

1 A review of methods for the reconstruction and analysis of integrated
2 genome-scale models of metabolism and regulation

3

4 Fernando Cruz¹, José P. Faria², Miguel Rocha¹, Isabel Rocha³, Oscar Dias^{1*}

5

6 ¹Centre of Biological Engineering, University of Minho, Braga 4710-057, Portugal

7 ²Data Science and Learning Division, Argonne National Laboratory, Argonne, IL, USA

8 ³Instituto de Tecnologia Química e Biológica António Xavier, Universidade Nova de
9 Lisboa, 2780-157 Oeiras, Portugal

10 * - Correspondence: odias@ceb.uminho.pt

11 (3243/2500-3500)

12 Abstract (89/250)

13 The current survey aims to describe the main methodologies for extending the
14 reconstruction and analysis of genome-scale metabolic models and phenotype
15 simulation with Flux Balance Analysis mathematical frameworks, via the integration of
16 Transcriptional Regulatory Networks and/or gene expression data. Although the
17 surveyed methods are aimed at improving phenotype simulations obtained from these
18 models, the perspective of reconstructing integrated genome-scale models of
19 metabolism and gene expression for diverse prokaryotes is still an open challenge.

20

21 Introduction

22 High-throughput large-scale omics experiments are nowadays disseminated in
23 biochemical research, supporting the study of the genomics, transcriptomics, and
24 metabolomics layers of the cellular's molecular machinery. Currently, the volume of
25 studies in the different omics fields has provided means for systems biology to thrive
26 (1). This interdisciplinary field proposes differentiated approaches, such as the
27 reconstruction of *in silico* networks and models, to provide quantitative and qualitative
28 descriptions of biological systems as a whole.

29

30 Reconstruction of GSMMs

31 Nowadays, the generation of Genome-Scale Metabolic Models (GSMMs) is a
32 common practice in systems biology. The reconstruction of these comprehensive
33 models, through modeling techniques and genomics data, allows predicting cells'
34 metabolic behavior (2–4).

35 A GSMM is an *in silico* representation of the biochemical reactions taking place within
36 a the metabolism of a given organism (5). A genome-wide functional annotation that
37 provides the required metabolic information over the organism of interest should be
38 performed to assemble this representation. This information is linked to existing
39 metabolic knowledge retrieved essentially from biochemical databases and literature.
40 These steps help to create the reaction set, upon which the metabolic network is
41 assembled.

42 The link from metabolic genes to proteins (mainly enzymes or membrane transporter
43 proteins), as well as from proteins to reactions, is established by Gene-Protein-Reaction
44 (GPR) associations. GPR associations must be cautiously defined during the
45 reconstruction, taking into account isoenzymes, protein complexes and cascade
46 reactions (3), through the use of *AND* or *OR* Boolean rules.

47 In the next iteration, biomass and organism-specific constraints are formulated from
48 the retrieved knowledge to assemble a final stoichiometric model. The final GSMM may
49 then be exported in a standard format, such as the Systems Biology Markup Language
50 (SBML) (6). Several platforms, such as *merlin* (7), ModelSEED (8), RAVEN (9) and
51 CarveMe (10), have been developed specifically for performing or assisting the
52 reconstruction of these models (11).

53 The classic principles of chemical engineering are used to infer the dynamic mass
54 balances of all metabolites in the metabolic network. A single ordinary differential
55 equation (ODE) is created for each metabolite, accounting for its stoichiometry in the
56 whole reaction set. Due to the lack of kinetic rates for all reactions in the ODE set, a
57 steady-state approximation is used to reduce the mass balances to a set of linear
58 equations. In a pseudo-steady-state paradigm, the concentration of a metabolite is
59 assumed to remain constant throughout time (4).

60 When used to determine flux values, the set of linear equations defines a linear
61 system, typically underdetermined, as the number of fluxes is much higher than the
62 number of mass balance constraints, also referred to as the null space of S (12).
63 Additional mass balance constraints can be added to the system to limit the flux that
64 each reaction can accommodate by the imposition of both lower and upper bounds.

65 The system can be solved mathematically transforming it into an optimization
66 problem, using several constraint-based approaches to predict the phenotypic behavior
67 of the organism on a wide variety of environmental and genetic conditions. One of the
68 most popular approaches is the Flux Balance Analysis (FBA) framework (12). FBA can
69 compute an optimal solution, out of the feasible space determined by both mass
70 balance and flux constraints using linear programming (LP). FBA requires the definition
71 of an objective function, which should be relevant to the underlying problem, which is
72 commonly defined as the maximization or minimization of a specific metabolic flux (e.g.,
73 biomass reaction), and quantitatively determines how much each reaction contributes
74 to a phenotype (4).

75 Parsimonious Flux Balance Analysis (pFBA) (13) and Flux Variability Analysis (FVA)
76 (14), are alternative mathematical frameworks that also employ LP to allow analyzing
77 *in silico* flux distributions. This set of tools is extremely helpful for validating a
78 reconstructed model using experimental data of the organisms of interest. COntstraint-
79 Based Reconstruction and Analysis (COBRA) Toolbox (15), COBRAPy (16), OptFlux (17),
80 and ReFramed (<https://github.com/cdanielmachado/reframed>) are prominent
81 computational tools that have implemented these methods.

82 Although GSMMs have proven to be valuable throughout the years (18–24), there
83 are limitations. Indeed, they are not yet capable of accounting for biological regulatory
84 phenomena, such as the control of gene expression (25). The lack of this additional layer
85 of information in these models can lead to erroneous *in silico* phenotype simulations,
86 due to the lack of constraints that allow reaching the most accurate flux distribution
87 according to experimental data.

88 Several methods have been proposed to improve phenotype simulations obtained
89 from GSMMs, which will be herein surveyed. Most of these new methodologies are
90 aimed at combining additional layers of omics data, namely transcriptomics, to limit the
91 cone of allowable flux distributions. Also, these methods often resort to the integration
92 of gene expression data and/or regulatory information obtained from Transcriptional
93 Regulatory Networks (TRN)s being, therefore, prominent efforts made towards the
94 reconstruction of integrated genome-scale models of metabolism and gene expression.
95 The utilization of these integrated models can be useful to improve phenotype
96 simulations or extend the analysis of regular GSMMs.

97

98 Reconstruction of TRNs

99 A TRN can be represented as a bipartite graph that comprehends vertices and edges.
100 Vertices are usually regulators and target genes, whereas edges determine how these
101 regulatory elements are connected and interact with each other, often under a causal
102 relationship.

103 Inferring TRNs is fundamentally an underdetermined problem associated with a large
104 search space where many solutions explain the data equally well (25,26). High-quality
105 transcriptional information is scarce, in databases or literature and focused on a few
106 well-studied organisms. Thus, the number of potential regulatory interactions between
107 a Transcription Factor (TF) and target genes is considerably larger than the actual true
108 biological interactions.

109 Although methods for reconstructing TRNs have been extensively reviewed in the
110 literature (25–31), there are a panoply of classification systems and procedures.
111 Moreover, as new methods are released each year, the complexity increases. Hence,
112 assigning classes to these new approaches can be a complicated task. More importantly,
113 this reveals that standard platforms and methodologies to assemble TRNs using diverse

114 sources of regulatory information, such as gene expression data (32,33) or
115 transcriptional information (34–36), are still missing.

116 Nevertheless, an integrative workflow for reconstructing bacterial TRNs has been
117 proposed by Faria *et al.* 2014 (25). The authors suggested that comparative-genomics
118 approaches, namely the inference of TRNs using template curated networks and the
119 prediction of cis-regulatory elements, can be integrated with the output of *de novo*
120 reverse engineering tools. This workflow addresses the possibility of reconstructing
121 TRNs for less described prokaryotic organisms using a variety of sources of regulatory
122 information.

123 Template-network methodologies are based on the conservation of prokaryotic TRNs
124 across evolution (37–39). As described by Faria *et al.* 2014 (25), template-network-based
125 methods usually perform a search for orthologous genes in the genome of the organism
126 of interest to propagate TRNs to strains of a well-characterized model organisms or
127 closely related ones.

128 *Cis*-regulatory elements detection rely on the assumption that a regulatory
129 interaction between a given TF and target gene can be inferred from the detection of
130 the Transcription Factor Binding Sites (TFBS). The prediction of these *cis*-regulatory
131 elements is a problem in which computational methods can assist (27,40). Although
132 these computational tools are unable to infer a complete TRN from TFBS data, they can
133 be integrated in the following workflow towards such a goal. The principles of this
134 methodology were implemented by Alkema *et al.* 2004 (41) in Regologger and the
135 RegPredict web-based platform (42).

136 *De novo* reverse engineering tools are widely used for inferring TRNs from gene
137 expression data. Indeed, a vast repertoire of computational tools based on the *de novo*
138 reverse engineering approach can be found in the literature, and consequently
139 numerous ways to classify these tools (28–30). Nevertheless, *de novo* reverse
140 engineering methods are usually classified by mathematical formulation. Hence, to
141 highlight the most used mathematical formulations, data-driven methods are usually
142 based on the following:

- 143 • Correlation (e.g. COREGNET (43)),
- 144 • Information-theoretic (e.g. (44)),
- 145 • Boolean algebra (making use of the widely known binary operators *AND*, *OR*,
146 and *NOT* to describe regulatory interactions (45), e.g. ModEnt (46)),
- 147 • Regression-based (e.g. GENIE3 (47)),
- 148 • ODEs (e.g. Inferelator (48)),
- 149 • Bayesian models (e.g. Gat-Viks *et al* 2007 (49)),

150 Available state-of-the-art TRNs' reconstructions for prokaryotic organisms include
151 well-known prokaryotic organisms such as *Escherichia coli* (50,51), *Bacillus subtilis*
152 (52,53) and *Mycobacterium tuberculosis* (54), having hundreds of regulators and
153 thousands of target genes. These TRNs can, therefore, be used as gold-standards in the
154 template-network-based approach or supervised methods. Interestingly, other TRNs for

155 prokaryotic organisms less described in the literature or having less amount of gene
156 expression data (44,55–61), are also available, though in less number.

157 The TRNs available in the literature are usually the result of a specific gene expression
158 data-driven analysis or the collection of regulatory information from literature and
159 databases. Although these TRNs can be used in the comparative genomics or data-
160 driven approaches, not all of them can be easily integrated in GSMMs, as only genome-
161 scale TRNs are actually useful for the integration and simulation of integrated models.

162

163 Integrated models

164 Combining regulatory elements with information on metabolic stoichiometry is a
165 complex task. There are many ways for controlling metabolism (62), which are well
166 represented in the large diversity of methods proposed to quantify such influence
167 (25,63–75). Nevertheless, the common denominator is that most methods start with
168 GSMMs.

169 In detail, several of these methods integrate complete functional TRNs (76–84) or
170 gene expression data (85–95) into GSMMs, whereas others impose additional
171 constraints using information on allosteric and post-translational modifications
172 (66,67,73). A different strategy is the combination of multiple layers of regulation
173 (63,65,72,74). For higher eukaryotes such as humans, the control of gene expression
174 also plays an essential role in the differentiation between different tissues or cell-types.
175 Thus, algorithms for tailoring a GSMM according to a specific cell-line or tissue,
176 commonly referred to as context-specific models, have been proposed (75,96–104).
177 These principles and their main implementations are depicted in the Figure 1.

178

179 **Figure 1:** Overview of several methods for integrating additional constraints into
180 GSMMs based on the regulation of metabolism. Whereas some methods integrate
181 complete functional TRNs or gene expression data into GSMMs, others impose further
182 constraints based on allosteric and post-translational modifications. Additionally, other
183 methods integrate multiple omics layers of regulation of metabolism. For higher
184 eukaryotes such as humans, context-specific models have also been based on tailoring
185 the flux cone of solutions.

186 Surveying all approaches is out of the scope of this review. The following sections will
187 cover the integration of TRNs or gene expression data into GSMMs, focusing on the
188 control of gene expression at the transcriptional level. Figure 2 highlights both
189 approaches, namely the integration of TRNs (Figure 2 A) and gene expression data
190 (Figure 2 B) into GSMMs.

191 The differences between the integration of TRNs and gene expression data into
192 GSMMs are associated with the type and amount of data that these sources can offer
193 to the metabolic landscape of GSMMs.

194 Methods capable of integrating TRNs into GSMMs provide comprehensive
195 knowledge regarding the metabolic and regulatory events occurring inside the cell at
196 the genome-scale. As a result, both regulatory and metabolic networks can be analyzed
197 together at the genome-scale, extending the range of applications of a regular GSMM.

198 On the other hand, gene expression data comprehend a set of snapshots of the
199 transcriptome for several experimental conditions. Thus, a gene expression dataset can
200 solely offer gene expression levels at a given experimental condition.

201 The group of methods aimed at integrating gene expression data with GSMM's
202 comprises methods using only transcriptomics data for tailoring the flux distributions,
203 so no structure or rules describing the regulatory interactions are observed in this class
204 of methods. Thus, the integration of gene expression data focuses on improving the
205 prediction of flux distributions, rather than the study and analysis of an additional
206 biochemical network at the genome-scale.

207 Methods have also been classified according to the main formulations, as previously
208 suggested by Machado *et al.* 2014 (70). Organizing methods into containers, according
209 to their main formulations and features, facilitates the decision process when selecting
210 an adequate method for the existing constraints and data sources.

211 Hence methods were classified into discrete (Figure 2 C) or continuous (Figure 2 D),
212 according to whether phenotype simulations were performed with discrete, namely
213 Boolean logic ("ON/OFF"), or continuous constraints.

214 Accordingly, a method is systematically classified as discrete if the result of the
215 integration is a Boolean value (e.g., 1 for "ON" and 0 for "OFF"), imposing additional
216 constraints on the system. These methods are also referred to as *switch*, since TRN or
217 gene expression data switch reactions on or off. The state of a given metabolic gene is
218 determined by evaluating either the Boolean regulatory rule or thresholding the gene
219 expression level. Then, metabolic reactions mapped to metabolic genes are accessed
220 according to the GPR rules to determine the resulting states. Thus, reactions having a
221 one-to-one direct GPR rule are active/inactive according to the state of the metabolic
222 gene. Reactions catalyzed by enzyme complexes, encoded by multiple yet mandatory
223 genes, are considered inactive if at least one metabolic gene is not available. In contrast,
224 reactions catalyzed by isoenzymes, namely multiple enzymes catalyzing the same
225 reaction, are considered active if at least one metabolic gene is active.

226 Alternatively, there are methods aimed at circumventing the rigid Boolean logic,
227 called *valve* methods, which impose continuous constraints to adjust a given flux
228 distribution gradually and according to penalties, expression scores, or normalized
229 expression levels obtained from the gene expression data. Typically, continuous
230 integration is performed through the implementation of slack variables in the
231 constraints' formulations, altering the reactions' bounds. The slack variable represents
232 penalties, expression scores, or normalized expression levels retrieved from gene
233 expression data for the metabolic genes associated with a given reaction. As before, the
234 state of the metabolic reactions mapped to metabolic genes is assessed through GPR

235 rules, through selecting the best penalty, expression score, or normalized expression
236 level for the slack variable.

237 The methodology for assigning a value to the slack variable, when a set of isozymes
238 catalyzes a given reaction, comprises several distinct approaches. These include:
239 methods in which the slack variable assumes the maximum expression score of the
240 associated genes; methods where the slack variable takes the sum of expression scores
241 of all genes encoding the isozymes catalyzing a single reaction; and, methods in which
242 the reaction is replicated, according to the number of isozymes, and each new reaction
243 is associated with one, and one only, gene.

244 Regarding reactions catalyzed by an enzyme complex, a group of methods establishes
245 that the minimum expression score of all encoding genes is assigned to the slack
246 variable. In contrast, other methods define the utilization of either the geometric or
247 arithmetic mean of the expression score of all genes associated with an enzyme complex
248 or isozymes, respectively.

249 Furthermore, methods capable of integrating gene expression data into GSMMs were
250 also divided into single-condition (Figure 2 A) or multi-condition (Figure 2 B). Notice that
251 this classification was not used to classify those methods aimed at integrating TRNs into
252 GSM models, as it will be explained next.

253 Methods were classified as single-condition (Figure 1 A) or multi-condition (Figure 1
254 B) according to whether phenotype simulations were performed for one or more
255 conditions/states in the gene expression dataset, respectively. For instance, a given
256 method is considered multi-condition if it adjusts the flux cone of solutions by
257 considering all conditions in the gene expression dataset or the gene differential
258 expression between two conditions. Otherwise, the methods are classified as single-
259 condition.

260 Notice that the latter classification was not used to classify those methods aimed at
261 integrating TRNs into GSMMs. Methods capable of integrating TRNs into GSMMs do not
262 require a gene expression dataset, thus classifying them into single- or multi-condition
263 would be meaningless. Other methods capable of assembling and integrating TRNs into
264 GSM models GSMMs often use the whole dataset and can then perform condition-
265 specific phenotype simulations. Hence, classifying these methods as single-condition
266 would be misleading.

267

268 **Figure 2:** Two examples of the integration of TRNs (A and C) or gene expression data
269 (B and D) into GSMMs. The integration of TRNs (A) does not require gene expression
270 data, while methods that integrate gene expression data (B) are capable of tailoring the
271 flux cone of solutions by accounting for one (single-condition) or more (multi-condition)
272 conditions in the gene expression dataset. Both types of integration can be mediated by
273 discrete (C) or continuous (D) variables.

274 An analysis of these methods, encompassing the year of publication, availability of a
275 tool with a user-friendly interface (namely a Graphical User Interface (GUI) without the
276 requirement of coding skills), type of reaction constraint formulation, as well as the
277 organism used for proof of concept has also been conducted. This information is
278 available at the supplementary material 1. Figure 3 provides, on the other hand, a
279 complete understanding of the methods described next, as well as their categorization
280 according to the classification axes described above.

281

282 Integrating TRNs

283 For simulation purposes, the first attempts to integrate TRNs within GSMMs, namely
284 Regulatory Flux Balance Analysis (rFBA) (76,105–107), Steady-state Regulatory Flux
285 Balance (SR-FBA) (77) and the method proposed by Herrgård *et al.* 2006 (79), are based
286 on the *switch* approach, to complement the metabolic system with additional constraints
287 outlining which genes are activated or silenced in the network for specific *stimuli*.

288 As proof of concept, rFBA was successfully used to create the very first integrated
289 genome-scale model of metabolism and gene expression for *E. coli* (106,107). In this
290 reconstruction, as well as in the integrated network of *S. cerevisiae* provided by Herrgård
291 *et al.* 2006 (79), a Boolean network collected from literature was integrated through a
292 set of GPR rules with the GSMM imposing regulatory events as additional time-
293 dependent constraints. On the other hand, SR-FBA performs steady-state simulations by
294 including all valid metabolic and regulatory constraints in the system in a single step,
295 through a Mixed-Integer Linear Programming (MILP) formulation. For that, nested
296 Boolean expressions are formulated as a set of linear constraints, by recursively iterating
297 over the structure of the regulatory layer and GPR rules, to add auxiliary variables
298 representing intermediate Boolean terms (77). As shown in Figure 3, these methods
299 have been classified as discrete, and none provides a user-friendly interface without the
300 requirement of coding skills.

301 Two platforms, namely Toolbox for Integrating Genome-scale Metabolism (TIGER)
302 (84) and FlexFlux (83), have been developed for integrating Boolean-based TRNs into
303 GSMMs. TIGER can convert a series of logic Boolean rules, which can be thought of as a
304 Boolean TRN, into a set of mixed-integer inequalities. Then, several algorithms for
305 integrating gene expression data into the metabolic model and simulating phenotypic
306 behavior can be implemented in the toolbox. Other implementations already available
307 in this toolbox, such as Metabolic Adjustment by Differential Expression (MADE) (95),
308 can be used for simulations.

309 FlexFlux differs from TIGER insofar as it is the only tool that provides a user-friendly
310 interface for the integration of TRNs into GSMMs. This computational tool developed in
311 Java® allows the input of Systems Biology Markup Language (SBML) (6) with the SBML
312 Qualitative (SBML-qual) extension. SBML-qual is the standard file format extension for
313 storing and sharing qualitative multi-state TRNs (108). In this way, a regular SBML file
314 can hold a computer representation of qualitative models of biological networks.

315 Qualitative multi-state regulatory networks can then be used to determine multi-state
316 qualitative constraints for metabolic flux analyses using FBA. Furthermore, FlexFlux
317 allows the translation of the discrete qualitative states into continuous intervals,
318 thereby constraining a reaction flux continuously or discretely (83).

319 Probabilistic Regulation of Metabolism (PROM) (78), PROM2.0 (109), and Integrated
320 Deduced REgulation And Metabolism (IDREAM) (82) are all based on a probabilistic
321 model for TRNs, which are integrated with a constraint-based model using a continuous
322 method. PROM and PROM2.0 were the first attempts to circumvent the previous rigid
323 discrete constraints added to a GSMM by setting the reactions' flux bounds proportional
324 to the probabilities of their associated metabolic genes. In turn, the probability of a
325 metabolic gene being activated in the whole set of conditions is defined together from
326 the TRN and gene expression dataset provided as input. In short, PROM approaches can
327 determine the probability of a given gene being or not activated, when the set of
328 regulating TFs is either activated or silenced. The probability is calculated according to
329 the frequency that each gene is active in the dataset (of either perturbed or over/under-
330 expressed TFs). Likewise, the effect of perturbations on the regulatory network can also
331 be robustly predicted.

332 Although PROM-based approaches are probably the best examples for integrating
333 both TRNs and gene expression data into a GSMM, the gene expression dataset must
334 have a large number of measurements per condition. PROM and PROM2.0 have been
335 validated with *E. coli* and *M. tuberculosis* experimental gene expression data and the
336 respective TRNs.

337 The IDREAM method resulted from the combination of Environment and Gene
338 Regulatory Influence Network (EGRIN) (55,110) and PROM frameworks to create an
339 enhanced genome-scale model of metabolism and gene expression for *Saccharomyces*
340 *cerevisiae* (82). Contrariwise to the previous approaches, this methodology has used a
341 *de novo* reverse engineering method called EGRIN to complement the yeast TRN
342 collected from the database YEAST Search for Transcriptional Regulators And Consensus
343 Tracking (YEASTRACT) (111). Then, the phenotype simulations are conducted in a similar
344 way as in the PROM-based approaches.

345 Transcriptional Regulated Flux Balance Analysis (TRFBA) (81) and CoRegFlux (80) also
346 provide a framework for the integration of gene expression data and TRNs in a
347 continuous manner. Whereas the former requires a TRN for the organism of interest,
348 the latter provides tools for inferring the regulatory network from gene expression data
349 using CoRegNet (43). Nevertheless, CoRegFlux allows us to use a curated TRN rather
350 than using the provided data-driven method.

351 Regarding the TRFBA methodology, this FBA-based approach considers gene
352 expression levels as two additional types of continuous constraints. The first is
353 represented by a constant parameter that converts the gene expression levels to the
354 upper bounds of the reactions. The second type of linear constraints to be added to the

355 system can be thought of as the linear regression of each target gene from the regulating
356 TFs.

357 CoRegFlux differs from TRFBA in that it uses a statistical reverse engineering method
358 to infer targets of a given set of regulators at the genome-scale. Then, the influence
359 score (similar to correlation scores for activation or repression) of each regulator in the
360 set of target genes is calculated with CoRegNet from a large gene expression training
361 dataset. Influence scores are used to train a linear model capable of predicting the gene
362 expression of metabolic genes using a new gene expression dataset. These predicted
363 levels of expression are then translated into flux bounds for the phenotype simulations
364 using FBA or Dynamic Flux Balance Analysis (dFBA) (112).

365 All methods surveyed here are listed in the supplementary material 1.

366

367 Integrating gene expression data

368 The method proposed by Åkesson *et al.* 2004 (87), followed by MADE (95), were the
369 earliest approaches for tailoring the flux cone of solutions using discrete variables
370 obtained solely from gene expression data. In the case of the method developed by
371 Åkesson *et al.* 2004 (87), a reaction is simply switched “off” with a zero flux bound if the
372 associated genes are found to be under-expressed in the corresponding condition
373 (single-condition method). MADE, on the other hand, tries to surpass the problem of
374 arbitrary thresholding under-expression by considering multiple conditions (multi-
375 condition method). Statistical significance between changes in gene expression levels
376 across sequential conditions is calculated to infer whether a gene is activated (95).

377 E-Flux (113) and the method proposed by Lee *et al.* 2012 (94) have introduced several
378 novelties when compared with the previous methodologies. These methods were the
379 first attempts to constraint an FBA-based model using continuous variables.
380 Nevertheless, these approaches are radically different. E-Flux directly maps gene
381 expression levels into flux bound constraints, assuming the maximum flux of a given
382 reaction to be a linear function of the expression of the associated genes in the same
383 condition (single-condition method). Lee and coworkers (94) do not introduce or alter
384 flux bound constraints directly into the GSMM. An alternative objective function that
385 minimizes the distance between flux distributions and gene expression data is applied
386 for each phenotype simulation (single-condition method).

387 The Transcriptional-controlled Flux Balance Analysis (TFBA) method, proposed by van
388 Berlo *et al.* 2011 (93), is aimed at overcoming the problem of setting an arbitrary
389 threshold to determine whether a gene is activated or not. The TFBA assumption is that
390 differential gene expression between two conditions should also be reflected in the flux
391 of the reactions associated with this gene. For that, the authors formulated constraints
392 defining upper and lower limits for fluxes according to the gene expression, though
393 assuming their transgression to be possible. The optimization problem (MILP

394 formulation) consists of finding the flux distribution that minimizes the number of
395 transgressions.

396 Likewise, the method developed by Fang *et al.* 2012 (92) is based on the differential
397 gene expression between two conditions, namely reference and perturbed conditions.
398 This method assumes that the flux distribution of a reference condition can be
399 determined using the FBA or FVA frameworks, while the differential gene expression
400 between the reference and perturbed conditions is used for tailoring the flux
401 distribution of the perturbed one. Also, this method considers the variation of the
402 biomass composition between reference and perturbed conditions.

403 Similarly to TFBA and the method proposed by Fang *et al.* 2012 (92), the Gene
404 Expression Flux Balance Analysis (GX-FBA) method (91) also determines the flux
405 distribution for the reference condition using FBA. Then, GX-FBA employs a new
406 objective function and new constraints derived from the difference between reference
407 and perturbed states to perform the *in-silico* phenotype simulation of the latter state. A
408 wide range of phenomena associated with temperature and known to induce virulence
409 in the gram-negative bacterium *Yersinia pestis* was used as proof of concept.

410 Temporal Expression-based Analysis of Metabolism (TEAM) (90) and Adaptation of
411 Metabolism (AdaM) (89) are the only methods developed for integrating time-series
412 gene expression data into constraint-based models. The former uses dFBA (112) to
413 predict time-series flux distributions based on temporal gene expression profiles. Using
414 a cost minimization scheme similar to the strategy proposed in the context-specific Gene
415 Inactivity Moderated by Metabolism and Expression (GIMME) method (98), TEAM is
416 capable of determining the flux distribution of a GSMM, constrained with gene
417 expression levels of each time step in the dataset. TEAM was tested with time-series
418 gene expression data from *Shewanella oneidensis*.

419 AdaM consists of a flux-based bilevel optimization problem that extracts minimal
420 operating networks from a given GSMM (89). This algorithm infers the minimal
421 operating networks in agreement with the differential gene expression pattern between
422 time-steps. Then, Elementary Flux Modes (EFM)s (114) are computed with these
423 minimal operating networks rather than computing the flux distributions at each time
424 step. Reactions are weighted according to the number of EFMs in which these are
425 present. The optimization problem consists of finding the minimal network having the
426 largest weight.

427 Angione *et al.* 2015 & 2016 (86,88) formulated methods, for example, the Metabolic
428 and Transcriptomics Adaptation Estimator (METRADE), aimed at measuring the
429 adaptability to a changing environmental condition over time. These approaches have
430 provided equally valid methodologies for integrating gene expression data in metabolic
431 networks. In short, these methods have modeled both upper and lower bounds of each
432 reaction as a continuous logarithmic function of the associated gene expression levels.

433 Reaction Inclusion by Parsimony and Transcript Distribution (RIPTiDe) (85) is aimed
434 at circumventing the assumption that reaction fluxes are directly related to the gene

435 expression levels for a given condition. The authors have proposed an unsupervised
436 method that assigns weights (continuous variable) to reactions according to the
437 normalized expression levels of associated genes over the entire dataset. Then, a pFBA
438 simulation considering these linear coefficients is performed. The novelty of this method
439 consists of its validation with precise transcript abundance obtained with RNA-
440 sequencing (RNA-seq).

441 The methods capable of integrating gene expression data into GSMMs addressed
442 herein are available in the supplementary material 1.

443

444 Synopsis

445 The reconstruction of GSMMs is common practice in systems biology nowadays. The
446 advent of the GSMM reconstruction for many organisms was facilitated by the adoption
447 of standard protocols (3), as well as the existence of user-friendly computational tools
448 (7,8), capable of assembling these models from different genomic, enzymatic and
449 stoichiometric data. Nevertheless, the simulation of GSMMs still presents today false-
450 positive phenotypes for several environmental conditions.

451 The reconstruction of TRNs is a well-known strategy in systems biology for
452 understanding the regulatory machinery of a given organism (26,28,30). Although there
453 are many methodologies for assembling a TRN, standard protocols and computational
454 platforms are yet missing to support the reconstruction of TRNs for less described
455 organisms using different data sources. The workflow suggested by Faria *et al.* 2014 (25)
456 highlighted several methodologies that can be combined to extend the reconstruction
457 of TRNs to more bacterial species. To the best of our knowledge little progress has been
458 made to provide a user-friendly platform capable of achieving such a goal. More
459 importantly, the reconstruction of genome-scale TRNs using such integrative workflow,
460 would be pivotal for the reconstruction and simulation of integrated models.

461 The integration of the control of gene expression into GSMMs has been surveyed in
462 this work. A systematic classification that grasps the difference between the several
463 methodologies, capable of integrating and simulating regulatory events into GSMMs
464 was proposed herein. Although part of the reviewed methods have already been
465 surveyed before (25,64,68–71,115), TIGER, FlexFlux, METRADE, IDREAM, TRFBA,
466 CoRegFlux, RIPTiDe and the method proposed by Angione *et al.* (2016) have never been
467 addressed elsewhere in reviews, to the best of our knowledge. Moreover, a detailed
468 categorization that highlights the methodologies used to perform the integration of the
469 regulatory layer into GSMMs has not been provided. This systematic categorization can
470 guide the decision process of selecting the most adequate method of integration and
471 simulation.

472 As shown in Figure 3, there are several methods and toolboxes capable of integrating
473 and simulating TRNs into GSMMs using a discrete approach (76,77,79,83,84). The TRNs
474 used by these methods and toolboxes were mainly reconstructed from literature, which

475 might be a time-consuming approach. The remaining methods allow to assemble TRNs
476 from gene expression data using *de novo* reverse engineering methods. The resulting
477 TRNs can be integrated and simulated with a given GSMM. FlexFlux is the prominent
478 exception as it can perform the integration of the TRN in the GSMM using either discrete
479 or continuous variables.

480 To date, only two prokaryotic organisms, *E. coli* (76–78,81,83) and *M. tuberculosis*
481 (78,109), and the yeast *S. cerevisiae* (79–82,84), have integrated genome-scale models
482 as a result of the integration of complete TRNs into a metabolic network. Nevertheless,
483 some of these reconstructions still require gene expression datasets, namely several
484 methods in the continuous sub-group.

485 Regarding the methods for integrating gene expression data, most of these have
486 provided means for integrating transcriptomics data as continuous constraints from one
487 or more conditions (Figure 3). Only the method proposed by Åkesson *et al.* 2004 (87), as
488 well as MADE (95), use discrete variables to simulate integrated models of metabolism
489 and gene expression.

490 Besides *E. coli*, *M. tuberculosis*, and *S. cerevisiae*, methods for integrating gene
491 expression data have also provided integrated models for *S. oneidensis* (90) and *Y. pestis*
492 (91).

493

494

495 **Figure 3:** Classification of methods aimed at the reconstruction of integrated
496 genome-scale models of metabolism and gene expression. These methods have been
497 divided according to the integration of TRNs (white boxes) or solely gene expression
498 data into GSMMs. Discrete and continuous categories were used to classify these
499 methods according to the usage of discrete, namely Boolean logic (“on/off”), or
500 continuous constraints. Methods capable of integrating gene expression data into
501 GSMMs have been further divided into single-condition (orange circles) and multi-
502 condition (blue ellipses) whether phenotype simulations were performed for one or
503 more conditions/states in the gene expression dataset, respectively. Each inner circle
504 stands for a prokaryotic organism, while the outer circle stands for the baker’s yeast
505 *Saccharomyces cerevisiae*.

506 A vast diversity of methods for the integration of gene expression data in GSMMs has
507 been found. Yet, most methods require large gene expression datasets to be robust,
508 which might not be the case for all organisms. Other methods resort to mapping levels
509 of gene expression directly with the reactions bounds, which again might not be the best
510 approach (70,115,116).

511 Furthermore, the methods for integrating gene expression data with metabolic
512 models previously evaluated by Machado and coworkers (70), namely E-Flux (113),
513 MADE (95), GX-FBA (91) and the method developed by Lee *et al.* 2012 (94) have shown
514 to perform poorly in the designed benchmark. None of the methods have outperformed

515 each other in the phenotype simulations nor pFBA, which indicates that the promising
516 results reported by these methods seem to be mere artifacts related to rigid constraints
517 created around the nature of the gene expression dataset.

518 The reconstruction of integrated models using TRNs is, in theory, more useful than
519 merely integrating gene expression data into GSMMs. Integrated models that result
520 from the integration of TRNs provide comprehensive knowledge regarding the
521 metabolic and regulatory events happening inside the cell, thus leading to a broader
522 range of applications when compared to a regular GSMM (117,118).

523 Moreover, the diversity of methods for reconstructing TRNs using different data
524 sources, such as gene expression, transcription factor binding site, or comparative
525 genomics analysis, eases the reconstruction of TRNs for most prokaryotic organisms
526 having a sequenced genome (25). However, the absence of a user-friendly
527 computational tool based on the ensemble of these different approaches is missing. In
528 contrast, the same strategy has yielded results in the reconstruction of GSMMs (7–
529 9,11,15,119).

530 In short, the existence of standardized protocols and easy to use computational tools
531 for the generation of GSMMs has eased its practice in systems biology to study the
532 metabolism of many organisms. In contrast, the absence of the computational tools that
533 ease the reconstruction of TRNs from different sources of regulatory data hindered a
534 similar approach.

535 The integration and analysis of regulatory events into GSMMs has been surveyed
536 herein. A systematic classification has been created to grasp the difference between the
537 several methodologies capable of integrating and simulating regulatory events into
538 GSMMs. As a result, two primary approaches have been determined, namely the
539 integration of TRNs and/or the integration of gene expression data.

540 The major obstacle when using the methods described in this survey to simulate
541 integrated genome-scale models of metabolism and gene expression is not reproducing
542 their results, but rather extending their implementations to other organisms and case
543 studies. This hurdle poses a stiff challenge for using these methods out of the scope they
544 were aimed at during development.

545 The requirement for large gene expression datasets with specific experimental
546 conditions, the usage of TRNs reconstructed solely from literature, and the output of
547 biased results strictly related to rigid constraints, are specific indicators of issues
548 preventing the scaling-up of the reconstruction and analysis of integrated models. In
549 short, there is a vast diversity of methods capable of integrating and simulating the
550 effect of regulation into the metabolism, though few approaches that ease the
551 reconstruction of these integrated models.

552 Hence, the perspective of reconstructing integrated genome-scale models of
553 metabolism and gene expression for diverse prokaryotes is still a complex endeavor.

554 The implementation of a user-friendly computational framework that does not
555 require coding skills, is capable of running a semi-automated pipeline for reconstructing
556 TRNs or analyzing gene expression data, and performs its integration into standard
557 GSMs, would be a clear breakthrough towards the reconstruction and simulation of
558 integrated genome-scale models of metabolism and gene expression. This hypothetical
559 computational tool should be able to combine different sources of regulatory
560 information which are seldom combined.

561

562 Perspectives

- 563 • The advent of high-throughput large-scale omics experiments has been supporting
564 the study of the genomics, transcriptomics, and metabolomics layers of the
565 cellular's molecular machinery. Systems biology can take advantage of the sheer
566 volume of studies in these different omics fields by proposing differentiated
567 approaches, such as the reconstruction of *in silico* networks and models
- 568 • The reconstruction and analysis of integrated models based on the integration of
569 TRNs into GSMs has not been conventional for non-model prokaryotic
570 organisms. Usually, these lack large gene expression datasets, or have few sources
571 of regulatory data. In addition to the absence of an established methodology and
572 of easy to use tools and algorithms, the reconstruction and integration of TRNs
573 into GSMs is almost impracticable.
- 574 • Hence, a user-friendly computational framework that facilitates the
575 reconstruction of TRNs and allows to integrate these into GSMs would be a step
576 towards facilitating the extension of integrated models to other prokaryotic
577 organisms.

578

579 Acknowledgements

580 This study was supported by the Portuguese Foundation for Science and Technology
581 (FCT) under the scope of the strategic funding of UIDB/04469/2020 unit and
582 BioTecNorte operation (NORTE-01-0145-FEDER-000004) funded by the European
583 Regional Development Fund under the scope of Norte2020 - Programa Operacional
584 Regional do Norte. Fernando Cruz holds a doctoral fellowship (SFRH/BD/139198/2018)
585 funded by the FCT. This study was supported by the European Commission through
586 project SHIKIFACTORY100 - Modular cell factories for the production of 100 compounds
587 from the shikimate pathway (Reference 814408). The submitted manuscript has been
588 created by UChicago Argonne, LLC as Operator of Argonne National Laboratory
589 ("Argonne") under Contract No. DE-AC02-06CH11357 with the U.S. Department of
590 Energy. The U.S. Government retains for itself, and others acting on its behalf, a paid-
591 up, nonexclusive, irrevocable worldwide license in said article to reproduce, prepare
592 derivative works, distribute copies to the public, and perform publicly and display
593 publicly, by or on behalf of the Government. The Department of Energy will provide

594 public access to these results of federally sponsored research in accordance with the
595 DOE Public Access Plan.

596 Competing Interests

597 The Authors declare that there are no competing interests associated with the
598 manuscript.

599

600 Contribution

601 FC and OD conceived and designed the study. FR assessed the methods and drafted the
602 manuscript. OD managed its coordination and helped to draft the manuscript. JPF, MR
603 and IR participated in the coordination of the study and helped to draft the manuscript.
604 All authors read and approved the final manuscript.

605 Bibliography

- 606 1. Gray AN, Koo B-M, Shiver AL, Peters JM, Osadnik H, Gross CA. High-throughput
607 bacterial functional genomics in the sequencing era. *Curr Opin Microbiol.* 2015
608 Oct 1;27:86–95.
- 609 2. Rocha I, Förster J, Nielsen J. Design and Application of Genome-Scale
610 Reconstructed Metabolic Models. *Methods Mol Biol vol 416 Microb Gene*
611 *Essentiality.* 2007;416:409–31.
- 612 3. Thiele I, Palsson BØ. A protocol for generating a high-quality genome-scale
613 metabolic reconstruction. *Nat Protoc.* 2010;5(1):93–121.
- 614 4. Dias O, Rocha I. Systems Biology in Fungi. In: Paterson R, editor. *Molecular*
615 *Biology of Food and Water Borne Mycotoxigenic and Mycotic Fungi.* Boca Raton:
616 CRC Press; 2015. p. 69–92.
- 617 5. Feist AM, Herrgård MJ, Thiele I, Reed JL, Palsson BØ. Reconstruction of
618 biochemical networks in microorganisms. *Nat Rev Microbiol.* 2009
619 Feb;7(2):129–43.
- 620 6. Hucka M, Finney A, Sauro HM, Bolouri H, Doyle JC, Kitano H, et al. The systems
621 biology markup language (SBML): a medium for representation and exchange of
622 biochemical network models. *Bioinformatics.* 2003 Mar 1;19(4):524–31.
- 623 7. Dias O, Rocha M, Ferreira EC, Rocha I. Reconstructing genome-scale metabolic
624 models with merlin. *Nucleic Acids Res.* 2015;43(8):3899–910.
- 625 8. Henry CS, DeJongh M, Best AA, Frybarger PM, Linsay B, Stevens RL. High-
626 throughput generation, optimization and analysis of genome-scale metabolic
627 models. *Nat Biotechnol.* 2010 Sep 29;28(9):977–82.
- 628 9. Agren R, Liu L, Shoaie S, Vongsangnak W, Nookaew I, Nielsen J. The RAVEN
629 Toolbox and Its Use for Generating a Genome-scale Metabolic Model for
630 *Penicillium chrysogenum.* Maranas CD, editor. *PLoS Comput Biol.* 2013 Mar
631 21;9(3):e1002980.

- 632 10. Machado D, Andrejev S, Tramontano M, Patil KR. Fast automated reconstruction
633 of genome-scale metabolic models for microbial species and communities.
634 *Nucleic Acids Res.* 2018 Sep 6;46(15):7542–53.
- 635 11. Faria JP, Rocha M, Rocha I, Henry CS. Methods for automated genome-scale
636 metabolic model reconstruction. Vol. 46, *Biochemical Society Transactions.*
637 Portland Press Ltd; 2018. p. 931–6.
- 638 12. Orth JD, Thiele I, Palsson BO. What is flux balance analysis? *Nat Publ Gr.*
639 2010;28(3):245–8.
- 640 13. Lewis NE, Hixson KK, Conrad TM, Lerman JA, Charusanti P, Polpitiya AD, et al.
641 Omic data from evolved *E. coli* are consistent with computed optimal growth
642 from genome-scale models. *Mol Syst Biol.* 2010 Jul 1;6(1):390.
- 643 14. Gudmundsson S, Thiele I. Computationally efficient flux variability analysis. *BMC*
644 *Bioinformatics.* 2010 Dec 29;11(1):489.
- 645 15. Heirendt L, Arreckx S, Pfau T, Mendoza SN, Richelle A, Heinken A, et al. Creation
646 and analysis of biochemical constraint-based models using the COBRA Toolbox
647 v.3.0. *Nat Protoc.* 2019 Mar 1;14(3):639–702.
- 648 16. Ebrahim A, Lerman JA, Palsson BO, Hyduke DR. COBRApy: COstraints-Based
649 Reconstruction and Analysis for Python. *BMC Syst Biol.* 2013 Aug 8;7(1):74.
- 650 17. Vilaça P, Rocha I, Rocha M. A computational tool for the simulation and
651 optimization of microbial strains accounting integrated metabolic/regulatory
652 information. *Biosystems.* 2011 Mar 1;103(3):435–41.
- 653 18. Edwards JS, Palsson BO. Systems properties of the *Haemophilus influenzae* Rd
654 metabolic genotype. *J Biol Chem.* 1999 Jun 18;274(25):17410–6.
- 655 19. Bro C, Regenberg B, Förster J, Nielsen J. In silico aided metabolic engineering of
656 *Saccharomyces cerevisiae* for improved bioethanol production. *Metab Eng.*
657 2006 Mar;8(2):102–11.
- 658 20. Henry CS, Zinner JF, Cohoon MP, Stevens RL. iBsu1103: A new genome-scale
659 metabolic model of *Bacillus subtilis* based on SEED annotations. *Genome Biol.*
660 2009 Jun 25;10(6).
- 661 21. Bordbar A, Palsson BO. Using the reconstructed genome-scale human metabolic
662 network to study physiology and pathology. *J Intern Med.* 2012 Feb;271(2):131–
663 41.
- 664 22. Dias O, Pereira R, Gombert AK, Ferreira EC, Rocha I. iOD907, the first genome-
665 scale metabolic model for the milk yeast *Kluyveromyces lactis*. *Biotechnol J.*
666 2014;9(6):776–90.
- 667 23. Dias O, Saraiva J, Faria C, Ramirez M, Pinto F, Rocha I. IDS372, a phenotypically
668 reconciled model for the metabolism of *Streptococcus pneumoniae* strain R6.
669 *Front Microbiol.* 2019;10(JUN).
- 670 24. Orth JD, Conrad TM, Na J, Lerman JA, Nam H, Feist AM, et al. A comprehensive
671 genome-scale reconstruction of *Escherichia coli* metabolism—2011. *Mol Syst*

- 672 Biol. 2011 Jan 11;7(1):535.
- 673 25. Faria JP, Overbeek R, Xia F, Rocha M, Rocha I, Henry CS. Genome-scale bacterial
674 transcriptional regulatory networks: reconstruction and integrated analysis with
675 metabolic models. *Brief Bioinform.* 2014 Jul 1;15(4):592–611.
- 676 26. De Smet R, Marchal K. Advantages and limitations of current network inference
677 methods. *Nature Reviews Microbiology.* 2010.
- 678 27. Thompson D, Regev A, Roy S. Comparative Analysis of Gene Regulatory
679 Networks: From Network Reconstruction to Evolution. *Annu Rev Cell Dev Biol.*
680 2015 Nov 13;31(1):399–428.
- 681 28. Barbosa S, Niebel B, Wolf S, Mauch K, Takors R. A guide to gene regulatory
682 network inference for obtaining predictive solutions: Underlying assumptions
683 and fundamental biological and data constraints. *Biosystems.* 2018 Dec
684 1;174:37–48.
- 685 29. Huynh-Thu VA, Sanguinetti G. Gene Regulatory Network Inference: An
686 Introductory Survey. In: *Methods in Molecular Biology.* Humana Press Inc.; 2019.
687 p. 1–23.
- 688 30. Karlebach G, Shamir R. Modelling and analysis of gene regulatory networks. *Nat*
689 *Rev Mol Cell Biol.* 2008 Oct 17;9(10):770–80.
- 690 31. Rodionov DA. Comparative Genomic Reconstruction of Transcriptional
691 Regulatory Networks in Bacteria. 2007;
- 692 32. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, et al.
693 NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Res.*
694 2012 Nov 26;41(D1):D991–5.
- 695 33. Kolesnikov N, Hastings E, Keays M, Melnichuk O, Tang YA, Williams E, et al.
696 ArrayExpress update—simplifying data submissions. *Nucleic Acids Res.* 2015 Jan
697 28;43(D1):D1113–6.
- 698 34. Novichkov PS, Brettin TS, Novichkova ES, Dehal PS, Arkin AP, Dubchak I, et al.
699 RegPrecise web services interface: programmatic access to the transcriptional
700 regulatory interactions in bacteria reconstructed by comparative genomics.
701 *Nucleic Acids Res.* 2012 Jul;40(W1):W604–8.
- 702 35. Santos-Zavaleta A, Sánchez-Pérez M, Salgado H, Velázquez-Ramírez DA, Gama-
703 Castro S, Tierrafría VH, et al. A unified resource for transcriptional regulation in
704 *Escherichia coli* K-12 incorporating high-throughput-generated binding data into
705 RegulonDB version 10.0. *BMC Biol.* 2018 Dec 16;16(1):91.
- 706 36. Sierra N, Makita Y, de Hoon M, Nakai K. DBTBS: a database of transcriptional
707 regulation in *Bacillus subtilis* containing upstream intergenic conservation
708 information. *Nucleic Acids Res.* 2008 Jan 1;36(suppl_1):D93–6.
- 709 37. Lozada-Chavez I, Janga SC, Collado-Vides J. Bacterial regulatory networks are
710 extremely flexible in evolution. *Nucleic Acids Res.* 2006 Jul 19;34(12):3434–45.
- 711 38. Gelfand MS. Evolution of transcriptional regulatory networks in microbial

- 712 genomes. Vol. 16, Current Opinion in Structural Biology. 2006. p. 420–9.
- 713 39. Madan Babu M, Teichmann SA, Aravind L. Evolutionary Dynamics of Prokaryotic
714 Transcriptional Regulatory Networks. *J Mol Biol.* 2006 Apr 28;358(2):614–33.
- 715 40. Tompa M, Li N, Bailey TL, Church GM, De Moor B, Eskin E, et al. Assessing
716 computational tools for the discovery of transcription factor binding sites.
717 *Nature Biotechnology.* 2005.
- 718 41. Alkema WBL, Lenhard B, Wasserman WW. Regulog analysis: Detection of
719 conserved regulatory networks across bacteria: Application to *Staphylococcus*
720 *aureus*. *Genome Res.* 2004;
- 721 42. Novichkov PS, Rodionov DA, Stavrovskaya ED, Novichkova ES, Kazakov AE,
722 Gelfand MS, et al. RegPredict: An integrated system for regulon inference in
723 prokaryotes by comparative genomics approach. *Nucleic Acids Res.* 2010;
- 724 43. Nicolle R, Radvanyi F, Elati M. COREGNET: reconstruction and integrated analysis
725 of co-regulatory networks. *Bioinformatics.* 2015 Sep 15;31(18):3066–8.
- 726 44. Faith JJ, Hayete B, Thaden JT, Mogno I, Wierzbowski J, Cottarel G, et al. Large-
727 scale mapping and validation of *Escherichia coli* transcriptional regulation from a
728 compendium of expression profiles. *PLoS Biol.* 2007;
- 729 45. Karlebach G, Shamir R. Constructing logical models of gene regulatory networks
730 by integrating transcription factor-DNA interactions with expression data: An
731 entropy-based approach. *J Comput Biol.* 2012 Jan 1;19(1):30–41.
- 732 46. Karlebach G, Shamir R. Constructing logical models of gene regulatory networks
733 by integrating transcription factor-DNA interactions with expression data: An
734 entropy-based approach. *J Comput Biol.* 2012 Jan 1;19(1):30–41.
- 735 47. Huynh-Thu VA, Irrthum A, Wehenkel L, Geurts P. Inferring regulatory networks
736 from expression data using tree-based methods. *PLoS One.* 2010;
- 737 48. Bonneau R, Reiss DJ, Shannon P, Facciotti M, Hood L, Baliga NS, et al. The
738 inferelator: An algorithm for learning parsimonious regulatory networks from
739 systems-biology data sets de novo. *Genome Biol.* 2006;
- 740 49. Gat-Viks I, Shamir R. Refinement and expansion of signaling pathways: The
741 osmotic response network in yeast. *Genome Res.* 2007 Mar;17(3):358–67.
- 742 50. Fang X, Sastry A, Mih N, Kim D, Tan J, Yurkovich JT, et al. Global transcriptional
743 regulatory network for *Escherichia coli* robustly connects gene expression to
744 transcription factor activities. *Proc Natl Acad Sci U S A.* 2017 Sep
745 19;114(38):10286–91.
- 746 51. Gao Y, Yurkovich JT, Seo SW, Kabimoldayev I, Dräger A, Chen K, et al. Systematic
747 discovery of uncharacterized transcription factors in *Escherichia coli* K-12
748 MG1655. *Nucleic Acids Res.* 2018 Aug 23;46(20):10682–96.
- 749 52. Faria JP, Overbeek R, Taylor RC, Conrad N, Vonstein V, Goelzer A, et al.
750 Reconstruction of the Regulatory Network for *Bacillus subtilis* and Reconciliation
751 with Gene Expression Data. *Front Microbiol.* 2016 Mar 18;7:275.

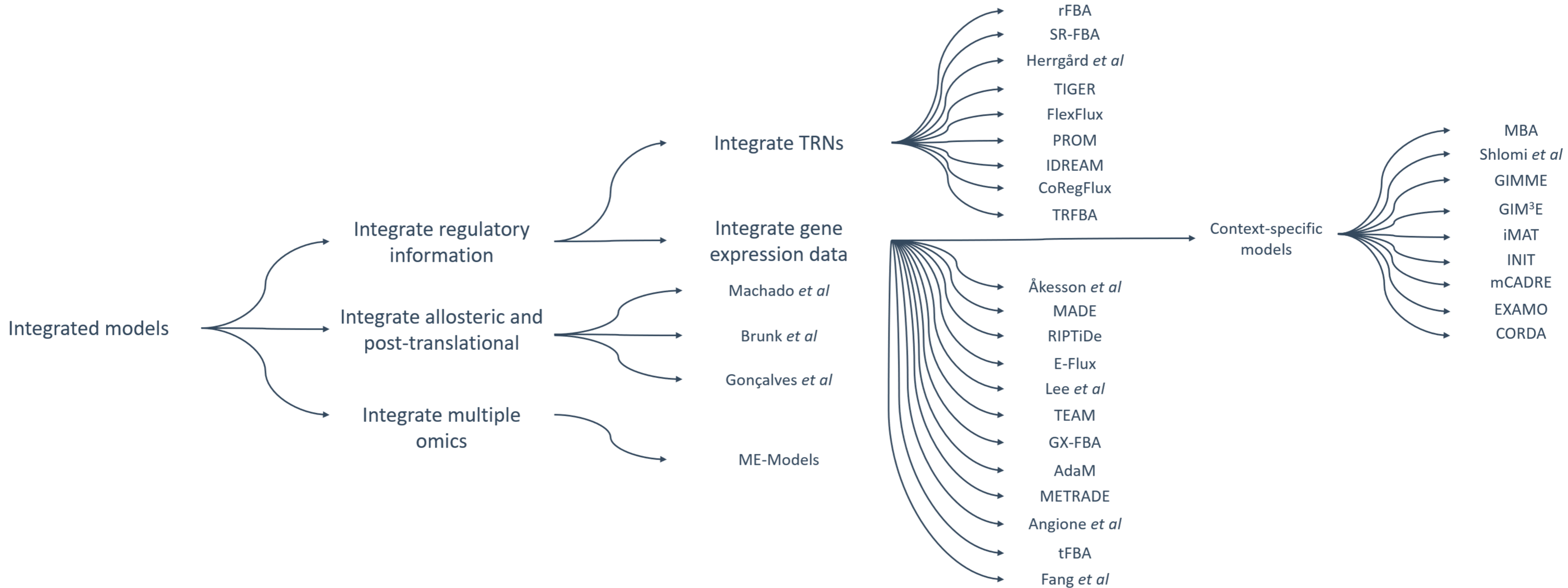
- 752 53. Arrieta-Ortiz ML, Hafemeister C, Bate AR, Chu T, Greenfield A, Shuster B, et al.
753 An experimentally supported model of the *Bacillus subtilis* global transcriptional
754 regulatory network. *Mol Syst Biol*. 2015 Nov 17;11(11):839.
- 755 54. Turkarslan S, Peterson EJR, Rustad TR, Minch KJ, Reiss DJ, Morrison R, et al. A
756 comprehensive map of genome-wide gene regulation in *Mycobacterium*
757 *tuberculosis*. *Sci Data*. 2015 Mar 31;2:150010.
- 758 55. Bonneau R, Facciotti MT, Reiss DJ, Schmid AK, Pan M, Kaur A, et al. A Predictive
759 Model for Transcriptional Control of Physiology in a Free Living Cell. *Cell*. 2007
760 Dec 28;131(7):1354–65.
- 761 56. Galán-Vázquez E, Luna B, Martínez-Antonio A. The Regulatory Network of
762 *Pseudomonas aeruginosa*. *Microb Inform Exp*. 2011;1(1):3.
- 763 57. Marbach D, Costello JC, Küffner R, Vega NMN, Prill RJ, Camacho DM, et al.
764 Wisdom of crowds for robust gene network inference. *Nat Methods*. 2012
765 Jul;9(8):796–804.
- 766 58. Fredrickson JK, Romine MF, Beliaev AS, Auchtung JM, Driscoll ME, Gardner TS, et
767 al. Towards environmental systems biology of *Shewanella*. Vol. 6, *Nature*
768 *Reviews Microbiology*. 2008. p. 592–603.
- 769 59. Rodionov DA, Rodionova IA, Li X, Ravcheev DA, Tarasova Y, Portnoy VA, et al.
770 Transcriptional regulation of the carbohydrate utilization network in
771 *Thermotoga maritima*. *Front Microbiol*. 2013;4(AUG).
- 772 60. de Jong A, Hansen ME, Kuipers OP, Kilstrup M, Kok J. The Transcriptional and
773 Gene Regulatory Network of *Lactococcus lactis* MG1363 during Growth in Milk.
774 *PLoS One*. 2013 Jan 17;8(1).
- 775 61. Schmitt WA, Raab RM, Stephanopoulos G. Elucidation of gene interaction
776 networks through time-lagged correlation analysis of transcriptional data.
777 *Genome Res*. 2004 Aug;14(8):1654–63.
- 778 62. Nelson DL, Cox MM. *Lehninger Principles of Biochemistry*. 6th edit. W.H.
779 Freeman; 2008. 1328 p.
- 780 63. Lerman JA, Hyduke DR, Latif H, Portnoy VA, Lewis NE, Orth JD, et al. In silico
781 method for modelling metabolism and gene product expression at genome
782 scale. *Nat Commun*. 2012 Jan 3;3(1):929.
- 783 64. Hao T, Wu D, Zhao L, Wang Q, Wang E, Sun J. The genome-scale integrated
784 networks in microorganisms. Vol. 9, *Frontiers in Microbiology*. Frontiers Media
785 S.A.; 2018.
- 786 65. Lloyd CJ, Ebrahim A, Yang L, King ZA, Catoiu E, O'Brien EJ, et al. COBRAME: A
787 computational framework for genome-scale models of metabolism and gene
788 expression. Darling AE, editor. *PLOS Comput Biol*. 2018 Jul 5;14(7):e1006302.
- 789 66. Brunk E, Chang RL, Xia J, Hefzi H, Yurkovich JT, Kim D, et al. Characterizing
790 posttranslational modifications in prokaryotic metabolism using a multiscale
791 workflow. *Proc Natl Acad Sci*. 2018 Oct 23;115(43):11096–101.

- 792 67. Gonçalves E, Sciacovelli M, Costa ASH, Tran MGB, Johnson TI, Machado D, et al.
793 Post-translational regulation of metabolism in fumarate hydratase deficient
794 cancer cells. *Metab Eng.* 2018 Jan 1;45:149–57.
- 795 68. Blazier AS, Papin JA. Integration of expression data in genome-scale metabolic
796 network reconstructions. *Front Physiol.* 2012 Aug 6;3:299.
- 797 69. Kim J, Reed JL. Refining metabolic models and accounting for regulatory effects.
798 *Curr Opin Biotechnol.* 2014 Mar;29C:34–8.
- 799 70. Machado D, Herrgård M. Systematic Evaluation of Methods for Integration of
800 Transcriptomic Data into Constraint-Based Models of Metabolism. Maranas CD,
801 editor. *PLoS Comput Biol.* 2014 Apr 24;10(4):e1003580.
- 802 71. Imam S, Schäuble S, Brooks AN, Baliga NS, Price ND. Data-driven integration of
803 genome-scale regulatory and metabolic network models. *Front Microbiol.*
804 2015;6:409.
- 805 72. O'Brien EJ, Palsson BO. Computing the functional proteome: recent progress
806 and future prospects for genome-scale models. *Curr Opin Biotechnol.* 2015 Aug
807 1;34:125–34.
- 808 73. Machado D, Herrgård MJ, Rocha I. Modeling the Contribution of Allosteric
809 Regulation for Flux Control in the Central Carbon Metabolism of *E. coli*. *Front*
810 *Bioeng Biotechnol.* 2015 Oct 8;3:154.
- 811 74. O'Brien EJ, Utrilla J, Palsson BO. Quantification and Classification of *E. coli*
812 Proteome Utilization and Unused Protein Costs across Environments. Maranas
813 CD, editor. *PLOS Comput Biol.* 2016 Jun 28;12(6):e1004998.
- 814 75. Opdam S, Richelle A, Kellman B, Li S, Zielinski DC, Lewis NE. A Systematic
815 Evaluation of Methods for Tailoring Genome-Scale Metabolic Models. *Cell Syst.*
816 2017 Mar 22;4(3):318-329.e6.
- 817 76. Covert MW, Knight EM, Reed JL, Herrgard MJ, Palsson BO. Integrating high-
818 throughput and computational data elucidates bacterial networks. *Nature.* 2004
819 May 6;429(6987):92–6.
- 820 77. Shlomi T, Eisenberg Y, Sharan R, Ruppin E. A genome-scale computational study
821 of the interplay between transcriptional regulation and metabolism. *Mol Syst*
822 *Biol.* 2007 Apr;3:101.
- 823 78. Chandrasekaran S, Price ND. Probabilistic integrative modeling of genome-scale
824 metabolic and regulatory networks in *Escherichia coli* and *Mycobacterium*
825 *tuberculosis*. *Proc Natl Acad Sci U S A.* 2010 Oct;107(41):17845–50.
- 826 79. Herrgård MJ, Lee B-S, Portnoy V, Palsson BØ. Integrated analysis of regulatory
827 and metabolic networks reveals novel regulatory mechanisms in *Saccharomyces*
828 *cerevisiae*. *Genome Res.* 2006 May;16(5):627–35.
- 829 80. Banos DT, Trébulle P, Elati M. Integrating transcriptional activity in genome-
830 scale models of metabolism. *BMC Syst Biol.* 2017 Dec 21;11(S7):134.
- 831 81. Motamedian E, Mohammadi M, Shojaosadati SA, Heydari M. TRFBA: an

- 832 algorithm to integrate genome-scale metabolic and transcriptional regulatory
833 networks with incorporation of expression data. *Bioinformatics*. 2017 Jan
834 8;33(7):btw772.
- 835 82. Wang Z, Danziger SA, Heavner BD, Ma S, Smith JJ, Li S, et al. Combining inferred
836 regulatory and reconstructed metabolic networks enhances phenotype
837 prediction in yeast. Nielsen J, editor. *PLOS Comput Biol*. 2017
838 May;13(5):e1005489.
- 839 83. Marmiesse L, Peyraud R, Cottret L. FlexFlux: combining metabolic flux and
840 regulatory network analyses. *BMC Syst Biol*. 2015 Dec 15;9(1):93.
- 841 84. Jensen PA, Lutz KA, Papin JA. TIGER: Toolbox for integrating genome-scale
842 metabolic models, expression data, and transcriptional regulatory networks.
843 *BMC Syst Biol*. 2011 Sep 23;5(1):147.
- 844 85. Jenior ML, Moutinho TJ, Papin JA. Parsimonious transcript data integration
845 improves context-specific predictions of bacterial metabolism in complex
846 environments. *bioRxiv*. 2019;
- 847 86. Angione C, Conway M, Lió P. Multiplex methods provide effective integration of
848 multi-omic data in genome-scale models. *BMC Bioinformatics*. 2016 Feb
849 2;17(S4):83.
- 850 87. Åkesson M, Förster J, Nielsen J. Integration of gene expression data into
851 genome-scale metabolic models. *Metab Eng*. 2004 Oct 1;6(4):285–93.
- 852 88. Angione C, Lió P. Predictive analytics of environmental adaptability in multi-omic
853 network models. *Sci Rep*. 2015 Oct 20;5.
- 854 89. Töpfer N, Jozefczuk S, Nikoloski Z. Integration of time-resolved transcriptomics
855 data with flux-based methods reveals stress-induced metabolic adaptation in
856 *Escherichia coli*. *BMC Syst Biol*. 2012 Nov 30;6.
- 857 90. Collins SB, Reznik E, Segrè D. Temporal Expression-based Analysis of
858 Metabolism. *PLoS Comput Biol*. 2012 Nov;8(11).
- 859 91. Navid A, Almaas E. Genome-level transcription data of *Yersinia pestis* analyzed
860 with a New metabolic constraint-based approach. *BMC Syst Biol*. 2012 Dec 6;6.
- 861 92. Fang X, Wallqvist A, Reifman J. Modeling Phenotypic Metabolic Adaptations of
862 *Mycobacterium tuberculosis* H37Rv under Hypoxia. *PLoS Comput Biol*. 2012
863 Sep;8(9).
- 864 93. van Berlo RJP, de Ridder D, Daran J-M, Daran-Lapujade PAS, Teusink B, Reinders
865 MJT. Predicting Metabolic Fluxes Using Gene Expression Differences As
866 Constraints. *IEEE/ACM Trans Comput Biol Bioinforma*. 2011 Jan;8(1):206–16.
- 867 94. Lee D, Smallbone K, Dunn WB, Murabito E, Winder CL, Kell DB, et al. Improving
868 metabolic flux predictions using absolute gene expression data. *BMC Syst Biol*.
869 2012 Jun 19;6(1):73.
- 870 95. Jensen PA, Papin JA. Functional integration of a metabolic network model and
871 expression data without arbitrary thresholding. *Bioinformatics*. 2011 Feb

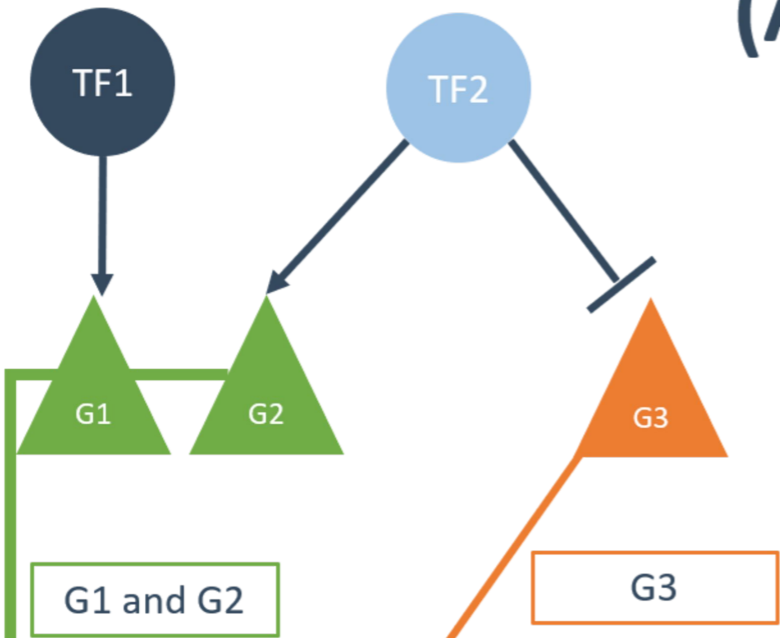
- 872 15;27(4):541–7.
- 873 96. Jerby L, Shlomi T, Ruppin E. Computational reconstruction of tissue-specific
874 metabolic models: Application to human liver metabolism. *Mol Syst Biol.*
875 2010;6.
- 876 97. Shlomi T, Cabili MN, Herrgård MJ, Palsson B, Ruppin E. Network-based
877 prediction of human tissue-specific metabolism. Vol. 26, *Nature Biotechnology.*
878 2008. p. 1003–10.
- 879 98. Becker SA, Palsson BO. Context-specific metabolic networks are consistent with
880 experiments. *PLoS Comput Biol.* 2008 May;4(5).
- 881 99. Zur H, Ruppin E, Shlomi T. iMAT: An integrative metabolic analysis tool.
882 *Bioinformatics.* 2010 Dec;26(24):3140–2.
- 883 100. Agren R, Bordel S, Mardinoglu A, Pornputtapong N, Nookaew I, Nielsen J.
884 Reconstruction of genome-scale active metabolic networks for 69 human cell
885 types and 16 cancer types using INIT. *PLoS Comput Biol.* 2012 May;8(5).
- 886 101. Wang Y, Eddy JA, Price ND. Reconstruction of genome-scale metabolic models
887 for 126 human tissues using mCADRE. *BMC Syst Biol.* 2012 Dec 13;6.
- 888 102. Schmidt BJ, Ebrahim A, Metz TO, Adkins JN, Palsson BØ, Hyduke DR. GIM3E:
889 condition-specific models of cellular metabolism developed from metabolomics
890 and expression data. *Bioinformatics.* 2013 Nov 15;29(22):2900–8.
- 891 103. Rossell S, Huynen MA, Notebaart RA. Inferring Metabolic States in
892 Uncharacterized Environments Using Gene-Expression Measurements. *PLoS*
893 *Comput Biol.* 2013;9(3).
- 894 104. Schultz A, Qutub AA. Reconstruction of Tissue-Specific Metabolic Networks
895 Using CORDA. *PLoS Comput Biol.* 2016 Mar 1;12(3).
- 896 105. Covert MW, Schilling CH, Palsson B. Regulation of Gene Expression in Flux
897 Balance Models of Metabolism. *J Theor Biol.* 2001 Nov 7;213(1):73–88.
- 898 106. Covert MW, Palsson BO. Constraints-based models: Regulation of Gene
899 Expression Reduces the Steady-state Solution Space. *J Theor Biol.* 2003 Apr
900 7;221(3):309–25.
- 901 107. Covert MW, Palsson BØ. Transcriptional regulation in constraints-based
902 metabolic models of *Escherichia coli*. *J Biol Chem.* 2002 Aug 2;277(31):28058–
903 64.
- 904 108. Chaouiya C, Bérenguier D, Keating SM, Naldi A, van Iersel MP, Rodriguez N, et al.
905 SBML qualitative models: a model representation format and infrastructure to
906 foster interactions between qualitative modelling formalisms and tools. *BMC*
907 *Syst Biol.* 2013 Dec 10;7(1):135.
- 908 109. Ma S, Minch KJ, Rustad TR, Hobbs S, Zhou S-L, Sherman DR, et al. Integrated
909 Modeling of Gene Regulatory and Metabolic Networks in *Mycobacterium*
910 *tuberculosis*. Stelling J, editor. *PLOS Comput Biol.* 2015 Nov 30;11(11):e1004543.

- 911 110. Brooks AN, Reiss DJ, Allard A, Wu W, Salvanha DM, Plaisier CL, et al. A system-
912 level model for the microbial regulatory genome. *Mol Syst Biol*. 2014
913 Jul;10(7):740.
- 914 111. Teixeira MC, Monteiro PT, Palma M, Costa C, Godinho CP, Pais P, et al.
915 YEASTRACT: An upgraded database for the analysis of transcription regulatory
916 networks in *Saccharomyces cerevisiae*. *Nucleic Acids Res*. 2018 Jan
917 1;46(D1):D348–53.
- 918 112. Mahadevan R, Edwards JS, Doyle FJ. Dynamic Flux Balance Analysis of diauxic
919 growth in *Escherichia coli*. *Biophys J*. 2002;83(3):1331–40.
- 920 113. Colijn C, Brandes A, Zucker J, Lun DS, Weiner B, Farhat MR, et al. Interpreting
921 Expression Data with Metabolic Flux Models: Predicting *Mycobacterium*
922 tuberculosis Mycolic Acid Production. Papin JA, editor. *PLoS Comput Biol*. 2009
923 Aug 28;5(8):e1000489.
- 924 114. Schuster S, Fell DA, Dandekar T. A general definition of metabolic pathways
925 useful for systematic organization and analysis of complex metabolic networks.
926 *Nat Biotechnol*. 2000;18(3):326–32.
- 927 115. Vivek-Ananth RP, Samal A. Advances in the integration of transcriptional
928 regulatory information into genome-scale metabolic models. *Biosystems*. 2016
929 Sep 1;147:1–10.
- 930 116. Kochanowski K, Sauer U, Chubukov V. Somewhat in control-the role of
931 transcription in regulating microbial metabolic fluxes. Vol. 24, *Current Opinion in*
932 *Biotechnology*. *Curr Opin Biotechnol*; 2013. p. 987–93.
- 933 117. Kim J, Reed JL. OptORF: Optimal metabolic and regulatory perturbations for
934 metabolic engineering of microbial strains. *BMC Syst Biol*. 2010 Apr;4(1):53.
- 935 118. Shen F, Sun R, Yao J, Li J, Liu Q, Price ND, et al. OptRAM: In-silico strain design
936 via integrative regulatory-metabolic network modeling. Ouzounis CA, editor.
937 *PLOS Comput Biol*. 2019 Mar 8;15(3):e1006835.
- 938 119. Monk J, Nogales J, Palsson BO. Optimizing genome-scale network
939 reconstructions. Vol. 32, *Nature Biotechnology*. Nature Publishing Group; 2014.
940 p. 447–52.
- 941



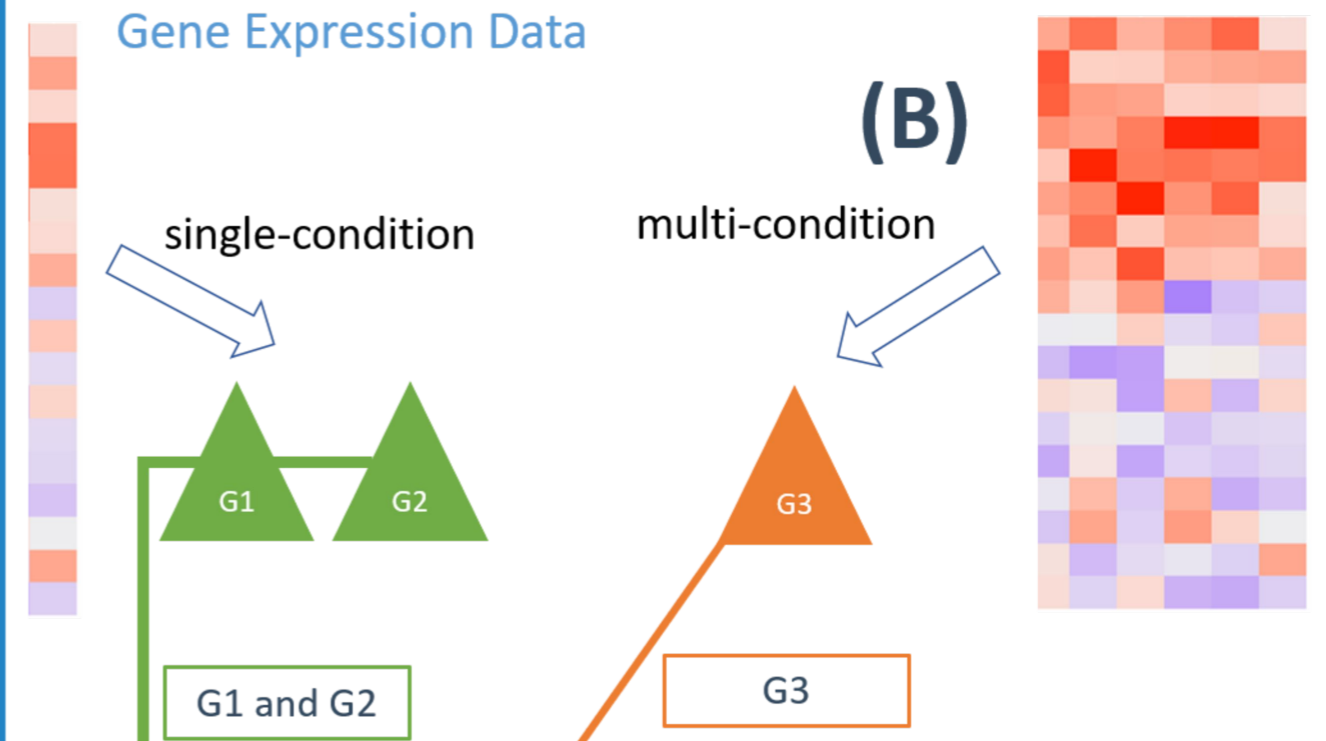
Regulatory Network

(A)



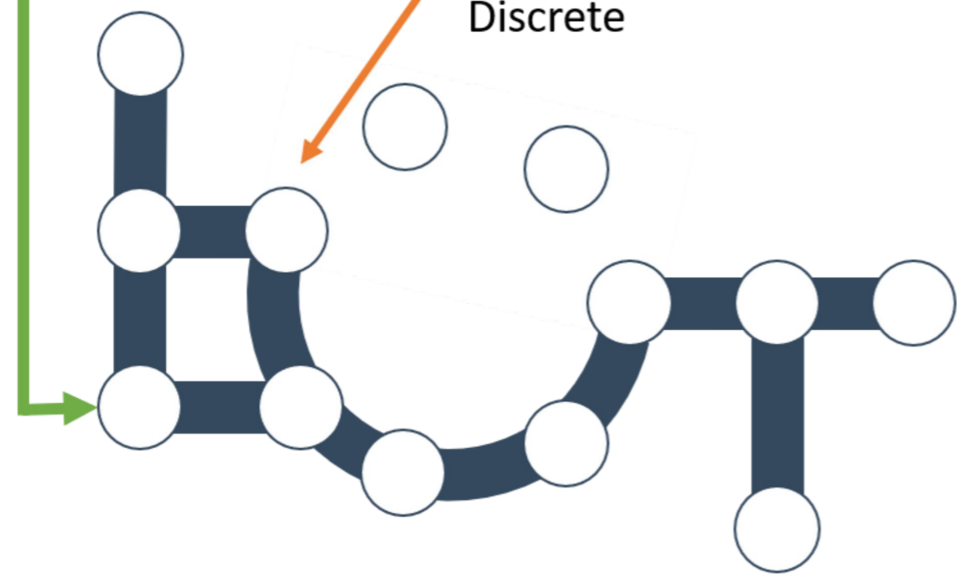
Gene Expression Data

(B)



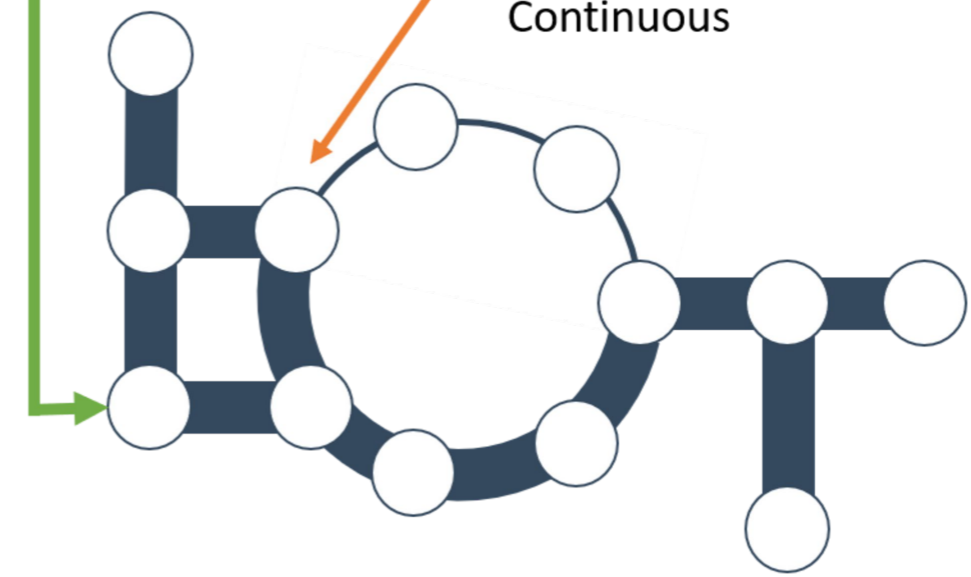
Absolute

Discrete



Metabolic Network

Continuous



Metabolic Network

