

1 Evolutionary algorithm for metabolic pathways synthesis

2 Matias F. Gerard*, Georgina Stegmayer, Diego H. Milone

3 *Research Institute for Signals, Systems and Computational Intelligence (sinc(i)),*
4 *FICH-UNL/CONICET, Argentina.*

5 Abstract

6 Metabolic pathway building is an active field of research, necessary to
7 understand and manipulate the metabolism of organisms. There are dif-
8 ferent approaches, mainly based on classical search methods, to find linear
9 sequences of reactions linking two compounds. However, an important lim-
10 itation of these methods is the exponential increase of search trees when
11 a large number of compounds and reactions is considered. Besides, such
12 models do not take into account all substrates for each reaction during the
13 search, leading to solutions that lack biological feasibility in many cases.
14 This work proposes a new evolutionary algorithm that allows searching not
15 only linear, but also branched metabolic pathways, formed by feasible reac-
16 tions that relate multiple compounds simultaneously. Tests performed using
17 several sets of reactions show that this algorithm is able to find feasible linear
18 and branched metabolic pathways.

19 *Keywords:* Evolutionary algorithms, search strategies, *de novo* pathway
20 building, reactions network, sets of compounds.

*Corresponding author. Tel: +54 (0342) 457 5234 int 118.

Email addresses: mgerard@sinc.unl.edu.ar (Matias F. Gerard),
gstegmayer@sinc.unl.edu.ar (Georgina Stegmayer), dmilone@sinc.unl.edu.ar (Diego
H. Milone)

21 **1. Introduction**

22 Systems biology has quickly progressed thanks to the technical advances
23 made in recent years to obtain quantitative and qualitative information of
24 biological systems at different scales. These developments, in addition to
25 contributions made by bioinformatics in several areas such as sequence anal-
26 ysis, modelling of protein structures, and building of interaction networks,
27 help to understand the functioning of living beings (Tenazinha and Vinga,
28 2011). However, the increasing volume of data produced in biological exper-
29 iments has led to the need to develop new computational tools capable of
30 manipulating and analyzing it to extract knowledge (Bordbar *et al.*, 2014;
31 Chen and Zhang, 2014).

32 In nature, metabolic processes do not occur in isolation, but rather
33 through complex networks made up of metabolic pathways that branch and
34 interconnect (Ravasz *et al.*, 2002; Lacroix *et al.*, 2008). They generate a large
35 variety of compounds that are used, for example, for structural purposes or
36 energy storage, or just as substrates for key reactions in other processes
37 (Jeong *et al.*, 2000). These networks are a natural way of organising rela-
38 tions (biochemical reactions) between compounds. Each reaction acts as a
39 rule that determines the compounds consumed (substrates) and produced
40 (products) in the process. These intricate relations are frequently modelled
41 employing different types of graphs (Arita, 2012). Determining the whole
42 sequence of reactions to produce a compound from another one consists in
43 searching for a path that links both compounds in the graph. This problem
44 is of particular interest in systems biology nowadays. The effort is focused
45 on developing tools that allow identification of metabolic pathways capable
46 of being manipulated to produce compounds of interest (Lee *et al.*, 2009;

47 Yim *et al.*, 2011).

48 There are different methods to automatically search for metabolic path-
49 ways between two compounds. They are mainly based on classical search
50 algorithms, such as breadth-first and depth-first search, and the A* algo-
51 rithm (Russell and Norvig, 2010). All of them start by transforming the
52 data into a type of graph appropriate for the search (Pey *et al.*, 2011). One
53 problem with these representations are the abundant compounds such as
54 water and Adenosine 5'-triphosphate (ATP), which have a high connectiv-
55 ity as they participate in a large number of reactions (Gerlee *et al.*, 2009).
56 Thus, frequently the solutions found by the search strategies do not make
57 biological sense since they use abundant compounds as intermediate steps
58 in the synthesis of the desired product, and the availability of the other
59 required substrates is not verified.

60 Different approaches to solve the problem of abundant compounds have
61 been proposed. Croes *et al.* (Croes *et al.*, 2005) propose a weighting scheme
62 to search a pathway between two compounds. They assign to each node
63 a weight equal to the number of reactions where it participates, and find
64 the lightest pathways between both ends. This approach was extended by
65 Faust *et al.* (Faust *et al.*, 2009), who applied the weighting scheme to a
66 graph where its edges indicates the transfer of atoms from one compound to
67 another one. Employing structural information of the compounds, McShan
68 *et al.* (McShan *et al.*, 2003) built vectors of characteristic for each compound
69 and performed the search by selecting the successive nodes using heuristics
70 based on the distance between vectors. Similarly, Rahman *et al.* (Rahman
71 *et al.*, 2005) generated a binary fingerprint for each compound and applied
72 similarity measures to guide the search process. Heath *et al.* (Heath *et al.*,
73 2010) proposed an approach based on tracking the flow of atoms, from the

74 starting to the ending compound, trying to preserve as many of these atoms
75 as possible. This allowed finding linear and branched pathways between
76 two compounds. Branched solutions contain several alternative mechanisms
77 to transfer atoms from the start to the end of the pathway. The main
78 problem faced by those methods is the exponential growth of the search
79 trees when a large number of highly connected reactions and compounds
80 are involved. Recently, a method based on evolutionary algorithms to search
81 metabolic pathways between two metabolites was developed (Gerard *et al.*,
82 2013), which avoids the problems of working with growing search trees.
83 These methods provide paths only between two compounds and take into
84 account the last synthesized product to select a new reaction.

85 Despite their characteristics, all these methods cannot find branching
86 metabolic pathways that relate more than two compounds. In an effort
87 to solve this issue, Faust *et al.* (Faust *et al.*, 2010, 2011) extended their
88 pathway search strategy to relate a set of compounds by means of a network
89 of reactions. Thus, solutions found consist of networks built as a combination
90 of linear pathways among all pairs of compounds specified. Even though
91 these solutions have ramifications, the feasibility of solutions is not taken
92 into account since the availability of all substrates is not guaranteed.

93 While all these proposals provide sequences of reactions that relate the
94 indicated compounds, the solutions found are often not biologically feasible.
95 This is due to the assumption that all substrates are available, thereby the
96 solution consists in finding a sequence of reactions to establish the relation.
97 Thus, the availability of the compounds is not taken into account to perform
98 the search and no restrictions are imposed on the possible reactions used to
99 generate the solutions. Furthermore, given that all the previously synthe-
100 sized compounds in the reactions chain are not taken into account to select a

101 new reaction, valuable information to guide the search is lost and not prop-
102 erly used. It is important to highlight that there are cases where a pathway
103 between two compounds needs a branching to be possible. For example, in
104 the case where a reaction needs two substrates, and each one of them should
105 be provided by independent reactions that must be carried out in parallel.
106 Supposing that only feasible solutions should be found, algorithms searching
107 lineal pathways could not find any solution in this case.

108 This work proposes a new approach based on the expanded set of com-
109 pounds concept (ESC), which allows to relate several compounds at the
110 same time by means of a network of feasible reactions. Given a set of avail-
111 able compounds and a feasible reaction from them, it is possible to expand
112 this set by adding the products of the reaction. In this way, it is possi-
113 ble for a higher number of reactions can take place from the new set of
114 compounds. Following this idea, our method only needs an initial set of
115 available compounds in order to search for a metabolic pathway that re-
116 lates the compounds of interest. To efficiently explore the search space, an
117 algorithm based on evolutionary computation is proposed. This family of
118 algorithms are inspired in biology and employ the principle of natural selec-
119 tion to evolve a population of potential solutions (Pal *et al.*, 2006; Affenzeller
120 *et al.*, 2009; Boussaïd *et al.*, 2013). These methods have been successfully
121 applied to solve a wide range of problems in bioinformatics (Lee and Hsiao,
122 2012; Kayaa and Şule Gündüz-Öğüdücü, 2013; de Magalhães *et al.*, 2014;
123 Garai and Chowdhury, 2015). The search is guided by the fitness of indi-
124 vidual in the population, which is evaluated using functions without formal
125 requirements. Each individual encodes a solution, evolved employing ge-
126 netic operators that combine the information of different individuals and
127 introduce small variations during the evolutionary process.

128 A web interface to the algorithm has appeared in (Gerard *et al.*, 2015).
129 That report simply described the software from a user point of view, without
130 details of the model and its functioning, mainly with a focus on the usability
131 of the tool and the visualizations provided. It has to be noticed that
132 this present contribution, instead, develops the main ideas behind the tool,
133 providing a detailed explanation of the evolutionary model, its internal parameters
134 and a wide experimental validation, with artificial as well as several
135 real data of increasing complexity. The analysis of sensibility to parameters
136 and robustness when facing a real problem is also included in the results.
137 Moreover, a real case study for a well-known metabolic pathway that relates
138 four biologically relevant compounds is presented, and two alternative
139 solutions found to the standard metabolic pathway are described.

140 The paper is organized as follows. Section 2 describes the model of sets
141 of compounds employed, the encoding in chromosomes, and the elements of
142 the evolutionary algorithm, analyzing in detail the proposed operators and
143 the measures that make up the fitness function. Section 3 describes the data
144 employed in the experiments, their processing, the measures used to evaluate
145 the algorithm performance, and several aspects of the searched networks.
146 Section 4 analyses the effect of the variation of different parameters of the
147 algorithm, the ability of the algorithm to scale to larger spaces, and a real
148 case study. Finally, Section 5 presents the conclusions and future work.

149 **2. Evolutionary algorithm based on expanded sets of compounds**

150 Metabolic networks are constituted by compounds and the biochemical
151 reactions r relating them (Lacroix *et al.*, 2008). These relations allow certain
152 groups of substrates to be modified in order to produce new products. For-

153 mally, reactions can be represented by means of the relation $S(r) \xrightarrow{r} P(r)$,
154 where $S(r)$ and $P(r)$ correspond to the substrates and products of the re-
155 action. Clearly, these relations require all substrates to be present in order
156 to take place. In some cases, substrates are available in the medium where
157 the reaction occurs. In other cases, they must be provided externally or
158 through a previous reaction. In any case, each reaction which takes place
159 can increase the available compounds so that new reactions can take place.
160 This idea can be employed to model a metabolic pathway by considering it
161 as a set of reactions carried out with a given order, that starts from a speci-
162 fied set of available compounds. Additionally, it is also possible to evaluate
163 the feasibility of each reaction in the pathway by analyzing the availability
164 of its substrates.

165 In an evolutionary algorithm, the linear structure of genes into a chromo-
166 some \mathbf{c} can be easily used to represent the sequence of reactions, considering
167 its order as indicative of the order that they take place in the pathway.
168 Besides, it is possible to evaluate the feasibility of the pathway by associ-
169 ating an initial set of available compounds C^0 to \mathbf{c} , and verifying whether
170 each reaction is possible based on this set and the products of all feasible
171 reactions that have been previously carried out. Additionally, the use of
172 an ESC enables to model branched metabolic pathways, where two or more
173 reactions must happen simultaneously in order to generate all the necessary
174 substrates for a subsequent reaction. Therefore, each chromosome encodes
175 a complete metabolic pathway, varying its size according to the number of
176 reactions the pathway has.

177 Figure 1 exemplifies a metabolic pathway encoded in a chromosome to-
178 gether with the ESC associated to each reaction. The substrates required
179 for the reaction r_k must be available in the ESC C^{k-1} , otherwise the re-

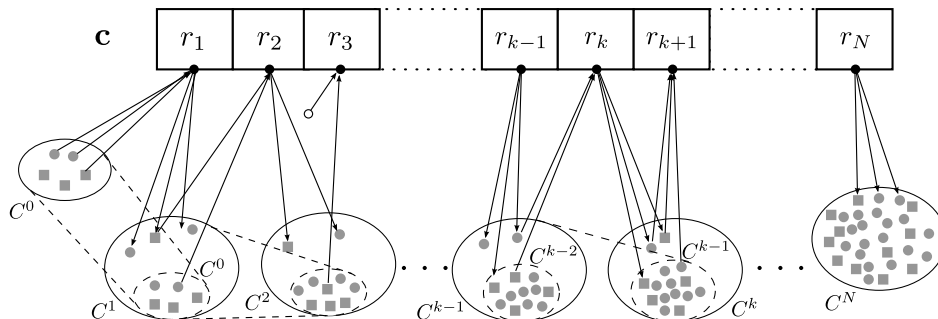


Figure 1: Representation of the ESC model in a chromosome. *Top*: Chromosome that encodes the reactions of a metabolic pathway. *Bottom*: The ESC for each reaction (solid lines) and previous sets (dash lines). Squares indicate available compounds. Filled circles correspond to new compounds generated in the metabolic pathway. The empty circle corresponds to a substrate required by the reaction r_3 that is not available in C^2 .

180 action will not be valid and the set of compounds will remain unmodified
 181 ($C^k = C^{k-1}$). Thus, if the substrates for the reaction r_k are available
 182 in the ESC $C^{k-1} = C^{k-2} \cup P(r_{k-1})$, this reaction produces the new set
 183 $C^k = C^{k-1} \cup P(r_k)$. Therefore, the ESC continues to be updated until the
 184 set C^N is reached.

185 2.1. Description of the algorithm

186 The proposed algorithm, named EvoMS (Evolutionary Metabolic Seeker),
 187 employs the sets of compounds model to search for feasible metabolic path-
 188 ways that relate a group D of specified compounds. In order to facilitate
 189 further explanations, the term *initial substrate* is introduced to denote the
 190 compound belonging to D used to find the pathway, and *final products* to
 191 indicate the remaining compounds in D after selecting the initial one. The

192 general structure of the algorithm and the selection operator are similar to
193 the ones used in genetic algorithms (Bäck *et al.*, 2000).

194 Briefly, the algorithm starts with the initialization and fitness evalua-
195 tion of the population in the first generation, $f(\mathbb{P}^0)$, which is subjected to
196 the evolution process until the stopping criterion is satisfied. This criterion
197 consists of two elements: a maximum allowed number of generations G_M
198 and a fitness value 1.0. The evolutionary process comprises six steps: ex-
199 tracting the best individual (chromosome \mathbf{c}^*), selecting the parents \mathbb{X}^G for
200 the new generation, creating the descendants \mathbb{C}^G through crossover of the
201 selected parents and mutation of their offspring, building the new popula-
202 tion $\mathbb{P}^{G+1} \leftarrow \{\mathbf{c}^*\} \cup \mathbb{X}^G \cup \mathbb{C}^G$ and evaluating the fitness $f(\mathbb{P}^{G+1})$ of the
203 new population. The solutions found by EvoMS correspond to networks
204 of feasible reactions that use C^0 to relate compounds in D . The feasible
205 reactions which are not part of these links are filtered later. The crossover
206 operator employed consists of a combination of one-point and two-point
207 crossover operators. Given two parents, this operator selects a portion of
208 genetic material from one parent and inserts it in a random position of the
209 other one, discarding the original genetic material in the second parent after
210 the point of insertion. The mutation operator and the initialization strategy
211 consider the use of sets of compounds. These will be explained in detail in
212 subsections 2.2 and 2.4.

213 *2.2. Initialization based on ESC*

214 The initialization of EvoMS is carried out employing a strategy based
215 on ESC and taking into account the validity of the reactions. The use of
216 this strategy has two objectives. On the one hand, it avoids using random
217 initialization, which could lead to very poor initial solutions. On the other

218 hand, it introduces the use of subpopulations. Each one is made up by a
219 set of individuals using the same initial substrate. It allows to overcome the
220 problem of selecting the initial one when there is no information to make
221 such decision. Thus, subpopulations will compete to determine the initial
222 substrate for the metabolic pathway searched. The initialization process is
223 carried out in two phases: identifying the number of subpopulations and
224 initializing the individuals. Algorithm 1 describes the steps of this process.
225 In order to initialize the population \mathbb{P} , it is necessary to define a set of
226 abundant compounds A , such as water and ATP, which will be available for
227 all reactions during the search. This set is automatically updated during
228 the initialization, incorporating the external compounds E to generate the
229 set $A' = A \cup E$. The set E is made up of all substrates that cannot be
230 synthesized by any reaction provided.

231 The first phase of the initialization consists in determining the number of
232 subpopulations to be generated (lines 7–12 of Algorithm 1). Each compound
233 $d \in D$ is evaluated in order to identify those which are used as substrate of
234 any reaction. Used compounds and substrates of those reactions are stored
235 in two lists, I and R , respectively. The amount of compounds in list I define
236 the number of subpopulations that should be created.

237 The second phase consists in the initialization of subpopulations, each
238 one containing equal number of individuals (lines 13–26). This process is
239 similar for all members. Firstly, the chromosome \mathbf{c} is initialized as an empty
240 list, and the number of genes N_I that it should contain is randomly selected.
241 Secondly, a set of available compounds C^0 associated to the chromosome is
242 built. It is made of the union of the abundant compounds (A) and the
243 external ones (E), plus all the substrates (Q_j) required by reactions that
244 use the initial substrate I_j . The initial reaction r_1 is randomly selected from

245 those using I_j as initial substrate, and its products update the set of available
246 compounds $C^1 = C^0 \cup P(r_1)$. Then, an iterative process is performed until
247 the specified number of genes N_I is reached, or there is no more reactions to
248 insert. In each step, a reaction r_k is selected at random, without repetition
249 from all reactions than can take place from the compounds present in the set
250 C^{k-1} . Afterwards, the set of accumulated compounds $C^k = C^{k-1} \cup P(r_k)$
251 is updated with products of the selected reaction. Finally, the individual
252 is incorporated to the population \mathbb{P} and the process is repeated. If the
253 final population has more than M individuals, some members are randomly
254 removed until the specified size is reached (lines 27–28).

Algorithm 1: Initialization strategy based on sets of compounds.

```

1  $A' \leftarrow A \cup E$ 
2  $N_M \leftarrow$  maximum pathway size allowed
3  $M \leftarrow$  population size
4  $N \leftarrow 0$ 
5  $Q, I \leftarrow$  empty list
6  $U \leftarrow \emptyset$ 
7 foreach  $d \in D$  do
8    $U \leftarrow \bigcup_{r/d \in S(r)} S(r)$ 
9   if  $U \neq \emptyset$  then
10      $N \leftarrow N + 1$ 
11      $Q_N \leftarrow U$ 
12      $I_N \leftarrow d$ 
13 for  $j \leftarrow 1$  to  $N$  do
14   for  $i \leftarrow 1$  to  $\lceil \frac{M}{N} \rceil$  do
15      $k \leftarrow 1$ 
16      $N_I \leftarrow$  select a random integer in  $[\frac{N_M}{2}, N_M]$ 
17      $\mathbf{c} \leftarrow$  empty list
18      $C^0 \leftarrow A' \cup (Q_j - D) \cup \{I_j\}$ 
19      $R \leftarrow \{r / |S(r) \cap C^{k-1}| = |S(r)| \wedge I_j \in S(r)\}$ 
20     while  $k \leq N_I$  and  $R \neq \emptyset$  do
21        $r_k \leftarrow$  select one reaction from  $R$  not included in  $\mathbf{c}$ 
22        $\mathbf{c} \leftarrow$  insert  $r_k$ 
23        $C^k \leftarrow C^{k-1} \cup P(r_k)$ 
24        $k \leftarrow k + 1$ 
25        $R \leftarrow \{r / |S(r) \cap C^{k-1}| = |S(r)|\}$ 
26      $\mathbb{P} \leftarrow$  insert  $\mathbf{c}$ 
27 if  $|\mathbb{P}| > M$  then
28    $\mathbb{P} \leftarrow$  randomly select  $M$  individuals from  $\mathbb{P}$ 
29 return  $\mathbb{P}$ 

```

256 *2.3. Fitness function*

257 The fitness $f(\mathbf{c})$ of the individuals in the population is evaluated em-
 258 ploying an additive function made up of four terms, each one focused on
 259 a specific property of the solution. The fitness function and its terms are
 260 normalized in $[0, 1]$, and the maximum fitness is reached when a solution

261 is found. A metabolic pathway is considered a solution when it meets two
 262 conditions: i) each reaction has the necessary substrates, and ii) there is a
 263 sequence of valid reactions that relate the initial substrate with each final
 264 product. Therefore, the fitness function is defined as

$$f(\mathbf{c}) = \frac{1}{4} [\mathcal{V}(\mathbf{c}) + \mathcal{L}(\mathbf{c}) + \mathcal{I}(\mathbf{c}) + \mathcal{C}(\mathbf{c})], \quad (1)$$

265 and the way of calculating the four measures is described below.

266 2.3.1. Validity

267 The term $\mathcal{V}(\cdot)$ evaluates the proportion of reactions in the metabolic
 268 pathway that have the required substrates. In this sense, the reaction r_k is
 269 valid if $S(r_k) \subseteq C^{k-1}$, which corresponds to the set of accumulated com-
 270 pounds until the reaction r_{k-1} . This measure is calculated as

$$\mathcal{V}(\mathbf{c}) = \frac{1}{|\mathbf{c}|} \sum_{k=1}^{|\mathbf{c}|} \mathbf{1}_{S(r_k) \subseteq C^{k-1}}, \quad (2)$$

271 where $|\mathbf{c}|$ is the number of genes of \mathbf{c} , and $\mathbf{1}_{A \subseteq B}$ is the indicator function,
 272 which takes the value 1 when $A \subseteq B$ and 0 in another case. The validity of
 273 a metabolic pathway is maximum when each reaction has the substrates it
 274 needs.

275 2.3.2. Linking

276 The term $\mathcal{L}(\cdot)$ in (1) evaluates two aspects of the metabolic pathway: i) if
 277 the initial substrate is used, at least, by one reaction, and ii) the proportion
 278 of the final products that are synthesized. This measure is calculated as

$$\mathcal{L}(\mathbf{c}) = \frac{1}{2} \left(|S^*(\mathbf{c}) \cap \{d\}| + \frac{|P^*(\mathbf{c}) \cap (D - \{d\})|}{|D - \{d\}|} \right), \quad (3)$$

279 where d denote the initial substrate of \mathbf{c} , $S^*(\mathbf{c}) = \bigcup_{r \in \mathbf{c}} S(r)$ and $P^*(\mathbf{c}) =$
 280 $\bigcup_{r \in \mathbf{c}} P(r)$ are the sets containing all substrates and products of the path-
 281 way, respectively. This measure reaches its maximum value when a reaction
 282 employs d as a substrate and all compounds $D - \{d\}$ are produced.

283 2.3.3. Innovation

284 The term $\mathcal{I}(\cdot)$ determines the proportion of reactions in the metabolic
 285 pathway that produce, at least, one compound that has not been previously
 286 generated in the sequence. Consequently, this term favours the incorporation
 287 of novel reactions that are not already present in the pathway and that
 288 produce new compounds. This measure is calculated as

$$\mathcal{I}(\mathbf{c}) = \frac{1}{|\mathbf{c}|} \sum_{k=1}^{|\mathbf{c}|} \mathbf{1}_{P(r_k) \not\subseteq C^{k-1}}. \quad (4)$$

289 The maximum value is reached when each reaction produces, at least, a new
 290 compound.

291 2.3.4. Connectivity

292 The term $\mathcal{C}(\cdot)$ in (1) evaluates the proportion of the final products for
 293 which there is a sequence of reactions that relates them with the initial
 294 substrate d . This measure is calculated in two steps. The first step consists
 295 in building a set of accumulated compounds Z , which is then used in the
 296 second step to calculate the connectivity. The set Z employed in the first
 297 step is built using Algorithm 2. From the initial set $Z = \{d\}$, the algorithm
 298 evaluates each reaction in the chromosome and verifies whether the reaction
 299 employs any of the compounds in Z as a substrate, updating this set with
 300 its products if the reaction is a valid one. The algorithm returns the set

Algorithm 2: Searching for compounds related to the initial substrate.

```

1  $Z \leftarrow$  initial substrate of  $\mathbf{c}$ 
2 for  $k \leftarrow 1$  to  $|\mathbf{c}|$  do
3   if  $|S(r_k) \cap Z| > 0$  then
4      $Z \leftarrow Z \cup (P(r_k) - C^0)$ 
5   if  $|S(r_k) \cap C^{k-1}| = |S(r_k)|$  then
6      $C^k \leftarrow P(r_k) \cup C^{k-1}$ 
7   else
8      $C^k \leftarrow C^{k-1}$ 
9 return  $Z$ 

```

301 of compounds that are employed to relate d with each member of $D - \{d\}$.
302 Then, connectivity is calculated from the set Z obtained according to

$$303 \quad \mathcal{C}(\mathbf{c}) = \frac{|Z \cap D| - 1}{|D| - 1}. \quad (5)$$

304 This measure takes its maximum value when there are sequences of re-
305 actions that relate the initial substrate with each final product.

306 2.4. Mutation based on ESC

307 The proposed mutation operator introduces changes based on the com-
308 position of the sets of accumulated compounds with a probability p_m . These
309 changes can be the deletion or insertion of one gene into the chromosome,
310 with probabilities p_e and $1 - p_e$, respectively. It introduces variations in
311 the pathway size, because deletions remove randomly one gene from the se-

312 quence, and each insertion adds one reaction that was not already in the
313 sequence.

314 Insertion starts by randomly selecting a position $k \in [1, N + 1]$, where
315 the gene will be inserted. Afterwards, two lists of reactions are built from
316 C^{k-1} . When $k = 1$, C^0 corresponds to the initial set of available compounds
317 associated to the chromosome \mathbf{c} . The list of valid reactions contains all the
318 possible reactions from the compounds in C^{k-1} , while the list of invalid
319 reactions has all the remaining reactions of the search space. The list of
320 valid or invalid reactions from which the reaction will be selected is chosen
321 with probabilities p_v and $1 - p_v$, respectively. When the chosen position is
322 in the interval $[1, N]$, the gene that is in that position and all genes coming
323 after in the sequence are moved one place forward to allow the insertion.

324 Figure 2 shows an example of the proposed mutation operator for a
325 chromosome containing $N = 10$ genes. In the case of insertion, the chosen
326 position is 4 and the set of accumulated compounds C^3 is built considering
327 the products of all valid reactions until the gene that contains r_3 . From
328 this set, the list of valid reactions, whose substrates are available in C^3 , is
329 generated, as well as the list of invalid reactions, which do not have all nec-
330 essary substrates in C^3 . A reaction from these lists is randomly extracted,
331 with probability p_v for the list of valid reactions. In the example of deletion,
332 the gene containing the reaction r_8 is eliminated from the sequence, and
333 adjacent reactions r_7 and r_9 are spliced. Clearly, in both cases the number
334 of genes in the chromosome is modified.

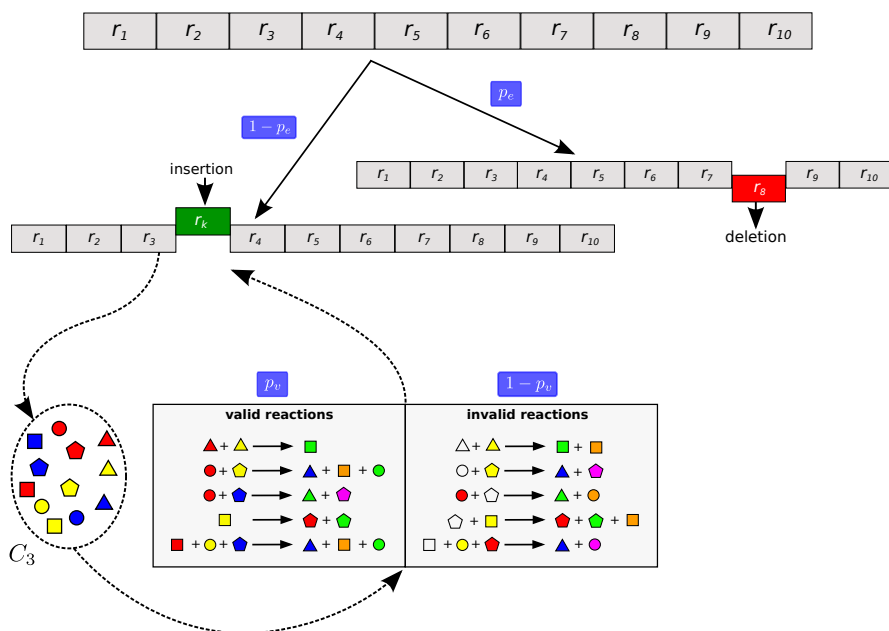


Figure 2: Diagram of the proposed mutation operator for a chromosome containing $N = 10$ reactions. *Left*: Example of gene insertion in position 4 of the chromosome. Available compounds in C^3 are indicated as filled polygons. *Right*: Example of gene deletion.

335 3. Data and evaluation measures

336 3.1. Reactions information

337 Reactions employed in the experiments were extracted from the KEGG
 338 database. Actually, reactions from other repositories, such as MetaCyc (Alt-
 339 man *et al.*, 2013), could be used as well. The direction for each reaction was
 340 assigned using the information contained in the KGML files associate to
 341 the reference maps (Ogata *et al.*, 1998; Goto *et al.*, 2002). Each reversible
 342 reaction was modelled as a pair of independent reactions with opposite di-
 343 rection. For example, the reaction $S(r) \longleftrightarrow P(r)$ was separated into the

Table 1: Abundant compounds employed to search for branched metabolic pathways. The table indicates the name of the compound and the corresponding KEGG code.

KEGG code	name	KEGG code	name	KEGG code	name	KEGG code	name
C00001	H ₂ O	C00005	NADPH	C00009	Phosphate	C00020	AMP
C00002	ATP	C00006	NADP ⁺	C00010	CoA	C00028	Hydrogen acceptor
C00003	NAD ⁺	C00007	O ₂	C00011	CO ₂	C00030	Hydrogen donor
C00004	NADH	C00008	ADP	C00014	Ammonia	C00080	H ⁺

344 semi-reactions $S(r) \rightarrow P(r)$ and $P(r) \rightarrow S(r)$. The set of abundant com-
345 pounds A employed in the experiments is shown in Table 1. The external
346 compounds E were extracted automatically for each set of reactions.

347 *3.2. Measures to evaluate the output network and the algorithm performance*

348 The algorithm performance was analyzed on the basis of 30 runs for
349 each combination of parameters. When results presented an asymmetric
350 distribution, the median and the median absolute deviation were employed
351 as robust estimators to characterize the distribution. The statistical analysis
352 was performed using the Wilcoxon signed-rank test, because it does not
353 assume normal distribution on the data and the outliers have less effect
354 than on a classical t -test (Derrac *et al.*, 2011).

355 Two measure groups were used to carry out the evaluations. The first
356 one includes measures that evaluate the algorithm performance such as:
357 N_G , the number of generations employed to find a solution; N_g , number of
358 generations required for all the population to be initialized with the same
359 compound, and F_S , the number of runs where compounds in D are linked
360 for a metabolic pathway. The second group corresponds to measures that
361 evaluate characteristics of the metabolic pathways found. The measures
362 employed to characterize the metabolic pathways are:

- 363 • Reactions (N_R): It provides information about the number of steps
 364 required to relate the compounds, counting the number of reactions
 365 the metabolic pathway has.
- 366 • Branching (ρ): It evaluates the relation between the pathway reac-
 367 tions by measuring the mean number of reactions that employ a non-
 368 abundant substrate. Compounds belonging to A' (abundant com-
 369 pounds) are not considered to calculate this measure, because the main
 370 interest is in the relationships among new compounds produced in the
 371 network. The pathway branching is calculated according to

$$\rho(\mathbf{c}) = \frac{1}{|S_f^*|} \sum_{i=1}^{|S_f^*|} \sum_{j=1}^{|\mathbf{c}|} \mathbf{1}_{s_i \subseteq S(r_j)}, \quad (6)$$

372 where S_f^* are the substrates of all reactions in \mathbf{c} after filtering the abun-
 373 dant compounds, $|\mathbf{c}|$ is the pathway size, $\mathbf{1}$ is the indicator function,
 374 and s_i is the i -th compound of the set $S(r_j)$ of substrates for reaction
 375 r_j .

- 376 • Leaves (λ): It counts the number of compounds produced by the
 377 metabolic pathway that are not employed as substrates by any re-
 378 action. This measure gives an idea of the degree of specificity the
 379 pathway has. A pathway with a high number of leaves indicates that
 380 it participates as an intermediary of a great variety of processes; a
 381 pathway with a low number of leaves indicates a high specificity for
 382 the synthesis of the indicated compounds. Let $S^*(\mathbf{c})$ and $P^*(\mathbf{c})$ be the
 383 sets of substrates and products of all reactions encoded in \mathbf{c} , respec-
 384 tively, the number of leaves λ is calculated as

$$\lambda(\mathbf{c}) = |P^*(\mathbf{c}) - (S^*(\mathbf{c}) \cup A')|. \quad (7)$$

- 385 • Difference between metabolic pathways (σ): This measure compares
 386 the sequence of compounds used to relate the elements in D and de-
 387 termines the proportion of compounds shared between two pathways.
 388 Let d_i and d_j be the initial substrates of the chromosomes \mathbf{c}_i and \mathbf{c}_j ,
 389 respectively, and let
 390 $Q_i = (P^*(\mathbf{c}_i) \cap S^*(\mathbf{c}_i) \cup \{d_i\}) - A'$ and $Q_j = (P^*(\mathbf{c}_j) \cap S^*(\mathbf{c}_j) \cup \{d_j\}) - A'$
 391 be the subsets of compounds belonging each pathway. The difference
 392 between both metabolic pathways is calculated as

$$\sigma(\mathbf{c}_i, \mathbf{c}_j) = 1 - \left[\frac{|Q_i \cap Q_j|}{\min(|Q_i|, |Q_j|)} \right]. \quad (8)$$

393 Two pathways have a difference $\sigma(\mathbf{c}_i, \mathbf{c}_j) = 0$ when they employ ex-
 394 actly the same compounds to relate the elements in D . This not im-
 395 plies that both are the same pathway, but rather one can be included
 396 in the other.

397 4. Results and discussion

398 In this section, the proposed algorithm performance is studied in three
 399 phases. The first one, presented in Section 4.1, studies the behavior of
 400 the algorithm for different parameters and operators. Section 4.2 analyzes
 401 the algorithm performance when the set of reactions previously employed
 402 is extended. Finally, Section 4.3 presents two case studies, where biological
 403 aspects of the solutions found are analyzed and discussed.

404 Experiments were conducted setting as finalization criteria a fitness equal
405 to 1.0 and a maximum of $G_M = 1000$ generations per search. Populations
406 were initialized with $M = 200$ individuals and a maximum size of chromo-
407 some $N_M = 100$ genes, to appropriately explore the solutions space. In
408 every case, the tournament selection strategy was employed with 5 individ-
409 uals and a crossover probability $p_x = 0.8$, since that value produced the best
410 results in preliminary experiments.

411 *4.1. Sensitivity to parameters and operators*

412 This section presents the performance measures for EvoMS. The effect
413 of the crossover type is analyzed and the influence of the different prob-
414 abilities that control the mutation operator is evaluated. In the experi-
415 ments, metabolic pathways relating the compounds L-Glutamate (C00025),
416 Fumarate (C00122), and L-Proline (C00763) were searched for. These par-
417 ticular compounds were selected given their importance in the metabolism,
418 and because only one (C00025) can be used to built a metabolic pathway
419 that links the three compounds. Thus, this situation allows to test the
420 method to determine the initial substrate. The search was carried out using
421 the set of reactions belonging to the arginine and proline reference metabolic
422 pathway (*apdata*)*. A total of 139 reactions were extracted, 24 of which are
423 reversible (broken down in 48 reactions) and 91 irreversible.

424 *4.1.1. Influence of the crossover type*

425 The EvoMS performance was compared using the standard one-point
426 crossover and the proposed crossover operator. The performance analysis
427 was evaluated in terms of the number of runs that produce a solution F_S ,

*<http://www.genome.jp/kegg/pathway/map/map00330.html>

Table 2: Effect of the crossover type on the evolutionary algorithm performance. The median and the median absolute deviation values are provided for N_G and N_g .

	one-point	modified
F_S	0.83	0.97
N_G	59 ± 27	57 ± 18
N_g	3 ± 0	3 ± 0

428 the number of generations required to find a solution N_G , and the number
429 of generations required to obtain a unique subpopulation N_g . Table 2 shows
430 the results obtained with each operator.

431 The most interesting fact observed in the table is the increase from 0.83
432 to 0.97 in the fraction of runs that lead to a solution when the modified oper-
433 ator is employed; there are not significant differences in the other measures.
434 This increase can be explained by the way in which metabolic pathways are
435 modelled. Since reactions are stored in the chromosome from left to right,
436 the ones located on the far right are more sensitive to the changes intro-
437 duced to the sequence, as they depend, to a greater extent, on the products
438 of previous reactions. On the other hand, since the algorithm requires all
439 reactions in the chromosome to be valid, incorporating a higher number
440 of reactions than the one needed to relate the compounds in D translates
441 into an additional effort the algorithm must make to meet this requirement.
442 Therefore, the insertion of only one portion of the genetic material from the
443 second parent decreases the number of reactions that do not probably meet
444 the validity requirement. At the same time, a lower number of generations
445 is required to find a solution.

Table 3: Generations required by EvoMS to find a metabolic pathway employing the initialization with a variable chromosome size. Results correspond to the median values and its deviations. † indicates experiments where a solution is found before 1000 generations and in more than 90% of the runs. The best results obtained with each mutation probability are highlighted in bold

N_G	$p_m = 0.02$			$p_m = 0.05$			$p_m = 0.08$		
	$p_e = 0.20$	$p_e = 0.50$	$p_e = 0.80$	$p_e = 0.20$	$p_e = 0.50$	$p_e = 0.80$	$p_e = 0.20$	$p_e = 0.50$	$p_e = 0.80$
$p_v = 0.20$	87 ± 37 [†]	72 ± 25	41 ± 10	164 ± 65	60 ± 19 [†]	36 ± 9 [†]	256 ± 130	69 ± 18	39 ± 13 [†]
$p_v = 0.50$	57 ± 18	56 ± 13 [†]	40 ± 11	102 ± 42	49 ± 10 [†]	35 ± 11	159 ± 80	70 ± 29 [†]	33 ± 8[†]
$p_v = 0.80$	72 ± 35 [†]	45 ± 6 [†]	41 ± 13	97 ± 48 [†]	47 ± 17	37 ± 12 [†]	162 ± 95 [†]	62 ± 28 [†]	36 ± 8 [†]

446 *4.1.2. Variation of mutation probabilities*

447 The proposed mutation operator plays an important role by introducing
 448 specific modifications that can change the branching of the metabolic path-
 449 way, and favour the exploration of new regions in the search space. Inserting
 450 new reactions can lead to the production of compounds necessary for other
 451 reactions to occur. Deletion allows to eliminate reactions that can be invalid
 452 or redundant. An appropriate balance of these operations can reduce the
 453 number of generations required to find the solution. To find the combination
 454 of probabilities leading to the best results, the values $p_m = \{0.02, 0.05, 0.08\}$,
 455 $p_e = \{0.20, 0.50, 0.80\}$ and $p_v = \{0.20, 0.50, 0.80\}$ were analyzed. Table 3
 456 shows the median and deviation values for the number of generations em-
 457 ployed in the runs for a specific set of parameters. The table has three
 458 blocks, each one corresponding to one mutation probability. For each block,
 459 valid insertion and deletion probabilities are shown in rows and columns,
 460 respectively. The combinations of probabilities with which solutions were
 461 obtained before 1000 generations and in more than 90% of the runs are

462 indicated with a mark ([†]).

463 In general terms, the combinations between deletion and valid insertion
464 probabilities lead to the same tendencies for the three mutation probabili-
465 ties evaluated. The increase in p_e is accompanied by the reduction in the
466 number of generations required to find a solution, as it is clearly seen when
467 $p_m = 0.05$ and $p_v = 0.5$, where there is a decrease from 102 to 35 gen-
468 erations when the value of p_e is increased. This is to be expected since,
469 during the initialization, a wide variety of reactions are incorporated, most
470 of which should be discarded during the evolution. Thus, the application of
471 mutations favoring the elimination of reactions will improve the algorithm
472 performance. Moreover, although no clear tendency is observed regarding
473 the effect produced by the valid insertion probability, in some cases, it is
474 seen that the increase in p_v is accompanied by a decrease in the number of
475 generations ($p_m = 0.02$ and $p_e = 0.5$).

476 As regards the mutation probability, it is possible to observe an increas-
477 ing tendency on cases in which a solution is obtained in more than 90% of
478 the runs with the raise of p_m . This trend might be explained considering
479 two effects produced by the mutation operator: increasing the genetic di-
480 versity and contributing to the consolidation of the validity of the sequence
481 of reactions. For that reason, there is an optimum number of insertions
482 that contributes to perform the search in the lowest number of generations.
483 Consequently, a low number of insertions makes the search slower, probably
484 because of the lack of genetic diversity; whereas an excess in the number of
485 insertions leads to the disproportionate increase of the pathways size and
486 makes it difficult to preserve the sequences validity. When the mutation
487 probability is low (few changes in the chromosome), the insertion of new
488 reactions has a more important contribution than deletion (low values for

489 p_e), probably collaborating to the generation of a sequence of valid reactions
490 and introducing genetic diversity. Nevertheless, when the mutation proba-
491 bility increases (a higher number of changes in the sequence), it is necessary
492 to increase the deletion of reactions in order to keep the balance between
493 insertions and the size of the pathways (containing unnecessary reactions).
494 In addition, it should be remembered that these results correspond to runs
495 in which the maximum number of generations is limited. Finally, it is ob-
496 served that the lowest number of generations employed with each mutation
497 probability (highlighted in bold) is obtained with $p_e = 0.80$ and $p_v = 0.50$,
498 being minimum for $p_m = 0.08$. Besides, this combination of probabilities
499 also provides solutions in 90% of the runs.

500 4.2. Scalability of the algorithm

501 In order to study the ability of the algorithm to perform similar searches
502 in spaces that scale in size, the search made in the previous section was
503 performed expanding the set of *apdata* reactions. The new dataset (*sdata*)
504 was built adding the reactions belonging to five reference metabolic path-
505 ways[†]. Thus, *sdata* has 443 one-way reactions, 132 of which are reversible
506 (broken down in 264 reactions) and 179 irreversible. Runs were carried out
507 employing the best parameters obtained in Section 4.1.

508 4.2.1. Algorithm performance and characteristics of the pathways

509 Table 4 shows the evaluation measures for the searches performed with
510 the two datasets. Blocks separate the performance measures (upper block)

[†]Glycolysis / Gluconeogenesis (map00010), Citrate cycle (map00020), Pentose phos-
phate pathway (map00030), Pentose and glucuronate interconversions (map00040) and
Alanine, aspartate and glutamate metabolism (map00250) in KEGG.

Table 4: Comparison of the algorithm performance employing the arginine and proline dataset (*apdata*) and extended dataset (*sdata*).

	<i>apdata</i>	<i>sdata</i>
F_S	1.00	0.97
N_G	33±8	29±7
N_g	3±0	4±1
N_R	8±1	6±1
ρ	1.2±0.1	1.3±0.1
λ	5±1	4±1

511 from the solutions quality measures (lower block). In general terms, no prac-
512 tical differences are observed in the algorithm performance. In both cases,
513 a solution is obtained in more than 90% of the runs ($F_S > 0.9$). More-
514 over, although the number of generations is lower when *sdata* is employed,
515 this behavior is only at a tendency level since the confidence intervals are
516 overlaped. Although the value of N_g is increased in one generation, from a
517 practical point of view this difference is not important, as in both cases the
518 winning subpopulation is quickly selected during the first generations.

519 Regarding the measures associated to the structure of the metabolic
520 pathways, a significant reduction is observed ($p < 0.0001$) in the size of the
521 pathways (N_R) found using *sdata*. This is to be expected since the number
522 of possible connections between compounds is higher and makes possible
523 the existence of smaller alternative paths that connect the compounds in D .
524 The branching ρ calculated for the solutions found with *sdata* supports this
525 explanation, as it experiences a significant increase ($p < 0.005$) from 1.2 to
526 1.3. The number of leaves λ , indicating that the pathways found with *sdata*

Table 5: Values of the difference between groups of equivalent solutions found with *apdata* and *sdata*. Difference values lower than 0.15 are highlighted in bold.

		<i>sdata</i>						
		Ib	IIb	IIIb	IVb	Vb	VIb	VIIb
	Ia	0.50	0.43	0.38	0.29	0.20	0.13	0.00
	IIa	0.50	0.43	0.25	0.29	0.20	0.22	0.00
<i>apdata</i>	IIIa	0.38	0.43	0.38	0.29	0.20	0.22	0.00
	IVa	0.25	0.29	0.38	0.29	0.00	0.22	0.00
	Va	0.25	0.29	0.38	0.29	0.00	0.11	0.00

527 include reactions that generate a lower number of unnecessary products
528 ($p < 0.0001$), is possibly due to the use of more specific process reactions.
529 It should be highlighted that, regardless of the branching differences, both
530 sets of reactions lead to solutions with values of ρ higher than the unit, since
531 some compounds in the networks found act as a substrate in more than one
532 reaction.

533 4.2.2. Difference between solutions

534 In order to measure if the proposed evolutionary algorithm is capable of
535 reproducing the searches in a solutions space extended by the incorporation
536 of additional reactions, the solutions found with both datasets were studied
537 and compared to determine the number of novel metabolic pathways in
538 common.

539 Typically, synthesizing a compound implies a number of steps until the
540 desired product is reached. Thus, a sequence of several intermediate com-
541 pounds linking the initial substrate and the final product is generated. How-
542 ever, in many cases, those intermediate compounds can be produced by
543 more than one reaction. This leads to metabolic pathways which are differ-
544 ent in terms of reactions, but equivalent in terms of the sequence of com-
545 pounds needed to perform the synthesis. According to (8), two metabolic
546 pathways \mathbf{c}_1 and \mathbf{c}_2 will be equivalent when they have a difference value
547 $\sigma(\mathbf{c}_1, \mathbf{c}_2) = 0.0$. Furthermore, this measure will increase when the number
548 of shared compounds decreases.

549 In a preliminary analysis, five *groups of equivalent solutions* were found
550 for *apdata* and seven for *sdata*. Table 5 shows the difference values between
551 the groups found with both sets of reactions. Rows and columns indicate
552 the group of equivalent solutions for *apdata* and *sdata*, respectively. The in-
553 tersection between a row and a column indicates the difference between the
554 groups considered. It can be seen that some groups of solutions are equiva-
555 lent, as it could be expected, since *apdata* and *sdata* share the mechanisms
556 to synthesize the three specified compounds. For instance, group VIIb does
557 not show differences with any of the solutions found with *apdata*. This is be-
558 cause the five groups of equal solutions found with *apdata* employ the same
559 sequence of compounds that the group VIIb, together with other additional
560 compounds. The group of solutions IVb also shows a similar behavior, pre-
561 senting a difference of 0.29 with all *apdata* solutions, which indicates that it
562 shares a portion of the sequence of compounds.

563 In order to analyze the differences found in more detail, two groups of
564 solutions with extreme difference values were selected and one metabolic
565 pathway representing each one was plotted. Figure 3 shows the metabolic

566 pathways corresponding to **VIa** and **Vb** groups, while Figure 4 contains the
567 pathways of **Ia** and **Ib** groups. In every case, the pathways must be inter-
568 preted in a descending manner, starting by the initial substrate C00025 (in
569 red), and descending through the sequence of reactions and until each one of
570 the final products (in yellow). Representations are simplified, not showing
571 abundant compounds.

572 Pathways representing solutions of groups **IVa** and **Vb** in Figure 3 do
573 not show any difference according to (8). Clearly, pathway from **IVa** em-
574 ploys almost twice the reactions as **Vb** to relate the same compounds. How-
575 ever, analyzing in detail the sequences of compounds used by both path-
576 ways, it is observed that the compounds used by **Vb** ($Q_{Vb} = \{C00025,$
577 $C03912, C00148, C00763, C00049, C00122\}$) are also employed by **IVa**
578 ($Q_{IVa} = \{C00025, C01165, C03912, C00148, C00763, C00169, C00077,$
579 $C00327, C03406, C00122\}$). Although the compound C00049 (which in
580 the sequence is indicated in italics) is not shared by the pathways, it should
581 not be considered to calculate the differences since it is part of the set of
582 abundant compounds. As a consequence, both solutions relate members of
583 D employing the same compounds.

584 When analyzing the solutions from **Ia** and **Ib** (Figure 4) it can be seen
585 that both contain the same number of reactions. However, the sequence of
586 compounds used by **Ia** ($Q_{Ia} = \{C00025, C00077, C00148, C00763, C00169,$
587 $C00327, C03406, C00122\}$) presents a large difference compared to the
588 one employed by **Ib** ($Q_{Ib} = \{C00025, C01165, C03912, C00148, C00763,$
589 $C00026, C00091, C00042, C00122\}$). On the one hand, the sequences of
590 compounds used to synthesise C00763 from C00025 only share the inter-
591 mediary compound C00148, which is produced through different reactions
592 in each solution. On the other hand, the C00122 synthesis is carried out

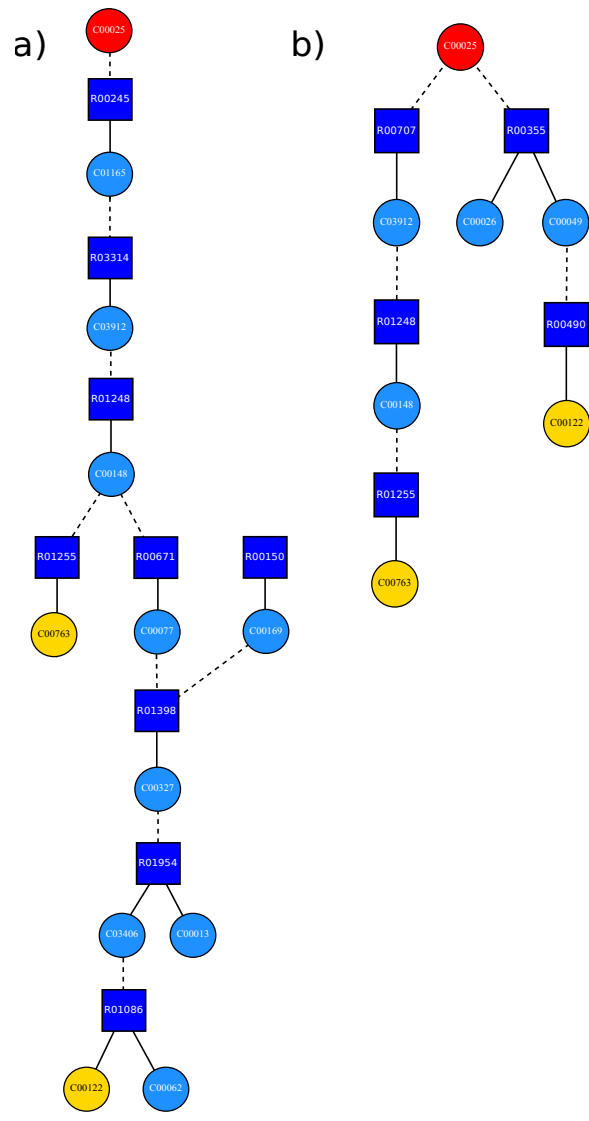


Figure 3: Pathways belonging to two groups of solutions found with *apdata* and *sdata*, respectively. a) Examples for: IVa, b) Vb. The initial compound is indicated in red (C00025), the compounds to be produced are indicated in yellow (C00122, C00763), and the compounds produced by the metabolic pathway are indicated in light blue. Reactions are indicated as blue squares. Available compounds are not included in the metabolic pathway.

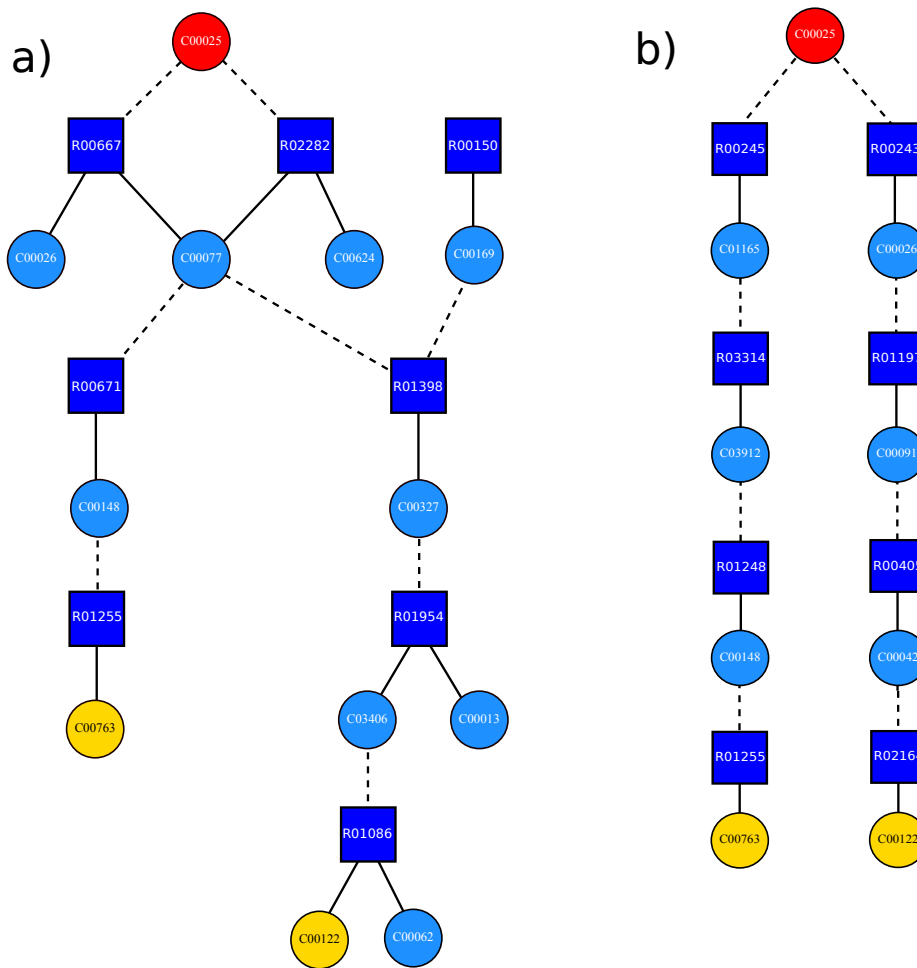


Figure 4: Pathways belonging to two groups of solutions found with *apdata* and *sdata*, respectively. a) Examples for: a) Ia, b) Ib. The initial compound is indicated in red (C00025), the compounds to be produced are indicated in yellow (C00122, C00763), and the compounds produced by the metabolic pathway are indicated in light blue. Reactions are indicated as blue squares. Available compounds are not included in the metabolic pathway.

593 using sequences of completely different compounds in both metabolic path-
594 ways, sharing only the compounds in the extremes. Thus, only 4 compounds
595 (C00148 and the three compounds to be related) are shared between Ia and
596 Ib, leading to a difference $\sigma = 0.5$ according to (8).

597 4.3. Case study: searching for relations between 4 compounds

598 The EvoMS performance in a more complex real problem was evalu-
599 ated and compared with a state-of-the-art algorithm (Faust *et al.*, 2011) for
600 searching a metabolic pathway relating 4 specific compounds. The search in-
601 volved the complete set of reactions stored in KEGG for the *Escherichia coli*
602 bacterium metabolism. After the pre-processing, the search space was made
603 up of 2354 reactions, 1061 of which were reversible (broken down in 2122
604 one way reactions) and 232 irreversible. The reference pathway for lysine,
605 threonine, and methionine biosynthesis (Figure 5) was taken as a case study
606 of a branched metabolic pathway. It synthesizes compounds C00047 (L-
607 Lysine), C00073 (L-Methionine), and C00188 (L-Threonine) from C00036
608 (Oxaloacetate).

609 The algorithm of Faust *et al.* (2011) combines several linear paths to
610 build a network of relationships among compounds. It performs the search
611 of the shortest path between each pair of compounds and combine all of
612 them into a network. With this approach, the authors were able to find a
613 pathway for the compounds using a high proportion (85%) of the reactions
614 belonging to the reference metabolic pathway. In comparison, EvoMS was
615 able to find a pathway with all the reactions (100%) of the reference pathway.
616 Furthermore, another important advantage is that feasibility of the solutions
617 found by EvoMS is guaranteed. EvoMS builds the pathway by verifying that
618 all reactions use available substrates, while the other algorithm does not even

619 takes into account that information during the search.

620 Besides the reference pathway, Figure 6 shows two other examples of
621 metabolic pathways synthesized by EvoMS, for the same search. In both
622 cases, solutions were fully feasible and allowed to relate the same 4 com-
623 pounds. Figure 6a shows the metabolic pathway found with C00036 as
624 initial substrate, containing four reactions also present in the reference path-
625 way (R00355, R03260, R01286 and R00946). It must be noted that reaction
626 R03260 plays a central function in the new pathway, producing two key
627 compounds (C01118 and C00097) needed to synthesize C00073 and C00188.
628 Also, it can be appreciated that the initial substrate has a key role in this
629 pathway, being a precursor to synthesize C00027 (needed for C00047), and
630 C00042 (needed for C00073 and C00188). Furthermore, the large number
631 of interconnections among reactions in this pathway shows an important
632 collaborative work to synthesize the final products.

633 Figure 6b presents another alternative metabolic pathway that is also
634 fully feasible and relates the same compounds. At a glance, it can be ob-
635 served that this novel pathway could be more efficient to link the 4 com-
636 pounds of interest than the previous one, because it requires a lower number
637 of reactions to relate them. This solution uses C00073 as initial substrate,
638 not sharing any reaction with the reference pathway. The novel pathway is
639 built by two main branches starting from C00073, one of which produces
640 C00047 and the other produces the remaining two products. As it can
641 be seen, C00022 plays a key role as precursor in the synthesis of C00036
642 and C00188. Similarly to the pathway in Figure 6a, C00022 in this novel
643 pathway allows to infer a relation between the glycolysis (a reference path-
644 way for many life forms) and the synthesis of both products. These exam-
645 ples evidence the natural interconnections present among metabolic path-

646 ways in living organisms. This also highlights the importance of developing
647 new algorithms for searching on large sets of reactions, providing branched
648 metabolic pathways of feasible reactions that relate multiple compounds
649 simultaneously.

650 **5. Conclusions and future work**

651 This work approached the problem of searching for metabolic pathways
652 that relate a set of compounds through networks of feasible reactions. A
653 model to build the pathways based on a set of compounds was proposed and
654 a new evolutionary algorithm, called EvoMS was developed to search for the
655 reactions required to build pathways between specific compounds. Also, new
656 operators and an initialization strategy that employ the set of compounds
657 model were developed. The fitness function was designed to evaluate essen-
658 tial characteristics required in the metabolic pathways searched, in order to
659 find feasible metabolic pathways. The tests carried out for a real problem
660 showed that EvoMS was capable of reproducing known metabolic pathways
661 and also finding alternative connections to synthesize the same final com-
662 pounds. In all searches, the algorithm found branched metabolic pathways
663 made up of feasible reactions from the initial compounds indicated. Besides,
664 in cases where reactions require compounds that do not belong to the abun-
665 dant ones, the algorithm was capable of previously incorporating reactions
666 to generate them. In summary, the possibility of generating a wide range of
667 connections between compounds, together with the ability to provide feasi-
668 ble solutions makes EvoMS a simple and powerful method to find feasible
669 networks connecting metabolic compounds. Moreover, flexibility of the eval-
670 uation function allows to easily extend it to incorporate new objectives to

671 optimize in the solution.

672 Future work will aim to improve the search process by adding informa-
673 tion to the evaluation function, for example, regarding the stoichiometry and
674 thermodynamics of the reactions, the degree of connectivity of compounds,
675 and/or the availability of enzymes. In addition, the crossover operator will
676 be modified to employ information of the compounds used by the metabolic
677 pathway, and mechanisms to automatically adjust the parameters of the
678 algorithm during the evolution will be studied.

679 The full source code for EvoMS algorithm is available for free aca-
680 demic use at <http://sourceforge.net/projects/sourcesinc/files/evoms/>. A
681 web-interface to run the evolutionary algorithm proposed in this work is
682 available online at <http://fich.unl.edu.ar/sinc/web-demo/evoms/>, whose main
683 inputs, outputs, features and ways of use are explained in (Gerard *et al.*,
684 2015).

685 **Acknowledgements**

686 This work was supported by National Scientific and Technical Research
687 Council (CONICET) [PIP 2013-2015 117], Universidad Nacional del Litoral
688 (UNL) [CAI+D 2011 548] and Agencia Nacional de Promoción Científica y
689 Tecnológica (ANPCyT) [PICT 2014 2627].

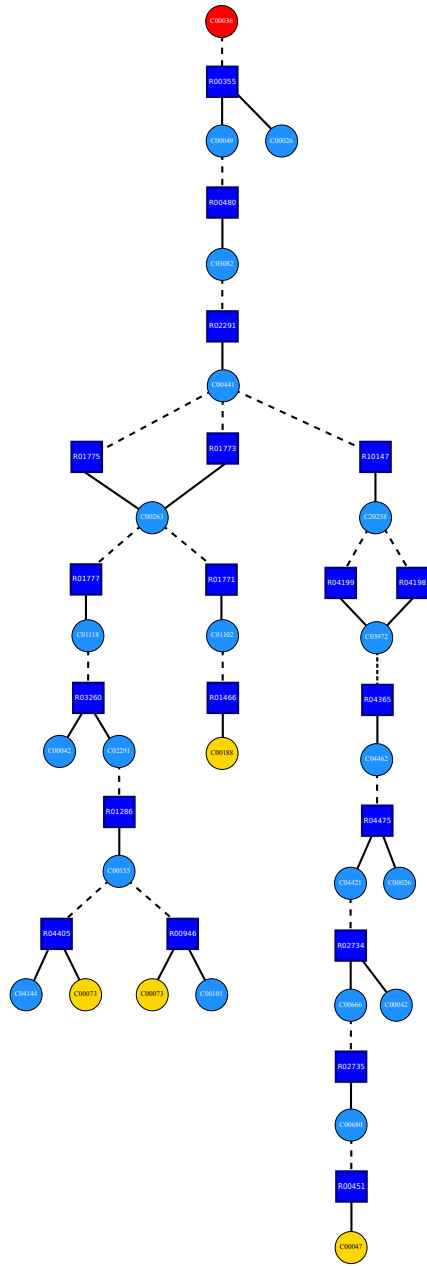


Figure 5: Reference metabolic pathway involving lysine, threonine and methionine biosynthesis. Note that reaction R00946 and R04405 produce the same compound C00073 in two different ways. Initial substrate is in red and the compounds to be produced are indicated in yellow. Reactions are indicated in blue, and their substrates and products are in dashed and solid lines, respectively. To provide a clearest view, only the more relevant compounds are shown.

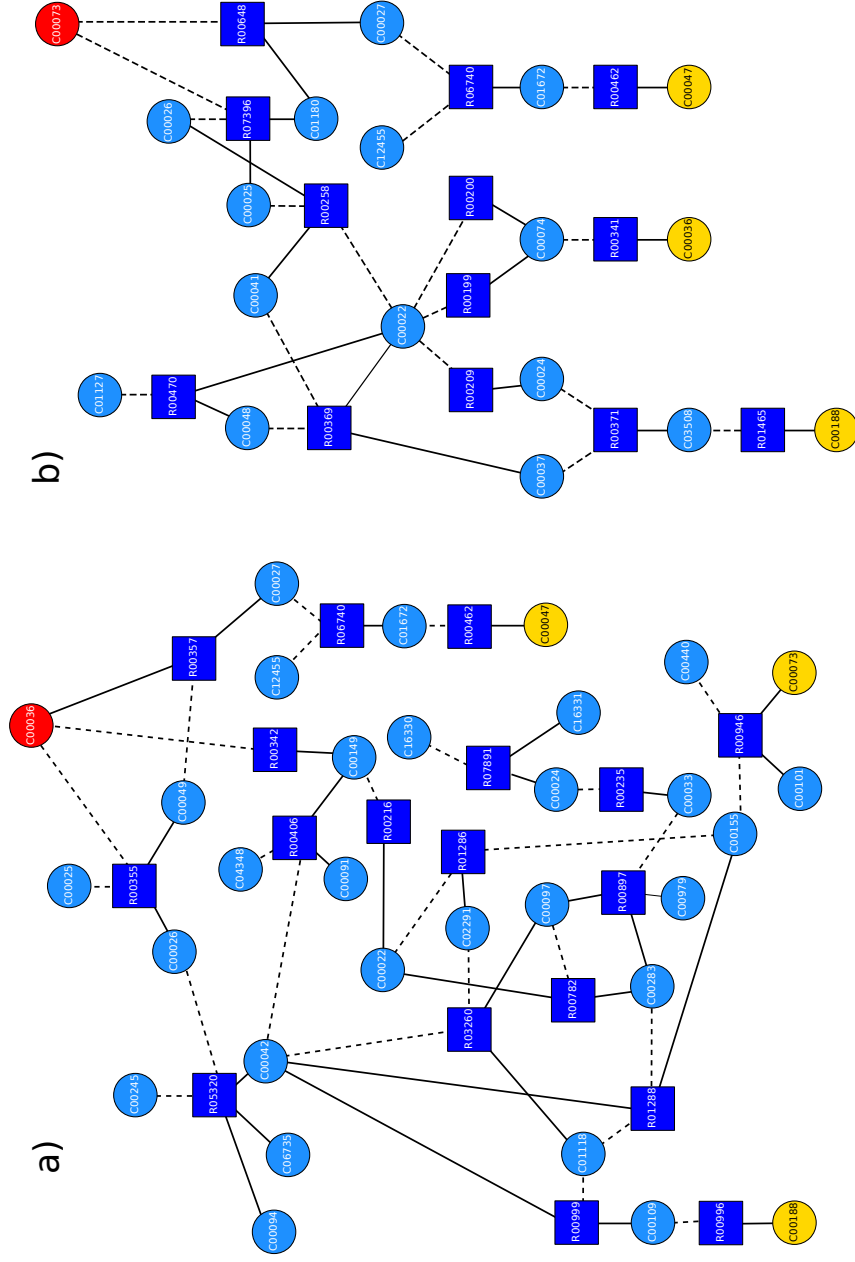


Figure 6: New metabolic pathways linking compounds C00036, C00047, C00073 and C00188. a) Metabolic pathway found by EvoMS with C00036 as initial substrate. b) Metabolic pathway found by EvoMS with C00073 as initial substrate. In every case, the initial substrate is indicated in red and the compounds to be produced are indicated in yellow. Reactions are indicated in blue. Substrates are connected with dashed lines and products with solid ones. To provide a clearest view, only the more relevant compounds are shown.

690 **References**

- 691 Affenzeller, M., Winkler, S., Wagner, S., and Beham, A. (2009). *Genetic Algorithms and Genetic*
692 *Programming: Modern Concepts and Practical Applications*. CRC Press.
- 693 Altman, T., Travers, M., Kothari, A., Caspi, R., and Karp, P. (2013). A systematic comparison of the
694 MetaCyc and KEGG pathway databases. *BMC Bioinformatics*, **14**, 112.
- 695 Arita, M. (2012). *Bacterial Molecular Networks : Methods and Protocols*, volume 804 of *Methods*
696 *in Molecular Biology*, chapter From Metabolic Reactions to Networks and Pathways, pages 93–106.
697 Springer.
- 698 Bäck, T., Fogel, D., and Michalewicz, Z. (2000). *Evolutionary Computation I: Basic Algorithms and*
699 *Operators*. Institute of Physics Publishing.
- 700 Bordbar, A., Monk, J. M., King, Z. A., and Palsson, B. O. (2014). Constraint-based models predict
701 metabolic and associated cellular functions. *Nature Reviews Genetics*, **15**, 107–120.
- 702 Boussaïd, I., Lepagnot, J., and Siarry, P. (2013). A survey on optimization metaheuristics. *Information*
703 *Sciences*, **237**, 82–117.
- 704 Chen, C. P. and Zhang, C.-Y. (2014). Data-intensive applications, challenges, techniques and technolo-
705 gies: A survey on Big Data. *Information Sciences*, **275**, 314–347.
- 706 Croes, D., Couche, F., Wodak, S., and van Helden, J. (2005). Metabolic Pathfinding: inferring relevant
707 pathways in biochemical networks. *Nucleic Acids Research*, **33**, W326–W330.
- 708 de Magalhães, C. S., Almeida, D. M., Barbosa, H. J. C., and Dardenne, L. E. (2014). A dynamic niching
709 genetic algorithm strategy for docking highly flexible ligands. *Information Sciences*, **289**, 206–224.
- 710 Derrac, J., García, S., Molina, D., and Herrera, F. (2011). A practical tutorial on the use of nonpara-
711 metric statistical tests as a methology for comparing evolutionary and swarm intelligence algorithms.
712 *Swarm and Evolutionary Computation*, **1**, 3–18.
- 713 Faust, K., Croes, D., and van Helden, J. (2009). Metabolic Pathfinding Using RPAIR Annotation.
714 *Journal of Molecular Biology*, **388**(2), 390–414.
- 715 Faust, K., Dupont, P., Callut, J., and van Helden, J. (2010). Pathway discovery in metabolic networks
716 by subgraph extraction. *Bioinformatics*, **26**, 1211–1218.
- 717 Faust, K., Croes, D., and van Helden, J. (2011). Prediction of metabolic pathways from genome-scale
718 metabolic networks. *BioSystems*, **105**, 109–121.
- 719 Garai, G. and Chowdhury, B. (2015). A cascaded pairwise biomolecular sequence alignment technique
720 using evolutionary algorithm. *Information Sciences*, **297**, 118–139.

- 721 Gerard, M., Stegmayer, G., and Milone, D. (2013). An evolutionary approach for searching metabolic
722 pathways. *Computers in Biology and Medicine*, **43**, 1704–1712.
- 723 Gerard, M., Stegmayer, G., and Milone, D. (2015). EvoMS: an evolutionary tool to find de novo
724 metabolic pathways. *BioSystems*, **134**, 43–47.
- 725 Gerlee, P., Lizana, L., and Sneppen, K. (2009). Pathway identification by network pruning in the
726 metabolic network of *Escherichia coli*. *Bioinformatics*, **25**, 3282–3288.
- 727 Goto, S., Okuno, Y., Hattori, M., Nishioka, T., and Kanehisa, M. (2002). LIGAND: database of chemical
728 compounds and reactions in biological pathways. *Nucleic Acids Research*, **30**, 402–404.
- 729 Heath, A., Bennett, G., and Kavraki, L. (2010). Finding metabolic pathways using atom tracking.
730 *Systems Biology*, **26**, 1548–1555.
- 731 Jeong, H., Tombor, B., Albert, R., Oltvai, Z., and Barabási, A. (2000). The large-scale organization of
732 metabolic networks. *Nature*, **407**, 651–654.
- 733 Kayaa, H. and Şule Gündüz-Öğüdücü (2013). SAGA: A novel signal alignment method based on genetic
734 algorithm. *Information Sciences*, **228**, 113–130.
- 735 Lacroix, V., Cottret, L., Thebault, P., and Sagot, M.-F. (2008). An Introduction to Metabolic Net-
736 works and Their Structural Analysis. *IEEE/ACM Transactions on Computational Biology and*
737 *Bioinformatics*, **5**(4), 594–617.
- 738 Lee, S. Y., Kim, H. U., Park, J. H., Park, J. M., and Kim, T. Y. (2009). Metabolic engineering of
739 microorganisms: general strategies and drug production. *Drug Discovery Today*, **14**, 78–88.
- 740 Lee, W.-P. and Hsiao, Y.-T. (2012). Inferring gene regulatory networks using a hybrid GA–PSO approach
741 with numerical constraints and network decomposition. *Information Sciences*, **188**, 80–99.
- 742 McShan, D., Rao, S., and Shah, I. (2003). PathMiner: predicting metabolic pathways by heuristic
743 search. *Bioinformatics*, **19**, 1692–1698.
- 744 Ogata, H., Goto, S., Fujibuchi, W., and Kanehisa, M. (1998). Computation with the KEGG pathway
745 database. *BioSystems*, **47**, 119–128.
- 746 Pal, S., Bandyopadhyay, S., and Ray, S. (2006). Evolutionary Computation in Bioinformatics: A Review.
747 *IEEE Transactions on Systems Man and Cybernetics*, **36**, 601–615.
- 748 Pey, J., Prada, J., Beasley, J., and Planes, F. (2011). Path finding methods accounting for stoichiometry
749 in metabolic networks. *Genome Biology*, **12**, R49.
- 750 Rahman, S., Advani, P., Schunk, R., Schrader, R., and Schomburg, D. (2005). Metabolic pathway
751 analysis web service (Pathway Hunter Tool at CUBIC). *Bioinformatics*, **21**, 1189–1193.

- 752 Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N., and Barabási, A.-L. (2002). Hierarchical
753 Organization of Modularity in Metabolic Networks. *Science*, **297**, 1551–1555.
- 754 Russell, S. and Norvig, P. (2010). *Artificial Intelligence: A Modern Approach (3 Ed.)*. Prentice Hall.
- 755 Tenazinha, N. and Vinga, S. (2011). A Survey on Methods for Modeling and Analyzing Integrated
756 Biological Networks. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, **8**,
757 943–958.
- 758 Yim, H., Haselbeck, R., Niu, W., Pujol-Baxley, C., Burgard, A., Boldt, J., Khandurina, J., Trawick,
759 J. D., Osterhout, R. E., Stephen, R., Estadilla, J., Teisan, S., Schreyer, H. B., Andrae, S., Yang,
760 T. H., Lee, S. Y., Burk, M. J., and Dien, S. V. (2011). Metabolic engineering of *Escherichia coli* for
761 direct production of 1,4-butanediol. *Nature Chemical Biology*, **7**, 445–452.