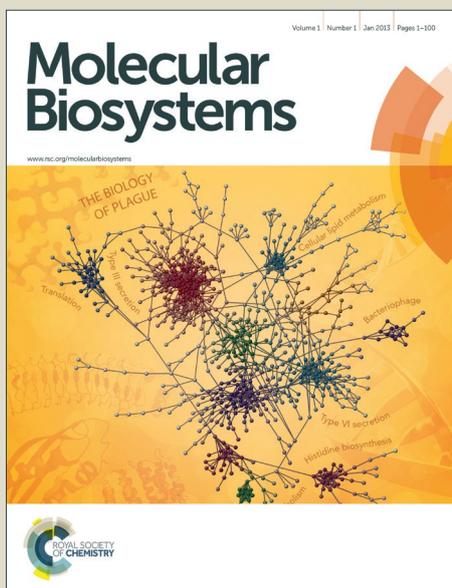


Molecular BioSystems

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.



www.rsc.org/molecularbiosystems

1 Improved metabolite profile
2 smoothing for flux estimation

3
4 Robert A. Dromms¹ and Mark P. Styczynski¹
5

6
7
8
9 ¹School of Chemical & Biomolecular Engineering, Georgia Institute of Technology, 311
10 Ferst Drive, Atlanta, GA 30332-0100
11

12
13 **Address for correspondence:**

14 Mark P. Styczynski

15
16
17 School of Chemical & Biomolecular Engineering, Georgia Institute of Technology, 311
18 Ferst Drive, Atlanta, GA 30332-0100. Tel: (404) 894-2825 Email:
19 mark.styczynski@chbe.gatech.edu
20

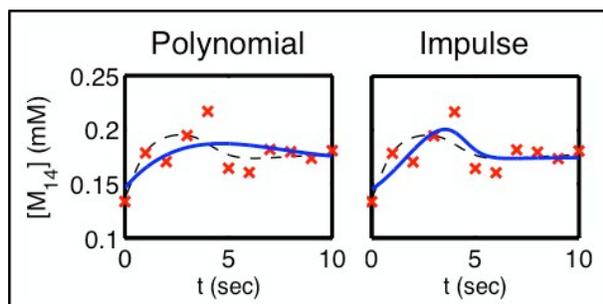
21 **Abstract** (250 words max):

22 As genome-scale metabolic models become more sophisticated and dynamic,
23 one significant challenge in using these models is to effectively integrate increasingly
24 prevalent systems-scale metabolite profiling data into them. One common data
25 processing step when integrating metabolite data is to smooth experimental time course
26 measurements: the smoothed profiles can be used to estimate metabolite accumulation
27 (derivatives), and thus the flux distribution of the metabolic model. However, this
28 smoothing step is susceptible to the (often significant) noise in experimental
29 measurements, limiting the accuracy of downstream model predictions. Here, we
30 present several improvements to current approaches for smoothing metabolite time
31 course data using defined functions. First, we use a biologically-inspired mathematical
32 model function taken from transcriptional profiling and clustering literature that captures
33 the dynamics of many biologically relevant transient processes. We demonstrate that it
34 is competitive with, and often superior to, previously described fitting schemas, and may
35 serve as an effective single option for data smoothing in metabolic flux applications. We
36 also implement a resampling-based approach to buffer out sensitivity to specific data
37 sets and allow for more accurate fitting of noisy data. We found that this method, as well
38 as the addition of parameter space constraints, yielded improved estimates of
39 concentrations and derivatives (fluxes) in previously described fitting functions. These
40 methods have the potential to improve the accuracy of existing and future dynamic
41 metabolic models by allowing for the more effective integration of metabolite profiling
42 data.

43 **Table of Contents entry**

44 We develop several methods to improve the estimation of metabolite concentrations
45 and accumulation fluxes from noisy time-course data, including use of a sigmoidal
46 impulse function and a resampling-based approach.

47



48

49 Introduction

50 Genome-scale metabolic modeling is an area of research with the potential for
51 significant impact on many biomedical and biotechnological applications. Such models
52 have been used to identify drug targets that specifically inhibit cancer proliferation¹, to
53 identify genomic manipulations that can facilitate production of valuable chemicals², and
54 to uncover and characterize metabolic pathways even in well-understood models³. This
55 approach entails using metabolic reconstructions that include all of the cataloged
56 metabolic reactions in an organism (i.e., genome-scale reconstructions) in a defined
57 mathematical modeling framework.

58 Effectively modeling biological systems at the genome scale calls for
59 measurements and data also at the genome scale. Metabolomics is the systems-scale
60 measurement of the small molecule intermediates in metabolism (the metabolites), a
61 field that has experienced rapid growth in the past decade. Modern analytical
62 technology enables the characterization of metabolic profiles in cells with increasingly
63 fine resolution; this provides relevant information to begin to replace steady state
64 assumptions on a genome-wide scale. However, to date, very few genome-scale
65 metabolic models have attempted to integrate metabolite profiling information, in
66 contrast to the prominent use of transcriptomic, fluxomic, and proteomic data in such
67 models⁴⁻⁸. In the few cases where metabolomics data have been integrated into these
68 models, the application of the data has typically been in setting thermodynamic
69 constraints and estimating free energies rather than in more direct applications^{9, 10}.

70 The primary reason for this omission is that most metabolic models using
71 genome-scale metabolic reconstructions assume the cell or organism to be at a steady
72 state, typically to simplify the model framework and associated computational
73 complexity. While models exploiting such an assumption have shown great utility, their
74 validity and potential for extrapolation have an intrinsic limit: while the steady state
75 assumption may be true over short time periods, it ultimately is violated once varying
76 forms of metabolic regulation begin to exert their influence.

77 The use of detailed ordinary differential equation (ODE) models would allow for
78 the capture of dynamic behaviors and regulation, but application of ODE models on a
79 genome-wide scale is not currently feasible due to (among other issues) the many
80 unknown reaction rate and thermodynamic parameters¹¹⁻¹³, each of which would require
81 extensive effort to be ascertained experimentally. As such, significant recent effort has
82 focused on softening the steady state assumption in genome-scale metabolic modeling
83 without requiring a full ODE model of the entire metabolic system^{5, 6, 14}. These efforts
84 hold great promise for future biotechnological applications, and they are the motivation
85 for the work presented here.

86 Use of metabolomics data is a promising approach for bridging the gap between
87 the steady state assumption and the dynamic intracellular reality. This data can be used
88 to estimate the accumulation or depletion “fluxes” of certain metabolites in a system,
89 which can then be used in place of the steady state assumption so common in genome-
90 scale metabolic modeling. This approach has been described and implemented in
91 multiple prior works¹⁵⁻¹⁹. The most common approach to estimating these accumulation

92 fluxes from metabolite data is to first smooth the data or fit it to a specific mathematical
93 function, and then use the resulting data or function to determine the flux of that
94 metabolite at any given time (potentially between measured time points). The accuracy
95 of these estimates has an obvious impact on the accuracy of the overall model, but
96 effective estimation of these fluxes is a non-trivial problem given the noise inherent to
97 measurement of metabolite levels and the limitations of the current methods for flux
98 estimation¹⁵.

99 One of the more thorough treatments of the problem of flux estimation from
100 metabolite data for metabolic modeling was included in work by Ishii *et al.*¹⁸ While the
101 main focus of that work was on developing a broader metabolic model, data smoothing
102 and flux estimation were integral parts of the data processing for the algorithm. They fit
103 a variety of polynomial and rational functions to simulated metabolite data and, on a
104 metabolite-wise basis, selected as the representative function the one that minimizes
105 the fitting error (accounting for the number of free parameters to minimize over-fitting).
106 Of note is that none of the candidate fitting functions are derived from or selected based
107 on biological insight. Additionally, as we show later, the fitting of an arbitrary dataset can
108 yield unphysical results. Splines, another common alternative, are sensitive to noise and
109 outliers—this is particularly problematic when the derivative of the concentration (the
110 accumulation flux) is the important quantity being estimated.

111 Here, we present two approaches for improving the estimation of accumulation
112 fluxes from metabolite time series data. First, we investigate the use of a biologically
113 reasonable and biologically-inspired sigmoidal impulse function^{20, 21} as an effective and

114 perhaps generalizable alternative to the fitting functions previously used. This functional
115 form emulates behavior observed in known biological systems, and our work represents
116 the first time that it has been applied in the context of metabolic modeling. Second, we
117 investigate whether a resampling-based approach to smoothing and fitting data might
118 yield more accurate concentration profile fits and derivative (flux) predictions than the
119 previously used approach. In the course of these investigations, we also identified the
120 importance of enforcing constraints on fitting equation parameter values to prevent the
121 selection of unphysical solutions. Each of these approaches improves the accuracy of
122 flux estimation from metabolite time series data, providing more reliable results to be
123 integrated into the larger metabolic modeling framework with reasonable computational
124 expense.

125

126 **Methods**

127 *Fitting functions*

128 Eight functions, shown in Table 1, were considered as candidates to best fit the
129 time series metabolite data. The first seven were used by Ishii *et al.*¹⁸. Four of these
130 were polynomials, of order two to five. The other three were rational functions,
131 composed of a first, second, or third order polynomial numerator and a first or second
132 order polynomial denominator. The eighth function was the sigmoidal impulse, which
133 was first presented in the context of filtering and clustering gene expression profiles^{20, 21};
134 it is here applied for the first time in the context of metabolic models. Unlike the other
135 functions, it has a biologically relevant interpretation: a two-phase transition from one

136 steady state to a (potentially new) steady state through an intermediate state. Its
137 parameters directly correspond to features of this trajectory, representing: transition time
138 delays; the initial, intermediate state, and steady-state metabolite levels; and the
139 sharpness of the transitions

140

141 *Synthetic Reference Data*

142 We tested our new methods using two different ODE models of central carbon
143 metabolism taken from the literature, which were used to generate noise-free “gold
144 standard” synthetic reference data for our analyses. These models were selected
145 because their dynamics are believed to reasonably represent *in vivo* metabolic
146 dynamics; the fact that they are not genome-scale does not detract from their relevance
147 as a model system, as the data smoothing/fitting step of flux estimation is independent
148 of the scale of the model.

149 The first model simulates central carbon metabolism in *E. coli*¹¹. While the model
150 includes 18 metabolites, only the 17 metabolites with substantial dynamics were
151 included in our analysis. (As implemented, metabolite 1 was a fixed value.) The second
152 model simulates central carbon metabolism in *S. cerevisiae*²², comprising 22
153 metabolites (21 of which had substantial dynamics, and were included in our analysis—
154 changes in metabolite 17 were several orders of magnitude smaller than the
155 concentration). While this model was initially presented in the context of stable
156 concentration oscillations, the initial conditions we used for our simulations do not

157 produce oscillatory behaviors. To validate our implementation of the model, we used it
158 to reproduce Fig. 6 from Hynne et al. (See Fig. S1)²².

159 We obtained curated SBML code for both models from the BioModels Database,
160 and solved systems of ODEs using the LSODA method in the Time Course module of
161 Copasi 4.14, Build 89, with the default tolerances and parameters^{23, 24}. For each model,
162 we solved the system of ODEs using the initial conditions specified in Table S1, derived
163 from those previously reported¹⁸, to simulate a perturbation in glucose concentration. As
164 previously described¹⁸, we used a perturbation from 0.0556 mM to 1.67 mM for
165 “Extracellular Glucose” in the *E. coli* model, and a perturbation from 2.5 mM to 5.0 mM
166 for “Mixed flow glucose” in the *S. cerevisiae* model. For the *E. coli* model, we fixed the
167 concentrations of ATP, ADP, AMP, NAD(H), and NADP(H) at their initial values, as was
168 done previously. The resulting gold-standard data contained concentrations at intervals
169 of 0.01 seconds for the *E. coli* model and 0.0025 minutes and for the *S. cerevisiae*
170 model.

171 To generate data for parameter estimation, simulated time points were sampled
172 at 1 second intervals from 0 seconds to 20 seconds for the *E. coli* model, and at 0.25
173 minute intervals from 0 minutes to 15 minutes for the *S. cerevisiae* model. The selection
174 of different sampling rates was to be consistent with the approach taken by Ishii *et al.*
175 for the *E. coli* model, but to account for the different time scales of the dynamics in the
176 two mathematical models as observed in the BioModels implementations while still
177 keeping the number of samples used for each respective model the same as that used
178 by Ishii *et al.* By keeping the number of samples the same as in previous work for each

179 respective model, our fitting results would be most directly comparable. We used a first-
180 order centered finite difference approximation on the ODE output to estimate the
181 derivatives in the synthetic reference data for each metabolite, C_i .

182

183 *Synthetic Noisy Data*

184 We generated sets of noisy metabolite time courses from this synthetic reference
185 data. For each metabolite C_i , we generated a noisy time course by adding noise at each
186 sampled time point, t_k , to the true value at that timepoint, $C_i(t_k)$, by drawing 5 simulated
187 measurements from a normal distribution, $N_{i,k} \sim (C_i(t_k), CoV \cdot C_i(t_k))$, and then taking
188 the mean of those 5 measurements, called $D_i(t_k)$. We refer to each individual noisy time
189 course as $D_{i,m}$. This approach paralleled the common experimental approach of taking
190 biological replicate measurements and then collapsing them into one value for analyses.
191 Here, we set the Coefficient of Variation (CoV) to 0.15, a reasonable value for many
192 mass spectrometry-based metabolite profiling approaches. The same noisy values were
193 used for all functions, allowing for direct comparison of the performance of each
194 function. In total, 500 noisy time courses were generated for each metabolite in each
195 model for the Direct Fit Method (described below), while an additional 50 time courses
196 were used as the base data for the Resampling Method (described below).

197

198 *Direct Fit Method*

199 We refer to a basic nonlinear least squares fitting of parameters as the “Direct
200 Fit” method for the purposes of this work. In this approach, we directly fitted each

201 function of interest to each noisy time course, $D_{i,m}$, to produce the smoothed time
202 course estimate, $f_{i,j,m}$. Best-fit parameters for a given function were selected by
203 minimizing the root-mean-square-displacement (RMSD) of the function to the data,
204 defined as

$$RMSD_{i,j,m} = \sqrt{\sum_k \frac{(D_{i,m}(t_k) - f_{i,j,m}(t_k))^2}{n - p_j}}$$

205 where i represents a specific metabolite, j represents a function being fitted, k
206 represents an individual time point, m represents the use of a specific noisy data set, n
207 is the number of sampled time points in the time course $D_{i,m}$, and p_j is the number of
208 parameters being fit for function f_j . The denominator reflects a penalty on the number of
209 parameters for a function, to help guard against over-fitting when comparing different
210 functions²⁵.

211 Polynomials were fit using the built-in `polyfit()` function in MATLAB. Rational
212 functions and the impulse function were fitted using `fmincon()` in MATLAB to allow for
213 bounds on the parameter space, as described in the Supplementary Methods (found in
214 Supplementary File 1). To improve the likelihood of finding globally optimal parameter
215 sets for the rational and impulse functions, we selected optimal parameters from 20
216 solver runs seeded with different sets of initial conditions (see Supplementary Methods).

217

218 *Resampling Method*

219 In an approach we refer to as the “Resampling Method”, we took advantage of
220 the stabilizing effect of calculating the median of fits to multiple noisy datasets to
221 produce more robust estimates of metabolite concentrations and derivatives.

222 Starting with the noisy time courses that model experimental data (described
223 above), we generated resampled time courses by repeating the procedure used to
224 produce the original noisy time courses, but using a noisy time course $D_{i,m}$ as input
225 rather than the true metabolite concentration C_i . We again used a fixed CoV of 15% for
226 this procedure; however, in practice, a dataset-specific and/or metabolite-specific CoV
227 could be estimated and use in place of the fixed CoV . We generated 250 such
228 resampled noisy time courses, $R_{i,m,w}$, for each initial noisy time course $D_{i,m}$.

229 We used the Direct Fit Method as described above to generate a nominal
230 parameter solution from each base noisy time course $D_{i,m}$. Then, for each resampled
231 time course $R_{i,m,w}$ derived from that noisy time course, we fit the function of interest
232 (once) using the parameter solution from the Direct Fit Method as the initial guess.
233 Parameter fitting was performed as described above.

234 We then used the resample-derived parameters to calculate concentration and
235 derivative trajectories for each resampled time course $R_{i,m,w}$, and calculated the median
236 value across all resampled time courses at the time points of interest (either the original
237 or interpolated time points, as described below). The output of the Resampling Method
238 was this list of concentration and derivative medians.

239

240 *Performance Calculations*

241 The performance of each fitting function using each method (direct and
 242 resampling) on both concentration and derivative predictions was quantified for each
 243 metabolite and for each base noisy time course, $D_{i,m}$. Concentration accuracy is useful
 244 for assessing the effectiveness of smoothing, while derivative accuracy is more relevant
 245 for downstream applications in estimating flux distributions¹⁷. Accuracy for each noisy
 246 time course $D_{i,m}$ was calculated using an adjusted RMSD between the synthetic
 247 reference data, C_i , and the predicted value for a given function, parameter set, and
 248 noisy data set, $f_{i,j,m}$. Specifically, we calculate accuracy as

$$RMSD_{i,j,m} = \frac{\sqrt{\sum_k (C_i(t_k) - f_{i,j,m}(t_k))^2}}{n_l \cdot S \cdot \mu}$$

249 where

$$S = \sqrt{\frac{\sum_k (f_{i,j,m}(t_k))^2}{n}}$$

$$\mu = \frac{n - p_j}{n}$$

250 and n_l is the number of time points used in assessing predictive accuracy, S is a scaling
 251 factor facilitating comparison and visualization by controlling for differences in the
 252 magnitude of different metabolites, and μ is a penalty factor scaling with the number of
 253 parameters in a function and the number of data points used to fit the function. For
 254 calculating derivative accuracy, the derivative values $f'_{i,j,m}(t_k)$ and $C'_i(t_k)$ are substituted
 255 in place of $f_{i,j,m}(t_k)$ and $C_i(t_k)$.

256 For these performance calculations, we more densely sampled metabolite
 257 concentration and derivative time courses to provide a more accurate representation of

258 interpolation performance, relevant to the general case of dynamic genome-scale
259 metabolic modeling. For each model, results were sampled at time steps a factor of ten
260 smaller than those used for the fitting data, resulting in $n_I = 201$ interpolated points for
261 the *E. coli* model and $n_I = 601$ interpolated points for the *S. cerevisiae* model (these
262 sets included the original sampled time points).

263 We ranked the functions' performance and averaged these ranks to provide a
264 quantitative overall comparison of each function. We ranked the performance of each
265 function for each noisy time course ($D_{i,m}$) of each metabolite and averaged the ranks for
266 each function across all of these time courses. In both cases, a harmonic mean was
267 used to average ranks, emphasizing the relative importance of comparing functions that
268 perform strongly in some cases; in this way, the difference between rank 1 and rank 2
269 was weighted more heavily than the difference between, for example, rank 4 and rank 5.

270 This averaged rank approach was used to compare performance of fitting
271 functions for the Direct Fit method only and for the Resampling Method only, as well as
272 to compare performance between these two methods for all of the different fitting
273 functions.

274 The MATLAB codes used to generate gold standard datasets, fit parameter
275 values, calculate metrics, and plot metrics, are collectively available in Supplementary
276 File 2.

277

278 **Results**

279 Two small-scale ODE metabolic models describing *E. coli* and *S. cerevisiae*
280 metabolism were used to generate synthetic reference data for the assessment of new
281 methods for concentration and flux inference from metabolite data. Using this synthetic
282 reference data as a basis, noisy time courses were generated to represent the noisy
283 data that typically result from metabolite profiling experiments. Eight different functions,
284 including four polynomials, three rational functions, and one impulse model function (as
285 described in the Methods section and in Table 1), were used as candidate fitting
286 functions for these noisy metabolite time course data. Two different approaches were
287 used to fit metabolite concentration curves to the noisy synthetic datasets generated
288 from the original ODE models.

289 The Direct Fit Method, described in the Methods section, was a standard fitting of
290 functions to given experimental data. The approach used to assess the effectiveness of
291 the Direct Fit Method for each of the candidate fitting functions is outlined in Fig. 1.
292 Briefly, after multiple noisy time courses were generated from the synthetic reference
293 data, each candidate function was fitted to each of the noisy time courses. Each of
294 these fits was then assessed for their performance at recapitulating and interpolating the
295 original data; these assessments were performed on both the fitted concentrations and
296 the derivative values that resulted from those fitted concentrations.

297 The Resampling Method, also described in the Methods section, involved fitting
298 multiple noisy datasets generated from a single experimental (or noisy synthetic)
299 dataset. By taking the median of these multiple fits, susceptibility to noise and outliers in
300 the original experimental data was reduced, providing more robust estimates of

301 metabolite concentrations and derivatives. The approach used to assess the
302 effectiveness of the Resampling Method for each of the candidate fitting functions is
303 outlined in Fig. 2. Briefly, multiple “base” noisy time courses were generated from the
304 original model to represent experimental measurements; these were fitted using the
305 Direct Fit Method for comparison. In parallel, additional noisy time course profiles were
306 generated (“resampled”) from each of these base noisy time courses and subsequently
307 fitted using the methods described for the Direct Fit Method—yielding a fitted
308 concentration for each resampled noisy time course for a given base noisy time course.
309 For each base noisy time course, the median per time point of the fitted profiles (or
310 profile derivatives) for the resampled noisy time courses was then used to determine the
311 overall fitted profile. This profile, along with the Direct Fit Method profile, was compared
312 to the original synthetic reference data to assess prediction accuracy.

313

314 *Parameter constraints improved the behavior of fitted results*

315 Fig. 3 provides representative examples of performance for different candidate
316 fitting functions using the Direct Fit Method and the *E. coli* model. Polynomial functions
317 provided computationally efficient data smoothing with little susceptibility to noise, but
318 had limited abilities to qualitatively capture the dynamics present in the *E. coli* model.
319 For certain sets of noisy data, the rational functions or the impulse function returned
320 unphysical or unreasonable results. This result highlighted a shortcoming in the basic
321 implementation of the rational functions and prompted the development of additional

322 constraints for use in the optimization step of fitting the rational functions and the
323 impulse function.

324 We observed that for approximately 29% of noisy datasets, the R_{22} rational
325 function produced asymptotic behavior, as shown in Fig. 3D. The frequency of
326 asymptote occurrence varied significantly across the different metabolites in the model,
327 as shown in Fig. S3A. The source of these asymptotes was selection of “optimal”
328 parameters such that the polynomial in the denominator of R_{22} had a root over the time
329 range of the data. Technically, such parameter selections would be optimal based on
330 the RMSD objective function, since the RMSD only considers the ability of the function
331 to match the data provided for fitting. However, such selections lead to clearly
332 unphysical profiles at interpolated points that would confound any efforts to use such
333 fitted functions in genome-scale metabolic simulations. Accordingly, we constrained the
334 RMSD optimization for all rational functions (as described in detail in the Supplementary
335 Methods, Fig. S3, and Table S4) such that parameters could not be selected that would
336 cause a zero in the denominator over the time range of the data. Fig. 3E shows the
337 trajectory of R_{22} after adding additional constraints to the allowed parameter values in
338 rational functions. However, this solution does not protect against near-asymptotic
339 behavior in R_{22} , where the denominator approaches but does not reach zero; Fig. 3F
340 depicts such a case using a different set of noisy data for the same metabolite.
341 Nonetheless, the results in Fig. 3E demonstrate significant improvement upon the
342 results from Fig. 3D with no parameter constraints.

343 The impulse function exhibited a similar phenomenon, insofar as it yielded results
344 that were technically correct based on the RMSD optimization function but were
345 physically unreasonable. As depicted in Fig. 3B, the impulse function sometimes
346 produced sharp shifts in concentration, which translated to sharp spikes in the derivative
347 trajectory. In addition, we noticed that our parameter-fitting solver was prone to getting
348 stuck in local minima when the resulting time delay parameters were outside the time
349 span of the data. These observations led us to implement an additional parameter
350 constraint strategy described in more detail in the Supplementary Methods.

351 Briefly, one fixed constraint and two new adjustable optimization parameters
352 were created that were used to constrain the possible parameter space. Since any
353 arbitrary dataset would not provide evidence for a sigmoidal shift outside of the time
354 range of the data, we constrained the possible sigmoidal response times to only be
355 within the time range of the data. We then defined two parameters, h_f and b_f , to further
356 constrain the parameter space based on the data. Since an arbitrary dataset would not
357 provide evidence for initial steady state, intermediate state, and final steady state levels
358 far outside of the range of the measured metabolite concentrations, the deviation of
359 function values above the maximum and below the minimum measured values was
360 constrained to be no more than h_f times the range of the metabolite data (with an
361 additional non-negativity constraint). Since an arbitrary dataset would not provide
362 evidence for concentration changes at a higher frequency than that of the sampling
363 frequency, sharp transitions between time points are unlikely to be realistic. Thus, the
364 steepness of the sigmoidal shift was constrained to be less than a value proportional to

365 the range of the data divided by the time difference between data points, with b_f as the
366 proportionality constant.

367 Using $h_f=0.1$ and $b_f=0.5$ resulted in more realistic profiles like those shown in Fig.
368 3C. Importantly, in addition to the direct physical interpretation of these, the results of
369 the parameter fitting are not highly sensitive to small changes in h_f and b_f (see Fig. S4),
370 and as a result the values of h_f and b_f that we used were generalizable to both model
371 systems even though they were selected only based on their performance for the *E. coli*
372 model.

373

374 *The impulse model consistently fits metabolite data with low error*

375 To quantitatively assess the effectiveness of the candidate fitting functions using
376 the Direct Fit Method in the *E. coli* model, we generated 500 noisy time course data sets
377 for each of the 17 metabolites. The parameters resulting from fitting each noisy time
378 course were used to calculate concentration and derivative trajectories, with the
379 corresponding performance accuracy calculated and averaged as described in the
380 Methods section. The results of these calculations are summarized in Table 2, which
381 presents the averaged ranks for each function and each metric. Fig. 4A and 4B provide
382 a detailed quantitative comparison of each fitting function. The impulse function, I,
383 showed the best rank averages for accuracy in both concentration and derivatives, and
384 was almost always the best-performing function across all of the metabolites.

385 The notable exceptions to the superior performance of the impulse function were
386 on Metabolites 12 and 18. Fig. 5 summarizes the performance of the impulse function

387 and an average fitting function, P_4 , for Metabolite 12, with representative fitted profiles in
388 Fig. 5A and 5B, and a direct comparison between the performance of P_4 and I in Fig.
389 5C. P_4 consistently performed better than I . However, as is clear from Fig. 5A and 5B,
390 the total change in metabolite level was smaller than the expected range of variability of
391 experimental measurements. Given the sparsity of samples, this metabolite's profile is
392 likely essentially unidentifiable, and so the performance of the different functions is likely
393 based only on general trends of the functional forms near the ends of the time range,
394 rather than any reliably accurate fitting.

395

396 *The Resampling Method can improve fitting and predictions in the E. coli*
397 *model*

398 To quantitatively assess the performance of the Resampling Method in the *E. coli*
399 model, we generated 50 noisy time courses from the synthetic reference data for each
400 of the 17 metabolites, and for each noisy time course, an additional 250 resampled
401 noisy time courses. For each noisy and resampled time course, each function was fitted
402 as described in the Methods, and the resulting Direct Fit or Resampling Method
403 trajectories used to calculate the performance metrics. The overall results are shown in
404 Table 3. Results jointly ranking the performance of functions across both the Direct Fit
405 Method and the Resampling Method are shown in Table 4. The Resampling Method
406 had the greatest impact on the ranking of the rational function R_{22} , resulting in it being
407 similar in accuracy and consistency to the impulse function, I . This consistently good

408 performance is also evident in Fig. 4C and 4D, which provide a detailed quantitative
409 comparison of each fitting function.

410 The impacts of the Resampling Method varied across the different types of
411 functions; representative graphs are presented in Fig. 6, with a complete summary
412 provided in Table 4. Polynomial functions showed little to no change in results from
413 using the Resampling Method, while rational functions show moderate to noticeable
414 benefit. The impulse function benefited in some cases as well. Across all functions, use
415 of the Resampling Method only infrequently caused decreased performance, and
416 typically with very small changes relative to the magnitude of the error.

417

418 *S. cerevisiae* model results show similar trends

419 We then quantitatively assessed the performance of all candidate fitting functions
420 using both the Direct Fit Method and the Resampling Method in the *S. cerevisiae* model.
421 We generated 500 noisy time courses for each of the 21 metabolites for use in the
422 Direct Fit method. For use in the Resampling Method we generated 50 base noisy time
423 courses for each of the 21 metabolites, along with an additional 250 resampled noisy
424 time courses for each base noisy time course. Parameters were fit for each method as
425 described in the Methods section. Tables 5 and 6 present the average ranks for the
426 Direct Method and Resampling Method, both separately and combined, respectively.
427 Fig. 7 provides a detailed quantitative comparison of each fitting function. For this
428 model, the R_{22} rational function and the impulse function, I , were usually among the

429 best-performing fitting functions, with R_{22} performing best for concentrations and I
430 performing best for derivatives.

431

432 **Discussion**

433 The goal of this work was to improve the prediction of concentration and
434 derivative time-course profiles derived from experimentally measured (or synthetic,
435 noisy) metabolite data. Two small-scale model metabolic systems were used as the
436 basis for assessing the performance of new methods to calculate and interpolate
437 concentration and flux values based on metabolite data. These two models have
438 different time scales and dynamics, which provided a broader assessment of the
439 potential utility of our approaches. These models were also used in previous work on
440 estimating flux distributions from metabolite data¹⁸, which allowed for direct comparison.
441 Integrating these systems numerically provided an exact reference dataset to which we
442 could compare fitted results. However, real metabolite concentration data contain
443 significant variability, so we only used noisy synthetic data derived from this reference
444 data to test the effectiveness of our approaches. In this way, we were able to generate
445 data of defined quality and arbitrary quantity with known underlying dynamics; this
446 allowed us to precisely and rigorously determine the performance of each approach
447 under study.

448 The approach of Ishii *et al.* was to fit all of the functions to the time course in
449 question and select the function with the lowest fitting error, once accounting for the
450 number of fitted parameters¹⁸. While this is certainly a viable approach that can be

451 extended to include the sigmoidal impulse model, here we have also investigated
452 whether this single, biologically reasonable function can be used instead of selecting the
453 best-fitting function from a list of arbitrary candidates. We consider the relative benefits
454 of each function type below.

455

456 *Polynomials are consistent but inaccurate*

457 The polynomial functions are computationally inexpensive to fit, use few
458 parameters (ranging from three to six), and are widely used for smoothing noisy data.
459 They are consistent and well-behaved, exhibiting very little sensitivity to noise. (As
460 described in Supplementary Methods and Tables S2 and S3, robustness of smoothed
461 profiles to noise was also assessed, but was found to closely depend on the number of
462 parameters used in a function and essentially represented a tradeoff between
463 consistency and accuracy of fitting.) As demonstrated by their ranks in Tables 2, 3, and
464 5, they can do a reasonable job in estimating concentrations and at times even in
465 estimating derivatives (ranking as low as 2.5 but often closer to 3.5 or 4). However, they
466 are ill-suited to capturing dynamics that include a terminal steady state, particularly
467 since their functional form requires them to be monotonically increasing or decreasing at
468 the ends of the time range; this also makes them a poor choice for even limited
469 extrapolation.

470

471 *Resampling improves rational function accuracy*

472 The rational functions (using three to five parameters) can exhibit a wider range
473 of behaviors than the polynomials with the same number of parameters, and it has been
474 reported that for many metabolite time courses, they yield better performance than the
475 polynomials¹⁸. Our parameter restriction strategy was largely effective in addressing
476 their potential to fit best with parameters that produce asymptotic behavior, though there
477 are still lingering issues with near-asymptotes that yield spurious behavior and even
478 negative concentrations for the R_{22} function (see Fig. 2F). However, as shown in Table
479 3, this effect is largely ameliorated by the use of the Resampling Method to filter out
480 asymptotic trajectories, making R_{22} one of the more effective functions we studied.

481

482 *The impulse function is a generally effective single fitting function model*

483 The last function, the sigmoidal impulse, is the product of two sigmoidal logistic
484 functions^{20, 21}. As previously stated, it recapitulates the dynamics of a common
485 biological process: a two-phase transition from one steady state to a (potentially new)
486 steady state through an intermediate state. Its parameters directly correspond to
487 features of this trajectory: the h parameters represent the initial, intermediate, and
488 steady-state metabolite levels; the τ parameters represent the timing of the on and off
489 transitions (accumulation and depletion driven by processes such as synthesis and
490 degradation) in response to a perturbation; and the β parameters represent how rapidly
491 those transition processes occur. In contrast with the work done by Chechik *et al.*, we
492 allowed the β parameters to vary independently to reflect the fact that the on and off

493 transitions can represent different biological processes (e.g., glucose uptake versus
494 metabolism), which one would reasonably expect to exhibit distinct dynamics²⁰.

495 While potentially exhibiting undesirable behaviors with unrestricted parameter
496 values, our parameter bounding strategies for avoiding broad local minima and overly
497 sharp curves were effective at preventing these undesirable behaviors (Fig. 3B and 3C).
498 Of particular note is that these parameters themselves typically exhibited broad local
499 optima in performance (Fig. S4), meaning that the fitting method was not very sensitive
500 to the specific values selected; additionally, the default parameters we selected for the
501 *E. coli* model generalized well to a completely separate model, meaning that while they
502 are technically adjustable parameters, they did not add significant risk of over-fitting to
503 the parameter selection process.

504 Using the Direct Fit Method for the *E. coli* model, the impulse function performed
505 consistently better than other functions (see Table 2) across all metabolites except for
506 two: metabolites 12 and 18. For these metabolites, the actual dynamic range of
507 metabolite concentrations in the synthetic reference data was substantially less than the
508 range of the random noise used to construct the noisy time courses (see Fig. 5). We
509 cannot realistically expect to recover the underlying concentration in this case without
510 either much more dense or much more accurate sampling. We suspect that the better
511 performance of the polynomials was due in part to their tendency to swing upwards or
512 downwards near the edges of the data, which captured the early time dynamics of each
513 of these metabolites well; we note that the other high-performing fitting function, R_{22} , did
514 poorly on these metabolites as well. The Resampling Method substantially improved the

515 performance of R_{22} and slightly improved the performance of the impulse function on
516 these metabolites (Fig. 4), leading to qualitative behavior where the derivative effectively
517 fluctuated around zero. Given the lack of statistically significant change over the time
518 course of these metabolites, we argue that this is the behavior we should not only
519 expect, but actually be seeking given the essentially unidentifiable change in metabolite
520 levels.

521

522 *The Resampling Method generally improves on Direct Fit Method results*

523 In general, the resampling method ranged from negligibly detrimental to highly
524 beneficial. In a few cases, a very minor loss of performance was observed.

525 Consistently, resampling provided no benefit to polynomials (Fig. 6A); this is to be
526 expected, since the polynomial functions are already insensitive to small changes in the
527 data. The R_{11} and R_{31} rational functions saw minor improvements in general, while the
528 impulse function saw improvements in cases where it performed most poorly (Fig 5C).
529 The Resampling Method had the biggest effect on R_{22} ; in the *E. coli* model, it moved
530 from one of the worst performers to one of the overall best (Fig. 4, Table 4). Generally
531 speaking, then, the Resampling Method seems to be an effective way to improve
532 accuracy at only a mild computational cost.

533 The Resampling Method appears to have an effect similar to parameter
534 regularization by avoiding over-fitting due to noisy data²⁶. However, we note that the
535 Resampling Method returns a median of multiple fits, rather than a single parameter set.
536 As a result, concentration and derivative values derived from this method need not

537 strictly adhere to the functional form of the smoothing function; this flexibility can allow
538 better approximation of the underlying data in cases where the form of the particular
539 function happens to be biased against the correct behavior.

540

541 *S. cerevisiae* model results generally recapitulate *E. coli* model results

542 The *S. cerevisiae* model generally recapitulated results from the *E. coli* model,
543 demonstrating the potential generalizability of the Resampling Method and the impulse
544 function (including the parameters used to restrict the fitting search space for the
545 impulse function). For both the Direct Fit and Resampling Methods, the impulse function
546 performed fairly well. One feature that distinguished the *S. cerevisiae* model from the *E.*
547 *coli* model was the wider range of time scales present in the model's dynamics. Several
548 metabolites (1-4,8-10,18-20) reached steady-state in several minutes, while others
549 (12,13,14) took tens of minutes, and as a result did not reach steady-state during the
550 time interval of the data. As the impulse function assumes long-term steady-state
551 behavior for the time course, it did not perform as strongly for the Direct Fit Method for
552 these metabolites. However, the Resampling Method did provide some improvement for
553 these metabolites.

554

555 *Selection of fitting functions should be driven by applications*

556 In this work we considered the problem of data smoothing specifically in the
557 context of genome-scale metabolic modeling. Two key factors in this application have
558 driven our assessment of function and method performance. First, we expect that we

559 may need to provide flux values at points other than those for which experimental
560 measurements are available (for instance, if a genome-scale model entails something
561 akin to a Runge-Kutta numerical integration). This means that function accuracy should
562 be assessed not only at the sampled points, but in between them as well. Without the
563 inclusion of such interpolated values, some differences can be seen in apparent
564 effectiveness; for example, previous work indicated that polynomials were more
565 frequently optimal for the *S. cerevisiae* model¹⁸, but in terms of practical applications
566 they are usually inferior to R_{22} and the impulse function. Second, the main application of
567 the metabolite concentration smoothing is for the estimation of metabolite fluxes; this
568 means that while recapitulating the concentration profile is important, the more directly
569 applicable metric is how accurate the derivative profile is. This distinction is most
570 relevant for the *S. cerevisiae* model, where R_{22} more accurately recapitulates
571 concentrations, but the impulse model more accurately recapitulates the derivatives that
572 will be used in downstream analyses.

573

574 *Single functions and biologically-inspired functions can be effective fitting*
575 *models*

576 While previous work selected the best-fitting of an essentially arbitrary set of
577 functions for each individual metabolite based on the experimental data, we suggest
578 that this may be a suboptimal approach. First, this increases the likelihood for over-
579 fitting; it is difficult to estimate the number of effective parameters that are introduced to
580 the system by allowing for the variable selection of seven different models, but it suffices

581 to say that the number of effective parameters is likely greater than the number of
582 explicit parameters in the highest-order polynomial. As such, restricting the fitting to one
583 function may be desirable from an information content perspective; both the R_{22} and
584 impulse functions seem like reasonable, viable candidates for universal fitting functions.
585 In fact, once the assessment metrics are based on a criterion more reasonable for the
586 application (i.e., inclusion of interpolated points), there are few if any cases where the
587 polynomials would be a desirable option. Second, there is inherent value in using
588 biologically-inspired fitting functions. These functions, by design, recapitulate behaviors
589 previously observed in biological systems; biasing the fit towards these results
590 integrates prior knowledge that may help ensure that the model is closer to the
591 underlying biology. Even though there are more parameters in these functions than the
592 polynomials, the space of characteristic curves that can be fit is more restrictive and
593 more relevant to expected biology, partially mitigating concerns about over-fitting due to
594 excess parameters. In this sense, the impulse function may be the most desirable
595 choice; either way, applying the Resampling Method ensures that the smoothing and
596 fitting is improved over previous approaches.

597

598 *Limitations*

599 There are a few limitations to our analyses that bear noting. First, the number of
600 variable parameters in the impulse function places a lower limit on the number of
601 samples needed to fit the function well, which could stretch the experimental feasibility
602 of acquiring a sufficient number of samples. However, our analyses have been

603 consistent with previous work in terms of the number of samples used, and considering
604 the possibility of using multiple biological replicates and multiple experiments to fit the
605 same data, obtaining one or two dozen samples is often reasonable for a metabolomics
606 experiment. Second, the impulse model assumes a steady state is reached at the end
607 of the experiment, which may not be valid for all datasets. However, this concern is
608 partially mitigated by the fact that many experiments would actually be continued until
609 something more closely resembling a steady state is reached, minimizing the number of
610 times significant non-zero derivatives were present at the end of the time range. There
611 is also an obvious computational cost to fitting non-linearizable functions (as opposed to
612 polynomials) and to applying the Resampling Method; however, since the data
613 smoothing task is ultimately performed just once, not many times, we believe that the
614 improvement in results is worth this computational cost, which is itself reasonable and
615 does not require parallelization or even particularly long runtimes. Finally, we have not
616 analyzed the ultimate downstream impacts in the genome-scale metabolic modeling
617 application of the improvements we have made to assess their magnitude. Based on the
618 tendency of functions like polynomials to have nonzero derivatives at the end of the time
619 range and the importance of being able to capture a steady state in a metabolic model,
620 we expect that these improvements may be important, but will be to some extent model-
621 specific and is thus beyond the scope of this work. Either way, it is often generally
622 accepted that optimization of each intervening analysis or data processing step is
623 desirable for complex modeling schema.

624

625 **Conclusions**

626 In this work, we have demonstrated two improvements to standard approaches to
627 smooth metabolite concentration data for application to genome-scale metabolic
628 modeling, including a Resampling Method to minimize susceptibility to experimental
629 noise and the establishment of a single, biologically-inspired fitting function that
630 performs well in almost all cases. In the course of this work, we also identified additional
631 constraints that should be applied to existing data smoothing fitting functions to increase
632 their robustness and activity. Taken together, these contributions have provided
633 consistent and substantial improvements in existing methods to smooth and fit
634 metabolite data for downstream applications, whether via a new fitting function or
635 improvements made to existing fitting functions. We have shown these results to be
636 generalizable across multiple models of metabolism, suggesting the potential for
637 general utility of these improved methods to improve the accuracy of flux distributions
638 calculated from the derivatives of their time courses.

639

640 **Acknowledgements**

641 RAD participated in the design of the study, carried out the computational experiments,
642 and helped to draft the manuscript. MPS conceived of the study, participated in its
643 design and coordination, and helped to draft the manuscript. All authors read and
644 approved the final manuscript. RAD was supported by NSF IGERT award # DGE
645 0965945, and MPS and RAD were supported by NSF award # 1254382. We would also
646 like to thank McKenzie Smith and Amy Su for their feedback on the manuscript draft.

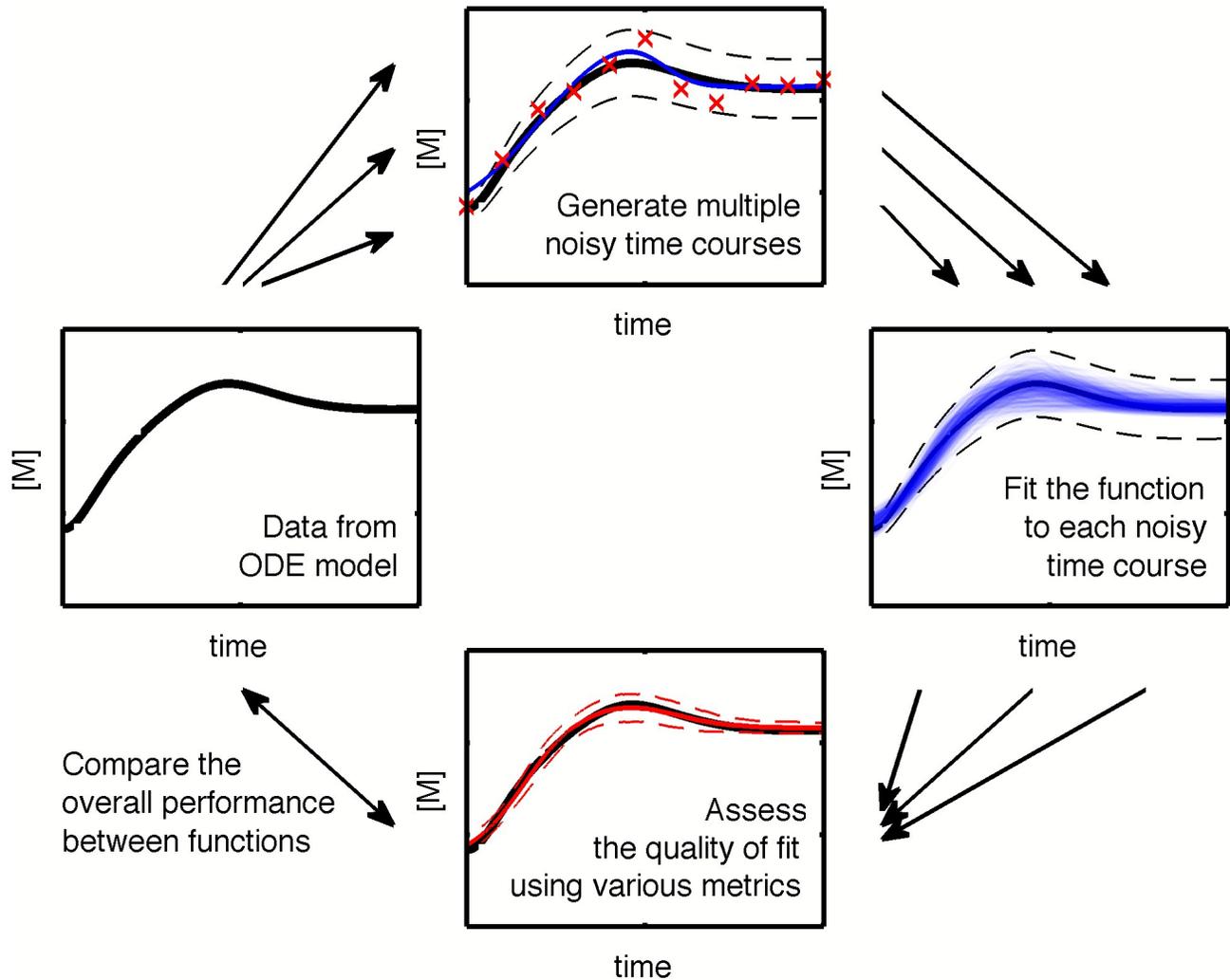
647

648 Table 1. Fitting functions evaluated in this work.
649

| Name | Formula |
|-----------------|---|
| P ₂ | $C(t) = p_1 \cdot t^2 + p_2 \cdot t + p_3$ |
| P ₂ | $C(t) = p_1 \cdot t^3 + p_2 \cdot t^2 + p_3 \cdot t + p_4$ |
| P ₄ | $C(t) = p_1 \cdot t^4 + p_2 \cdot t^2 + p_3 \cdot t^2 + p_4 \cdot t + p_5$ |
| P ₅ | $C(t) = p_1 \cdot t^5 + p_2 \cdot t^4 + p_3 \cdot t^3 + p_4 \cdot t^2 + p_5 \cdot t + p_6$ |
| R ₁₁ | $C(t) = \frac{p_1 \cdot t + p_2}{t + p_3}$ |
| R ₂₂ | $C(t) = \frac{p_1 \cdot t^2 + p_2 \cdot t + p_3}{t^2 + p_4 \cdot t + p_5}$ |
| R ₃₁ | $C(t) = \frac{p_1 \cdot t^3 + p_2 \cdot t^2 + p_3 \cdot t + p_4}{t + p_5}$ |
| I | $C(t) = \frac{1}{h_1} \cdot s(t, \tau_1, h_0, \beta_1) \cdot s(t, \tau_2, h_2, \beta_2)$ $s(t, \tau, h, \beta) = h + \frac{(h_1 - h)}{1 + e^{-4\beta(t-\tau)}}$ |

650

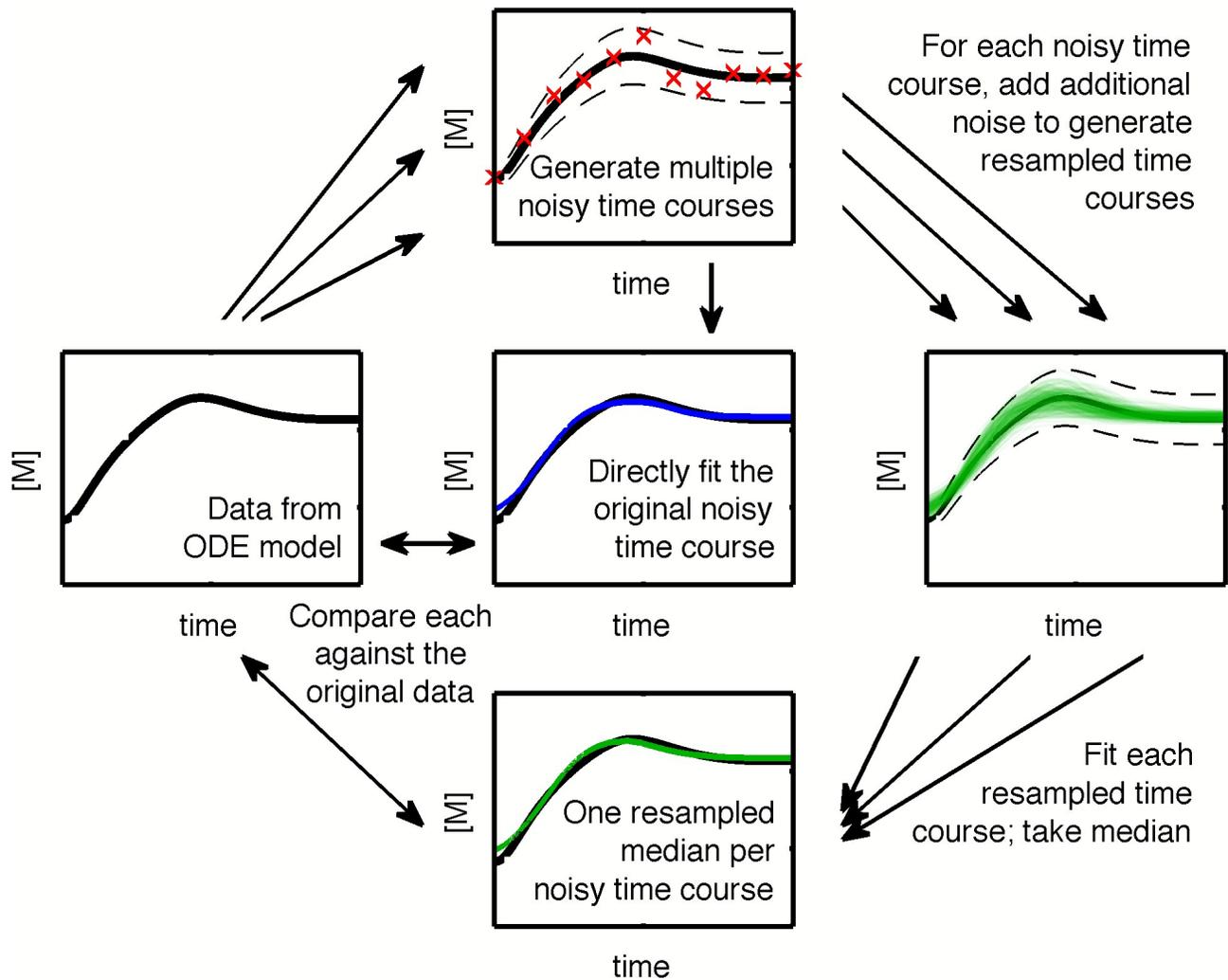
651 Fig. 1. Schematic of the Direct Fit Method.
 652 Synthetic gold standard data are generated by simulating a system of ODEs over the
 653 time interval of interest. From the synthetic data, noisy time courses are generated by
 654 adding Gaussian noise with a 15% coefficient of variation to the synthetic data, to
 655 simulate experimental sources of variation in measurements. Multiple such noisy time
 656 courses are generated. A smoothing function is fit directly to a noisy time course, and
 657 the resulting fit (or its derivative) is compared against the synthetic data to determine
 658 how closely they match. The performance of each function can then be compared
 659 based on their performance relative to the initial synthetic data.
 660



661
 662

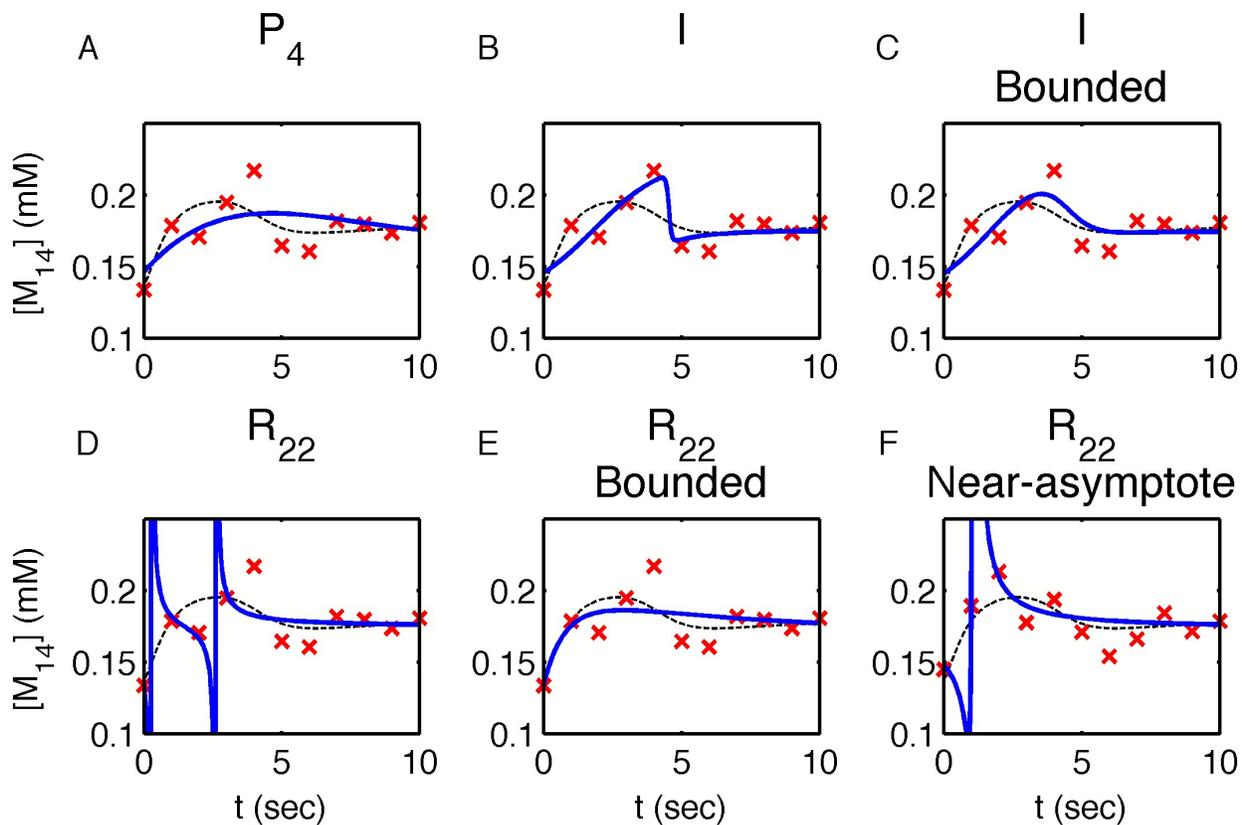
663 Fig. 2. Schematic of the Resampling Method.

664 As in the Direct Fit method, synthetic data and base noisy time courses are generated
 665 from a system of ODEs. In the Resampling Method, each base noisy time course is then
 666 used to generate a set of "Resampled" time courses, by using the same process used to
 667 generate the base noisy time courses from the synthetic data, only now with the base
 668 noisy time course as the input. The function of interest is fit to each of these resampled
 669 time courses, and the median of these functions (or their derivatives) is used to
 670 generate the resulting smoothed time course corresponding to the specific base noisy
 671 time course. As in the Direct Fit method, these median profiles can be assessed to
 672 determine accuracy and performance of the function.
 673



674
 675

676 Fig. 3. Performance of different fitting functions for fitting concentration trajectories.
 677 Thin, dotted black lines are the original synthetic data. Red crosses are the noisy time
 678 course data used to fit the functions. Solid blue lines are the function fitted to the data.
 679 A) Polynomial curves were consistent but typically not very accurate. B) The sigmoidal
 680 impulse function performed well but sometimes exhibited steep derivatives. C)
 681 Constraining the parameter space for the impulse function prevented this behavior. D)
 682 The rational function R_{22} can exhibit unphysical asymptotes in the time interval of the
 683 data due to a polynomial term in the denominator. E) Constraining the parameter space
 684 for R_{22} prevents such asymptotes. F) However, near-asymptote behavior can still occur
 685 in the rational functions, despite the parameter restrictions, when the value of the
 686 denominator polynomial becomes sufficiently small. Note: A-E all use the same noisy
 687 data set.
 688



----- ODE Time Course × Noisy Data — Fitted Time Course

689
690

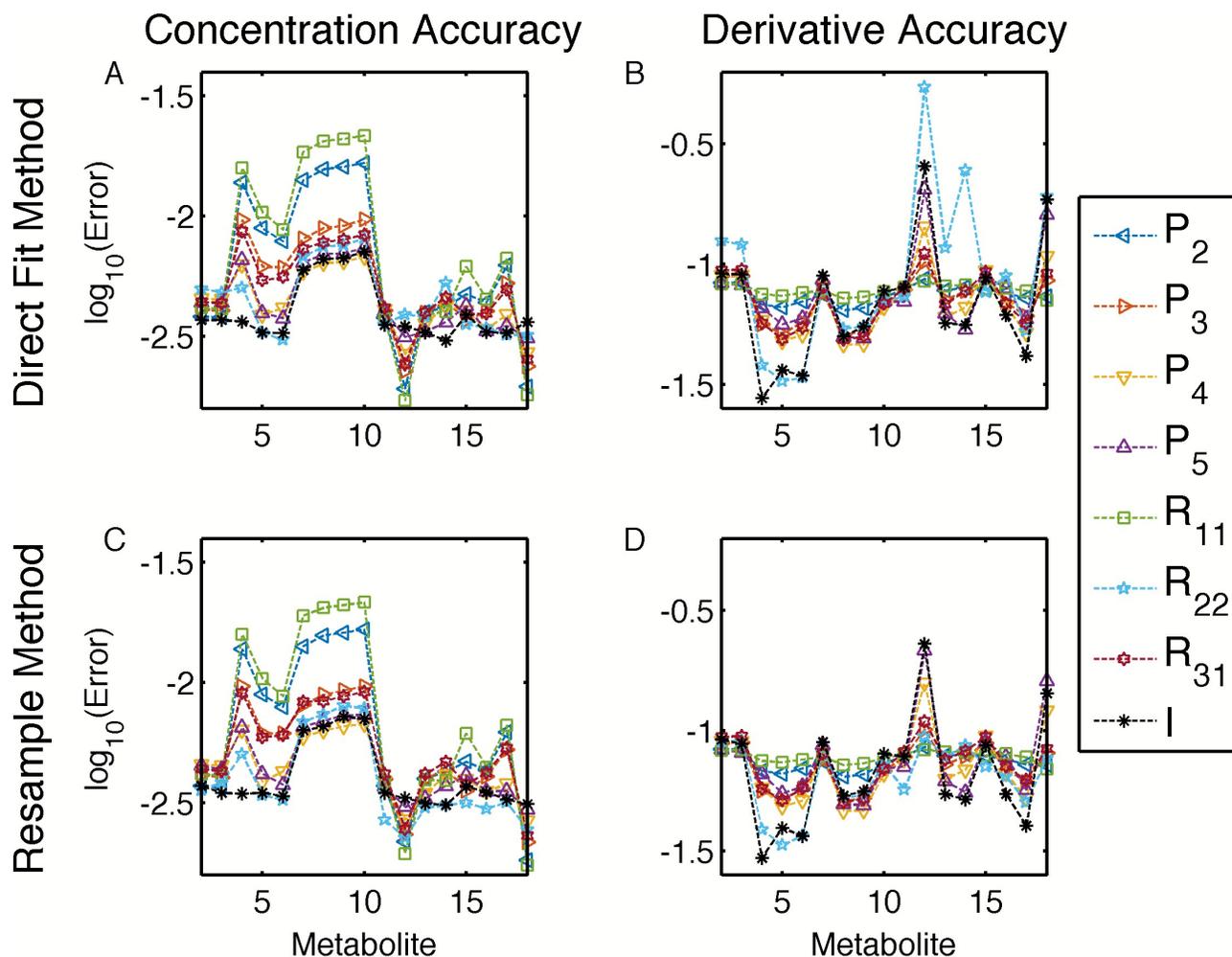
691 Table 2. Average rank of function accuracy using the Direct Fit method on the *E. coli*
692 model.
693

| Average Rank of Metric | P ₂ | P ₃ | P ₄ | P ₅ | R ₁₁ | R ₂₂ | R ₃₁ | I |
|------------------------|----------------|----------------|----------------|----------------|-----------------|-----------------|-----------------|------|
| Concentration Accuracy | 3.68 | 4.13 | 2.50 | 2.94 | 3.94 | 2.33 | 4.83 | 1.74 |
| Derivative Accuracy | 3.18 | 3.45 | 2.48 | 3.08 | 3.58 | 2.61 | 3.77 | 2.18 |

694

695 Fig. 4. Quantitative assessment of function accuracy across metabolites in the *E. coli*
 696 model.

697 The impulse function performs consistently well across most metabolites for both (A)
 698 concentration and (B) derivative accuracy. The resampling method improves the
 699 performance of a number of functions for both (C) concentration and (D) derivative
 700 accuracy. Error metrics are normalized to average metabolite concentrations (see
 701 Methods) for easier visualization and are presented in log-transformed format.
 702



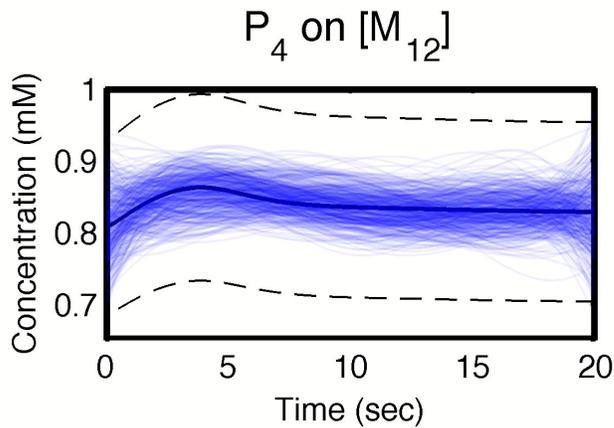
703
 704

705 Fig. 5. Comparison of the Impulse and P_4 on Metabolite 12 (6-Phosphogluconate) over
 706 500 random noisy time courses.

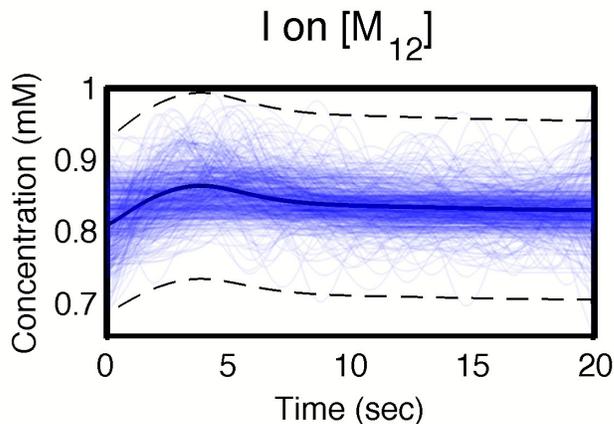
707 A) The P_4 polynomial function intrinsically curves upwards or downwards at the ends of
 708 the interval, which helps match the early slope in the synthetic data. B) The impulse
 709 function exhibits greater variability across different noisy replicates due to the small
 710 dynamic concentration range in the synthetic data relative to the noise introduced. Solid
 711 black lines indicate the synthetic data. Dashed black lines indicate the 15% coefficient of
 712 variation envelope, used to generate the noisy time course data. Blue lines indicate the
 713 concentration trajectory of functional fits to individual noisy time courses. C) As a result,
 714 the P_4 polynomial consistently fits the synthetic data concentration with lower error than
 715 the impulse. Blue dots indicate the error of each function in recapitulating the synthetic
 716 data when fit to a particular noisy time course. The red star indicates the average error
 717 of the blue dots.

718

A

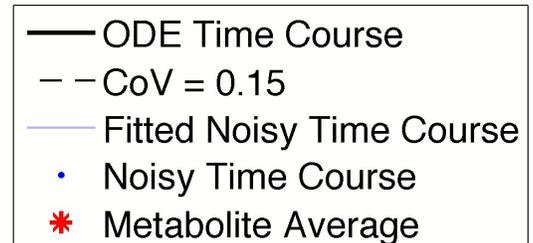
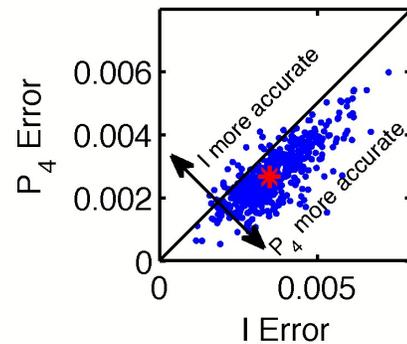


B



C

M_{12} Concentration Accuracy



719

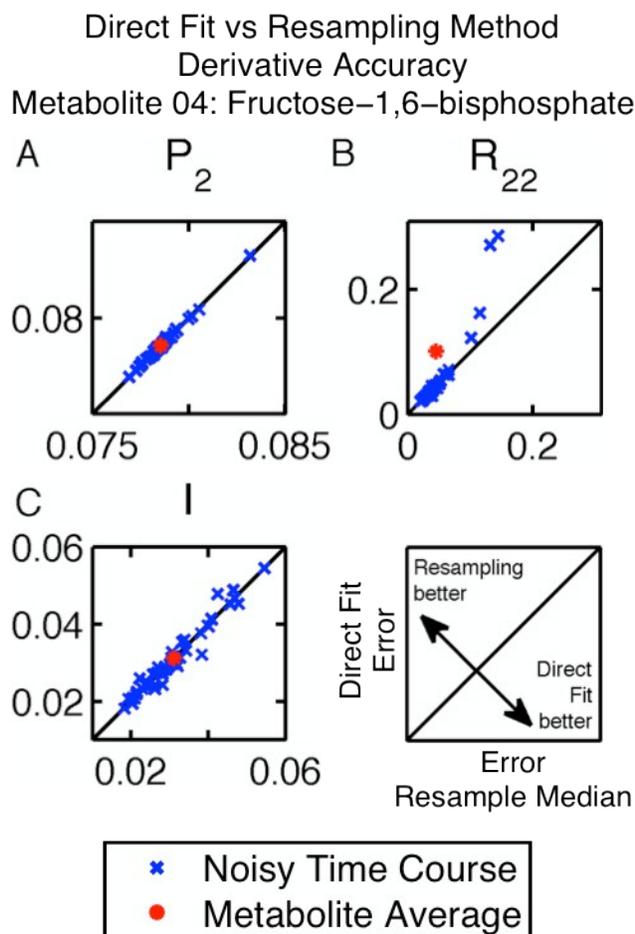
720

721 Table 3. Average rank of function accuracy using the Resampling Method on the *E. coli*
722 model.
723

| Average Rank of Metric | P ₂ | P ₃ | P ₄ | P ₅ | R ₁₁ | R ₂₂ | R ₃₁ | I |
|------------------------|----------------|----------------|----------------|----------------|-----------------|-----------------|-----------------|------|
| Concentration Accuracy | 4.02 | 4.16 | 2.44 | 3.11 | 4.22 | 1.83 | 5.32 | 1.90 |
| Derivative Accuracy | 3.38 | 3.40 | 2.50 | 3.07 | 3.68 | 2.16 | 4.66 | 2.20 |

724

725 Fig. 6. The effect of the Resampling Method on the derivative accuracy of three
 726 representative functions.
 727 The error for fitted concentration profiles was determined for both the Direct Fit and
 728 Resampling Methods and directly compared. A) For polynomial functions the
 729 Resampling Method produces results nearly identical to the Direct Fit method. B) The
 730 R_{22} rational function can produce derivative errors several orders of magnitude greater
 731 using the Direct Fit method (not shown on these axes) than when using the Resampling
 732 Method, making the Resampling Method more accurate on average. C) The impulse
 733 function is generally consistent between the Direct Fit and Resampling Methods, but
 734 does show some variability. Other metabolites exhibit modest benefits from the
 735 Resampling Method relative to the Direct Fit Method.
 736



737
 738

739 Table 4. Average rank of function and method accuracy using the *E.coli* model. Results
 740 from both the Direct Fit (DF) and Resampling (RM) methods are all ranked together to
 741 facilitate direct comparison of their performance.
 742

| Average Rank of Metric | P ₂ | | P ₃ | | P ₄ | | P ₅ | | R ₁₁ | | R ₂₂ | | R ₃₁ | | I | |
|------------------------|----------------|------|----------------|------|----------------|------|----------------|------|-----------------|------|-----------------|------|-----------------|-------|------|------|
| | DF | RM | DF | RM | DF | RM | DF | RM | DF | RM | DF | RM | DF | RM | DF | RM |
| Concentration Accuracy | 6.62 | 6.70 | 7.36 | 7.35 | 3.76 | 3.94 | 5.34 | 5.35 | 7.17 | 6.62 | 3.48 | 2.55 | 8.77 | 10.17 | 2.60 | 2.88 |
| Derivative Accuracy | 5.40 | 5.50 | 6.20 | 6.21 | 3.98 | 4.02 | 5.12 | 5.09 | 6.49 | 5.85 | 3.76 | 3.12 | 6.33 | 8.96 | 3.30 | 3.17 |

743

744 Table 5. Average rank of function accuracy using the *S. cerevisiae* model. Here, the
 745 Direct Fit and Resampling Methods are ranked and averaged separately.
 746

| Average Rank of Metric | Direct Fit Method | | | | | | | | Resampling Method | | | | | | | |
|------------------------|-------------------|----------------|----------------|----------------|-----------------|-----------------|-----------------|------|-------------------|----------------|----------------|----------------|-----------------|-----------------|-----------------|------|
| | P ₂ | P ₃ | P ₄ | P ₅ | R ₁₁ | R ₂₂ | R ₃₁ | I | P ₂ | P ₃ | P ₄ | P ₅ | R ₁₁ | R ₂₂ | R ₃₁ | I |
| Concentration Accuracy | 4.28 | 4.00 | 3.83 | 3.22 | 4.81 | 1.34 | 4.45 | 2.07 | 4.48 | 4.15 | 3.90 | 3.33 | 4.82 | 1.24 | 4.79 | 2.10 |
| Derivative Accuracy | 3.99 | 3.65 | 3.55 | 2.77 | 4.80 | 1.95 | 4.44 | 1.66 | 4.39 | 4.00 | 3.81 | 2.92 | 4.81 | 1.61 | 5.06 | 1.64 |

747

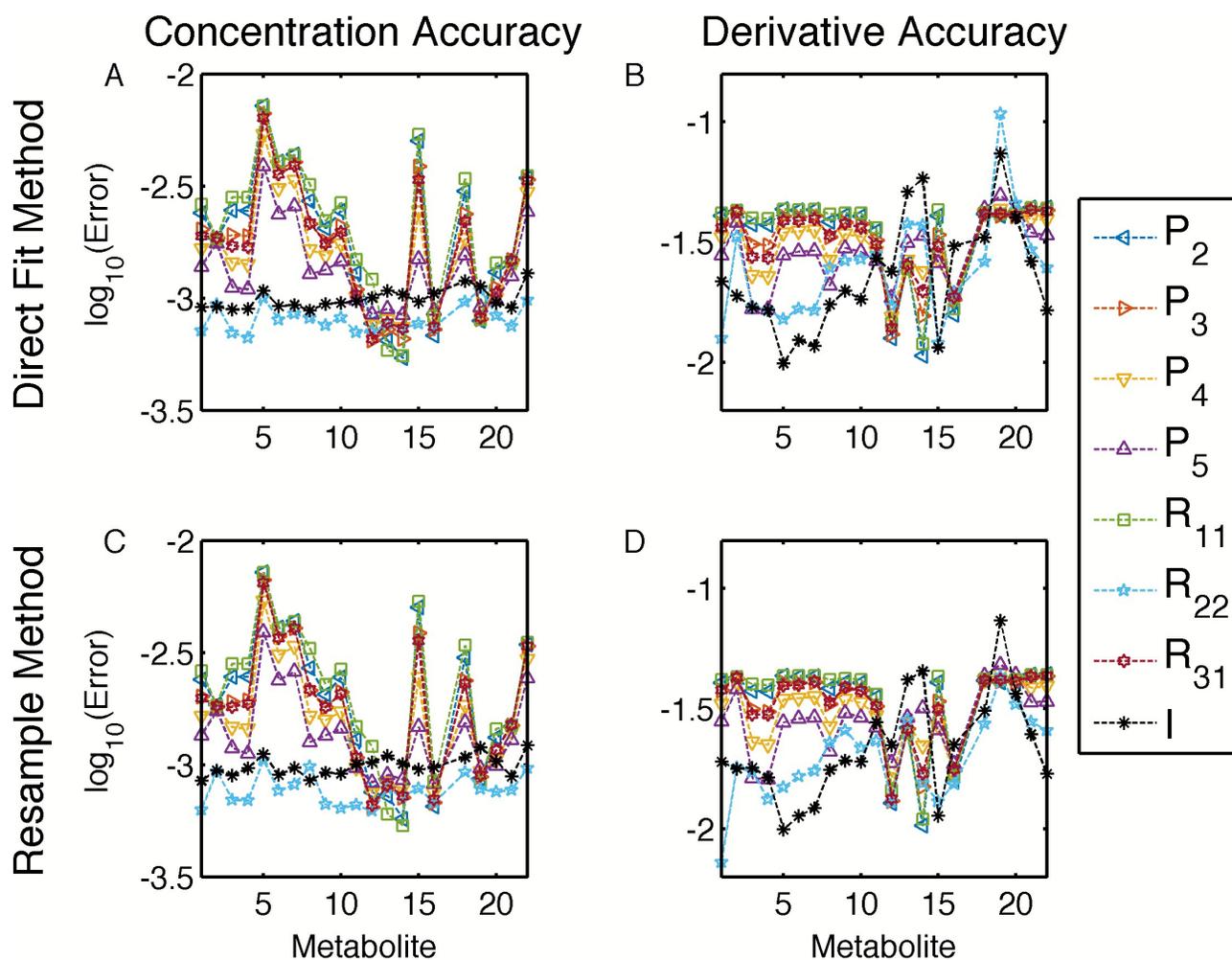
748 Table 6. Average rank of function and method accuracy using the *S. cerevisiae* model.
 749 Results from both the Direct Fit (DF) and Resampling (RM) methods are all ranked
 750 together to facilitate direct comparison of their performance.

751

| Average Rank of Metric | P ₂ | | P ₃ | | P ₄ | | P ₅ | | R ₁₁ | | R ₂₂ | | R ₃₁ | | I | |
|------------------------|----------------|------|----------------|------|----------------|------|----------------|------|-----------------|------|-----------------|------|-----------------|------|------|------|
| | DF | RM | DF | RM | DF | RM | DF | RM | DF | RM | DF | RM | DF | RM | DF | RM |
| Concentration Accuracy | 7.37 | 7.82 | 7.05 | 7.55 | 7.14 | 7.17 | 5.86 | 6.02 | 7.92 | 7.98 | 1.85 | 1.65 | 7.85 | 8.98 | 3.59 | 3.22 |
| Derivative Accuracy | 7.52 | 7.41 | 7.16 | 6.75 | 6.64 | 6.74 | 4.79 | 4.85 | 8.23 | 8.10 | 2.95 | 2.14 | 8.34 | 9.43 | 2.72 | 2.15 |

752

753 Fig. 7. Quantitative assessment of function accuracy across metabolites in the *S.*
 754 *cerevisiae* model.
 755 Results by metric are presented for the Direct Fit Method for (A) concentration accuracy
 756 and (B) derivative accuracy, and for the Resampling Method for (C) concentration
 757 accuracy and (D) derivative accuracy. Error metrics are normalized to average
 758 metabolite concentrations (see Methods) for easier visualization and are presented in
 759 log-transformed format.
 760



761
 762

763 **Notes and References**

764 “Supplementary File 1.pdf” contains Fig. S1-S5, Tables S1-S4, and Supplementary
765 Methods. [PDF, 4.9MB]

766

767 “Supplementary File 2.zip” contains an archive of the code used to generate datasets,
768 fit parameter values, calculate metrics, and plot metrics; and descriptions of file contents
769 and directions on use. [ZIP, 8.2MB]

770

- 771 1. K. Yizhak, E. Gaude, S. Le Dévédec, Y. Y. Waldman, G. Y. Stein, B. van de
772 Water, C. Frezza and E. Ruppín, *eLife*, 2014, **3**, e03641.
- 773 2. A. P. Burgard, P. Pharkya and C. D. Maranas, *Biotechnol Bioeng*, 2003, **84**, 647-
774 657.
- 775 3. K. Nakahigashi, Y. Toya, N. Ishii, T. Soga, M. Hasegawa, H. Watanabe, Y. Takai,
776 M. Honma, H. Mori and M. Tomita, *Mol Sys Biol*, 2009, **5**, n/a-n/a.
- 777 4. M. W. Covert, C. H. Schilling and B. Palsson, *J Theor Biol*, 2001, **213**, 73-88.
- 778 5. M. W. Covert, N. Xiao, T. J. Chen and J. R. Karr, *Bioinformatics*, 2008, **24**, 2044-
779 2050.
- 780 6. J. Min Lee, E. P. Gianchandani, J. A. Eddy and J. A. Papin, *PLoS Comput Biol*,
781 2008, **4**, e1000086.
- 782 7. C. Cotten and J. Reed, *BMC Bioinformatics*, 2013, **14**, 32.
- 783 8. D. McCloskey, B. Ø. Palsson and A. M. Feist, *Mol Sys Biol*, 2013, **9**, 661.
- 784 9. C. S. Henry, L. J. Broadbelt and V. Hatzimanikatis, *Biophysical J*, 2007, **92**,
785 1792-1805.
- 786 10. A. Kümmel, S. Panke and M. Heinemann, *Mol Sys Biol*, 2006, **2**, 2006.0034.
- 787 11. C. Chassagnole, N. Noisommit-Rizzi, J. W. Schmid, K. Mauch and M. Reuss,
788 *Biotechnol Bioeng*, 2002, **79**, 53-73.
- 789 12. R. N. Gutenkunst, J. J. Waterfall, F. P. Casey, K. S. Brown, C. R. Myers and J. P.
790 Sethna, *PLoS Computat Biol*, 2007, **3**, e189.
- 791 13. K. van Eunen, J. Bouwman, P. Daran-Lapujade, J. Postmus, A. B. Canelas, F. I.
792 C. Mensorides, R. Orij, I. Tuzun, J. van den Brink, G. J. Smits, W. M. van Gulik,
793 S. Brul, J. J. Heijnen, J. H. de Winde, M. J. Teixeira de Mattos, C. Kettner, J.
794 Nielsen, H. V. Westerhoff and B. M. Bakker, *FEBS J*, 2010, **277**, 749-760.
- 795 14. R. Mahadevan, J. S. Edwards and F. J. Doyle, *Biophysical J*, 2002, **83**, 1331-
796 1340.
- 797 15. E. O. Voit, G. Goel, I. C. Chou and L. L. Fonseca, in *IET Syst Biol*, Institution of
798 Engineering and Technology, Editon edn., 2009, vol. 3, pp. 513-522.

- 799 16. I.-C. Chou and E. Voit, *BMC Syst Biol*, 2012, **6**, 84.
800 17. G. Goel, I.-C. Chou and E. O. Voit, *Bioinformatics*, 2008, **24**, 2505-2511.
801 18. N. Ishii, Y. Nakayama and M. Tomita, *Theor Biol Med Model*, 2007, **4**, 19.
802 19. K. Yugi, Y. Nakayama, A. Kinoshita and M. Tomita, *Theor Biol Med Model.*,
803 2005, **2**, 42.
804 20. G. Chechik and D. Koller, *J Comput Biol*, 2009, **16**, 279-290.
805 21. J. Sivriver, N. Habib and N. Friedman, *Bioinformatics*, 2011, **27**, i392-i400.
806 22. F. Hynne, S. Danø and P. G. Sørensen, *Biophys Chem*, 2001, **94**, 121-163.
807 23. N. Le Novère, B. Bornstein, A. Broicher, M. Courtot, M. Donizelli, H. Dharuri, L.
808 Li, H. Sauro, M. Schilstra, B. Shapiro, J. L. Snoep and M. Hucka, *Nucleic Acids*
809 *Res*, 2006, **34**, D689-D691.
810 24. S. Hoops, S. Sahle, R. Gauges, C. Lee, J. Pahle, N. Simus, M. Singhal, L. Xu, P.
811 Mendes and U. Kummer, *Bioinformatics*, 2006, **22**, 3067-3074.
812 25. R. G. D. Steel and J. H. Torrie, *Principles and procedures of statistics, with*
813 *special reference to the biological sciences*, McGraw-Hill, New York, 1960.
814 26. T. Evgeniou, T. Poggio, M. Pontil and A. Verri, *Comput. Stat. Data Anal.*, 2002,
815 **38**, 421-432.
816
817