r^l} + 1$$

Hence, this point is a minimum of **PB** and corresponds to a vertex.

**Case 3:**  $u_2 = 0, u_1 \neq 0$ . From Eqs. (26) and (27), we get:

$$v_r^* = v_r^{l,p}; u_1 = \frac{1}{v_r^l} - \frac{1}{v_r^{l,p}} \geq 0 \quad OF = \ln \left( \frac{v_r^{l,p}}{v_r^l} \right) - \frac{v_r^{l,p}}{v_r^l} + 1$$

This point (again a vertex) is another minimum of **PB**. Hence, the solution of **PB** must correspond to a vertex, and the proof is complete.  $\square$

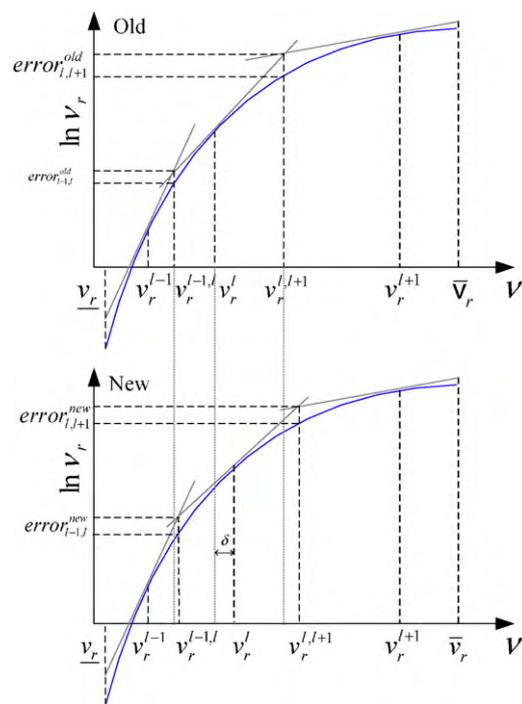


Fig. 7. Illustration of the decrease in  $error_{max}$  by moving *SH<sub>l</sub>* a distance  $\delta$  towards the vertex  $v_r^{l,l+1}$ .

**Proposition 1.** The allocation of *L* hyper-planes that minimizes  $error_{max}$  is that in which the error is the same in all the *L*+1 vertexes.

**Proof.** The proof is by contradiction. From Lemma 1, we know that the maximum error between the hyper-planes and the logarithmic function is attained in a vertex. Assume that in the optimal allocation there is at least one vertex  $v_r^{l,l+1}$  with a different error. Now, identify a supporting hyperplane *SH<sub>l</sub>* with different errors in its extreme vertexes ( $error_{l-1,l}$  in  $v_r^{l-1,l}$  and  $error_{l,l+1}$  in  $v_r^{l,l+1}$ ). Assume, without losing generality, that  $error_{l,l+1} \geq error_{l-1,l}$ . Now, we consider two cases:

**Case 1:** the maximum  $error_{max} = \max_{l \neq l'} \{error_{l,l'}\}$  corresponds to the right vertex  $v_r^{l,l+1}$  of the selected hyperplane, that is,  $error_{max} = error_{l,l+1}$ . Now, define  $error_{l,l+1}^{old} = error_{l,l+1}^{old}$  and move the hyperplane *SH<sub>l</sub>* a small distance  $\delta$  towards  $v_r^{l,l+1}$ , that is, make  $v_r^{l,l+1,new} = v_r^{l,l+1,old} + \delta$ , thus decreasing the slope of *SH<sub>l</sub>*. This move decreases  $error_{l,l+1}$  at the expense of increasing  $error_{l-1,l}$ . Since the logarithmic function is continuous, it is possible to find  $\delta$  such that  $error_{l-1,l}^{old} < error_{l-1,l}^{new} = error_{l,l+1}^{new} < error_{l,l+1}^{old}$  (Fig. 7), that is, a new solution with a smaller error in the right vertex of *SH<sub>l</sub>*, and hence with a smaller  $error_{max}$ . This contradicts the fact that in the optimal solution there are vertexes with different errors.

**Case 2:**  $error_{max}$  is placed in another hyper-plane *SH<sub>l'</sub>* ( $l' \neq l, l+1$ ). In this case, the same procedure can be repeated recursively to the rest of the hyper-planes until no more hyper-planes containing different errors in their vertexes remain. It is straightforward to see that this would lead to a solution with lower  $error_{max}$ , which contradicts the assumption that the optimal allocation implies the existence of at least one hyperplane with different errors in its extreme vertexes.  $\square$

**Lemma 2.** Consider an underestimation of the  $\ln v_r$  function with a linear piecewise (PW) section (Fig. 8). The maximum error,  $error_{max}$ , defined as the  $L_1$  norm (i.e.,  $error(v_r) = \ln v_r - av_r - b$ ) between the  $\ln v_r$  and the PW linear function occurs at  $v_r^* = 1/a$

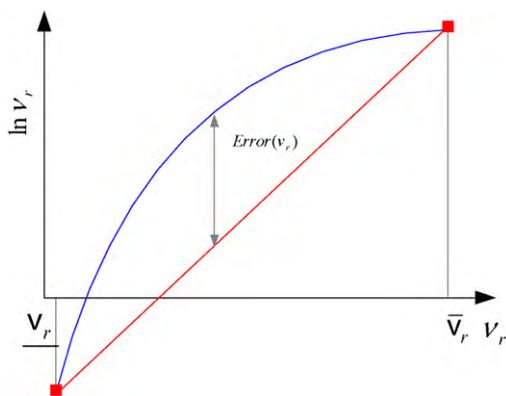


Fig. 8. Approximation of the  $\ln v_r$  function by a one interval linear function.

**Proof.** The error is a concave function that depends on a single variable, hence in the optimal solution we get:

$$error' = \frac{1}{v_r} - a = 0 \Leftrightarrow v_r^* = \frac{1}{a} \quad error'' = -\frac{1}{v_r^2} \leq 0$$

Therefore,  $v_r^* = \frac{1}{a}$  is a maximum of the function error.  $\square$

**Proposition 2.** Consider a piece-wise approximation of the function  $\ln v_r$  with  $H$  intervals. Let  $error_h$  be the maximum error in the  $h^{th}$  interval of the PW function and  $error_{max}$  the maximum error among the different intervals ( $error_{max} = \max_h \{error_h\}$ ). The piece-wise approximation that minimizes  $error_{max}$  is that in which  $error_h = error_{h'}$ ,  $\forall h, h' (h \neq h')$ .

**Proof.** The proof is by contradiction. From Lemma 2 we know that the maximum error in a piecewise section  $PW_h$  is attained at  $1/a$ . Assume that the optimal distribution of the domain is that where there is at least one section  $h$  with a different error,  $error_h > error_{h+1}$ . Consider the following cases:

**Case 1:**  $error_h = error_{max}$ . Move the grid point  $\bar{v}_r^h = v_r^{h+1}$  a distance  $\delta$  towards  $\bar{v}_r^h$ , that is, make  $\bar{v}_r^{h, new} = \bar{v}_r^{h, old} - \delta$ , thus increasing  $a_h$  and decreasing  $\bar{a}_{h+1}$ . This decreases  $error_h$  and increases  $error_{h+1}$ . Since the logarithmic function is continuous, we can define a  $\delta$  such that  $error_{h+1}^{old} < error_{h+1}^{new} = error_h^{new} < error_h^{old}$  (Fig. 9), that is, a new solution with a smaller  $error_{max}$ . This contradicts the original statement that in the optimal solution there are sections with different errors.

**Case 2:**  $error_{max} = error_{h'} \neq error_h (h' \neq h)$ . It is straightforward to see that we can follow the same strategy described before until the error in every couple of adjacent sections is the same and smaller than  $error_{max}^{old}$ .  $\square$

**Proposition 3.** The distribution of the  $H$  intervals where  $error_h = error_{h'}, \forall h, h' (h \neq h')$  corresponds to that where all PW sections are of equal width  $Q$  in the logarithmic space.

**Proof.** Consider two piecewise sections  $^2PW_h$  and  $PW_{h+1}$  defined by grid point  $v_r^{h+1}$ . When  $error_h = error_{h+1}$  the following relationship holds:

$$\ln\left(\frac{1}{a_h}\right) - 1 - b_h = \ln\left(\frac{1}{a_{h+1}}\right) - 1 - b_{h+1} \quad (31)$$

$$\rightarrow \ln\left(\frac{a_{h+1}}{a_h}\right) - (b_h - b_{h+1}) = 0$$

As the piecewise functions take the same value at the common grid point  $v_r^h = v_r^{h+1}$ , we can rewrite Eq. (31) in terms of  $a_h, a_{h+1}$  and

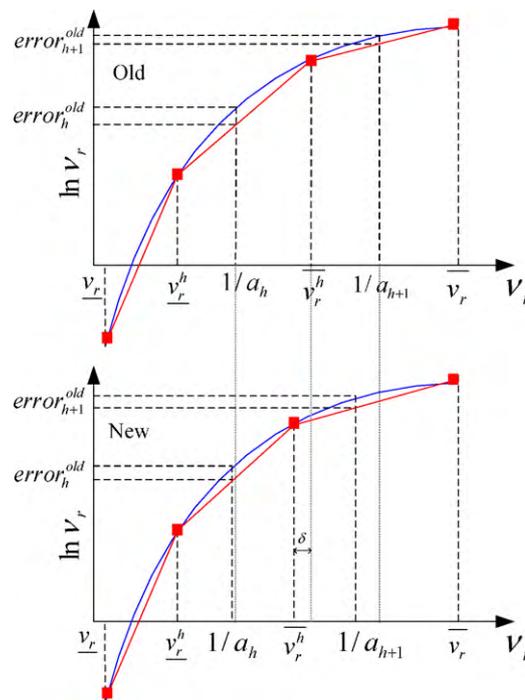


Fig. 9. Illustration of the decrease in  $error_{max}$  by moving the grid point  $\bar{v}_r^h$  a distance  $\delta$ .

$$\bar{v}_r^{h+1}:$$

$$\ln\left(\frac{a_{h+1}}{a_h}\right) - v_{r-h+1}(a_{h+1} - a_h) = 0 \quad (32)$$

Now, we introduce two new variables  $Q$  and  $Q'$  defined as:

$$Q = \ln \bar{v}_r^h - \ln \underline{v}_r^h \quad (33)$$

$$Q' = \ln \bar{v}_r^{h+1} - \ln \underline{v}_r^{h+1} \quad (34)$$

Hence, we can express the width of each interval in the cartesian space (i.e.,  $\bar{v}_r^h - \underline{v}_r^h$  and  $\bar{v}_r^{h+1} - \underline{v}_r^{h+1}$ ) in terms of  $Q$  and  $Q'$ :

$$\underline{v}_r^h = \frac{\bar{v}_r^h}{\exp Q} \rightarrow \bar{v}_r^h - \underline{v}_r^h = \frac{\bar{v}_r^h}{\exp Q} (\exp Q - 1) \quad (35)$$

$$\underline{v}_r^{h+1} = \frac{\bar{v}_r^{h+1}}{\exp Q'} \rightarrow \bar{v}_r^{h+1} - \underline{v}_r^{h+1} = \frac{\bar{v}_r^{h+1}}{\exp Q'} (\exp Q' - 1) \quad (36)$$

Similarly, we can redefine the slope of each of the linear piecewise functions,  $a_h$  and  $a_{h+1}$ , in terms of  $Q, Q'$  and  $\bar{v}_r^{h+1}$ :

$$a_h = \frac{\ln \bar{v}_r^h - \ln \underline{v}_r^h}{\bar{v}_r^h - \underline{v}_r^h} = \frac{Q (\exp Q)}{\bar{v}_r^{h+1} (\exp Q - 1)} \quad (37)$$

$$a_{h+1} = \frac{\ln \bar{v}_r^{h+1} - \ln \underline{v}_r^{h+1}}{\bar{v}_r^{h+1} - \underline{v}_r^{h+1}} = \frac{Q'}{\bar{v}_r^{h+1} (\exp Q' - 1)} \quad (38)$$

By introducing Eqs. (37) and (38) into Eq. (32), the following equality is obtained:

$$\ln\left(\frac{Q' (\exp Q - 1)}{(\exp Q' - 1) Q (\exp Q)}\right) - \left(\frac{Q'}{(\exp Q' - 1)} - \frac{Q (\exp Q)}{(\exp Q - 1)}\right) = 0 \quad (39)$$

When  $Q' = Q$  (i.e., when the intervals are of equal width in the logarithmic space) this equation is satisfied.  $\square$

<sup>2</sup> Note that more intervals could be considered and generality would not be lost.

## References

- Alvarez-Vasquez, F., Canovas, M., Iborra, J., & Torres, N. (2002). Modeling, optimization and experimental assessment of continuous l-(-)-carnitine production by *Escherichia coli* cultures. *Biotechnology and Bioengineering*, 80(7), 794–805.
- Alves, R., Vilaprinyo, E., Hernandez-Bermejo, B., & Sorribas, A. (2009). Mathematical formalisms based on approximated kinetic representations for modeling genetic and metabolic pathways. *Biotechnology & Genetic Engineering Reviews*, 25, 1–40.
- Bailey, J. (1991). Toward a science of metabolic engineering. *Science*, 252, 1668–1675.
- Bailey, J. (1999). Lessons from metabolic engineering for functional genomics and drug discovery. *Nature Biotechnology*, 17, 616–618.
- Bailey, J., Birnbaum, S., Galazzo, J., Khosla, C., & Shanks, J. (1990). Strategies and challenges in metabolic engineering. *Annals of the New York Academy of Sciences*, 589, 1–15.
- Banga, J. (2008). Optimization in computational systems biology. *BMC Systems Biology*, 2–47.
- Bergamini, M., Aguirre, P., & Grossmann, I. (2005). Logic-based outer approximation for globally optimal synthesis of process networks. *Computers and Chemical Engineering*, 29, 1914–1933.
- Cameron, D., & Chaplen, F. (1997). Developments in metabolic engineering. *Current Opinion in Biotechnology*, 8, 175–180.
- Cameron, D., & Tong, J. (1993). Cellular and metabolic engineering: An overview. *Applied Biochemistry and Biotechnology*, 38, 105–140.
- Chang, Y., & Sahinidis, N. (2005). Optimization of metabolic pathways under stability considerations. *Computers and Chemical Engineering*, 29, 467–479.
- Chou, I., & Voit, E. (2009). Comparative characterization of the fermentation pathway of *Saccharomyces cerevisiae* using biochemical systems theory and metabolic control analysis: Model definition and nomenclature. *Mathematical Biosciences*, 219, 57–83.
- Curto, R., Sorribas, A., & Cascante, M. (1995). Comparative characterization of the fermentation pathway of *Saccharomyces cerevisiae* using biochemical systems theory and metabolic control analysis: Model definition and nomenclature. *Mathematical Biosciences*, 130, 25–50.
- Floudas, C., & Gounaris, C. (2009). A review of recent advances in global optimization. *Journal of Global Optimization*, 45, 3–38.
- Gavalas, G. (1968). *Nonlinear differential equations of chemical reacting systems*. Berlin: Springer-Verlag.
- Grossmann, I., & Bigler, L. (2004). Part ii. future perspective on optimization. *Computers and Chemical Engineering*, 28, 1193–1218.
- Guillén-Gosálbez, G., & Sorribas, A. (2009). Identifying quantitative operation principles in metabolic pathways: A systematic method for searching feasible enzyme activity patterns leading to cellular adaptive responses. *BMC Bioinformatics*, 10(386).
- Hatzimanikatis, V., Floudas, C., & Bailey, J. (1996). Optimization of regulatory architectures in metabolic reaction networks. *Biotechnology and Bioengineering*, 52, 485–500.
- Heinrich, R., & Schuster, S. (1996). *The regulation of cellular systems*. New York: Chapman and Hall.
- Marin-Sanguino, A., & Torres, N. (2003). Optimization of biochemical systems by linear programming and general mass action model representations. *Mathematical Biosciences*, 184(2), 187–200.
- Marin-Sanguino, A., Voit, E., Gonzalez-Alzon, C., & Torres, N. (2007). Optimization of biotechnological systems through geometric programming. *Theoretical Biology & Medical Modelling*, 4–38.
- Mendes, P., & Kell, D. (1996). Making cells work – metabolic engineering for everyone. *Trends in Biotechnology*, 15, 6–7.
- Polisetty, P., Gatzke, E., & Voit, E. (2008). Yield optimization of regulated metabolic systems using deterministic branch-and-reduce methods. *Biotechnology and Bioengineering*, 99(5), 1154–1169.
- Savageau, M. (1969a). Biochemical systems analysis. i. Some mathematical properties of the rate law for the component enzymatic reactions. *Journal of Theoretical Biology*, 25, 365–369.
- Savageau, M. (1969b). Biochemical systems analysis. ii. The steady-state solutions for an n-pool system using a power-law approximation. *Journal of Theoretical Biology*, 25, 370–379.
- Sorribas, A., Hernandez-Bermejo, B., Vilaprinyo, E., & Alves, R. (2007). Cooperativity and saturation in biochemical networks: A saturable formalism using Taylor series approximations. *Biotechnology and Bioengineering*, 97, 1259–1277.
- Sorribas, A., Pozo, C., Vilaprinyo, E., Guillén-Gosálbez, G., Jiménez, L., & Alves, R. (2010). Global optimization techniques in Generalized Mass Action models. *J. Biotechnol.*, doi:10.1016/j.jbiotec.2010.01.026
- Torres, N., & Voit, E. (2002). *Pathway analysis and optimization in metabolic engineering*. Cambridge: Cambridge University Press.
- Vecchietti, A., Sangbum, L., & Grossmann, I. (2003). Modeling of discrete/continuous optimization problems: Characterization and formulation of disjunctions and their relaxations. *Computers and Chemical Engineering*, 27, 433–448.
- Vera, J., de Atauri, P., Cascante, M., & Torres, N. (2003). Multicriteria optimization of biochemical systems by linear programming: Application to production of ethanol by *Saccharomyces cerevisiae*. *Biotechnology and Bioengineering*, 83(3), 335–343.
- Voit, E. (1992). Optimization in integrated biochemical systems. *Biotechnology and Bioengineering*, 40(5), 572–582.
- Voit, E. (2000). *Computational analysis of biochemical systems. A practical guide for biochemists and molecular biologists*. Cambridge: Cambridge University Press.
- Voit, E. (2003). Design principles and operating principles: the yin and yang of optimal functioning. *Mathematical Biosciences*, 182, 81–92.

# A Spatial Branch-and-Bound Framework for the Global Optimization of Kinetic Models of Metabolic Networks

C. Pozo,<sup>†</sup> G. Guillén-Gosálbez,<sup>\*,†</sup> A. Sorribas,<sup>‡</sup> and L. Jiménez<sup>†</sup>

<sup>†</sup>Departament d'Enginyeria Química (EQ), Escola Tècnica Superior d'Enginyeria Química (ETSEQ), Universitat Rovira i Virgili (URV), Campus Sescelades, Avinguda Països Catalans, 26, 43007 Tarragona, Spain

<sup>‡</sup>Departament de Ciències Mèdiques Bàsiques, Institut de Recerca Biomèdica de Lleida (IRBLLEIDA), Universitat de Lleida, Montserrat Roig 2, 25008 Lleida, Spain

**ABSTRACT:** The identification of the enzymatic profile that achieves a maximal production rate of a given metabolite is an important problem in the biotechnological industry, especially if there is a limit on the number of enzymatic modulations allowed. The intrinsic nonlinear behavior of metabolic processes enforces the use of kinetic models, such as the generalized mass action (GMA) models, giving rise to nonconvex MINLP formulations with multiple local solutions. In this paper, we introduce a customized spatial branch-and-bound strategy devised to solve efficiently these particular problems to global optimality. A tight MILP-based relaxation of the original nonconvex MINLP is constructed by means of supporting hyperplanes and piecewise linear underestimators. The overall solution procedure is expedited through the use of bound tightening techniques and a special type of cutting plane. The capabilities of the proposed strategy are tested through its application to the maximization of the citric acid production in *Aspergillus niger*. We also provide a numerical comparison of our algorithm with the commercial package BARON and an outer approximation-based method earlier proposed by the authors.

## 1. INTRODUCTION

Cellular and molecular biology has experienced a dramatic paradigm switch driven by the introduction of new technological and computational tools. This change has led to a wide acceptance of networks and their emergent properties as a central subject for understanding the evolution of cell metabolism. The emergence of systems biology as a discipline based on high throughput experimental techniques, bioinformatics methods, and mathematical modeling is the result of these advances. One of the consequences of this activity is a renewed interest in biotechnological applications that ranges from industrial products based on modified organisms to the possibility of designing new organisms.<sup>1–4</sup>

Cellular metabolism is a complex system that involves a huge number of components interacting in a dynamic way through nonlinear processes. This makes biological systems much more challenging than human designed factories and industrial products. In most problems, appropriate simplifications are required to grasp part of this complexity and to obtain practical results both in understanding the evolution of emergent properties and in predicting systems responses to experimental manipulation.<sup>5–7</sup>

Advances in molecular biology techniques have made it possible to modulate the expression of genes in a given organism in order to obtain strains with enhanced phenotypes.<sup>8,9</sup> Being able to improve the yield through modified strains is a crucial aspect for successful biotechnological applications. However, the intrinsic complexity of metabolic networks makes an intuitive inference of the most promising genetic changes a highly difficult (if not impossible) task. Henceforth, systematic optimization tools are required for improving metabolic engineering so that biotechnological applications can be made useful and affordable.

Optimization is not at all a new concept in biology.<sup>10–13</sup> It is clear that mathematical programming approaches offer a

promising framework for analyzing mathematical models of biological systems in a systematic way, shedding light on the strategies that must be followed in order to improve their properties.<sup>9,12,14–16</sup> In particular, one of the areas in which systematic tools based on mathematical programming hold good promise is the analysis and manipulation of metabolic networks through gene expression modification.<sup>17–20</sup> From the point of view of industrial applications, the use of optimization methods in systems biology applications has gained wider interest.<sup>9</sup> Besides their application in increasing the yield of specific products, these techniques have also been used to explain the current adaptive responses of organisms and to predict the properties of new designs.<sup>9,21,22</sup>

While existing optimization techniques may be of some help, the complexity of cellular metabolism requires the development of global optimization methods that could be applied to these kinds of nonlinear problems. With these techniques, one expects that actual biological processes could be further improved by identifying quantitative operation principles that would help in deciding which genes should be modified and which is the optimal profile for obtaining a given goal. The fact that biological experiments are expensive and time-consuming<sup>9</sup> coupled with the usefulness of computational techniques when modeling metabolic networks<sup>23</sup> contributes to increasing the attractiveness of developing appropriate optimization approaches to address these problems.

**Special Issue:** Puigjaner Issue

**Received:** June 28, 2010

**Accepted:** October 19, 2010

**Revised:** October 11, 2010

**Published:** December 01, 2010

Flux balance analysis attempts this prediction by the optimization of stoichiometric models.<sup>24</sup> This approach leads to mixed-integer linear problems (MILP) that can be effectively solved by standard branch-and-bound techniques. This has been the main key of their success in different applications.<sup>25–30</sup> Unfortunately, this technique fails to capture the regulatory relationships that commonly exist between processes in metabolic networks.<sup>31</sup> These limitations can be overcome by resorting to kinetic metabolic models that account for the relationship between the concentration of metabolites and the fluxes in the network. Specifically, nonlinear kinetic expressions are preferred, as linear estimations have been found to be only valid for a narrow range around the approximation point.<sup>8</sup>

Among the available formalisms, models based on the power-law formalism in the variant form known as generalized mass action (GMA from here on) exhibit some particular advantages that make their application rather convenient.<sup>32–37</sup> For instance, as will be explained in detail later in the paper, they can adequately capture the nonlinear behavior of the metabolic regulations while exhibiting some linear properties when expressed in the logarithmic space. Furthermore, they are able to describe any particular metabolic network<sup>37</sup> what grants the generality of the framework presented herein. On the other hand, this approach gives rise to nonconvex models and, hence, to multimodality (i.e., existence of multiple solutions).<sup>9</sup> It should be emphasized that guaranteeing global optimality is of paramount importance in this type of problem, as a local optimal solution may lead to a completely different physical interpretation and objective function value than that associated with the global optimum, thus hampering the entire biological analysis.<sup>38</sup>

Global optimization addresses the computation and characterization of global optima (i.e., minima and maxima) of nonconvex functions constrained in a specified domain.<sup>39</sup> It has been the object of intense research during the past 15 years, but it is expected to continue as a major challenge in nonlinear optimization in the upcoming years.<sup>40</sup>

Global optimization approaches can be classified into stochastic or deterministic ones. Stochastic methods are nondeterministic approaches (i.e., they cannot guarantee global optimality) that make use of meta-heuristics in order to guide the search for “good” solutions from a series of pseudorandom generated points. These methods are often based on physical and biological analogies. On the other hand, deterministic methods are rigorous and, thus, can guarantee global optimality within a desired optimality gap. These methods rely on the calculation of a series of valid upper and lower bounds for the global optimum of the problem that approach each other during the execution of the algorithm until the optimality gap is reduced below a predefined tolerance. Among the different methods that may be included in this group, the most commonly used are the outer-approximation (OA)<sup>41</sup> and the spatial branch-and-bound (B&B) methods.<sup>42–46</sup>

In OA, the original problem is decomposed into two different subproblems at two different hierarchical levels: a master lower bounding problem and a slave upper bounding problem. The former is a relaxation of the original problem (i.e., it overestimates the feasible region of the original problem) that provides lower bounds on its global optimum. The latter entails the solution of the original problem in a reduced search space. In each iteration, the solution of the master problem is used as a starting point to solve locally the slave problem in a reduced search space (i.e., bounds are provided to some variables according to the solution of the master problem). If the

optimality gap is found to be within a given tolerance, the algorithm terminates. Otherwise, the relaxation of the master problem is improved (i.e., is tightened) at the expense of introducing more variables.

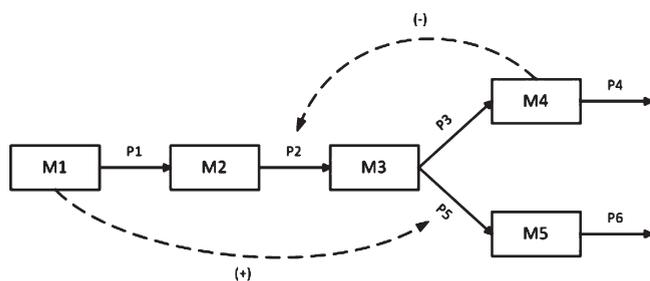
On the other hand, in the spatial branch-and-bound (sBB here on; do not confuse with the MINLP solver sBB that implements a nonlinear branch and bound) method, the original problem is allocated in the root node of an exploration tree. Lower and upper bounds for the problem are compared, and if the desired tolerance is not met, the problem is split into two smaller subproblems (descendants) by partitioning the feasible space of a continuous variable (branching variable). Then, the two new problems are solved, if required, by recursive partitioning. If a node is proved not to contain the global optimum, then the associated branch in the sBB tree can be pruned. At the end, the global optimal solution is to be found in one of the subproblems derived during the process. This method is based on the idea of “divide and conquer” as each of the subproblems is smaller, and thus easier to solve, than the original one.

Multiple methods have been devised so far as variations from the original sBB. These methods include branch-and-reduce,<sup>47,48</sup>  $\alpha$ BB,<sup>49–54</sup> symbolic reformulation,<sup>38,55,56</sup> reduced-space branch-and-bound,<sup>57</sup> branch-and-contract,<sup>58</sup> and the branch-and-cut framework proposed by Barton.<sup>59</sup> Some interval arithmetic global optimization methods<sup>60–62</sup> are sBB-like methods.<sup>63</sup> It has been observed that the performance of global optimization methods is highly dependent on the type of nonlinearities.<sup>8</sup> Henceforth, by exploiting the special mathematical structure of the problem under investigation,<sup>64,65</sup> it is possible to devise tighter relaxations that lead to faster algorithms.<sup>66</sup>

The application of global optimization methods to the analysis of metabolic networks that are described through nonlinear models (e.g., GMA formalism) has been scarce. Polisetty et al.<sup>67</sup> were the first ones to address this problem. In their work, they present a B&B procedure to identify the enzymes to be modified for efficiency in yield and cost. Later, Pozo et al.<sup>68</sup> proposed an outer-approximation algorithm that improved the method by Polisetty in terms of quality of the solutions provided (i.e., significantly smaller optimality gaps) and CPU time. The authors also presented a rigorous theoretical analysis on the construction of tight piecewise approximations and supporting hyperplanes. This method was also used to study the evolution of the cellular metabolism.<sup>21,22</sup>

In this work, we present a novel sBB method for the global optimization of metabolic networks that are modeled via the GMA formalism. Our computational procedure exploits the specific structure of the GMA models in order to construct tight MILP-based relaxations of the original nonconvex formulation. These linear relaxations are tightened through the use of a special type of cutting planes that are derived from some equations of the model. The sBB method is further expedited by tailored-made branching rules and bound contraction procedures based on interval analysis. The capabilities of this customized sBB are tested through a case study that addresses the optimization of citric acid production by *Aspegillus niger*. The results produced by our algorithm are compared with those generated by an outer approximation-based method introduced by the authors in previous works and also with the commercial global optimization package BARON.<sup>21,22,68</sup>

The paper is organized as follows. The problem is presented in section 2, and its mathematical formulation is proposed in section 3. The customized sBB is described in detail in section



**Figure 1.** Example of a generic metabolic network, where processes are represented by arrows and metabolites by boxes.

4, whereas section 5 contains some numerical results. Finally, in section 6, we discuss some particular issues about the performance of the proposed methodology and its implementation.

## 2. PROBLEM STATEMENT

A metabolic network (Figure 1) is composed of a set of reactions and transportation processes (represented by arrows in the figure), generally ruled by enzymes, which transform organic substrates into metabolic intermediates and energy compounds (i.e., metabolites, in general). Some of these metabolites (represented by boxes in the figure) can also inhibit or facilitate some processes in the metabolic network. For instance, M4 inhibits P2 and M1 facilitates P5 in the figure.

The problem under study is the determination of the levels of the enzymes activities that maximize the synthesis rate of a particular metabolite in a metabolic pathway. The GMA representation is used to model the metabolic network behavior assuming steady state conditions. It is considered that all model parameters are deterministic in nature (i.e., perfectly known in advance without any variability). These parameters include the stoichiometric coefficients of the chemical reactions and the transportation processes, as well as the rate constants and kinetic orders of the power-law formalism describing these processes.

Under these conditions, we aim to customize a sBB global optimization method that may improve our previous results for this class of models. Due to the canonical representation provided by the GMA modeling strategy, this goal is of paramount importance for practical biotechnological applications.

## 3. MATHEMATICAL FORMULATION

The optimization problem is mathematically formulated as a MINLP, in which continuous variables denote metabolite concentrations and velocities, and binary variables model the changes in the enzyme levels. We first present the GMA formalism and then introduce the overall MINLP formulation.

**3.1. GMA Representation.** The concentration  $X$  of every single metabolite  $i$  present in the metabolic network can be determined at a particular time  $t$  from the  $p$  flows of the network:

$$\frac{dX_i}{dt} = \sum_{r=1}^p \mu_{ir} v_r \quad i = 1, \dots, n \quad (1)$$

In eq 1, the stoichiometric coefficient,  $\mu_{ir}$ , accounts for the number of molecules of metabolite  $i$  that are involved in process  $r$ . Hence, it is an integer parameter that is positive if process  $r$  contributes to the synthesis of metabolite  $i$ , negative if it depletes the concentration of  $i$ , and zero if process  $r$  does not directly influence the concentration of metabolite  $i$ . The velocity at which process  $r$  occurs, which is denoted by  $v_r$ , is described by a kinetic

equation. In GMA models, the so-called power-law formalism<sup>69–71</sup> is the kinetic equation of choice (eq 2).

$$v_r = \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \quad r = 1, \dots, p \quad (2)$$

Here,  $\gamma_r$  is the basal state activity of the enzyme governing process  $r$ , whereas  $f_{rj}$  is the kinetic order of metabolite  $j$  in process  $r$ . This representation accounts for the  $m$  external (i.e., independent) metabolites, whose concentration is constant throughout the process ( $X_j = \text{constant}$ ,  $j = n + 1, \dots, m$ ). By introducing eq 2 into eq 1 and removing the time dependence (we are interested in solving the steady state for which  $dX_i/dt = 0$  applies), a complete GMA model as in eq 3 is obtained.

$$\sum_{r=1}^p (\mu_{ir} \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}}) = 0 \quad i = 1, \dots, n \quad (3)$$

**3.2. MINLP Formulation.** Since genetic manipulations will take place on an unmodified strain (i.e., at its basal state), it is convenient to express the optimal enzyme activities as a fold-change  $K_r$  over their basal state levels  $\gamma_r$ . According to this, we can rewrite eq 2 as follows:

$$v_r = K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \quad r = 1, \dots, p \quad (4)$$

Here,  $K_r$  is a positive continuous variable that will take the value of 1 at the basal state (i.e., when the enzyme levels are not manipulated). Furthermore,  $K_r > 1$  indicates overexpression of enzyme  $r$ , and  $K_r < 1$  denotes its inhibition. This variable is allowed to change between given bounds,  $K_r^{\text{LB}}$  and  $K_r^{\text{UB}}$  as stated in eq 5.

$$K_r^{\text{LB}} \leq K_r \leq K_r^{\text{UB}} \quad r = 1, \dots, p \quad (5)$$

The number of enzymes that can be modified at a time is constrained to be lower than an upper limit. The motivation for this is that a large number of genetic manipulations might be impractical. This is modeled through a disjunction that determines whether a specific enzyme is modified or not:

$$\left[ K_r^{\text{LB}} \leq K_r \leq 1 - \delta \right] \vee \left[ 1 - \delta \leq K_r \leq 1 + \delta \right] \\ \vee \left[ 1 + \delta \leq K_r \leq K_r^{\text{UB}} \right] \quad Y_{r1}, Y_{r2}, Y_{r3} \in \{\text{True}, \text{False}\} \\ r = 1, \dots, p \quad (6)$$

Here,  $\delta$  is a sufficiently small parameter (i.e., numerical results shown in this work were obtained using a value of  $5 \times 10^{-7}$ ), and  $Y_r$  is a Boolean variable that is true if the associated term of the disjunction is satisfied and false otherwise. The disjunction in eq 6 can be reformulated into linear inequalities by applying either the Big-M or convex hull reformulations.<sup>72,73</sup> The latter, known to provide a relaxation at least as tight as the former,<sup>72</sup> gives rise to eqs 7–11.

$$K_r = K_{r1} + K_{r2} + K_{r3} \quad r = 1, \dots, p \quad (7)$$

$$K_r^{\text{LB}} y_{r1} \leq K_{r1} \leq (1 - \delta) y_{r1} \quad r = 1, \dots, p \quad (8)$$

$$(1 - \delta)y_{r2} \leq K_{r2} \leq (1 + \delta)y_{r2} \quad r = 1, \dots, p \quad (9)$$

$$(1 + \delta)y_{r3} \leq K_{r3} \leq K_r^{UB} y_{r3} \quad r = 1, \dots, p \quad (10)$$

$$y_{r1} + y_{r2} + y_{r3} = 1 \quad r = 1, \dots, p \quad (11)$$

These equations enforce the definition of the binary variables  $y_{r1}$ ,  $y_{r2}$ , and  $y_{r3}$ , which take the value of one if the corresponding term of the disjunction holds true and zero otherwise. These binary variables are then used to define an upper bound ME on the total number of enzymes that can be modified as follows:

$$\sum_{r=1}^p y_{r1} + \sum_{r=1}^p y_{r3} \leq \text{ME} \quad (12)$$

Typically, metabolite concentrations will be allowed to change within given bounds ( $X_i^{LB}$  and  $X_i^{UB}$ , respectively):

$$X_i^{LB} \leq X_i \leq X_i^{UB} \quad i = 1, \dots, n \quad (13)$$

Generally, the objective of these problems is to maximize the synthesis rate of the desired product (note that any other objective function could be evaluated if required). For the sake of simplicity, we pose the problem as a minimization one by reversing the sign of the objective function:

$$\min - \sum_{r=1}^p \mu_{ir} v_r \quad (14)$$

Recall that only the velocities involved in the production of the desired metabolite must be considered in eq 14. The resulting MINLP that embeds the GMA equations can be expressed in compact form as follows:

$$\begin{aligned} (\text{OMINLP}) \min & - \sum_{r=1}^p \mu_{ir} v_r \\ \text{s.t.} & \text{eqs 1, 4, and 7-13} \end{aligned}$$

Model OMINLP [note that the authors have uploaded a similar model to ref 74] seeks the appropriate changes in the enzyme activities (continuous variables) that maximize the synthesis rate of the desired product. The enzyme activities calculated by the model can be implemented in the real system by tuning the expressions of the corresponding genes. Note that when the number of simultaneous modifications is not limited (recall that in our case it is), we can drop the binary variables, which gives rise to a nonconvex NLP problem.

Constraints in OMINLP define a nonconvex search space where multiple local optima may exist. Hence, in order to solve OMINLP to global optimality, we must resort to global optimization techniques.

#### 4. SOLUTION STRATEGY

In this section, we present our customized sBB method for solving problem OMINLP to global optimality. This method makes use of a MILP-based linear relaxation of the nonlinear equations present in the MINLP formulation. We first describe in detail the way in which this relaxation is constructed before presenting the particularities of the sBB algorithm.

**4.1. Relaxed Subproblem.** In order to build a linear relaxation of OMINLP, we introduce two new auxiliary variables,  $k_r$  and  $x_i$ , that are defined by an exponential transformation as follows:

$$K_r = \exp k_r \quad r = 1, \dots, p \quad (15)$$

$$X_i = \exp x_i \quad i = 1, \dots, n \quad (16)$$

These variables replace the original ones,  $K_r$  and  $X_i$ , appearing in eq 4, thus giving rise to eq 17.

$$v_r = (\exp k_r) \gamma_r \prod_{j=1}^{n+m} (\exp x_j)^{f_{rj}} \quad r = 1, \dots, p \quad (17)$$

Let  $p(i)$  denote the number of velocity terms explicitly expressed in the mass balance of metabolite  $i$ , that is, for which  $\mu_{ir} \neq 0$ . Velocities  $v_r$  appearing only in instances of eq 1 with  $p(i) = 2$  are next transferred to linear constraints by introducing eq 17 into eq 1 and taking logarithms as follows:

$$\begin{aligned} 0 &= \mu_{ir} v_r + \mu_{ir'} v_{r'} \\ \mu_{ir} (\exp k_r) \gamma_r \prod_{j=1}^{n+m} (\exp x_j)^{f_{rj}} &= -\mu_{ir'} (\exp k_{r'}) \gamma_{r'} \prod_{j=1}^{n+m} (\exp x_j)^{f_{r'j}} \\ \ln(\mu_{ir}) + k_r + \ln(\gamma_r) + \sum_{j=1}^{n+m} f_{rj} x_j &= \ln(-\mu_{ir'}) + k_{r'} \\ \ln(\gamma_{r'}) + \sum_{j=1}^{n+m} f_{r'j} x_j &\quad \forall i | p(i) = 2 \end{aligned} \quad (18)$$

Recall that when the concentration of a metabolite is only determined by two processes, the stoichiometric coefficient of one of them must be negative, and hence, no domain violation for the logarithmic function can occur in eq 18.

On the other hand, when  $v_r$  appears in at least one instance of eq 1 with more than two terms (i.e.,  $p(i) \geq 3$ ), we make the following changes. We reformulate eq 4 by taking logarithms in both sides of the constraint:

$$\ln(v_r) = k_r + \ln(\gamma_r) + \sum_{j=1}^{n+m} f_{rj} x_j \quad \forall r \in r_{\text{lin}} \quad (19)$$

In this equation,  $r_{\text{lin}}$  denotes the set of velocities  $r$  that are linearized by this process. In mathematical terms,  $r \in r_{\text{lin}} \subset \{r\} \Leftrightarrow \exists i | \mu_{ir} \neq 0 \wedge p(i) \geq 3$ .

The right-hand side of eq 19 is now linear, but the logarithm in the left-hand side gives rise to a nonconvex search space. To linearize this nonconvex term, we reformulate the equation into two inequalities (eqs 20 and 21) and replace their left-hand sides with linear estimators.<sup>75</sup>

$$\ln(v_r) \geq k_r + \ln(\gamma_r) + \sum_{j=1}^{n+m} f_{rj} x_j \quad \forall r \in r_{\text{lin}} \quad (20)$$

$$\ln(v_r) \leq k_r + \ln(\gamma_r) + \sum_{j=1}^{n+m} f_{rj} x_j \quad \forall r \in r_{\text{lin}} \quad (21)$$

The left-hand side of equation eq 20 can be overestimated by  $L$  supporting hyperplanes, which are first-order Taylor expansions of the natural logarithm defined at  $L$  linearization points  $v_r^l$  within the domain  $[v_r^{LB}, v_r^{UB}]$ .

$$\ln v_r \leq \ln v_r^l + \frac{1}{v_r^l} (v_r - v_r^l) \quad \forall r \in r_{\text{lin}} \quad l = 1, \dots, L \quad (22)$$

By combining eq 22 with eq 20, we obtain the following linear constraint (eq 23):

$$\ln v_r^l + \frac{1}{v_r^l}(v_r - v_r^l) \geq k_r + \ln(\gamma_r) + \sum_{j=1}^{n+m} f_{rj}x_j \quad \forall r \in r_{\text{lin}} \quad l = 1, \dots, L \quad (23)$$

Note that the quality of the relaxation depends on the number of linearizations added to the model.

On the other hand, the logarithmic term  $\ln v_r$  in eq 21 is underestimated by means of a piecewise linear function<sup>76–78</sup> defined over  $H$  subintervals within the domain  $[v_r^{\text{LB}}, v_r^{\text{UB}}]$  as follows:

$$\ln v_r \geq \begin{cases} a_r^1 v_r + b_r^1 & v_r^1 \leq v_r \leq v_r^2 \\ a_r^2 v_r + b_r^2 & v_r^2 \leq v_r \leq v_r^3 \\ \dots & \dots \\ a_r^h v_r + b_r^h & v_r^h \leq v_r \leq v_r^{h+1} \\ \dots & \dots \\ a_r^H v_r + b_r^H & v_r^H \leq v_r \leq v_r^{H+1} \end{cases} \quad (24)$$

where  $a_r^h$  and  $b_r^h$  are the coefficients of the straight line that is active in the  $h$ th interval defined by the limits  $v_r^1 = v_r^{\text{LB}}$  and  $v_r^{H+1} = v_r^{\text{UB}}$ . This can be modeled as a disjunction with  $h$  terms as follows:

$$\bigvee_{h=1}^H \left[ \begin{array}{c} z_r^h \\ v_r^h \leq v_r \leq v_r^{h+1} \\ a_r^h v_r + b_r^h \leq k_r + \ln(\gamma_r) + \sum_{j=1}^{n+m} f_{rj}x_j \end{array} \right] \quad \forall r \in r_{\text{lin}} \quad z_r^h \in \{\text{True}, \text{False}\} \quad (25)$$

Here, the Boolean variable  $z_r^h$  indicates whether the  $h$ th interval of the  $r$ th velocity is active or not. The last equation inside the disjunction is obtained by combining eq 21 and eq 24. The disjunction in eq 25 can be translated into linear equations through the convex hull reformulation.

$$v_r = \sum_{h=1}^H v_r^h z_r^h \quad \forall r \in r_{\text{lin}} \quad (26)$$

$$v_r^h z_r^h \leq v_r^h \leq v_r^{h+1} z_r^h \quad \forall r \in r_{\text{lin}} \quad h = 1, \dots, H \quad (27)$$

$$\sum_{h=1}^H z_r^h = 1 \quad \forall r \in r_{\text{lin}} \quad (28)$$

$$\sum_{h=1}^H (a_r^h v_r^h z_r^h + b_r^h z_r^h) \leq k_r + \ln(\gamma_r) + \sum_{j=1}^{n+m} f_{rj}x_j \quad \forall r \in r_{\text{lin}} \quad (29)$$

where  $v_r^h z_r^h$  is a disaggregated variable and  $z_r^h$  is a binary variable that takes the value of 1 if the  $h$ th interval of the  $r$ th velocity is active and 0 otherwise. Note that, in contrast with the supporting hyperplanes, the piecewise formulation does require the definition of binary variables. Hence, a proper balance should be found between the number of intervals and the quality of the relaxation,

so that the computational burden of the model does not explode with the addition of a large number of binary variables.

Finally, eqs 7–10 are rewritten as follows:

$$k_r = k_{r1} + k_{r2} + k_{r3} \quad r = 1, \dots, p \quad (30)$$

$$\ln(K_r^{\text{LB}})y_{r1} \leq k_{r1} \leq \ln(1 - \delta)y_{r1} \quad r = 1, \dots, p \quad (31)$$

$$\ln(1 - \delta)y_{r2} \leq k_{r2} \leq \ln(1 + \delta)y_{r2} \quad r = 1, \dots, p \quad (32)$$

$$\ln(1 + \delta)y_{r3} \leq k_{r3} \leq \ln(K_r^{\text{UB}})y_{r3} \quad r = 1, \dots, p \quad (33)$$

Recall that bounds on variable  $X_i$  need to be expressed in the space of variables  $x_i$  as shown in eq 34.

$$\ln(X_i^{\text{LB}}) \leq x_i \leq \ln(X_i^{\text{UB}}) \quad i = 1, \dots, n \quad (34)$$

The lower bounding problem can be expressed in compact form as follows:

$$\begin{aligned} (\text{CMILP}) \quad & \min - \sum_{r=1}^p \mu_{ir} v_r \\ \text{s.t.} \quad & \text{eqs 1, 11, 12, 18, 23, and 26–34} \end{aligned}$$

It should be clarified that the reformulation presented here is an opt-reformulation since all local and global optima of the original problem are mapped into local and global optima of the reformulated model.<sup>63</sup> Problem CMILP can be solved via standard methods for MILP problems such as the B&B.<sup>46</sup>

**4.2. Customized Spatial Branch-and-Bound.** The spatial branch-and-bound algorithm we propose to solve problem OMINLP exploits the particular features of the GMA model. The method is based on sequentially solving subproblems obtained by partitioning the original domain. A spatial branch-and-bound tree (sBB tree from here on) is used to represent the hierarchy of nodes.

Let OMINLP<sup>*k*</sup> and CMILP<sup>*k*</sup> denote the OMINLP and CMILP subproblems associated with node *k* of the sBB tree. The original problem, OMINLP, is allocated in the root node (*k* = 0). A convex relaxation of the original problem (model CMILP<sup>0</sup>) is solved in order to obtain a valid lower bound on the global optimum of the original formulation.<sup>42,44,79,80</sup> An upper bound for the node can also be computed by optimizing locally the original model OMINLP<sup>0</sup> using the solution provided by CMILP<sup>0</sup> as starting point. If the optimality gap of the node is above the tolerance, then we generate subproblems OMINLP<sup>1</sup> and OMINLP<sup>2</sup> by splitting (branching) the domain of one of the *p* velocities  $v_r$ . This is equivalent to creating two descendant nodes in the sBB tree. Every time a subproblem OMINLP<sup>*k*</sup> is created, it is added to a list *T* containing all of the active (i.e., yet to explore) nodes in the sBB tree. Each of these subproblems is then solved exactly in the same manner as OMINLP<sup>0</sup>, in order to produce lower and upper bounds for each of the nodes. Recall that in these subproblems, we impose lower and upper limits on the variables according to the selected branching scheme. Every time a node is evaluated, the associated OMINLP<sup>*k*</sup> problem is eliminated from *T*.

If at any node *k* of the sBB tree CMILP<sup>*k*</sup> is infeasible, the node can be pruned, as it does not contain any feasible solution to OMINLP. If this happens at node 0, then OMINLP is infeasible. Similarly, if the optimal solution to CMILP<sup>*k*</sup>, denoted by *rOF\**, is above the overall upper bound OUB (i.e., the best bound considering all the nodes of the sBB tree), we can prune this

node, as proceeding in this branch will only lead to worse solutions (note that as we go deeper in the tree, subproblems are more restricted). After updating OUB, we can prune those nodes in the active list with a lower bound greater than OUB.

Search trees are only finite for an  $\varepsilon$ -tolerance.<sup>40</sup> Hence, a node can be fathomed when the difference between the upper and lower bounds is smaller than the tolerance. We update OUB whenever the upper bound of the node is lower than the current OUB. The overall lower bound (OLB) corresponds to the lowest among the lower bounds of the active nodes in the sBB tree. The algorithm terminates when the gap between OUB and OLB is reduced below the  $\varepsilon$ -tolerance.

In the next few sections, we highlight some particular features of our sBB strategy.

**4.2.1. Branching Strategy.** An effective branching technique<sup>50,81</sup> aims at minimizing the size of the sBB tree and, thus, can strongly affect the performance of the algorithm.<sup>63</sup> In contrast with the application of B&B to MILP optimization, where the optimal solution of the relaxation is only infeasible in the original problem when integer variables take fractional values, in nonlinear optimization, infeasibilities may also be due to continuous variables violating constraints that have been relaxed. We must keep in mind that the termination criterion for the proposed strategy is achieving a sufficiently small optimality gap. A tight CMILP formulation capable of providing high-quality lower bounds plays a major role in the performance of the algorithm. Recall that eq 4 is the only equation of OMINLP that is relaxed to build CMILP. Hence, by deriving a tight approximation of the logarithmic function therein, it is possible to determine tight bounds on the global optimal solution of OMINLP. The proposed method branches on the velocities  $v_r | r \in r_{\text{lin}}$ . This is a common feature with the reduced space B&B<sup>57</sup> that only branches on a subset of variables. With this strategy, the linear estimators (i.e., piecewise linear functions and hyperplanes) concentrate on the lower region of the branching velocity in the left-hand descendant subproblem and in the upper region of the velocity in the right-hand one. This improves the quality of the relaxation without increasing the number of variables and the associated complexity.

At each node, the algorithm branches on one single velocity. Our branching strategy consists of branching on the velocity term with the worst relaxation (i.e., the one for which the difference between the solutions of the relaxed and original problem takes a maximum value). Let  $v^{\text{CMILP}^k}$  be the vector containing the value of the  $p$  velocities  $v_r$  in the optimal solution of subproblem CMILP<sup>k</sup> and  $v^{\text{OMINLP}^k}$  be the equivalent vector for subproblem OMINLP<sup>k</sup>. The branching velocity in node  $k$  is that with the largest distance between its optimal value in the original problem and the relaxation:

$$r^k = \arg \max_{r \in r_{\text{lin}}} (\text{abs}(v_r^{\text{CMILP}^k} - v_r^{\text{OMINLP}^k})) \quad (35)$$

If no optimal solution to OMINLP<sup>k</sup> is available (i.e., OMINLP<sup>k</sup> was found infeasible in a local search), the branching velocity is selected with the same equation but  $v_r^{\text{OMINLP}^k}$  is then calculated as a function  $\phi$  of the optimal values of  $k_r^{\text{CMILP}^k}$  and  $x^{\text{CMILP}^k}$ :

$$\begin{aligned} v_r^{\text{OMINLP}^k} &= \phi(k_r^{\text{CMILP}^k}, x^{\text{CMILP}^k}) \\ &= \exp(k_r^{\text{CMILP}^k}) \gamma_r \prod_{j=1}^{n+m} \exp(x_j^{\text{CMILP}^k})^{f_j} \end{aligned} \quad (36)$$

Another important consideration when branching is the selection of the branching point, that is, the point in which the

domain of the branching velocity will be split. One possible strategy consists of using the optimal solution to CMILP<sup>k</sup>,  $v_r^{\text{CMILP}^k}$ , as the branching point. From numerical examples, we found that this strategy usually led to large CPU times, mainly because it produces the same solutions in both descendant nodes. In contrast, allocating this point close to one of the extreme points of  $[v_r^{\text{LB},k}, v_r^{\text{UB},k}]$  is likely to produce a very easy subproblem and a very hard one. The same applies to the rule presented in ref 57, where the branching point is selected as  $v_r^{\text{br},k} = 0.9v_r^{\text{LB},k} + 0.1v_r^{\text{UB},k}$  if  $v_r^{\text{OMINLP}^k} \leq v_r^{\text{mid},k}$  (with  $v_r^{\text{mid},k} = (v_r^{\text{LB},k} + v_r^{\text{UB},k})/2$ ) and  $v_r^{\text{br},k} = 0.1v_r^{\text{LB},k} + 0.9v_r^{\text{UB},k}$  otherwise. Another alternative, perhaps the most intuitive one, is using the bisecting rule, in which the interval is divided by its mid point,  $v_r^{\text{mid},k}$ . Particularly, we have obtained the best performance of the algorithm by applying one of the strategies presented in ref 82. This strategy relies on using a convex combination between the optimal solution  $v_r^{\text{OMINLP}^k}$  and the midpoint of the interval  $v_r^{\text{mid},k}$ , as illustrated by eq 37:

$$v_r^{\text{br},k} = 0.5v_r^{\text{OMINLP}^k} + 0.5v_r^{\text{mid},k} \quad (37)$$

Again, if  $v_r^{\text{OMINLP}^k}$  is not available, it is calculated as in eq 36. With this strategy, we concentrate the efforts around the optimal solution without compromising the balance between the complexity of the two subproblems.

**4.2.2. Bound Contraction and Interval Analysis.** The quality of the OLB strongly depends on the bounds imposed on the variables.<sup>40</sup> These bounds can be tightened during the performance of the algorithm using bound contraction techniques. In general, we can distinguish between two lines of bound tightening procedures: optimality-based bounds tightening (OBBT<sup>55,58,83-87</sup>) and feasibility-based bounds tightening (FBBT<sup>55,66,83,87-91</sup>).

OBBT derives tight bounds for  $n$  variables by solving  $2n$  optimization problems, where each of the  $n$  variables is minimized and maximized subjected to the problem constraints. When  $n$  is large, this procedure becomes time-consuming. Consequently, OBBT is typically performed only in the root node prior to the global optimization procedure.<sup>84</sup> We implement the same strategy, using OBBT to improve the bounds of the  $p$  velocities  $v_r$  by solving subproblems OBLB and OBUB:

$$\begin{aligned} \text{(OBLB) for every } r : & \min v_r \\ & \text{s.t. eqs 1, 11, 12, 18, 23 and 26-34} \end{aligned}$$

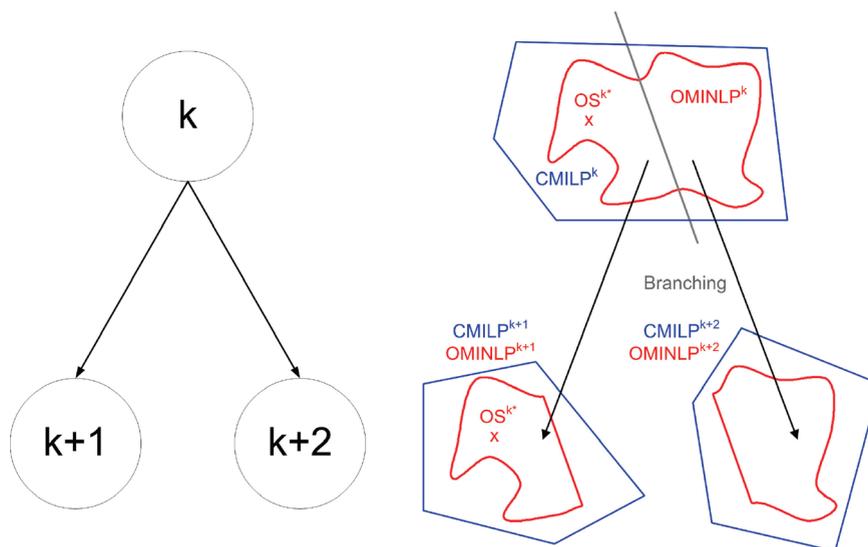
$$\begin{aligned} \text{(OBUB) for every } r : & \max v_r \\ & \text{s.t. eqs 1, 11, 12, 18, 23 and 26-34} \end{aligned}$$

To avoid cutting off feasible values of  $v_r$ , we use the linear relaxation CMILP to generate bounds on the variables. For those cases in which the computational burden of model OBLB/UB is large or the number of velocities is particularly high, we can relax the integer variables in these subproblems in order to expedite their solution. Note that this is done at the expense of obtaining weaker bounds for the velocities.

On the other hand, FBBT inherits the knowledge from recursive arithmetic intervals<sup>47</sup> in order to infer new bounds for the variables from the information provided by the problem constraints. Every time we branch in a node  $k$ , we modify the bounds for the branching velocity in the descendant subproblems as follows:

$$v_r^{\text{UB},k+1} = v_r^{\text{br},k} \quad (38)$$

$$v_r^{\text{LB},k+2} = v_r^{\text{br},k} \quad (39)$$



**Figure 2.** Scheme of sBB partitioning procedure. Solution  $OS^{k^*}$  belongs to the feasible space of subproblem  $OMINLP^{k+1}$ .

where  $k + 1$  and  $k + 2$  denote the left-hand side and right-hand side subproblems, respectively. Consider a hypothetical metabolite  $X_a$  for which the mass balance is described as follows:

$$\frac{dX_a}{dt} = 0 = 2v_1 + v_2 - 3v_3 \quad (40)$$

From this equation, we know that  $v_1^{LB} \geq (3v_3^{LB} - v_2^{UB})/2$  and  $v_1^{UB} \leq (3v_3^{UB} - v_2^{LB})/2$ . Similar expressions can be derived to get bounds on  $v_2$  and  $v_3$ . These equations improve the effect of the branching strategy by generating tighter bounds for variables others than the one on which we have branched. In general, the following expressions hold:

$$v_r^{LB,i} = \sum_{r' \left| \left( \frac{\mu_{ir'}}{-\mu_{ir}} \right) > 0 \right.} \frac{\mu_{ir'}}{-\mu_{ir}} v_r^{LB} + \sum_{r' \left| \left( \frac{\mu_{ir'}}{-\mu_{ir}} \right) < 0 \right.} \frac{\mu_{ir'}}{-\mu_{ir}} v_r^{UB} \quad (41)$$

$$r' \neq r \quad i = 1, \dots, n$$

$$v_r^{UB,i} = \sum_{r' \left| \left( \frac{\mu_{ir'}}{-\mu_{ir}} \right) > 0 \right.} \frac{\mu_{ir'}}{-\mu_{ir}} v_r^{UB} + \sum_{r' \left| \left( \frac{\mu_{ir'}}{-\mu_{ir}} \right) < 0 \right.} \frac{\mu_{ir'}}{-\mu_{ir}} v_r^{LB} \quad (42)$$

$$r' \neq r \quad i = 1, \dots, n$$

Note that each mass balance equation in which velocity  $r$  participates can potentially lead to new tighter bounds. To account for this, we introduce the index  $i$  in the bounds  $v_r^{LB,i}$  and  $v_r^{UB,i}$ . The bounds obtained in each equation are finally compared in order to keep the tightest one:

$$v_r^{LB} = \max(v_r^{LB,old}, \max_{i=1}^n(v_r^{LB,i})) \quad (43)$$

$$v_r^{UB} = \min(v_r^{UB,old}, \min_{i=1}^n(v_r^{UB,i})) \quad (44)$$

Since during the FBBT procedure bounds may be updated, it may be worth it to repeat the process recursively in order to obtain tighter bounds. It is convenient to consider an iteration

limit on the number of times that the procedure is performed. More sophisticated criteria (e.g., repeating the FBBT until the best improvement falls below a given tolerance) can also be used.

It is known that FBBT provides weaker bounds than OBBT.<sup>63</sup> However, it tends to be faster. One of the main advantages of FBBT is that it can detect infeasible subproblems prior to their optimization. A subproblem  $k$  is infeasible when  $v_r^{LB,k} > v_r^{UB,k}$  for at least one  $r$ :

$$\exists r | v_r^{LB,k} > v_r^{UB,k} \rightarrow OMINLP^k = \phi \quad (45)$$

OBBT and FBBT are thus valuable techniques for expediting the overall performance of the algorithm.

**4.2.3. Strengthening Cuts.** A special type of linear cuts that tighten the relaxation of OMINLP can be derived from the stoichiometric coefficients that relate the  $p$  velocities in the mass balance of every dependent metabolite  $i$ . Let us consider the example introduced in the previous section. Two cuts can be deduced from eq 40 as follows:

$$v_3 \geq \frac{2v_1}{3} \quad (46)$$

$$v_3 \geq \frac{v_2}{3} \quad (47)$$

In general, from any mass balance equation associated with metabolite  $i$  with  $p(i)$  velocities in which only one  $\mu_{ir}$  has a different sign than the remaining ones (i.e.,  $\exists r | \mu_{ir} \mu_{ir'} < 0 \forall r \neq r' \wedge \mu_{ir'} \mu_{ir''} > 0 \forall r', r'' \neq r, r' \neq r''$ ), it is possible to generate  $p(i) - 1$  strengthening cuts according to eq 48:

$$v_r \geq \frac{\mu_{ir'}}{-\mu_{ir}} v_{r'} \quad \forall i, r' | \exists r | \mu_{ir} \mu_{ir'} < 0 \forall r \neq r' \wedge \mu_{ir'} \mu_{ir''} > 0 \quad (48)$$

$$\forall r'' \neq r, r'$$

These inequalities can be obtained offline and added to CMILP before the optimization takes place. A major advantage of these cuts is that we can easily linearize them by applying the exponential transformation described before. Particularly, if we introduce eq 4 into eq 48 and replace the original  $X_i$  and  $K_r$  as



**Table 1. Size of Citric Acid Models after Preprocessing<sup>a</sup>**

case	OMINLP			PW <sub>0</sub>	CMILP		
	equations	CV	IV		equations	CV	IV
B1	692	448	3	12	5339	1072	924
B2	692	448	3	10	5111	958	810
C1	692	445	6	10	5111	958	810
C2	692	445	6	30	7451	2098	1890
D1	692	442	9	15	5681	1243	1095
D2	692	442	9	10	5111	958	810
E1	692	436	15	12	5339	1072	924
E2	692	436	15	15	5681	1243	1095

<sup>a</sup> OMINLP: full-space problem. CMILP: MIP relaxation of OMINLP. CV: number of continuous variables. IV: number of integer variables. PW<sub>0</sub>: number of piecewise sections in the initial iteration of the algorithm.

**Table 2. Enzymes That Can Be Modified in Each of the Instances of OMINLP**

case	ME	subcase	modifiable			case	ME	subcase	modifiable		
			enzymes						enzymes		
B	1	1	[40]	D	3	1	[1, 40, 60]	2	[1, 40, 59]		
		2	[59]								
C	2	1	[40, 59]	E	5	1	[1, 39, 40, 59, 60]	2	[1, 28, 40, 59, 60]		
		2	[1, 40]								

OS<sup>k\*</sup> in the space of variables of the linear relaxation. We accomplish this by applying a logarithmic transformation on the continuous variables, and by fixing the values of the binary variables according to the intervals of the piecewise approximation in which the original solution has fallen. This provides an integer feasible solution that is used as a starting value for the B&B solvers, thereby expediting the solution of the lower bounding problem in node  $k + 1$ . Note that this initialization scheme is only applicable to one descendant node.

## 5. COMPUTATIONAL RESULTS

The problem selected for testing the capabilities of our customized sBB algorithm is the maximization of the citric acid production in *Aspergillus niger* (see Figure 3). On the basis of the results of Polisetty et al.<sup>67</sup> and Pozo et al.,<sup>68</sup> we solve several instances of OMINLP, which differ in the number of reactions (ME) allowed for simultaneous modification. We assume that the 60 reactions included in the model can be modified by genetic manipulation. Note, however, that any practical solution should consider only a limited number of changes. In this specific case, Polisetty et al.<sup>67</sup> showed that by manipulating only 5 enzymes, it is possible to attain a solution close to the one found when all of the enzymes can be modified.

Here, we take the results from Pozo et al.<sup>68</sup> as a reference for comparison purposes. We focus on optimizing the system when only one, two, three, or five enzymes can be modified (case B, ME = 1; case C, ME = 2; case D, ME = 3; and case E, ME = 5). The nomenclature is the same used in Pozo et al.<sup>68</sup> The cases discussed in Table 8 of that paper are used to test the performance of the novel sBB method. A total of eight instances are solved with the customized sBB approach (see Tables 1 and 2). In all these cases, those enzymes that are not allowed for modification are fixed to their basal state. The maximum change

**Table 3. Parameters Setting in the sBB Algorithm**

parameter	configuration
node selection	highest LB
CPLEX tolerance	0.00%
FBBT stop criterion	10 iterations
number of hyperplanes	50
branching point selection	see eq 37

for each enzyme is 5 fold over its basal state. The optimization constraints are the same as in the referenced paper.

Our results are compared with those obtained by the OA technique introduced by the authors in an earlier work<sup>21,68</sup> and also with the global optimization package BARON. With regard to the sBB and OA methods, it should be noted that both of them solve iteratively the same subproblems: the MILP-based relaxation CMILP and the bounded OMINLP. From numerical examples, we observed that both algorithms worked better when the binary variables associated with the genetic manipulations of the enzyme levels are fixed in the original problem according to the output of the linear MILP relaxation. For this reason, the lower bounds are generated by solving a bounded NLP instead of a bounded MINLP. [Note that the optimization task is posed as a maximization problem, so CMILP predicts upper bounds on the global optimum of OMINLP.]

In all of the examples, we used CPLEX 11.2.1 as MILP solver and CONOPT 3.14s for the NLPs, whereas BARON v.8.1.5<sup>92</sup> was employed to solve the full-space OMINLP problems. The algorithms were implemented in GAMS 23.0.2 on an Intel 1.2 GHz machine. An optimality tolerance of 2.00% was fixed in all of the cases.

The performance of the sBB algorithm depends on a series of factors that can be configured at will. The ones with the highest influence are the branching rule, the CPLEX tolerance, the stop criterion for the FBBT procedure, the selection of the branching point, the number of supporting hyperplanes, and the number of piecewise intervals. The particular configuration of the algorithm chosen to perform the calculations is given in Table 3. The only parameter that was particularly tuned for every single instance being solved was the number of piecewise intervals. The results obtained with the aforementioned sBB configuration are shown in Table 4, which also summarizes the performance of the other methods. Recall that the optimal solution reported corresponds to the best solution provided by the lower bounding problem when the termination criterion was met. Note also that all of the algorithms are compared on the basis of the CPU time required to attain a solution with an optimality gap of 2.00%. This is the same comparison criterion used in Pozo et al.<sup>68</sup> Other criteria could have been used instead. Nevertheless, we observed that the conclusions of the analysis are very similar for all of the cases.

As can be seen, the proposed methodology can solve all of the instances within the required tolerance. The same occurs with the OA, whereas BARON was not able to improve the starting point (which corresponds to the basal state solution) even after 3600 s of CPU time. This might be due to the use of generic techniques for building the relaxed upper bounding problem (when maximizing) that do not benefit from the particular structure of the GMA formalism.

The number of nodes explored in the sBB tree varies from one example to another without any clear tendency. However, the node in which the optimal solution is found is generally very close to the root node (20 first nodes) in all the instances except C1

Table 4. Comparison between the Best Results Obtained with the OA and the Customized sBB for Each Instance<sup>a</sup>

case	sBB						BARON			OA				
	PW <sub>0</sub> <sup>b</sup>	nodes	NO	LB	UB	CPU	LB	UB	CPU	PW <sub>0</sub>	PW <sub>f</sub>	LB	UB	CPU
B1	12	34	10	25.82	26.05	36	12.36	— <sup>c</sup>	17	3	12	25.82	26.33	17
B2	10	198	12	12.37	12.57	229	12.36	—	9	3	15	12.35	12.51	45
C1	10	298	244	25.83	26.34	368	12.36	—	13	3	16	25.78	26.14	53
C2	30	26	15	25.82	26.34	61	12.36	—	39	3	12	25.82	26.33	18
D1	15	340	1	40.88	41.59	1155	12.36	—	52	3	18	40.88	41.42	167
D2	10	42	18	176.8	179.87	35	12.36	176.79	3600	3	12	176.79	180	18
E1	12	1	1	347.92	353.22	1	12.36	347.93	3600	7	9	347.93	353.35	6
E2	15	32	1	256.59	261.68	83	12.36	256.68	3600	3	18	256.59	261.32	1093

<sup>a</sup> PW<sub>0</sub>: number of piecewise sections in the initial iteration of the algorithm. PW<sub>f</sub>: number of piecewise sections in the last iteration of the algorithm. NO: node in which the optimal solution was found. LB: lower bound on the global optimum in mM min<sup>-1</sup>. UB: upper bound on the global optimum in mM min<sup>-1</sup>. CPU: CPU time in seconds. <sup>b</sup> Note that for the sBB algorithm, PW<sub>0</sub> = PW<sub>f</sub> as the number of piecewise sections is not modified throughout the algorithm. <sup>c</sup> BARON failed to provide a rigorous upper bound in cases with —.

(244). This is in consonance with the common observation that B&B algorithms may take a long time to verify optimality, although good (sometimes optimal) solutions are usually found in the early stages of the search.<sup>48</sup>

Regarding the optimal number of piecewise intervals, they range between 10 and 15. In only one case (C2) did the customized sBB perform better with a larger number of piecewise terms (i.e., 30). In contrast, in the OA, the optimal number of initial piecewise terms is small on average (i.e., 4), and even in the worst case, the algorithm works better with fewer intervals (i.e., 9) than in any instance of the sBB. Recall that, in each iteration of the OA algorithm, the total number of piecewise sections is increased by one, as the interval containing the optimal solution of CMILP is split into two subintervals. Since there is a binary variable associated with each of these new intervals, starting with too many sections is likely to lead to large instances that are hard to solve in short CPU times.

On the other hand, the final number of piecewise sections required by the OA exceeds the piecewise terms used in the sBB in all of the cases except two (i.e., C2 and E1). Let us note that in the OA method, the number of piecewise intervals is progressively increased to construct tighter relaxations, whereas in the customized sBB, bound tightening techniques and tailored branching rules allow reducing of the search space while keeping the number of piecewise terms constant in each subproblem. Henceforth, the sBB can produce a relaxation as tight as that of the OA with fewer piecewise sections.

With regard to the CPU time, the OA proved to be faster on average (177s compared to 253s). Specifically, it performed better than the sBB in more instances (6 vs 2). The reason for this to happen might be that at each iteration of the OA, we tighten the relaxation of the logarithm of all the velocities  $|r_{\text{lin}}|$ , whereas in the sBB, only one velocity is tightened at each branching point of the tree. This leads to more nodes and hence large CPU times. This finally results in a faster convergence of the lower and upper bounds in the case of the OA.

The customized sBB algorithm proved to be significantly faster than the OA in the two most difficult instances (i.e., those with a higher number of manipulations allowed). It should be noted that these results may vary according to the settings of each algorithm.

In order to better explore the advantages of the proposed method, we also solved instances of cases B, E, and F (the latter

corresponding to ME = 4), where the best combination of enzymatic manipulations is searched. In these computations, no binary variables are fixed prior to the optimization. This exercise attempts to mimic the search for an optimal genetic modification in a biotechnological application.

As can be seen in Table 5, the solution identified for case B by both algorithms corresponds to the one obtained in case B1 (see Table 4). The same occurs with the solution obtained by the sBB for case E, which is the same as that of case E1. In contrast, the solution of case E found with the OA corresponds to a different combination of enzymes that leads to a lower (i.e., worse) value of the objective function. Solutions found for ME = 1 and ME = 5 did not improve those already reported in Polisetty et al.<sup>67</sup> and Pozo et al.<sup>68</sup> On the other hand, the solution obtained for case F, which implies modifying four enzymes, is indeed very close to the one found when all the enzymes are allowed to change. In addition, this solution can be attained modifying different combinations of enzymes.

Regarding the performance of the algorithms, both of them were capable of finding solutions with low optimality gaps in all of the instances, showing the sBB method the best performance. These results are partly due to the bound tightening techniques discussed in section 4.2.2, which were not included in the OA proposed in Pozo et al.<sup>68</sup> For instance, note that the customized sBB identifies the global optimum of cases E and F in just one node using 11 and 12 piecewise intervals respectively. One could think that the same result could be obtained with the OA using the same number of piecewise sections. However, when no bound contraction is performed, this setup does not allow for attaining the specified tolerance in one iteration of the OA algorithm. In fact, with this number of initial piecewise terms, the algorithm shows a worse performance than with 4 intervals (i.e., the final CPU times exceed those reported in Table 5 for four initial sections).

These results also indicate that good/optimal solutions can be found in the early stages of the search. Note that the global optimum of case B is identified in node 238 with sBB, but a solution with an optimality gap of 3% is found in node 1. This can be better seen in Figure 4: solutions very close to the global optimum are identified in the very first seconds of the execution of the algorithms, while the remaining time is spent in reducing the optimality gap. Particularly, the OA provides smaller optimality gaps than the sBB in the beginning of the search (i.e., in the

Table 5. Comparison between the Best Results Obtained with the Customized sBB and the OA for Each Instance<sup>a</sup>

case	sBB							OA							
	PW <sub>0</sub> <sup>b</sup>	nodes	NO	OE	K <sub>r</sub>	LB	UB	CPU	PW <sub>0</sub>	PW <sub>f</sub>	OE	K <sub>r</sub>	LB	UB	CPU
B	3	296	238	[40]	[5.00]	25.82	26.25	484	2	12	[40]	[5.00]	25.82	26.19	864
E	11	1	1	[1, 39, 40, 59, 60]	[1.40, 0.92, 5.00, 5.00, 1.07]	347.92	353.81	91	4	7	[28, 39, 40, 59, 60]	[1.13, 1.45, 5.00, 5.00, 1.05]	347.26	353.75	203
F	12	1	1	[1, 39, 40, 59]	[1.23, 0.88, 5.00, 5.00]	347.25	351.41	50	4	7	[39, 40, 55, 59]	[1.53, 5.00, 1.10, 5.00]	347.26	351.71	78

<sup>a</sup> PW<sub>0</sub>: number of piecewise sections in the initial iteration of the algorithm. NO: node in which the optimal solution was found. OE: enzymes *r* being modified in the optimal solution of the instance. K<sub>r</sub>: optimal fold-change in activity of enzyme *r*. LB: lower bound on the global optimum in mM min<sup>-1</sup>. UB: upper bound on the global optimum in mM min<sup>-1</sup>. CPU: CPU time in seconds. PW<sub>f</sub>: number of piecewise sections in the last iteration of the algorithm. <sup>b</sup> Note that for the sBB algorithm, PW<sub>0</sub> = PW<sub>f</sub> as the number of piecewise sections is not modified throughout the algorithm.

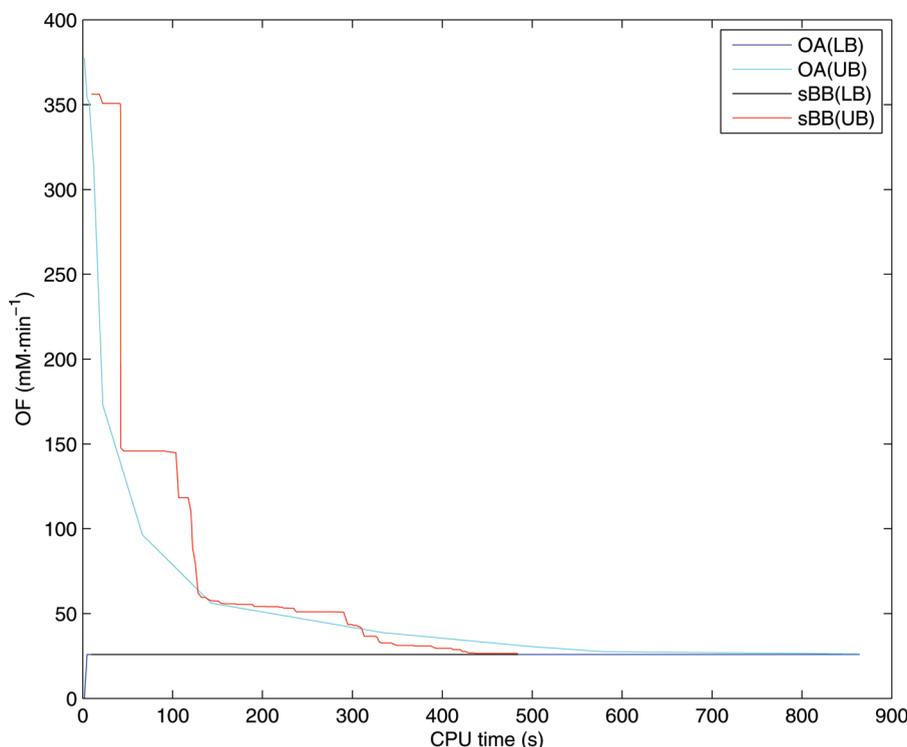


Figure 4. Evolution of the lower and the upper bounds of the global optimum of case B for the OA and the sBB algorithms.

first 300s). At this point, the tendency changes and the sBB shows better performance. This is due to the increase in the number of binary variables and hence in the complexity of the MILP subproblems solved by the OA. In contrast, the size of the MILP subproblems calculated in the sBB is kept constant in the nodes of the tree.

## 6. CONCLUSIONS

This paper has addressed the global optimization of metabolic networks described through the GMA formalism. A customized sBB algorithm that benefits from the specific structure of this type of model has been presented for this purpose. The optimization task was posed as a nonconvex MINLP in which integer variables denote the number of manipulations allowed. Tight bounds on the global optimum were obtained by constructing a linear MILP-based relaxation that exploits the mathematical structure of the GMA formalism. The method incorporates branching rules and bound contraction strategies devised to expedite the overall solution procedure.

Our strategy was compared against an outer approximation (OA) algorithm and the global optimization package BARON.

Numerical results showed that the first two methods outperformed BARON in all of the instances under study. This is due to the quality of the MILP-based relaxation that is obtained by performing a logarithmic transformation on the power-law equations and approximating them by under- and overestimators. We also observed that none of these two methods (sBB and OA) proved to be superior in all of the cases. Nevertheless, the sBB showed a better performance in the most complicated instances, which is probably due to the ability of this strategy to reduce the problem domain without increasing the number of variables. Problems with a similar structure (i.e., with a large number of sigmoidal terms) may also benefit from the proposed strategy. Future work will focus on devising systematic tuning strategies that will improve the performance of the customized sBB algorithm.

The results obtained clearly show that we can tackle problems of moderate complexity when expressed as GMA models. One difficulty encountered when addressing the global optimization of complex metabolic networks is the current limited biological knowledge of some of these systems. While stoichiometric models can easily be constructed, GMA models require

additional information that may not be available for large models. Although detailed GMA genome-wide models are far in the future, our results show that it is worth it to collect the required information, as we are able to obtain optimization results that go beyond those possible with stoichiometric models.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: gonzalo.guillen@urv.cat.

## ACKNOWLEDGMENT

The authors wish to acknowledge support of this research work from the Spanish Ministry of Education and Science (projects DPI2008-04099, PHB2008-0090-PC, BFU2008-00196 and CTQ2009-14420-C02), the Spanish Ministry of External Affairs (projects HS2007-0006 and A/023551/09), and the Generalitat de Catalunya (FI programs).

## NOMENCLATURE

### Indexes

$h$  = interval of the piecewise underestimation of the logarithm function  
 $i$  = dependent metabolite  
 $j$  = metabolite (dependent or independent)  
 $l$  = supporting hyperplane  
 $r$  = flow, process, velocity

### Sets

$r_{\text{lin}}$  = set of processes  $r$  whose kinetic equations are linearized

### Variables

$K_r$  = fold-change in the basal state activity of enzyme governing process  $r$   
 $k_r$  = logarithm of the fold change in the basal state activity of enzyme governing process  $r$   
 $K_{r,1}$  = auxiliary disaggregated variable associated with process  $r$   
 $K_{r,2}$  = auxiliary disaggregated variable associated with process  $r$   
 $K_{r,3}$  = auxiliary disaggregated variable associated with process  $r$   
 $t$  = time  
 $v_r$  = velocity of process  $r$   
 $v_r^h$  = disaggregated variable associated with the  $h$ th term of the convex hull reformulation of the piecewise underestimator of velocity  $r$   
 $X_i$  = concentration of metabolite  $i$   
 $x_i$  = logarithm of the concentration of metabolite  $i$   
 $y_{r,1}$  = binary variable associated with the first term of the convex hull of the disjunction of process  $r$   
 $y_{r,2}$  = binary variable associated with the second term of the convex hull of the disjunction of process  $r$   
 $y_{r,3}$  = binary variable associated with the third term of the convex hull of the disjunction of process  $r$   
 $z_r^h$  = binary variable associated with the  $h$ th term of the convex hull of the piecewise underestimation of velocity  $r$

### Parameters

$\delta$  = sufficiently small parameter  
 $\gamma_r$  = basal state activity of enzyme governing process  $r$   
 $\mu_{ir}$  = stoichiometric coefficient of process  $r$  in the mass balance of metabolite  $i$   
 $a_r^h$  = slope of the segment used in interval  $h$  of the piecewise

approximation of velocity  $r$   
 $b_r^h$  = vertical axis intercept of the segment used in interval  $h$  of the piecewise approximation of velocity  $r$   
 $f_{rj}$  = kinetic order of metabolite  $j$  in process  $r$   
 $H$  = total number of intervals in the piecewise underestimator of the logarithmic function  
 $K_r^{\text{LB}}$  = lower bound on the fold change in the basal state activity of enzyme governing process  $r$   
 $K_r^{\text{UB}}$  = upper bound on the fold change in the basal state activity of enzyme governing process  $r$   
 $L$  = total number of supporting hyperplanes (linearization points)  
 $m$  = total number of independent metabolites  
 $ME$  = maximum number of enzymes allowed for modification  
 $n$  = total number of dependent metabolites  
 $p$  = total number of flows (processes) involved in the metabolic network under study  
 $p(i)$  = total number of flows (processes) involved in the mass balance of metabolite  $i$   
 $v_r^h$  = lower limit of interval  $h$  in the piecewise underestimator of velocity  $r$   
 $v_r^{h+1}$  = upper limit of interval  $h$  in the piecewise underestimator of velocity  $r$   
 $v_r^{\text{LB}}$  = lower bound on velocity  $r$   
 $v_r^{\text{UB}}$  = upper bound on velocity  $r$   
 $X_i^{\text{LB}}$  = lower bound on the concentration of metabolite  $i$   
 $X_i^{\text{UB}}$  = upper bound on the concentration of metabolite  $i$

## REFERENCES

- (1) Guell, M.; van Noort, V.; Yus, E.; Chen, W. H.; Leigh-Bell, J.; Michalodimitrakis, K.; Ya-mada, T.; Arumugam, M.; Doerks, T.; Kuhner, S.; Rode, M.; Suyama, M.; Schmidt, S.; Gavin, A. C.; Bork, P.; et al. Transcriptome complexity in a genome-reduced bacterium. *Science (Washington, DC, U.S.)* **2009**, *326*, 1268–1271.
- (2) Kuhner, S.; van Noort, V.; Betts, M. J.; Leo-Macias, A.; Batisse, C.; Rode, M.; Yamada, T.; Maier, T.; Bader, S.; Beltran-Alvarez, P.; Castano-Diez, D.; Chen, W. H.; Devos, D.; Guell, M.; Norambuena, T.; et al. Proteome organization in a genome-reduced bacterium. *Science (Washington, DC, U.S.)* **2009**, *326*, 1235–1240.
- (3) Yus, E.; Maier, T.; Michalodimitrakis, K.; van Noort, V.; Yamada, T.; Chen, W. H.; Wodke, J. A.; Guell, M.; Martinez, S.; Bourgeois, R.; Kuhner, S.; Raineri, E.; Letunic, I.; Kalinina, O. V.; Rode, M.; et al. Impact of genome reduction on bacterial metabolism and its regulation. *Science (Washington, DC, U.S.)* **2009**, *326*, 1263–1268.
- (4) Gibson, D. G.; Glass, J. I.; Lartigue, C.; Noskov, V. N.; Chuang, R. Y.; Algire, M. A.; Benders, G. A.; Montague, M. G.; Ma, L.; Moodie, M. M.; Merryman, C.; Vashee, S.; Krishnakumar, R.; Assad-Garcia, N.; Andrews-Pfannkoch, C.; et al. Creation of a Bacterial Cell Controlled by a Chemically Synthesized Genome. *Science (New York, N.Y.)* **2010**, .
- (5) Vilaprinyo, E.; Alves, R.; Sorribas, A. Use of physiological constraints to identify quantitative design principles for gene expression in yeast adaptation to heat shock. *BMC Bioinf.* **2006**, *7*, 184.
- (6) Feist, A. M.; Palsson, B. O. The growing scope of applications of genome-scale metabolic reconstructions using *Escherichia coli*. *Nat. Biotechnol.* **2008**, *26*, 659–667.
- (7) Oberhardt, M. A.; Palsson, B. O.; Papin, J. A. Applications of genome-scale metabolic reconstructions. *Mol. Syst. Biol.* **2009**, *5*, 320.
- (8) Vital-López, F.; Armaou, A.; Nikolaev, E.; Maranas, C. A. Computational Procedure for Optimal Engineering Interventions Using Kinetic Models of Metabolism. *Biotechnol. Prog.* **2006**, *22*, 1507–1517.
- (9) Banga, J. Optimization in computational systems biology. *BMC Syst. Biol.* **2008**, *2*–47.
- (10) Bailey, J. Toward a science of metabolic engineering. *Science* **1991**, *252*, 1668–1675.

- (11) Hatzimanikatis, V.; Floudas, C.; Bailey, J. Optimization of regulatory architectures in metabolic reaction networks. *Biotechnol. Bioeng.* **1996**, *52*, 485–500.
- (12) Voit, E. Optimization in integrated biochemical systems. *Biotechnol. Bioeng.* **1992**, *40* (5), 572–582.
- (13) Bailey, J. Lessons from metabolic engineering for functional genomics and drug discovery. *Nat. Biotechnol.* **1999**, *17*, 616–618.
- (14) Marin-Sanguino, A.; Torres, N. Optimization of biochemical systems by linear programming and general mass action model representations. *Math. Biosci.* **2003**, *184* (2), 187–200.
- (15) Marin-Sanguino, A.; Voit, E.; Gonzalez-Alzon, C.; Torres, N. Optimization of biotechnological systems through geometric programming. *Theor. Biol. Med. Model.* **2007**, *4*, 4–38.
- (16) Alvarez-Vasquez, F.; Canovas, M.; Iborra, J.; Torres, N. Modeling, optimization and experimental assessment of continuous L-(–)-carnitine production by *Escherichia coli* cultures. *Biotechnol. Bioeng.* **2002**, *80* (7), 794–805.
- (17) Torres, N. V.; Voit, E. O.; Gonzalez-Alcon, C. Optimization of nonlinear biotechnological processes with linear programming: Application to citric acid production by *Aspergillus niger*. *Biotechnol. Bioeng.* **1996**, *49*, 247–258.
- (18) Alvarez-Vasquez, F.; Gonzalez-Alcon, C.; Torres, N. V. Metabolism of citric acid production by *Aspergillus niger*: model definition, steady-state analysis and constrained optimization of citric acid production rate. *Biotechnol. Bioeng.* **2000**, *70*, 82–108.
- (19) Lin, H.; Bennett, G. N.; San, K. Y. Metabolic engineering of aerobic succinate production systems in *Escherichia coli* to improve process productivity and achieve the maximum theoretical succinate yield. *Metab. Eng.* **2005**, *7*, 116–127.
- (20) Chang, Y.; Sahinidis, N. Optimization of metabolic pathways under stability considerations. *Comput. Chem. Eng.* **2005**, *29*, 467–479.
- (21) Guillén-Gosálbez, G.; Sorribas, A. Identifying quantitative operation principles in metabolic pathways: a systematic method for searching feasible enzyme activity patterns leading to cellular adaptive responses. *BMC Bioinf.* **2009**, *10*, 386.
- (22) Sorribas, A.; Pozo, C.; Vilaprinyo, E.; Guillén-Gosálbez, G.; Jiménez, L.; Alves, R. Optimization and evolution in metabolic pathways: Global optimization techniques in Generalized Mass Action models. *J. Biotechnol.* **2010**, DOI: 10.1016/j.jbiotec.2010.01.026.
- (23) Bower, J.; Bolouri, H. *Computational Modeling of Genetic and Biochemical Networks*; MIT Press: London, 2004.
- (24) Orth, J. D.; Thiele, I.; Palsson, B. O. What is flux balance analysis? *Nat. Biotechnol.* **2010**, *28*, 245–248.
- (25) Edwards, J.; Ibarra, R.; Palsson, B. In silico predictions of *Escherichia coli* metabolic capabilities are consistent with experimental data. *Nat. Biotechnol.* **2001**, *19* (2), 125–130.
- (26) Forster, J.; Famili, I.; Fu, P.; Palsson, B.; Nielsen, J. Genome-scale reconstruction of the *Saccharomyces cerevisiae* metabolic network. *Genome Res.* **2003**, *13* (2), 244–253.
- (27) Alper, H.; Jin, Y.; Moxley, J.; Stephanopoulos, G. Identifying gene targets for the metabolic engineering of lycopene biosynthesis in *Escherichia coli*. *Metab. Eng.* **2005**, *7* (3), 155–164.
- (28) Cox, S.; Levanon, S.; Sanchez, A.; Lin, H.; Peercy, B.; Bennett, G.; San, K. Development of a metabolic network design and optimization framework incorporating implementation constraints: A succinate production case study. *Metab. Eng.* **2006**, *8* (1), 46–57.
- (29) Pramanik, J.; Keasling, J. Stoichiometric model of *Escherichia coli* metabolism: Incorporation of growth-rate dependent biomass composition and mechanistic energy requirements. *Biotechnol. Bioeng.* **1997**, *56* (4), 398–421.
- (30) Fong, S.; Burgard, A.; Herring, C.; Knight, E.; Blattner, F.; Maranas, C.; Palsson, B. In silico design and adaptive evolution of *E. coli* for production of lactic acid. *Biotechnol. Bioeng.* **2005**, *91* (5), 643–648.
- (31) Voit, E. Design principles and operating principles: the yin and yang of optimal functioning. *Math. Biosci.* **2003**, *182*, 81–92.
- (32) Voit, E. O.; Savageau, M. A. Accuracy of alternative representations for integrated biochemical systems. *Biochemistry* **1987**, *26*, 6869–6880.
- (33) Sorribas, A.; Curto, R.; Cascante, M. Comparative characterization of the fermentation pathway of *Saccharomyces cerevisiae* using biochemical systems theory and metabolic control analysis: model validation and dynamic behavior. *Math. Biosci.* **1995**, *130*, 71–84.
- (34) Cascante, M.; Curto, R.; Sorribas, A. Comparative characterization of the fermentation pathway of *Saccharomyces cerevisiae* using biochemical systems theory and metabolic control analysis: steady-state analysis. *Math. Biosci.* **1995**, *130*, 51–69.
- (35) Curto, R.; Sorribas, A.; Cascante, M. Comparative characterization of the fermentation pathway of *Saccharomyces cerevisiae* using biochemical systems theory and metabolic control analysis: Model definition and nomenclature. *Math. Biosci.* **1995**, *130*, 25–50.
- (36) Alves, R.; Vilaprinyo, E.; Hernandez-Bermejo, B.; Sorribas, A. Mathematical formalisms based on approximated kinetic representations for modeling genetic and metabolic pathways. *Biotechnol. Genet. Eng. Rev.* **2008**, *25*, 1–40.
- (37) Alves, R.; Vilaprinyo, E.; Sorribas, A. Integrating Bioinformatics and Computational Biology: Perspectives and Possibilities for In Silico Network Reconstruction in Molecular Systems Biology. *Curr. Bioinf.* **2008**, *3*, 98–129.
- (38) Smith, E.; Pantelides, C. A symbolic reformulation/spatial branch-and-bound algorithm for the global optimization of nonconvex MINLPs. *Comput. Chem. Eng.* **1999**, *23*, 457–478.
- (39) Floudas, C.; Akrotirianakis, I.; Caratzoulas, S.; Meyera, C.; Kallrath, J. Global optimization in the 21st century: Advances and challenges. *Comput. Chem. Eng.* **2005**, *29*, 1185–1202.
- (40) Grossmann, I.; Bigler, L. Part II. Future perspective on optimization. *Comput. Chem. Eng.* **2004**, *28*, 1193–1218.
- (41) Horst, R.; Thoai, N.; Vries, L. D. A new simplicial cover technique in constrained global optimization. *J. Global Optimiz.* **1992**, *2*, 1–19.
- (42) Falk, J.; Soland, R. An algorithm for separable nonconvex programming problems. *Manage. Sci.* **1969**, *15*, 550–569.
- (43) Al-Khayyal, F. Generalized bilinear programming. Part I. Models, applications and linear programming relaxation. *Eur. J. Operat. Res.* **1992**, *60*, 306–314.
- (44) Al-Khayyal, F.; Falk, F. Jointly constrained biconvex programming. *Math. Operat. Res.* **1983**, *8*, 273–286.
- (45) R, H.; Tuy, H. On the convergence of global methods in multiextremal optimization. *J. Optimiz. Theory Appl.* **1987**, *54*, 253.
- (46) Horst, R.; Tuy, H. *Global Optimization: Deterministic Approaches*, 2nd ed.; Springer-Verlag: Berlin, 1993.
- (47) Ryoo, H.; Sahinidis, N. V. Global optimization of nonconvex NLPs and MINLPs with applications in process design. *Comput. Chem. Eng.* **1995**, *19* (5), 551–566.
- (48) Ryoo, H.; Sahinidis, N. V. A branch-and-reduce approach to global optimization. *J. Global Optimiz.* **1996**, *8* (2), 107–138.
- (49) Adjiman, C.; Androulakis, I.; Floudas, C. A global optimization method,  $\alpha$ bb, for general twice-differentiable constrained NLPs: II. Implementation and computational results. *Comput. Chem. Eng.* **1998**, *22* (9), 1159–1179.
- (50) Adjiman, C.; Androulakis, I.; Floudas, C. Global optimization of MINLP problems in process synthesis and design. *Comput. Chem. Eng.* **1997**, *21*, S445–S450.
- (51) Adjiman, C.; Androulakis, I.; Maranas, C.; Floudas, C. A global optimization method,  $\alpha$ BB, for process design. *Comput. Chem. Eng.* **1996**, *20*, S419–S424.
- (52) Adjiman, C.; Dallwig, S.; Floudas, C.; Neumaier, A. A global optimization method,  $\alpha$ BB, for general twice-differentiable constrained NLPs: I. Theoretical advances. *Comput. Chem. Eng.* **1998**, *22* (9), 1137–1158.
- (53) Adjiman, C.; Floudas, C. Rigorous convex underestimators for general twice-differentiable problems. *J. Global Optimiz.* **1996**, *9* (1), 23–40.
- (54) Adjiman, C.; Schweiger, C.; Floudas, C. Mixed-integer nonlinear optimization in process synthesis. *Handbook of Combinatorial Optimization*; Kluwer Academic Publishers: Norwell, MA, 1998.
- (55) Smith, E. *On the Optimal Design of Continuous Processes*; Ph.D. thesis, Imperial College of Science, Technology and Medicine, University of London, 1996.

- (56) Smith, E.; Pantelides, C. Global optimization of nonconvex MINLPs. *Comput. Chem. Eng.* **1997**, *21*, S791–S796.
- (57) Epperly, T.; Pistikopoulos, E. A reduced space branch and bound algorithm for global optimization. *J. Global Optimiz.* **1997**, *11*, 287–311.
- (58) Zamora, J.; Grossmann, I. A branch and contract algorithm for problems with concave univariate, bilinear and linear fractional terms. *J. Global Optimiz.* **1999**, *14*, 217–249.
- (59) Kesavan, P.; Barton, P. Generalized branch-and-cut framework for mixed-integer nonlinear optimization problems. *Comput. Chem. Eng.* **2000**, *24*, 1361–1366.
- (60) O'Grady, A.; Bogle, I.; Fraga, E. Interval analysis in automated design for bounded solutions. *Chem. Zvesti* **2001**, *55* (6), 376–381.
- (61) Vaidyanathan, R.; El-Halwagi, M. Global optimization of non-convex MINLPs by interval analysis. In *Global Optimization in Engineering Design*; Grossmann, I., Ed.; Kluwer Academic Publishers: Dordrecht, The Netherlands, 1996; pp 175–193.
- (62) Zilinskas, J.; Bogle, I. Evaluation ranges of functions using balanced random interval arithmetic. *Informatica* **2003**, *14* (3), 403–416.
- (63) Belotti, P.; Lee, J.; Liberti, L.; Margot, F.; Wächter, A. Branching and bounds tightening techniques for non-convex MINLP. *Optimiz. Methods Software* **2009**, *24* (4–5), 597–634.
- (64) Sahinidis, N.; Grossmann, I. MINLP model for cyclic multi-product scheduling on continuous parallel production lines. *Comput. Chem. Eng.* **1991**, *15*, 85–103.
- (65) Sahinidis, N.; Grossmann, I. Reformulation of multi-period MINLP model for capacity expansion of chemical process. *Oper. Res.* **1992**, *40*, S127–S144.
- (66) Carrizosa, E.; Hansen, P.; Messine, F. Improving interval analysis bounds by translations. *J. Global Optimiz.* **2004**, *29* (2), 157–172.
- (67) Polisetty, P.; Gatzke, E.; Voit, E. Yield optimization of regulated metabolic systems using deterministic branch-and-reduce methods. *Biotechnol. Bioeng.* **2008**, *99* (5), 1154–1169.
- (68) Pozo, C.; Guillén-Gosálbez, G.; Sorribas, A.; Jiménez, L. Outer approximation-based algorithm for biotechnology studies in systems biology. *Comput. Chem. Eng.* **2010**, DOI: 10.1016/j.compchemeng.2010.03.001.
- (69) Savageau, M. Biochemical systems analysis. I. Some mathematical properties of the rate law for the component enzymatic reactions. *J. Theor. Biol.* **1969**, *25*, 365–369.
- (70) Savageau, M. Biochemical systems analysis. II. The steady-state solutions for an n-pool system using a power-law approximation. *J. Theor. Biol.* **1969**, *25*, 370–379.
- (71) Voit, E. *Computational Analysis of Biochemical Systems. A Practical Guide for Biochemists and Molecular Biologists*; Cambridge University Press: Cambridge, U.K., 2000.
- (72) Vecchiotti, A.; Sangbum, L.; Grossmann, I. Modeling of discrete/continuous optimization problems: characterization and formulation of disjunctions and their relaxations. *Comput. Chem. Eng.* **2003**, *27*, 433–448.
- (73) Lee, S.; Grossmann, I. Global optimization of nonlinear generalized disjunctive programming with bilinear equality constraints: applications to process networks. *Comput. Chem. Eng.* **2003**, *27*, 1557–1575.
- (74) Biegler, L.; Grossmann, I.; Margot, F.; Sahinidis, N.; Lee, J.; Waechter, A.; Belotti, P.; Castro, P.; Ruiz, J. CMU-IBM Cyber-Infrastructure for MINLP. <http://www.minlp.org/> (accessed Oct 2010).
- (75) Lu, H.-C.; Li, H.-L.; Gounaris, C.; Floudas, C. Convex relaxation for solving posynomial programs. *J. Global Optimiz.* **2010**, *46* (1), 147–154.
- (76) Bergamini, M.; Scenna, N.; Aguirre, P. Global Optimal Structures of Heat Exchanger Networks by Piecewise Relaxation. *Ind. Eng. Chem. Res.* **2007**, *46*, 1752–1763.
- (77) Bergamini, M.; Grossmann, I.; Scenna, N.; Aguirre, P. An improved piecewise outer-approximation algorithm for the global optimization of MINLP models involving concave and bilinear terms. *Comput. Chem. Eng.* **2008**, *32*, 477–493.
- (78) Karuppiah, R.; Grossmann, I. Global optimization for the synthesis of integrated water systems in chemical processes. *Comput. Chem. Eng.* **2006**, *30*, 650–673.
- (79) McCormick, G. *Nonlinear Programming, Theory, Algorithms, and Applications*; John Wiley & Sons: New York, 1983.
- (80) Murtagh, B.; Saunders, M. *MINOS 5.1 user's guide. Technical Report SOL 83-20R*, Systems Optimization Laboratory, Stanford University: Palo Alto, CA, 1987.
- (81) Schilling, G.; Pantelides, C. A simple continuous-time process scheduling formulation and a novel solution algorithm. *Comput. Chem. Eng.* **1996**, *20*, S1221–S1226.
- (82) Tawarmalani, M.; Sahinidis, N. *Convexification and global optimization in continuous and mixed-integer nonlinear programming: Theory, algorithms, software and applications; Non-convex Optimization and Its Applications*; Kluwer Academic Publishers: Dordrecht, The Netherlands, 2002; Vol. 65.
- (83) Hansen, P.; Jaumard, B.; Lu, S. An analytical approach to global optimization. *Math. Programming* **1991**, *52*, 227–254.
- (84) Quesada, I.; Grossmann, I. Global optimization algorithm for heat exchanger networks. *Ind. Eng. Chem. Res.* **1993**, *32*, 487–499.
- (85) Quesada, I.; Grossmann, I. A global optimization algorithm for linear fractional and bilinear programs. *J. Global Optimiz.* **1995**, *6*, 39–76.
- (86) Quesada, I.; Grossmann, I. Global optimization of bilinear process networks and multicomponent flows. *Comput. Chem. Eng.* **1995**, *19* (12), 1219–1242.
- (87) Liberti, L. Writing global optimization software. In *Global Optimization: From Theory to Implementation*; Liberti, L., Maculan, N., Eds.; Springer: Berlin, 2006; pp 211–262.
- (88) Shtectman, J.; Sahinidis, N. A finite algorithm for global minimization of separable concave programs. *J. Global Optimiz.* **1998**, *12*, 1–36.
- (89) Hentenryck, P. V.; Michel, L.; Deville, Y. *Numerica, a Modeling Language for Global Optimization*; MIT Press: Cambridge, MA, 1997; Vol. 65.
- (90) Shtectman, J.; Sahinidis, N. A finite algorithm for global minimization of separable concave programs. *J. Global Optimiz.* **1998**, *12*, 1–36.
- (91) Messine, F. Deterministic global optimization using interval constraint propagation techniques. *RAIRO-RO* **2004**, *38* (4), 277–294.
- (92) Sahinidis, N. V. *BARON: a general purpose global optimization software package. J. Global Optimiz.* **1996**, *8*, 201–205.



Contents lists available at [ScienceDirect](#)

Journal of Biotechnology

journal homepage: [www.elsevier.com/locate/jbiotec](http://www.elsevier.com/locate/jbiotec)



## Optimization and evolution in metabolic pathways: Global optimization techniques in Generalized Mass Action models

Albert Sorribas<sup>a,\*</sup>, Carlos Pozo<sup>b</sup>, Ester Vilaprinyo<sup>c</sup>, Gonzalo Guillén-Gosálbez<sup>b</sup>, Laureano Jiménez<sup>b</sup>, Rui Alves<sup>a</sup>

<sup>a</sup> Departament de Ciències Mèdiques Bàsiques, Institut de Recerca Biomèdica de Lleida (IRBLLEIDA), University of Lleida, Montserrat Roig, 2, 25008 Lleida, Spain

<sup>b</sup> Departament d'Enginyeria Química, Universitat Rovira i Virgili, Tarragona, Spain

<sup>c</sup> Fundation Dr. Ferran, Hospital Verge de la Cinta, Tortosa, Spain

### ARTICLE INFO

#### Article history:

Received 15 September 2009

Received in revised form 23 January 2010

Accepted 29 January 2010

#### Keywords:

Optimization

Power-law formalism

Evolution

Design principles

Operation principles

### ABSTRACT

Cells are natural factories that can adapt to changes in external conditions. Their adaptive responses to specific stress situations are a result of evolution. In theory, many alternative sets of coordinated changes in the activity of the enzymes of each pathway could allow for an appropriate adaptive readjustment of metabolism in response to stress. However, experimental and theoretical observations show that actual responses to specific changes follow fairly well defined patterns that suggest an evolutionary optimization of that response. Thus, it is important to identify functional effectiveness criteria that may explain why certain patterns of change in cellular components and activities during adaptive response have been preferably maintained over evolutionary time. Those functional effectiveness criteria define sets of physiological requirements that constrain the possible adaptive changes and lead to different *operation principles* that could explain the observed response. Understanding such operation principles can also facilitate biotechnological and metabolic engineering applications. Thus, developing methods that enable the analysis of cellular responses from the perspective of identifying operation principles may have strong theoretical and practical implications. In this paper we present one such method that was designed based on nonlinear global optimization techniques. Our methodology can be used with a special class of nonlinear kinetic models known as GMA models and it allows for a systematic characterization of the physiological requirements that may underlie the evolution of adaptive strategies.

© 2010 Elsevier B.V. All rights reserved.

### 1. Introduction

Cells are natural factories that can adapt to changes in external conditions (Causton et al., 2001; Gasch et al., 2000; Mitchell et al., 2009). Their adaptive responses are a result of evolution through different mechanisms that include random mutation, gene duplication, gene transfer, etc. (Koonin, 2009). During steady-state growth conditions, the cell works within normal operating ranges that are characterized by fluxes and metabolite levels moving within more or less narrow ranges (Watson, 1970; Wiebe et al., 2008).

As the conditions in the medium change, the operating range of cells may also change. If environmental changes are spurious, there are internal control mechanisms that play a fundamental role in maintaining the operating range of cells about its initial value. However, when the environmental changes are relevant or sustained, an adaptive response is mounted by the cells. Such adaptive responses occur during heat shock, oxidative stress, or other

stresses. If those situations are prevalent in the evolutionary history of the cell, specific behaviors and mechanisms that facilitate cell adaptation through changes in gene expression and protein activity and assure cell viability are selected for. Such behaviors lead to a fine tuning of metabolic fluxes and concentrations (Vilaprinyo et al., 2006). The specificity of the adaptive response mounted by each cell type in response to a given stress depends both on the challenges it responds to and on the evolutionary history of the cell or organism (Bedford and Hartl, 2009; Kashiwagi et al., 2006; Teusink et al., 2009; Wilkins, 2007).

For example, the heat shock caused by a sudden rise in the temperature of the growing media triggers an ordered response in yeast that causes an arrest in cell cycle and specific changes in the coordinated activity of several metabolic pathways (Trotter et al., 2001). These changes help the cell to synthesize protective molecules that permit its adaptation and survival (Causton et al., 2001; Eisen et al., 1998; Gasch et al., 2000; Jenkins, 2003). In principle, many alternative sets of coordinated changes in the activity of pathways could allow for an appropriate adaptive readjustment of metabolism. However, experimental and theoretical measurements of the actual responses show that these follow fairly

\* Corresponding author. Tel.: +34 608533249; fax: +973702426.

E-mail address: [albert.sorribas@cmb.udl.cat](mailto:albert.sorribas@cmb.udl.cat) (A. Sorribas).

well defined patterns that are consistent with an evolutionary optimization of this response with respect to different physiological and functional effectiveness criteria (El-Samad et al., 2005; Kurata et al., 2006; Molina-Navarro et al., 2008; Vilaprinyo et al., 2006). Thus, it is important to identify functional effectiveness criteria that may explain why certain patterns of change in cellular components and activities during adaptive response have been preferably maintained over evolutionary time (Coelho et al., 2009; Han, 2008; Salvador and Savageau, 2003, 2006; Savageau, 1971, 1974a,b, 1976; Savageau et al., 2009). Such criteria are necessarily derived from the analysis of systemic properties that emerge from the integrated molecular behavior of the cellular components, and they may include robustness, dynamic stability, minimization of intermediates, minimization of biosynthetic cost, temporal responsiveness, etc. (Chang and Sahinidis, 2005; Coelho et al., 2009; Salvador and Savageau, 2003, 2006; Savageau et al., 2009). The functional effectiveness criteria for a response define sets of physiological constraints that shape that response and lead to different *operation principles* that could explain why the cells adapt in a certain way at the molecular level (Bedford and Hartl, 2009; Braunstein et al., 2008; Vilaprinyo et al., 2006; Voit and Radivoyevitch, 2000; Voit, 2003).

Although the operational principles of cellular responses are a result of evolution, they can be applied to and validated in biotechnological applications. Metabolic engineering manipulates naturally evolved organisms in order to obtain increased amount of new products (Bailey et al., 1990, 1996; Bailey, 1991, 1999, 2001; Hatzimanikatis and Liao, 2002). This manipulation often involves a process of optimization that searches for the best modified strain with respect to the initial optimization criteria (Goodman, 2008). Thus, developing methods that permit analysis of cellular responses from the perspective of identifying operational principles may have strong theoretical and practical implications. Often, this goal can only be achieved through methods that involved the creation, analysis and comparison of mathematical models for the processes and responses one is interested in studying (Alvarez-Vasquez et al., 2004; Klipp et al., 2005; Sims et al., 2004; Voit, 2003).

In this work, we discuss and extend a method that can be used to identify and study the operation principles of cellular response at the molecular level, by characterizing feasibility regions for those responses (Guillén-Gosálbez and Sorribas, 2009; Pozo et al., submitted for publication). Such feasibility regions encompass all possible ranges in enzyme activity that allow for an appropriate response by the cell after an environmental challenge. This method may help in both, understanding the evolution of such responses and guiding manipulations of gene expression in metabolic engineering applications. The proposed method for identifying feasibility regions uses a recently developed global optimization method (Guillén-Gosálbez and Sorribas, 2009; Pozo et al., submitted for publication). Here, the capabilities of that optimization method are enhanced through an iterative and systematic search strategy that identifies all the parameter regions containing admissible solutions that are compatible with the considered physiological constraints. The general framework presented here has the potential for solving problems of great interest in systems biology studies. As an example we analyze a mathematical model created to represent the heat shock response of the yeast *Saccharomyces cerevisiae*.

## 2. Methods

### 2.1. Generalized Mass Action models

Generalized Mass Action (GMA) models are a special class of models defined within the general framework of Biochemical Systems Theory (BST) (Voit, 2000). These models use the power-law

formalism to obtain a representation of the different processes involved in the target system. For a network with  $p$  processes (enzyme reactions, transport systems, etc.),  $n$  internal metabolites, and  $m$  external parameters or independent variables, a GMA model is defined as follows:

$$\frac{dX_i}{dt} = \sum_{r=1}^p \mu_{ir} v_r = \sum_{r=1}^p \mu_{ir} \left( \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \right) \quad i = 1, \dots, n \quad (1)$$

In Eq. (1), the parameters  $\mu_{ir}$  account for the stoichiometry of the process, i.e. the number of molecules of  $X_i$  produced by or used in reaction  $v_r$  (for instance +1, +2 for production, or -1, -2, etc., for degradation). The parameters of the power-law representation of each reaction are the apparent rate-constant  $\gamma_r$  and the kinetic-order  $f_{rj}$ , defined as (Savageau, 1969a,b, 1976):

$$f_{ir} = \left( \frac{\partial v_r}{\partial X_i} \right)_0 \frac{X_{i0}}{v_{r0}} \quad (2)$$

The subscript 0 stands for the operating point where the power-law representation is derived. Appropriate parameter values for a given system can be estimated using different procedures. As this is a broad subject, the reader is referred to the recent review by Chou and Voit (2009). In the following, we shall assume that a parameter set has been obtained and that the GMA model can be used for characterizing the properties of the system.

GMA representations integrate information about network stoichiometry and regulation (kinetic-orders) into a dynamic mathematical model. These models can be used for computing both the transient and steady-state responses of metabolites and fluxes to changes in the environment of the model. Due to their structure and to the available methods, GMA models are well suited for evaluating parameter sensitivities and for developing optimization techniques (Chang and Sahinidis, 2005; Marin-Sanguino et al., 2007; Polisetty et al., 2008; Torres et al., 1996, 1997; Voit, 1992). Thus, that representation is especially useful as a framework for systems biology applications and provides a description of processes that is more accurate than the one provided by other techniques based on the stoichiometric matrix alone, such as Flux Balance Analysis (FBA) (Lee et al., 2006). This added accurateness comes at the price of needing more information to estimate parameter values for GMA models.

### 2.2. Characterization of the effect of changes in enzyme activities

Given a GMA model, changes in enzyme activities can be implemented by changing the value of the rate-constant for the processes in which the enzymes are involved.<sup>1</sup> For simplicity, we can write

$$\frac{dX_i}{dt} = \sum_{r=1}^p \mu_{ir} v_r = \sum_{r=1}^p \mu_{ir} \left( k_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \right) \quad i = 1, \dots, n \quad (3)$$

where  $k_r$  indicates the change-fold over the original enzyme activity (which is  $\gamma_r$ ). Thus, in the reference state,  $k_r = 1$ . Accordingly, a vector  $(k_1, k_2, \dots, k_p)$  would correspond to a specific pattern of fold changes in enzyme activities. For this vector, the change in the

<sup>1</sup> Enzyme activities can be explicitly included in the model as independent variables. However, for constant levels of enzyme activity, doing so is equivalent to changing the rate-constant directly. If the model includes gene regulation and modulatory changes in protein activity, enzymes should be explicitly included as internal variables in the model. Mimicking changes in the medium can be done either by changing the values of an external variable or by changing the values of rate constants for the processes that are responsible for sensing those changes.

system steady-state can be easily computed numerically by solving the steady-state equation

$$0 = \sum_{r=1}^p \mu_{ir} \left( k_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \right) \quad i = 1, \dots, n \quad (4)$$

We shall use Eqs. (3) and (4) to analyze the effect of different activity patterns on the systemic performance of the model and evaluate how this performance influences the overall physiological outcome of the response.

### 2.3. Criteria for functional effectiveness in cellular metabolism

Changes in the reference steady-state as a consequence of a change in the enzyme activity pattern can be compared to a series of functional effectiveness criteria (Vilaprinyo et al., 2006; Voit and Radivoyevitch, 2000). Those criteria, which define the boundaries of internal change that the cell must go through in order to adapt and survive, are matched against the internal changes that are caused by the changes in enzyme activity of a given response profile. While some of those criteria may be quite general, others may be case-specific and may have different quantitative thresholds in different cases (Salvador and Savageau, 2003, 2006; Vilaprinyo et al., 2006). We now briefly discuss some of the criteria that have been used in the literature. These are useful for discussing operative ranges, evolution, and optimization of metabolic processes.

#### 2.3.1. Change in metabolic fluxes

Changes in the rate of synthesis for key metabolites are important indices of functional effectiveness. For example, if a system regulates production of a metabolite in response to the cellular demand for that metabolite, an increase in the demand should lead to an increase in the production (e.g. Alves and Savageau, 2000 and references therein). The specific flux criteria are dependent on the system one is interested in and should be considered in the context of the whole system and not as isolated processes. For example, in the adaptive response to heat shock an increase in ATP production that causes depletion of NADPH or a dramatic decrease in glycolytic flux may be inappropriate in the general context of the adaptive response (Vilaprinyo et al., 2006).

In GMA models, steady-state fluxes can be easily computed for each condition using the following equation

$$v_{r_{ss}} = k_r \gamma_r \prod_{j=1}^{n+m} X_{j_{ss}}^{f_{rj}} \quad r = 1, \dots, p \quad (5)$$

where subscript *ss* indicates the relevant steady-state values corresponding to the new conditions. As stated before, the steady-state solution for metabolites can be obtained by numerically solving Eq. (4).

In larger networks that involve different branch points and regulatory effects, it is possible to obtain similar increases in a given set of fluxes with different patterns of modified enzyme activities. Thus, this criterion, by itself, will seldom be enough to assess the adaptive value of a set of changes and fully explain the observed operation principles for the system.

#### 2.3.2. Metabolite accumulation

Changes in steady-state fluxes may often lead to changes in metabolite levels. From a practical point of view, either in biotechnological applications or in natural systems, one may argue that accumulation of intermediary metabolites may cause undesirable cross regulation side effects and tax the finite solvability capacity of the cell (see Alves and Savageau, 2000 and references therein). Thus, minimization of intermediate metabolite accumulation will

be typically regarded as an important effectiveness criterion of an adaptive response, except for those cases in which metabolite accumulation might play an important role (for instance accumulation of trehalose in the heat shock response). Changes in metabolite levels are given by the steady-state solution to Eq. (4).

#### 2.3.3. Overall changes in enzyme activities

Changes in enzyme activity are easy to simulate. However, it is often difficult to assess in a real situation whether those changes are indirect and due to the modulation of either gene expression or stability of mRNA (Garcia-Martinez et al., 2007; Romero-Santacreu et al., 2009), or direct and due to modulator effects on the activity of the protein. The later can arise via reversible covalent modification of specific residues or via changes in the conformation of the protein in response to a new set of physical chemical parameters in the medium. Changes in gene expression are costly in terms of metabolic resources (Wagner, 2005). They lead to mRNA and protein synthesis, which are among the most expensive metabolic activities of a cell. Thus, minimization of fold change can be considered an important functional effectiveness criterion (Raiford et al., 2008). If one assumes that changes in protein activity during the long term adaptive response of a cell are mostly due to changes in gene expression then, to a first approximation, one can estimate the cost of a given set of changes in enzyme activity by adding up all the  $k_r$  values. One possible way to account for both up- and down-regulations consists of defining a “biological” cost of a response that is mathematically given by  $\sum v_i |\ln(k_i)|$ .

#### 2.3.4. Parameter robustness

Parameter robustness is an important criterion as it refers to the system's sensitivity to slight differences in parameter values (Aldana et al., 2007; Coelho et al., 2009; Kitano, 2004; Kitano, 2007; Morohashi et al., 2002; Savageau, 1971). Systems with large parameter sensitivities may indicate the existence of processes that are more responsive to noise. Thus, they could be considered as less well adapted than systems that are more sensitive to parameter changes. Although low parameter sensitivities may arise from poorly identified parameters, one can argue that, in most cases, low sensitivity is a desirable property in well-adapted systems. This criterion has been extensively used in identifying design principles and in evaluating model adequacy and behavior (Cascante et al., 1995; Coelho et al., 2009; Curto et al., 1997; de Atauri et al., 2000; Voit, 2000).

#### 2.3.5. Temporal responsiveness

Temporal responsiveness is another criterion that is important for systemic performance. Systems with inadequate temporal responsiveness may not survive to reach a new steady state, independently of the adequacy of their steady-state responses. In general, evaluating this criterion requires numerical simulations, except for the case where we are only interested in studying the dynamics in the neighborhood of a steady-state solution. In such a case one can linearize the system of equations about the steady state and obtain analytical solutions for the transient behavior of the dependent variables (Hlavacek and Savageau, 1998).

Unlike the other criteria that were discussed so far, using temporal responsiveness as a criterion for optimization poses many problems. In the context of globally optimizing metabolic systems, there are indeed very few methods capable of handling the dynamic constraints required to assess the temporal responsiveness. In fact, the strategies proposed so far are only applicable to specific types of models, and usually optimization uses large amounts of CPU time, even when tackling small problems with few variables and constraints (Chachuat et al., 2006; Chang and Sahinidis, 2005; Esposito and Floudas, 2000; Papamichail and Adjiman, 2002, 2004; Singer and Barton, 2006). This limitation can be overcome by performing

the assessment of the temporal responsiveness in the post-optimal analysis of the solutions found. Hence, once a feasible solution is identified, the evaluation of its temporal responsiveness can add an extra criterion for deciding the relevance of such solution. In terms of evolution, this may be important as a given optimum can involve dynamic properties that will make the solution unfeasible in practice.

### 2.3.6. Steady-state stability

Dynamic stability is a criterion that evaluates the ability of a given system for returning to a steady-state after a perturbation. A stable system can accommodate fluctuations and will be able of maintain a reference state. Evaluation of steady-state stability should be a complementary criterion for testing the appropriateness of a proposed change in the system (Savageau, 1974a, 1975, 1998). In the optimization of metabolic systems, this criterion can be included in the optimization model itself (Chang and Sahinidis, 2005) or it can be assessed in the post-optimal analysis of the solutions for Eq. (4).

## 3. Feasibility regions in biochemical pathways: definition and their practical significance

### 3.1. Definitions

A feasibility region is a region in parameter space whose internal membership is defined by the sets of all parameter values that are compatible with specific physiological constraints (Dayarian et al., 2009; Guillén-Gosálbez and Sorribas, 2009). Here, without loss of generality, we shall concentrate on the special case of feasibility regions defined by changes in enzyme activities, that is the set of vectors representing the fold change in enzyme activities:  $(k_1, k_2, \dots, k_p)$ , that are compatible with a set of functional effectiveness criteria (constraints). These functional criteria must be assessed through mathematical models, such as the GMA representation, that allow predicting the biological performance of a system in a specific environment. In mathematical terms, performing a feasibility analysis entails conducting a systematic search for determining the set of values of some variables of the biological model for which the overall formulation remains feasible. In this context, linear models usually fail to capture the whole complexity of the biological system, so it is necessary to use nonlinear formulations.

Hence, finding the boundaries for this class of feasibility regions requires obtaining global optimal solutions for nonlinear optimization problems. One of the important limitations of standard nonlinear optimization techniques is that they cannot guarantee the global optimality of the solutions found when they are applied to nonlinear problems that have non-convexities. Non-convexities, such as bilinear terms, fractional terms, etc., are very common in many engineering problems. In the context of our analysis, these non-convexities arise from the kinetic equations required to link the concentration of the metabolites with the velocities of the reactions that take place in the metabolic network.

There are currently several global optimization methods that can handle non-convex problems and provide solutions that are globally optimal within a desired tolerance (Tawarmalani and Sahinidis, 2002). Most of these methods are general purpose, that is to say, they can be applied to a wide range of problems regardless of the type of non-convexities embedded in the model. However, their performance can change drastically from one application to another depending on the specific structure of the problem to be solved (for a detailed review of these methods see Grossman and Biegler, 2004). A possible way of expediting the search for global solutions for nonlinear non-convex problems consists of exploiting

the structure of the involved non-convexities. The major classes of non-convex problems studied so far include concave minimization (Hansen et al., 1992) and problems with linear fractional and bilinear terms (Quesada and Grossman, 1995), and a method for problems with signomial parts (Porn et al., 2008). Different optimization strategies have also been suggested for S-system and GMA models within BST (Chang and Sahinidis, 2005; Hatzimanikatis et al., 1996; Marin-Sanguino et al., 2007; Polisetty et al., 2008; Voit, 1992). Recently, a highly efficient global optimization technique for GMA models has been developed by our group. Technical aspects of this optimization are discussed elsewhere (Guillén-Gosálbez and Sorribas, 2009; Pozo et al., submitted for publication). We shall use this technique in the feasibility analysis presented here.

### 3.2. Characterization of feasibility regions in GMA models

The method for finding the feasibility regions was first introduced by Guillén-Gosálbez and Sorribas (2009). Here, we briefly review it and discuss the different steps and their importance. Mainly, steps 2–3 are critical for reducing the search space and obtain a useful result. After reviewing the method, we shall apply it to two practical cases showing its utility both for optimization and evolutionary studies. Finally, we shall stress the role of the set of physiological constraints in defining the feasibility region. As stated before, the strategy presented relies on the use of global optimization methods that are customized for this particular application.

A feasibility region for a particular problem can be identified through the following steps:

1. Define a set of constraints that must be fulfilled by any solution (limits for fluxes, concentrations, gene expression, etc.). At this point, collaboration with experts in the biological problem is fundamental.
2. Define the search space for the fold change of each enzyme. (i.e., the lower and upper bounds,  $k_i^{LO}$  and  $k_i^{UP}$  that define the interval within which the fold change must fall). Based on experimental information, one can restrict the search space for practical purposes. Thus, if measurements in microarray experiments show that during the studied response a given gene is over expressed between 5 and 8-fold, we could consider allowing changes from 1 up to 20-fold from the basal condition for that gene. By making the range so large, one covers for other plausible values that may also be linked to alternative adaptive solutions.
3. Find the maximum and minimum (*bound contraction*) values for changing each enzyme that are compatible with the set of constraints defined in 1 and with the limits established in 2. This is achieved by defining  $k_i$  as the objective function. Note that these optimizations provide bounds for all the variables  $k_i$  (i.e.,  $k_i^{LO*}$  and  $k_i^{UP*}$ ) that will fall within those first considered in step 2. Thus, we can assure that outside the obtained bounds for  $k_i$  no other combination of changes in the enzymes produces a valid solution. In mathematical terms, we have that a solution will be unfeasible if there exists at least one  $k_i$  that satisfies that  $k_i \notin [k_i^{LO*}, k_i^{UP*}]$ .
4. Define a grid of values for each  $k_i$ :  $(k_i^1, k_i^2, \dots, k_i^{n_i})$ , using the minimum and maximum values obtained in the previous step. Typically, we have divided the allowed range in 10 sections. The bound contraction step shortens the search region, which results in a more efficient search of the feasibility region.
5. Consider a set of the hyper-rectangles, each of which is defined by lower and upper limits imposed on the values of each  $k_i$ . For instance, a particular hyper-rectangle would be defined by fold changes that are between 2.5 and 3.7 for enzyme 1, between 10 and 12.3 for enzyme 2, and so on.

6. Find the global optimum using any of the velocities as the objective function. This will give a set of fold changes ( $k_{11}, k_{21}, \dots, k_{p1}$ ) for which this optimum is attained. At this stage, as the goal is to find admissible solutions, we can select any of the velocities as the objective function. The results of the feasibility analysis will be the same independently of this choice. At this point, all we need is to be sure that at least a solution exist that is compatible with the set of constraints. Hence, it is not strictly required to globally optimize the model in each hyper-rectangle, since a feasible solution suffices for the purpose of the analysis. One can wrongly conclude from this observation that it is possible to conduct the feasibility analysis using a local optimization method. This is not true, as the task of the algorithm is not only to identify feasible solutions in each hyper-rectangle, but also to discard regions in which no feasible point exists. Standard local optimization methods cannot accomplish the latter task, as they can fail even in solving convex problems (Tawarmalani and Sahinidis, 2002), in such a way that one will never be sure if the convergence problems that will arise when attempting to optimize empty hyper-rectangles will really indicate the absence of feasible solutions. One possible way to circumvent this issue is to rely on a lower bounding problem, which is one of the main ingredients of any global optimization technique, capable of providing a valid lower bound on the global optimum of the model. Particularly, the feasibility analysis requires the use of a linear lower bounding problem, since linear and mixed-integer linear programming techniques (LP and MILP, respectively) can indeed identify problems with no feasible solutions. Hence, our method exploits the fact that LP and MILP techniques will only fail when attempting to solve models of small/medium size that are really unfeasible (i.e., do not contain any feasible point). Note that in this context the main task of the lower bounding problem is not to provide a tight bound on the global solution of the model, as is the case in standard global optimization methods, but to detect empty hyper-rectangles from unfeasible models. Thus, the particular features of our feasibility analysis justify the need for a customized global optimization method.

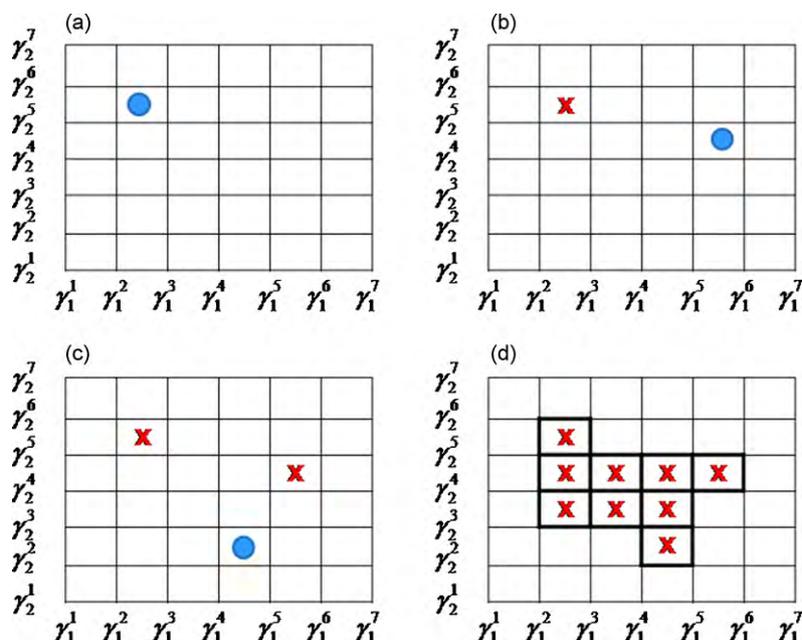
7. Identify and annotate the hyper-rectangle that contains this solution. This will be the one whose lower and upper limits contain the values of  $k_i$  associated with the optimal solution identified in the current iteration.
8. Repeat steps 4–7 by excluding the hyper-rectangle containing the optimal solution obtained in step 5 by adding an integer cut to the lower bounding problem. This is repeated until no further solution is found to be compatible with the remaining hyper-rectangles (i.e., until the lower bounding problem turns out to be unfeasible).
9. Analyze the obtained results and compare the feasible region with actual experimental data. If the feasibility region contains the observed data, this is an indication that the considered set of constraints may explain the adaptive response. Alternative constraints can be introduced and a new feasibility region can be obtained by starting again the analysis at point 2. In the next section we will discuss the interpretation of results obtained with different sets of constraints.

This procedure is illustrated in Fig. 1. For simplicity, we show results for two enzymes only. However, at each optimization, all the enzymes are allowed to change values (see below). The constraints considered in each optimization are those defined in step 1 plus the limits on the values of  $k_i$  that define each hyper-rectangle. By following the procedure described above, a region of feasibility is bounded and defined by a set of feasible hyper-rectangles. These regions can be further refined by increasing the granularity of the hyper-rectangles within the region(s) of feasibility. Note that any of the hyper-rectangles that define the feasible regions contains at least one admissible solution.

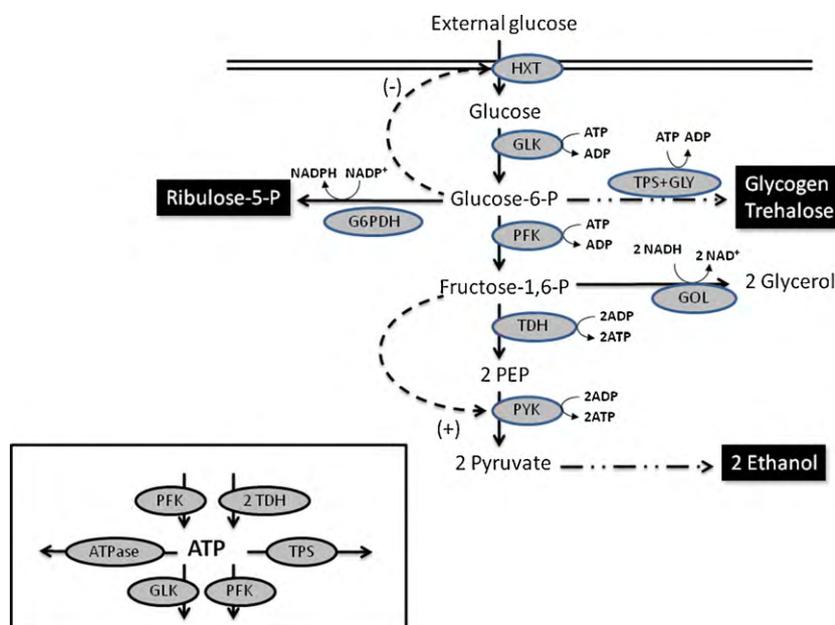
### 3.3. Utility of feasibility regions characterization

The feasibility regions determined through the previously described methodology can be most useful in two situations:

1. *In biotechnological applications.* In practical applications, it may be impossible to attain the optimal solution that would



**Fig. 1.** Strategy for finding a feasible region. In the first step, the global optimum is identified and the hyper-rectangle where it occurs is annotated (a). In the next step, this hyper-rectangle is discarded and the new optimum is located (b). The process is repeated (c) until no new optimum is obtained (d). Here we show results for two of the hypothetical enzymes but the search is done for all simultaneously.



**Fig. 2.** Schematic representation of the central metabolism of yeast. Details are discussed in previous papers (Polisetty et al., 2008; Vilaprinyo et al., 2006; Voit and Radivoyevitch, 2000). Basal values for enzyme activities and the resulting steady state are given in Tables 1 and 2.

correspond to an optimization analysis. Thus, one can define a minimum percentage of the optimum that would identify a practical cost-beneficial strategy. Feasible regions that are compatible with that threshold can be obtained using the method proposed in this work. The desired minimal limit can be mathematically represented by a simple inequality constraint. Once the feasibility region is determined the user can select the most appropriate values for practical implementation. The feasible regions will contain a global optimum for attaining this practical threshold as well as many other suboptimal solutions. Note that all the identified solutions, including the sub-optimal ones, would be guaranteed to attain the minimum increase in the objective function considered in the analysis.

2. *In evolutionary studies.* Feasible regions that are compatible with physiological requirements can be identified in studies about evolution of responses. If the model captures the features of the system that are important for the response, one would expect to find the actual adaptive response within this region. An iterative analysis considering different physiological constraints may help in identifying which of these constraints are more important as selective pressures for evolving an appropriate response, avoiding the spandrel effects. Furthermore, comparison of actual data with optimal solutions can help in understanding the selective pressures in a given case.

### 3.4. Examples

#### 3.4.1. Metabolic model

As an example for showing the applicability of the method described above, we shall consider a simplified conceptual model for the basal metabolism of yeast that is derived from previous models of the same pathways (Curto et al., 1995; Polisetty et al., 2008; Voit and Radivoyevitch, 2000).

This model, summarized in Fig. 2, accounts for glycolysis, the synthesis of glycogen and trehalose, the branching from fructose-1,6-P to glycerol, and the branching from glycolysis to the pentose phosphate metabolism. For convenience, we consider simplified reactions by lumping together a number of processes. For example, we consider an aggregated process leading to trehalose and glycogen. Numerically, we shall consider that the flux into trehalose is

a fraction of the total flux for this branch (Vilaprinyo et al., 2006; Voit and Radivoyevitch, 2000). For more details on the simplifications, assumptions, and experimental evidences used to build this model the reader is referred to the paper by Voit and Radivoyevitch (2000). The different processes are modeled using the power-law formalism as:

Process	Velocity	Power-law representation	Steady-state rate
HXT	$v_1$	$0.9023X_2^{-0.2344}X_6$	17.73
GLK	$v_2$	$3.1847X_1^{0.7464}X_5^{0.0253}X_7$	17.73
PFK	$v_3$	$0.5232X_2^{0.7318}X_5^{-0.3941}X_8$	15.946
TDH	$v_4$	$0.011X_3^{0.6159}X_5^{0.1308}X_9X_{14}^{-0.6088}$	15.06
PYK	$v_5$	$0.0947X_3^{0.05}X_4^{0.533}X_5^{-0.0822}X_{10}$	30.00
TPS + GLY	$v_6$	$0.0009X_2^{0.7318}X_{11}$	0.014
G6PDH	$v_7$	$1.76898X_2^{0.0526}X_{15}^{0.9646}$	1.77
GOL	$v_8$	$0.103209X_3^{0.05}X_4^{0.533}X_5^{-0.0822}X_{12}$	1.772
ATPase	$v_9$	$0.937905X_5^1X_{13}$	26.55

The stoichiometric matrix corresponding to the model in Fig. 2 is given by:

$$N = \begin{pmatrix} -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & -1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 & 0 & 0 & -1/2 & 0 \\ 0 & 0 & 0 & 2 & -1 & 0 & 0 & 0 & 0 \\ 0 & -1 & -1 & 2 & 1 & -1 & 0 & 0 & -1 \end{pmatrix} \quad (7)$$

Multiplying the stoichiometric matrix by the vector of velocities  $\mathbf{V} = (v_1, \dots, v_9)$ , we would obtain the set of differential equations for the model in GMA form<sup>2</sup>:

$$\dot{\mathbf{X}} = \mathbf{N} \cdot \mathbf{V} \quad (8)$$

<sup>2</sup> Here we use the notation  $\dot{X}_i = dX_i/dt$ .

**Table 1**  
 Basal enzyme activities.

Symbol	Name	Value
$X_6$	Glucose uptake (HXT)	19.7 mM min <sup>-1</sup>
$X_7$	Hexokinase (GLK)	68.5 mM min <sup>-1</sup>
$X_8$	Phosphofruktokinase (PFK)	31.7 mM min <sup>-1</sup>
$X_9$	Glyceraldehyde-3-phosphate dehydrogenase (GAPD or, as alternative name, TDH)	49.9 mM min <sup>-1</sup>
$X_{10}$	Pyruvate kinase (PYK)	3440 mM min <sup>-1</sup>
$X_{11}$	Polysaccharide production (glycogen + trehalose)	14.31 mM min <sup>-1</sup>
$X_{12}$	Glycerol production (GOL)	203 mM min <sup>-1</sup>
$X_{13}$	ATPase	25.1 mM min <sup>-1</sup>
$X_{14}$	NAD <sup>+</sup> /NADH ratio	0.042

The complete mathematical model is given by:

$$\begin{aligned}
 \dot{X}_1 &= 0.9023X_2^{-0.2344}X_6 - 3.1847X_1^{0.7464}X_5^{0.0253}X_7 \\
 \dot{X}_2 &= 3.1847X_1^{0.7464}X_5^{0.0253}X_7 - 0.5232X_2^{0.7318}X_5^{-0.3941}X_8 - 0.0009X_2^{0.7318}X_{11} - 1.76898X_2^{0.0526}X_{15}^{0.9646} \\
 \dot{X}_3 &= 0.5232X_2^{0.7318}X_5^{-0.3941}X_8 - 0.011X_3^{0.6159}X_5^{0.1308}X_9X_{14}^{-0.6088} - 0.0516X_3^{0.05}X_4^{0.533}X_5^{-0.0822}X_{12} \\
 \dot{X}_4 &= 2 \times (0.011X_3^{0.6159}X_5^{0.1308}X_9X_{14}^{-0.6088}) - 0.0947X_3^{0.05}X_4^{0.533}X_5^{-0.0822}X_{10} \\
 \dot{X}_5 &= 2 \times (0.011X_3^{0.6159}X_5^{0.1308}X_9X_{14}^{-0.6088}) + 0.0947X_3^{0.05}X_4^{0.533}X_5^{-0.0822}X_{10} - 3.1847X_1^{0.7464}X_5^{0.0243}X_7 - \\
 &0.0009X_2^{0.7318}X_{11} - 0.5232X_2^{0.7318}X_5^{-0.3941}X_8 - 0.937905X_5^1X_{13}
 \end{aligned} \tag{9}$$

The basal enzyme activities used in the models are shown in Table 1. The steady-state calculated from these values and the model parameters given in Eq. (9) is shown in Table 2.

### 3.4.2. Feasible regions for a significant increase in ethanol production

Optimization of cellular processes is an important goal in biotechnology. However, optimal solutions obtained with a model will seldom be practically realizable. In most cases, significant increases in flux would imply modifying many enzymes at the same

**Table 2**  
 Steady-state values of the considered model at the basal conditions.

Symbol	Name	Basal concentration (mM)
$X_1$	Internal glucose	0.0345
$X_2$	Glucose-6-phosphate	1.011
$X_3$	Fructose-1,6-diphosphate	9.144
$X_4$	Phosphoenolpyruvate (PEP)	0.0095
$X_5$	ATP	1.1278

time, which can be unpractical. One possible application of the feasibility method proposed here is to explore possible changes in enzyme activity leading to acceptable solutions, say a given percentage over the basal value or a given percentage below the optimum value. As an example, we will consider the optimization of ethanol production using the reference model. Because of the

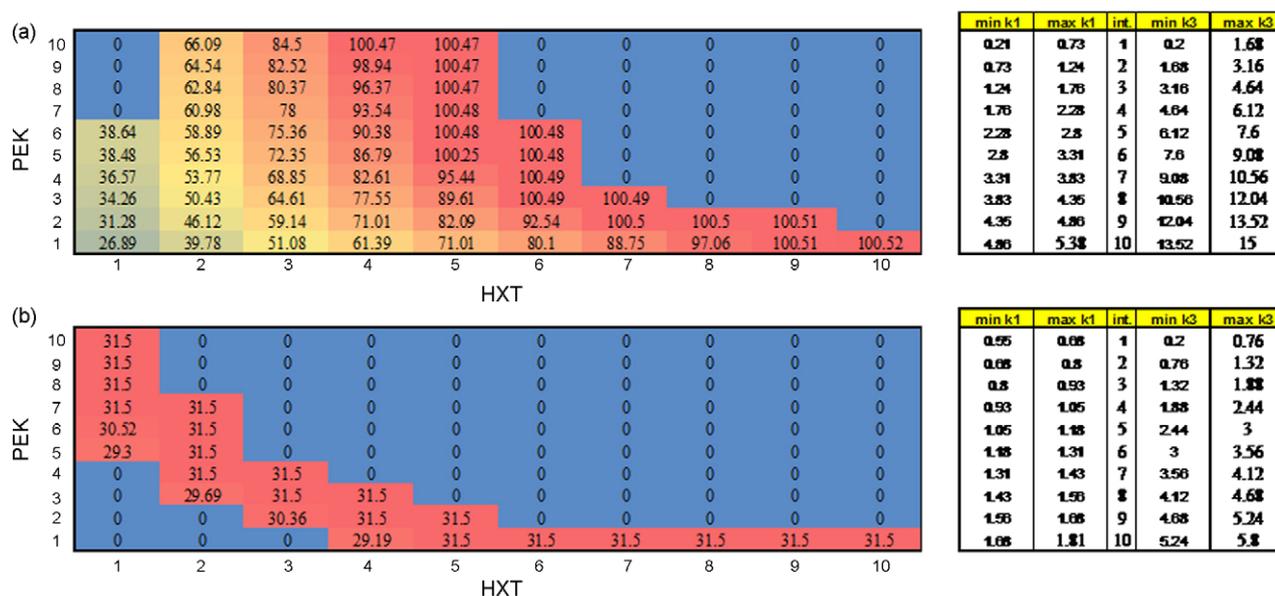
simplifications introduced in the model, the rate of synthesis of ethanol is the same as that for the synthesis of pyruvate.

First, as a reference for comparison, we explore the maximum rate of ethanol production that can be achieved if changes are allowed in all enzyme activities. The results of the optimization analysis using the method described elsewhere (Guillén-Gosálbez and Sorribas, 2009) are shown in Table 3. For comparative purposes, we obtain the optimal solution with different allowed ranges for enzyme activity changes. A nearly linear increase in ethanol production is achieved as we allow higher increases in the enzymes.

**Table 3**  
 Maximization of ethanol production.

	No constraints Maximum fold change in any enzyme				VNADPH (5% maximum variation) Maximum fold change in any enzyme				VNADPH, VATP (5% maximum variation) Maximum fold change in any enzyme			
	5	10	15	20	5	10	15	20	5	10	15	20
Fold change values at the different optimum												
HXT	5	10	15	20	5	10	15	20	1.01	1.01	1.01	1.01
GLK	1.16	10	10.56	16.03	1.16	2.16	7.6	20	0.27	0.27	0.27	0.27
PFK	5	10	15	20	5	10	15	20	1.14	1.14	1.14	1.14
TDH	5	10	15	5.27	5	9.99	15	20	2.58	2.58	2.58	2.58
PYK	5	10	15	20	5	10	15	20	0.34	0.34	0.34	0.34
TPS	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2
G6PDH	0.2	0.2	0.2	0.2	0.94	0.94	0.94	0.94	0.95	0.95	0.95	0.95
GOL	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.22	0.22	0.22	0.22
ATPase	5	10	15	20	5	10	15	20	0.94	0.94	0.94	0.94
Steady-state values for metabolites (mM) corresponding to the different optimum												
Glu	0.23	0.03	0.05	0.04	0.23	0.25	0.08	0.03	0.21	0.21	0.21	0.21
Glu-6-P	1.23	1.24	1.24	1.24	1.21	1.22	1.23	1.23	0.95	0.95	0.95	0.95
F-1,6-P	10.41	10.43	10.44	91.01	10.29	10.38	10.40	10.42	2.05	2.05	2.05	2.05
PEP	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.09	0.09	0.09	0.09
ATP	1.42	1.43	1.43	1.43	1.40	1.41	1.42	1.42	1.29	1.29	1.29	1.29
Steady-state values for fluxes (mM min <sup>-1</sup> ) corresponding to the different optimum												
$V_{ATP}$	336.1	673.8	1011.5	1348.9	332.8	670.2	1008.2	1345.9	63.4	63.4	63.4	63.4
$V_{TRE}$	0.00027	0.00027	0.00028	0.00028	0.00027	0.00027	0.00027	0.00027	0.00022	0.00022	0.00022	0.00022
$V_{NADPH}$	0.36	0.36	0.36	0.36	1.68	1.68	1.68	1.68	1.68	1.68	1.68	1.68
$V_{GLY}$	0.43	0.43	0.43	0.43	0.43	0.43	0.43	0.43	1.33	1.33	1.33	1.33
$V_{ETHANOL}$	168.0	336.9	505.8	674.6	166.3	335.2	504.1	672.9	31.9	31.9	31.9	31.9

The different scenarios are defined by allowing a maximum fold change increase for any of the enzymes of 5, 10, 15, and 20-fold. Optimal enzyme patterns are obtained without any other restriction (left) and with a maximum allowable change in the rate of NADPH synthesis of 5% about its basal value (center) and a maximum allowable change in the rate of NADPH and ATP synthesis of 5% about its basal value (right). Steady-state values of metabolites and relevant fluxes resulting from the optimal change are also shown for comparison.



**Fig. 3.** Feasibility analysis of the maximum ethanol production when only HXT ( $k_1$ ) and PFK ( $k_3$ ) are allowed to change. Values inside the left tables indicates optimum ethanol production within each cell. Cells are defined by the values of  $k_1$  and  $k_3$  indicated in the right tables. In each case, the minimum and maximum value defining each cell are shown in those tables. Color code shows the decreasing ethanol production that can be attained in different conditions. Blue color and 0 values indicate unfeasible combinations. (a) Optimization constrained to prevent an accumulation of intermediary metabolites that is over 10 times their basal value. (b) Optimization constrained to prevent an accumulation of intermediary metabolites that is larger than 5% about the basal value. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of the article.)

When no restriction is considered (Table 3 left), with a 20-fold change one can reach a velocity of  $674.6 \text{ mM min}^{-1}$  that is much higher than the basal value of  $30.0 \text{ mM min}^{-1}$ . While HXT, PFK, PYK, and ATPase should increase 20-fold, GLK and TDH require a lower increase. In all four scenarios, optimization of ethanol production should require lowering TPS, G6PDH, and GOL activities. Imposing limits on the changes of NADPH production leads to a similar result, but now the activity of G6PDH is almost unchanged (Table 3 center). It is important to stress that in all the cases the synthesis of ATP increases by a large amount, from a basal value of  $46.07 \text{ mM min}^{-1}$  to  $1348.9 \text{ mM min}^{-1}$  when the maximum fold change allowed is 20. If we restrict both the increase in NADPH and ATP production, then the maximum attainable ethanol production drastically decreases (Table 3 right).

In the previous examples, all the enzymes were allowed to change. A more realistic approach that could be translated into wet lab experiments should analyze the practical possibilities of increasing ethanol production when only a small number of enzymes are changed. For illustrative purposes, based on the previous analyses of this problem (Guillén-Gosálbez and Sorribas, 2009; Polisetty et al., 2008; Vilaprinyo et al., 2006), we select HXT and PFK. We shall maintain all the other enzymes fixed at their basal activity. As we are looking for solutions that do not compromise cell viability, we shall enforce the condition that all the internal metabolites should not change more than 10-fold from their basal values (Polisetty et al., 2008). Under such constrains, the maximum rate of ethanol production that can be obtained in the model is  $100.52 \text{ mM min}^{-1}$  (Fig. 3a).

Now, we will obtain the feasibility region under the same constrains. This is an alternative to just finding the optimal solution and it may help in discussing the changes that can be implemented in practice. First, we obtain the feasibility region without limitation in additional fluxes. In this case, the feasibility region has admissible ranges between 0.21 and 5.38 for changing HXT, and a range between 0.2 and 15 for PFK. Outside these limits, no feasible solution can be obtained (Fig. 3a). While it is reasonable to expect that ethanol production would increase by increasing HXT, our results show that increasing simultaneously HXT above 4.86-

fold and PFK above 3.16-fold (cell number 9 for HXT and cell 3 for PFK in Fig. 3a) leads to unfeasible solutions. This is so because intermediary metabolites would accumulate and the fitness of cells would decrease. Thus, a biotechnological implementation of a 5-fold change in HXT and a 5-fold change in PFK is expected to result in a failure in producing a viable strain. Our results also show that near optimum increases in ethanol production can be obtained in different conditions. As far as HXT activity is increased a minimum of 6.12-fold, we can obtain an almost optimal ethanol increase with different increases in PFK activity. For instance, we could decide a 3-fold increase in PFK and an 8-fold increase in HXT to obtain an almost optimal solution.

Planning a biotechnological strategy for increasing the production of a given metabolite must consider all the implications of the planned changes in the overall cellular metabolism. As a second scenario, we have determined the feasible solutions by imposing the additional constraint of maintaining the rate of NADPH within a 5% of its basal level (Fig. 3b). Now, the feasible region has been drastically reduced and the possible increase in ethanol production is almost minimal when compared to its basal value. Thus, in those cases in which maintaining the rate of NADPH unchanged is an important limitation, it is impossible to find a strategy involving changes in HXT and PFK capable of producing a significant increase in ethanol production (Fig. 3b).

These results show the potential application of our feasibility analysis in practical applications. Following this procedure, we can efficiently obtain an overall picture of the attainable values and a clear estimation of the unfeasible changes. This may help in discussing the most convenient implementation and the expected increase one would obtain in the objective function.

### 3.4.3. Feasible regions for the adaptive response to heat shock in yeast

Understanding the evolution of adaptive responses was the main motivation for developing the feasibility method (Guillén-Gosálbez and Sorribas, 2009). As stated before, the goal of the feasibility analysis is to obtain admissible values for enzyme activity changes that drive the model to a new state in which a set of

constraints are satisfied. Here, two fundamental ingredients are required. First, we need a mathematical model that is accurate enough to represent the biological problem at hand. Second, a set of constraints must be defined, so that changes in enzyme activity can be evaluated for compatibility.

Both issues pose significant challenges. Useful mathematical models are hard to build and, in most cases, parameters values for those models are difficult to obtain. This limitation is common to any application as discussed in the optimization section. Finding a set of constraints for the response is not an easy task either, and a sound biological understanding of the problem is required. As an example of the potential use of the methodology presented here, we consider the model presented in Fig. 2. First, we shall perform a feasibility analysis taking into account the set of constraints  $C_1$ – $C_6$  suggested by Vilaprinyo et al. (2006) (see Table 4). These constraints were identified and used for a previous analysis of operating principles in the adaptive response of yeast to heat shock.

Taking these constraints into account, we first find the upper and lower admissible values for changing each enzyme. Mathematically, this corresponds to performing a *bound contraction* on some continuous variables of the non-convex problem. Specifically, those limits are obtained by solving an optimization problem that finds the maximum and minimum values of a given  $k_i$  for which admissible solutions are found. Results are shown in Fig. 4a. These results are a generalization of those in Fig. 2(D) in the paper of Vilaprinyo et al. (2006). As we are now using a systematic search, our results include those obtained previously by intensive compu-

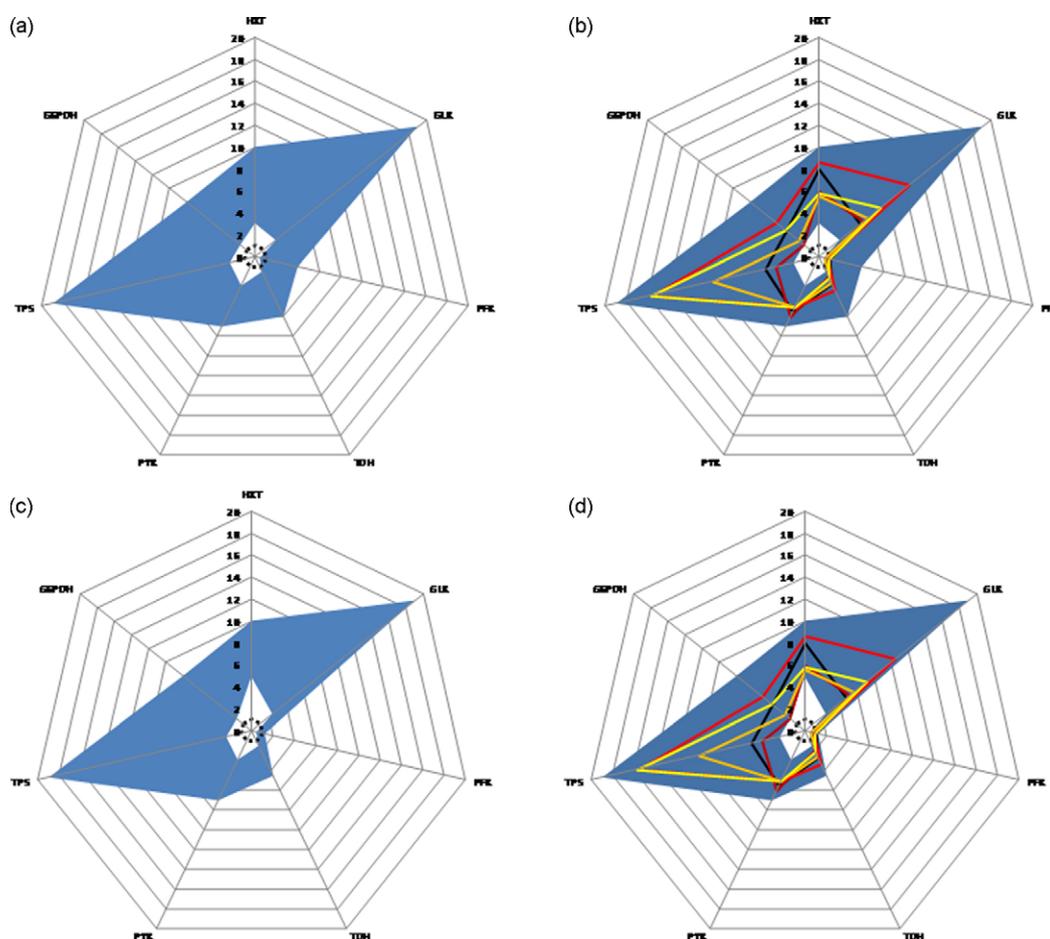
**Table 4**  
 Physiological constraints for the feasibility analysis (see Vilaprinyo et al., 2006 for details).

Constraint	Value
$C_1$	$V_{ATP} > 180.6 \text{ mM min}^{-1}$
$C_2$	$V_{TRE} > 0.03 \text{ mM min}^{-1}$
$C_3$	$V_{NADPH} > 3.54 \text{ mM min}^{-1}$
$C_4$	Internal glucose $> 0.04 \text{ mM}$
	G6P $< 20.22 \text{ mM}$
	F16P $< 22.86 \text{ mM}$
	PEP $< 0.01 \text{ mM}$
$C_5$	Cost $< 12.06$
$C_6$	$V_{Glycerol} > 0.39 \text{ mM min}^{-1}$

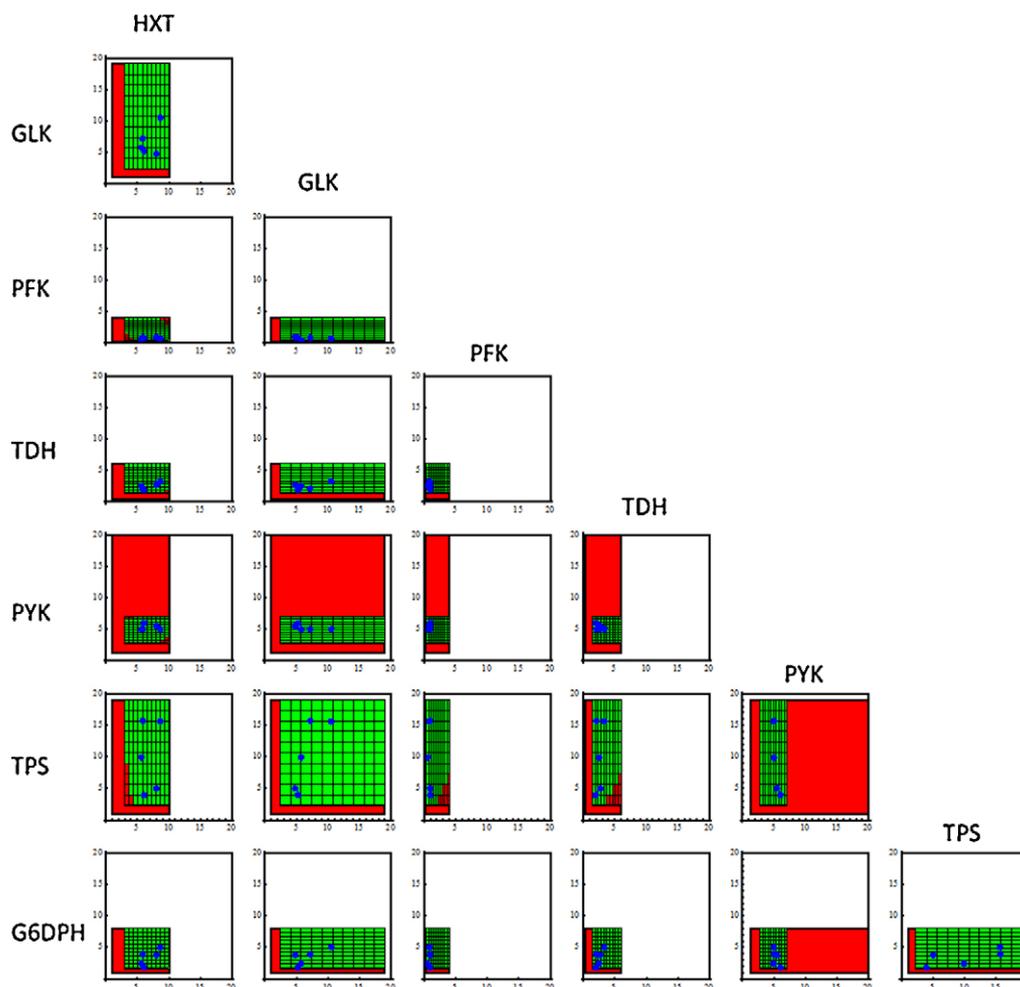
See Vilaprinyo et al. (2006) for details.

tations and are slightly wider as the previous analysis was done by considering only a set of discrete values. Furthermore, by using the new procedure, computational time is dramatically reduced to seconds.

In Fig. 4b, we plot the activity profiles corresponding to different experimental measurements (see details in Table 1 of Vilaprinyo et al., 2006). Note that all the experimental results are within the predicted values. Imposing two extra constraints ( $C_7$ – $C_8$ ) on the changes in PFK and TPS relative to the rate of trehalose synthesis ( $\Psi = (\Delta PFK \times \Delta TPS) / V_{TRE}$ ,  $\Psi < 100$ ), and a minimum value of F-1,6-P of 8.16, (criteria  $C_7$ ,  $C_8$  in Vilaprinyo et al., 2006), the limits for PFK are drastically reduced (Fig. 4c), although the resulting



**Fig. 4.** Result of the bound contraction procedure. In each case, the maximum and minimum admissible change folds for each enzyme are indicated. (a) Bounds with  $C_1$ – $C_6$  (Table 5), (b) experimental data plotted to show they are located in the admissible region found in (a), (c) Bounds with  $C_1$ – $C_8$  (see text for the definition of criteria  $C_7$  and  $C_8$ ), (d) experimental data plotted to show they are located in the admissible region found in (c). Experimental data are those of Tables 1 and 2 in Vilaprinyo et al. (2006). The search regions allowed for the change-fold in each enzyme are shown in Table 6.



**Fig. 5.** Feasibility regions for a simultaneous change in two enzymes. In each case, all the other enzymes can change to compensate and make the changes compatible with the constraints  $C_1$ – $C_6$ . Red rectangle identifies the limit for changing a given enzyme. For instance, in the case of PFK we have considered changes between 0.2 and 4. These conditions are the same considered in Vilaprinyo et al. (2006) and are maintained here for comparison. Red rectangles indicate admissible solutions. White regions are unexplored in that example. Blue points indicate experimental values described in Table 1 of Vilaprinyo et al. (2006). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of the article.)

bounds still contain the observed results (Fig. 4d, see also Fig. 2(E) in Vilaprinyo et al., 2006).

The complete results of performing the feasibility analysis are presented in Fig. 5. As indicated above, we first obtain the limits for admissible solutions taking into account criteria  $C_1$ – $C_6$  and specific limits imposed on the allowable fold changes in each enzyme based of experimental results (Table 5). For clarity, in Fig. 5 we show one-by-one figures that show the simultaneous feasibility regions for two particular enzymes. It can be seen that some enzyme activities

can take a wide range of values within their allowable boundaries, while still fulfilling the imposed constraints. This is the case for TPS and GLK. For other enzymes, feasible changes are more restricted. For example, PFK and TDH cannot increase by more than 5-fold their basal values. Outside this range, the system cannot compensate the changes and the constraints are not met. This is also the case for PFK and PYK. Feasible solutions for changes in both enzymes are obtained only in a relatively narrow margin. Interestingly, experimentally measured changes from different experiments are found to be within the feasibility regions identified by our method (see Vilaprinyo et al., 2006 for details). This is consistent with the notion that the set of constraints defined for the response are relevant for the physiological adaptation of yeast.

**Table 5**

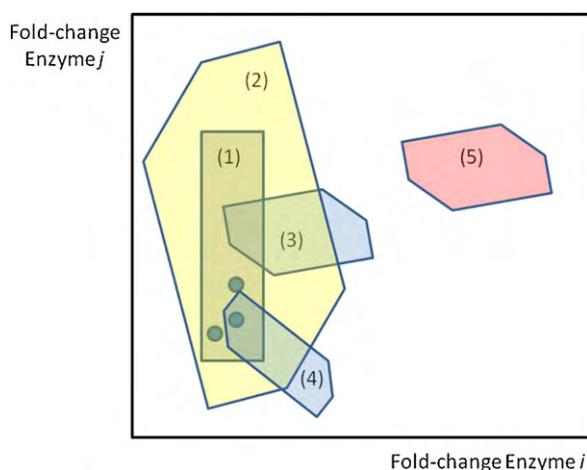
Limits for the fold change in the different enzymes in the feasibility analysis of Fig. 5.

Enzyme	Explored fold change
HXT	$1 < k_1 < 10$
GLK	$1 < k_2 < 19$
PFK	$0.25 < k_3 < 4$
TDH	$0.25 < k_4 < 6$
PYK	$0.25 < k_5 < 20$
TPS	$1 < k_6 < 19$
G6PDH	$1 < k_7 < 8$

These limits were defined considering experimental results. In each, a wide region around the values observed after heat shock are selected (see Vilaprinyo et al., 2006 for details).

### 3.4.4. On the importance of an appropriate set of constraints

The set of initial physiological constraints that are applied to the optimization procedure play a fundamental role in the feasibility analysis. Different sets of constraints are likely to produce different feasibility regions. The situation is exemplified in Fig. 6. Each of the represented regions would correspond to different sets of constraints. In this hypothetical situation, regions (1), (2), and (4) contain experimental results, while (3) and (5) do not. Thus, the constraint sets leading to regions (3) and (5) could be discarded as explanatory physiological constraints for that case. Constraints



**Fig. 6.** Hypothetical feasibility regions obtained from five different sets of constraints. Points represent experimental results on fold change for enzymes  $E_i$  and  $E_j$  in a given adaptive response (see text for details).

producing region (1) should be considered more restrictive than those of (2), although both explain the observed result. Finally, the set of constraints that produce region (3) partially explain some results but not others. In principle, this set of constraints would be less appropriate to describe the physiological requirements of the response than sets (1) and (2).

How would changing constraint sets affect the results in our analysis of yeast heat shock response? As an example, we have considered an alternative set of constraints to those used above (compare Tables 4 and 6). For illustrative purposes, in feasibility analysis using constraints from either Table 4 or Table 6, the activity of any enzyme is allowed to change between 0.2 and 20-fold. For simplicity, only the results for PFK, TDH and PYK are shown in Fig. 7. The two sets of constraints result in different feasibility regions that share some common values. Interestingly, the feasibility region obtained with the new set of constraints does not contain all the experimental values (see Table 1 in Vilaprinyo et al., 2006). This suggests that this second set of constraints does not adequately describe the physiological requirements that may have shaped the adaptive response of yeast to heat shock.

**Table 6**

Alternative set of constraints for evaluating heat shock response in yeast.

Constraint	Value
$C_1$	$V_{ATP} > 5B$
$C_2$	$V_{FRE} > 30B$
$C_3$	$V_{NADPH} > 5B$
$C_4$	$3B < \text{internal glucose} < 5B$ $15B < G6P < 20B$ $2B < F16P < 5B$ $2B < PEP < 5B$ $3B < ATP < 6B$
$C_5$	Cost $< 20$
$C_6$	$V_{Glycerol} > B$

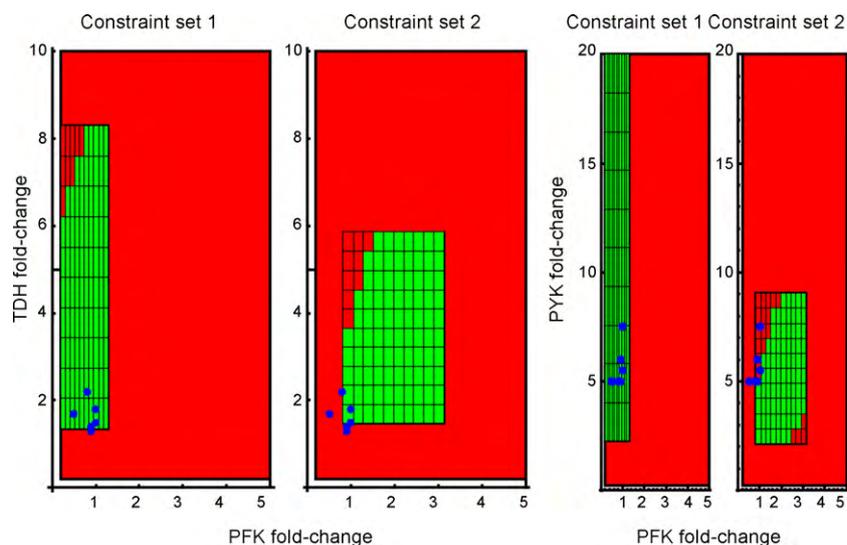
$B$  indicates de corresponding basal value (Tables 1 and 2) for the flux or metabolite considered in each criterion.

#### 4. Discussion

Understanding why metabolic pathways evolved to be as they are and how to optimize them are two closely related subjects. Studies in either field often use similar tools to compute the response of the whole system to changing conditions.

In optimization problems, control variables are manipulated by the experimenter and a predefined goal is pursued. This is often the case in metabolic engineering studies, where the general goal is that of modifying cells so that specific production targets can be reached (Hatzimanikatis et al., 1998). Typically, one considers optimizing the yield of a given process, maximizing flux through a pathway, etc. Then, optimization procedures are used on a mathematical model for the relevant processes in order to analyze which changes are the most likely to produce the desired result (Gianchandani et al., 2008). There is a wide scope of optimization methods that can be used for this task, based on different optimization strategies (Banga, 2008; Chang and Sahinidis, 2005; Marin-Sanguino et al., 2007; Nielsen, 2007; Polisetty et al., 2008; Schuetz et al., 2007; Vital-Lopez et al., 2006).

In evolutionary studies, however, we are faced with conserved changes that can appear in organisms by random mutations, by gene transfer, gene duplication, gene deletion, and other mechanisms. Natural selection may operate as a purifying mechanism that acts upon the systemic effect of these changes on the overall fitness



**Fig. 7.** Feasibility analysis obtained with two different sets of constraints. Constraint set 1 correspond to Table 5, constraint set 2 to Table 6. In both cases, enzymes are allowed to change between 0.2 and 20-fold over basal. Blue points indicate experimental measurements as presented in Table 1 of Vilaprinyo et al. (2006).

and leads to the fixation of new designs and operative patterns in a population, due to the differential reproduction of individuals. As a result, organisms often evolve towards some quasi-optimal regime under the conditions they live in. Such regime however may become quite sub-optimal if conditions change drastically. Those changes could lead to a new round of natural selection, this time with different physiological constraints. Thus, evolution in natural systems can be seen as a perpetual optimization-like process, with the parameter conditions that maximize survival and reproduction shifting over time.

In fact, one of the biggest current problems in this area is how to establish a connection between what researchers see as the actual functioning conditions of the molecular pathways that allow a cell to perform appropriately and the fitness of that cell. Causative genotype phenotype models (Martens et al., 2009; O'Connor and Mundy, 2009) are but a start in connecting the optimization of the molecular determinants of life and the fitness of organisms. We hypothesize that adaptive responses are to be found within feasible regions that allow the system to meet a set of physiological constraints that are required for cell survival (Guillén-Gosálbez and Sorribas, 2009; Vilaprinyo et al., 2006). The numerical thresholds considered in these constraints would shape the admissible changes in the system parameters so that the effect on global fitness can be sensed by natural selection. As a result, a specific adaptive response would evolve. Future work should deal with connecting the molecular aspects of the adaptive response to the direct survival ability.

From a practical point of view, there is a set of considerations that should be taken into account in optimization related studies of biological problems at the pathway level: (1) a model that can be used to compute fluxes, metabolite levels, the effect of changes in parameters, dynamic response, etc., is required; (2) stoichiometry-based models, such as Flux Balance Analysis models, are not sufficiently accurate to be used for characterizing quantitative changes. This is so because they do not account for regulatory interactions within the network and cannot be used to accurately calculate metabolite levels, dynamic changes, and other quantitative information (Nikolaev, 2009); (3) models that include information about the regulatory signals are essential for an accurate analysis; (4) kinetic information, even if it is only approximated, is required to define a quantitative model that may help in the analysis. Because of (4), at present we are still unable to create genome-wide models for metabolism, because not enough information is available. However, the obtained results show the importance of developing GMA-like models as a basis for a more complete analysis of system optimization and evolution. In this paper we have presented a methodology designed to address important practical questions, both in metabolic engineering applications and in studies of pathway evolution, through the use of a global nonlinear optimization technique and the characterization of feasibility regions. Although linear global optimization methods had been used before to search for survivability regions in Flux Balance Analysis models, those studies have the limitations described in points (1)–(4) of the previous paragraph. Our methodology overcomes those limitations and it can be applied to a special class of nonlinear differential equations models known as GMA models.<sup>3</sup> If such a model is defined for a given metabolic problem, then our method allows for an exhaustive exploration of different evolutionary strategies and a systematic characterization of the physiological requirements that may underlie the evolution of adaptive strategies.

<sup>3</sup> It should be noted that ODE models written using other mathematical forms can be recasted into GMA models, increasing the generality of the method presented here.

## Acknowledgements

AS acknowledges the financial support from grant BFU2008-0196 (Ministerio de Ciencia e Innovación (MICINN, Spain)). GG-G expresses his gratitude for the financial support from the MICINN (projects DPI2008-04099, PHB2008-0090-PC, BFU2008-754 00196 and CTQ2009-14420-C02-01) and the Spanish Ministry of External Affairs (projects A/8502/07, HS2007-0006 and A/020104/08). RA is supported by the MICINN through the Ramon y Cajal program and grant BFU2007-62772/BMC. EV is supported by grant PI06/1649 from the Spanish Ministry of Health.

## References

- Aldana, M., Balleza, E., Kauffman, S., Resendiz, O., 2007. Robustness and evolvability in genetic regulatory networks. *J. Theor. Biol.* 245, 433–448.
- Alvarez-Vasquez, F., Sims, K.J., Hannun, Y.A., Voit, E.O., 2004. Integration of kinetic information on yeast sphingolipid metabolism in dynamical pathway models. *J. Theor. Biol.* 226, 265–291.
- Alves, R., Savageau, M.A., 2000. Extending the method of mathematically controlled comparison to include numerical comparisons. *Bioinformatics* 16, 786–798.
- Bailey, J.E., 2001. Reflections on the scope and the future of metabolic engineering and its connections to functional genomics and drug discovery. *Metab. Eng.* 3, 111–114.
- Bailey, J.E., 1999. Lessons from metabolic engineering for functional genomics and drug discovery. *Nat. Biotechnol.* 17, 616–618.
- Bailey, J.E., Sburlati, A., Hatzimanikatis, V., Lee, K., Renner, W.A., Tsai, P.S., 1996. Inverse metabolic engineering: a strategy for directed genetic engineering of useful phenotypes. *Biotechnol. Bioeng.* 52, 109–121.
- Bailey, J.E., 1991. Toward a science of metabolic engineering. *Science* 252, 1668–1675.
- Bailey, J.E., Birnbaum, S., Galazzo, J.L., Khosla, C., Shanks, J.V., 1990. Strategies and challenges in metabolic engineering. *Ann. N. Y. Acad. Sci.* 589, 1–15.
- Banga, J.R., 2008. Optimization in computational systems biology. *BMC Syst. Biol.* 2, 47.
- Bedford, T., Hartl, D.L., 2009. Optimization of gene expression by natural selection. *Proc. Natl. Acad. Sci. U.S.A.* 106, 1133–1138.
- Braunstein, A., Mulet, R., Pagnani, A., 2008. Estimating the size of the solution space of metabolic networks. *BMC Bioinformatics* 9, 240.
- Cascante, M., Curto, R., Sorribas, A., 1995. Comparative characterization of the fermentation pathway of *Saccharomyces cerevisiae* using biochemical systems theory and metabolic control analysis: steady-state analysis. *Math. Biosci.* 130, 51–69.
- Causton, H.C., Ren, B., Koh, S.S., Harbison, C.T., Kanin, E., Jennings, E.G., Lee, T.I., True, H.L., Lander, E.S., Young, R.A., 2001. Remodeling of yeast genome expression in response to environmental changes. *Mol. Biol. Cell* 12, 323–337.
- Chachuat, B., Singer, A.B., Barton, P.I., 2006. Global methods for dynamic optimization and mixed-integer dynamic optimization. *Ind. Eng. Chem. Res.* 45, 8373–8392.
- Chang, Y.J., Sahinidis, N.V., 2005. Optimization of metabolic pathways under stability considerations. *Comput. Chem. Eng.* 29, 467–479.
- Chou, I.C., Voit, E.O., 2009. Recent developments in parameter estimation and structure identification of biochemical and genomic systems. *Math. Biosci.* 219, 57–83.
- Coelho, P.M., Salvador, A., Savageau, M.A., 2009. Quantifying global tolerance of biochemical systems: Design implications for moiety-transfer cycles. *PLoS Comput. Biol.* 5, e1000319.
- Curto, R., Voit, E.O., Sorribas, A., Cascante, M., 1997. Validation and steady-state analysis of a power-law model of purine metabolism in man. *Biochem. J.* 324 (Pt 3), 761–775.
- Curto, R., Sorribas, A., Cascante, M., 1995. Comparative characterization of the fermentation pathway of *Saccharomyces cerevisiae* using biochemical systems theory and metabolic control analysis: model definition and nomenclature. *Math. Biosci.* 130, 25–50.
- Dayarian, A., Chaves, M., Sontag, E.D., Sengupta, A.M., 2009. Shape, size, and robustness: feasible regions in the parameter space of biochemical networks. *PLoS Comput. Biol.* 5, e1000256.
- de Atauri, P., Sorribas, A., Cascante, M., 2000. Analysis and prediction of the effect of uncertain boundary values in modeling a metabolic pathway. *Biotechnol. Bioeng.* 68, 18–30.
- Eisen, M.B., Spellman, P.T., Brown, P.O., Botstein, D., 1998. Cluster analysis and display of genome-wide expression patterns. *Proc. Natl. Acad. Sci. U.S.A.* 95, 14863–14868.
- El-Samad, H., Kurata, H., Doyle, J.C., Gross, C.A., Khammash, M., 2005. Surviving heat shock: control strategies for robustness and performance. *Proc. Natl. Acad. Sci. U.S.A.* 102, 2736–2741.
- Espósito, W.R., Floudas, C.A., 2000. Deterministic global optimization in nonlinear optimal control problems. *J. Global Optim.* 17, 97–126.
- García-Martínez, J., González-Candelas, F., Pérez-Ortín, J.E., 2007. Common gene expression strategies revealed by genome-wide analysis in yeast. *Genome Biol.* 8, R222.

- Gasch, A.P., Spellman, P.T., Kao, C.M., Carmel-Harel, O., Eisen, M.B., Storz, G., Botstein, D., Brown, P.O., 2000. Genomic expression programs in the response of yeast cells to environmental changes. *Mol. Biol. Cell* 11, 4241–4257.
- Gianchandani, E.P., Oberhardt, M.A., Burgard, A.P., Maranas, C.D., Papin, J.A., 2008. Predicting biological system objectives de novo from internal state measurements. *BMC Bioinformatics* 9, 43.
- Goodman, C., 2008. Engineering ingenuity at iGEM. *Nat. Chem. Biol.* 4, 13.
- Grossman, I.E., Biegler, L.T., 2004. Part II. Future perspective on optimization. *Comput. Chem. Eng.* 28, 1193–1218.
- Guillén-Gosálbez, G., Sorribas, A., 2009. Identifying quantitative operation principles in metabolic pathways: a systematic method for searching feasible enzyme activity patterns leading to cellular adaptive responses. *BMC Bioinformatics* 10, 386.
- Han, J.D., 2008. Understanding biological functions through molecular networks. *Cell Res.* 18, 224–237.
- Hansen, P., Jaumard, B., Lu, S., 1992. Global optimization of univariate Lipschitz functions. I. *Surrey and properties. Math. Program.* 55, 251–272.
- Hatzimanikatis, V., Liao, J.C., 2002. A memorial review of Jay Bailey's contribution in prokaryotic metabolic engineering. *Biotechnol. Bioeng.* 79, 504–508.
- Hatzimanikatis, V., Emmerling, M., Sauer, U., Bailey, J.E., 1998. Application of mathematical tools for metabolic design of microbial ethanol production. *Biotechnol. Bioeng.* 58, 154–161.
- Hatzimanikatis, V., Floudas, C.A., Bailey, J.E., 1996. Optimization of regulatory architectures in metabolic reaction networks. *Biotechnol. Bioeng.* 52, 485–500.
- Hlavacek, W.S., Savageau, M.A., 1998. Method for determining natural design principles of biological control circuits. *J. Intell. Fuzzy Syst.* 6, 87.
- Jenkins, G.M., 2003. The emerging role for sphingolipids in the eukaryotic heat shock response. *Cell Mol. Life Sci.* 60, 701–710.
- Kashiwagi, A., Urabe, I., Kaneko, K., Yomo, T., 2006. Adaptive response of a gene network to environmental changes by fitness-induced attractor selection. *PLoS One* 1, e49.
- Kitano, H., 2007. Towards a theory of biological robustness. *Mol. Syst. Biol.* 3, 137.
- Kitano, H., 2004. Biological robustness. *Nat. Rev. Genet.* 5, 826–837.
- Klipp, E., Nordlander, B., Kruger, R., Gennemark, P., Hohmann, S., 2005. Integrative model of the response of yeast to osmotic shock. *Nat. Biotechnol.* 23, 975–982.
- Koonin, E.V., 2009. Darwinian evolution in the light of genomics. *Nucleic Acids Res.* 37, 1011–1034.
- Kurata, H., El-Samad, H., Iwasaki, R., Ohtake, H., Doyle, J.C., Grigoro, I., Gross, C.A., Khammash, M., 2006. Module-based analysis of robustness tradeoffs in the heat shock response system. *PLoS Comput. Biol.* 2, e59.
- Lee, J.M., Gianchandani, E.P., Papin, J.A., 2006. Flux balance analysis in the era of metabolomics. *Brief Bioinform.* 7, 140–150.
- Marin-Sanguino, A., Voit, E.O., Gonzalez-Alcon, C., Torres, N.V., 2007. Optimization of biotechnological systems through geometric programming. *Theor. Biol. Med. Model.* 4, 38.
- Martens, H., Veflingstad, S.R., Plahte, E., Martens, M., Bertrand, D., Omholt, S.W., 2009. The genotype-phenotype relationship in multicellular pattern-generating models—the neglected role of pattern descriptors. *BMC Syst. Biol.* 3, 87.
- Mitchell, A., Romano, G.H., Groisman, B., Yona, A., Dekel, E., Kupiec, M., Dahan, O., Pilpel, Y., 2009. Adaptive prediction of environmental changes by microorganisms. *Nature* 460, 220–224.
- Molina-Navarro, M.M., Castells-Roca, L., Belli, G., Garcia-Martinez, J., Marin-Navarro, J., Moreno, J., Perez-Ortin, J.E., Herrero, E., 2008. Comprehensive transcriptional analysis of the oxidative response in yeast. *J. Biol. Chem.* 283, 17908–17918.
- Morohashi, M., Winn, A.E., Borisuk, M.T., Bolouri, H., Doyle, J., Kitano, H., 2002. Robustness as a measure of plausibility in models of biochemical networks. *J. Theor. Biol.* 216, 19–30.
- Nielsen, J., 2007. Principles of optimal metabolic network operation. *Mol. Syst. Biol.* 3, 126.
- Nikolaev, E.V., 2009. The elucidation of metabolic pathways and their improvements using stable optimization of large-scale kinetic models of cellular systems. *Metab. Eng.*
- O'Connor, T.D., Mundy, N.I., 2009. Genotype-phenotype associations: Substitution models to detect evolutionary associations between phenotypic variables and genotypic evolutionary rate. *Bioinformatics* 25, i94–i100.
- Papamichail, I., Adjiman, C.S., 2004. Global optimization of dynamic systems. *Comput. Chem. Eng.* 28, 403–415.
- Papamichail, I., Adjiman, C.S., 2002. A rigorous global optimization algorithm for problems with ordinary differential equations. *J. Global Optim.* 24, 1–33.
- Polisetty, P.K., Gatzke, E.P., Voit, E.O., 2008. Yield optimization of regulated metabolic systems using deterministic branch-and-reduce methods. *Biotechnol. Bioeng.* 99, 1154–1169.
- Porn, R., Bjork, K.M., Westerlund, T., 2008. Global solution of optimization problems with signomial parts. *Discrete Optim.* 5, 108–120.
- Pozo, C., Sorribas, A., Jiménez, L., Guillén-Gosálbez, G., submitted for publication. Outer approximation-based algorithm for biotechnology studies in systems biology. *Comput. Chem. Eng.*
- Quesada, I., Grossman, I.E., 1995. A global optimization algorithm for linear fractional and bilinear programs. *J. Global Optim.* 6, 39–76.
- Raiford, D.W., Heizer Jr., E.M., Miller, R.V., Akashi, H., Raymer, M.L., Krane, D.E., 2008. Do amino acid biosynthetic costs constrain protein evolution in *Saccharomyces cerevisiae*? *J. Mol. Evol.* 67, 621–630.
- Romero-Santacreu, L., Moreno, J., Perez-Ortin, J.E., Alepuz, P., 2009. Specific and global regulation of mRNA stability during osmotic stress in *Saccharomyces cerevisiae*. *RNA*.
- Salvador, A., Savageau, M.A., 2006. Evolution of enzymes in a series is driven by dissimilar functional demands. *Proc. Natl. Acad. Sci. U.S.A.* 103, 2226–2231.
- Salvador, A., Savageau, M.A., 2003. Quantitative evolutionary design of glucose 6-phosphate dehydrogenase expression in human erythrocytes. *Proc. Natl. Acad. Sci. U.S.A.* 100, 14463–14468.
- Savageau, M.A., Coelho, P.M., Fasani, R.A., Tolla, D.A., Salvador, A., 2009. Phenotypes and tolerances in the design space of biochemical systems. *Proc. Natl. Acad. Sci. U.S.A.* 106, 6435–6440.
- Savageau, M.A., 1998. Demand theory of gene regulation. I. Quantitative development of the theory. *Genetics* 149, 1665–1676.
- Savageau, M.A., 1976. *Biochemical Systems Analysis: A Study of Function and Design in Molecular Biology*. Addison-Wesley, Reading, MA.
- Savageau, M.A., 1975. Optimal design of feedback control by inhibition: dynamic considerations. *J. Mol. Evol.* 5, 199–222.
- Savageau, M.A., 1974a. Optimal design of feedback control by inhibition. *J. Mol. Evol.* 4, 139–156.
- Savageau, M.A., 1974b. Genetic regulatory mechanisms and the ecological niche of *Escherichia coli*. *Proc. Natl. Acad. Sci. U.S.A.* 71, 2453–2455.
- Savageau, M.A., 1971. Parameter sensitivity as a criterion for evaluating and comparing the performance of biochemical systems. *Nature* 229, 542–544.
- Savageau, M.A., 1969a. Biochemical systems analysis. II. The steady-state solutions for an n-pool system using a power-law approximation. *J. Theor. Biol.* 25, 370–379.
- Savageau, M.A., 1969b. Biochemical systems analysis. I. Some mathematical properties of the rate law for the component enzymatic reactions. *J. Theor. Biol.* 25, 365–369.
- Schuetz, R., Kuepfer, L., Sauer, U., 2007. Systematic evaluation of objective functions for predicting intracellular fluxes in *Escherichia coli*. *Mol. Syst. Biol.* 3, 119.
- Sims, K.J., Spassieva, S.D., Voit, E.O., Obeid, L.M., 2004. Yeast sphingolipid metabolism: clues and connections. *Biochem. Cell Biol.* 82, 45–61.
- Singer, A.B., Barton, P.L., 2006. Global optimization with nonlinear ordinary differential equations. *J. Global Optim.* 34, 159–190.
- Tawarmalani, M., Sahinidis, N.V., 2002. Convexification and global optimization. In: Anonymous (Ed.), *Continuous and Mixed-Integer Nonlinear Programming: Theory, Algorithms, Software, and Applications (Nonconvex Optimization and its Applications)*. Springer.
- Teusink, B., Wiersma, A., Jacobs, L., Notebaart, R.A., Smid, E.J., 2009. Understanding the adaptive growth strategy of *Lactobacillus plantarum* by in silico optimisation. *PLoS Comput. Biol.* 5, e1000410.
- Torres, N.V., Voit, E.O., Glez-Alcon, C., Rodriguez, F., 1997. An indirect optimization method for biochemical systems: description of method and application to the maximization of the rate of ethanol, glycerol, and carbohydrate production in *Saccharomyces cerevisiae*. *Biotechnol. Bioeng.* 55, 758–772.
- Torres, N.V., Voit, E.O., Gonzalez-Alcon, C., 1996. Optimization of nonlinear biotechnological processes with linear programming: Application to citric acid production by *Aspergillus niger*. *Biotechnol. Bioeng.* 49, 247–258.
- Trotter, E.W., Berenfeld, L., Krause, S.A., Petsko, G.A., Gray, J.V., 2001. Protein misfolding and temperature up-shift cause G1 arrest via a common mechanism dependent on heat shock factor in *Saccharomyces cerevisiae*. *Proc. Natl. Acad. Sci. U.S.A.* 98, 7313–7318.
- Vilaprinyo, E., Alves, R., Sorribas, A., 2006. Use of physiological constraints to identify quantitative design principles for gene expression in yeast adaptation to heat shock. *BMC Bioinformatics* 7, 184.
- Vital-Lopez, F.G., Armaou, A., Nikolaev, E.V., Maranas, C.D., 2006. A computational procedure for optimal engineering interventions using kinetic models of metabolism. *Biotechnol. Prog.* 22, 1507–1517.
- Voit, E.O., 2003. Design principles and operating principles: the yin and yang of optimal functioning. *Math. Biosci.* 182, 81–92.
- Voit, E.O., 2000. *Computational Analysis of Biochemical Systems. A Practical Guide for Biochemists and Molecular Biologists*. Cambridge University Press, Cambridge, U.K.
- Voit, E.O., Radivoyevitch, T., 2000. Biochemical systems analysis of genome-wide expression data. *Bioinformatics* 16, 1023–1037.
- Voit, E.O., 1992. Optimization in integrated biochemical systems. *Biotechnol. Bioeng.* 40, 572–582.
- Wagner, A., 2005. Energy constraints on the evolution of gene expression. *Mol. Biol. Evol.* 22, 1365–1374.
- Watson, T.G., 1970. Effects of sodium chloride on steady-state growth and metabolism of *Saccharomyces cerevisiae*. *J. Gen. Microbiol.* 64, 91–99.
- Wiebe, M.G., Rintala, E., Tamminen, A., Simolin, H., Salusjarvi, L., Toivari, M., Kokkonen, J.T., Kiuru, J., Ketola, R.A., Jouhten, P., Huuskonen, A., Maheimo, H., Ruohonen, L., Penttila, M., 2008. Central carbon metabolism of *Saccharomyces cerevisiae* in anaerobic, oxygen-limited and fully aerobic steady-state conditions and following a shift to anaerobic conditions. *FEMS Yeast Res.* 8, 140–154.
- Wilkins, A.S., 2007. Between “design” and “bricolage”: genetic networks, levels of selection, and adaptive evolution. *Proc. Natl. Acad. Sci. U.S.A.* 104 (Suppl. 1), 8590–8596.

METHODOLOGY ARTICLE

Open Access

# Steady-state global optimization of metabolic non-linear dynamic models through recasting into power-law canonical models

Carlos Pozo<sup>1</sup>, Alberto Marín-Sanguino<sup>2</sup>, Rui Alves<sup>3</sup>, Gonzalo Guillén-Gosálbez<sup>1</sup>, Laureano Jiménez<sup>1</sup> and Albert Sorribas<sup>3\*</sup>

## Abstract

**Background:** Design of newly engineered microbial strains for biotechnological purposes would greatly benefit from the development of realistic mathematical models for the processes to be optimized. Such models can then be analyzed and, with the development and application of appropriate optimization techniques, one could identify the modifications that need to be made to the organism in order to achieve the desired biotechnological goal. As appropriate models to perform such an analysis are necessarily non-linear and typically non-convex, finding their global optimum is a challenging task. Canonical modeling techniques, such as Generalized Mass Action (GMA) models based on the power-law formalism, offer a possible solution to this problem because they have a mathematical structure that enables the development of specific algorithms for global optimization.

**Results:** Based on the GMA canonical representation, we have developed in previous works a highly efficient optimization algorithm and a set of related strategies for understanding the evolution of adaptive responses in cellular metabolism. Here, we explore the possibility of recasting kinetic non-linear models into an equivalent GMA model, so that global optimization on the recast GMA model can be performed. With this technique, optimization is greatly facilitated and the results are transposable to the original non-linear problem. This procedure is straightforward for a particular class of non-linear models known as Saturable and Cooperative (SC) models that extend the power-law formalism to deal with saturation and cooperativity.

**Conclusions:** Our results show that recasting non-linear kinetic models into GMA models is indeed an appropriate strategy that helps overcoming some of the numerical difficulties that arise during the global optimization task.

## 1 Background

Identifying optimization strategies for increasing strain productivity should be possible by applying optimization methods to detailed kinetic models of the target metabolism. Thus, a rational approach would pinpoint the changes to be done - e.g. by modulating gene expression - in order to achieve the desired biotechnological goals [1-4]. To build such models we can either start from a detailed description of the underlying processes (bottom-up strategy) or we can fit kinetic models to experimental data obtained *in vivo* (top-down strategy).

The bottom-up approach was the original strategy for model building in the biological sciences. Bottom-up kinetic models require information that is seldom available, despite the increasing amount of kinetic data contained in a growing set of databases (for example see [5,6] and the webpage <http://kinetics.nist.gov/kinetics/index.jsp>). Even in the best case scenarios where kinetic data are available, the data have often been obtained in different labs and under *in vitro* conditions that are hardly ever comparable or representative of the situation *in vivo*. In addition, models built using this strategy often fail to adequately reproduce the known behavior of the target system [7-10]. With the accumulation of time-series data, which were originally generated from the accurate measurement of transient responses, top-down modeling became viable as an alternative to the

\* Correspondence: [albert.sorribas@cmb.udl.cat](mailto:albert.sorribas@cmb.udl.cat)

<sup>3</sup>Departament de Ciències Mèdiques Bàsiques, Institut de Recerca Biomèdica de Lleida (IRBLLEIDA), Universitat de Lleida, Montserrat Roig 2, 25008 Lleida, Spain

Full list of author information is available at the end of the article

bottom-up strategy [11]. However, top-down modeling also faces important difficulties. For example, regulatory interactions between metabolites and enzymes are very poorly characterized and most metabolic maps lack such crucial information. Therefore, for a given network structure (i.e. a stoichiometric description) obtained from databases, a large number of alternative regulatory patterns may exist that account for the observed experimental data [12]. Model discrimination among the alternative regulatory patterns requires appropriate experimental design. However, this is seldom considered when performing the time series measurements. Last, but not the least, parameter identifiability in highly non-linear models can be problematic (for a review see [13]).

An additional issue that is common to models built using both strategies is that such detailed kinetic models include non-convexities that lead to the existence of multiple local optima in which standard non-linear optimization algorithms may get trapped during the search. Several stochastic and deterministic global optimization methods have been proposed to overcome this limitation [14]. Deterministic methods, which are the only ones that can rigorously guarantee global optimality, rely on the use of convex envelopes or underestimators to formulate lower-bounding convex problems that are typically combined with spatial branch and bound strategies. Most of these methods are general purpose and assume special structures in the continuous terms of the mathematical model. Because of this, they can encounter numerical difficulties in specific metabolic engineering systems that require the optimization of kinetic models with a large number of non-convexities of different nature.

Given all these issues, it is hardly surprising that linear stoichiometric models have emerged as the most popular tool to analyze genome-wide metabolic networks using optimization techniques. Linear optimization problems can be solved using very fast and efficient algorithms [15,16] that are implemented in almost every kind of computer, ranging from laptops to cloud computing centers. In addition, such models require a relatively small amount of information.

The possibility of condensing information about a very large network in a compact form enabled stoichiometric models to provide interesting insights in many different cases. However, the apparent simplicity in building and analyzing stoichiometric models comes at the cost of neglecting regulatory signals, metabolite levels and dynamic constraints. Accounting for these features in a dynamic way requires using more detailed, non-linear, mathematical models [17,18].

These models go a step further than stoichiometric models by incorporating regulatory influences through a set of ordinary differential equations that can account

for the system's dynamics. Building such models is often impossible because the appropriate functional form that needs to be used to describe the dynamical behavior of specific processes is in general unknown. Modeling strategies based on systematic approximated kinetic representations, such as power-laws [19-22], Saturating and Cooperative [23], or convenience kinetics [24], side-step this issue by providing uniform forms that are guaranteed to be accurate over a range of conditions and reduce the amount of information required to build the models. Because of the regularity in the form of the mathematical function, models based on approximate formalisms can be automatically built from the reaction scheme of the target system. The model parameters can subsequently be estimated from experimental data using different procedures [13,25].

Although building and analyzing of comprehensive genome-wide detailed models is still not viable in most cases (see however [26,27]), developing ways to extend large scale optimization analysis to larger and more realistic non-linear kinetic models is an important part of the future of systems biology [18]. In fact, the optimization of certain types of non-linear problems can already be solved very efficiently and geometric programming problems with up to 1,000 variables and 10,000 constraints can be solved in minutes on a personal computer.

Efficient global optimization techniques are available for power-law models [1,28-30], either in S-system form or in Generalized Mass Action (GMA) form (for a review see [31]). In the case of S-system models, a simple logarithmic transformation brings the model to a linear form [1]. In the case of GMA models, the problem can be efficiently solved using branch-and-bound [28,32] and outer-approximation techniques [29,30].

The usefulness of the global optimization techniques developed for GMA models has been shown in the analysis of the adaptive response of yeast to heat shock [29,33]. In essence, starting with a GMA model and considering a set of constraints on flux and metabolite values, we can obtain: (i) The pattern of enzyme activities that maximizes a given objective, (ii) The region of feasible changes in enzyme activities so that the model fulfills a set of constraints on fluxes, metabolites, maximum allowable change in activity, etc., and (iii) A heat map of how the objective function changes within the feasible region. These results share some similarities with those produced with stoichiometric models, but incorporate many additional features.

Based on ideas similar to those that led to the development of the power-law formalism, Sorribas et al. [23] proposed a new Saturable and Cooperative (SC) formalism, that extends the power-law representation to include cooperativity and saturation. Although models

built using this new formalism loses some of the simplicity inherent to the analysis of S-systems and GMA models, they tend to be accurate over a wider range of conditions than both the S-System and GMA representations [23]. Thus, it is important to enlarge the scope of global optimization methods developed for power-law models in order to deal with the SC formalism and analyze under which situations the later models behave better than the former.

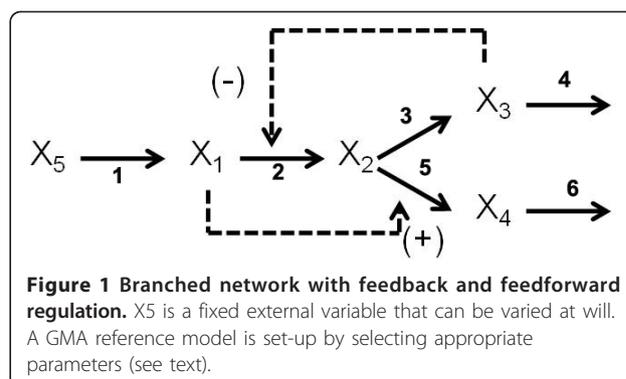
Optimization of SC models faces a number of practical problems common to kinetic non-linear models [34,35]. To sidestep these problems, and in order to be able to use global optimization methods developed for power-law models, we will use a technique called recasting. Recasting permits the exact transformation of a continuous non-linear model with an arbitrary form into a canonical GMA model [36,37]. This transformation is typically performed by increasing the number of variables of the original model. Through this technique, arbitrary non-linear models can be represented using a canonical form such as GMA or S-system that can be used for simulation and optimization purposes, which opens the door for effectively extending the optimization and feasibility analysis originally devised for GMA models to other detailed kinetic models.

In this paper, and as a first step to define a framework for optimization of non-linear models with arbitrary form and extend FBA and related approaches to detailed kinetic models, we shall show the practical utility of recasting SC models into GMA models for optimization purposes. This technique is similar to the symbolic reformulation algorithm proposed by Smith and Panteleides [38]. Our method, however, focuses on obtaining a power-law representation that greatly facilitates global optimization, instead of continuing with the recasting until converting the model to a standard form containing linear constraints and a set of nonlinearities corresponding to bilinear product, linear fractional, simple exponentiation and univariate function terms. After recasting the model to the canonical form, we can apply any of the optimization strategies we have presented for GMA models [29,32] to obtain the global optimum of the original SC problem.

## 2 Results

### 2.1 Global optimization of non-linear models through recasting

For a proof of concept of the difficulties of global optimizing non-linear models and of the use of recasting for attaining practical solutions, we shall start by defining a reference biochemical network that corresponds to the reaction scheme in Figure 1. This hypothetical system has a source metabolite  $X_5$  and four internal metabolites. The network includes six reactions and a branch



point.  $X_3$  acts as a feed-back inhibitor of the synthesis of  $X_2$ , while  $X_1$  is an activator of the synthesis of  $X_4$ .

The generic model for this system is:

$$\begin{aligned} \dot{X}_1 &= v_1 - v_2 \\ \dot{X}_2 &= v_2 - v_3 - v_5 \\ \dot{X}_3 &= v_3 - v_4 \\ \dot{X}_4 &= v_5 - v_6 \end{aligned} \quad (1)$$

Each of the velocities is a non-linear function of the involved metabolites. The SC representation, provides a systematic way for defining a functional model of this pathway. As a demonstrative example, let us suppose that the numerical model is:

$$\begin{aligned} \frac{dX_1}{dt} &= \frac{20k_1X_5^1}{X_5^1 + 1} - \frac{40k_2X_1^1}{(X_1^1 + 1)X_3^2 \left(1 + \frac{1}{X_3^2}\right)} \\ \frac{dX_2}{dt} &= \frac{40k_2X_1^1}{(X_1^1 + 1)X_3^2 \left(1 + \frac{1}{X_3^2}\right)} - \frac{7.5k_3X_2^{2.5}}{X_2^{2.5} + 0.25} \\ &\quad - \frac{16k_5X_1^1X_2^1}{(X_1^1 + 1)(X_2^1 + 1)} \\ \frac{dX_3}{dt} &= \frac{7.5k_3X_2^{2.5}}{X_2^{2.5} + 0.25} - \frac{12k_4X_3^1}{X_3^1 + 1} \\ \frac{dX_4}{dt} &= \frac{16k_5X_1^1X_2^1}{(X_1^1 + 1)(X_2^1 + 1)} - \frac{8k_6X_4^1}{X_4^1 + 1} \end{aligned} \quad (2)$$

In these equations,  $k_r$ ,  $r = 1, \dots, 6$  is an auxiliary variable used to model changes in the enzyme activity. At the basal level,  $k_r = 1$  for all the reactions. During the optimization tasks, it is possible to limit the maximum change in gene expression by imposing a maximum allowable change in  $k_r$ .

We shall now address the following questions:

(i) To what extent can general purpose global optimization methods be applied to SC models?, (ii) Given that a SC model can be recast as a GMA (rGMA), is

this useful for optimization of the original SC model?, (iii) Are the results obtained with the rGMA equivalent to the results of the original SC model?, and (iv) What are the practical advantages of optimizing a rGMA model?.

## 2.2 Optimization goals

In order to address the questions posed at the end of the previous section we shall define the following optimizations tasks (note that changes in enzyme activities and metabolite concentrations are constrained between  $0.2 \leq k_r \leq 5.0$  and  $0.1 \leq X_i \leq 10.0$  respectively in all the instances unless otherwise specified):

- O1: What is the optimal pattern of changes in enzyme activities that maximizes the objective function in the new steady-state for a fixed value of  $X_5$ ?
- O2: What is the optimal pattern of changes in enzyme activities that maximizes the objective function in the new steady-state for a fixed value of  $X_5$  considering a maximum allowable variation of 10% in the steady-state values of the intermediaries?
- O3: What is the optimal pattern of changes in enzyme activities that maximizes the objective function in the new steady-state for a fixed value of  $X_5$  considering changes in the output flux from  $X_4$  of less than 10% with respect to its reference value?
- O4: What is the best set of changes, assuming that we can only manipulate three enzymes, that maximizes the objective function in the new steady-state for a fixed value of  $X_5$  considering a maximum variation of 10% in the steady-state values of the intermediaries?

Two different objective functions (OF), steady-state concentration of  $X_3$  and flux  $v_4$ , have been considered for each optimization case, except for O3. This latter case has been optimized in terms only of the first objective (i.e., steady-state concentration of  $X_3$ ), because limits on  $v_4$  are already included in the formulation of the optimization problem.

## 2.3 Global optimization of SC models using BARON

We first address the optimization of the aforementioned model in their original SC form using state of the art global optimization techniques. The model was coded in the algebraic modeling system GAMS 23.0.2 and solved with the commercial global optimization package BARON v.8.1.5. on an Intel 1.2 GHz machine. An optimality gap (i.e., tolerance) of 0.2% was set in all the instances. As can be seen in Table 1, BARON produce results with an optimality gap (OG) below the specified tolerance.

**Table 1 Results for the maximization of  $X_3$  and  $v_4$  and optimization goals O1-O4 using BARON v.8.1.5. for a tolerance of 0.2%**

O	$k_1$	$k_2$	$k_3$	$k_4$	$k_5$	$k_6$	$X_3$	OG (%)	CPU (s)
1	0.26	5.00	4.97	0.20	0.20	0.54	8.30	0.20	136.17
2	0.20	0.24	0.22	0.20	0.21	0.20	1.10	0.00	0.06
3	0.60	5.00	5.00	0.53	0.20	0.27	5.39	0.20	96.39
4	0.99	1.15	1.00	0.96	1.00	1.00	1.10	0.00	1.42
O	$k_1$	$k_2$	$k_3$	$k_4$	$k_5$	$k_6$	$v_4$	OG (%)	CPU (s)
1	4.61	5.00	5.00	5.00	0.72	1.20	37.40	0.20	157.83
2	3.22	3.73	5.00	4.99	0.21	0.22	31.33	0.00	1.67
3	0.88	0.94	0.88	0.96	0.23	3.00	6.60	0.00	10.53
4	1.16	1.00	1.34	1.34	1.00	1.00	7.61	0.00	3.61

Table 1 only shows one solution for each particular instance. However, BARON identified in each case a set of equivalent optima (i.e, solutions with the same objective function value) involving different changes in enzyme activities, which indicates that the optimization problem is somehow degenerated. This redundancy is a consequence of the system's structure and has practical implications. As an example, we have calculated some of these equivalent points for case O1- $v_4$  using the *NumSol* option of BARON (see Figure 2). In particular, a well defined triangular region containing the changes in  $k_2$  and  $k_5$ , and  $k_1$  and  $k_2$  that lead to the same objective function value is identified. Within these regions, one can decide which combination of changes should be selected based on additional cost arguments, as they all show the same performance in terms of the predefined objective function. This region could be further reduced by imposing additional constraints to the optimization.

## 2.4 Recasting SC models into GMA models

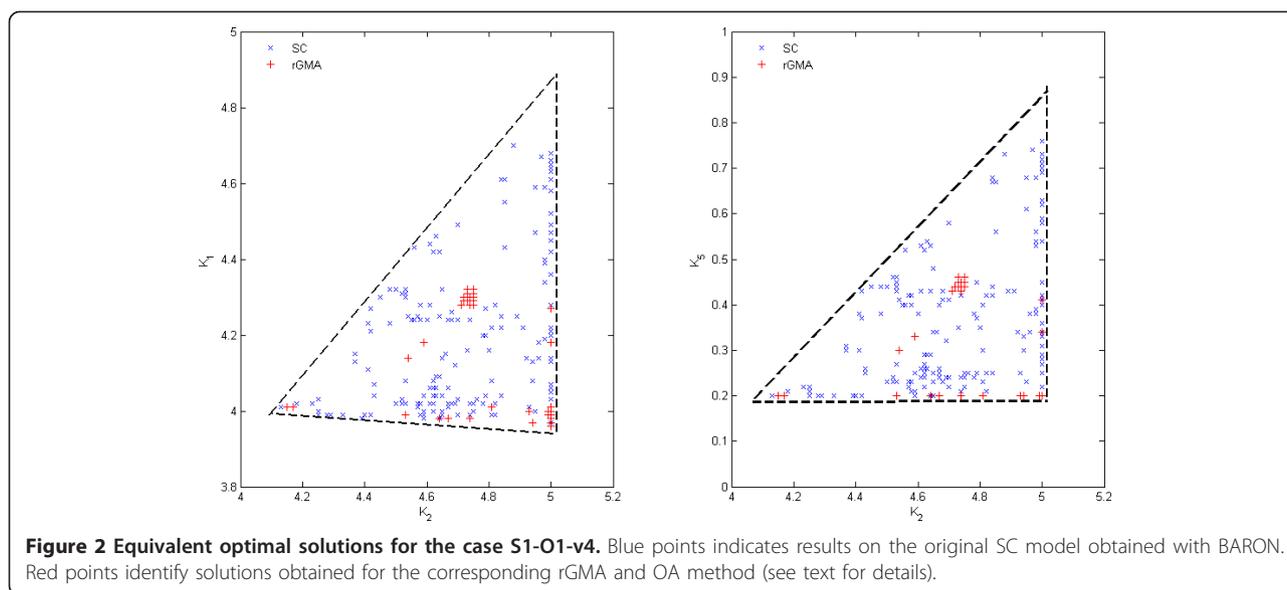
Any SC model can be recast into a GMA canonical model by introducing the auxiliary variables  $z_{rj} = K_{rj} + X_j^{n_{rj}}$ . Substitution and differentiation generates the following recast GMA (rGMA) model:

$$\dot{X}_i = \sum_{r=1}^p \mu_{ir} V_r \prod_{j=1}^{n+m} X_j^{n_{rj}} z_{rj}^{-1} \quad i = 1, \dots, n \quad (3a)$$

$$\dot{z}_{rj} = n_{rj} X_j^{n_{rj}-1} \dot{X}_j \quad \begin{matrix} r = 1, \dots, p \\ j = 1, \dots, n+m \end{matrix} \quad (3b)$$

with appropriate initial conditions  $X_j(0) = X_{j_0}$  and  $z_{rj}(0) = K_{rj} + X_{j_0}^{n_{rj}}$ .

For simulation purposes, model (3) is equivalent to the original SC model. As discussed in [36], a model recast into a GMA model has the same steady-state that the original non-linear model. The steady-state equations of the rGMA model can be expressed as:



$$\sum_{r=1}^p \mu_{ir} V_r \prod_{j=1}^{n+m} X_j^{n_{rj}} z_{rj}^{-1} = 0 \quad i = 1, \dots, n \quad (4a)$$

$$n_{rj} X_j^{n_{rj}-1} \dot{X}_j = 0 \quad \begin{matrix} r = 1, \dots, p \\ j = 1, \dots, n+m \end{matrix} \quad (4b)$$

### 2.5 Steady-state optimization of SC models through recasting

The steady-state solutions of Eqn. (4b) satisfy also Eqn. (4a). Thus, for optimization purposes, the steady-state constraints of interest are:

$$\sum_{r=1}^p \mu_{ir} V_r \prod_{j=1}^{n+m} X_j^{n_{rj}} z_{rj}^{-1} = 0 \quad i = 1, \dots, n \quad (5a)$$

$$K_{rj} + X_{j_0}^{n_{rj}} = z_{rj_0} \quad \begin{matrix} r = 1, \dots, p \\ j = 1, \dots, n+m \end{matrix} \quad (5b)$$

According to these results, the optimization problem can be stated as:

$$\begin{aligned} \min - OF & \quad OF = \{X_i \text{ or } v_r\} \\ \text{s.t.} & \\ \sum_{r=1}^p \mu_{ir} k_r v_r \prod_{j=1}^{n+m} X_j^{n_{rj}} z_{rj}^{-1} = 0 & \quad i = 1, \dots, n \\ K_{rj} + X_{j_0}^{n_{rj}} = z_{rj_0} & \quad r = 1, \dots, p \\ & \quad j = 1, \dots, n+m \\ X_{iL} \leq X_i \leq X_{iU} & \quad i = 1, \dots, n \\ k_{rL} \leq k_r \leq k_{rU} & \quad r = 1, \dots, p \\ \dots \text{ additional constraints} \dots & \end{aligned} \quad (6)$$

In our reference model, we shall consider the following constraints:

$$\begin{aligned} \min - OF & \quad OF = \{X_3, v_4\} \\ \text{s.t.} & \\ \sum_{r=1}^p \mu_{ir} k_r v_r \prod_{j=1}^{n+m} X_j^{n_{rj}} z_{rj}^{-1} = 0 & \quad i = 1, \dots, n \\ K_{rj} + X_{j_0}^{n_{rj}} = z_{rj_0} & \quad r = 1, \dots, p \\ & \quad j = 1, \dots, n+m \\ \text{Specific constraints for each optimization task} & \\ \text{(O1, O3 only)} & \quad 0.1 \leq X_i \leq 10 \quad i = 1, \dots, n \\ \text{(O1, O2, O3 only)} & \quad 0.2 \leq k_r \leq 5 \quad r = 1, \dots, p \\ \text{(O2, O4 only)} & \quad 0.9 X_i^{BAS} \leq X_i \leq 1.1 X_i^{BAS} \quad i = 1, \dots, n \\ \text{(O3 only) and (OF : } X_3 \text{ only)} & \quad 0.9 v_4^{BAS} \leq v_4 \leq 1.1 v_4^{BAS} \\ \text{(O4 only)} & \\ k_r = k_{r1} + k_{r2} + k_{r3} & \quad r = 1, \dots, p \\ k_r^{LB} \gamma_{r1} \leq k_{r1} \leq (1 - \delta) \gamma_{r1} & \quad r = 1, \dots, p \\ (1 - \delta) \gamma_{r2} \leq k_{r2} \leq (1 + \delta) \gamma_{r2} & \quad r = 1, \dots, p \\ (1 + \delta) \gamma_{r3} \leq k_{r3} \leq k_r^{LB} \gamma_{r3} & \quad r = 1, \dots, p \\ \gamma_{r1} + \gamma_{r2} + \gamma_{r3} = 1 & \quad r = 1, \dots, p \\ \sum_{r=1}^p \gamma_{r1} + \sum_{r=1}^p \gamma_{r3} \leq ME = 3 & \end{aligned} \quad (7)$$

Once the problem has been recast into a rGMA, its mathematical structure can be exploited in order to improve the efficiency of the solution procedure, as demonstrated by the authors in previous works. This problem has a GMA form except for the auxiliary constraint 5b, which is required to recast the SC into the rGMA. This constraint can be easily handled by means of relaxation techniques and exponential transformations similar to those used by the authors in their global optimization algorithms for pure GMA models [32,33]. In particular, two algorithms were developed for the global optimization of GMA models: a customized outer-approximation (OA, [30]) and a tailored spatial branch-

and-bound (sBB, [32]). The authors showed that the numerical performance of these methods depends on the specific problem being solved, and that none of them is clearly better than the other one. Here, we use the OA algorithm to solve 6, as this method proved to be faster than sBB for problems of smaller size ([32]). Again, the main body of the algorithm was coded in GAMS 23.0.2, using CPLEX 11.2.1 as MILP solver for the master subproblems and CONOPT 3.14 s as NLP solver for the slave subproblems of the algorithm. For a fair comparison, we also set a tolerance of 0.2%, the same as when using BARON.

As can be seen in Table 2, the optimization of the rGMA formulation using our customized OA yields similar results to those obtained when BARON is applied to the original SC model. In some cases, significant reductions in computational time are attained with our OA algorithm. While BARON took a total time of 407.68 CPU seconds for solving the 8 instances, the customized OA algorithm solved the same problems in 8.5 CPU seconds.

Note that the objective function values obtained with the SC and rGMA models only differ within the tolerance imposed. In some cases, discrepancies regarding the enzymatic profiles calculated are observed mainly due to the system's structure, that is, to the fact that the problem contains multiple solutions attaining the same performance in terms of objective function value but involving different enzymatic configurations, as discussed in section 2.3.

To further investigate this issue, we apply the multi-solution capability of BARON to the rGMA model (Figure 2). Again, different equivalent optima are obtained, but this time the dispersion of the equivalent solutions generated for a given case tend to concentrate either in the center or in the extremes of the region containing the solutions with the same objective function value calculated with the SC model.

**Table 2 Results for the maximization of  $X_3$  and  $v_4$  using the rGMA model and optimization goals O1-O4 using the customized OA for a tolerance of 0.2%**

O	$k_1$	$k_2$	$k_3$	$k_4$	$k_5$	$k_6$	$X_3$	OG (%)	CPU (s)
1	0.26	5.00	5.00	0.20	0.20	0.20	8.30	0.20	2.94
2	0.21	0.22	0.21	0.20	0.20	0.20	1.10	0.00	0.06
3	0.60	5.00	5.00	0.53	0.20	0.24	5.40	0.13	2.35
4	1.00	1.05	0.97	0.92	1.00	1.00	1.10	0.00	0.23
O	$k_1$	$k_2$	$k_3$	$k_4$	$k_5$	$k_6$	$v_4$	OG (%)	CPU (s)
1	3.96	5.00	5.00	5.00	0.20	2.99	37.47	0.00	0.16
2	3.22	3.55	5.00	4.99	0.20	0.21	31.33	0.17	0.66
3	0.68	1.79	1.12	1.27	0.20	0.21	6.60	0.00	0.12
4	1.16	1.00	1.34	1.34	1.00	1.00	7.61	0.11	1.98

**Table 3 Results (objective function) of the optimization of case O1-  $v_4$  for specific regions of  $k_2$  and  $k_5$  obtained with BARON for the SC model**

$k_5/k_2$	1	2	3	4	5	6	7	8
8	36.50	36.71	36.90	37.08	37.24	37.37	37.47	37.47
7	36.62	36.83	37.02	37.19	37.34	37.46	37.47	37.47
6	36.75	36.95	37.14	37.31	37.44	37.47	37.47	37.47
5	36.88	37.08	37.26	37.41	37.47	37.47	37.47	37.47
4	37.02	37.21	37.38	37.47	37.47	37.47	37.47	37.47
3	37.15	37.34	37.47	37.47	37.47	37.47	37.47	37.47
2	37.29	37.46	37.47	37.47	37.47	37.47	37.47	37.47
1	37.43	37.47	37.47	37.47	37.47	37.47	37.47	37.47

Domain of each  $k_i$  ( $4 \leq k_2 \leq 5$ ;  $0.2 \leq k_5 \leq 0.8$ ) has been split into 8 intervals with equal width.

The region illustrated in Figure 2 should not be misunderstood as a feasibility region. In fact, solutions do exist outside this region, but they lead to worse objective function values. To further clarify this issue, we consider a grid of values for  $k_2$  and  $k_5$  in the region defined by constraints  $4 \leq k_2 \leq 5$  and  $0.2 \leq k_5 \leq 0.8$ , and solve the optimization problem within each cell applying BARON to the SC model, and our OA to the rGMA model. Recall that these linear constraints define a region that contains that in Figure 2. The results obtained in this optimization are illustrated in Tables 3 and 4, and are exactly equal for both methods. However, the CPU time is much lower when using our OA algorithm applied to rGMA (11,811 CPU seconds for generating all the points with BARON applied to the SC model vs 17 CPU seconds with the customized OA applied to the rGMA model; as shown in Tables 5 and 6).

## 2.6 Difficult optimization tasks can be solved via recasting

The reference model can be optimized either by general purpose techniques or by rGMA specific methods such as the customized OA. However, even with this simple

**Table 4 Results (objective function) of the optimization of case O1- $v_4$  for specific regions of  $k_2$  and  $k_5$  obtained with the customized OA for the rGMA model**

$k_5-k_2$	1	2	3	4	5	6	7	8
8	36.50	36.71	36.90	37.08	37.24	37.37	37.47	37.47
7	36.62	36.83	37.02	37.19	37.34	37.46	37.47	37.47
6	36.75	36.95	37.14	37.31	37.44	37.47	37.47	37.47
5	36.88	37.08	37.26	37.41	37.47	37.47	37.47	37.47
4	37.02	37.21	37.38	37.47	37.47	37.47	37.47	37.47
3	37.15	37.34	37.47	37.47	37.47	37.47	37.47	37.47
2	37.29	37.46	37.47	37.47	37.47	37.47	37.47	37.47
1	37.43	37.47	37.47	37.47	37.47	37.47	37.47	37.47

Domain of each  $k_i$  ( $4 \leq k_2 \leq 5$ ;  $0.2 \leq k_5 \leq 0.8$ ) has been split into 8 intervals with equal width.

example, we may encounter instances that are hard to solve using standard techniques. Consider, for instance, the same reaction scheme as before but this time with the alternative parameters indicated in the following model:

$$\begin{aligned} \frac{dX_1}{dt} &= \frac{11.11k_1X_5^{2.86}}{X_5^{2.86} + 0.81} \\ &\quad - \frac{12.35k_2X_1^{1.54}}{(X_1^{1.54} + 0.61)X_3^{6.81} \left(0.11 + \frac{1}{X_3^{6.81}}\right)} \\ \frac{dX_2}{dt} &= \frac{12.35k_2X_1^{1.54}}{(X_1^{1.54} + 0.61)X_3^{6.81} \left(0.11 + \frac{1}{X_3^{6.81}}\right)} \\ &\quad - \frac{4.44k_3X_2^{4.14}}{X_2^{4.14} + 0.11} \\ &\quad - \frac{7.41k_5X_1^{0.51}X_2^{26.51}}{(X_1^{0.51} + 0.19)(X_2^{26.51} + 0.11)} \\ \frac{dX_3}{dt} &= \frac{4.44k_3X_2^{4.14}}{X_2^{4.14} + 0.11} - \frac{4.44k_4X_3^{4.14}}{X_3^{4.14} + 0.11} \\ \frac{dX_4}{dt} &= \frac{7.41k_5X_1^{0.51}X_2^{26.51}}{(X_1^{0.51} + 0.19)(X_2^{26.51} + 0.11)} \\ &\quad - \frac{6.67k_6X_4^{1.57}}{X_4^{1.57} + 1.40} \end{aligned} \tag{8}$$

The optimization task of interest being:

- O5: Which is the optimal pattern of changes in enzyme activities that maximize  $v_6$  in the new steady-state for a fixed value of  $X_5$  and considering the following constraints?

$$\begin{aligned} 0.3 &\leq X_1 \leq 30 \\ 0.1 &\leq X_2 \leq 10 \\ 0.1 &\leq X_3 \leq 10 \\ 0.6 &\leq X_4 \leq 50 \\ 0.1 &\leq k_r \leq 20 \quad r = 1, \dots, p \end{aligned} \tag{9}$$

When BARON is employed to solve this case using the native SC form, it cannot reduce the optimality gap

**Table 5 Results (CPU time in seconds) of the optimization of case O1- v<sub>4</sub> for specific regions of k<sub>2</sub> and k<sub>5</sub> obtained with BARON for the SC model**

k <sub>5</sub> /k <sub>2</sub>	1	2	3	4	5	6	7	8
8	212.53	308.53	185.64	201.80	222.30	201.53	139.16	178.31
7	194.81	161.16	215.80	196.81	344.73	243.02	0.03	174.81
6	234.30	203.75	147.08	180.69	328.34	254.42	304.11	280.53
5	212.08	282.41	329.33	237.34	208.02	292.27	200.00	154.62
4	288.00	160.14	92.94	235.80	172.69	147.14	56.11	150.28
3	125.56	111.17	150.27	187.52	337.97	158.16	112.66	264.12
2	239.70	190.59	100.03	138.47	106.38	205.14	119.39	246.34
1	140.42	102.12	80.45	21.69	73.12	96.61	89.94	80.03

Domain of each k<sub>r</sub>(4 ≤ k<sub>2</sub> ≤ 5; 0.2 ≤ k<sub>5</sub> ≤ 0.8) has been split into 8 intervals with equal width.

**Table 6 Results (CPU time in seconds) of the optimization of case O1-v<sub>4</sub> for specific regions of k<sub>2</sub> and k<sub>5</sub> obtained with the customized OA for the rGMA model**

k <sub>5</sub> /k <sub>2</sub>	1	2	3	4	5	6	7	8
8	0.13	0.27	0.23	0.18	0.17	0.19	0.28	0.28
7	0.26	0.28	0.28	0.26	0.28	0.23	0.32	0.25
6	0.32	0.30	0.28	0.28	0.27	0.23	0.19	0.25
5	0.31	0.21	0.25	0.25	0.26	0.28	0.27	0.29
4	0.25	0.27	0.32	0.30	0.25	0.27	0.26	0.28
3	0.20	0.22	0.28	0.28	0.29	0.30	0.19	0.53
2	0.28	0.25	0.19	0.19	0.22	0.17	0.30	0.25
1	0.23	0.24	0.26	0.27	0.23	0.21	0.24	0.31

Domain of each k<sub>r</sub>(4 ≤ k<sub>2</sub> ≤ 5; 0.2 ≤ k<sub>5</sub> ≤ 0.8) has been split into 8 intervals with equal width.

below the specified tolerance after 1 hour of CPU time. In contrast, when the model is recast into its rGMA form and our OA method is applied, the global optimum can be determined with an optimality gap of 2% in 10.95 seconds (see Table 7). This illustrates both, the utility of using the rGMA as a canonical form for dealing with the optimization of SC models, and the computational efficiency of our global optimization methods specifically designed to take advantage of the mathematical structure of the GMA.

### 3 Discussion

While experimental tools to manipulate gene expression are already available, there is no established set of guidelines on how these tools can be used to achieve a certain goal. So far, two main difficulties have prevented model driven optimization from becoming a standard in providing such guidelines: (i) the lack of information to build detailed kinetic models and (ii) the computational difficulties that arise upon the optimization of such models. The latter can be exemplified by the application of mixed integer non-linear optimization techniques (MINLP) in the context of kinetic models presented in [34,35]. In such cases, the optimization task showed to be computationally very demanding and global optimality could not be guaranteed in many cases. We propose that using models with a standardized structure may offer a solution to both problems. On one hand, approximate kinetics, such as the SC formalism, can provide very accurate approximations and retain key features of the system like saturation and cooperativity. On the other hand, these formalisms can be automatically recast

**Table 7 Results of the optimization of model 8 with BARON (SC model) and the customized OA (rGMA model)**

Solver	k <sub>1</sub>	k <sub>2</sub>	k <sub>3</sub>	k <sub>4</sub>	k <sub>5</sub>	k <sub>6</sub>	OF	OG (%)	CPU (s)
BARON (SC)	6.24	5.16	0.46	0.6	8.46	9.09	60.36	45.18	3600
OA (rGMA)	6.25	5.17	0.45	0.6	8.44	9.1	60.46	2.18	10.95

into GMA form and using efficient global optimization methods developed specifically for this canonical representation. Although this technique will certainly have limitations, our previous results indicate that it can be applied to models of moderate complexity without major problems [32]. Optimization of GMA models comprising up to 60 reactions and 40 metabolites offer no limitation to our technique. We have shown how these methods can be easily used to optimize SC via recasting into rGMA models while still being quite efficient.

Our results can be of particular interest for dealing with multicriteria optimization on realistic models. This kind of problems are relevant when exploring the adaptive response to changing conditions, where conflicting goals may be at play [39,40]. Particularly, we should notice that several multi-objective optimization techniques, such as the weighted sum or epsilon constraint methods [41] are based on solving a set of auxiliary single-objective problems. These approaches could directly benefit from the numerical advances presented in this work. This kind of problems are relevant when exploring the adaptive response to changing conditions, where conflicting goals may be on play [39,40]. The highly efficient OA algorithm applied to rGMA models provide a practical way for extending multicriteria optimization methods, for instance as used in [39], to non-linear kinetic models. It is in principle possible to make use of methods such as ours to analyze the optimality of large scale dynamic systems much in the same way that Flux Balance Analysis can be applied to analyze the stoichiometry of an organism on a genomic scale. To make this possible, however, extensive experimental and modeling efforts would be required to characterize the most important properties of the involved processes. In fact, we anticipate that practical limitations to apply the techniques presented here in solving larger problems will be dominated by the lack of information about the component processes and metabolites rather than by the technical capacity of the optimization technique presented here. Although a complete kinetic characterization of the processes in a complete metabolic network may yet be far, information on kinetic orders and saturation fractions is easier to obtain. In this context, the SC formalism provides a sound approximation that results in a mathematical representation useful for simulation and optimization through recasting.

#### 4 Conclusions

We expect that the possibility of building models using non-linear approximate formalisms and of subsequently optimizing these models will trigger interest in the experimental characterization of the components of cellular metabolism. After the genomic explosion, we need

to step back and begin to measure enzyme activities, metabolite levels, and regulatory signals on a larger scale than we used to do before, if we want to understand the emergence of the dynamic properties of biological systems and to be able to develop successful biotechnological applications.

## 5 Methods

### 5.1 Modelling strategies

The process of model building and optimization can be used to understand how a system should be changed in order to achieve specific biotechnological goals or how the same system has evolved in order to more efficiently execute a given biological function. Different trade-offs are considered during the modeling process. On the one hand, one wants to use models that are as simple as possible to guarantee numerical tractability. Unfortunately simplifications may lead to models whose accuracy is only ensured for a limited range of physiological conditions. On the other hand, models that are very detailed and accurate over a wide range of physiological conditions are typically more difficult to analyze and optimize. Needless to say, the type of modeling strategy and the model one chooses to implement have a large impact on the results of the analysis. The most widely used strategies in the context of optimization are: (1) Stoichiometric models, (2) Kinetic models, and (3) Approximated models.

The three strategies have as a starting point a set of ordinary differential equations, in which the dependent variables or nodes are the chemical species whose dynamical behavior one is interested in studying. For a system with  $n$  dependent variables,  $p$  processes and  $m$  independent variables, the node equations are written as follows:

$$\dot{X}_i = \sum_{r=1}^p \mu_{ir} v_r \quad i = 1, \dots, n \quad (10)$$

$\mu_{ir}$  stands for the stoichiometry of each metabolite  $X_i$  in each reaction  $r$  with respect to metabolite  $i$  and can be derived from the reaction scheme.

At this stage, the various strategies begin to differ in the way that they implement and analyze the equations. Typically, Flux balance analysis (FBA) and related techniques consider only the steady state behavior of the system, and treat  $v_r$  as a variable whose value can be changed in order to optimize specific steady state constraints. To accomplish this, FBA-like methods attempt to find solutions for the following system of linear equations:

$$0 = \sum_{r=1}^p \mu_{ir} v_r \quad i = 1, \dots, n \quad (11)$$

This system of equations is solved under different assumptions. A typical problem is that of understanding the effect of knocking out different genes from the system. This analysis can be performed by setting  $v_r = 0$  for the process(es) that depend on the product of the genes that are knocked out. Once these constraints are set, linear optimization techniques can be used to identify the region of the variable space that satisfies the steady state and optimizes at the same time a set of specific measurable aspects of the systems [42-44]. It must be noted that FBA analysis of Eqn. (11) does not account for the regulatory effects that can result from gene knockout and it cannot be used to predict changes in metabolic concentrations and temporal responses. Thus, optimization constraints are limited to steady-state fluxes [15].

To overcome these limitations, we must use more complex kinetic models where the effect of changing the values of the variables on the fluxes is taken into account. This requires defining a functional form for each  $v_r$  in Eqn. (10). Often, this functional form is drawn from a number of classical enzyme kinetic rate-laws. As a result, we use an approximate expression for the kinetic behavior of each elementary process whose form depends on the underlying mechanism of the process. The reason for this is that the classical rate laws are rational functions of the variables and they are built upon different types of simplifying assumptions on the detailed mechanism of the reactions. Such assumptions range from considering that the elementary chemical steps of the catalytic process occur at very different timescales to assuming that the concentration of the catalyst and of the reactants differ in orders of magnitude. Thus, rate laws such as the popular Michaelis-Menten are approximations to the actual mechanism in specific conditions. However, more often than not, one does not have enough information to judge if such conditions meet those one is trying to model. Thus, using rational enzyme kinetics in models lacks a sound theoretical ground. In fact, within the complex architecture of the intracellular milieu, many of the assumptions that justify these classical rate-laws may not hold [45-47]. Even in the best case scenario where a detailed kinetic model using classical enzyme kinetics can be derived and numerically identified, it may be hard to globally optimize that model using general purpose algorithms. As we will show here, available optimization techniques may fail to solve fairly trivial optimization tasks even in simple models. These numerical difficulties can be overcome by defining reformulated models based on canonical representations that are easier to handle using customized global optimization algorithms devised for specific canonical functional forms.

As an alternative, theoretically well supported canonical representations can be derived using approximation theory. One type of such representations are power-law models. In a power-law model, each  $v_r$  in Eqn. (10) is approximated as [19,21]:

$$v_r(X_1, \dots, X_n, \dots, X_{n+m}) \approx \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \quad r = 1, \dots, p \quad (12)$$

This approximation is derived at a given operating point  $(X_{1,0}, X_{2,0}, \dots, X_{(n+m),0})$  as a first-order Taylor series representation of the target function in log-log space. This approximation can generate models with different representations. The two that are most commonly used are the S-system representation and the GMA representation. The S-system representation is obtained by lumping the various processes that contribute to the synthesis of a given metabolite into a global process of synthesis  $V_i^+$  and those that contribute to the utilization of a given metabolite into a global degradation process

$$\begin{aligned} \dot{X}_i &= \sum_{r=1}^p \mu_{ir} v_r \\ &= \sum_{r=1}^p \mu_{ir}^+ v_r - \sum_{r=1}^p \mu_{ir}^- v_r \quad ; \\ &= V_i^+ - V_i^- \quad i = 1, \dots, n \\ \dot{X}_i &= \sum_{r=1}^p \mu_{ir} v_r \\ &= \sum_{r=1}^p \mu_{ir}^+ v_r - \sum_{r=1}^p \mu_{ir}^- v_r \\ &= V_i^+ - V_i^- \quad i = 1, \dots, n \end{aligned} \quad (13)$$

Then, the aggregated processes are represented by power-law functions:

$$\begin{aligned} \dot{X}_i &= \alpha_i \prod_{j=1}^{n+m} X_j^{g_{ij}} - \\ &\quad \beta_i \prod_{j=1}^{n+m} X_j^{h_{ij}} \quad i = 1, \dots, n \end{aligned} \quad (14)$$

Alternatively, the GMA form is obtained representing each individual  $v_r$  as a power-law:

$$\begin{aligned} \dot{X}_i &= \sum_{r=1}^p \mu_{ir} v_r \\ &= \sum_{r=1}^p \left( \mu_{ir} \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \right) \quad i = 1, \dots, n \end{aligned} \quad (15)$$

The parameters in these representations have a clear physical interpretation. Kinetic orders, the exponents in the power-laws, are local sensitivities of the fluxes, either individual ( $f_{rj}$  for  $v_r$ ) or aggregated ( $g_{ij}$  for  $V_i^+$  and  $h_{ij}$  for  $V_i^-$ ), with respect to  $X_j$ . Rate-constants ( $\alpha_i$ ,  $\beta_i$  and  $\gamma_r$ ) are parameters that are computed so that the flux in the model at steady state is equal to the operating flux at the operating point for the metabolites. Parameter estimation techniques have been developed so that power-law parameters can be calculated from experimental measurements [13]. It should also be noted that the use of estimation procedures (i.e., least-squares), alternate regression or similar procedures to estimate power-law parameters from dynamic curves lead to a power-law representation that is no longer local according to the classical definition [48-50]. Those models may, by definition, slightly improve their accuracy over strictly local models.

To complement the power-law approach, the Saturable and Cooperative (SC) formalism was introduced by Sorribas et al. [23] as an extension of the ideas that led to the power-law formalism. The SC representation of a given velocity is:

$$v_r(X_1, \dots, X_n, \dots, X_{n+m}) \approx \frac{V_r \prod_{j=1}^{n+m} X_j^{n_{rj}}}{\prod_{j=1}^{n+m} (K_{rj} + X_j^{n_{rj}})} \quad (16)$$

This representation can be obtained from a power-law model defined at a given operating point  $X_0 = (X_{10}, \dots, X_{(n+m)0})$  through the following relationships:

$$n_{rj} = \frac{f_{rj}}{1 - p_{rj}} \quad r = 1, \dots, p \quad j = 1, \dots, n + m \quad (17)$$

$$K_{rj} = \frac{1 - p_{rj}}{p_{rj}} X_{j0}^{n_{rj}} \quad r = 1, \dots, p \quad j = 1, \dots, n + m \quad (18)$$

Thus SC uses the same information as the power-law except for the new parameters  $p_{rj}$  (saturation fractions), which are defined as:

$$p_{rj} = v_{r0} / V_{rj} \quad r = 1, \dots, p \quad j = 1, \dots, n + m \quad (19)$$

where  $v_{r0} = v_r(X_{10}, \dots, X_{n0}, \dots, X_{(n+m)0})$  and  $V_{rj}$  is either the limit velocity (saturation) when  $X_j \rightarrow \infty$  if  $n_{rj} > 0$ , or the limit velocity when  $X_j \rightarrow 0$  if  $n_{rj} < 0$ .

Using SC models for global optimization can raise some numerical issues. These difficulties can be avoided to a large extent by recasting SC models into a

canonical GMA model, through the introduction of auxiliary variables, as will be shown in the next section.

## 5.2 Recasting non-linear models into power-law canonical models by increasing the number of variables

Non-linear models can be *exactly* recast into GMA or S-system models through the use of auxiliary variables [36]. As a result, the final model is an exact representation of the original model, written in a canonical form. In other words, the resulting GMA model is not an approximation to the original model: it is an exact replica of it. To avoid confusion, hereafter, we refer to a GMA model that exactly recasts another as an rGMA model.

As a very simple introductory example, consider a linear pathway with two internal metabolites  $X_1$  and  $X_2$  and a source metabolite  $X_3$  (Figure 3). In this pathway,  $X_2$  is a competitive inhibitor of the synthesis of  $X_1$  from the source metabolite. A generic model using Michaelis-Menten kinetic functions, assuming a competitive inhibition of the first reaction by  $X_2$ , can be written as:

$$\begin{aligned} \dot{X}_1 &= \frac{V_1 X_3}{K_1(1 + K_i/X_2) + X_3} - \frac{V_2 X_1}{K_2 + X_1} \\ \dot{X}_2 &= \frac{V_2 X_1}{K_2 + X_1} - \frac{V_3 X_2^2}{K_3^2 + X_2^2} \end{aligned} \quad (20)$$

in which  $X_3$  is an externally fixed variable.

Recasting this model as a rGMA can be done as follows. First, let us define three new variables:

$$\begin{aligned} X_4 &= K_1(1 + K_i/X_2) + X_3 \\ X_5 &= K_2 + X_1 \\ X_6 &= K_3^2 + X_2^2 \end{aligned} \quad (21)$$

We can now write the model in 20 as:

$$\begin{aligned} \dot{X}_1 &= V_1 X_3 X_4^{-1} - V_2 X_1 X_5^{-1} \\ \dot{X}_2 &= V_2 X_1 X_5^{-1} - V_3 X_2^2 X_6^{-1} \end{aligned} \quad (22)$$

with initial conditions  $X_1(0) = X_{10}$  and  $X_2(0) = X_{20}$ .

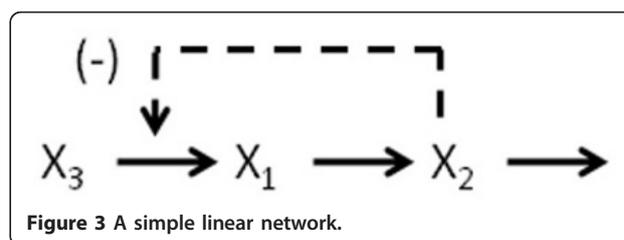


Figure 3 A simple linear network.

To complete the recasting we must now provide the equations that follow the change in the new variables over time. These are given by the following equations:

$$\begin{aligned}\dot{X}_4 &= -\frac{K_1 K_i \dot{X}_2}{X_2^2} \\ &= V_3 K_1 K_i X_6^{-1} - V_2 K_1 K_i X_1 X_5^{-1} \\ X_2^{-2} \dot{X}_5 &= \dot{X}_1 \\ &= V_1 X_3 X_4^{-1} - V_2 X_1 X_5^{-1} \\ \dot{X}_6 &= 2X_2 \dot{X}_2 \\ &= 2V_2 X_1 X_2 X_5^{-1} - 2V_3 X_2^3 X_6^{-1}\end{aligned}\quad (23)$$

with initial conditions  $X_4(0) = K_1(1 + K_i/X_{2_0}) + X_{3_0}$ ,  $X_5(0) = K_2 + X_{1_0}$ , and  $X_6(0) = K_3^2 + X_{2_0}^2$ .

The resulting rGMA model (22-23) is an exact representation of model in (20). Hence, for a set of appropriate initial conditions, the simulation of the dynamic response using either the model recast as a rGMA or the original model will produce the same trajectory. In principle, any non-linear model can be recast into a rGMA following a similar procedure [36]. This can be extremely useful, because it allows for the application of tailored global optimization procedures originally devised for GMA models [28-30,32,51,52] to generic non-linear models.

#### Acknowledgements

AS is funded by MICINN (Spain) (BFU2008-0196). RA is partially supported by MICINN (Spain) through Grants BFU2007-62772/BMC and BFU2010-17704. AS and RA are members of the 2009SGR809 research group of the Generalitat de Catalunya. GG-G and CP acknowledges support from the Spanish Ministry of Science and Innovation (Projects DPI2008-04099 and CTQ2009-14420-C02-01) and the Spanish Ministry of External Affairs and Cooperation (Projects A/023551/09 and A/031707/10).

#### Author details

<sup>1</sup>Departament d'Enginyeria Química, Universitat Rovira i Virgili, Avinguda Països Catalans 26, 43007-Tarragona, Spain. <sup>2</sup>Technische Universität München. Fachgebiet für Systembiotechnologie, Boltzmannstr. 15 85748 Garching, Germany. <sup>3</sup>Departament de Ciències Mèdiques Bàsiques, Institut de Recerca Biomèdica de Lleida (IRBLLEIDA), Universitat de Lleida, Montserrat Roig 2, 25008 Lleida, Spain.

#### Authors' contributions

AM-S suggested the potential utility of recasting for optimizing non-linear kinetic models. AS and AM-S elaborate on the recasting of SC models and planned the work. CP, GG-G and LJ implemented the OA algorithm and worked out the technical solution for applying it to a rGMA model. CP and GG-G performed the optimization tasks. AS and RA defined the reference model and obtained the numerical parameters used in the paper. All authors read and approved the final manuscript.

#### Competing interests

The authors declare that they have no competing interests.

Received: 5 April 2011 Accepted: 25 August 2011  
Published: 25 August 2011

#### References

1. Voit EO: Optimization in integrated biochemical systems. *Biotechnol Bioeng* 1992, **40**(5):572-82.
2. Alvarez-Vasquez F, Gonzalez-Alcon C, Torres NV: Metabolism of citric acid production by *Aspergillus niger*: model definition, steady-state analysis and constrained optimization of citric acid production rate. *Biotechnol Bioeng* 2000, **70**:82-108.
3. Banga JR: Optimization in computational systems biology. *BMC Syst Biol* 2008, **2**:47.
4. Rodriguez-Prados JC, de Atauri P, Maury J, Ortega F, Portais JC, Chassagnole C, Akerenza L, Lindley ND, Cascante M: In silico strategy to rationally engineer metabolite production: A case study for threonine in *Escherichia coli*. *Biotechnology and bioengineering* 2009, **103**(3):609-620.
5. Scheer M, Grote A, Chang A, Schomburg I, Munaretto C, Rother aSC M, Stelzer M, JT, DS: BRENDA, the enzyme information system in 2011. *Nucleic Acids Res* 2011, **39**:670-676.
6. Rojas I, Golebiewski M, Kania R, Krebs O, Mir S, Weidemann A, Wittig U: SABIO-RK: a database for biochemical reactions and their kinetics. *BMC Systems Biology* 2007, **1**:56.
7. Shiraiishi F, Savageau MA: The tricarboxylic acid cycle in *Dictyostelium discoideum*. I. Formulation of alternative kinetic representations. *The Journal of biological chemistry* 1992, **267**(32):22912-22918.
8. Shiraiishi F, Savageau MA: The tricarboxylic acid cycle in *Dictyostelium discoideum*. II. Evaluation of model consistency and robustness. *The Journal of biological chemistry* 1992, **267**(32):22919-22925.
9. Shiraiishi F, Savageau MA: The tricarboxylic acid cycle in *Dictyostelium discoideum*. III. Analysis of steady state and dynamic behavior. *The Journal of biological chemistry* 1992, **267**(32):22926-22933.
10. Shiraiishi F, Savageau MA: The tricarboxylic acid cycle in *Dictyostelium discoideum*. IV. Resolution of discrepancies between alternative methods of analysis. *The Journal of biological chemistry* 1992, **267**(32):22934-22943.
11. Fonseca LL, Sanchez C, Santos H, Voit EO: Complex coordination of multi-scale cellular responses to environmental stress. *Molecular bioSystems* 2010.
12. Sorribas A, Cascante M: Structure identifiability in metabolic pathways: parameter estimation in models based on the power-law formalism. *The Biochemical journal* 1994, **298**(Pt 2):303-311.
13. Chou IC, Voit EO: Recent developments in parameter estimation and structure identification of biochemical and genomic systems. *Mathematical Biosciences* 2009, **219**(2):57-83.
14. Grossmann I, Biegler LT: Part II. Future perspective on optimization. *Computers and Chemical Engineering* 2004, **28**:1193-1218.
15. Terzer M, Maynard ND, Covert MW, Stelling J: Genome-scale metabolic networks. *Wiley interdisciplinary reviews. Systems biology and medicine* 2009, **1**(3):285-297.
16. Gianchandani EP, Chavali AK, Papin JA: The application of flux balance analysis in systems biology. *Wiley interdisciplinary reviews. Systems biology and medicine* 2010, **2**(3):372-382, [UID: 101516550; ppublish].
17. Voit EO: Design principles and operating principles: the yin and yang of optimal functioning. *Math Biosci* 2003, **182**:81-92.
18. Jamshidi N, Palsson BO: Formulating genome-scale kinetic models in the post-genome era. *Molecular systems biology* 2008, **4**:171.
19. Savageau MA: Biochemical systems analysis. I. Some mathematical properties of the rate law for the component enzymatic reactions. *Journal of theoretical biology* 1969, **25**(3):365-369.
20. Savageau MA: Biochemical systems analysis. II. The steady-state solutions for an n-pool system using a power-law approximation. *Journal of theoretical biology* 1969, **25**(3):370-379.
21. Savageau MA: Biochemical systems analysis. 3. Dynamic solutions using a power-law approximation. *Journal of theoretical biology* 1970, **26**(2):215-226.
22. Savageau MA: *Biochemical Systems Analysis: A Study of Function and Design in Molecular Biology* Reading, Mass.: Addison-Wesley; 1976.
23. Sorribas A, Hernandez-Bermejo B, Vilaprinyo E, Alves R: Cooperativity and saturation in biochemical networks: a saturable formalism using Taylor series approximations. *Biotechnology and bioengineering* 2007, **97**(5):1259-1277.

24. Liebermeister W, Klipp E: **Bringing metabolic networks to life: convenience rate law and thermodynamic constraints.** *Theoretical biology & medical modelling* 2006, **3**:41.
25. Goel G, Chou IC, Voit EO: **System Estimation from Metabolic Time Series Data.** *Bioinformatics (Oxford, England)* 2008.
26. Ni T, Savageau M: **Model assessment and refinement using strategies from biochemical systems theory: Application to metabolism in human red blood cells.** *Journal of Theoretical Biology* 1996, **179**(4):329-368.
27. Arkin A, Ross J, McAdams H: **Stochastic kinetic analysis of developmental pathway bifurcation in phage lambda-infected Escherichia coli cells.** *Genetics* 1998, **149**(4):1633-1648.
28. Polisetty PK, Gatzke EP, Voit EO: **Yield optimization of regulated metabolic systems using deterministic branch-and-reduce methods.** *Biotechnol Bioeng* 2008, **99**(5):1154-69.
29. Guillén-Gosálbez G, Sorribas A: **Identifying quantitative operation principles in metabolic pathways: A systematic method for searching feasible enzyme activity patterns leading to cellular adaptive responses.** *BMC Bioinformatics* 2009, **10**.
30. Guillén-Gosálbez G, Pozo C, Jiménez L, Sorribas A: **A global optimization strategy to identify quantitative design principles for gene expression in yeast adaptation to heat shock.** *Computer Aided Chemical Engineering* 2009, **26**:1045-1050.
31. Voit EO: *Computational Analysis of Biochemical Systems. A Practical Guide for Biochemists and Molecular Biologists* Cambridge, U.K.: Cambridge University Press; 2000.
32. Pozo C, Guillén-Gosálbez G, Sorribas A, Jiménez L: **A Spatial Branch-and-Bound Framework for the Global Optimization of Kinetic Models of Metabolic Networks.** *Industrial and Engineering Chemistry Research* 2010.
33. Sorribas A, Pozo C, Vilaprinyo E, Guillén-Gosálbez G, Jiménez L, Alves R: **Optimization and evolution in metabolic pathways: global optimization techniques in Generalized Mass Action models.** *Journal of Biotechnology* 2010, **149**(3):141-153.
34. Chassagnole C, Noisommit-Rizzi N, Schmid J, Mauch K, Reuss M: **Dynamic modeling of the central carbon metabolism of Escherichia coli.** *Biotechnol Bioeng* 2002, **79**:53-73.
35. Nikolaev E: **The elucidation of metabolic pathways and their improvements using stable optimization of large-scale kinetic models of cellular systems.** *Metab Eng* 2010, **12**:26-38.
36. Savageau MA, Voit EO: **Recasting nonlinear differential equations as S-systems: a canonical nonlinear form.** *Mathematical Biosciences* 1987, **87**:83-115.
37. Voit EO: **Recasting nonlinear models as S-systems.** *Mathematical and Computer Modelling* 1988, **11**(C):140-145.
38. Smith EMB, Pantelides CC: **A symbolic reformulation/spatial branch-and-bound algorithm for the global optimisation of nonconvex MINLPs.** *Computers and Chemical Engineering* 1999, **23**(4-5):457-478.
39. Vera J, de Atauri P, Cascante M, Torres NV: **Multicriteria optimization of biochemical systems by linear programming: application to production of ethanol by *Saccharomyces cerevisiae*.** *Biotechnol Bioeng* 2003, **83**(3):335-43.
40. Vilaprinyo E, Alves R, Sorribas A: **Use of physiological constraints to identify quantitative design principles for gene expression in yeast adaptation to heat shock.** *BMC Bioinformatics* 2006, **7**:184.
41. Ehrgott M: *Multicriteria Optimization* Springer; 2005.
42. Famili I, Forster J, Nielsen J, Palsson BO: ***Saccharomyces cerevisiae* phenotypes can be predicted by using constraint-based analysis of a genome-scale reconstructed metabolic network.** *Proc Natl Acad Sci USA* 2003, **100**(23):13134-9.
43. Famili I, Mahadevan R, Palsson BO: **k-Cone analysis: determining all candidate values for kinetic parameters on a network scale.** *Biophys J* 2005, **88**(3):1616-25.
44. Price ND, Papin JA, Schilling CH, Palsson BO: **Genome-scale microbial in silico models: the constraints-based approach.** *Trends Biotechnol* 2003, **21**(4):162-9.
45. Savageau MA: **Influence of fractal kinetics on molecular recognition.** *Journal of Molecular Recognition: JMR* 1993, **6**(4):149-157.
46. Savageau MA: **Michaelis-Menten mechanism reconsidered: implications of fractal kinetics.** *Journal of theoretical biology* 1995, **176**:115-124.
47. Savageau MA: **Development of fractal kinetic theory for enzyme-catalysed reactions and implications for the design of biochemical pathways.** *Bio Systems* 1998, **47**(1-2):9-36.
48. Hernandez-Bermejo B, Fairen V, Sorribas A: **Power-law modeling based on least-squares minimization criteria.** *Mathematical biosciences* 1999, **161**(1-2):83-94.
49. Hernandez-Bermejo B, Fairen V, Sorribas A: **Power-law modeling based on least-squares criteria: consequences for system analysis and simulation.** *Mathematical biosciences* 2000, **167**(2):87-107.
50. Alves R, Vilaprinyo E, Hernandez-Bermejo B, Sorribas A: **Mathematical formalisms based on approximated kinetic representations for modeling genetic and metabolic pathways.** *Biotechnology and Genetic Engineering Reviews* 2008, **25**:1-40.
51. Marin-Sanguino A, Torres NV: **Optimization of biochemical systems by linear programming and general mass action model representations.** *Math Biosci* 2003, **184**(2):187-200.
52. Marin-Sanguino A, Voit EO, Gonzalez-Alcon C, Torres NV: **Optimization of biotechnological systems through geometric programming.** *Theor Biol Med Model* 2007, **4**:38.

doi:10.1186/1752-0509-5-137

Cite this article as: Pozo et al.: Steady-state global optimization of metabolic non-linear dynamic models through recasting into power-law canonical models. *BMC Systems Biology* 2011 5:137.

**Submit your next manuscript to BioMed Central and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
www.biomedcentral.com/submit



# Identifying the preferred subset of enzymatic profiles in nonlinear kinetic metabolic models via multiobjective global optimization and Pareto filters

C. Pozo<sup>1</sup>, G. Guillén-Gosálbez<sup>1,\*</sup>, A. Sorribas<sup>2</sup> L. Jiménez<sup>1</sup>

**1** Departament d'Enginyeria Química (EQ), Escola Tècnica Superior d'Enginyeria Química (ETSEQ), Universitat Rovira i Virgili (URV), Tarragona, Spain

**2** Departament de Ciències Mèdiques Bàsiques, Institut de Recerca Biomèdica de Lleida (IRBLLEIDA), Universitat de Lleida, Lleida, Spain

\* E-mail: gonzalo.guillen@urv.cat

## Abstract

Optimization models in metabolic engineering and systems biology focus typically on optimizing a unique criterion, usually the synthesis rate of a metabolite of interest or the rate of growth. Connectivity and nonlinear regulatory effects, however, make it necessary to consider multiple objectives in order to identify useful strategies that balance out different metabolic issues. This is a fundamental aspect, as optimization of maximum yield in a given condition may involve unrealistic values in other key processes. Due to the difficulties associated with detailed non-linear models, analysis using stoichiometric descriptions and linear optimization methods have become rather popular in systems biology. However, despite being useful, these approaches fail in capturing the intrinsic nonlinear nature of the underlying metabolic systems and the regulatory signals involved. Targeting more complex biological systems requires the application of global optimization methods to non-linear representations. In this work we address the multi-objective global optimization of metabolic networks that are described by a special class of models based on the power-law formalism: the generalized mass action (GMA) representation. Our goal is to develop global optimization methods capable of efficiently dealing with several biological criteria simultaneously. In order to overcome the numerical difficulties of dealing with multiple criteria in the optimization, we propose a heuristic approach based on the epsilon constraint method that reduces the computational burden of generating a set of Pareto optimal alternatives, each achieving a unique combination of objectives values. To facilitate the post-optimal analysis of these solutions and narrow down their number prior to being tested in the laboratory, we explore the use of Pareto filters that identify the preferred subset of enzymatic

profiles. We demonstrate the usefulness of our approach by means of a case study that optimizes the ethanol production in the fermentation of *Saccharomyces cerevisiae*.

## Introduction

Genetic manipulation of microorganisms for obtaining improved strains involves expensive and time consuming experiments that have typically relied on trial-and-error mutagenesis and selection of promising variants. Nowadays, mathematical models of cell metabolism and gene regulation circuits have become reliable enough for metabolic engineering applications [1–3]. These models can be coupled with optimization algorithms in order to identify the most promising genetic manipulations leading to an enhanced phenotype in a given microorganism. This approach requires defining a suitable objective function, for instance the maximum yield or flux of interest. Optimization is then performed by considering the model equations describing the microorganisms' metabolism and a set of constraints relevant for cell viability [4–8]. This method provides, a sound theoretical basis for experimentalists on the best strategies for manipulating the biological system, either by changing enzyme levels through genetic manipulations or by altering environmental conditions [9].

The selection of an appropriate mathematical model is a crucial step towards success in this field. Two main strategies can be followed at this stage. On the one hand, one can choose mathematical simplicity and a genome-wide scope. In this context, flux balance analysis (FBA) provides an appropriate solution. This method makes use of stoichiometric models to represent the metabolic networks, which gives rise to mixed-integer linear formulations (MILP) that are easy to solve with standard techniques [10]. This MILP approach, however, fails at capturing the regulatory loops existing in metabolic networks [11]. On the other hand, one can choose a kinetic detailed description, which necessarily will be limited to relatively few pathways at a time. Detailed kinetic models can deal with all kind of regulatory signals and reaction mechanisms, but involve nonlinear equations (e.g., Michaelis-Menten, Hill or power-law, etc.) required to appropriately represent the reaction rates as a function of the involved metabolite concentrations. These nonlinearities give rise to nonconvexities which in turn lead to the potential existence of multiple local optima (i.e., multimodality). This may prevent standard algorithms from identifying the global optimum, as they can get trapped in local wells during the search. Global optimization strategies overcome this limitation, guaranteeing convergence to the global optimum within a desired tolerance. It should be

emphasized that global optimization is of paramount importance in these theoretical biological studies since misidentifying a local optimum as the global one may lead to spurious conclusions [12, 13].

For S-Systems models, a particular class of power-law models, Voit [4] proposed a reformulation strategy based on a logarithmic transformation that brings the model to an LP/MILP form, making it possible to apply standard optimization methods that ensure global optimality. This reformulation cannot be applied to other non-linear models, such as GMA models or detailed kinetic models. These last models must be tackled though using global optimization methods. One such method for GMA models based on an outer approximation algorithm was proposed by Polisetty et al. [8]. Guillén-Gosálbez and Sorribas [12] presented further developments using an outer approximation-based algorithm [14] and related advanced strategies [12, 15] to globally optimize GMA models. These methods have been recently extended further to deal with detailed kinetic models through a mathematical reformulation framework termed recasting that converts them into GMA models [13].

Biotechnology studies typically seek optimizing a single flux in the metabolic network as unique criterion. In practice, however, there are other criteria of interest for experimentalists, such as minimizing the number of enzymatic changes, metabolic concentration of intermediates [16] or transient times [17]. Despite the importance of such additional criteria, the majority of works in metabolic engineering are based on single-objective formulations. Although some of these functional criteria can be treated as constraints ensuring cell viability, they should be treated as additional objectives [18]. This would eventually allow for the identification of solutions in which cell viability is further improved at the expense of marginal reductions in other objectives such as growth.

The importance of multiobjective optimization in metabolic studies has been pointed out by several authors [19–21]. Technically, the solution of a multiobjective optimization (MOO) problem is given by a set of points known as the Pareto set. All these solutions feature the property that it is not possible to find another one that improves any of them in one objective without worsening at least one of the others (see Figure 1). Because of the presence of continuous variables, optimization problems arising in metabolic engineering may have an infinite number of Pareto-optimal solutions. Clearly, testing all these alternatives in the laboratory would be prohibitive in terms of time and resources. Multi-criteria decision-making (MCDM) can be of great help at this stage to rank and/or screen alternatives, ruling out the less promising and keeping the best. Unfortunately, the complexity of both, MOO and MCDM, increases with the number of objectives. In practice, the visualization and analysis of the Pareto set

becomes highly difficult in problems with more than three objectives. The need for advanced methods to support these tasks in biochemical systems has already been acknowledged [21,22].

Several approaches have been proposed for identifying a subset of Pareto solutions of special interest for decision-makers. For instance, Branke et al. [23] and later Deb [24] suggested either to specify the extreme pair-wise trade-off information about objectives or to attach relative weights to them, in order to concentrate the search in a particular region of the Pareto set. Branke and Deb [25] proposed a projection-based method to obtain a biased distribution of Pareto solutions. Farina and Amato [26] introduced a more restrictive dominance concept that produces less number of Pareto solutions. Branke et al. [27] introduced a method for obtaining those Pareto solutions with a significantly different slope (i.e., "knee" solutions). Deb and Gupta [28] focused on identifying robust (i.e. less sensitive to parameter changes) solutions. The concept of Pareto filter was also employed by several authors for eliminating non-Pareto or locally optimal Pareto solutions [29–33].

MOO and MCDM have been extensively studied in the context of a wide variety of engineering problems (for instance, refer to [34]). In contrast, their application to metabolic engineering has been quite scarce [35]. In this work, we address the MOO of metabolic networks. Our study assumes a GMA model of the target metabolic network where all model parameters are known. These include the stoichiometric coefficients of the reactions involved in the production/consumption of each internal metabolite; and the parameters of the power-law formalism that model the kinetics of each reaction at the basal state. Then, we will seek the optimization of a given flux assuming two important complementary objectives: (i) We assume that any increment in gene expression is a limiting factor for the cell as it involves an important metabolic burden; (ii) We also consider that an excessive increment in intermediate concentrations compromises cell viability. These two criteria will be used as complementary objectives that should be minimized when possible.

Under these conditions, we aim to develop a systematic framework to (i) calculate the Pareto front of the kinetic metabolic model in this multi-objective problem and (ii) identify from it a small enough set of the most promising changes in enzyme activity to be tested in the laboratory. In other words, the goal of this analysis is to determine a set containing the preferred enzymatic profiles that optimize the synthesis rate of a metabolite at minimum cost (minimum number of changes in these activities, i.e. minimum change in gene expression) and minimum increase in the concentration of intermediate metabolites.

Note that there are two main difficulties associated with the identification of such set. First, we need

to solve a high dimensional non-convex multiobjective optimization problem in which several criteria must be simultaneously minimized. This problem is challenging not only because of the high number of objectives, but also due to the existence of non-convexities. Second, even if a sufficiently large number of Pareto solutions can be identified, there is still the issue of analyzing and interpreting them, in order to keep the most promising for further evaluation in the laboratory. Deb and Saxena [36] reviewed the main difficulties associated with the calculation and analysis of the Pareto solutions of MOO problems with large number of objectives, like those arising in metabolic engineering. As will be shown later in the paper, our systematic approach allows overcoming some of these difficulties.

In particular, our strategy relies on the combined use of multiobjective global optimization and Pareto filters, which are both applied to metabolic networks described using the GMA formalism. The method presented builds upon our global optimization framework for single-objective models of metabolic networks [14, 37], which is adequately modified herein to handle multiple objectives. This method is based on an outer approximation algorithm that decomposes the target problem into a master MILP and a slave NLP, which respectively provide lower bounds (LB) and upper bounds (UB) on the global optimum. These bounds tend to approach as iterations proceed until a given tolerance is satisfied.

Note that our methodology shares some common features with that presented by [35] for S-Systems models. However, while the former strategy ends with the generation of the Pareto optimal front, ours goes one step beyond by suggesting a subset of preferred alternatives that are identified using Pareto filters. Hence, this work presents advances in two main fronts: (i) the generation of Pareto optimal solutions for multiobjective GMA models, and (ii) the identification of the most promising alternatives using systematic filters.

The capabilities of the proposed methodology are illustrated in the optimization of the fermentation of *Saccharomyces cerevisiae* considering 14 objectives. This process has been already studied in the past by several authors. For instance, Sendín et al. [35] used an ad-hoc model of this metabolic pathway to address by means of different MOO methods a 6-objective MOO problem considering the ethanol synthesis rate and the concentration of 5 dependent metabolites. Most of the approaches compared therein show some limitations, as they either rely on local solvers (this is the case of weighted sum, attainment goal and NBI) or employ stochastic optimization methods (MOEA) that are unable to guarantee convergence to the global optimum in a finite number of iterations, which may result in a spurious Pareto front. The other method studied in that work (MIOM) requires the transformation of the original model into

an S-Systems representation, which is something unnecessary when relying directly on GMA models. Furthermore, we address here a more complex problem that accounts for 14 objectives (the fold-change in 8 different enzyme activities, expressed as the absolute value of the natural logarithm of the enzyme activity fold-change; the concentration of 5 dependent metabolites; and the ethanol synthesis rate). This represents a significant advance compared to traditional biotechnological approaches that maximize the ethanol yield and impose biological constraints for maintaining metabolites and enzymes levels around their basal state so as to preserve cell homeostasis [9].

## Results

In order to illustrate the capabilities of our approach we solved a case study that optimizes the ethanol production in the fermentation of *Saccharomyces cerevisiae*. For this, steps 2 and 4 of the algorithm proposed (refer to the Methods Section for further details) were coded in GAMS 23.2.0, while the normalization step 3 was implemented off-line using Microsoft Excel. Numerical experiments were performed on an Intel 1.2 GHz. The GMA model (Step 1) was retrieved from [8]. The reader is referred to this paper for further technical details. Bounds on metabolite concentrations and changes in enzyme activities were the same as those reported in [14].

Note that we assume that the GMA model is given. If this was not the case, a previous step would be necessary to construct such a model from dynamic profiles using parameter estimation methods. We should note also that the modeling software GAMS is a versatile tool that allows implementation of all the framework's steps, offering standard coding capabilities and interfacing with powerful optimization solvers.

### Obtention of the Pareto set

The MOO problem was solved using the epsilon constraint method, which was enhanced through a heuristic procedure based on generating solutions for all possible bi-criteria subproblems. We defined 10 epsilon parameters for each objective, which gave rise to 910 single iterations (note that the same number of objectives and epsilon intervals would lead to more than  $1 \cdot 10^{14}$  instances using the traditional epsilon constraint approach). The outer approximation-based algorithm [14,37] was then employed to solve these instances to global optimality. CPLEX 11.2.1 was used as MILP solver for the lower bounding master

problem, and CONOPT 3.14s for the slave NLPs. All the sub-problems of the algorithm were solved to global optimality within a tolerance of 0.2%, which is the same tolerance that we used in [14] for the analogous single objective problem. A set of Pareto optimal solutions was finally obtained through the above commented procedure. Figure 2 shows the 2D Pareto set for the maximization of  $V_{ethanol}$  vs minimization of hexose transporters (i.e.  $K_1$ ) changes. As observed, as we increase the value of  $K_1$  (recall that we are representing  $|\ln(K_1)|$ ), the ethanol synthesis rate increases. In the same Figure, we have also projected the points resulting from the other bi-criteria optimizations, that is, in Figure 2 we have included also the points obtained from the optimization of  $V_{ethanol}-K_2$ ,  $V_{ethanol}-K_3$ , ...,  $V_{ethanol}-K_8$ . As observed, while there is a clear tendency in the points coming from one bi-criteria optimization, the same is not true when we consider the remaining solutions generated by the other bi-criteria results. Hence, while we can "easily" analyze the trade-off between two single objectives, it is difficult to perform the same analysis when several criteria come into play.

The Pareto set was next normalized (see the Section "Normalization of the Pareto optimal solutions" in Methods) assuming a normal distribution for all objectives. We further assumed that the mean and standard deviation are the same as those of the samples (i.e., the solutions generated with the epsilon constraint method). Note that this brings the data to the [0,1] range. Figure 3 shows the box plot associated with the normalized Pareto solutions. As seen, objective  $K_2$  shows a very small variability (the 25<sup>th</sup> and 75<sup>th</sup> percentiles correspond to the same value, around 0.34, as the median). This implies in turn that it is easy to obtain a good (i.e. small) value for this objective. The same happens in the case of objectives  $K_5$ ,  $K_7$ ,  $X_1$ ,  $X_3$  and  $X_4$ , for which the median and 25<sup>th</sup> percentile are also rather close, indicating that the solutions are concentrated around their minimum values. On the contrary, most solutions are allocated at high (i.e., poor) values of objectives  $K_4$ ,  $K_8$  and  $V_{ethanol}$ , while very few are close to their minimum values.

## Selection of preferred subset of solutions

The Smart filter was applied next in order to remove indistinguishable solutions from the pool. The application of this algorithm has also the effect of providing a more uniform spread of points. Note that choosing larger values of tolerance  $\Delta t$  will allow discarding more solutions from the pool, but this may come at the expense of losing valuable solutions (i.e., promising enzymatic profiles). To illustrate this, we performed the calculations for two different values of  $\Delta t$ . In particular, selecting a  $\Delta t = 0.01$  allowed

to reduce the size of the Pareto set from 910 to 611 solutions, whereas only 321 solutions were retained for a  $\Delta t = 5.00$ . We found that using a  $\Delta t = 5.01$  resulted in an excessive loss of information in this case study, and hence, kept the results obtained with a  $\Delta t = 0.01$ .

We next resort to the second type of Pareto filter: the order of efficiency filter. We started by imposing a  $Q = 13$  (i.e.,  $Q = NO - 1$ ), and searched for nondominate solutions in any of the  $Q$ -elements subsets of objectives. This narrowed down the number of Pareto solutions from 611 to 214 alternatives. The procedure was repeated for decreasing values of  $Q$  until an empty set of solutions was identified, which occurred for a value of  $Q = 10$ . In particular, 14 solutions were found to be efficient of order 12, while only 1 solution was efficient of order 11.

Figure 4 shows the minimum and maximum objective values among those solutions retained for a given  $Q$ . This plot provides valuable insight on how much quality is lost as we decrease efficiency order. The closer the lower bound curve of a set of solutions is to the lower bound curve of the original set, the better is the quality of the set, as this implies that such set contains solutions with objective function values close to the best possible performance that can be attained in each criterion.

Particularly, the lower and upper limits of the 214 solutions efficient of order 13 are quite close to the bounds corresponding to the 611 solutions of the Pareto set obtained using the Smart filter, showing a small decrease (about 2%) in the ethanol synthesis rate with respect to the maximum possible value. There are 14 solutions efficient of order 12 with a curve rather close in most objectives to that of the 611 original solutions. In this set, however, the ethanol synthesis rate drops by an additional 69%, which is consistent with the trend observed in Figure 3. We should clarify that it is possible to artificially add in the final pool of solutions any other alternative for further consideration, with special interest on those with good performance in one criterion and poor in the others that are not efficient of order 12.

Remarkably, the only solution efficient of order 11 (which is not included in Figure 4) is not the closest to the utopia point, that is, it is not the one with the minimum Euclidean distance to the utopia point, which is a common criterion for selecting a single final candidate from a Pareto set.

Table 1 shows the values obtained for the 14 objectives in the solutions efficient of order 12. It can be seen that some of the solutions are very close to the ethanol production rate of the basal solution (i.e., solution with the  $K_r$  values fixed to one), which turns out to be  $30.11 \text{ mM min}^{-1}$  [8]. The best solution comprising only three changes in enzyme activity achieves a ethanol production rate of  $37.68 \text{ mM min}^{-1}$  and involves a 2.3 fold increase in  $E_3$  (which corresponds to a  $|\ln(K_3)|=0.84$ ), and about a 5 fold increase

in  $E_5$  and  $E_7$ . A ethanol production rate of  $42.88 \text{ mM min}^{-1}$  can be achieved by changing four enzymes. This leads to a 42% increase over the basal production rate. In this case,  $E_3$  must be modified by a factor of 3.5,  $E_4$  5 times,  $E_7$  2.1 times, and  $E_8$  1.7 times approximately. Further increases in ethanol production would require manipulating a larger set of enzymes. Single objective optimization focusing on maximizing the ethanol production would obtain better yields, but would entail higher (costly) enzyme changes and probably higher metabolic concentrations that would compromise the cell viability.

## Discussion

In this paper, we have introduced a systematic framework for the multiobjective deterministic global optimization of metabolic networks modeled through the GMA formalism. The proposed strategy integrates the epsilon constraint method, deterministic global optimization tools, and a set of Pareto filters that narrow down the final number of candidate solutions to be tested in the laboratory. The method presented does not rely on any visualization procedure, being therefore suitable for problems with a large number of objectives. The capabilities of the proposed approach were illustrated by means of a benchmark problem that addressed the optimization of the ethanol synthesis rate in *Saccharomyces cerevisiae*.

Biological objectives, such as the concentration of intermediate metabolites and the enzymatic changes were considered in addition to the ethanol synthesis rate. By selecting the auxiliary problems of the epsilon constraint method in a smart way, we could reduce the computational burden considerably. Furthermore, the Pareto filters allowed reducing the number of promising alternatives significantly from 910 to 14 (i.e., 98% reduction), illustrating the usefulness of the approach in the post optimal analysis of the candidate solutions. In different test problems, the outer approximation algorithm integrated in our systematic framework efficiently solved problems with up to 30 independent metabolites and 60 reactions in short CPU times (i.e., few minutes). Hence, we expect the method to scale up smoothly when tackling more complex models, even though we have yet to explore its limits. Note, however, that genome-wide scale problems are still beyond the capabilities of current deterministic global optimization methods. First, there is a lack of kinetic data to build realistic genome scale models. Second, assuming the existence of a detailed enough kinetic model, there is still the issue of solving it to global optimality in short CPU time. For these reasons, genome scale models are usually solve via FBA, despite the known limitations of this method. Nevertheless, we think that advances in deterministic global optimization theory and

software applications will pave the way for more efficient algorithms leading to significant CPU savings, which will make it possible to tackle complex genome scale kinetic models.

In summary, our approach allows for the global optimization of metabolic networks on different objectives simultaneously. The method presented reduces the computational burden associated with the generation of solutions, and facilitates the post-optimal analysis of these alternatives by systematically identifying the best ones (i.e., more balanced) for subsequent experiments in the laboratory. Hence, our method is particularly suited for problems of moderate size. Larger kinetic models could be tackled with stochastic methods, but even if they are the method of choice, it will be still possible to use the Pareto filters introduced in our work. However, we will not have any information on the quality of the solution found. Finally, for genome-wide scale models, FBA might be the method of choice, despite having some limitations already discussed in the literature.

## Methods

Our systematic framework comprises the following steps (see Figure 5):

1. Model building and parameter estimation (optional): construct a GMA model for the targeted metabolic network.
2. Global optimization of the GMA model on several biological criteria.
3. Normalization of the solutions obtained in step 2.
4. Application of Pareto filters to identify the preferred subset of alternatives.

The sections that follow describe in detail each of these steps.

### **Mathematical model: GMA representation**

The optimization of the metabolic network is posed in mathematical terms as a multiobjective NLP (i.e., moNLP) that embeds GMA equations. Note that there are different possible ways to obtain this GMA model. Particularly, we can follow a top-down approach, that is, find the parameters of a GMA model that make it consistent with dynamic data by solving a parameter estimation problem. On the contrary, we might be interested in following a bottom-up strategy and acquire the GMA model of interest from

the literature. In what follows, we describe briefly the GMA formalism before presenting the details of the moNLP.

We assume that the concentration  $X_i$  of every metabolite  $i$  present in a metabolic network varies with time  $t$  as a result of the action of  $p$  flows:

$$\frac{dX_i}{dt} = \sum_{r=1}^p \mu_{ir} v_r \quad i = 1, \dots, n \quad (1)$$

The stoichiometric coefficient,  $\mu_{ir}$ , appearing in Eq. 1 is an integer parameter accounting for the number of molecules of metabolite  $X_i$  that are involved in the process  $r$ . It is positive when the reaction  $r$  produces metabolite  $X_i$  and negative when  $r$  consumes  $X_i$ . Note that not all the  $p$  processes in the metabolic network are directly involved in the production of every single metabolite  $X_i$ , which implies that some parameters  $\mu_{ir}$  are zero ( $\mu_{ir} = 0$ ) for some particular combinations of  $i$  and  $r$ . The velocity at which process  $r$  occurs, is represented using the so-called power-law formalism [38–40] as in Eq. 2.

$$v_r = \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \quad r = 1, \dots, p \quad (2)$$

Here,  $\gamma_r$  is a parameter denoting the basal state activity of the enzyme governing process  $r$ , whereas  $f_{rj}$  is the kinetic order of metabolite  $X_j$  in process  $r$ . This representation accounts for the  $n$  internal dependent and  $m$  external (i.e., independent) metabolites. At this point, the concentration of the external metabolites will be considered fixed. Thus, the term  $X_j^{f_{rj}}$  behaves as a variable for  $i = 1, \dots, n$  and as a parameter for  $i = n + 1, \dots, n + m$ . By combining Eq. 2 and Eq. 1, we obtain a GMA model (Eq. 3).

$$\frac{dX_i}{dt} = \sum_{r=1}^p \left( \mu_{ir} \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \right) \quad i = 1, \dots, n \quad (3)$$

To model the effect of genetic manipulations performed on the strain, we introduce an auxiliary continuous variable,  $K_r$  that accounts for the fold-change over the basal state enzymatic level  $\gamma_r$  as follows:

$$v_r = K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \quad r = 1, \dots, p \quad (4)$$

Recall that, in Eq. 4, the product  $K_r \gamma_r$  denotes the actual enzyme activity. Hence, the values of  $K_r$  in the optimal solution will dictate the modification to be performed in the strain:  $K_r > 1$  indicates

overexpression of enzyme  $r$ ,  $K_r < 1$  denotes its downregulation, and a value of 1 means that enzyme  $r$  is not manipulated. Furthermore, bounds  $K_r^{LB}$  and  $K_r^{UB}$  are imposed on this variables as stated in Eq. 5.

$$K_r^{LB} \leq K_r \leq K_r^{UB} \quad r = 1, \dots, p \quad (5)$$

Similarly, metabolite concentrations are allowed to change within given bounds ( $X_i^{LB}$  and  $X_i^{UB}$ , respectively):

$$X_i^{LB} \leq X_i \leq X_i^{UB} \quad i = 1, \dots, n \quad (6)$$

Since we are interested in solving the steady state, the time dependence can be dropped from the formulation:

$$\frac{dX_i}{dt} = \sum_{r=1}^p \left( \mu_{ir} K_r \gamma_r \prod_{j=1}^{n+m} X_j^{f_{rj}} \right) = 0 \quad i = 1, \dots, n \quad (7)$$

For demonstrative purposes, we assume that the main objective is to maximize the synthesis rate of a desired product. This rate is calculated by summing up the velocities of those processes contributing to its synthesis, as illustrated in Eq. 8.

$$\min f_1 = - \sum_{r \in FP_i} \mu_{ir} v_r \quad i \in FP \quad (8)$$

Here,  $FP$  is the set of metabolites  $i$  that are regarded as final products and  $FP_i$  is the set of processes  $r$  contributing to the synthesis of metabolite  $i$  (i.e., those processes for which  $\mu_{ir} > 0$ ). Note that, for simplicity, we have posed the problem as a minimization one by reversing the sign of the objective function.

Two additional criteria are appended to the objective function. The first is the minimization of the metabolites concentrations, proposed as an optimality principle for metabolic networks [16]. Genetic manipulation of many genes at once may be costly and technically difficult. To take this into account, the model seeks to minimize the individual changes in enzyme activities. The resulting MOO problem that embeds the GMA equations can be expressed in compact form as follows:

$$\begin{aligned} (moGMA) \quad \min \quad & (f_1, \dots, f_{NO}) \\ \text{s.t.} \quad & \text{Eqs. 1, 4 - 6} \end{aligned} \quad (9)$$

Thus, model *moGMA* seeks the appropriate changes in enzyme activities (continuous variable  $K_r$ ) that maximize simultaneously the synthesis rate of the desired product and minimize the concentration of the metabolites and changes in enzyme activities. Objective  $f_1$  represents the synthesis rate targeted, while  $f_2$  to  $f_{NO}$  denote the metabolites concentrations  $X_i$  and individual changes in enzyme activities. To quantify deviations in enzyme activities from the basal state, we use the absolute value of the natural logarithm of the fold-change in enzyme activities. The enzyme activities calculated by the model can be later implemented in the real system by tuning the expressions of the corresponding genes.

The optimization problem takes the form of a nonconvex NLP, in which multiple local optima may exist. We employ global optimization techniques to ensure global optimality within a desired tolerance.

## Multiobjective global optimization of metabolic networks described by a GMA model

In general, the Pareto set of a GMA model may be nonconvex due to the nonlinear kinetic equations. Different MOO algorithms could be used to calculate this set (i.e., NBI [41], NNC [42]). We use herein the epsilon-constraint method because unlike other methods, such as the weighted sum one, it can identify points located in the nonconvex part of the Pareto set. Note that this property is also shared by the more complex NBI and NNC methods, which also offer the appealing property of providing a uniform spread of Pareto points. However, this limitation of the epsilon constraint is alleviated by coupling it with a Smart filter (refer to Section “Smart filter” in Methods). We should clarify, however, that our global optimization approach could work with other deterministic MOO algorithms, such as the NBI or NNC.

In the epsilon constraint technique, one objective is regarded as main objective, while the rest are transferred to auxiliary constraints that impose upper bounds  $\epsilon_b^e$  on their values:

$$\begin{aligned}
 (ecGMA) \quad & \min f_b \quad b = 1 \\
 & s.t. \quad f_{b'} \leq \epsilon_{b'}^e \quad e = 1, \dots, E + 1 \quad b' = 2, \dots, NO \\
 & \quad \quad \quad \text{Eqs.1, 4 - 6}
 \end{aligned} \tag{10}$$

The  $\epsilon_b^e$  values appearing in the auxiliary constraints are commonly obtained as follows:

1. Solve problem *moGMA* for each individual objective separately.

2. Store the best ( $\underline{f}_b$ ) and worst ( $\overline{f}_b$ ) values obtained in step 1 for each objective. These values are the limits within which the auxiliary epsilon parameters must fall (i.e.,  $\epsilon_b^e \in [\underline{f}_b, \overline{f}_b] \forall e$ ).
3. Split the epsilon interval into  $E$  subintervals to generate parameters  $\epsilon_b^e$  (i.e.,  $\epsilon_b^e = \underline{f}_b + (e-1) \cdot \frac{(\overline{f}_b - \underline{f}_b)}{(E)}$ ).

Note that step 1 provides the so-called anchor points, that is, the extreme solutions of the Pareto frontier.

In the traditional epsilon constraint approach, problem *ecGMA* is solved for all possible combinations of  $\epsilon_b^e$ , which leads to a total of  $(E+1)^{NO}$  instances. The complexity of this approach grows exponentially with the number of objectives. As an example, for 3 objectives and 4 sub-intervals, we have 125 iterations; for 4 objectives and the same number of sub-intervals, we have 625, and for 5 objectives and identical number of sub-intervals, we have 3125 iterations.

Here, we follow a heuristic approach for generating Pareto solutions that consists of solving a set of bi-criteria problems corresponding to all possible combinations of any two objectives. This strategy presents some advantageous features. First, the Pareto points generated in the two-dimensional space are also Pareto optimal in higher dimensional spaces [34], and hence in the original  $NO$ -dimensional space. Second, this approach requires solving  $\binom{NO}{2} \cdot (E+1)$  single-objective models, rather than  $(E+1)^{NO}$ , which dramatically reduces the computational effort. For instance, it would reduce the number of iterations required in the previous example from 125 to 15, from 625 to 30 and from 3125 to 60, respectively.

The epsilon constraint method transforms the MOO problem into a set of single-objective problems. This is very convenient, since it makes it possible to apply our global optimization methods devised for single-objective GMA models [14, 37] to multiobjective problems. In particular, in this work we use the outer-approximation-based algorithm we developed in [14], which was inspired by the works of Polisetty et al. [8] and Bergamini et al. [43].

Following this approach, the original problem (i.e., *ecGMA* in this case) is divided into two subproblems at two different hierarchical levels. A master problem consisting of a linear relaxation of *ecGMA* is solved in the upper level to predict a LB on the global optimum. A slave problem based on the original model is then solved locally in the lower level using the solution of the master problem as starting point in order to predict an UB. The solutions computed during the first iteration are used to tighten the relaxation of the master problem, which will produce better LBs in subsequent iterations. The algorithm proceeds in this manner until the optimality gap (OG, defined as the relative difference between the UB

and the LB) is reduced below a given tolerance.

The most important step of the outer approximation is the construction of the master MILP problem. This MILP is built by applying an exponential transformation that brings the model into a canonical form that can be relaxed in a straightforward manner using piecewise linear approximations and supporting hyper-planes. For the sake of brevity, technical details about this procedure are omitted herein. The interested reader is referred to the original works by Pozo et al. [14,37] for further details.

## Normalization of the Pareto optimal solutions

A normalization procedure is applied to the Pareto set of solutions in order to bring them to the same scale and units, so they become readily comparable. A plethora of alternative methodologies are available for this purpose. One of the main drawbacks of normalization methods is that they tend to concentrate the points in some regions of the feasible domain.

In a recent work, Cloquell et al. [44] presented a normalization methodology previously proposed in another work [45] that aims at overcoming this limitation. According to this strategy, the normalized value of a given solution  $s$  is calculated as follows:

$$fn_{s,b} = P_{DF(f_b)}(f_b \leq f_{s,b}) \quad (11)$$

Where  $fn_{s,b}$  is the normalized value associated with the non-normalized value  $f_{s,b}$ , and  $DF(f_b)$  is the probability distribution function of the objective variable  $f_b$ . The form of this distribution is assumed beforehand, with the normal distribution being the common choice.

## Pareto filters

The previous steps provide as output a set of normalized Pareto points. As mentioned previously, an infinite number of such points may exist for problems involving continuous variables. Testing all of them in the laboratory would be highly expensive and time consuming. Hence, a method is required for screening and ranking them, narrowing down their total number. We explore the application of two different Pareto filters. A Smart filter [46] is applied first to remove indistinguishable alternatives from the pool. A second filter based on the order of efficiency of the Pareto solutions [47] is then employed to identify solutions that are well-balanced, that is, they show "good" performance simultaneously in all of

the objectives.

### Smart filter

Two arbitrary solutions that are rather close in the objective space might be equally appealing for decision-makers, despite representing completely different experimental manipulations. If any of these is preferred over the other because of differences in any of the required changes, this differentiating feature should then be regarded as an additional objective [46]. A possible way to reduce the size of the Pareto set is to eliminate solutions which are within a given tolerance in the objectives space, that is, solutions which entail insignificant differences compared to others. Figure 6 illustrates the underlying idea behind this filter. As seen in Figure 6a, a region is defined around each normalized solution  $FN_s$ . Any other solution  $FN_{s'}$  falling inside this region is said to be indistinguishable from  $FN_s$ , and automatically removed from the pool. Consider for instance the example presented in Figure 6b where a small set of solutions is presented. We start by comparing solution  $FN_1$  with the rest, and then removing those contained inside the shaded region defined around the reference point. After comparing all the points, we pick the next candidate solution and repeat the procedure again. In this particular example, solution  $FN_2$  is found within the specified tolerance of  $FN_1$  and  $FN_5$  is within the region defined by  $FN_4$ .

To this end, we use the following algorithm, which is based on that presented by Mattson et al. [46]:

Let  $FN_s$  be one of the  $NS$  normalized solutions of the normalized Pareto set (i.e.,  $FN_s = fn_{s,1}, \dots, fn_{s,NO}$ ) obtained through steps 2 and 3 of the solution approach, and let  $SOS$  be the set containing all these solutions. The application of the filter comprises the following steps.

1. Define tolerance  $\Delta t$ , a set of rejected solutions  $SOR = \emptyset$ , a set of candidate solutions  $SOC = \emptyset$  and start iteration counters  $s = 0$  and  $ss = 0$ .
2. While  $s < NS$ ,
  - (a)  $s = s + 1$
  - (b) If  $\nexists FN_s | FN_s \in SOS$ , return to 2.a. Else:
  - (c) While  $ss < NS$ ,
    - i.  $ss = ss + 1$
    - ii. If  $\nexists FN_{ss} | FN_{ss} \in SOS$ , return to 2.c.i. Else:

iii. If  $s = ss$ , return to 2.c.i. Else, if  $f_{n_{s,b}} - f_{n_{ss,b}} \leq \Delta t \forall b$ , let  $SOR = SOR \cup FN_{ss}$  and  $SOS = SOS \setminus SOR$ .

(d) End while

(e) Restart iteration counter  $ss = 0$ .

3. End while

4. Make  $SOC = SOS$

We should clarify that this algorithm is a special case of the one proposed by Mattson et al. [46], in which the original  $\Delta t$  and  $\Delta r$  are assumed to be equal to  $\Delta t$ . Furthermore, note that the value of this control parameter is the same in all of the objectives, since the Pareto points are normalized prior to the application of the filter.

This filter is particularly useful when coupled with the epsilon constraint method, as it alleviates its tendency to concentrate points in given regions of the Pareto front, thus giving rise to a more uniform spread of points.

### Order of efficiency filter

The filter described above allows reducing the number of Pareto solutions. Further reductions can be attained by applying the concept of order of efficiency, as introduced by Das [47]. A solution is said to be efficient of order  $Q$  if it is not dominated by any other solution in any of the possible  $Q$ -elements subsets of objectives. In mathematical terms, a solution  $F_s$  is said to be efficient of order  $Q$ , if and only if,  $\nexists F_{s'} | F_{s'} \prec F_s$  for any subset of objectives of cardinality  $Q$ . In this definition, we consider that a solution  $F_s$  dominates  $F_{s'}$  (i.e.,  $F_s \prec F_{s'}$ ) if and only if,  $f_{s,b} \leq f_{s',b} \forall b$  with at least one  $b$  in which  $f_{s,b} < f_{s',b}$ .

Figure 7 provides an illustrative example of the concept of Pareto efficiency of order  $Q$ . Consider we have a MOO problem with 5 biotechnological criteria: final product yield, aggregated cost of changing the enzyme activities via gene expression, and concentration of 3 different metabolites ( $X_1 - X_3$ ). Assume that the values of 3 different solutions (blue, red and green) have already been normalized as described previously, so that the minimum value of each of the 5 objectives represents their individual optima. As seen, the three solutions plotted are Pareto optimal since none of them can improve any of the others simultaneously in all of the objectives. At this point, one can start eliminating solutions which are not efficient of order  $Q = 4$  by identifying sets of 4 objectives in which a given solution is dominated. For

instance, the blue solution is dominated by the green and red ones in the set  $\{yield, X_1, X_2, X_3\}$ . On the other hand, the red solution is not efficient of order 3, since it is in turn dominated by the green one in  $\{yield, X_1, X_2\}$ . Hence, the green solution is the only one that is efficient of order 3, while none of them is efficient of order 2 (i.e., the green solution is dominated by the red one in  $\{cost, X_3\}$ ).

According to the definitions previously exposed, if a solution is efficient of order  $Q$ , it is also efficient of order  $Q + L$  with  $L = 1, \dots, NO - Q$  (see [47] for proofs). Note that the concept of efficiency of order  $Q$  is stronger than the Pareto optimality condition [47], and can thus be used to discern between efficient alternatives. Furthermore, this concept avoids the use of any arbitrary “criterion of merit” or visualization technique, making it suitable for high-dimensionality problems [47].

We propose to apply this filter for searching efficient solutions of order  $NO - 1$ , and then repeat the process recursively for successively inferior orders of efficiency until either an empty set is found or the number of solutions retained is sufficiently small. As pointed out by Das [47], solutions with lower order of efficiency are expected to be well-balanced. This is because solutions behaving well in some objectives but poorly in others are expected to be dominated at least in the subsets including the latter criteria [47].

## Acknowledgments

The authors wish to acknowledge support of this research work from the Spanish Ministry of Education and Science (projects DPI2008-04099, PHB2008-0090-PC, BFU2008-00196 and CTQ2009-14420-C02), the Spanish Ministry of External Affairs (projects HS2007-0006 and A/023551/09), the Spanish Ministry of Science and Innovation (ENE2011-28269-C03-03), the Generalitat de Catalunya (FI programs) and the Spanish Ministry of Education, Culture and Sport (FPU programs).

## References

1. Kim IK, Roldao A, Siewers V, Nielsen J (2012) A systems-level approach for metabolic engineering of yeast cell factories. *FEMS yeast research* 12(2): 228-248.
2. Parachin NS, Bergdahl B, van Niel EW, Gorwa-Grauslund MF (2011) Kinetic modelling reveals current limitations in the production of ethanol from xylose by recombinant *saccharomyces cerevisiae*. *Metabolic engineering* 13: 508-517.

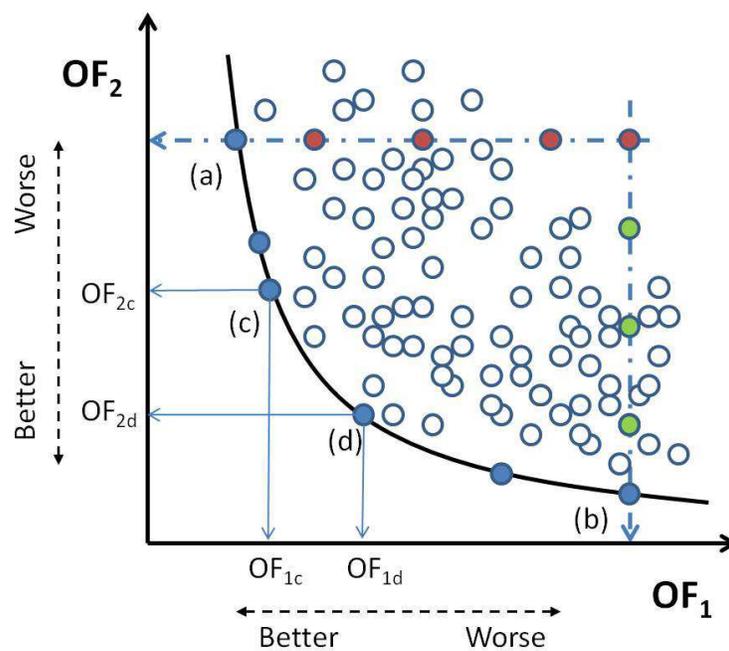
3. Li RD, Li YY, Lu LY, Ren C, Li YX, et al. (2011) An improved kinetic model for the acetone-butanol-ethanol pathway of *Clostridium acetobutylicum* and model-based perturbation analysis. *BMC systems biology* 5 Suppl 1: S12.
4. Voit EO (1992) Optimization in integrated biochemical systems. *Biotechnol Bioeng* 40: 572-82.
5. Hatzimanikatis V, Floudas CA, Bailey JE (1996) Optimization of regulatory architectures in metabolic reaction networks. *Biotechnology and bioengineering* 52: 485-500.
6. Alvarez-Vasquez F, González-Alcón C, Torres NV (2000) Metabolism of citric acid production by *Aspergillus niger*: model definition, steady-state analysis and constrained optimization of citric acid production rate. *Biotechnol Bioeng* 70: 82-108.
7. Marín-Sanguino A, Torres NV (2003) Optimization of biochemical systems by linear programming and general mass action model representations. *Math Biosci* 184: 187-200.
8. Polisetty P, Gatzke E, Voit E (2008) Yield optimization of regulated metabolic systems using deterministic branch-and-reduce methods. *Biotechnology and Bioengineering* 99(5): 1154-1169.
9. Marín-Sanguino A, Torres N, Mendoza E, Oesterhelt D (2009) Metabolic engineering with power-law and linear-logarithmic systems. *Mathematical Biosciences* 218(1): 50-58.
10. Grossmann I, Biegler L (2004) Retrospective on optimization. *Computers and Chemical Engineering* 28: 1169-1192.
11. Voit E (2003) Design principles and operating principles: the yin and yang of optimal functioning. *Math Bioscience* 182: 81-92.
12. Guillén-Gosálbez G, Sorribas A (2009) Identifying quantitative operation principles in metabolic pathways: a systematic method for searching feasible enzyme activity patterns leading to cellular adaptive responses. *BMC Bioinformatics* 10(386).
13. Pozo C, Marín-Sanguino A, Alves R, Guillén-Gosálbez G, Jiménez L, et al. (2011) Steady-state global optimization of metabolic non-linear dynamic models through recasting into power-law canonical models. *BMC systems biology* 5: 137-0509-5-137.

14. Pozo C, Guillén-Gosálbez G, Sorribas A, Jiménez L (2010) Outer approximation-based algorithm for biotechnology studies in systems biology. *Computers and Chemical Engineering* 34(10): 1719-1730.
15. Sorribas A, Pozo C, Vilaprinyo E, Guillén-Gosálbez G, Jiménez L, et al. (2010) Optimization and evolution in metabolic pathways: Global optimization techniques in generalized mass action models. *Journal of Biotechnology* 149(3): 141-153.
16. Schuster S, Schuster R, Heinrich R (1991) Minimization of intermediate concentrations as a suggested optimality principle for biochemical networks. ii. time hierarchy, enzymatic rate laws, and erythrocyte metabolism. *J Math Biol* 29(5): 443-455.
17. Heinrich R, Schuster S (2005) *The regulation of cellular systems*. Chapman and Hall, New York.
18. Sendín J, Exler O, Banga J (2010) Multi-objective mixed integer strategy for the optimization of biological networks. *IET Systems Biology* 4(3): 236-248.
19. Vera J, de Atauri P, Cascante M, Torres NV (2003) Multicriteria optimization of biochemical systems by linear programming: application to production of ethanol by *saccharomyces cerevisiae*. *Biotechnol Bioeng* 83: 335-43.
20. Liu PK, Wang FS (2008) Inference of biochemical network models in s-system using multiobjective optimization approach. *Bioinformatics* 24: 1085-92.
21. Wu WH, Wang FS, Chang MS (2011) Multi-objective optimization of enzyme manipulations in metabolic networks considering resilience effects. *BMC systems biology* 5: 145-0509-5-145.
22. Handl J, Kell D, Knowles J (2007) Multiobjective optimization in bioinformatics and computational biology. *IEEE/ACM Trans Comput Biol Bioinf* 4(2): 279-291.
23. Branke J, Kaussler T, Schmeck T (2001) Guidance in evolutionary multi-objective optimization. *Advances in Engineering Software* 32: 499-507.
24. Deb K (2003) Multi-objective evolutionary algorithms: Introducing bias among pareto-optimal solutions. In: Ghosh A, Tsutsui S, editors, *Advances in Evolutionary Computing: Theory and Applications*. London, England: Springer-Verlag, pp. 263-292.

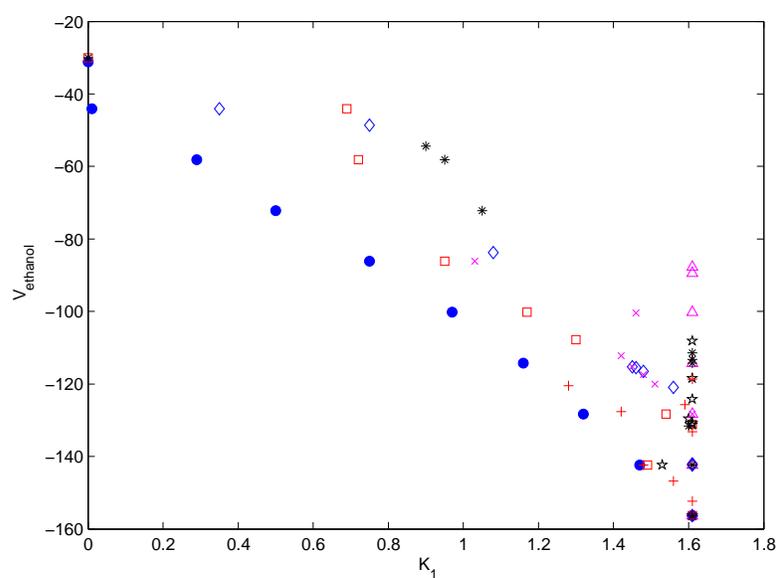
25. Branke J, Deb K (2004) Integrating user preferences into evolutionary multi-objective optimization. Knowledge Incorporation in Evolutionary Computation : 461-477.
26. Farina M, Amato P (2004) A fuzzy definition of "optimality" for many-criteria optimization problems. IEEE Transactions on Systems, Man, and Cybernetics Part A: Systems and Humans 34(3): 315-326.
27. Branke J, Deb K, Dierolf H, Osswald M (2004) Finding knees in multi-objective optimization. Parallel Problem Solving from Nature PPSN-VIII: 722-731.
28. Deb K, Gupta H (2005) Searching for robust pareto-optimal solutions in multi-objective optimization. In: Third Evolutionary Multi-Criteria Optimization (EMO-05) Conference. pp. 150-164.
29. Mattson C, Messac A (2003) Concept selection using s-pareto frontiers. AIAA Journal 41(6): 1190-1198.
30. Messac A, Ismail-Yahaya A, Mattson C (2003) The normalized normal constraint method for generating the pareto frontier. Structural and Multidisciplinary Optimization 25(2): 86-98.
31. Montusiewicz J, Osyczka A (1990) A decomposition strategy for multicriteria optimization with application to machine tool design. Engineering Costs and Production Economics 20: 191-202.
32. Abraham S, Rau B, Schriber R (2000) Fast design space exploration through validity and quality filtering of subsystem designs. Technical report, HPL-2000-98, Hewlett Packard.
33. Cheng FY, Li D (1998) Genetic algorithm development for multiobjective optimization of structures. AIAA Journal 36(6): 1105-1112.
34. Ehrgott M (1998) Multicriteria Optimization. Berlin: Springer.
35. Sendín J, Vera J, Torres N, Banga J (2006) Model-based optimization of biochemical systems using multiple objectives: a comparison of several solution strategies. Math Comput Model Dyn Syst 12(5): 469-487.
36. Deb K, Saxena D (2005) On finding pareto-optimal solutions through dimensionality reduction for certain large-dimensional multi-objective optimization problems. Tech rep, Kanpur Genetic Algorithms Laboratory (KanGAL), Indian Institute of Technology Kanpur .

37. Pozo C, Guillén-Gosálbez G, Sorribas A, Jiménez L (2011) A spatial branch-and-bound framework for the global optimization of kinetic models of metabolic networks. *Industrial and Engineering Chemistry Research* 50(9): 5225-5238.
38. Savageau M (1969) Biochemical systems analysis. i. some mathematical properties of the rate law for the component enzymatic reactions. *J Theor Biol* 25: 365-369.
39. Savageau M (1969) Biochemical systems analysis. ii. the steady-state solutions for an n-pool system using a power-law approximation. *J Theor Biol* 25: 370-379.
40. Voit E (2000) *Computational Analysis of Biochemical Systems. A Practical Guide for Biochemists and Molecular Biologists*. Cambridge University Press: Cambridge, U.K.
41. Das I, Dennis J (1998) Normal-boundary intersection: A new method for generating the pareto surface in nonlinear multicriteria optimization problems. *SIAM Journal on Optimization* 8(3): 631-657.
42. Messac A, Mattson C (2004) Normal constraint method with guarantee of even representation of complete pareto frontier. *AIAA Journal* 42(10): 2101-2111.
43. Bergamini M, Aguirre P, Grossmann I (2005) Logic-based outer approximation for globally optimal synthesis of process networks. *Computers and Chemical Engineering* 29: 1914-1933.
44. Cloquell V, Santamarina M, Hospitaler A (2001) Nuevo procedimiento para la normalización de valores numéricos en la toma de decisiones. In: XVII Congreso Nacional de Ingeniería de Proyectos - Murcia 2001.
45. Cloquell V (1999) Contribución al desarrollo de un modelo generalizado y sistemático de localización de actividades económicas. Ph.D. thesis, Universidad Politécnica de Valencia.
46. Mattson C, Mullur A, Messac A (2004) Smart pareto filter: obtaining a minimal representation of multiobjective design space. *Engineering Optimization* 36(6): 721-740.
47. Das I (1999) A preference ordering among various pareto optimal alternatives. *Structural Optimization* 18: 30-35.

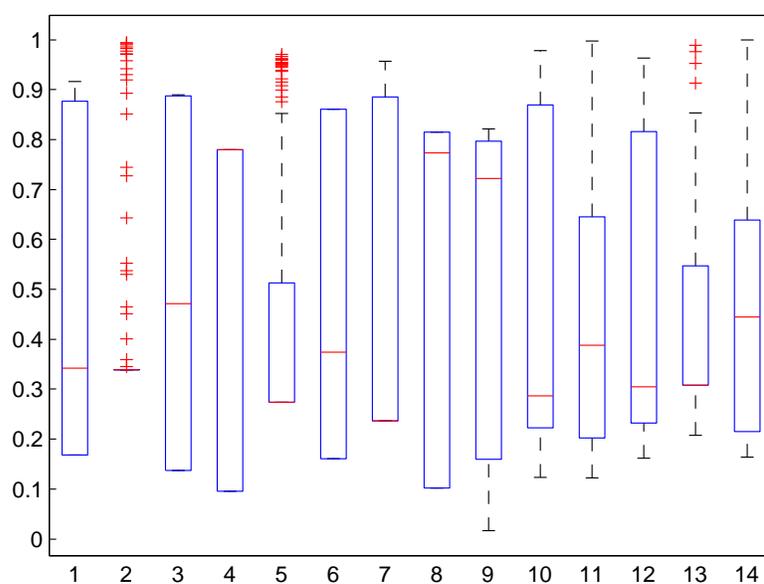
## Figure Legends



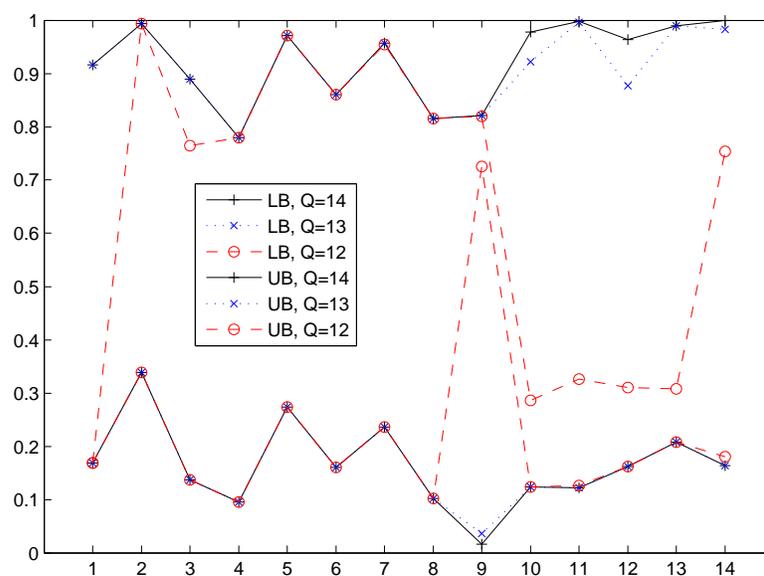
**Figure 1. Generic Pareto front.** Full blue points indicate members of the Pareto set. Point (a) is the optimum for objective function  $OF_1$  for a given value of  $OF_2$  (red points). Point (b) minimizes  $OF_2$  for another value of  $OF_1$  (compared to green points). For a member of the Pareto set, say (c), any attempt to improve a goal involves worsening the other, point (d) for comparison. Empty blue points are other possible solutions that are worse than those in the Pareto set.



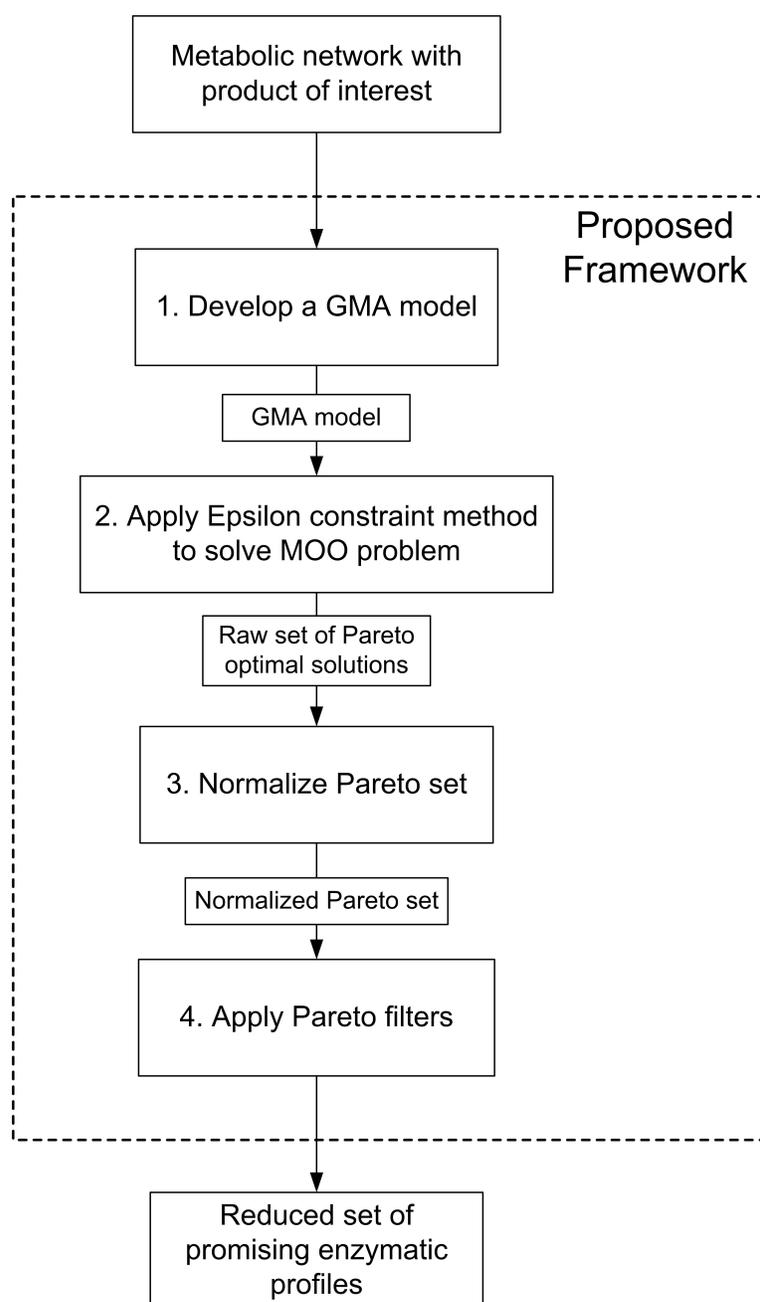
**Figure 2. Pareto curve (blue circles) of the bi-criteria problem considering  $V_{ethanol}$  and  $K_1$  (Hexose transporters).** The other points represent projections of the same variables obtained during other bi-criteria optimization problems:  $V_{ethanol}-K_2$  (red squares),  $V_{ethanol}-K_3$  (magenta triangles),  $V_{ethanol}-K_4$  (black stars),  $V_{ethanol}-K_5$  (blue diamonds),  $V_{ethanol}-K_6$  (red plus signs),  $V_{ethanol}-K_7$  (magenta cross signs) and  $V_{ethanol}-K_8$  (black asterisks). Fold-Change factors correspond to:  $K_1$ : Hexose transporters,  $K_2$ : Glucokinase/Hexokinase,  $K_3$ : Phosphofruktokinase,  $K_4$ : Trehalose 6-phosphate syntase complex (+Glycogen production),  $K_5$ : Glyceraldehyde-3-phosphate dehydrogenase,  $K_6$ : GOL (Glycerol production),  $K_7$ : Pyruvate kynase,  $K_8$ : ATPase.



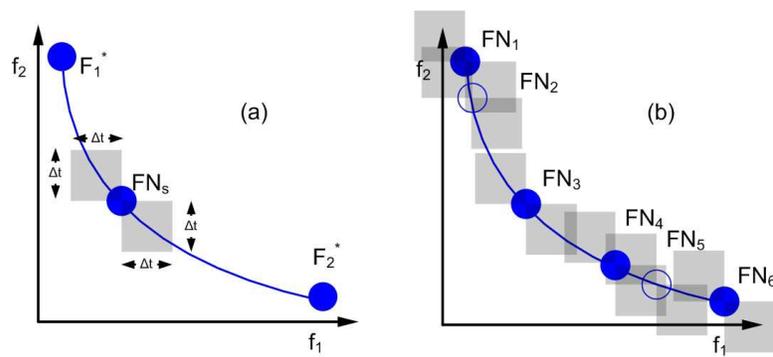
**Figure 3. Box plot for the normalized Pareto set.** In the bottom axis the fourteen objectives are represented. Objectives 1-8 correspond to  $K_1$ - $K_8$  (see legend in Figure 2), objective 9 is indeed  $V_{ethanol}$  whereas the remaining 5 objectives represent  $X_1$ - $X_5$ .  $X_1$ : Internal glucose,  $X_2$ : Glucose-6-phosphate,  $X_3$ : Fructose-1,6-diphosphate,  $X_4$ : Phosphoenolpyruvate,  $X_5$ : Adenosine triphosphate.



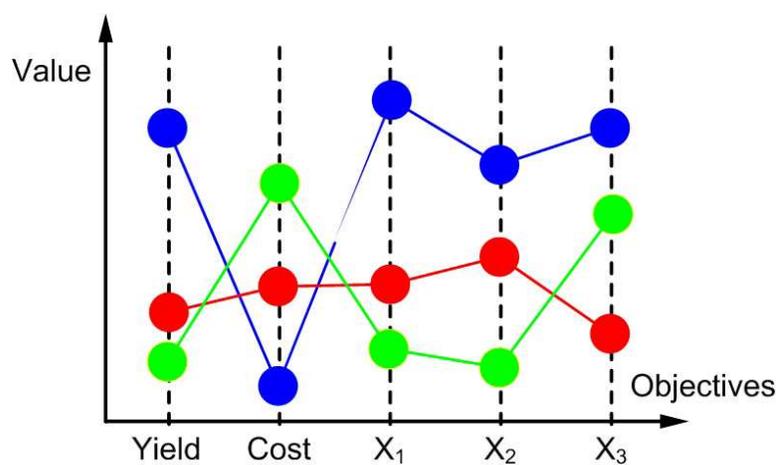
**Figure 4.** Lower and upper bounds for objectives among the values attained by the set of **Pareto solutions of order  $Q$** . In particular, 611 solutions are efficient of order 14 (i.e., these are indeed the solutions obtained after applying the Smart filter); 214 solutions are efficient of order 13; and 14 solutions are efficient of order 12. Objectives are ordered as in Figure 3. See legends in Figure 2 and 3.



**Figure 5. Proposed algorithm for the multiobjective global optimization of metabolic networks.** This method allows not only to generate a Pareto set, but also to systematically select the most promising subset of enzymatic profiles embedded therein.



**Figure 6. Illustration of the smart Pareto filter.** a) Indistinguishability region. b) Algorithm performance example.



**Figure 7. Illustrative example for the Pareto order of efficiency concept.** Blue solution is efficient of order 5, whereas red solution is efficient of order 4 and green solution is efficient of order 3.

## Tables

**Table 1.** 14 solutions efficient of order 12 in decreasing order of  $V_{ethanol}$ .

$K_1$	$K_2$	$K_3$	$K_4$	$K_5$	$K_6$	$K_7$	$K_8$	$V_{ethanol}$	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
0.00	0.00	0.86	1.61	1.16	0.00	1.16	1.61	43.27	0.06	0.26	3.27	<0.01	0.34
0.00	0.00	1.26	1.61	0.00	0.00	0.75	0.56	42.88	0.05	0.26	16.93	<0.01	0.94
0.00	0.00	1.14	1.61	1.52	0.00	0.15	0.82	41.95	0.05	0.26	1.49	0.01	0.70
0.00	0.00	0.97	1.61	0.00	1.31	0.39	1.07	38.36	0.05	0.26	16.59	0.01	0.44
0.00	0.00	0.84	0.00	1.61	0.00	1.59	0.00	37.68	0.04	0.47	0.91	<0.01	1.48
0.00	0.00	0.59	1.61	0.00	0.00	0.81	0.20	35.83	0.04	0.55	12.15	<0.01	1.14
0.00	1.61	0.56	1.61	0.00	1.57	0.25	1.58	34.97	0.01	0.29	16.53	0.01	0.22
0.00	1.00	1.17	1.61	1.61	1.61	0.05	0.28	34.43	0.01	0.26	0.91	0.01	0.74
0.00	1.18	0.00	0.00	1.48	0.00	0.00	1.22	33.53	0.01	0.64	1.25	0.01	0.37
0.00	0.00	0.00	0.00	0.00	1.30	0.53	1.35	32.20	0.04	0.58	13.66	<0.01	0.29
0.00	0.00	0.00	0.00	0.00	1.29	0.55	1.30	32.17	0.04	0.59	13.50	<0.01	0.31
0.00	0.00	0.57	0.00	1.61	1.61	0.00	0.86	31.46	0.05	0.36	0.91	0.01	0.37
0.00	1.61	0.45	1.61	0.00	1.29	0.16	0.02	30.54	<0.01	0.60	9.69	0.01	0.98
0.00	0.00	0.44	1.61	0.00	1.61	0.44	0.00	30.24	0.04	0.61	9.54	<0.01	0.98

Recall that columns labeled as  $K_r$  represent indeed  $|\ln(K_r)|$ . Enzyme 1: Hexose transporters, enzyme 2: Glucokinase/Hexokinase, enzyme 3: Phosphofructokinase, enzyme 4: Trehalose 6-phosphate syntase complex (+Glycogen production), enzyme 5: Glyceraldehyde-3-phosphate dehydrogenase, enzyme 6: GOL (Glycerol production), enzyme 7: Pyruvate kynase, enzyme 8: ATPase, metabolite 1: Internal glucose, metabolite 2: Glucose-6-phosphate, metabolite 3: Fructose-1,6-diphosphate, metabolite 4: Phosphoenolpyruvate, metabolite 5: Adenosine triphosphate.



Contents lists available at [SciVerse ScienceDirect](#)

# Chemical Engineering Science

journal homepage: [www.elsevier.com/locate/ces](http://www.elsevier.com/locate/ces)



## On the use of Principal Component Analysis for reducing the number of environmental objectives in multi-objective optimization: Application to the design of chemical supply chains

C. Pozo<sup>a</sup>, R. Ruíz-Femenia<sup>b</sup>, J. Caballero<sup>b</sup>, G. Guillén-Gosálbez<sup>a,\*</sup>, L. Jiménez<sup>a</sup>

<sup>a</sup> *Departament d'Enginyeria Química (EQ), Escola Tècnica Superior d'Enginyeria Química (ETSEQ), Universitat Rovira i Virgili (URV), Campus Sescelades, Avinguda Països Catalans, 26, 43007 Tarragona, Spain*

<sup>b</sup> *Department of Chemical Engineering, University of Alicante, Apdo. 99, 03080 Alicante, Spain*

### ARTICLE INFO

#### Article history:

Received 19 July 2011

Received in revised form

14 September 2011

Accepted 6 October 2011

Available online 18 October 2011

#### Keywords:

Multi-objective optimization

Principal component analysis

Life cycle assessment

Supply chain management

Dimensionality reduction

Mixed-integer linear programming

### ABSTRACT

Multi-objective optimization (MOO) has recently attracted an increasing interest in environmental engineering. One major limitation of the existing solution methods for MOO is that their computational burden tends to grow rapidly in size with the number of environmental objectives. In this paper, we study the use of Principal Component Analysis (PCA) to identify redundant environmental metrics in MOO that can be omitted without disturbing the main features of the problem, thereby reducing the associated complexity. We show that, besides its numerical usefulness, the use of PCA coupled with MOO provides valuable insights on the relationships between environmental indicators of concern for decision-makers. The capabilities of the proposed approach are illustrated through its application to the design of environmentally conscious chemical supply chains (SCs).

© 2011 Elsevier Ltd. All rights reserved.

### 1. Introduction

In the recent past, multi-objective optimization (MOO) has gained wider interest in environmental applications. This technique treats the environmental concerns as individual objectives rather than as constraints imposed on the system. This allows identifying alternatives that achieve significant environmental savings at a marginal increase in cost.

The definition of a suitable environmental metric to support objective environmental assessments is an open issue in the literature. The scientific community has not yet reached an agreement on the use of a universal environmental metric, which has motivated the development of a plethora of environmental assessment methodologies. For instance, Il Lim et al. (1999) proposed a global pollution index function for the design of chemical processes. Dantus and High (2000) proposed to quantify the environmental impact using hazardous values of different compounds generated in a process. Alternatively, Zhou et al. (2000) suggested to minimize the resources consumption and waste generation in supply chain management applications.

\* Corresponding author. Tel.: +34 977 558618.

E-mail address: [gonzalo.guillen@urv.cat](mailto:gonzalo.guillen@urv.cat) (G. Guillén-Gosálbez).

Among the environmental indicators that are available, those based on Life Cycle Assessment (LCA) principles (Guinée et al., 1992) are nowadays becoming the prevalent approach. The combined use of LCA and MOO was first proposed by Stefanis et al. (1995), whereas the theoretical foundations of the framework were in fact established by Azapagic and Clift (1999). Since then, several authors have applied this general approach to engineering problems of different nature and scale. For instance, the design and planning of chemical supply chains was addressed in the works by Hugo and Pistikopoulos (2003), Mele et al. (2005), Puigjaner and Guillén-Gosálbez (2008), Guillén-Gosálbez and Grossmann (2009, 2010) and Bojarski et al. (2009). At the single-site level, Stefanis et al. (1995); Chakraborty and Linninger (2002) and Guillén-Gosálbez et al. (2008) incorporated LCA principles in the design of chemical plants, whereas Stefanis et al. (1997) and Cavin et al. (2004) addressed the minimization of the LCA impact of batch facilities. Chang and Hwang (1996) and Papandreou and Shang (2008) introduced some models for the optimal design of utility systems that incorporated LCA-based metrics. Hugo et al. (2005) and Guillén-Gosálbez et al. (2010) addressed the strategic planning of hydrogen supply chains, whereas Pistikopoulos and Stefanis (1998) applied LCA principles in the selection of solvents in mass separating agent driven technologies. More recently, Gebreslassie et al. (2009, 2010) introduced MOO models that incorporated LCA metrics for the design of absorption cooling cycles.

The computational burden of the standard MOO solution methods that are employed in these approaches tends to grow rapidly with the number of objectives (Deb and Saxena, 2005). Further, the visualization and interpretation of the solution space of these problems become harder as one increases the number of objectives, which hampers the decision-making process. These limitations are critical in environmental applications, and particularly in those based on the combined use of MOO and LCA, in which the simultaneous analysis of a wide range of environmental metrics is sought. Previous attempts to ameliorate this difficulty focused on aggregation schemes based on the use of indicators calculated by attaching weights to the single environmental metrics considered in the analysis (see Guillén-Gosálbez, 2011).

This approach has two major drawbacks. First, the weights used, which are typically defined by a panel of experts, may not represent the decision-makers' interests. Second, aggregation methods modify the Pareto structure of the problem, in a manner such that some parts of the search space might be left out of the analysis. An alternative approach that surmounts these difficulties is to use dimensionality reduction methods that allow omitting redundant metrics from the problem in order to keep it in a manageable size. In this paper, we investigate the use of Principal Component Analysis (PCA) to reduce the dimensionality of multi-objective optimization problems arising in environmental applications. Our final goal is to reduce the problem complexity while still preserving its Pareto structure to the extent possible. An additional objective of our study is to gain valuable insights on the Pareto structure of a given problem in order to facilitate the decision-making process. The capabilities of the combined use of PCA and MOO are highlighted in the discussion of two case studies that address the design of environmentally conscious chemical SCs.

## 2. Materials and methods

The next sections review some general concepts about PCA and dimensionality reduction in MOO with emphasis on their application in environmental engineering. In what follows, we will consider MOO problems of the following form:

$$\begin{aligned} MO(x,y) = \min & \quad (f_1(x,y), \dots, f_m(x,y)) \\ \text{s.t.} & \quad h(x,y) = 0 \\ & \quad g(x,y) \leq 0 \\ & \quad x \in \mathfrak{R}, y \in \{0, 1\} \end{aligned} \quad (1)$$

where  $m$  objective functions are to be minimized (one economic metric and  $m-1$  environmental indicators). In supply chain management (SCM) problems, like the ones addressed in this paper, variables  $x$  denote mass flow rates, capacities of SC entities and economic and environmental performance indicators, whereas binary variables  $y$  are used to model the execution of capacity expansions and establishment of transportation links. Likewise, equality constraints represent mass balances and economic and environmental calculations, whereas inequality constraints express capacity limitations (see the Appendix for further details).

Problem  $MO$  can be solved by any MOO solution method (see Ehrgott, 1998). As discussed by Deb and Saxena (2005), the complexity of these techniques grows rapidly with the number of objectives. Handling a large number of objectives (i.e., more than three) in MOO causes additional difficulties related to the visualization and analysis of the Pareto set (Deb and Saxena, 2005).

In this paper, we explore the use of PCA (Johnson and Wichner, 1998) to ameliorate these limitations. As will be shown later in the article, PCA can be effectively employed to reduce the

dimension of MOO problems arising in environmental engineering, providing valuable insights on their structure.

### 2.1. Dimensionality reduction

Consider the set of objectives  $F := (f_1(x,y), \dots, f_m(x,y))$  in model  $MO$ . The goal of dimensionality reduction methods is to determine a subset  $F_0$  of  $F$  ( $F_0 \subseteq F$ ) such that the Pareto structure of the model is preserved to the extent possible when the problem is solved in this reduced domain. Zitzler et al. (2003) were the first to introduce a general notion of conflict between objective sets as a theoretical foundation for objective reduction. They also introduced the concept of delta error to quantify the change in the dominance structure resulting from removing objectives originally contained in  $F$ . An exact and an approximated algorithm were proposed by these authors for minimizing such an error. More recently, Guillén-Gosálbez (2011) developed an MILP-based method for dimensionality reduction that was applied to two environmental engineering problems. This approach employs branch and cut techniques to identify those environmental objectives that minimize the delta error of the approximation obtained after omitting them. Despite being rigorous, all of these algorithms based on the delta error concept have the limitation of being very sensitive to the number of solutions and objectives (i.e., their computation burden grows rapidly with the number of solutions and objectives).

In this paper, we will focus on a special type of methods for dimensionality reduction based on PCA. This technique presents the advantage of being faster than those based on the delta error. We next review some generalities on PCA before describing its use in the context of our problem.

### 2.2. Principal Component Analysis

Principal Component Analysis (PCA) is a multivariate technique that allows to transform a series of correlated variables  $c_1, c_2, \dots, c_m$  into a set of uncorrelated variables  $u_1, u_2, \dots, u_m$  known as principal components (PCs). PCs consist of convex combinations of the original variables. As a result, each of their components (sometimes referred to as loadings) represents the contribution of an original variable towards that PC. The PCs of a data set  $x$  are obtained by solving an eigenvalue–eigenvector problem. The eigenvectors of the covariance matrix  $\Sigma$  of  $x$  correspond to the PCs themselves. The correlation matrix  $R$  is preferred instead of  $\Sigma$  when the variables are expressed in different units. In this case

$$RP = LP \quad (2)$$

where  $P$  is a matrix containing the PCs arranged in descending order of their magnitudes (i.e., the  $j$ th column of  $P$  corresponds to the PC with the  $j$ th largest eigenvalue  $\lambda_j$ ) and  $L$  is a diagonal matrix storing the associated eigenvalues, also listed in the same way. These eigenvalues, which are obtained from a positive semi-definite matrix and are therefore always positive or null (Abdi and Williams, 2010), provide the amount of variance explained by the associated PC. Further, because of the symmetry of matrix  $R$ , PCs are pairwise orthogonal when their eigenvalues are different (for proofs see Strang, 2009). Henceforth, the first PC explains the largest portion of the problem's variance, followed by the second PC, and so on. All the PCs are constrained to be orthogonal between them. Hence, the percentage of total variance explained by the first  $k$  PCs is defined as follows:

$$CVAR_k(\%) = 100 \frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^m \lambda_j} \quad (3)$$

In practice, most of the problem's variance is explained by only a few PCs. The percentage of variance explained by the PCs can be

used as a criterion to select a subset of PCs for further analysis (Jolliffe, 2002).

In the context of environmental engineering, PCA can be used to analyze relationships between environmental criteria. Gutiérrez et al. (2010) explored the combined use of LCA and PCA to uncover and visualize the structure of large multidimensional data sets in the context of waste water treatment plants. After applying PCA, they observed strong correlations between some LCA metrics. They applied the same approach to mussel cultivation and waste electrical and electronic equipment (Gutiérrez et al., 2010).

PCA can be effectively employed for selecting a few variables among a wider set of correlated variables. In the context of MOO, this property is certainly useful as it allows identifying redundant objectives that can be removed from the model, thereby alleviating its complexity.

To the best of our knowledge, Deb and Saxena (2005) were the first to use PCA in MOO. They proposed some heuristic rules based on PCA that allow identifying redundant objectives from a set of feasible points of the MOO. These rules are based on the analysis of the components of the eigenvectors of the correlation matrix. The authors coupled this technique with evolutionary algorithms (i.e., elitist non-dominated sorting genetic algorithm, NSGA-II) in order to improve their performance.

In contrast to the work of Deb and Saxena (2005), we propose herein to integrate PCA with a rigorous MOO solution method (i.e., epsilon constraint coupled with branch and cut) that guarantees the global optimality of the solutions found. As will be discussed later in the paper, this approach has the advantage of providing solutions that are globally optimal within a desired tolerance in shorter CPU times, as opposed to the approach by Deb and Saxena (2005) that is unable to guarantee the optimality of the solutions found. An additional novelty of this work is that, to the best of our knowledge, it is the first contribution that makes use of PCA for reducing the number of objectives in an environmental engineering MOO problem.

### 2.3. Combined used of PCA and MOO

We next present a PCA based epsilon constraint procedure to solve MOO problems arising in environmental engineering. The method is applied in an iterative manner as follows. First, some Pareto solutions of the problem with all the original objectives are generated using the epsilon constraint method. PCA is then applied to these high-dimensional points in order to identify redundant environmental metrics that can be omitted. After removing redundant environmental objectives, the epsilon constraint method is applied again to the reduced space model, generating new solutions that will be analyzed using PCA. The algorithm proceeds in this manner until no further reductions in the dimension of the MOO problem are possible. The detailed steps of the algorithm are as follows:

#### 1. Initialization.

- (a) Set a threshold cut  $TC$  above which no more PCs will be considered, and a number of iterations  $it$  for the epsilon constraint method.
- (b) Optimize each individual objective  $f_i \in F$  separately. Store the best and worst values ( $\underline{f}_i$  and  $\bar{f}_i$ , respectively) of each objective function obtained in these optimizations.

#### 2. Apply the epsilon constraint method to model $MO$ .

- (a) Choose an objective  $f_i$  as main objective and transfer the remaining ones (i.e.,  $f_r \neq f_i$ ) to auxiliary constraints, giving rise to problem  $MOEC$ :

$$MOEC(x,y) = \min f_i(x,y)$$

$$\begin{aligned} \text{s.t. } & f_r(x,y) \leq \varepsilon_r^n \quad \forall f_r \in F \setminus f_i \quad n = 1, \dots, N \\ & h(x,y) = 0 \\ & g(x,y) \leq 0 \\ & x \in \mathfrak{R}, y \in \{0, 1\} \end{aligned} \quad (4)$$

- (b) Calculate the epsilon parameters ( $\varepsilon_r^n$ ) for each objective  $f_r \neq f_i$  by splitting the interval  $[\underline{f}_r, \bar{f}_r]$  into  $N-1$  sub-intervals.
  - (c) Solve  $MOEC$  for each set of epsilon parameters (a total of  $N^{|F|-1}$  problems must be solved).
  - (d) Generate matrix  $M$  (with dimension  $N^{|F|-1} \times |F|$ ) by storing in each row the values of the  $|F|$  objectives associated with each solution  $n$  of problem  $MOEC$ .
3. Compute the PCs of matrix  $M$ .
- (a) Filter matrix  $M$  by eliminating repeated, dominated or infeasible solutions.
  - (b) Standardize the filtered matrix  $M$  so as to make its centroid equal to 0. For this, we subtract the mean of each column  $\mu_i$  from each data point in the matrix.
  - (c) If the magnitudes of the  $|F|$  objectives are different, divide each data point in the matrix by the standard deviation of the corresponding column  $\sigma_i$  and compute the correlation matrix  $R$  of the standardized matrix  $M$ . Otherwise, compute the covariance matrix  $\Sigma$  of  $M$ .
  - (d) Determine the eigenvectors and eigenvalues of  $R$  (or  $\Sigma$ ). Order them in decreasing order of their eigenvalues and assign PCs recursively. That is, let the first eigenvector be the first PC; the second eigenvector be the second PC and so on so forth. Let the variance explained by each PC be equal to the associated eigenvalue ( $VAR_j = \lambda_j$ ).
4. Apply the objective reduction heuristic proposed by Deb and Saxena (2005).
- (a) Define an alternative set of objectives  $F_0 = \emptyset$ .
  - (b) Calculate the portion of the total variance explained by the first  $k$  PCs ( $CVAR_k$ ) for  $k = 1, \dots, |F|$ . Retain for further analysis the minimal subset of PCs with a total cumulative variance above  $TC$ . Fathom the rest of PCs.
  - (c) Add to  $F_0$  the objectives with the most positive and most negative contribution to the eigenvector of the first PC.
  - (d) For the remaining PCs, proceed as follows. If  $\lambda_j \leq 0.1$ , add to  $F_0$  the objective with the highest contribution in absolute value. If  $\lambda_j > 0.1$ :
    - i. If all the components of the PC are positive, add to  $F_0$  the objective corresponding to the highest component of the eigenvector.
    - ii. If all the components of the PC are negative, add all the objectives to  $F_0$ .
    - iii. If none of the previous two cases apply, proceed as follows. Let  $mp_j$  be the most positive component of the PC under consideration and  $mn_j$  its most negative one.
      - A. If  $mp_j < 0.9|mn_j|$ , add the objective corresponding to  $mn_j$  to  $F_0$ .
      - B. If  $0.9|mn_j| \leq mp_j < |mn_j|$ , add to  $F_0$  the objectives associated with both  $mn_j$  and  $mp_j$ .
      - C. If  $0.8|mp_j| \leq mn_j < |mp_j|$ , add to  $F_0$  the objectives associated with both  $mn_j$  and  $mp_j$ .
      - D. If none of the previous three cases apply, add to  $F_0$  the objective associated with  $mp_j$ .
5. If  $F_0 \neq F$ , make  $F = F_0$  and go to step 2. Otherwise stop.

#### 2.3.1. Remarks

- In step 1.a, it is recommended to use a large value of  $TC$  (refer to Deb and Saxena (2005) for further discussion on this issue).

- Steps 1.a and 1.b depend on the PCA metrics used to identify redundant objectives and the solution procedure employed to solve MO, respectively. Hence, the initialization step may vary according to the methods of choice.
- In step 2, we can use any MOO solution procedure coupled with PCA. Even when the epsilon constraint method is employed, the original set of Pareto solutions used in the PCA analysis can be generated in different ways. One possible manner is to run the epsilon constraint method in the original search space, as mentioned above. An alternative approach is to execute the epsilon constraint method for an arbitrary reduced set of objectives.
- Although the number of iterations of the epsilon constraint method can be modified at will during the algorithm, it is recommended to increase the value of  $N$  (step 2.b) in successive iterations in which the number of objectives diminishes. This will increase the granularity of the data for a similar computational burden (see numerical examples for further clarification on this issue).
- The solutions generated in the reduced spaces obtained during the execution of the algorithm are all guaranteed to be weakly Pareto optimal in the original search space.
- Recall that the final step of the original heuristic by Deb and Saxena (2005), which is described in Section 5.2.3 of the referenced work, is not included in the framework proposed herein. Note also that, in the proposed strategy, there is no need to proceed with the objective reduction procedure once  $|F| = 2$ , since the first step of the heuristic will lead to retaining the two objectives which are already in the set.
- In step 4, different criteria based on the outcome of the PCA can be employed for identifying redundant objectives, as discussed by Gutiérrez et al. (2010).

### 3. Results and discussion

The design of chemical SCs is taken as a test-bed to illustrate the capabilities of the use of PCA in the multi-objective optimization of industrial processes. Note that the interest here is on the identification of redundancies between LCA metrics, and not on the analysis of the main features of the SC configurations obtained from the multi-objective optimization. For further details on the latter topic, the reader is referred to the original works by Guillén-Gosálbez and Grossmann (2009, 2010).

The environmentally conscious design of chemical SCs has been the focus of an increasing interest in the last years. In a seminar paper, Mele et al. (2005) addressed the optimization of SCs with economic and LCA-based environmental concerns through a combined simulation-optimization approach. Hugo and Pistikopoulos (2005) proposed a MILP formulation for the long-range planning and design of SCs, in which the environmental performance was measured via the Eco-indicator 99 (Eco99). Bojarski et al. (2009) introduced a MILP formulation for the design and planning of SCs considering economical and environmental issues, which incorporated the CML 2001 methodology to assess their environmental performance. Guillén-Gosálbez and Grossmann (2009) (see also Grossmann and Guillén-Gosálbez, 2010) proposed two MINLP formulations for the design of chemical SCs under uncertainty that explicitly consider the variability of the life cycle inventory of emissions and damage assessment model, respectively.

The problem of interest can be formally stated as follows. We are given a demand of chemicals to be satisfied in a set of final markets, a set of available manufacturing technologies and potential locations for plants and warehouses, and cost and environmental data associated with the SC operation. The goal is to determine the set of Pareto optimal SC configurations that maximize the net present value (NPV) and minimize a set of environmental metrics.

A detailed mathematical formulation for the problem described above can be found in the works by Guillén-Gosálbez and Grossmann (2009, 2010). For the sake of simplicity, we have omitted the treatment of the uncertainties associated with the life cycle inventory of emissions and damage assessment model. We therefore assume herein that all model parameters can be perfectly known in advance and do not show any variability. The constraints of the model can be roughly classified into three main blocks: mass balances, capacity limitations and objective function equations. Due to space limitations, we provide the complete formulation in the Appendix (the reader is referred to the works by Guillén-Gosálbez and Grossmann (2009) for further details).

#### 3.1. Numerical examples

Two case studies are presented to illustrate the usefulness of the proposed approach. These examples were first introduced by Guillén-Gosálbez and Grossmann (2009). We consider a superstructure based on a three-echelon SC (production-storage-market) with different available production technologies for plants, potential locations for SC entities and transportation links (see Fig. 1). The goal of the analysis is to determine the SC configuration along with the associated planning decisions that maximize the NPV and minimize the associated environmental impact.

In the original work by Guillén-Gosálbez and Grossmann (2009), the environmental performance was assessed through the Eco99. In contrast, the mathematical formulation considered in this work seeks to minimize simultaneously the three damage categories included in the Eco99: ecosystem quality (EQ), human health (HH), and damage to natural resources (NR), which gives rise to a 4-objective optimization problem. The environmental data have been updated with the latest version of the Eco-invent database (see Tables 1 and 5). On the other hand, the process data (cost, demand, yields, etc.) are those reported in Guillén-Gosálbez and Grossmann (2009).

The epsilon-constraint method was coded in GAMS 23.0.2 using CPLEX 11.2.1 as MILP solver. The objective reduction technique was implemented in MATLAB R2009a. A threshold cut  $TC = 95\%$  was used in all the cases. The numerical experiments were performed on an Intel 1.2 GHz machine.

##### 3.1.1. Case study 1

The first case study considers an existing SC (see Fig. 8 in Guillén-Gosálbez and Grossmann, 2009) comprising one plant, with a capacity of 100 kt/year, and one warehouse, with a capacity of 100 kt and an initial inventory of 0 kt. Both of them, the plant and the warehouse, are located in Tarragona (Spain). The demand of four different markets (Leuna in Germany, Neratovice in Czech Republic, Sines in Portugal and Tarragona in

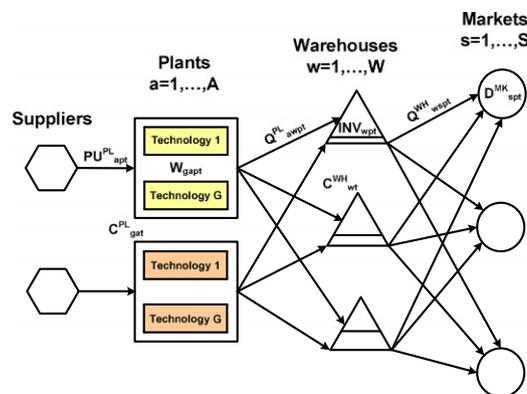


Fig. 1. Generic 3-echelon supply chain topology. Taken from Guillén-Gosálbez and Grossmann (2009).

**Table 1**

Life cycle impact assessment data (i.e., product of the life cycle inventory of emissions with the corresponding damage factors) used in the calculations of case study 1. All the damages are expressed in ecopoints per unit of reference flow.

Item	Reference units	EQ	HH	NR
Heat, heavy fuel oil, at 1 MW furnace	kg	0.0264598	0.1069444	0.1805763
Ammonia	kg	0.022189	0.085481	0.14725
Benzene	kg	0.0058944	0.026526	0.2073
Ethylene	kg	0.0012775	0.019804	0.20742
Hydrochloric acid	kg	0.0075906	0.0464	0.030403
Hydrogen cyanide	kg	0.0047149	0.083475	0.34225
Oxygen	kg	0.0013857	0.011564	0.0066515
Propylene	kg	0.0012435	0.02012	0.2117
Sodium hydroxide	kg	0.0056171	0.036556	0.020385
Sulfuric acid	kg	0.0029373	0.030643	0.0051743
Transport, lorry 20–28 t, fleet average	t km	0.0013865	0.0070944	0.0094519

Spain) must be satisfied, at least up to a minimum level of 85%, in every period in which the total time horizon of 3 years is divided. For this purpose, we consider six different technologies (see Fig. 9 in Guillén-Gosálbez and Grossmann, 2009). Table 1 displays the main environmental data, whereas the remaining process and cost data can be found in Tables 2–7 of the original work by Guillén-Gosálbez and Grossmann (2009). A detailed formulation of the model (which will be referred to as  $P1^1$  from here on) can also be found either in the aforementioned reference or in the Appendix.

Following our solution procedure, we first generated a set of Pareto solutions of  $P1^1$  in the original search space. For this, we optimized each single scalar objective separately (i.e., initialization step). This provided the lower and upper limits for each epsilon parameter. The epsilon intervals were then split into six subintervals, which led to 216 (i.e.,  $6^3$ ) single iterations. The 216 solutions were next filtered in order to remove infeasible and suboptimal solutions. We finally identified 16 non-dominated points, which were used in the PCA. These solutions were stored in matrix  $M1^1$ , which is presented in Table 2.

We next calculated the correlation matrix (see Table 3) in order to reveal whether the objectives were correlated. This additional step, not included as such in our algorithm, provides valuable information on the underlying relationships between the LCA metrics prior to the application of PCA.

As seen, the economic objective (NPV) is conflicting with the three environmental metrics (EQ, HH and NR). This was expected, given the traditional trade-off existing between economic and environmental criteria in many engineering applications. This observation suggests in turn that the three impact categories are somehow equivalent in this problem. The correlation between the environmental metrics is particularly strong. Henceforth, it was expected that the application of PCA would allow for significant reductions in the number of objectives.

Matrix  $M1^1$  was next standardized. This was done by subtracting from each measurement (i.e., Pareto point) the mean value of the corresponding column,  $\mu_j^1$  ( $\forall j$ ), so that the centroid of the data set became 0, and by dividing each value by the corresponding standard deviation,  $\sigma_j^1$  ( $\forall j$ ).

At this point, the principal components were computed, together with the associated eigenvalues, by solving an eigenvalue–eigenvector problem. For this, we used the correlation matrix  $R1^1$  (see step 3 of the algorithm), since the magnitudes of the four objectives were not the same. Note that the correlation matrix  $R1^1$  used to obtain the PCs is not the same as the correlation matrix shown in Table 3, which was obtained before the standardization of matrix  $M1^1$ . The results of the PCA are summarized in Table 4. Recall that the components of a PC denote the contribution of each individual objective toward that PC. A positive value indicates that the

**Table 2**

Solutions obtained for case study 1 after filtering the original 216 points generated by applying the epsilon constraint method to the 4-dimensional problem.

Solution	NPV (\$)	EQ	HH	NR
1	103,772	22,530,227	104,242,733	359,642,956
2	113,850	22,889,709	106,029,789	365,276,235
3	114,954	22,928,460	106,233,138	369,279,030
4	114,980	22,930,280	106,242,154	369,309,474
5	117,401	23,180,056	107,474,994	370,909,514
6	119,897	23,330,333	108,220,586	376,260,912
7	119,903	23,332,064	108,223,542	376,391,893
8	120,158	23,450,603	108,787,926	376,542,793
9	122,652	23,730,386	110,109,661	382,176,072
10	122,681	23,730,386	110,103,519	382,300,057
11	122,679	23,749,505	110,213,947	382,176,072
12	122,781	23,752,498	110,213,947	382,505,703
13	122,682	23,752,871	110,232,303	382,176,072
14	123,716	24,130,438	112,147,778	385,082,154
15	123,739	24,141,495	112,204,351	385,157,527
16	124,284	24,530,491	114,194,755	387,809,352

**Table 3**

Correlation matrix  $R1^1$  of matrix  $M1^1$ .

	NPV	EQ	HH	NR
NPV	1.0000	−0.8928	−0.8903	−0.9568
EQ	−0.8928	1.0000	0.9999	0.9674
HH	−0.8903	0.9999	1.0000	0.9649
NR	−0.9568	0.9674	0.9649	1.0000

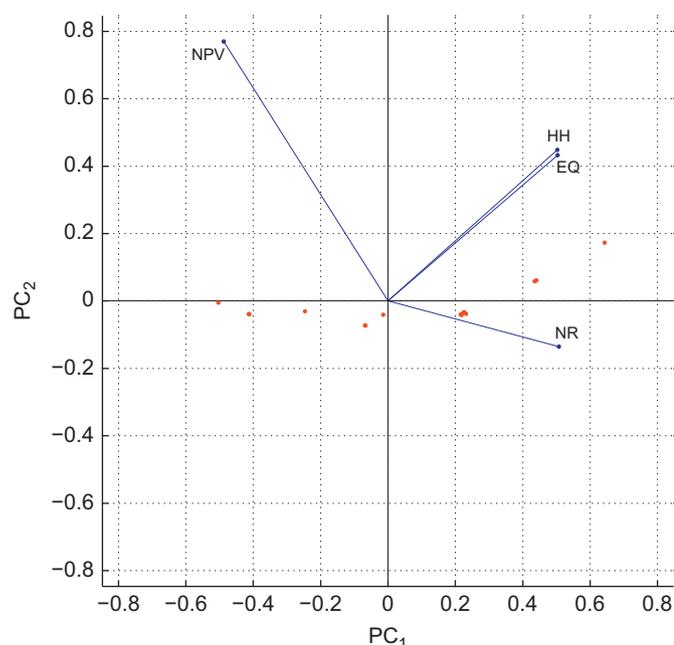
**Table 4**

PCs (presented in decreasing order of their eigenvalues) and other relevant data computed from the correlation matrix of the results obtained by solving the 4-dimensional multi-objective optimization in case study 1. Columns NPV, EQ, HH and NR show the contribution of the corresponding objective to each PC, whereas the amount of variance explained (i.e., the associated eigenvalue) is given in column  $VAR_j$ . Column %Exp. denotes the same quantity but expressed as a percentage of the total variance of the problem, whereas  $CVAR_k$  accounts for the accumulated percentage of variance explained by the first  $k$  PCs.

PC	NPV	EQ	HH	NR	$VAR_j$	%Exp.	$CVAR_k$ (%)
1	−0.4871	0.5033	0.5026	0.5068	3.8369	95.9236	95.9236
2	−0.7700	−0.4331	−0.4486	0.1350	0.1468	3.6704	99.5940
3	−0.4120	0.1941	0.2638	−0.8503	0.0162	0.4045	99.9985
4	−0.0116	0.7221	−0.6903	−0.0437	0.0001	0.0015	100.0000

objective is increased as we move along the PC axes, whereas a negative value indicates the opposite.

Fig. 2 depicts the objectives (vectors) and Pareto points (dots) in the space of the two first PCs, which explain 95.9 and 3.7% of



**Fig. 2.** The figure depicts the objectives (vectors) and Pareto points (dots) in the space of the two first PCs for case study 1.

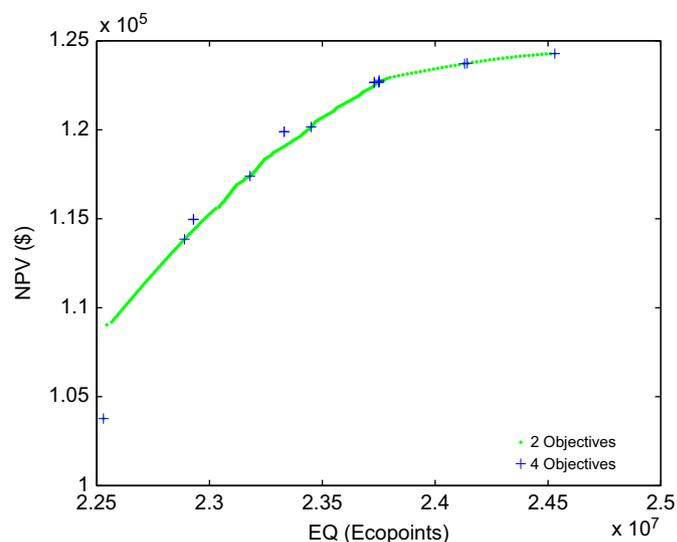
the total variance, respectively. Note that the projections of the vectors representing the environmental impact categories onto the  $PC_1$  axis are all positive, whereas that of the profit is negative. This is due to the existence of a clear trade-off between the economic and environmental performance of the network. In addition, the vectors associated with HH and EQ point towards the same direction, which indicates a very strong correlation between both metrics.

We next applied the heuristic proposed by Deb and Saxena (2005) to the PCA results. By selecting the most positive (NR) and most negative (NPV) objectives of the first PC, we obtained the two most conflicting objectives (those kept in the problem formulation). No further PCs were analyzed, since the variance explained by this first PC was above the  $TC$ .

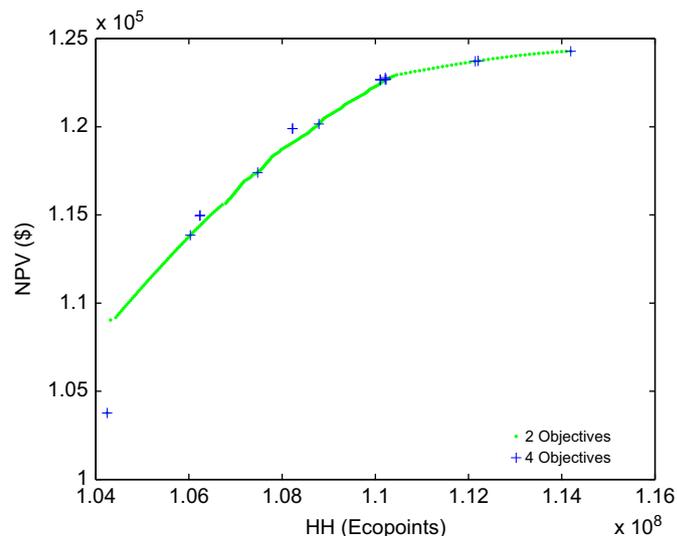
Observing Fig. 2, one could wrongly conclude that the data set is indeed 3-dimensional, as vectors representing the objectives point towards three different directions. It should be noted, however, that these three clusters appear only when the second PC, which explains 3.67% of the total variance and can thus be neglected, is considered in the analysis. In contrast, when the vectors are projected onto the  $PC_1$  axes, the 2-dimensionality of the problem is uncovered, as the three objectives, EQ, HH and NR, practically fall in the same spot.

We next removed the two redundant objectives (i.e., EQ and HH) from the pool, and generated a new set of Pareto solutions, this time considering only NPV and NR. We will refer to this problem as  $P1^2$ . Again, 216 iterations were run, but this time we concentrated on optimizing only the aforementioned two objectives. One single epsilon interval was therefore defined for NR, which was split into 216 subintervals giving rise again to 216 iterations.

The results obtained after applying the epsilon constraint method in this new reduced space are presented in Figs. 3–5. More precisely, these figures show the projections of the Pareto points onto 2-dimensional spaces. In dots, we have represented the solutions obtained when optimizing two objectives, whereas plus signs denote those determined when four objectives are considered. Note that the computational effort associated with the generation of both sets of solutions is indeed very similar, as



**Fig. 3.** Projections of the points obtained from the multi-objective optimization of case study 1 in the subspace of objectives NPV and EQ. Plus signs correspond to the results obtained with four objectives in the pool whereas dots represent solutions obtained with two objectives in the pool.



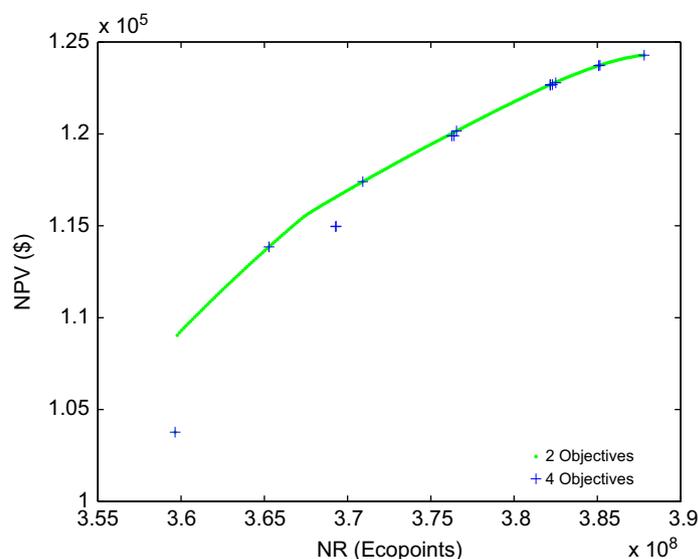
**Fig. 4.** Projections of the points obtained from the multi-objective optimization of case study 1 in the space of objectives NPV and HH. Plus signs correspond to the results obtained with four objectives in the pool whereas dots represent solutions obtained with two objectives in the pool.

we have run the same number of iterations (i.e., 216) in both cases. Hence, it is clear that for a very similar CPU time, the reduced model provides more Pareto points.

We should clarify that the computational savings of reducing objectives may vary from one solution method to another. Nevertheless, as pointed out by Deb and Saxena (2005), all of the MOO methods available tend to be very sensitive indeed to the number of objectives.

As can be seen in the figures, most of the information of the original 4-dimensional problem is still retained when two objectives are considered. That is, almost all the points generated when four objectives are optimized are identified after removing two of them from the optimization. Furthermore, the granularity of the data is dramatically improved after omitting the redundant objectives.

Let us clarify that all the points depicted in the figures are indeed Pareto optimal in the original 4-dimensional search space.



**Fig. 5.** Projections of the points obtained from the multi-objective optimization of case study 1 in the space of objectives NPV and NR. Plus signs correspond to the results obtained with four objectives in the pool whereas dots represent solutions obtained with two objectives in the pool.

**Table 5**

Life cycle impact assessment data (i.e., product of the life cycle inventory of emissions with the corresponding damage factors) used in the calculations of case study 2. All the damages are expressed in ecopoints/unit of reference flow.

Item	Reference units	EQ	HH	NR
Heat, heavy fuel oil, at 1 MW furnace	kg	0.0264598	0.1069444	0.1805763
Ammonia	kg	0.022189	0.085481	0.14725
Benzene	kg	0.0058944	0.026526	0.2073
Carbon monoxide	kg	0.0085804	0.050513	0.13189
Chlorine	kg	0.005501	0.035715	0.019729
Ethylene	kg	0.0012775	0.019804	0.20742
Hydrogen chloride	kg	0.0075906	0.0464	0.030403
Hydrogen cyanide	kg	0.0047149	0.083475	0.34225
Methane	kg	0.0496408	0.031777	0.017093
Methanol	kg	0.02341	0.014411	0.0099561
Methyl acetate	kg	0.0047723	0.025554	0.13676
Nitrogen	kg	0.0014686	0.012256	0.0070495
Oleum	kg	0.0029373	0.030643	0.0051743
Oxygen	kg	0.0013857	0.011564	0.0066515
Propylene	kg	0.0012435	0.02012	0.2117
Sodium hydroxide	kg	0.0056171	0.036556	0.020385
Sulfuric acid	kg	0.0029373	0.030643	0.0051743
Transport, lorry 20–28 t, fleet average	t km	0.0013865	0.007094	0.0094519

This is because, as pointed out by Brockhoff and Zitzler (2009), any Pareto solution of a MOO problem defined for any subset  $F_0 \subset F$ , is also Pareto optimal in the higher dimensional space  $F$ .

Recall that solutions depicted as dots in Fig. 5 belong to the Pareto optimal frontier of the 2-dimensional problem. The remaining points are projections of the Pareto solutions generated in a higher dimensional space onto the corresponding axes. For instance, plus sign points in Figs. 3–5 correspond to projections of the 4-dimensional Pareto set, as commented above. Similarly, dot points in Figs. 3 and 4 are projections of the Pareto frontier depicted in Fig. 5 onto the same axis. As observed in Fig. 5, some points of the original search space cannot be identified after reducing the dimensionality of the problem. This is because these points are suboptimal when NPV and NR are the only objectives considered in the analysis.

The application of the PCA algorithm presented in this paper allows omitting two objectives. Henceforth, no further objective

reduction is possible and no more iterations are required in this case (see Remark 6).

### 3.1.2. Case study 2

This second case study considers another SC to be set in Europe (see Fig. 18 in Guillén-Gosálbez and Grossmann, 2009). In this case, there are no existing facilities already under operation. Twenty technologies manufacturing 14 different products, some of which can be recycled and used as raw materials for other processes, are available (refer to Fig. 19 in Guillén-Gosálbez and Grossmann, 2009). No limits on the purchases of raw materials and intermediates are considered, while no outsourcing is allowed for final products. Four final markets are considered (Kazinbarcika in Hungary, Wloclawek in Poland, Neratovice and Tarragona), in which a minimum demand satisfaction level of 97.5% must be attained. The problem data can be found in Tables 9–15 of

the work by Guillén-Gosálbez and Grossmann (2009). Refer also to the original work by Guillén-Gosálbez and Grossmann (2009) and to the Appendix for a detailed formulation of the model ( $P2^1$  from here on).

The same solution strategy was applied to this example. Each scalar objective was first optimized separately in order to identify the extreme limits for each single objective. The associated intervals were then split into six subintervals, resulting in 216 iterations of the epsilon-constraint method. Only 11 non-dominated solutions were identified after filtering the set of 216 points. These points were stored in matrix  $M2^1$ , which is shown in Table 6.

We next constructed the correlation matrix (see Table 7) of matrix  $M2^1$  in order to identify redundancies between objectives. The analysis of this matrix reveals the existence of a conflict between the economic (NPV) and environmental objectives (EQ, HH and NR). Further, the correlation between the three environmental objectives is very strong (i.e., equal to 1), which indicates the existence of redundant metrics.

We next standardized matrix  $M2^1$  in order to compute its correlation matrix  $R2^1$  and calculate the PCs (see Table 8). Note

**Table 6**

Solutions obtained for case study 2 after filtering the original 216 points generated by applying the epsilon constraint method to the 4-dimensional problem.

Solution	NPV (\$)	EQ	HH	NR
1	357,919	24,149,124	97,871,816	164,805,370
2	376,676	24,867,509	10,088,811	169,707,241
3	380,078	25,620,368	103,921,340	174,845,248
4	380,161	25,636,325	103,994,502	174,954,090
5	380,161	25,636,400	103,994,804	174,954,599
6	382,139	26,370,472	106,946,101	179,964,426
7	382,163	26,379,926	106,984,314	180,028,948
8	382,163	26,379,965	106,984,471	180,029,213
9	384,481	27,118,769	109,970,863	185,071,226
10	384,486	27,121,328	109,981,274	185,088,692
11	386,750	27,867,127	112,995,625	190,178,443

**Table 7**

Correlation matrix  $R2^1$  of matrix  $M2^1$ .

	NPV	EQ	HH	NR
NPV	1.0000	-0.8457	-0.8491	-0.8457
EQ	-0.8457	1.0000	1.0000	1.0000
HH	-0.8491	1.0000	1.0000	1.0000
R	-0.8457	1.0000	1.0000	1.0000

**Table 8**

PCs (presented in decreasing order of their eigenvalues) and other relevant data computed from the correlation matrix of the results obtained by solving the 4-dimensional multi-objective optimization in case study 2. Columns NPV, EQ, HH and R show the contribution of the corresponding objective to each PC, whereas the amount of variance explained (i.e., the associated eigenvalue) is given in column  $VAR_j$ . Column %Exp. denotes the same quantity but expressed as a percentage of the total variance of the problem, whereas  $CVAR_k$  accounts for the accumulated percentage of variance explained by the first  $k$  PCs

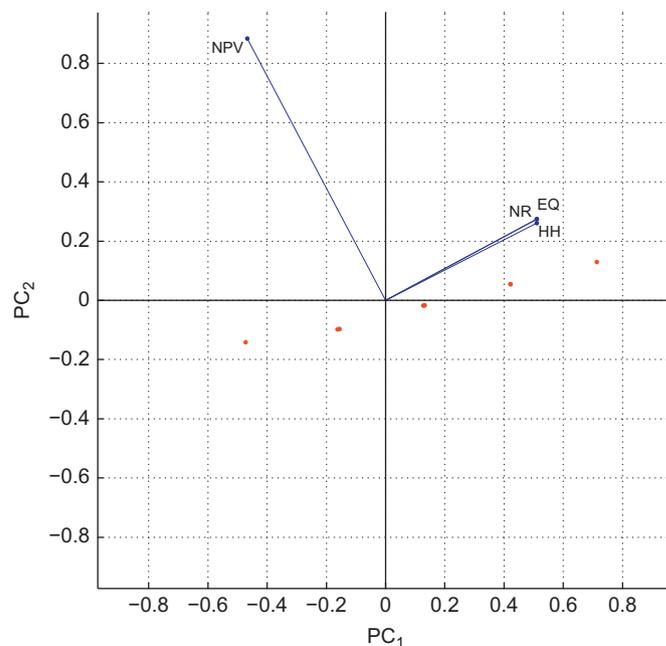
PC	NPV	EQ	HH	NR	$VAR_j$	%Exp.	$CVAR_k$ (%)
1	-0.4673	0.5103	0.5107	0.5103	3.7752	94.3806	94.3806
2	-0.8841	-0.2742	-0.2608	-0.2743	0.2248	5.6193	99.9999
3	-0.0099	0.4029	-0.8192	0.4080	0.0000	0.0000	99.9999
4	0.0000	0.7086	-0.0029	-0.7056	0.0000	0.0000	99.9999

that this correlation matrix is not the same as the one presented in Table 7, which provided preliminary information on the relationships between objectives.

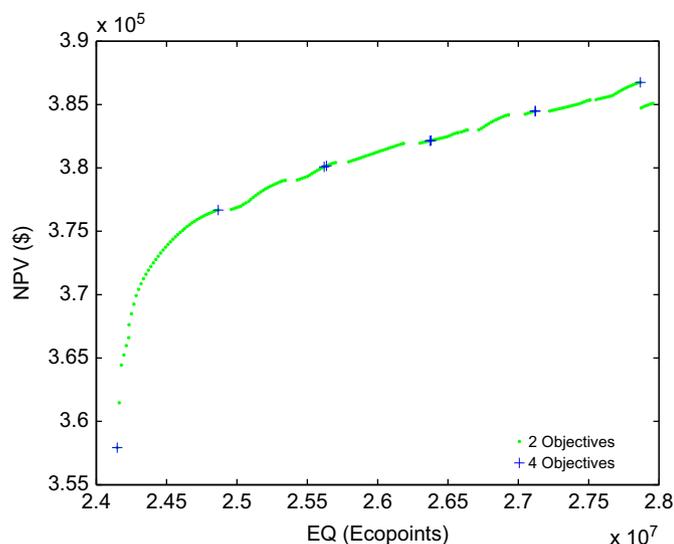
Following the proposed algorithm, the objectives with the most positive and most negative contribution to the first PC (HH and NPV, respectively) were selected and added to the reduced set of objectives  $F_0$ . In this case, the first PC, which explains 94.38% of the variance, is not enough to attain the proposed TC of 95%. Note however, that the problem's variance explained by the second PC is rather small compared to that explained by the first PC. This indicates that the data set in  $M2^1$  is essentially one-dimensional regarding the PCs and bi-dimensional regarding the objectives (see Fig. 6). Henceforth, no more PCs, and consequently objectives, were considered. Objectives EQ and NR were then removed from the pool by making  $F = F_0$ .

After omitting the redundant objectives, we generated a new set of Pareto optimal solutions for problem  $P2^2$  considering the reduced objective set  $F = \{NPV, HH\}$ . Since the dimensionality of the problem was reduced to 2 in the first iteration of the algorithm, it was possible to concentrate all the iterations of the epsilon constraint method on a single objective. Hence, the single epsilon parameter was divided into the same number of subintervals as iterations were performed with 4 objectives (i.e., 216). The results of these calculations were used to construct matrix  $M2^2$ . The new Pareto solutions are depicted in Figs. 7–9 (dots), which also show the points obtained from the minimization of four objectives (plus signs).

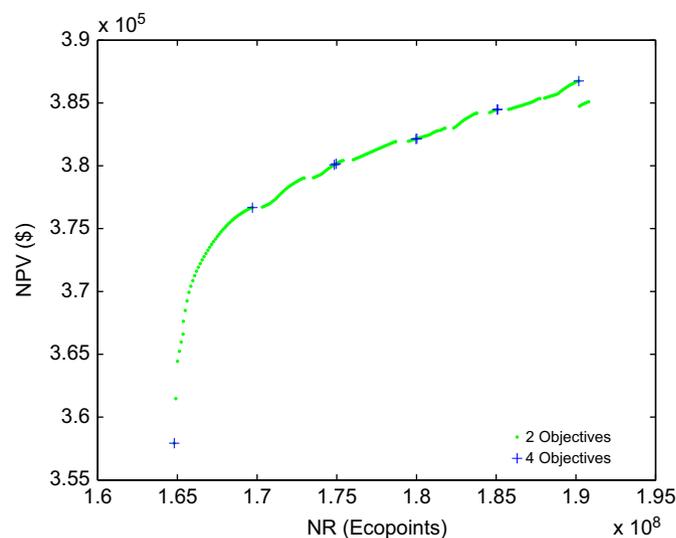
As seen in the figures, for a similar computational effort, it is possible to generate more non-dominated solutions when only two objectives are considered. Moreover, the Pareto set generated in the reduced space contains most of the points calculated with the original 4-dimensional formulation. Hence, it seems clear that the proposed approach provides a better representation of the Pareto set for the same CPU time. Note also that it would be possible to generate the 11 solutions obtained in the high dimensional space by running the epsilon constraint method in the reduced space with only 11 iterations, which would lead to significant CPU savings. At this point  $|F| = 2$ , so there is no need to perform further reductions in the number of objectives and the algorithm stops.



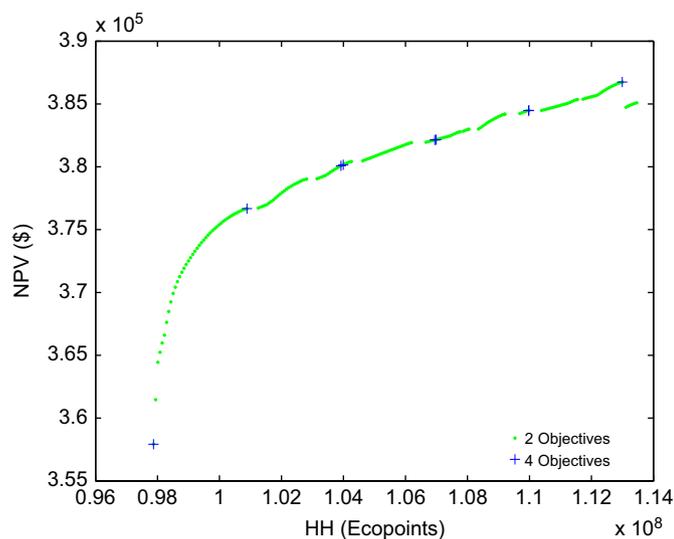
**Fig. 6.** The figure depicts the objectives (vectors) and Pareto points (dots) in the space of the two first PCs for case study 2.



**Fig. 7.** Projections of the points obtained from the multi-objective optimization of case study 2 in the subspace of objectives NPV and EQ. Plus signs correspond to the results obtained with four objectives in the pool whereas dots represent solutions obtained with two objectives in the pool.



**Fig. 9.** Projections of the points obtained from the multi-objective optimization of case study 2 in the space of objectives NPV and NR. Plus signs correspond to the results obtained with four objectives in the pool whereas dots represent solutions obtained with two objectives in the pool.



**Fig. 8.** Projections of the points obtained from the multi-objective optimization of case study 2 in the space of objectives NPV and HH. Plus signs correspond to the results obtained with four objectives in the pool whereas dots represent solutions obtained with two objectives in the pool.

#### 4. Conclusions

This paper has investigated the use of PCA for dimensionality reduction in environmental engineering problems. We have presented a novel method for solving MOO problems that integrates PCA into the standard epsilon constraint method in order to overcome the numerical difficulties that arise when dealing with a large number of environmental objectives.

The capabilities of the proposed approach have been illustrated through its application to the design of chemical supply chains. It has been clearly shown how PCA can be effectively employed to identify redundant environmental metrics that can be left out of the analysis while still preserving to a large extent its Pareto structure. In addition, the use of PCA coupled with MOO provides valuable insights on the relationships between the

environmental metrics considered in the analysis, facilitating the task of decision-makers during the analysis and interpretation of the Pareto set of an environmental problem.

#### Nomenclature

##### Abbreviations

LCA	Life Cycle Assessment
MOO	Multi-Objective Optimization
PC	Principal Component
PCA	Principal Component Analysis
SC	Supply Chain
SCM	Supply Chain Management

##### Set/indices

$A$	set of plants indexed by $a$
$B$	set of environmental burdens indexed by $b$
$D$	set of damage categories indexed by $d$
$F$	set of objectives indexed by $i$
$G$	set of manufacturing technologies indexed by $g$
$J$	set of principal components indexed by $j$
$N$	set of epsilon constraint method subintervals indexed by $n$
$P$	set of products indexed by $p$
$S$	set of markets indexed by $s$
$T$	set of time periods indexed by $t$
$W$	set of warehouses indexed by $w$

##### Subsets

$F_0$	subset of objectives $F_0 \subset F$
$IN(p)$	set of manufacturing technologies that consume $p$
$MP(g)$	set of main products $p$ of technology $g$
$OUT(p)$	set of manufacturing technologies that produce $p$

##### Parameters

$\overline{CE}_{gat}^{PL}$	upper bound on the capacity expansion of manufacturing technology $g$ at plant $a$ in time period $t$
$\underline{CE}_{gat}^{PL}$	lower bound on the capacity expansion of manufacturing technology $g$ at plant $a$ in time period $t$

$\overline{C}_{wt}^{WH}$	upper bound on the capacity expansion of warehouse $w$ in time period $t$	$\beta_{wst}^{TWH}$	fixed investment term associated with the establishment of a transport link between warehouse $w$ and market $s$ in time period $t$
$\underline{C}_{wt}^{WH}$	lower bound on the capacity expansion of warehouse $w$ in time period $t$	$\omega_{bp}^{PU}$	emissions/feedstock requirements of chemical $b$ per unit of raw material $p$ generated
$\overline{D}_{spt}^{MK}$	maximum demand of product $p$ sold at market $s$ in period $t$	$\omega_{bp}^{PR}$	emissions/feedstock requirements of chemical $b$ per unit of intermediate/final product $p$ generated
$\underline{D}_{spt}^{MK}$	minimum demand of product $p$ to be satisfied at market $s$ in period $t$	$\omega_b^{EN}$	emissions/feedstock requirements of chemical $b$ per unit of energy consumed
$ir$	interest rate	$\omega_b^{TR}$	emissions/feedstock requirements of chemical $b$ per unit of mass transported one unit of distance
$\overline{f}_i$	upper bound on objective function $f_i$	$\eta_{gap}^{EN}$	energy consumed per unit of chemical $p$ produced with manufacturing technology $g$ at plant $a$
$\underline{f}_i$	lower bound on objective function $f_i$	$\delta_{aw}^{PL}$	distance between plant $a$ and warehouse $w$
$\overline{FCI}$	upper limit on the total capital investment	$\delta_{ws}^{WH}$	distance between warehouse $w$ and market $s$
$k$	number of PCs selected ( $k \leq m$ ) to calculate the cumulative variance they explain	$\theta_{bd}$	damage factor of chemical $b$ contributing to damage category $d$
$m$	number of original objectives $i$	$\tau$	minimum desired percentage of the available installed capacity that must be utilized
$N$	number of subintervals $n$ in the epsilon constraint method		
$NEXD_{ga}^{PL}$	maximum number of capacity expansions for technology $g$ available at plant $a$		
$NEXD_w^{WH}$	maximum number of capacity expansions for warehouse $w$		
$NT$	number of time periods		
$\overline{PU}_{apt}$	upper bound on the purchases of product $p$ at plant $a$ in period $t$		
$\underline{PU}_{apt}$	lower bound on the purchases of product $p$ at plant $a$ in period $t$		
$\overline{Q}_{awt}^{PL}$	upper bound on the flow of materials between plant $a$ and warehouse $w$ in time period $t$		
$\underline{Q}_{awt}^{PL}$	lower bound on the flow of materials between plant $a$ and warehouse $w$ in time period $t$		
$\overline{Q}_{wst}^{WH}$	upper bound on the flow of materials between warehouse $w$ and market $s$ in time period $t$		
$\underline{Q}_{wst}^{WH}$	lower bound on the flow of materials between warehouse $w$ and market $s$ in time period $t$		
$SV$	salvage value		
$TC$	threshold cut in the heuristic proposed by Deb and Saxena (2005)		
$TOR_w$	turnover ratio of warehouse $w$		
$\epsilon_i^n$	epsilon constraint parameter for subinterval $n$ of objective $i$		
$\xi_{gp}$	mass balance coefficient associated with product $p$ and manufacturing technology $g$		
$\varphi$	tax rate		
$\gamma_{spt}^{FP}$	price of final product $p$ sold at market $s$ in time period $t$		
$\gamma_{apt}^{RM}$	price of raw material $p$ purchased at plant $a$ in time period $t$		
$\nu_{gapt}$	operating cost of manufacturing technology $g$ available at plant $a$ per unit of main product $p$ in time period $t$		
$\pi_{wt}$	inventory cost at warehouse $w$ in period $t$		
$\psi_{awpt}^{PL}$	unit transportation cost of product $p$ sent from plant $a$ to warehouse $w$ in time period $t$		
$\psi_{wspt}^{WH}$	unit transportation cost of product $p$ sent from warehouse $w$ to market $s$ in time period $t$		
$\alpha_{gat}^{PL}$	variable investment term associated with technology $g$ at plant $a$ in time period $t$		
$\alpha_{wt}^{WH}$	variable investment term associated with warehouse $w$ in time period $t$		
$\beta_{gat}^{PL}$	fixed investment term associated with technology $g$ at plant $a$ in time period $t$		
$\beta_{wt}^{WH}$	fixed investment term associated with warehouse $w$ in time period $t$		
$\beta_{awt}^{TPL}$	fixed investment term associated with the establishment of a transport link between the plant $a$ and warehouse $w$ in time period $t$		
		<b>Variables</b>	
		$C_{gat}^{PL}$	capacity of manufacturing technology $g$ at plant $a$ in time period $t$
		$CE_{gat}^{PL}$	capacity expansion of manufacturing technology $g$ at plant $a$ in time period $t$
		$C_{wt}^{WH}$	capacity of warehouse $w$ in time period $t$
		$CE_{wt}^{WH}$	capacity expansion of warehouse $w$ in time period $t$
		$CF_t$	cash flow in period $t$
		$CVAR_k$	cumulative variance explained by the first $k$ PCs expressed as a percentage
		$DAM_d$	impact in damage category $d$
		$DEP_t$	depreciation term in period $t$
		$EQ$	damage to the ecosystem quality
		$f_i$	objective function $i$
		$FCI$	fixed capital investment
		$FTDC_t$	fraction of the total depreciable capital that must be paid in period $t$
		$HH$	damage to human health
		$IL_{wt}$	average inventory level at warehouse $w$ in time period $t$
		$INV_{wpt}$	inventory of product $p$ kept at warehouse $w$ in period $t$
		$L$	diagonal $m \times m$ matrix storing the eigenvalues $\lambda_j$ arranged in descending order of their magnitudes
		$LCI_b$	life cycle inventory entry (i.e., emissions/feedstock requirements) associated with chemical $b$
		$M$	matrix storing solutions generated with the epsilon constraint method
		$mp_j$	most positive component of $PC_j$
		$mn_j$	most negative component of $PC_j$
		$NE_t$	net earnings in period $t$
		$NPV$	net present value
		$P$	$m \times m$ matrix storing the PC components arranged in descending order of their eigenvalues $\lambda_j$
		$PU_{apt}$	purchases of product $p$ made by plant $a$ in period $t$
		$Q_{awpt}^{PL}$	flow of product $p$ sent from plant $a$ to warehouse $w$ in period $t$
		$Q_{wspt}^{WH}$	flow of product $p$ sent from warehouse $w$ to market $s$ in period $t$
		$R$	correlation matrix
		$NR$	damage to natural resources
		$SA_{spt}$	sales of product $p$ at market $s$ in time period $t$
		$VAR_j$	amount of variance explained by $PC_j$
		$W_{gapt}$	input/output flow of product $p$ associated with technology $g$ at plant $a$ in $t$

$X_{gat}^{PL}$	binary variable (1 if the capacity of manufacturing technology $g$ at plant $a$ is expanded in time period $t$ , 0 otherwise)
$X_{wt}^{WH}$	binary variable (1 if the capacity of warehouse $w$ is expanded in time period $t$ , 0 otherwise)
$Y_{awt}^{PL}$	binary variable (1 if a transportation link between plant $a$ and warehouse $w$ is established in time period $t$ , 0 otherwise)
$Y_{wst}^{WH}$	binary variable (1 if a transportation link between warehouse $w$ and market $s$ is established in time period $t$ , 0 otherwise)
$\lambda_j$	eigenvalue associated to $PC_j$
$\mu_i$	mean value of objective $i$
$\sigma_i$	standard deviation of objective $i$
$\Sigma$	covariance matrix

$$\sum_w Q_{wspt}^{WH} = SA_{spt} \quad \forall s, p, t \quad (10)$$

### A.1.2. Capacity constraints

**Plants:** Eq. (11) bounds the capacity expansion in each time period ( $CE_{gat}^{PL}$ ), whereas Eq. (12) determines the total capacity in period  $t$  ( $C_{gat}^{PL}$ ) from the previous capacity and the capacity expansion carried out in the current period. Eq. (13) limits the number of times that the capacity of technology  $g$  available at plant  $a$  can be expanded during the entire planning horizon. In constraints (11) and (13),  $X_{gat}^{PL}$  is a binary variable that takes the value of 1 if the capacity of technology  $g$  established in plant  $a$  is expanded in period  $t$  and 0 otherwise. Further,  $NEXP_{ga}^{PL}$  represents the maximum number of capacity expansions of technology  $g$  that can be executed in plant  $a$  during the entire time horizon. Finally, Eq. (14) imposes lower and upper production limits based on the existing capacities. Here,  $\tau$  is a parameter that denotes the minimum percentage of the capacity that must be utilized. Note that this parameter should always be in the range [0,1]

$$\underline{CE}_{gat}^{PL} X_{gat}^{PL} \leq CE_{gat}^{PL} \leq \overline{CE}_{gat}^{PL} X_{gat}^{PL} \quad \forall g, a, t \quad (11)$$

$$C_{gat}^{PL} = C_{gat-1}^{PL} + CE_{gat}^{PL} \quad \forall g, a, t \quad (12)$$

$$\sum_t X_{gat}^{PL} \leq NEXP_{ga}^{PL} \quad \forall g, a \quad (13)$$

$$\tau C_{gat}^{PL} \leq W_{gap} \leq C_{gat}^{PL} \quad \forall g, a, t \quad (14)$$

**Warehouses:** Constraints (15)–(17) are equivalent to Eqs. (11)–(13), but apply to warehouses

$$\underline{CE}_{wt}^{WH} X_{wt}^{WH} \leq CE_{wt}^{WH} \leq \overline{CE}_{wt}^{WH} X_{wt}^{WH} \quad \forall w, t \quad (15)$$

$$C_{wt}^{WH} = C_{wt-1}^{WH} + CE_{wt}^{WH} \quad \forall w, t \quad (16)$$

$$\sum_t X_{wt}^{WH} \leq NEXP_w^{WH} \quad \forall w \quad (17)$$

In these equations  $C_{wt}^{WH}$  and  $CE_{wt}^{WH}$  denote the total capacity and capacity expansion associated with warehouse  $w$  in period  $t$ , respectively. On the other hand,  $X_{wt}^{WH}$  is a binary variable that takes the value of 1 if the capacity of a warehouse is expanded in period  $t$  and 0 otherwise, whereas  $NEXP_w^{WH}$  represents the maximum allowable number of times that the capacity of a warehouse can be expanded. Eqs. (18) and (19) impose limits on the inventory kept at each warehouse at the end of period  $t$  ( $INV_{wpt}$ ) and also on the average inventory level ( $IL_{wt}$ ), which is calculated via Eq. (20). Note that Eq. (18) accounts for possible fluctuations in the demand by assuming that the capacity required to handle a given amount of products, considering regular shipment and delivery schedule, is twice the average storage inventory level kept at the warehouse (Shapiro, 2001)

$$\sum_p INV_{wpt} \leq C_{wt}^{WH} \quad \forall w, t \quad (18)$$

$$2IL_{wt} \leq C_{wt}^{WH} \quad \forall w, t \quad (19)$$

$$IL_{wt} = \frac{\sum_s \sum_p Q_{wspt}^{WH}}{TOR_w} \quad \forall w, t \quad (20)$$

**Transportation links:** The amount of products sent from plants to warehouses ( $Q_{awt}^{PL}$ ) and from warehouses to plants ( $Q_{wst}^{WH}$ ) must lie between upper and lower limits, provided a transportation link between the corresponding nodes of the network is established, as stated in Eqs. (21) and (22). In these equations,  $Y_{awt}^{PL}$  and  $Y_{wst}^{WH}$  are binary variables that represent the existence of transportation links between plants and warehouses and between warehouses

## Acknowledgments

The authors wish to acknowledge support of this research work from the Spanish Ministry of Education and Science (projects DPI2008-04099, PHB2008-0090-PC, BFU2008-00196 and CTQ2009-14420-C02), the Spanish Ministry of External Affairs (projects HS2007-0006 and A/023551/09), the Spanish Ministry of Science and Innovation (ENE2011-28269-C03-03) and the Generalitat de Catalunya (FI programs).

## Appendix A

### A.1. Mathematical formulation

For the sake of completeness of this work, we next present an outline of the mathematical formulation used in the case studies. A more detailed model can be found in the work by Guillén-Gosálbez and Grossmann (2009). Similar formulations are also available in the works by Guillén et al. (2007) and Láinez et al. (2007).

#### A.1.1. Mass balances

The mass balances associated with the manufacturing plants and warehouses are expressed via constraints (5) and (6) (for plants), and (7) (for warehouses)

$$PU_{apt} + \sum_{g \in OUT(p)} W_{gap} = \sum_w Q_{awpt}^{PL} + \sum_{g \in IN(p)} W_{gap} \quad \forall a, p, t \quad (5)$$

$$W_{gap} = \xi_{gp} W_{gap} \quad \forall g, a, p, t \quad \forall p' \in MP(g) \quad (6)$$

$$INV_{wpt-1} + \sum_a Q_{awpt}^{PL} = \sum_s Q_{wspt}^{WH} + INV_{wpt} \quad \forall w, p, t \quad (7)$$

In these equations  $PU_{apt}$  represents the amount of product  $p$  purchased by plant  $a$  in period  $t$ ,  $W_{gap}$  is the input/output flow of  $p$  associated with technology  $g$  at plant  $a$  in  $t$ ,  $Q_{awpt}^{PL}$  and  $Q_{wspt}^{WH}$  are the flows of  $p$  between plant  $a$  and warehouse  $w$  and between warehouse  $w$  and market  $s$ , respectively, in  $t$ ,  $\xi_{gp}$  is a material balance coefficient and  $INV_{wpt}$  is the inventory of  $p$  kept at warehouse  $w$  at the end of period  $t$ . Constraints (8) and (9) impose lower and upper limits on the purchases of raw materials ( $PU_{apt}$ ) and sales of products ( $SA_{spt}$ ), respectively. The sales of products are determined from the flows of materials between the warehouses and the final markets, as shown in Eq. (10).

$$\underline{PU}_{apt}^{PL} \leq PU_{apt}^{PL} \leq \overline{PU}_{apt}^{PL} \quad \forall a, p, t \quad (8)$$

$$\underline{D}_{spt}^{MK} \leq SA_{spt} \leq \overline{D}_{spt}^{MK} \quad \forall s, p, t \quad (9)$$

and markets, respectively

$$Q_{awt}^{PL} Y_{awt}^{PL} \leq Q_{awt}^{PL} \leq \overline{Q_{awt}^{PL}} Y_{awt}^{PL} \quad \forall a, w, t \quad (21)$$

$$Q_{wst}^{WH} Y_{wst}^{WH} \leq Q_{wst}^{WH} \leq \overline{Q_{wst}^{WH}} Y_{wst}^{WH} \quad \forall w, s, t \quad (22)$$

### A.1.3. Objective function

NPV: Eqs. (23)–(30) allow to calculate the NPV, which is determined from the cash flow in each period  $t$  ( $CF_t$ )

$$NPV = \sum_t \frac{CF_t}{(1+ir)^{t-1}} \quad (23)$$

$$CF_t = NE_t - FTDC_t, \quad t = 1, \dots, NT - 1 \quad (24)$$

$$CF_t = NE_t - FTDC_t + SVFCI, \quad t = NT \quad (25)$$

$$NE_t = (1-\phi) \left[ \sum_s \sum_p \gamma_{spt}^{FP} SA_{spt} - \sum_a \sum_p \gamma_{apt}^{RM} PU_{apt} - \sum_g \sum_a \sum_{p \in MP(g)} v_{gapt} W_{gapt} - \sum_w \tau_{wt} IL_{wt} - \sum_a \sum_w \sum_p \psi_{awpt}^{PL} Q_{awpt}^{PL} - \sum_w \sum_s \sum_p \psi_{wsp}^{WH} Q_{wsp}^{WH} \right] + \phi DEP_t \quad \forall t \quad (26)$$

$$DEP_t = \frac{(1-SV)FCI}{NT} \quad \forall t \quad (27)$$

$$FCI = \sum_g \sum_a \sum_t (\alpha_{gat}^{PL} CE_{gat}^{PL} + \beta_{gat}^{PL} X_{gat}^{PL}) + \sum_w \sum_t (\alpha_{wt}^{WH} CE_{wt}^{WH} + \beta_{wt}^{WH} X_{wt}^{WH}) + \sum_a \sum_w \sum_t (\beta_{awt}^{TPL} v_{awt}^{PL}) + \sum_w \sum_s \sum_t (\beta_{wst}^{TWH} Y_{wst}^{WH}) \quad (28)$$

$$FCI \leq \overline{FCI} \quad (29)$$

$$FTDC_t = \frac{FCI}{NT} \quad \forall t \quad (30)$$

The cash flow is calculated from the net earnings ( $NE_t$ ), and the fixed investment term in period  $t$  ( $FTDC_t$ ), as stated in Eqs. (24) and (25). The net earnings are given by the difference between the incomes (i.e., sales of products) and the total cost, which accounts for the raw materials cost, the operating and inventory costs and the transportation expenses, as shown in Eq. (26). Eq. (27) is used to determine the depreciation of the capital invested ( $DEP_t$ ) assuming the straight-line method. The total fixed cost investment FCI appearing in Eq. (27) is determined from the capacity expansions made in plants and warehouses and the establishment of transportation links within the time horizon, as stated in Eq. (28). Eq. (29) imposes an upper limit on the total capital investment. Finally, Eq. (30) determines the payments of the total fixed capital investment, which are assumed to be equally distributed over time.

**Environmental impact:** the environmental performance of the network is quantified by the three damage categories following the Eco-indicator 99 methodology (PRé-Consultants, 2000). This requires the calculation of the life cycle inventory associated with the SC operation (Eq. (31)), which includes the emissions released and feedstock requirements ( $LCI_b$ ) associated with the manufacturing and distribution tasks carried out in the SC entities. These entries of the life cycle inventory, which are given by the production of raw materials ( $PU_{apt}$ ), the manufacture of final products ( $W_{gapt}$ ) and the transportation of materials between plants and warehouses ( $Q_{awpt}^{PL}$ ) and warehouses and final markets ( $Q_{wsp}^{WH}$ ), are then translated into a set of environmental damages

( $DAM_d$ ) by using Eq. (32). In the latter equation,  $\theta_{bd}$  represents the impact in damage category  $d$  per unit of chemical  $b$  released to/extracted from the environment.

$$LCI_b = \sum_a \sum_p \sum_t \omega_{bp}^{PU} PU_{apt} + \sum_g \sum_a \sum_{p \in MP(g)} \sum_t \omega_{bp}^{PR} W_{gapt} + \sum_g \sum_a \sum_{p \in MP(g)} \sum_t \omega_b^{EN} \eta_{gap}^{EN} W_{gapt} + \sum_a \sum_w \sum_p \sum_t \omega_b^{TR} \zeta_{aw}^{PL} Q_{awpt}^{PL} \times \sum_w \sum_s \sum_p \sum_t \omega_b^{TR} \zeta_{ws}^{WH} Q_{wsp}^{WH} \quad \forall b \quad (31)$$

$$DAM_d = \sum_b \theta_{bd} LCI_b \quad \forall d \quad (32)$$

## References

- Abdi, H., Williams, L., 2010. Principal component analysis. Wiley Interdiscipl. Rev.: Comput. Stat. 2 (4), 433–459.
- Azapagic, A., Clift, R., 1999. The application of life cycle assessment to process optimisation. Comput. Chem. Eng. 23 (10), 1509–1526.
- Bojarski, A., Láinez, J., Espuna, A., Puigjaner, L., 2009. Incorporating environmental impacts and regulations in a holistic supply chains modeling: an lca approach. Comput. Chem. Eng. 33 (10), 1747–1759.
- Brockhoff, D., Zitzler, E., 2009. Are all objectives necessary? On dimensionality reduction in evolutionary multiobjective optimization. Parallel Problem Solving from Nature, PPSN, I.X., Proceedings, vol. 4193. Springer-Verlag, Berlin, Germany, pp. 1523–1534.
- Cavin, L., Fischer, U., Glover, F., Hungerbühler, K., 2004. Multi-objective process design in multi-purpose batch plants using a tabu search optimization algorithm. Comput. Chem. Eng. 28 (4), 459–478.
- Chakraborty, A., Linninger, A., 2002. Plant-wide wastewaters management, 1. Synthesis and multiobjective design. Ind. Eng. Chem. Res. 41 (18), 4591–4604.
- Chang, C., Hwang, J., 1996. A multiobjective programming approach to waste minimization in the utility systems of chemical processes. Chem. Eng. Sci. 51 (16), 3951–3965.
- Dantus, M., High, K., 2000. Evaluation of waste minimization alternatives under uncertainty: a multiobjective optimization approach. Comput. Chem. Eng. 23 (10), 1493–1508.
- Deb, K., Saxena, D., 2005. On finding pareto-optimal solutions through dimensionality reduction for certain large-dimensional multi-objective optimization problems. Technical Report. Kanpur Genetic Algorithms Laboratory (KanGAL), Indian Institute of Technology, Kanpur.
- Ehrgott, M., 1998. Multicriteria Optimization. Springer, Berlin.
- Gebreslassie, B., Guillén-Gosálbez, G., Jiménez, L., Boer, D., 2009. Design of environmentally conscious absorption cooling systems via multi-objective optimization and life cycle assessment. Appl. Energy 86 (9), 1712–1722.
- Gebreslassie, B., Guillén-Gosálbez, G., Jiménez, L., Boer, D., 2010. A systematic tool for the minimization of the life cycle impact of solar assisted absorption cooling systems. Energy 39 (9), 3849–3862.
- Grossmann, I., Guillén-Gosálbez, G., 2010. Scope for the application of mathematical programming techniques in the synthesis and planning of sustainable processes. Comput. Chem. Eng. 34 (9), 1365–1376.
- Guillén, G., Badell, M., Puigjaner, L., 2007. A holistic framework for short-term supply chain management integrating production and corporate financial planning. International Journal of Production Economics 106 (1), 288–306.
- Guillén-Gosálbez, G., 2011. A novel mip-based objective reduction method for multi-objective optimization: application to environmental problems. Comput. Chem. Eng. 35 (8), 1469–1477.
- Guillén-Gosálbez, G., Caballero, J., Jiménez, L., 2008. Application of life cycle assessment to the structural optimization of process flowsheets. Ind. Eng. Chem. Res. 47 (3), 777–789.
- Guillén-Gosálbez, G., Grossmann, I., 2009. Optimal design and planning of sustainable chemical supply chains under uncertainty. AIChE J. 55 (1), 99–121.
- Guillén-Gosálbez, G., Grossmann, I., 2010. A global optimization strategy for the environmentally conscious design of chemical supply chains under uncertainty in the damage assessment model. Comput. Chem. Eng. 34 (1), 42–58.
- Guillén-Gosálbez, G., Mele, F., Grossmann, I., 2010. A bi-criterion optimization approach for the design and planning of hydrogen supply chains for vehicle use. AIChE J. 56 (3), 650–667.
- Guinée, J., Gorree, M., Heijungs, R., Huppes, G., Kleijn, R., de Konig, A., 1992. Environmental Life Cycle Assessment of Products: Background and Guide. Multi-Copy, Leiden.
- Gutiérrez, E., Lozano, S., Moreira, T., Feijoo, G., 2010. Assessing relationships among life-cycle environmental impacts with dimension reduction techniques. J. Environ. Manage. 91, 1002–1011.
- Gutiérrez, E., Lozano, S., Adenso-Díaz, B., 2010. Dimensionality reduction and visualization of the environmental impacts of domestic appliances. J. Ind. Ecol. 14 (6), 878–889.
- Hugo, A., Pistikopoulos, E., 2003. Environmentally conscious planning and design of supply chain networks. In: Chen, B., Westerberg, A.W. (Eds.), Process Systems Engineering 2003, Elsevier, Amsterdam, pp. 214–219.

- Hugo, A., Pistikopoulos, E., 2005. Environmentally conscious long-range planning and design of supply chain networks. *J. Cleaner Prod.* 13 (15), 1428–1448.
- Hugo, A., Rutter, P., Pistikopoulos, S., Amorelli, A., Zoia, G., 2005. Hydrogen infrastructure strategic planning using multi-objective optimization. *J. Hydrogen Energy* 30 (15), 1523–1534.
- Il Lim, Y., Floquet, P., Joulia, X., Kim, S., 1999. Multiobjective optimization in terms of economics and potential environment impact for process design and analysis in a chemical process simulator. *Ind. Eng. Chem. Res.* 38 (12), 4729–4741.
- Johnson, R., Wichner, D., 1998. *Applied Multivariate Statistical Analysis*, fourth ed. Prentice-Hall, Englewood Cliffs, NJ.
- Jolliffe, I., 2002. *Principal component analysis*, second ed. Springer Series in Statistics Springer, New York.
- Laínez, J., Guillén-Gosálbez, G., Badell, M., Espuña, A., Puigjaner, L., 2007. Enhancing corporate value in the optimal design of chemical supply chains. *Ind. Eng. Chem. Res.* 46 (23), 7739–7757.
- Mele, F., Espuna, A., Puigjaner, L., 2005. Environmental impact considerations into supply chain management based on life-cycle assessment. In: *Innovation by Life Cycle Management 2003*. LCM International Conference.
- Papandreou, V., Shang, Z., 2008. A multi-criteria optimisation approach for the design of sustainable utility systems. *Comput. Chem. Eng.* 32 (7), 1589–1602.
- Pistikopoulos, E., Stefanis, S., 1998. Optimal solvent design for environmental impact minimization. *Comput. Chem. Eng.* 22 (6), 717–733.
- PRÉ-Consultants, 2000. *The eco-indicator 99, a damage oriented method for life cycle impact assessment. methodology report and manual for designers*. Technical Report. Amersfoort, The Netherlands: PRÉ Consultants.
- Puigjaner, L., Guillén-Gosálbez, G., 2008. Towards an integrated framework for supply chain management in the batch chemical process industry. *Comput. Chem. Eng.* 32 (4–5), 650–670.
- Shapiro, J., 2001. *Modeling the Supply Chain*. MIT Press, Cambridge, MA.
- Stefanis, S., Livingston, A., Pistikopoulos, E., 1995. Minimizing the environmental impact of process plants: a process systems methodology. *Comput. Chem. Eng.* 19 (Suppl. 1), 39–44.
- Stefanis, S., Livingston, A., Pistikopoulos, E., 1997. Environmental impact considerations in the optimal design and scheduling of batch processes. *Comput. Chem. Eng.* 21, 1073–1094.
- Strang, G., 2009. *Introduction to Linear Algebra*, fourth ed. Springer Series in Statistics Springer, New York.
- Zhou, Z., Cheng, S., Hua, B., 2000. Supply chain optimization of continuous process industries with sustainability considerations. *Comput. Chem. Eng.* 24 (2–7), 1151–1158.
- Zitzler, E., Thiele, L., Laumanns, M., Fonseca, C., DaFonseca, V., 2003. Performance assessment of multiobjective optimizers: an analysis and review. *IEEE Trans. Evol. Comput.* 7 (2), 117–132.

# Appendices

## 1. List of publications

### 1.1. Research articles

Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: Identifying the preferred subset of enzymatic profiles in nonlinear kinetic metabolic models via multiobjective global optimization and Pareto filters. Accepted for publication in *PLoS ONE*.

Pozo, C., Ruíz-Femenia, R., Caballero, J., Guillén-Gosálbez, G., Jiménez, L.: On the use of Principal Component Analysis for reducing the number of environmental objectives in multi-objective optimization: Application to the design of chemical supply chains. *Chemical Engineering Science* 2012, 69(1), 146-158.

Miró, A., Pozo, C., Guillén-Gosálbez, G., Egea, J.A., Jiménez, L.: Deterministic global optimization algorithm based on outer approximation for the parameter estimation of nonlinear dynamic biological systems. *BMC Bioinformatics*, DOI: 10.1186/1471-2105-13-90.

Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: Outer approximation-based algorithm for biotechnology studies in systems biology. *Computers and Chemical Engineering* 2010, 34(10), 1719-1730.

Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: A spatial branch-and-bound framework for the global optimization of kinetic models of metabolic networks. *Industrial and Engineering Chemistry Research* 2011, 50(9), 5225-5238

Sorribas, A., Pozo, C., Vilaprinyo, E., Guillén-Gosálbez, G., Jiménez, L., Alves, R.: Optimization and evolution in metabolic pathways: Global optimization techniques in Generalized Mass Action models. *Journal of Biotechnology* 2010, 149(3), 141-153.

Pozo, C., Marín-Sanguino, A., Alves, R., Guillén-Gosálbez, G., Jiménez, L., Sorribas, A.: Steady-state global optimization of metabolic non-linear dynamic models through recasting into power-law canonical models. *BMC Systems Biology* 2011, 5, art. no. 137.

### 1.2. Book chapters

Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: Automatic selection of the most promising enzymatic modulations for metabolic engineering: a multi-objective optimization approach. 22nd European Symposium on *Computer Aided Chemical Engineering* 2012, 30. Elsevier, B.V. ISBN: 978-0-444-59431-0.

Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: Deterministic global optimization of kinetic models of metabolic networks: Outer approximation vs. Spatial branch and bound. 21st European Symposium on *Computer Aided Chemical Engineering*, 2011, 29, 583-586. Elsevier, B.V. ISBN: 978-0-444-53895-6.

Guillén-Gosálbez, G., Pozo, C., Jiménez, L., Sorribas, A.: An Outer Approximation Algorithm for the Global Optimization of Regulated Metabolic Systems. *10th International Symposium on Process Systems Engineering: Part A*. 2009, 1707-1712. Elsevier, B.V. ISBN: 978-0-444-53435-4.

Guillén-Gosálbez, G., Pozo, C., Jiménez, L., Sorribas, A.: A global optimization strategy to identify quantitative design principles for gene expression in yeast adaptation to heat shock. 19th European Symposium on *Computer Aided Chemical Engineering* 2009, 28, 1045-1050. Elsevier, B.V. ISBN: 978-0-444-53433-0.

## **2. Congress contributions**

### **2.1. Keynotes**

Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: Multi-Objective global optimization framework for non-linear kinetic models of metabolic networks. *ANQUE International Congress of Chemical Engineering 2012, Sevilla (SPAIN)*, 2012.

Sorribas, A., Guillén-Gosálbez, G., Pozo, C., Jiménez, L., Marín-Sanguino, A., Vilaprinyo, E., Alves, R.: Optimization and evolution of metabolic processes: is there life beyond flux balance analysis? *The XII International Congress on Molecular Systems Biology, Lleida (SPAIN)*, 2011.

### **2.2. Oral presentations**

Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: Deterministic global optimization of kinetic models of metabolic networks. *European Symposium on Computer Aided Process Engineering - 21, Chaldiki (GREECE)*, 2011.

Sabio, N., Pozo, C., Guillén-Gosálbez, G., Jiménez, L., Karuppiah, R., Vasudevan, V., Sawaya, N., Farrell, J.T.: Systematic Methods for the Elimination of Redundant Life Cycle Assessment Metrics in the Multi-Objective Optimization of Industrial Processes. *American Institute of Chemical Engineering 2011 Annual Meeting, Minneapolis (USA)*, 2011.

Pozo, C., Guillén-Gosálbez, G., Jiménez, L., Sorribas, A.: Global Optimization of Detailed Non-linear Kinetic Models of Metabolic Networks Through Recasting into Canonical Power-Law Representations. *American Institute of Chemical Engineering 2011 Annual Meeting, Minneapolis (USA)*, 2011.

Pozo, C., Guillén-Gosálbez, G., Jiménez, L., Sorribas, A.: Computational Strategies for the Global Optimization of Kinetic Models of Metabolic Networks. *American Institute of Chemical Engineering 2010 Annual Meeting, Salt Lake City (USA)*, 2010.

Sabio, N., Pozo, C., Guillén-Gosálbez, G., Jiménez, L., Vasudevan, V., Karuppiah, R., Sawaya, N., Farrell, J.T.: Improving the Environmental and Economic Performance of Industrial Processes Using a Multi-Objective Optimization Framework. *American Institute of Chemical Engineering 2010 Annual Meeting, Salt Lake City (USA)*, 2010.

Guillén-Gosálbez, G., Pozo, C., Jiménez, L., Sorribas, A.: An outer approximation algorithm for the global optimization of regulated metabolic systems. *10th International Symposium on Process Engineering, Salvador de Bahia (BRAZIL)*, 2009.

Pozo, C., Guillén-Gosálbez, G., Jiménez, L., Sorribas, A.: A Novel Outer Approximation Algorithm for the Global Optimization of Metabolic Networks. *American Institute of Chemical Engineering 2009 Annual Meeting, Nashville (USA)*, 2009.

### 2.3. Posters

Pozo, C., Guillén-Gosálbez, G., Sorribas, A., Jiménez, L.: Automatic selection of the most promising enzymatic modulations for metabolic engineering: a multi-objective optimization approach. *European Symposium on Computer Aided Process Engineering – 22, London (UK)*, 2012.

Pozo, C., Guillén-Gosálbez, G., Jiménez, L. , Sorribas, A.: Multi-Objective Global Optimization for Metabolic Engineering. *American Institute of Chemical Engineering 2012 Annual Meeting, Pittsburgh (USA)*, 2012.

Pozo, C., Marín-Sanguino, A., Alves, R., Guillén-Gosálbez, G., Jiménez, L. , Sorribas, A.: Steady-state global optimization of metabolic non-linear dynamic models through recasting into power-law canonical models. *The XII International Congress on Molecular Systems Biology, Lleida (SPAIN)*, 2011.

Pozo, C., Marín-Sanguino, A., Guillén-Gosálbez, G., Jiménez, L. , Sorribas, A.: Global Optimization of Metabolic Kinetic Non-Linear Models through Recasting. *12th Mediterranean Congress of Chemical Engineering, Barcelona (SPAIN)*, 2011.

Guillén-Gosálbez, G., Pozo, C., Jiménez, L. , Sorribas, A.: A global optimization strategy to identify quantitative design principles for gene expression in yeast adaptation to heat shock. *European Symposium on Computer Aided Process Engineering - 19, Crakow (POLAND)*, 2009.

Guillén-Gosálbez, G., Pozo, C., Jiménez, L. , Sorribas, A.: A systematic method for searching feasible enzyme activity patterns. *8th World Congress of Chemical Engineering, Montreal (CANADA)*, 2009.