

Document downloaded from:

<http://hdl.handle.net/10251/50289>

This paper must be cited as:

González Martínez, JM.; Abel Folch-Fortuny; Llaneras Estrada, F.; Tortajada Serra, M.; Picó Marco, JA.; Ferrer, A. (2014). Metabolic flux understanding of *Pichia pastoris* grown on heterogenous culture media. *Chemometrics and Intelligent Laboratory Systems*. 134:89-99. doi:10.1016/j.chemolab.2014.02.003.



The final publication is available at

<http://dx.doi.org/10.1016/j.chemolab.2014.02.003>

Copyright Elsevier

# Metabolic Flux Understanding of *Pichia pastoris* Grown on Heterogenous Culture Media

J. M. González-Martínez<sup>a,b,\*</sup>, A. Folch-Fortuny<sup>b</sup>, F. Llaneras<sup>c</sup>, M. Tortajada<sup>d</sup>, J. Picó<sup>e</sup>, A. Ferrer<sup>b</sup>

<sup>a</sup>Shell Global Solutions International B.V., Shell Technology Centre Amsterdam, PO Box 38000, 1030 BN Amsterdam, The Netherlands

<sup>b</sup>Departamento de Estadística e Investigación Operativa Aplicadas y Calidad, Universitat Politècnica de València, Camino de Vera s/n, Edificio 7A, 46022 Valencia, Spain

<sup>c</sup>Department of Electrical, Electronic and Automatic Engineering, Universitat de Girona, Campus de Montilivi, 17071 Girona, Spain

<sup>d</sup>Biopolis, S.L., Parc Científic Universitat de València, c/ Catedrático Agustín Escardino Benlloch 9, 46980 Paterna, Spain

<sup>e</sup>Institut Universitari d' Automàtica i Informàtica Industrial, Universitat Politècnica de València, Camino de Vera s/n, Edificio 5C, 46022 Valencia, Spain

---

## Abstract

Within the emergent field of Systems Biology, mathematical models obtained from physical-chemical laws (the so-called first principles-based models) of microbial systems are employed to discern the principles that govern cellular behaviour and achieve a predictive understanding of cellular functions. The reliance on this biochemical knowledge has the drawback that some of the assumptions (specific kinetics of the reaction system, unknown dynamics and values of the model parameters) may not be valid for all the metabolic possible states of the network. In this uncertainty context, the combined use of fundamental knowledge and data measured in the fermentation that describe the behavior of the microorganism in the manufacturing process is paramount to overcome this problem. In this paper, a grey modelling approach is presented combining data-driven and first principles information at different scales, developed for *Pichia pastoris* cultures grown on different carbon sources. This approach will allow us to relate patterns of recombinant protein production to intracellular metabolic states and correlate intra and extracellular reactions in order to understand how the internal state of the cells determines the observed behaviour in *P. pastoris* cultivations.

**Keywords:** *Metabolic network, Possibilistic consistency analysis, Monte Carlo sampling, Principal Component Analysis, Missing-data methods for Exploratory Data Analysis*

---

## 1. Introduction

Currently, biotechnological industries are devoted to the production of economically important enzymes and proteins, generally using genetically modified microorganisms. The main goal of these industries is to maximize protein yields and productivity. The production of these high-added value products is governed by highly correlated factors that require a multidisciplinary approach to process optimization (biochemistry, molecular biology, process engineering, biotechnology, *etc*).

The development of accurate monitoring schemes to control the manufacturing process becomes a challenging task due to the scarcity of measurements and the high complexity of the biochemical synthesis process. Only few process variables can be measured in industrial microbial fermentations, such as pH, temperature, and oxygen consumption. Others, such as substrate consumption, can be inferred depending on the operational strategy. In this kind of processes, measurements corresponding to biological process variables (such as intracellular specific reaction rates) are promising to achieve a more accurate process

---

\*Corresponding author.

Email address: jgonmar@gmail.com (J. M. González-Martínez)

28 control. In order to develop novel monitoring schemes, the study of different cellular behaviours is crucial  
29 for the biotechnological production of high-added value biochemicals. For this purpose, the modelling of the  
30 available data is needed to know which key variables control the main metabolic pathways and, possibly,  
31 their regulation mechanisms.

32 First principles-based models of microbial systems can be developed to describe the principles that gov-  
33 ern cellular behaviour and achieve a predictive understanding of cellular functions [1]. Typically, networks  
34 of biochemical reactions are used to approach an organism microbial metabolism and growth [2, 3]. These  
35 networks are modelled assuming that certain constraints operate at steady-state, such as environmental con-  
36 straints [4], regulatory constraints [5, 6], gene expression data [7], mass balances or reactions irreversibilities  
37 [8] (the so-called *constraint-based perspective*) [9, 10, 11]. The imposed constraints define a solution space  
38 that encloses all the possible states of the network (*i.e.* flux distributions through the reactions). The  
39 development of this type of models based solely on fundamental or knowledge information has the drawback  
40 that the unknown part of the process is not represented as well as some of the underlying assumptions  
41 (*e.g.* specific kinetics of the reaction system, unknown dynamics, values of the model parameters, objec-  
42 tive functions) may not be valid for all the metabolic possible states of the network [12, 13]. To address  
43 this problem, hybrid models that combine knowledge-based models, which fit the theoretical behavior, and  
44 empirical models, which fit any remaining systematic variation, can be used [14].

45 In the context of grey modelling, there are different approaches to decompose the data into the three  
46 types of variation (known causes, unknown causes and residuals) [15], which be roughly classified into three  
47 categories. The first category are the models based on known constraints. There exist general frameworks  
48 that enable to impose very specific constraints on each type of information, *e.g.* observed experimental  
49 information [16] or transformations on the original variables [17]. These methods are based on the projection  
50 of a data matrix, followed by multivariate model decomposition. Principal Component Analysis (PCA) [18]  
51 is one of the most applied multivariate statistical projection methods to reveal the internal structure of the  
52 cell. This analysis is commonly preceded by a Monte Carlo sampling in order to produce a data set of  
53 possible states or feasible solutions from which the PCA elucidates the meaningful principal components  
54 (PCs) [19, 20, 21]. PCA has also been compared to other multivariate techniques, such as Multivariate  
55 Linear Regression [21] and Parallel Factor Analysis (PARAFAC) [22, 23] in the field of Systems Biology.  
56 Partial Least Squares regression (PLS) [24] has been applied directly [25, 26] and combined with Hierarchical  
57 Clustering (HC-PLSR) [27] to deal with situations where the input-output relations (*e.g.* the effect of the  
58 substrates consumption of the cell or the environmental conditions in the production of a particular protein)  
59 are highly nonlinear or non-monotone. Recently, Grey Component Analysis (GCA) has been proposed using  
60 a cost function to maximise the interpretability of the solutions by forcing the decomposition towards the  
61 direction of the prior information - a chemically or biologically meaningful solution - [28]. A second strategy  
62 is formed by methods based on introducing *a priori* knowledge by means of mathematical relations that  
63 describe the system behaviour or dynamics. The starting point is some specific structure based on first  
64 principles mathematical relations, where some functions must be estimated. Different tools can be used to  
65 calculate these functions, such as artificial Neural Networks (NNs) [29] or Kalman filters [30, 31]. Finally,  
66 a third category are the methods based on incorporating the fundamental knowledge through constraints  
67 on the modelling algorithms. For instance, some model parameters can be forced to have values within  
68 certain regions in the parameters space [32]. Projection to Latent Pathways (PLP) [33] has been recently  
69 formulated as a modification of the PLS regression algorithm by using the concept of Elementary Modes (*i.e.*  
70 thermodynamically feasible pathways through the metabolic network). This method is devoted to obtain a  
71 more biologically explanatory set of latent variables (LVs) relating the observed behaviour of the cell and  
72 its initial conditions.

73 The complexity of data available from microbial systems requires the design of sophisticated grey mod-  
74 els that combine data-driven and knowledge-based information at different scales for biochemical process  
75 understanding. The main goal of this paper is to use this hybrid framework to analyse the behaviour of  
76 the methylophilic yeast *P. pastoris* [3], as a first step to analysing which conditions and through which  
77 reactions the cell achieves an optimal state for our interests. Several scenarios corresponding to different  
78 chemostat runs are collected from the literature [34, 35, 36, 37, 38, 39, 40, 41, 42] with the aim of starting  
79 the analysis with a rich data set of different culture conditions. A recently developed adaptation [3, 43]

80 of the possibilistic theory [44] is applied in order to check the consistency between model and data. For  
81 the completion of the unmeasured data, a Monte Carlo sampling method is applied to produce feasible flux  
82 solutions for the microbial system under study. At this point, a PCA is performed to obtain a reduced  
83 number of PCs explaining most of the variance of the collected and sampled data. Finally, the Missing-data  
84 method for Exploratory Data Analysis (MEDA) [45] is applied to obtain a better interpretation of the PCs  
85 derived from the PCA model.

86 This paper is organized as follows. Section 2 presents the metabolic network reconstruction of the yeast  
87 *P. pastoris* and the different scenarios used in the study. Section 3 describes the grey modelling approach  
88 proposed in detail. This procedure is applied to the available data from *P. pastoris* in Section 4. Finally,  
89 some conclusions on the grey modelling approach presented in this paper and how it may be applied to  
90 improve the understanding of microbial cultures are drawn in Section 5.

## 91 2. Materials

### 92 *Metabolic network reconstruction*

93 The methylotrophic *P. pastoris* has become one of the most widely used yeasts for heterologous protein  
94 production since its development, in the early 1970s [46]. This system is of particular industrial interest due  
95 to its powerful and tightly regulated methanol-inducible alcohol oxidase 1 promoter (pAOX1), its capacity  
96 for foreign protein secretion, its ability to perform post-translational modifications (including glycosylation  
97 and disulfide bond formation) and its capability to grow on defined media at high cell densities [47, 48].

98 The constraint-based model, whose corresponding metabolic network is shown in Figure 1, has been used  
99 throughout this work. The model is a simplified representation of the whole metabolism of the yeast *P.*  
100 *pastoris*, meaning that only a reduced number of biochemical reactions has been included (45), from the  
101 larger amount available from genomic information (more than 1200). The reactions were selected on the  
102 basis of previous models found in literature, as lumped equivalents of more complex pathways. This model  
103 was previously validated by the authors for this organism and the experimental conditions studied [3, 49]  
104 and is the only one used in the referred experiments throughout this work. The model represents the most  
105 significant features of *P. pastoris* metabolism, including the main catabolic pathways of the yeast, such as  
106 glycolysis, the citric acid cycle, glycerol and methanol oxidation and fermentative pathways [49]. Anabolism  
107 is introduced through the pentose phosphate pathway and a general lumped biomass equation according to  
108 which growth is assumed to depend exclusively on key biochemical precursors. Branch-point metabolites,  
109 such as NADH, NADPH, AcCoA, oxalacetate and pyruvate, are considered in compartmentalized cytosolic  
110 and mitochondrial pools [34].

### 111 *P. pastoris* experimental data set

112 In this work, experimental data from several fermentation runs with different *P. pastoris* strains have  
113 been taken from the available literature, building the different scenarios considered for the subsequent  
114 statistical analysis. For the sake of visualization, the 40 scenarios under study have been grouped attending  
115 to the experimental substrates (*i.e.* glucose, glycerol, methanol, and glycerol and methanol mixtures) (see  
116 Figure 2). Scenario *A1* is taken from the strain expressing the Fab fragment of the human anti-HIV  
117 antibody 3H6 [34]. Scenarios from *B1* to *B7*, and *C1* and *C2*, are from a strain expressing a *Rhizopus*  
118 *oryzae* lipase (ROL) [35, 36]. Scenarios from *D1* to *D10* come from a *P. pastoris* strain expressing and  
119 secreting recombinant avidin [37]. Scenario *E1* has been obtained from a macrokinetic model for *P. pastoris*  
120 expressing recombinant human serum albumin (HSA) [38]. Scenarios from *F1* to *F7* are from a *P. pastoris*  
121 strain genetically modified to produce sea raven antifreeze protein [39]. Scenarios from *G1* to *G10* are  
122 obtained from a *P. pastoris* strain producing recombinant human chymotrypsinogen B [40]. Scenario *H1*  
123 has been obtained from the continuous fermentation of a *P. pastoris* strain for the extracellular production  
124 of a recombinant ovine interferon protein [41]. Finally, scenario *I1* comes from the expression of recombinant  
125 chitinase with a genetically modified *P. pastoris* strain [42]. The data for all these scenarios are detailed in  
126 Figure 2.

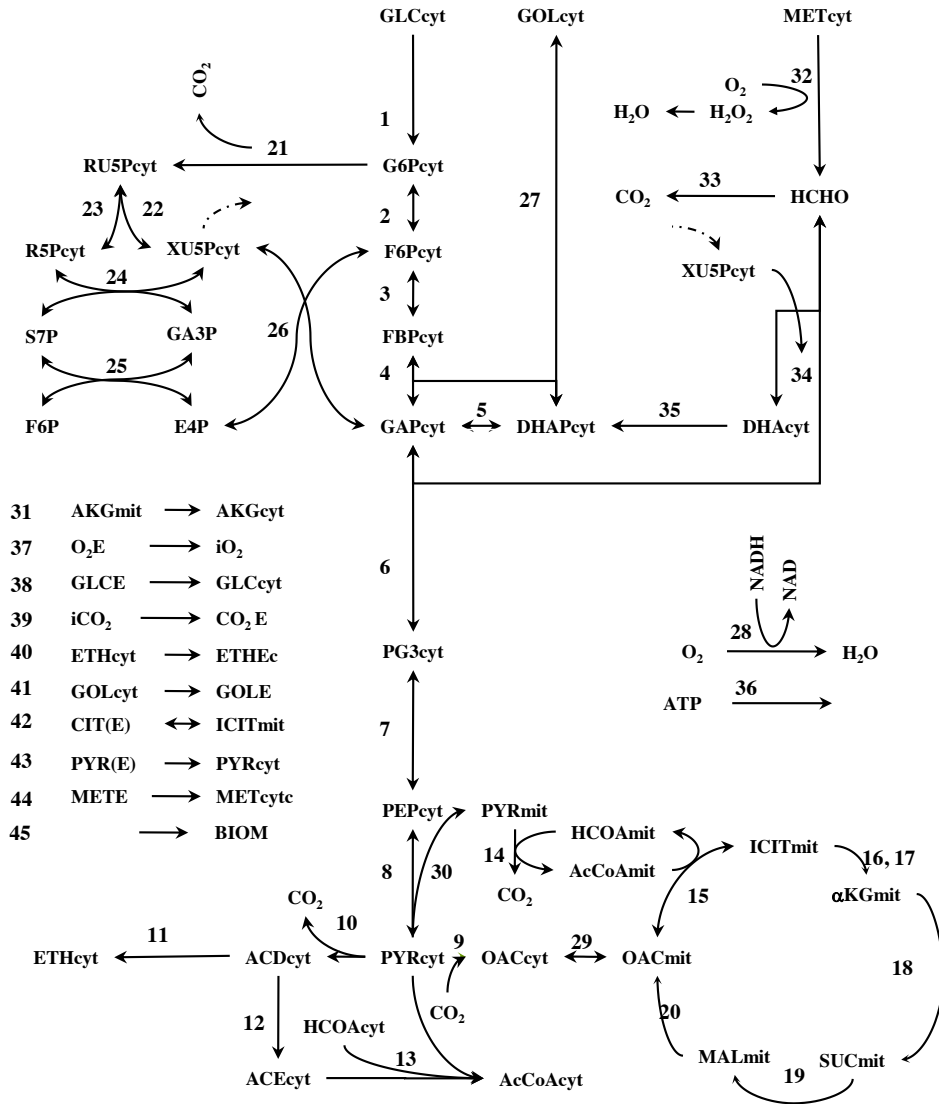


Figure 1: Summarized representation of the metabolic network of *P. pastoris*, representing the central carbon metabolism of the yeast during growth on glucose, glycerol and methanol. For the purpose of clarity, the biomass equation is not represented in the figure. Please refer to the stoichiometric matrix for details about each reaction and the involved metabolites.

127 At this point, there is a paramount comment that is in due. Batch effects, which are defined as systematic  
 128 non-biological variation between groups of samples (or batches) due to experimental artifacts [50, 51, 52, 53],  
 129 can be present in data collected from different cultures. In case that replicates of the same scenario are  
 130 collected (*i.e.* same strain and same quantities of initial substrates) and the presence of batch effects is  
 131 statistically confirmed, this artificial variation must be removed. Otherwise, the bias introduced by the  
 132 non-biological nature of this kind of effects may confound true biological differences [52], affecting the  
 133 results of statistical analysis. In this study, the scenarios within a single strain of *P. pastoris* have different  
 134 initial substrate quantities (see Figure 2). Hence, the variation observed across scenarios can be due to  
 135 these different initial conditions, which were applied with the aim of obtaining different flux values. This

136 fact jointly with the scarcity of information about the experimentation conditions disable the possibility to  
 137 straightforwardly confirm actual batch effects in data.

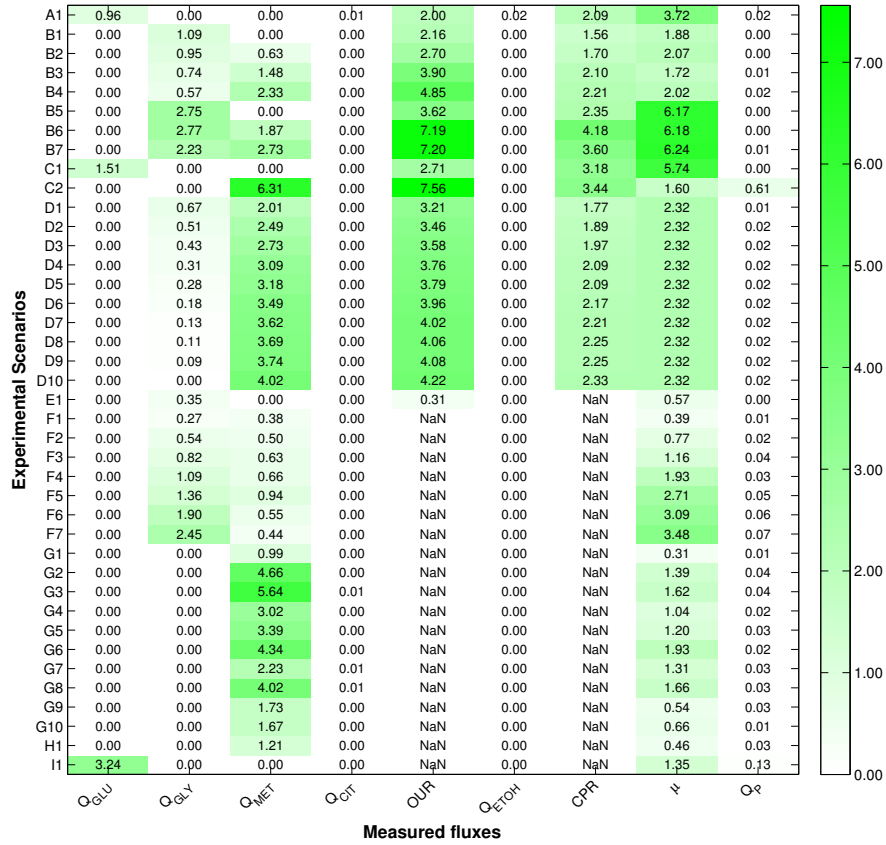


Figure 2: Set of 40 experimental scenarios corresponding to *P. pastoris* chemostat cultures grown on glucose, glycerol and methanol mixtures. For each scenario, the values of measured fluxes belonging to substrate and product specific consumption and production are shown. The substrates are glucose ( $Q_{GLU}$ ), glycerol ( $Q_{GLYC}$ ), methanol ( $Q_{MET}$ ), citrate ( $Q_{CIT}$ ) and oxygen ( $OUR$ ). The products are ethanol ( $Q_{ETOH}$ ), carbon dioxide ( $CPR$ ), biomass ( $\mu$ ) and protein ( $Q_P$ ). Note that NaN values stand for missing measured external fluxes.

### 138 3. Methods

139 The grey modelling approach proposed is composed of several steps (see Figure 3). Firstly, the constraint-  
 140 based model is built by transforming the network in a mathematical form (the stoichiometric matrix  $\mathbf{S}$  and  
 141 the flux irreversibilities are detailed in the Supplementary Information). At the same time, different ex-  
 142 perimental fermentation scenarios are collected from the literature. The combination of these two steps  
 143 represents the novel grey modelling approach, detailed in the previous sections. Then, a Possibilistic consis-  
 144 tency analysis is performed to elucidate which scenarios are not consistent with the model. On the consistent  
 145 scenarios, a Monte Carlo sampling is applied to obtain a hundred different feasible solutions satisfying the  
 146 proposed model. With the feasible flux solutions matrix, a Principal Component Analysis (PCA) is per-  
 147 formed with the aim of getting an insight of the metabolic structure of the yeast. All the outliers are detected

148 by using two Shewhart-type control charts based on the Hotelling- $T^2$  and Square Prediction Error (SPE)  
 149 statistics and, later on, root causes are diagnosed. Once the outliers are isolated, PCA is computed again.  
 150 This procedure is repeated until the percentage of outliers is consistent with the confidence limits (99%  
 151 confidence level)). Finally, the Missing-data method for Exploratory Data Analysis (MEDA) is applied in  
 152 order to attain more informative components. The theory behind these steps are described in the following  
 153 subsections.

#### 154 *Stoichiometric modelling*

155 To build a constraint-based model, the stoichiometric information embedded in the metabolic network  
 156 (*i.e.* metabolites or cofactors involved in each reaction) must be arranged into a  $I \times J$  matrix  $\mathbf{S}$  (the so-called  
 157 stoichiometric matrix). Rows of this matrix represent the  $I$  metabolites, columns the  $J$  metabolic reactions  
 158 and each element  $(i, j)$  the stoichiometric coefficient  $S_{i,j}$  of the  $i$ th metabolite in the  $j$ th reaction. A value  
 159 of  $S_{i,j} = -1$  indicates that the  $i$ th metabolite is consumed by the  $j$ th reaction. In contrast, a  $S_{i,j} = 1$   
 160 indicates the  $i$ th metabolite is produced by the  $j$ th reaction. Finally, a value of  $S_{i,j} = 0$  stands for the  $i$ th  
 161 metabolite is not involved in the  $j$ th reaction.

162 The stoichiometric matrix is used, in combination with the flux vector  $\mathbf{v} = (v_1, \dots, v_J)$ , the intracellular  
 163 metabolites concentration  $\mathbf{c} = (c_1, \dots, c_I)$  and the specific growth rate of the cell  $\mu$ , to represent the mass  
 164 balances through the metabolic network. The ordinary differential equation describing this process is as  
 165 follows:

$$\frac{d\mathbf{c}}{dt} = \mathbf{S} \cdot \mathbf{v} - \mu \cdot \mathbf{c} \quad (1)$$

166 This equation is called the dynamic mass balance equation, and describes the evolution of the concentra-  
 167 tion of each metabolite over time [11]. In stoichiometric modelling, the dynamic intracellular behaviour is  
 168 disregarded on the basis assumption of pseudosteady state for the internal metabolites [8]. This assumption  
 169 is supported by the observation that intracellular dynamics are much faster than extracellular dynamics.  
 170 Therefore, it is sensible to assume that these compounds reach the steady state instantaneously and, hence,  
 171 its transient behavior can be omitted. In addition, the dilution term  $\mu \cdot \mathbf{c}$  is also discarded because it is gen-  
 172 erally much smaller than the fluxes affecting the same metabolite. Under these considerations, the general  
 173 equation can be expressed as:

$$\mathbf{S} \cdot \mathbf{v} = \mathbf{0} \quad (2)$$

174 This equation constrains the  $J$ -dimensional space of feasible solutions. An extra constraint is added,  
 175 assuming that some of the fluxes of the metabolic network flow only in one direction:

$$\mathbf{D} \cdot \mathbf{v} \geq \mathbf{0} \quad (3)$$

176 where  $\mathbf{D}$  is a  $J \times J$  diagonal matrix with binary values: 1 for the irreversible fluxes and 0 for the reversible  
 177 ones.

178 Finally, a maximum value for each flux value is computed:

$$v_j \leq v_{j,max} \quad \forall j \in 1, \dots, J \quad (4)$$

179 The combination of the constraints imposed by Equations 2-4 define a space (a bounded convex cone)  
 180 of feasible steady-state flux distributions: only flux vectors that fulfill Equation 2-4 are considered valid  
 181 cellular states. In this way, Equations 2-4 define our model of *Pichia pastoris*, following a constraint-based  
 182 modelling approach [9, 10, 11].

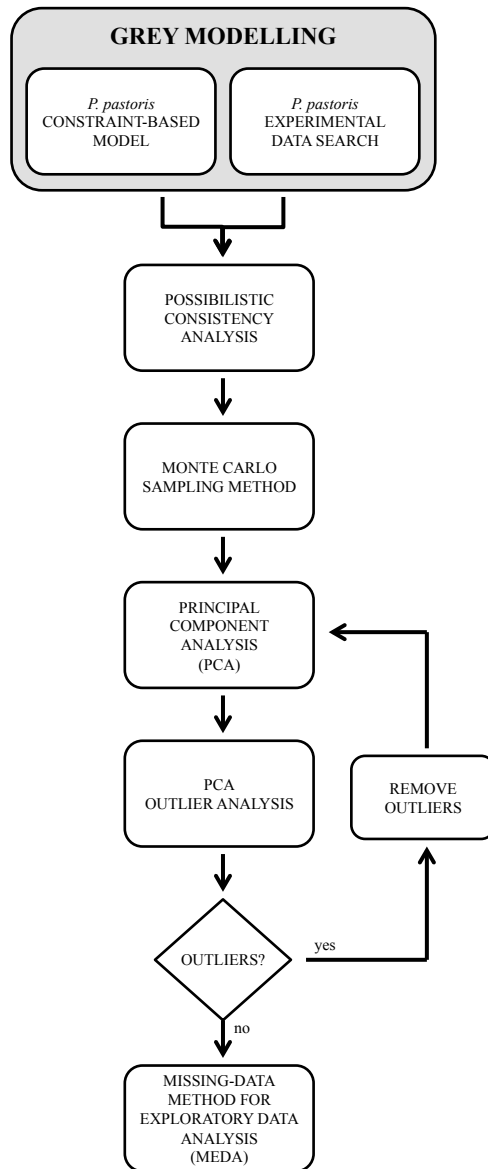


Figure 3: Flow diagram of the grey modelling approach of *P. pastoris*.



184 The simplest consistency analysis could be performed by checking that the flux states shown by cells  
 185 fulfill the constraints imposed by the model (see Equations 2-4) [3]. However, this simple approach would be  
 186 impractical because measurements are imprecise and do not exactly satisfy the constraints. Such difficulty  
 187 is overcome by taking into account uncertainty as follows:

$$w = v_m + e \quad (5)$$

188 where  $e$  represents the deviation error between an actual flux  $v_m$  and its measured value  $w$ .

189 The consistency analysis can be also formulated as a possibilistic constraint satisfaction problem [43].  
 190 The basic idea is that a flux vector fulfilling Equations 2 and 3, and compatible with the measurements  
 191 will be considered as “possible”, otherwise as “impossible”. This can be refined to cope with measurements  
 192 errors by introducing the notion of “degree of possibility” [44].

193 This degree of possibility provides an indication of the consistency between the model and the measure-  
 194 ments. A possibility equal to one must be interpreted as complete agreement between the model and the  
 195 original measurements. Lower values of possibility imply that certain error in the measurements is needed  
 196 to find a flux vector fulfilling the model constraints. For further details readers are referred to [3, 43].

197 The main formulation of Possibilistic consistency analysis is summarized in this section.

198 **Model and measurements constraints.** Firstly we consider the constraints conforming the model (Equa-  
 199 tion 2-4). Then measures of (some) extracellular fluxes are incorporated as additional constraints (Equation  
 200 6):

$$\begin{cases} \mathbf{v}_m = \mathbf{w} + \boldsymbol{\varepsilon}_1 - \boldsymbol{\mu}_1 + \boldsymbol{\varepsilon}_2 - \boldsymbol{\mu}_2 \\ \boldsymbol{\varepsilon}_1, \boldsymbol{\mu}_1, \boldsymbol{\varepsilon}_2, \boldsymbol{\mu}_2 \geq \mathbf{0} \\ \boldsymbol{\varepsilon}_2 \leq \boldsymbol{\varepsilon}_{2,max} \\ \boldsymbol{\mu}_2 \leq \boldsymbol{\mu}_{2,max} \end{cases} \quad (6)$$

201 where vector  $\mathbf{v}_m$  represents the actual metabolite concentrations and  $\mathbf{w}$  the measured values, which differ  
 202 due to errors and imprecision (uncertainty). This uncertainty is represented by the vectors of slack variables  
 203  $\boldsymbol{\varepsilon}$ 's and  $\boldsymbol{\mu}$ 's.

204 **Possibility.** Let us denote each candidate solution of Equation 6 as  $\delta = \{\mathbf{v}, \mathbf{w}, \boldsymbol{\varepsilon}, \boldsymbol{\mu}\}$  in  $\Delta$ . The basic  
 205 building block of possibility theory is a user-defined possibility distribution  $\pi(\delta) : \Delta \rightarrow [0, 1]$ . This function  
 206 defines the possibility of each solution  $\delta$  in  $\Delta$ , ranging between impossible ( $\pi = 0$ ) and fully possible ( $\pi = 1$ ).  
 207 Among different possible choices, a simple -yet sensible- way to define possibility is using a linear cost index  
 208 such as Equation :

$$J(\delta) = \boldsymbol{\alpha}^T \cdot \boldsymbol{\varepsilon}_1 + \boldsymbol{\beta}^T \boldsymbol{\mu}_1 \quad (7)$$

209 and define the possibility of each solution  $\delta$  as follows:

$$\pi(\delta) = e^{-J(\delta)}, \quad \delta \in \Delta \quad (8)$$

210 where  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  are row vectors of user-defined, sensor accuracy coefficients.

211 The interpretation of Equations 6-8 may be:  $\mathbf{v} = \mathbf{w}$  is fully possible; the more  $\mathbf{v}$  and  $\mathbf{w}$  differ, the less  
 212 possible such situation is

213 **Representing uncertainty.** Two pairs of vectors of slack variables have been chosen to represent the  
 214 uncertainty of each measurement:  $\boldsymbol{\varepsilon}_2$  and  $\boldsymbol{\mu}_2$  define an interval of fully possible values, and  $\boldsymbol{\varepsilon}_1$  and  $\boldsymbol{\mu}_1$   
 215 penalise values out of it (with weights  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$ ). This is achieved choosing two vectors of bounds. Hence,  
 216 in all computations the uncertainty of each measurement has been represented as follows:

- 217 (i) Full possibility ( $\pi = 1$ ) is assigned to values with less than  $\pm 5\%$  of deviation.
- 218 (ii) Larger deviations are penalised, so values with a deviation equal to  $\pm 20\%$  have a possibility of  $\pi = 0.1$ ,  
 219 and those with a deviation equal to  $\pm 10\%$  have a  $\pi \approx 0.5$ .

220 (iii) Uncertainty is considered as symmetric, and thus  $\alpha = \beta$ .

221 This is achieved choosing bounds  $\varepsilon_{2,max}$  and  $\mu_{2,max}$  and weights  $\alpha$  and  $\beta$  for each measurement: (i)  
 222 implies that  $\varepsilon_{2,max} = \mu_{2,max} = 0.05 \cdot \mathbf{w}$ , and (ii) defines  $\alpha$ , noticing that  $0.2 \cdot \mathbf{w} = \mu_{1,20\%} + \mu_{2,max}$ , then  
 223  $\alpha = -\log(0.1)/(0.2 - 0.05)/\mathbf{w}$ .

224 **Possibilistic consistency evaluation.** This method can be applied to evaluate the degree of consistency  
 225 between a given model and a set of experimental measurements. Notice that the most possible solution  
 226 of the constraint-satisfaction problem is the maximum possibility (minimum-cost) solution, which can be  
 227 obtained solving a linear programming problem (LP):

$$J^{min} = \min_{\varepsilon, \mu, \mathbf{v}} J \quad (9)$$

228 subject to Equations 2-4 and the experimental measurements. This solution has an associated degree of  
 229 possibility:

$$\pi^{mp} = e^{-J^{min}} \quad (10)$$

230 This value,  $\pi^{mp}$  in  $[0,1]$ , grades the consistency between model and measurements. A possibility equal  
 231 to one must be interpreted as complete agreement, while lower values imply that there is some error in the  
 232 measurements, the model or both, which severity depends on how the uncertainty has been defined (see  
 233 above). More details on Possibilistic consistency analysis are given in a previous work, where the model of  
 234 *P. pastoris* was validated [3, 49].

#### 235 *Monte Carlo sampling method*

236 Metabolic Flux Analysis was designed with the aim of obtaining the flux values of all reactions based on  
 237 the known fluxes, typically extracellular, which are easier to measure. Assuming the  $J_1$  measured fluxes of  
 238 the  $J = J_1 + J_2$  fluxes of the metabolic network, the  $J_2$  unmeasured fluxes can be derived from the general  
 239 equation of the stoichiometric modelling (see Equation 2):

$$\mathbf{S}_{J_1} \cdot \mathbf{v}_{J_1} = -\mathbf{S}_{J_2} \cdot \mathbf{v}_{J_2} \quad (11)$$

240 where  $\mathbf{S}_{J_1}$  and  $\mathbf{v}_{J_1}$  involves the measured fluxes (the external ones), and  $\mathbf{S}_{J_2}$  and  $\mathbf{v}_{J_2}$  are related to the  
 241 unmeasured fluxes (the internal ones). The problem with this formulation is that the number of internal  
 242 fluxes often remain high compared to the number of external fluxes. Thus, the system shown in Equation  
 243 11 is undetermined, *i.e.* there are different flux distributions compatible with the known flux values.

244 In this context, Monte Carlo sampling methods can be used to produce feasible flux distributions across  
 245 the cell [19, 20, 21, 54, 55]. This way, the available experimental data and the first-principles knowledge  
 246 captured by the model are coupled together, providing a new richer data-set amenable to further analysis  
 247 with a statistical multivariate projection method. To randomly generate possible values for the unmeasured  
 248 fluxes (internal fluxes) for each cultivation, stoichiometry (see Equation 2), irreversibility (see Equation 3)  
 249 and measured fluxes (see Equation 4 and experimental data on Figure 2) are taken into account.

250 In order to deal with experimental errors, external fluxes are allowed to vary within a defined range  
 251 of values centered on the original measured value. The upper (lower) bound of this range is the sum  
 252 (subtraction) of the measured value and the maximum value between 0.001 and the 10% of the measured  
 253 value:

$$(LB, UB)_j = (v_j - \max(0.001, 0.1 \times v_j), v_j + \max(0.001, 0.1 \times v_j)) \quad \forall j \in 1, \dots, J \quad (12)$$

254 where  $v_j$  is a measured flux, and  $LB$  and  $UB$  are the lower and upper bounds for the Monte Carlo sampling  
 255 method.

256 At this point, it is worth commenting that the feasible solutions for each scenario are obtained by  
 257 sampling within the *slice* of the cone defined by Equations 2-4 and the experimental data, *i.e.* the measured  
 258 fluxes reduce the feasible solution space from the initial cone, which is bounded only by the constraint-based  
 259 model, to the portion of it fulfilling these specific experimental measurements. The complete procedure can  
 260 be visualized in Figure 4.

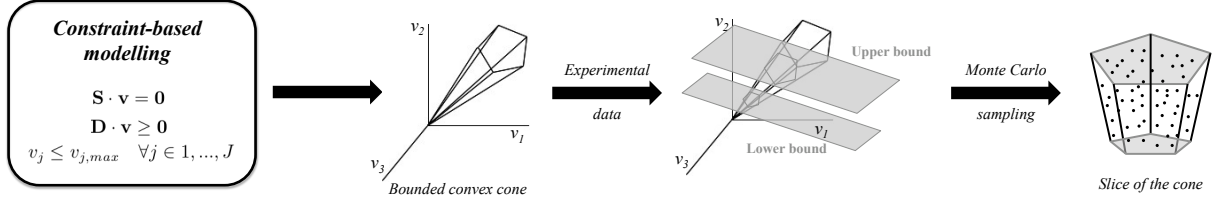


Figure 4: Constraint-based modelling and allowed flux states sampling. The convex cone is obtained by Equations 2-4. The experimental measurements constrain the cone through Equation 12. Finally, the sampling is performed on the resulting *slice* of the cone.

### 261 *Principal Component Analysis*

262 The aim of Principal Component Analysis (PCA) is to find the subspace in the space of the variables  
 263 where data mostly vary [56]. The original variables, commonly correlated, are linearly transformed into  
 264 a lower number of uncorrelated variables (the so-called principal components, PCs). PCA follows the  
 265 expression:

$$\mathbf{X} = \mathbf{T}_A \cdot \mathbf{P}_A^t + \mathbf{E}_A \quad (13)$$

266 where  $\mathbf{X}$  is a  $N \times M$  matrix of data,  $\mathbf{T}_A$  is the  $N \times A$  scores matrix containing the projection of the objects  
 267 in the  $A$  PCs subspace,  $\mathbf{P}_A$  is the  $M \times A$  loadings matrix containing the linear combination of the variables  
 268 represented in each of the PCs, and  $\mathbf{E}_A$  is the  $N \times M$  matrix of residuals.

269 As a previous step of PCA, the data matrix is autoscaled, i.e. each variable (flux) is centered and  
 270 divided by its standard deviation, making all variables have a variance equal to 1. In the present work,  
 271 since the components obtained by PCA are linear combinations of different fluxes, the more positive the  
 272 coefficient of a flux is, the more positive correlated is with this particular component, in the sense that the  
 273 flux is higher than its mean value in this component. As well, the more negative its coefficient is, the more  
 274 negative correlated is with this component, in the sense that the flux is lower than its mean value in this  
 275 component. In other words, fluxes with positive coefficients in a component are overused, and fluxes with  
 276 negative coefficients are underused.

### 277 *PCA outlier detection*

278 Square Prediction Error (SPE) and Hotelling- $T^2$  are two statistics widely used to detect outliers on a  
 279 given data. SPE is the orthogonal distance of a particular object to the  $A$ -dimensional subspace of latent  
 280 variables defined by PCA. It is expressed as:

$$SPE_n = \mathbf{e}_n^t \cdot \mathbf{e}_n \quad \forall n \in 1, \dots, N \quad (14)$$

281 where  $\mathbf{e}_n$  is the  $n$ th row of the residual matrix  $\mathbf{E} = \mathbf{X} - \mathbf{T}_A \mathbf{P}_A^t$ . By taking the eigenvalues of the covariance  
 282 matrix of the residual matrix ( $\lambda_{A+1}, \dots, \lambda_M$ ), the control limit of the SPE [57] is computed as follows:

$$SPE_\alpha = \theta_1 \left[ \frac{z_\alpha \sqrt{2\theta_2 h_0^2}}{\theta_1} + 1 + \frac{\theta_2 h_0 (h_0 - 1)}{\theta_1^2} \right]^{1/h_0} \quad (15)$$

283 where  $\theta_m = \sum_{j=A+1}^M (\lambda_j)^m$ ,  $h_0 = 1 - \frac{2\theta_1 \theta_3}{3\theta_2^2}$  and  $z_\alpha$  is the  $100 \times (1 - \alpha)$  percentile of a standard Normal  
 284 distribution.

285 Hotelling- $T^2$  is a statistic based on the Mahalanobis distance [58]. This statistic is used in multivariate  
 286 monitoring to compute the distance between one object and model's centre according to the covariance  
 287 structure [59]. When the data is centered (the mean of each column of  $\mathbf{X}$  is equal to zero), the distance  
 288 between an observation  $\mathbf{x}_n$  (a row from the data matrix  $\mathbf{X}$ ) and the centre of the original  $M$ -dimensional  
 289 variable space is:

$$\chi_n^2 = \mathbf{x}_n^t \cdot \boldsymbol{\Sigma}^{-1} \cdot \mathbf{x}_n \quad \forall n \in 1, \dots, N \quad (16)$$

290 where  $\boldsymbol{\Sigma}$  is the real covariance matrix of the original  $M$ -dimensional variable space, and  $\chi_n^2$  follows a  $\chi^2$   
 291 distribution with  $M$  degrees of freedom. In practice, the mean and the covariance matrix are estimated  
 292 by the data matrix  $\mathbf{X}$  as  $\mathbf{S} = \mathbf{X}^t \mathbf{X} / (N - 1)$ . So the approximation to the Mahalanobis distance is the  
 293 Hotelling- $T^2$ :

$$T_n^2 = \mathbf{x}_n^t \cdot \mathbf{S}^{-1} \cdot \mathbf{x}_n \quad \forall n \in 1, \dots, N \quad (17)$$

294 where  $\mathbf{x}_n$  is the  $n$ th row of the data matrix  $\mathbf{X}$ , corresponding to a concrete object. The control limit for the  
 295 Hotelling- $T^2$  [60] is computed as :

$$T_\alpha^2 = \frac{(N^2 - 1)A}{N(N - A)} F_\alpha(A, N - A) \quad (18)$$

296 where  $A$  is the number of PCs of the model, and  $F_\alpha(A, N - A)$  is the  $100 \times (1 - \alpha)$  percentile of a Snedecor's  
 297  $F$  distribution with  $(A, N - A)$  degrees of freedom.

#### 298 *MEDA*

299 The Missing-data methods for Exploratory Data Analysis (MEDA) [61] can be seen as a substitute of  
 300 rotation methods with better properties. First of all, it is more accurate than rotation methods in the  
 301 detection of relations between pairs of variables. Also, it is robust to the overestimation of the number of  
 302 PCs and it does not depend on the normalization of the loadings.

303 Once the PCA has been performed, the MEDA approach consists of the following steps for each variable  
 304  $m$  ( $m = 1, \dots, M$ ):

- 305 1. Construct matrix  $\tilde{\mathbf{X}}_m$ , which is a  $N \times M$  matrix full with zeros except in the  $m$ th column where it  
 306 contains the  $m$ th column of matrix  $\mathbf{X}$ ,  $\mathbf{x}_m$ .

$$\tilde{\mathbf{X}}_m = [\mathbf{0} \dots \mathbf{0} \quad \mathbf{x}_m \quad \mathbf{0} \dots \mathbf{0}] \quad (19)$$

- 307 2. Estimate the scores from  $\tilde{\mathbf{X}}_m$  using a missing-data method. In this case, the Known-Data Regression  
 308 (KDR) method is applied, which has been proved to be statistically superior to other missing data  
 309 imputation techniques in [62].

$$\hat{\mathbf{T}}_A = MD(\tilde{\mathbf{X}}_m) \quad (20)$$

- 310 3. Estimate the reconstruction of the original measurements with  $A$  latent variables and compute the  
 311 estimation error

$$\hat{\mathbf{X}}_A = \hat{\mathbf{T}}_A \cdot \mathbf{P}_A^t \quad (21)$$

$$\hat{\mathbf{E}}_A = \mathbf{X} - \hat{\mathbf{X}}_A \quad (22)$$

312 where  $\mathbf{P}$  is the estimated loadings matrix from  $\mathbf{X}$  (the complete  $N \times M$  matrix of data),  $\hat{\mathbf{X}}_A$  is the  
 313 estimation matrix and  $\hat{\mathbf{E}}_A$  the estimation error matrix.

- 314 4. Compute an index of goodness of prediction [63] in all columns but the  $m$ th one

$$Q_{A,(m,l)}^2 = 1 - \frac{\sum_{n=1}^N (\hat{E}_{A,(n,l)})^2}{\sum_{n=1}^N (X_{n,l})^2}, \quad \forall l \neq m \quad (23)$$

316 where  $X_{n,l}$  is the element located at the  $n$ th row and the  $l$ th column of  $\mathbf{X}$ , and  $\hat{E}_{A,(n,l)}$  is its estimation  
317 error. The closer  $Q_{A,(m,l)}^2$  is to 1, the more related variables  $m$  and  $l$  are.

318 Once the values of  $Q_{A,(m,l)}^2$  for all possible combinations of  $m$  and  $l$  are computed, a matrix  $\mathbf{Q}_A^2$  can  
319 be constructed so that  $Q_{A,(m,l)}^2$  is located at row  $m$  and column  $l$ . This matrix is similar in nature to the  
320 element-wise squared correlation matrix. Structural relations between variables are detected as high values  
321 in  $\mathbf{Q}_A^2$ , but the direct/inverse pair-wise relation is not represented on the matrix because of the squared  
322 values. To avoid obvious relations, the values of principal diagonal of  $\mathbf{Q}_A^2$  matrices are set to zero. When  
323 the number of variables is large, matrix  $\mathbf{Q}_A^2$  can be shown as a grey map to improve interpretability.

324 The  $\mathbf{Q}_A^2$  matrices have been built in a cumulative way, *i.e.* they have the variability of the first  $A$  PCs.  
325 These matrices can also be constructed by taking the information of a single PC. For this purpose, the  
326 method previously detailed has to be changed in Equations 20-23. The new equations have to consider only  
327 the  $a$ -th component for estimation. Finally, this kind of MEDA matrices, which have been used in this  
328 work, are written as  $\mathbf{Q}_{(a)}^2$ , where  $a = 1, \dots, A$ .

#### 329 *Software*

330 All methods commented in this Section have been computed in Matlab environment. The Monte Carlo  
331 sampling method has been applied using the COBRA toolbox [64]. PCA has been performed on MATLAB's  
332 Statistical Toolbox. Finally, MEDA has been performed using Explanatory Data Analysis Toolbox [65].

## 333 4. Results and Discussion

334 In this section, the grey modelling approach proposed is applied to the methylotrophic yeast *P. pastoris*  
335 to discover patterns of heterologous protein production and correlate intra and extracellular reactions in  
336 order to understand how the internal state of the cells determines their observed behavior.

#### 337 *Possibilistic consistency analysis*

338 The different scenarios collected from the literature are combined with the proposed model in order  
339 to validate which ones are consistent and which ones are not. From each one of the 40 scenarios, the  
340 flux values through the external reactions, which are different depending on the initial conditions of each  
341 experiment, are validated against the stoichiometric modelling of the *P. pastoris*. As explained in Section 3,  
342 the most possible solution for each scenario (*i.e.* experimental dataset) is computed to perform a Possibilistic  
343 consistency analysis. The corresponding possibility values ( $\pi$ ) are shown in Table 1. The majority of datasets  
344 are highly consistent with the model (65% are fully possible, and 87% have a possibility higher than 0.5).  
345 There are, however, 4 out of 40 datasets with a possibility lower than 0.25 (*i.e.* a possibility that is equivalent  
346 to an error of 14% in one measurement, or to an error of 8% in three measurements). These scenarios (B3,  
347 B4, C2, and E1) are not fully consistent with the model. The inconsistency can be due to (a) limitations of  
348 the model, which may be unable to capture phenomena occurring in those experiments, (b) larger errors than  
349 expected in the data measured in those scenarios, or (c) the two previous reasons acting simultaneously. For  
350 this reason, we decided to remove these scenarios (B3, B4, C2, and E1) so they are not considered in the  
351 following analysis.

#### 352 *Monte Carlo sampling*

353 The previous analysis concludes that 36 out of the 40 scenarios are consistent with the model. However,  
354 only the external fluxes of each solution have been measured. Due to the complexity of measuring the  
355 internal fluxes, the Monte Carlo sampling method is proposed to simulate different possible flux solutions,  
356 consistent with the proposed model and the measured subset of fluxes, in order to get enough complete flux  
357 solutions to be analysed.

358 Once the sampling has been performed, a feasible flux solution matrix  $\mathbf{X}$  is built.  $\mathbf{X}$  has the complete  
359 3600 sampled flux solutions in its rows (36 scenarios  $\times$  100 samples) and the corresponding 45 flux values  
360 and the protein production for each scenario in its columns (see Additional file 2).

Scenario	Group	$\pi$
A1	<i>glucose</i>	1,000
B1	<i>glycerol</i>	1,000
B2	<i>glycerol + methanol</i>	0,739
B3	<i>glycerol + methanol</i>	0,246(*)
B4	<i>glycerol + methanol</i>	0,082(*)
B5	<i>glycerol</i>	1,000
B6	<i>glycerol + methanol</i>	0,819
B7	<i>glycerol + methanol</i>	0,319
C1	<i>glucose</i>	0,658
C2	<i>methanol</i>	0,052(*)
D1	<i>glycerol + methanol</i>	1,000
D2	<i>glycerol + methanol</i>	1,000
D3	<i>glycerol + methanol</i>	1,000
D4	<i>glycerol + methanol</i>	1,000
D5	<i>glycerol + methanol</i>	1,000
D6	<i>glycerol + methanol</i>	0,908
D7	<i>glycerol + methanol</i>	0,709
D8	<i>glycerol + methanol</i>	0,637
D9	<i>glycerol + methanol</i>	0,614
D10	<i>methanol</i>	0,500
E1	<i>glycerol</i>	0,065(*)
F1	<i>glycerol + methanol</i>	1,000
F2	<i>glycerol + methanol</i>	1,000
F3	<i>glycerol + methanol</i>	1,000
F4	<i>glycerol + methanol</i>	1,000
F5	<i>glycerol + methanol</i>	1,000
F6	<i>glycerol + methanol</i>	1,000
F7	<i>glycerol + methanol</i>	1,000
G1	<i>methanol</i>	1,000
G2	<i>methanol</i>	1,000
G3	<i>methanol</i>	1,000
G4	<i>methanol</i>	1,000
G5	<i>methanol</i>	1,000
G6	<i>methanol</i>	1,000
G7	<i>methanol</i>	1,000
G8	<i>methanol</i>	1,000
G9	<i>methanol</i>	1,000
G10	<i>methanol</i>	1,000
H1	<i>methanol</i>	1,000
I1	<i>glucose</i>	1,000

Table 1: Possibility values ( $\pi$ ) for each scenario. Those scenarios that are not consistent (i.e.  $\pi < 0.25$ ) with the constrained-based model are signaled with (\*).

361 *Principal Component Analysis*

362 A PCA is performed to the feasible flux solutions matrix  $\mathbf{X}$  to obtain a low number of principal compo-  
363 nents (PCs) explaining a high variance percentage of the complete data set. These PCs are linear combina-  
364 tions of the actual 46 variables (45 flux values of the reactions and the protein production).

365 After the PCA has been fitted, explaining 94.3% of total variance with five PCs, the outlier analysis is  
366 applied to the scores and the residuals of the different scenarios. The results on the statistic SPE show that  
367 scenario *C1*, classified on group *glucose* widely exceed the control limits. Thus, the hundred observations  
368 generated by the Monte Carlo sampling for this scenario are knocked-out. Afterwards, a new PCA is fitted.

369 The results with the second analysis are that the first five PCs capture 95.9% of total variance in data:  
370 42.4% for the first component, 24.2% for the second one, 19.7% for the third one, 7.0% for the fourth, and  
371 2.5% for the last one. In these results no outliers are detected.

372 *Enhancement of model interpretation*

373 The MEDA is applied to the first five components obtained by the PCA, and the  $\mathbf{Q}_{(a)}^2$ ,  $a = 1, \dots, 5$   
374 matrices are obtained. Looking at the values in each matrix, the first three PCs are sufficient to explain the  
375 behaviour of the yeast, which capture 86,3% variance in data. Fourth and fifth PCs are classified as noise.  
376 The first three MEDA matrices can be seen in Figure 5.

377 If analysed from a biological standpoint, the first principal component relates protein production rate to  
378 reactions 5-8 (glycolysis), 14-16 and 18 (TCA cycle), 19-20, 28, 30, 36 and 37. In Figure 5a these reactions  
379 are rounded by the solid line rectangle. It can be seen that this relations are indeed strongly correlated,  
380 having  $Q_{(1),(m,l)}^2$  coefficients close to 1. As can be seen in the stoichiometric matrix, each of these groups  
381 is directly connected to NADH and ATP metabolism: ATP is formed in reactions 6, 8, 18 and 28, whereas  
382 NADH is formed in reactions 6, 14, 16 and 18-20. Finally, reactions 28, 30 and 36 represent the electronic  
383 transport chain, oxygen consumption and ATP dissimulation. The first PC can be then understood as  
384 the main pathway for ATP formation and dissimulation, this is, energy generation. Interestingly, protein  
385 productivity and ATP generation have been previously related in a first-principles based approach to predict  
386 recombinant protein production [49].

387 The second principal component is related to the biomass growth rate, which involves reactions 9-13  
388 (fermentative pathways), 17, 21, 29 and 41 (relations shown by dashed line rectangles in Figure 5b). Except  
389 for reaction 41, corresponding to the glycerol consumption rate, reactions 12 (around which reactions 9, 10,  
390 11, 13 and 29 are connected), 17 and 21 share NADPH (either mitochondrial or cytosolic) production, which  
391 is, in fact, one of the major contributing precursors to biomass formation. It is worth noting that reaction 17  
392 (corresponding to NADPH-requiring form) and not 16 (corresponding to the isoenzyme NADH-requiring)  
393 is identified.

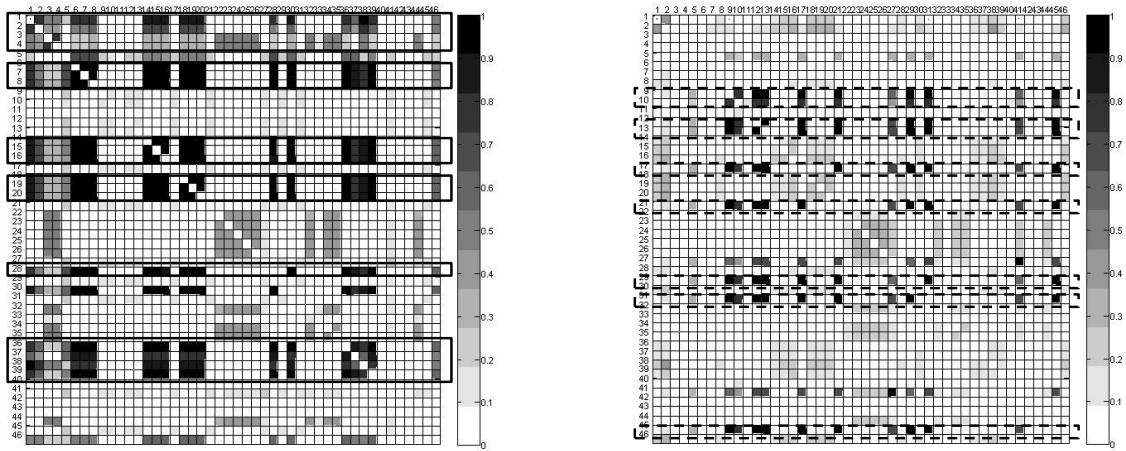
394 Finally, the third principal component relates methanol consumption rate to the pentose phosphate  
395 pathway, strongly connected by reaction 34 (reactions correlated are rounded by dotted rectangles in Figure  
396 5c). Reactions 3-4, 22-26, 32, 35 and 44 are also related with this component.

397 The first three principal pathways are depicted in Figure 6. In this way, the reactions involved by the  
398 three first principal components seem to pinpoint specific metabolic indicators (cofactors NADH, NADPH  
399 and ATP) and their relation with protein, biomass and substrate (glycerol and methanol) consumption.

400 It is worth pointing out that the fit of a PCA model on the available experimental data is not feasible  
401 due to two main reasons: i) only seven out of nine external fluxes are measured for all scenarios under  
402 study, of which three have zero values mostly (see Figure 2), ii) the flux distributions across the metabolic  
403 network cannot be represented since no internal fluxes are considered. Actually, a PCA does not clearly  
404 relate substrates consumption to biomass and protein production, so this model is not meaningful (results  
405 not shown).

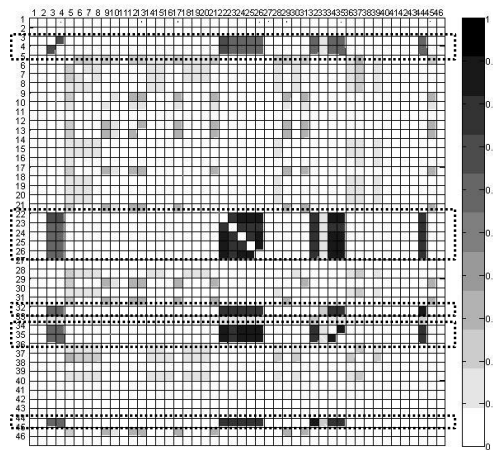
406 **5. Conclusions**

407 In this paper, a grey modelling strategy that combines data-driven and knowledge-based information at  
408 different scales is presented to analyse the behaviour of the methylotrophic yeast *P. pastoris*. This strategy



(a)  $Q^2_{(1)}$

(b)  $Q^2_{(2)}$



(c)  $Q^2_{(3)}$

Figure 5: MEDAs plots for the first (a), second (b) and third (c) PC. Solid line rectangles marks reactions related to the first PC, dashed line rectangles round reactions associated to the second PC and, finally, reactions related to the third PC are rounded by dotted line rectangles.

409 is composed of five main steps. Firstly, the available flux measurements, mainly external, are coupled with  
 410 a model-based estimation of the unmeasured fluxes, mainly internal. Secondly, a possibilistic analysis  
 411 is applied to check the consistency between the constraint-based model and data. Thirdly, a Monte Carlo  
 412 sampling is performed to produce feasible flux solutions for the microbial system under study. As a result,  
 413 a large solution data set -a mixture of experimental data, data-based estimations, and variability resulting  
 414 from uncertainty- is obtained. Fourthly, a PCA model is fitted on the sampled data to reveal its internal  
 415 biochemical structure. Finally, MEDA analysis is performed to enhance the interpretation of the Principal  
 416 Components (PCs) derived from the PCA analysis.

417 The grey modelling of the methylotrophic yeast *P. pastoris* yielded three meaningful PCs from the  
 418 biological point of view, which are sufficient to explain most of the variance of the sampled data. The first



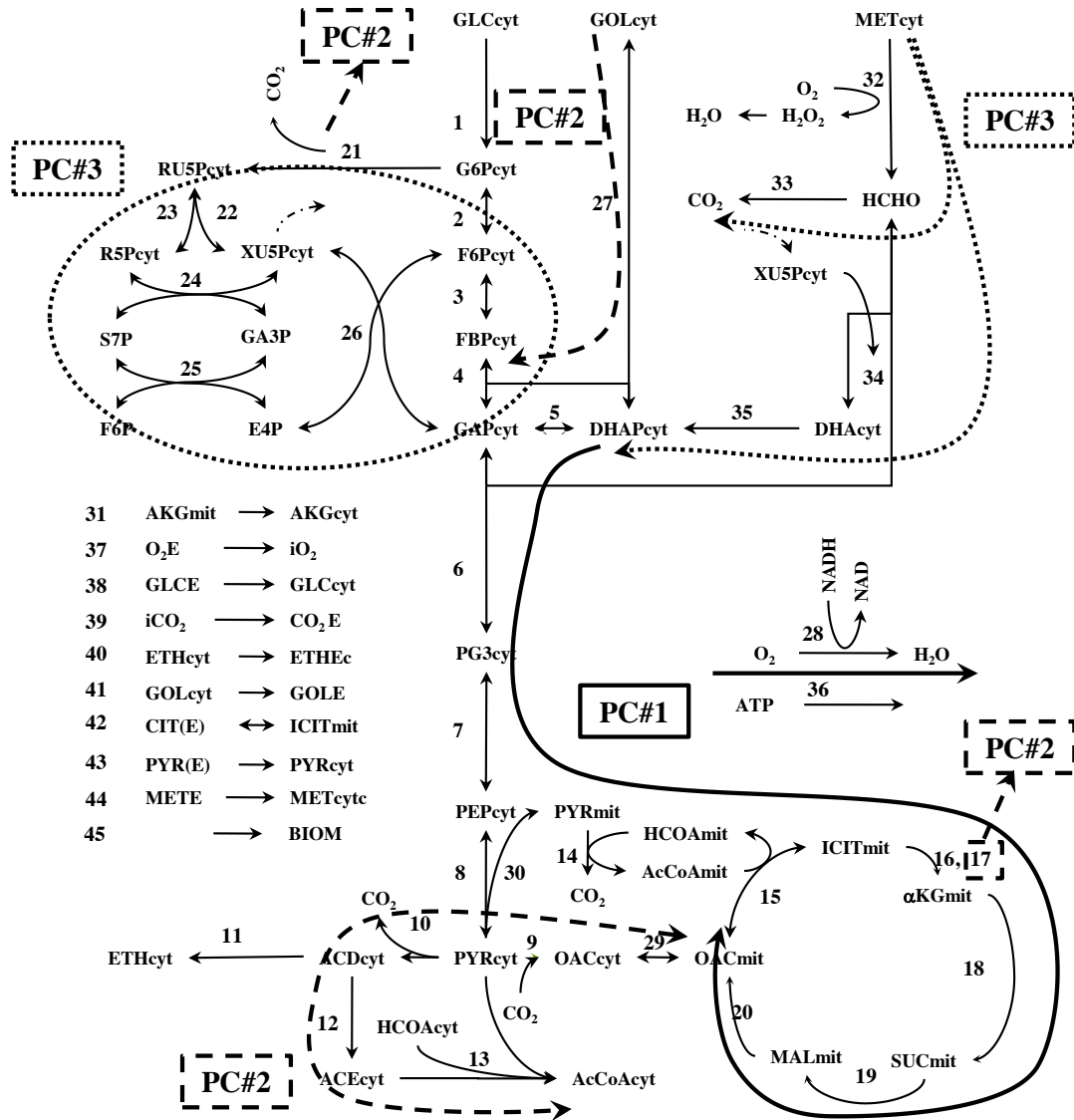


Figure 6: The first three PCs represent the main metabolic pathways through the yeast *P. pastoris*.

419 PC is related to protein productivity, the second one corresponds to the biomass growth rate, and the third  
 420 one represents methanol consumption rate. Note that the experimental data taken from literature do not  
 421 describe the whole space of behaviours that *P. pastoris* could exhibit. In order to explore all the feasible  
 422 solution space, future work should address this issue in two ways: (a) a generalisation: incorporating more  
 423 datasets to explore in a wider scope, and (b) a particularization: perform a similar study only with similar  
 424 datasets (*e.g.* mixed cultures of *glycerol* and *methanol*), with the aim of getting a deeper insight in their  
 425 differences and their impact of the process performance.

426 An important benefit of the grey modelling and analysis approach presented in this paper is its scalability.  
 427 New knowledge, *e.g.* metabolomics or gene regulation, can be incorporated in the form of extra constraints,  
 428 new flux data can be added via new scenarios, and other data pieces -not fluxes- could be incorporated  
 429 directly into the data matrix before performing the multivariate methods. This joint with its capability to  
 430 describe the most important biochemical processes makes this strategy promising for the design of real-time

431 monitoring systems.

## 432 6. Acknowledgements

433 Research in this study was partially supported by the Spanish Ministry of Science and Innovation and  
434 FEDER funds from the European Union through grants DPI2011-28112-C04-01 and DPI2011-28112-C04-02.  
435 The authors are also grateful to Biopolis SL for supporting this research. We also gratefully acknowledge  
436 Associate Professor José Camacho for providing the Exploratory Data Analysis Toolbox.

## 437 References

- 438 [1] F. Llaneras, Interval and possibilistic methods for constraint-based metabolic models, Master's thesis, Universidad Politécnica de Valencia (2010).  
439  
440 [2] A. Kayser, J. Weber, V. Hecht, U. Rinas, Metabolic flux analysis of *Escherichia coli* in glucose-limited continuous culture. i. growth-rate-dependent metabolic efficiency at steady state, *Microbiology* 151 (3) (2005) 693–706.  
441  
442 [3] M. Tortajada, F. Llaneras, J. Picó, Validation of a constraint-based model of *Pichia pastoris* metabolism under data scarcity, *BMC Systems Biology* 4 (115) (2010) 1–11.  
443  
444 [4] T. Benyamini, O. Folger, E. Ruppin, T. Shlomi, Flux balance analysis accounting for metabolite dilution, *Genome Biology* 11 (4) (2010) 1–9.  
445  
446 [5] M. Covert, E. Knight, J. Reed, M. Herrgard, B. Palsson, Integrating high-throughput and computational data elucidates bacterial networks, *Nature* 429 (2004) 92–96.  
447  
448 [6] M. W. Covert, C. H. Schilling, B. Palsson, Regulation of gene expression in flux balance models of metabolism, *Journal of Theoretical Biology* 213 (1) (2001) 73 – 88.  
449  
450 [7] M. Åkesson, J. Förster, J. Nielsen, Integration of gene expression data into genome-scale metabolic models, *Metabolic Engineering* 6 (4) (2004) 285 – 293.  
451  
452 [8] G. N. Stephanopoulos, A. A. Aristidou, J. Nielsen, *Metabolic engineering: principles and methodologies*, Academic Press, San Diego, 1998.  
453  
454 [9] B. Palsson, *Systems Biology: Properties of Reconstructed Networks*, Cambridge University Press, New York, USA, 2006.  
455 [10] N. D. Price, J. A. Papin, C. H. Schilling, B. O. Palsson, Genome-scale microbial in silico models: the constraints-based approach, *Trends in Biotechnology* 21 (4) (2003) 162 – 169.  
456  
457 [11] F. Llaneras, J. Picó, Stoichiometric modelling of cell metabolism, *Journal of Bioscience and Bioengineering* 105 (1) (2008) 1–11.  
458  
459 [12] M. D. Mesarovic, S. N. Sreenath, J. D. Keene, Search for organising principles: understanding in systems biology, *Systems Biology, IEE* 1 (1) (2004) 19–27.  
460  
461 [13] K. J. Kauffman, P. Prakash, J. S. Edwards, Advances in flux balance analysis, *Current Opinion in Biotechnology* 14 (5) (2003) 491 – 496.  
462  
463 [14] S. Feyoazevedo, B. Dahm, F. Oliveira, Hybrid modelling of biochemical processes: A comparison with the conventional approach, *Computers & Chemical Engineering* 21 (1997) S751–S756.  
464  
465 [15] H. Ramaker, E. van Sprang, S. Gurden, J. Westerhuis, A. Smilde, Improved monitoring of batch processes by incorporating external information, *Journal of Process Control* 12 (4) (2002) 569 – 576.  
466  
467 [16] Y. Takane, T. Shibayama, Principal component analysis with external information on both subjects and variables, *Psychometrika* 56 (1) (1991) 97–120.  
468  
469 [17] Y. Takane, H. Kiers, J. Leeuw, Component analysis with different sets of constraints on different dimensions, *Psychometrika* 60 (2) (1995) 259–280.  
470  
471 [18] J. E. Jackson, *A User's Guide to Principal Components*, Wiley Series in Probability and Statistics, 1991.  
472 [19] B. Sariyar, S. Perk, U. Akman, A. Hortasu, Monte carlo sampling and principal component analysis of flux distributions yield topological and modular information on metabolic networks, *Journal of Theoretical Biology* 242 (2) (2006) 389–400.  
473  
474 [20] C. Barrett, M. Herrgard, B. Palsson, Decomposing complex reaction networks using random sampling, principal component analysis and basis rotation, *BMC Systems Biology* 3 (30) (2009) 1–8.  
475  
476 [21] S. Van Dien, S. Iwatani, Y. Usuda, K. Matsui, Theoretical analysis of amino acid-producing *Escherichia coli* using a stoichiometric model and multivariate linear regression, *Journal of Bioscience and Bioengineering* 102 (1) (2006) 34–40.  
477  
478 [22] R. Bro, PARAFAC. tutorial and applications, *Chemometrics and Intelligent Laboratory Systems* 38 (2) (1997) 149–171.  
479 [23] M. Verouden, R. Notebaart, J. Westerhuis, M. van der Werf, B. Teusink, A. Smilde, Multi-way analysis of flux distributions across multiple conditions, *Journal of Chemometrics* 23 (7-8) (2009) 406–420.  
480  
481 [24] S. Wold, M. Sjöström, L. Eriksson, PLS-regression: a basic tool of chemometrics, *Chemometrics and Intelligent Laboratory Systems* 58 (2) (2001) 109–130.  
482  
483 [25] N. Carinhas, V. Bernal, A. Teixeira, M. Carrondo, P. Alves, R. Oliveira, Hybrid metabolic flux analysis: Combining stoichiometric and statistical constraints to model the formation of complex recombinant products, *BMC Systems Biology* 5 (34) (2011) 1–13.  
484  
485 [26] F. Dieterle, A. Ross, G. Schlotterbeck, H. Senn, Metabolite projection analysis for fast identification of metabolites in metabolomics. application in an amiodarone study, *Analytical Chemistry* 78 (11) (2006) 3551–3561.  
486  
487

- 488 [27] K. Tondel, U. Indahl, A. Gjuvslund, J. Vik, P. Hunter, S. Omholt, H. Martens, Hierarchical cluster-based partial least  
489 squares regression (HC-PLSR) is an efficient tool for metamodelling of nonlinear dynamic models, *BMC Systems Biology*  
490 5 (90) (2011) 1–17.
- 491 [28] J. Westerhuis, E. Derks, H. Hoefsloot, A. Smilde, Grey component analysis, *Journal of Chemometrics* 21 (10-11) (2007)  
492 474–485.
- 493 [29] C. Ng, M. Hussain, Hybrid neural network - prior knowledge model in temperature control of a semi-batch polymerization  
494 process, *Chemical Engineering and Processing: Process Intensification* 43 (4) (2004) 559–570.
- 495 [30] H. Bechmann, H. Madsen, N. K. Poulsen, M. K. Nielsen, Grey box modeling of first flush and incoming wastewater at a  
496 wastewater treatment plant, *Environmetrics* 11 (1) (2000) 1–12.
- 497 [31] B. Sohlberg, Grey box modelling for model predictive control of a heating process, *Journal of Process Control* 13 (3)  
498 (2003) 225 – 238.
- 499 [32] J. M. F. ten Berge, A. K. Smilde, Non-triviality and identification of a constrained tucker3 analysis, *Journal of Chemo-*  
500 *metrics* 16 (12) (2002) 609–612.
- 501 [33] A. Teixeira, J. Dias, N. Carinhas, M. Sousa, J. Clemente, A. Cunha, M. von Stosch, P. Alves, M. Carrondo, R. Oliveira,  
502 Cell functional enomics: Unravelling the function of environmental factors, *BMC Systems Biology* 5 (92) (2011) 1–16.
- 503 [34] M. Dragosits, J. Stadlmann, J. Albiol, K. Baumann, M. Maurer, B. Gasser, M. Sauer, F. Altmann, P. Ferrer, D. Mat-  
504 tanovich, The effect of temperature on the proteome of recombinant *pichia pastoris*, *Journal of Proteome Research* 8 (3)  
505 (2009) 1380–1392.
- 506 [35] A. Solà, P. Jouhten, H. Maaheimo, F. Sánchez-Ferrando, T. Szyperski, P. Ferrer, Metabolic flux profiling of *pichia pastoris*  
507 grown on glycerol/methanol mixtures in chemostat cultures at low and high dilution rates, *Microbiology* 153 (1) (2007)  
508 281–290.
- 509 [36] A. Solà, Estudi del metabolisme central del carboni de *pichia pastoris*, Ph.D. thesis, Universitat Autònoma de Barcelona  
510 (2004).
- 511 [37] C. Jungo, I. Marison, U. von Stockar, Mixed feeds of glycerol and methanol can improve the performance of *pichia pastoris*  
512 cultures: A quantitative study based on concentration gradients in transient continuous cultures, *Journal of Biotechnology*  
513 128 (4) (2007) 824–837.
- 514 [38] H. Ren, J. Yuan, K.-H. Bellgardt, Macrokinetic model for methylotrophic *pichia pastoris* based on stoichiometric balance,  
515 *Journal of Biotechnology* 106 (1) (2003) 53–68.
- 516 [39] M. C. d’Anjou, A. J. Daugulis, A rational approach to improving productivity in recombinant *pichia pastoris* fermentation,  
517 *Biotechnology and bioengineering* 72 (1) (2001) 1–11.
- 518 [40] S. Curvers, J. Linnemann, T. Klauser, C. Wandrey, R. Takors, Recombinant protein production with *pichia pastoris* in  
519 continuous fermentation - kinetic analysis of growth and product formation, *Chemical Engineering and Technology* 25 (8)  
520 (2002) 229–235.
- 521 [41] W. Zhang, C.-P. Liu, M. Inan, M. Meagher, Optimization of cell density and dilution rate in *pichia pastoris* continuous  
522 fermentations for production of recombinant proteins, *Journal of Industrial Microbiology and Biotechnology* 31 (7) (2004)  
523 330–334.
- 524 [42] B. Schilling, J. Goodrick, N. Wan, Scale-up of a high cell-density continuous culture with *pichia pastoris* x-33 for the  
525 constitutive expression of rh-chitinase, *Biotechnology Progress* 17 (4) (2001) 629–633.
- 526 [43] F. Llaneras, A. Sala, J. Picó, A possibilistic framework for constraint-based metabolic flux analysis, *BMC Systems Biology*  
527 3 (79) (2009) 1–22.
- 528 [44] L. Zadeh, Fuzzy sets as a basis for a theory of possibility, *Fuzzy Sets and Systems* 1 (1) (1978) 3–28.
- 529 [45] J. Camacho, Missing-data theory in the context of exploratory data analysis, *Chemometrics and Intelligent Laboratory*  
530 *Systems* 103 (1) (2010) 8–18.
- 531 [46] G. Potvin, A. Ahmad, Z. Zhang, Bioprocess engineering aspects of heterologous protein production in *pichia pastoris*: A  
532 review, *Biochemical Engineering Journal* 64 (2012) 91–105.
- 533 [47] S. Macauley-Patrick, M. Fazenda, B. McNeil, L. Harvey, Heterologous protein production using the *pichia pastoris* ex-  
534 pression system, *Yeast* 22 (4) (2005) 249–270.
- 535 [48] O. Cos, R. Ramón, J. Montesinos, F. Valero, Operational strategies, monitoring and control of heterologous protein  
536 production in the methylotrophic yeast *pichia pastoris* under different promoters: A review, *Microbial Cell Factories*  
537 5 (17) (2006) 1–20.
- 538 [49] M. Tortajada, F. Llaneras, D. Ramón, J. Picó, Estimation of recombinant protein production in *pichia pastoris* based on  
539 a constraint-based model, *Journal of Process Control* 22 (6) (2012) 1139–1151.
- 540 [50] M. Benito, J. Parker, Q. Du, J. Wu, D. Xiang, C. M. Perou, J. S. Marron, Adjustment of systematic microarray data  
541 biases, *Bioinformatics* 20 (1) (2004) 105–114.
- 542 [51] W. E. Johnson, C. Li, A. Rabinovic, Adjusting batch effects in microarray expression data using empirical bayes methods,  
543 *Biostatistics* 8 (1) (2007) 118–127.
- 544 [52] J. Luo, M. Schumacher, A. Scherer, D. Sanoudou, D. Megherbi, T. Davison, T. Shi, W. Tong, L. Shi, H. Hong, A  
545 comparison of batch effect removal methods for enhancement of prediction performance using maqc-ii microarray gene  
546 expression data, *The Pharmacogenomics Journal* (4) (2010) 278291.
- 547 [53] S. E. Reese, K. J. Archer, T. M. Therneau, E. J. Atkinson, C. M. Vachon, M. de Andrade, J.-P. A. Kocher, J. E.  
548 Eckel-Passow, A new statistic for identifying batch effects in high-throughput genomic data that uses guided principal  
549 components analysis, *Bioinformatics*doi:10.1093/bioinformatics/btt480.
- 550 [54] F. Hadlich, K. Nöh, W. Wiechert, Determination of flux directions by thermodynamic network analysis: Computing  
551 informative metabolite pools, *Mathematics and Computers in Simulation* 82 (3) (2011) 460–470.
- 552 [55] D. Machado, R. Costa, E. Ferreira, I. Rocha, B. Tidor, Exploring the gap between dynamic and constraint-based models

- 553 of metabolism, *Metabolic Engineering* 14 (2) (2012) 112–119.
- 554 [56] J. Camacho, J. Picó, A. Ferrer, Data understanding with PCA: structural and variance information plots, *Chemometrics*  
555 *and Intelligent Laboratory Systems* 100 (1) (2010) 48–56.
- 556 [57] J. Jackson, G. S. Mudholkar, Control procedures for residuals associated with principal component analysis., *Technometrics*  
557 21 (3) (1979) 341–349.
- 558 [58] P. C. Mahalanobis, On the generalised distance in statistics, *Proceedings of the National Institute of Sciences of India*  
559 2 (1) (1936) 49–55.
- 560 [59] F. Arteaga, Control estadístico multivariante de procesos con datos faltantes mediante análisis de componentes principales,  
561 Ph.D. thesis, Universidad Politécnica de Valencia (2003).
- 562 [60] N. D. Tracy, J. C. Young, R. L. Mason, Multivariate control charts for individual observations, *Journal of Quality*  
563 *Technology* 24 (2) (1992) 88–95.
- 564 [61] J. Camacho, Missing-data theory in the context of exploratory data analysis, *Chemometrics and Intelligent Laboratory*  
565 *Systems* 103 (1) (2010) 8 – 18.
- 566 [62] F. Arteaga, A. Ferrer, Dealing with missing data in MSPC: several methods, different interpretations, some examples,  
567 *Journal of Chemometrics* 16 (8-10) (2002) 408–418.
- 568 [63] S. Wold, Cross-validatory estimation of the number of components in factor and principal components models, *Techno-*  
569 *metrics* 20 (4) (1978) 397–405.
- 570 [64] S. Becker, A. Feist, M. Mo, G. Hannum, B. Palsson, M. Herrgard, Quantitative prediction of cellular metabolism with  
571 constraint-based models: The COBRA toolbox, *Nature Protocols* 2 (3) (2007) 727–738.
- 572 [65] J. Camacho[link].  
573 URL <http://wdb.ugr.es/~josecamacho/>