

1 *Conference Proceedings Paper*

2 **A poly-omics machine learning method to predict** 3 **metabolite production in CHO cells**

4 **Guido Zampieri¹, Macauley Coggins², Giorgio Valle¹ and Claudio Angione^{2,*}**

5 ¹ CRIBI Biotechnology Centre, University of Padova, viale G. Colombo 3, 35131 Padova, Italy;
6 guido.zampieri@phd.unipd.it

7 ² Department of Computer Science and Information Systems, Teesside University, Borough road, TS1 3BA
8 Middlesbrough, UK; c.angione@tees.ac.uk

9 * Correspondence: c.angione@tees.ac.uk; Tel.: +4401642342681

10 Academic Editor: name

11 Published: date

12 **Abstract:** The success of biopharmaceuticals as highly effective clinical drugs has
13 recently led industrial biotechnology towards their large-scale production. The
14 ovary cells of the Chinese hamster (CHO cells) are one of the most common
15 production cell line. However, they are very inefficient in producing desired
16 compounds. This limitation can be tackled by culture bioengineering, but
17 identifying the optimal interventions is usually expensive and time-consuming. In
18 this study, we combined machine learning techniques with metabolic modelling
19 to estimate lactate production in CHO cell cultures. We trained our poly-omics
20 method using gene expression data from varying conditions and associated
21 reaction rates in metabolic pathways, reconstructed *in silico*. The poly-omics
22 reconstruction is performed by generating a set of condition-specific metabolic
23 models, specifically optimised for lactate export estimation. To validate our
24 approach, we compared predicted lactate production with experimentally
25 measured yields in a cross-validation setting. Importantly, we observe that
26 integration of metabolic predictions significantly improves the predictive ability
27 of our machine learning pipeline when compared to the same pipeline based on
28 gene expression alone. Our results suggest that, compared to transcriptomic-only
29 studies, combining metabolic modelling with data-driven methods vastly
30 improves the automatised design of cultures, by accurately identifying optimal
31 growth conditions for producing target therapeutic compounds.

32 **Keywords:** CHO cell; Biopharmaceutical; Metabolic modelling; Machine learning;
33 Flux balance analysis.

34

35 **1. Introduction**

36 Chinese hamster ovary (CHO) cells are widely regarded as one of the most reliable
37 cell types for industrial-scale mammalian protein production. As compared to
38 bacterial cell lines such as those of *Escherichia coli*, CHO cultured cells are less

39 productive, much fragile and grow slowly. In turn, this means that the
40 manufacturing methods that facilitate protein production using CHO cell lines are
41 much more expensive and time-consuming. However, heavy interest is put in
42 optimising CHO cell lines as they are required to produce mammalian
43 recombinant proteins.

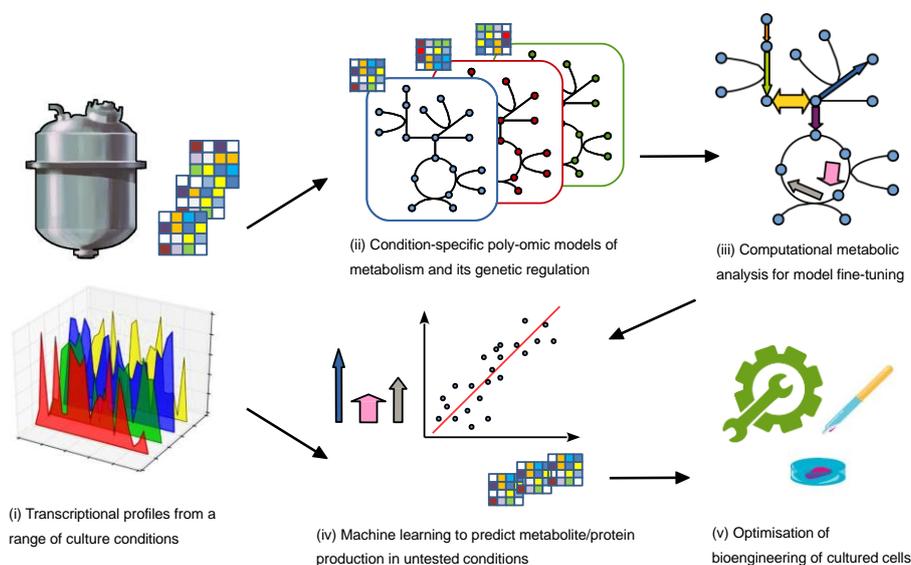
44 Recent advances in this context have focused on unraveling the complex
45 biological machinery controlling desirable characteristics of protein synthesis and
46 secretion [1]. While gene expression profiling has proved helpful in past studies,
47 there have been recent efforts to combine genetic data with knowledge of
48 metabolic pathways through the reconstruction of genome-scale metabolic models
49 (GSMMs). GSMMs attempt to describe cellular metabolism *in silico* through gene
50 annotation and stoichiometry associated with reactions and metabolites, as well as
51 with constraints such as upper or lower bounding of metabolic flux rates. Flux
52 balance analysis (FBA) allows to predict the configuration of metabolic reaction
53 fluxes within GSMMs under general growth conditions [2]. Condition-specific
54 GSMMs can be built using a variety of methods and extended FBA pipelines. The
55 idea is to use omic-data available in each condition, and a set of rules to constrain
56 the flux rates of the general-purpose GSMM [20,21].

57 Metabolic models have recently been reconstructed for CHO-K1, CHO-S, and
58 CHO-DG44 cell lines, along with a general consensus model [3]. These models
59 were useful in quantifying the protein synthesis capacity of these cell lines and
60 revealed that bioprocessing treatments such as histone deacetylase inhibitors' lead
61 to an inefficiency in increasing product yield. FBA can thus reveal the impact of
62 various media and culture conditions on growth and yield of cultured cells, aiding
63 CHO cells bioengineering [3-6]. Moreover, computational estimation of metabolic
64 fluxes can be an asset when experimental data is not available [7].

65 However, the precision of GSMMs strongly depends on available pathway and
66 biochemical knowledge. Especially when dealing with the complexity of
67 mammalian cells, more advanced computational techniques may be necessary for
68 an effective application to real problems within the bio-processing industry. In
69 particular, machine learning coupled with computational modelling of CHO cells
70 has the potential to effectively elucidate optimal bioengineering steps towards
71 improved production of therapeutic metabolites and proteins [8].

72 Here we present a new approach integrating machine learning and metabolic
73 modelling for the computational prediction of protein production in CHO cells. We
74 propose to integrate experimental data on the gene level with data generated *in*
75 *silico* via a GSMM of CHO cells metabolism within an integrated data-driven
76 framework (Figure 1). We evaluated this approach by a computational validation,
77 estimating the average prediction error in general settings. Importantly, we
78 observe that metabolic predictions coupled with gene expression data can
79 significantly improve estimations of lactate production based solely on gene
80 expression.

81
82



83

84 **Figure 1.** Workflow of the proposed approach for the prediction of metabolite and protein prediction in
85 CHO cells. Steps (i)-(iv) are presented in the Methods section of this work. They serve the final goal of
86 optimising culture bioengineering, depicted in step (v). Integrating transcriptomics data, machine learning
87 methods and metabolic modelling improves the predictive ability of transcriptomic-only methods.

88 2. Materials and Methods

89 2.1 Publicly available gene expression data

90 As a first data source, a large-scale gene expression dataset from two different
91 CHO cell lines was used [9]. This dataset contains 295 microarray profiles with
92 expression values for 3592 genes from 121 CHO cell cultures of varying conditions
93 in terms of including cell density, growth rate, viability, lactate and ammonium
94 accumulation and cell productivity. We extracted the 127 profiles with available
95 quantification of lactate accumulation.

96 2.2 Genome scale reconstruction of CHO metabolism

97 We used a recently developed GSMM of CHO cell metabolism, previously
98 used to accurately predict growth phenotypes [3]. This model is the largest
99 reconstruction of CHO metabolism to date, with 1766 genes and 6663 reactions,
100 aggregating community knowledge from various sources. Being a consensus
101 model, it provides general mechanistic relationships that can be refined depending
102 on the particular task or cell line of interest.

103 2.3 Building condition-specific poly-omics models of CHO cells

104 To create condition and cell line-specific poly-omics models the genome-scale
105 model of CHO cell metabolism was combined with the gene expression data from
106 CHO cell cultures in varying conditions. In this step, data accessible via the BIGG

107 repository was employed to match gene identifiers [10]. A model for each
108 condition was created by computing gene set effective expressions Θ for each
109 reaction, following previous investigations [11,12]. The effective expression at
110 reaction level is thereby determined by gene expressions $\theta(g)$ and by gene-protein-
111 reaction rules, properly converted to min/max rules depending on the type of gene
112 set. In particular, we define $\Theta(g) = \theta(g)$ for single genes, $\Theta(g_1 \wedge g_2) = \min\{\theta(g_1),$
113 $\theta(g_2)\}$ for enzymatic complexes and $\Theta(g_1 \vee g_2) = \max\{\theta(g_1), \theta(g_2)\}$ for isozymes.
114 Lower bounds and upper bounds for each reaction were obtained by applying the
115 following multiplicative coefficient to its native bounds:

$$\phi(\Theta) = [1 + \gamma |\log(\Theta)|]^{\text{sgn}(\Theta-1)}, \quad (1)$$

116 where γ is a parameter controlling the impact of gene expression on reaction
117 bounds.

118 *2.4 Extraction of metabolic features*

119 After a model for each condition was created, flux distributions were
120 computed using FBA by maximising the biomass for producing cell lines included
121 in the CHO model [3]. To perform FBA we employed the COBRA toolbox and a
122 multi-level linear program structure [13,24]. All simulations were carried out in
123 Matlab R2014b with the Gurobi solver.

124 *2.5 Feature processing and selection*

125 Principle Component Analysis (PCA) is a very effective statistical tool that uses
126 an orthogonal transformation to reduce a set of variables to a smaller set of linearly
127 uncorrelated variables, known as the principle components [14]. Here PCA was
128 used to process metabolic flux features in order to extract informative metabolic
129 features.

130 Moreover, elastic net was applied to select relevant features, both at a gene
131 expression and metabolic level [15]. Given an α in the interval]0, 1] and a non-
132 negative λ , elastic net solves the following optimisation problem:

$$\min_{\beta_0, \beta} \left(\frac{1}{2N} \sum_{i=1}^N (y_i - \beta_0 - x_i^T \beta)^2 + \lambda P_\alpha(\beta) \right). \quad (2)$$

133 In this formula, x represents the gene expression and metabolic flux rates variables,
134 y corresponds to measured metabolite yield and N is the total number of training
135 conditions. $P_\alpha(\beta)$ is a regularisation term depending on a vector of linear
136 coefficients β and on parameter α . Non-null entries of β resulting from this
137 minimisation correspond to relevant features selected by elastic net.

138 *2.6 Training generalised linear models to predict metabolite/protein production*

139 Generalised linear models (GLM) were trained to predict lactate yield starting
140 from poly-omics information [16]. A GLM gives an estimate of metabolite
141 production y_i^{pred} calculated as follows:

$$y_i^{pred} = \beta_0 + x_i^T \beta. \quad (3)$$

142 GLM accuracy was assessed by nested cross-validation, consisting of two
 143 cross-validation loops which together evaluate a selected model based on training
 144 data [17]. The nested loop selects the values of α and λ of elastic net on 5 training
 145 and test folds. The outer loop is used for model evaluation and is ran over 10 folds.
 146 GLM accuracy for each test fold was evaluated by computing the root-mean-
 147 square error (RMSE) defined by the following formula:

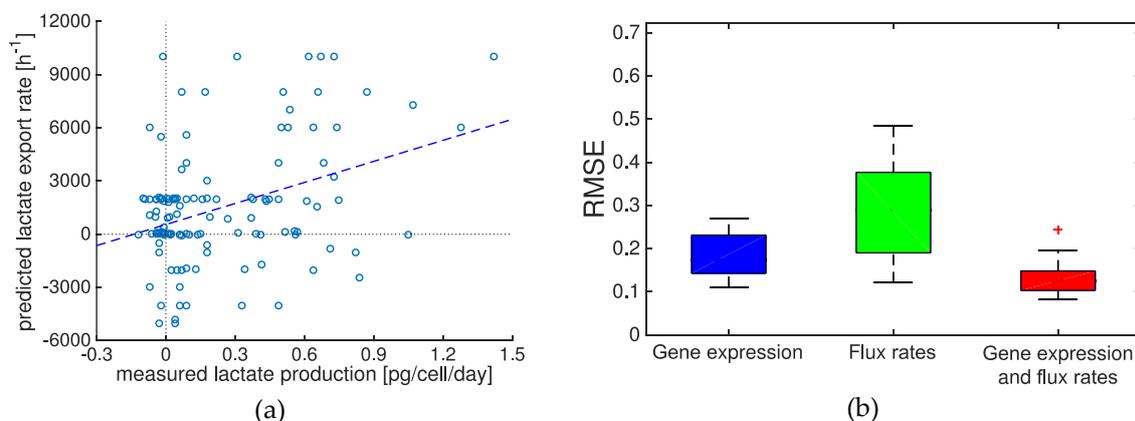
$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i^{pred} - y_i)^2}{n}}, \quad (4)$$

148 where n is the number of test conditions in the fold.

149 3. Results

150 3.1. Metabolic model optimisation

151 We validated our proposed approach on the prediction of lactate production,
 152 resorting to experimental data from the study of Clarke et al. [9]. We selected the
 153 conditions with both microarray and measured lactate production, obtaining 127
 154 conditions. In order to optimise metabolic flux information, we performed a
 155 sensitivity analysis on the gene expression mapping parameter γ in Equation (1).
 156 Specifically, we studied the Pearson correlation r between measured lactate
 157 accumulation in culture media and simulated lactate export rates for varying
 158 values of γ across several orders of magnitude. The maximum correlation
 159 coefficient obtained was $r = 0.36$ (p-value = $2.6 \cdot 10^{-5}$). The relationships between
 160 these two quantities can be visualised in Figure 2a. We thus employed condition-
 161 specific models with the optimal γ to generate fluxes for the following analysis.
 162



163 **Figure 2.** Validation results of the proposed approach on lactate production prediction: (a) comparison
 164 between simulated lactate export through condition-specific GSMMs and measured lactate production; this
 165 step enables GSMMs optimisation for the target metabolite in the following step; (b) RMSE distribution plots
 166 for lactate production predictions as a function of employed data sources. Two outliers for the green box lie
 167 outside of the current scale.

168 3.2 Predictions of lactate production

169 To accurately predict lactate production in CHO cells, we employed elastic net
170 and GLMs as described in the Methods section. We estimated the generalised
171 prediction error by means of a 10-fold cross-validation, repeatedly swapping
172 conditions used in training and in tests [17]. We calculated the RMSE of predicted
173 lactate yield across the test conditions in each fold, which quantifies the average
174 difference between predicted and experimentally measured lactate yield. We
175 repeated this procedure under three data sources scenarios, where gene
176 expression, metabolic fluxes and their combination was evaluated separately. The
177 results are shown in Figure 2b and summarised in Table 1. Interestingly, although
178 flux rates alone lead to poor predictions, if combined with gene expression they
179 achieve the minor average and median RMSE across the 10 folds. In the latter case,
180 associated RMSE distribution is significantly different to that obtained from gene
181 expression alone on the basis of a one-tailed Wilcoxon rank sum test at a 5%
182 threshold (p-value = 0.027) [18].

183

	Gene expression	Flux rates	Gene expression and flux rates
Mean RMSE	0.19	1.08	0.14
Median RMSE	0.17	0.26	0.13
RMSE standard deviation	0.06	2.41	0.05

184

185 **Table 1.** Comparison of 10-fold cross-validation RMSE statistics for the prediction of lactate production from
186 different data sources. Combining gene expression and metabolic flux data leads to best values for all
187 statistical measures. These results correspond to those shown in Figure 2b.

188

189 4. Discussion

190 The growing demand for natural products in global healthcare requires
191 advanced automation of CHO cell culture design for biotechnological industry to
192 reach commercial-scale production levels. Notably, recent advances in metabolic
193 modelling and in data-driven prediction algorithms have not been yet exploited in
194 combination for this purpose. In this study, we started to explore this research line:
195 the overall goal of the work was to develop a poly-omics approach capable of
196 predicting metabolite/protein production in CHO cells. The approach comprises a
197 GLM trained on gene expression data originating from cultures in varying
198 conditions and on metabolic flux rates obtained *in silico* from FBA on a GSMM of
199 CHO metabolism. The accuracy of our approach was evaluated in comparison to
200 GLMs employing only a single type of data. This allowed us to show that

201 combining gene expression and metabolic fluxes improves accuracy compared to
202 just using gene expression or metabolic fluxes separately.

203 Generation of condition-specific metabolic information can in principle be
204 achieved through various types of computational analysis. In this study, we used
205 FBA as this is the most widely used technique to capture flux configurations in a
206 growth steady state [2]. In principle, different techniques could potentially extract
207 even more useful information, further improving final data-driven predictions. For
208 instance, in a preliminary evaluation we tested also a modified version of
209 parsimonious enzyme usage FBA minimising the norm-2 of reaction fluxes [22,23].
210 However, we observed that normal FBA achieved best results (data not shown).

211 The main limitation of this work is represented by a scarce availability of large-
212 scale public data on CHO cells and by the prototypical state of present GSMMs.
213 Proposed strategies for model refining are expected to lead to further prediction
214 improvements [19]. With more comprehensive datasets, both in terms of number
215 of samples and in terms of metabolic gene coverage, we expect our pipeline to
216 vastly improve its predictive ability. Moreover, although our validation focussed
217 on lactate production, the proposed methodological framework can be
218 straightforwardly implemented around any target metabolite or protein.

219 Despite the above-mentioned limitations, our results show that metabolism-
220 based machine learning methods can significantly improve the predictive power of
221 common transcriptomic-only methods. This is due to the introduction of metabolic
222 features coupled with transcriptomic features. The present study therefore
223 represents a preliminary assessment that we plan to extend in future
224 investigations.

225

226 **Acknowledgments:** This work was partially supported by funding from BBSRC/EPSRC BioProNET. We thank
227 Jonathan Welsh from CPI-NBMC for helpful discussions about CHO cell products.

228 **Author Contributions:** C.A. and G.Z. conceived and designed the experiments; G.Z. and M.C. performed the
229 experiments; G.Z. analysed the data; C.A. G.Z. and G.V. contributed analysis tools; G.Z., M.C. and C.A. wrote
230 the paper. All authors read and approved the final version of the paper.

231 **Conflicts of Interest:** The authors declare no conflict of interest. The founding sponsors had no role in the
232 design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, and
233 in the decision to publish the results.

234 References

- 235 1. Richelle A. and Lewis N.E. Improvements in protein production in mammalian cells from targeted
236 metabolic engineering. *Curr. Opin. Syst. Bio.* **2017**, 6, 1-6, doi:10.1016/j.coisb.2017.05.019.
- 237 2. Orth J.D., Thiele I. and Palsson B.O. What is flux balance analysis? *Nat. Biotech.* **2010**, 28, 245-248,
238 doi:10.1038/nbt.1614.
- 239 3. Hefzi H., Ang K.S., Hanscho M., et al. A Consensus Genome-scale Reconstruction of Chinese
240 Hamster Ovary Cell Metabolism. *Cell Syst.* **2016**, 3(5), 434-443.e8, doi:10.1016/j.cels.2016.10.020.

- 241 4. Martínez V.S., Dietmair S., Quek L.E., Hodson M.P., Gray P. and Nielsen L.K. Flux balance analysis
242 of CHO cells before and after a metabolic switch from lactate production to consumption. *Biotechnol.*
243 *Bioeng.* **2013**, 10(2), 660-6, doi:10.1002/bit.24728.
- 244 5. Rejc Ž., Magdevska L., Tršelič T., Osolin T., Vodopivec R., Mraz J., Pavliha E., Zimic N., Cvitanović
245 T., Rozman D., Moškon M. and Mraz M. Computational modelling of genome-scale metabolic
246 networks and its application to CHO cell cultures. *Comp. Biol. Med.* **2017**, 88, 150-160,
247 doi:10.1016/j.combiomed.2017.07.005.
- 248 6. Pan X., Dalm C., Wijffels R.H. and Martens D.E. Metabolic characterization of a CHO cell size
249 increase phase in fed-batch cultures. *Appl. Microbiol. Biotechnol.* **2017**, 101(22), 8101-8113,
250 doi:10.1007/s00253-017-8531-y.
- 251 7. Sengupta N., Rose S.T. and Morgan J.A. Metabolic flux analysis of CHO cell metabolism in the late
252 non-growth phase. *Biotechnol. Bioeng.* **2011**, 108(1), 82-92, doi:10.1002/bit.22890.
- 253 8. Galleguillos S.N., Ruckerbauer D., Gerstl M.P., Borth N., Hanscho M. and Zanghellini J. What can
254 mathematical modelling say about CHO metabolism and protein glycosylation? *Comput. Struct.*
255 *Biotechnol. J.* **2017**, 15, 212-221, doi:10.1016/j.csbj.2017.01.005.
- 256 9. Clarke C., Doolan P., Barron N., Meleady P., O'Sullivan F., Gammell P., Melville M., Leonard M. and
257 Clynes M. Large scale microarray profiling and coexpression network analysis of CHO cells
258 identifies transcriptional modules associated with growth and productivity. *J. Biotechnol.* **2011**, 155(3),
259 350-359, doi:10.1016/j.jbiotec.2011.07.011.
- 260 10. King Z.A., Lu J.S., Dräger A., Miller P.C., Federowicz S., Lerman J.A., Ebrahim A., Palsson B.O. and
261 Lewis N.E. BiGG Models: A platform for integrating, standardizing, and sharing genome-scale
262 models. *Nucleic Acid Res.* **2016**, 44(D1), D515-D522, doi:10.1093/nar/gkv1049.
- 263 11. Angione C. and Lió P. Predictive analytics of environmental adaptability in multi-omics network
264 models. *Sci. Rep.* **2015**, 5, 15147, doi:10.1038/srep15147.
- 265 12. Angione C. Integrating splice-isoform expression into genome-scale models characterizes breast
266 cancer metabolism. *Bioinformatics* **2017**, btx562, doi:10.1093/bioinformatics/btx562.
- 267 13. Schellenberger J., Que R., Fleming R.M.T., Thiele I., Orth, J.D., Feist, A.M., Zielinski D.C., Bordbar,
268 A., Lewis, N.E., Rahmanian S., Kang J., Hyduke D.R. and Palsson B.O. Quantitative prediction of
269 cellular metabolism with constraint-based models: the COBRA Toolbox v2.0. *Nat. Prot.* **2011**, 6(9),
270 1290-1307, doi:10.1038/nprot.2011.308.
- 271 14. Jolliffe I.T. *Principal component analysis*, Series: Springer Series in Statistics, 2nd ed.; Springer, New
272 York, United states, 2002.
- 273 15. Zou H. and Hastie T. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B*
274 *(Stat. Methodol.)* **2005**, 67, 301-320, doi:10.1111/j.1467-9868.2005.00503.x.
- 275 16. McCullagh P. and Nelder J.A. *Generalized linear models*, 2nd ed.; Chapman and Hall, London, United
276 Kingdom, 1989.
- 277 17. Devijver P.A. and Kittler J. *Pattern Recognition: A Statistical Approach*; Prentice-Hall, London, United
278 Kingdom, 1982.
- 279 18. Hollander M. and Wolfe D.A. *Nonparametric Statistical Methods*; John Wiley & Sons, Inc., Hoboken,
280 United States, 1999.

- 281 19. Chowdhury R., Chowdury A. and Maranas C.D. Using Gene Essentiality and Synthetic Lethality
282 Information to Correct Yeast and CHO Cell Genome-Scale Models. *Metabolites* **2015**, 29;5(4), 536-70,
283 doi:10.3390/metabo5040536.
- 284 20. Vijayakumar S., Conway M., Lió P. and Angione, C. Seeing the wood for the trees: a forest of
285 methods for optimization and omic-network integration in metabolic modelling. *Briefings in*
286 *Bioinformatics* **2017**, bbx053, doi: 10.1093/bib/bbx053
- 287 21. Opdam S., Richelle A., Kellman B., Li S., Zielinski D.C., Lewis N.E.. A Systematic Evaluation of
288 Methods for Tailoring Genome-Scale Metabolic Models. *Cell Systems* **2017**, 22;4(3), 318-29, doi:
289 10.1016/j.cels.2017.01.010
- 290 22. Lewis N.E., Hixson K.K., Conrad T.M., Lerman J.A., Charusanti P., Polpitiya A.D., Adkins J.N.,
291 Schramm G., Purvine S.O., Lopez-Ferrer D., Weitz K.K.. Omic data from evolved E. coli are
292 consistent with computed optimal growth from genome-scale models. *Molecular systems biology* **2010**,
293 6(1):390, doi: 10.1038/msb.2010.47
- 294 23. Kim M.K., Lane A., Kelley J.J., Lun D.S. E-Flux2 and SPOT: validated methods for inferring
295 intracellular metabolic flux distributions from transcriptomic data. *PloS one* **2016**, 11(6):e0157101, doi:
296 10.1371/journal.pone.0157101
- 297 24. Angione C., Conway M., Lió P. Multiplex methods provide effective integration of multi-omic data
298 in genome-scale models. *BMC bioinformatics* **2016**, 17(4):83, doi: 10.1186/s12859-016-0912-1



© 2017 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC-BY) license (<http://creativecommons.org/licenses/by/4.0/>).