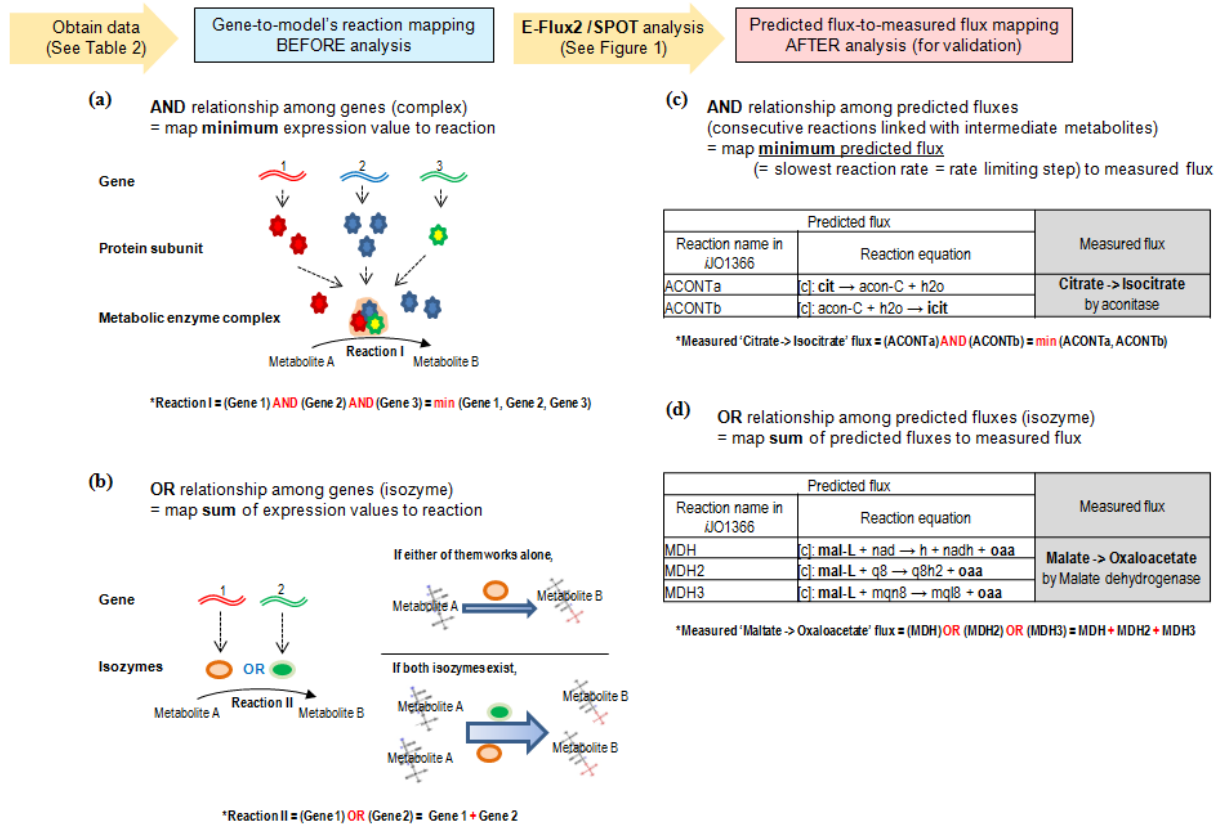


Supplementary figures and methods

E-Flux2 and SPOT: Validated methods for inferring intracellular
metabolic flux distributions from transcriptomic data

Min Kyung Kim, Anatoliy Lane, James J. Kelley, and Desmond S. Lun

1 Supplementary figure 1. Schematic overview of this study.



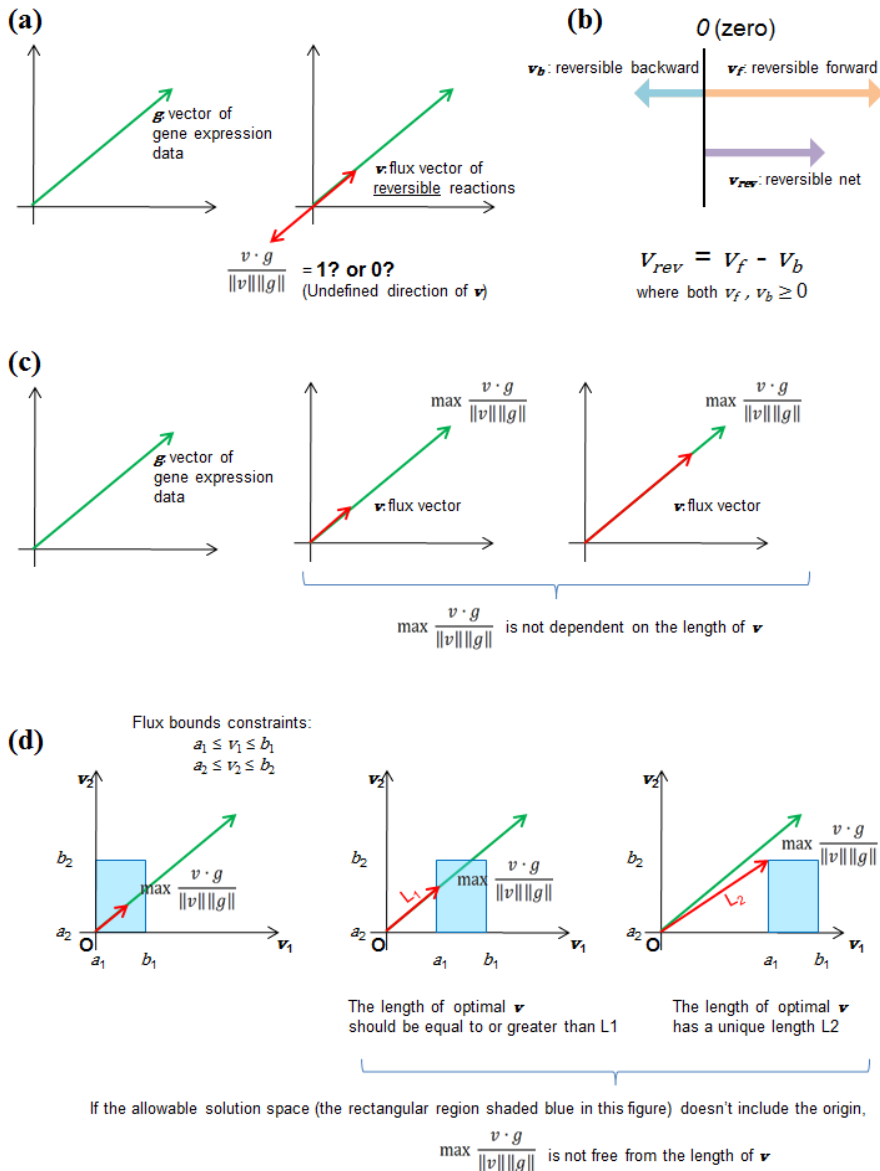
2

3 The process for our research can be classified into five steps, which are: 1) obtaining both transcriptomic and
 4 fluxomic data measured under the same conditions (See Table 2); 2) mapping gene expression data onto
 5 corresponding reactions in the model based on Gene-Protein-Reaction association relationships; 3) creating one of
 6 the two template metabolic models depending on carbon source information; 4) solving an optimization problem
 7 with one of the two algorithms depending on the availability of information on biomass objective (See Figure 1);
 8 and 5) calculating the correlation between predicted and measured fluxes.

9 Supplementary figure 1a and Supplementary figure 1b illustrate how transcriptomic data were mapped onto
 10 corresponding reactions in the model based on Gene-Protein-Reaction association relationships. Supplementary
 11 figure 1c and Supplementary figure 1d shows how predicted fluxes were matched with a corresponding measured
 12 flux to calculate correlation between them. (a) In the case where an enzyme complex mediates a certain metabolic
 13 reaction (AND relationship), we mapped the minimum value of the expression level of the associated genes
 14 encoding its subunits to the corresponding reaction since the least-expressed components is likely to determine the
 15 final concentration of the complete enzyme complex. (b) If a reaction is catalyzed by isozymes (OR relationship),

1 we took the sum of the expression values of the associated genes for mapping since the total capacity of the reaction
2 is given by the sum of the capacities of its isozymes. (c) If a certain measured reaction corresponds to the set of
3 consecutive reactions in the model that share intermediate metabolites (AND relationship), the slowest reaction rate
4 also known as the rate-determining step (the minimum flux) among those predicted fluxes was used to match with
5 the corresponding measured flux. (d) In the case where a measured flux corresponds to multiple reactions in the
6 model that mediate an identical chemical conversion independently with each other (OR relationship), the sum of
7 those predicted fluxes was used to match with the corresponding measured flux.

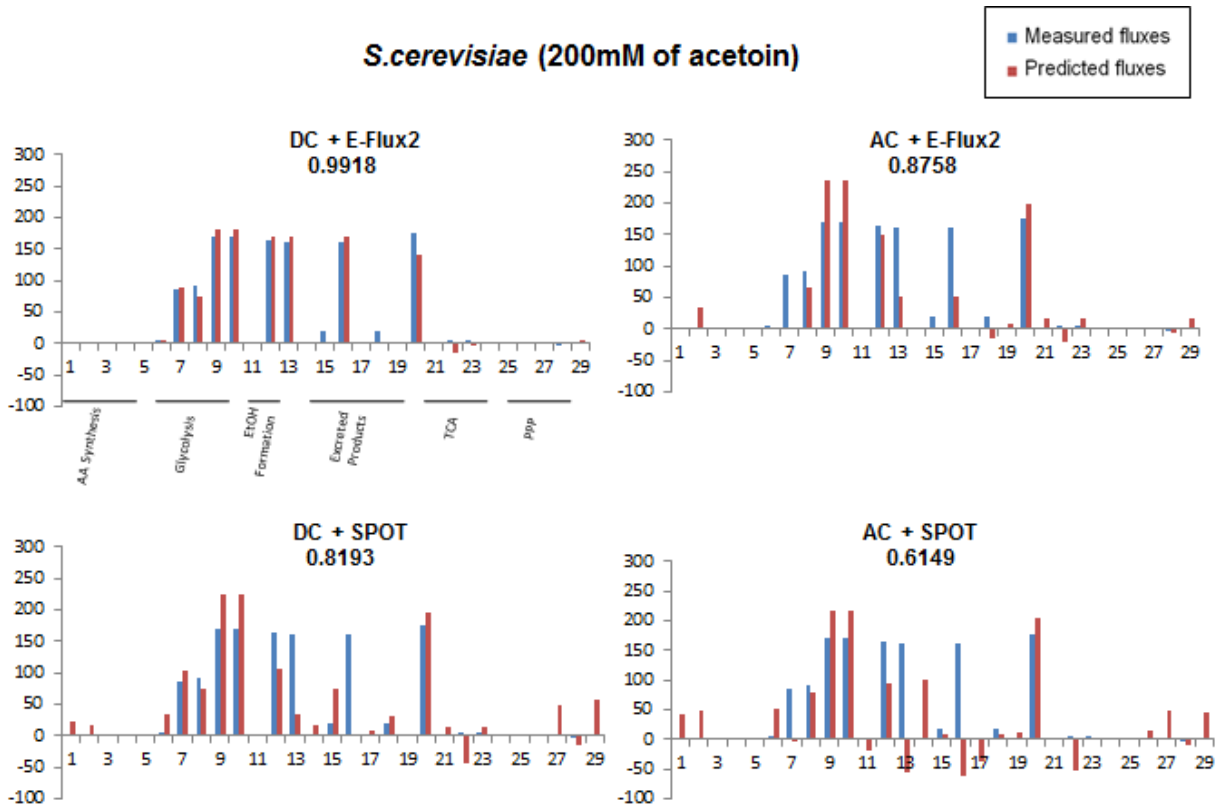
1 **Supplementary figure 2. Rationale for the SPOT method.**



- 2 (a) We noticed that the dot product between a flux vector (v , denoted as red arrows in the figure) and its
- 3 corresponding gene expression data (g , denoted as green arrows in the figure) in the numerator of the objective
- 4 function cannot be calculated for the set of reversible reactions since the directions of reversible reactions are
- 5 undefined whereas gene expression data values are always positive. (b) So we decomposed every reversible reaction
- 6 in the model into two positive irreversible reactions, the forward reaction, v^f , and the backward reaction, v^b , where
- 7 $v^{rev} = v^f - v^b$, and $v^f, v^b \geq 0$ (c) Usually, the maximum Pearson product-moment correlation is not dependent
- 8 on the length of v but dependent on the angle between v and g . (d) However, if any of the flux bounds does not

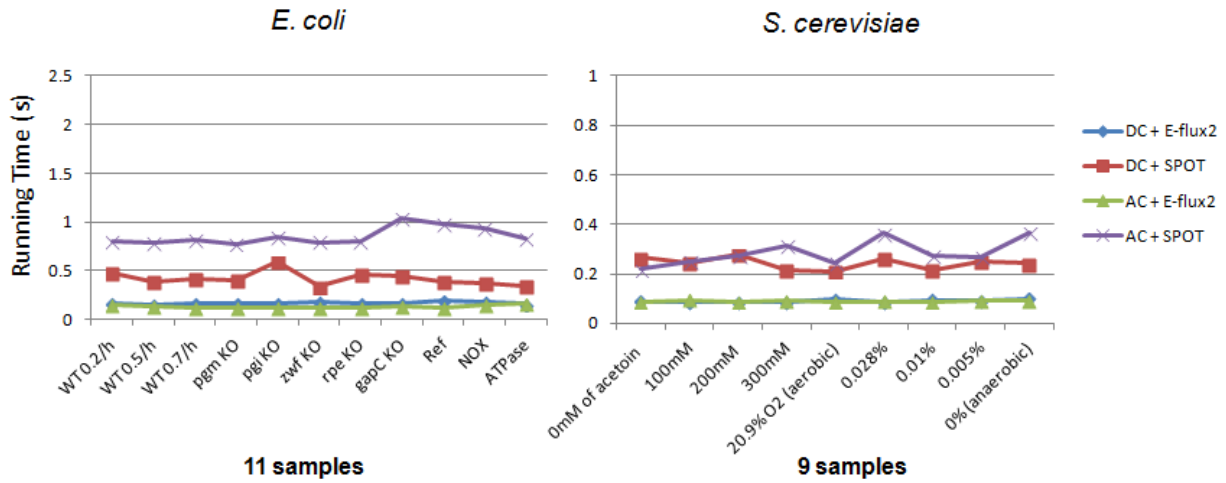
- 1 include zero, the origin in the graph, the maximum correlation is no longer independent of the length of the flux
- 2 vector, v . Thus, it is a prerequisite for using SPOT method to make sure that the allowable solution space includes
- 3 the origin. The light blue-colored rectangular space shows the allowable solution space that is determined by flux
- 4 bounds (i.e. a_1 , b_1 and a_2 , b_2 in the figure) of all reactions (i.e. v_1 and v_2).

- 1 **Supplementary figure 3. Comparison of the predicted fluxes with the measured fluxes of *S.***
- 2 ***cerevisiae* data (200mM of acetoin-treated sample of Celton *et al.*).**



- 3
- 4 The x-axis represents metabolic reactions used to calculate correlation between the measured (blue bars in the
- 5 figure) and the predicted fluxes (red bars in the figure), and the y-axis indicates flux value. The scale and the units
- 6 on the y-axis are based on those of the measured flux.

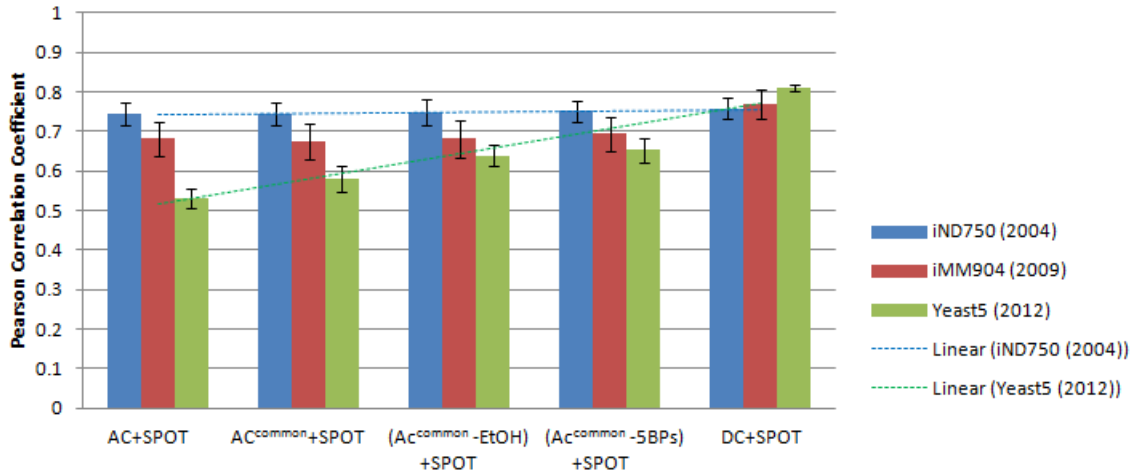
1 **Supplementary figure 4. Average running time of our algorithm.**



2

3 We measured the running time of our algorithm implemented using MATLAB (The Mathworks, Inc., Natick, Mass.)
4 and Gurobi Optimizer (Gurobi Optimization, Inc., Houston, Texas) for all 80 samples (4 simulation methods and 20
5 samples per simulation method) using the built-in MATLAB function, profile. Regardless of which simulation
6 method is used, our method completes within one second for both *E. coli* and *S. cerevisiae*. Computations were
7 carried out on the Window 8 OS platform using a personal computer with an Intel Core i5 3.10 GHz processor with
8 8 GB of RAM.

1 **Supplementary figure 5. Exploration of the way to improve a poor performance of Yeast 5**
 2 **in AC+SPOT.**



3 As described in Methods, we built the AC (All possible Carbon sources) model, which has lower bounds of
 4 negative infinity for all exchange reactions of possible carbon sources (i.e. external metabolites containing carbon)
 5 to simulate the unknown carbon source situation. As listed in Additional file 5, a total of 108, 154, and 158 potential
 6 carbon sources were allowed to be taken up by the cell for *iND750*, *iMM904*, and Yeast 5 models of *S. cerevisiae*,
 7 respectively. To test whether Yeast 5 performs worse than the older models because it has more carbon source
 8 uptake reactions (leading to more incorrect carbon sources to confound the prediction method), we performed SPOT
 9 again after blocking the uptake of model-specific carbon sources, leaving 106 exchange reactions that are common
 10 across all three models (see Additional file 7 for details). In other words, we updated the three yeast AC models so
 11 that they all have the same set of 106 possible carbon source exchange reactions, which we call the AC^{common} (All
 12 possible *common* Carbon sources) model in the figure to distinguish it from the original AC model.

13 As shown in the AC^{common}+SPOT case of the figure above, although the average correlation of Yeast 5 was
 14 improved from 0.5313 to 0.5791 after reducing the number of carbon source uptake reactions to 106, the older yeast
 15 models still outperform Yeast 5, which suggests that the lower correlation achieved by SPOT on the Yeast 5 AC
 16 model is not simply due to having a greater number of possible carbon sources.

17 Interestingly, unlike the older models, the performance of Yeast 5 in AC+SPOT was much improved by limiting the
 18 uptake of well-known by-products of yeast such as ethanol and glycerol [1]. As shown in the (AC^{common}-

1 EtOH)+SPOT and (AC^{common}-5BPs)+SPOT case, the correlation of Yeast 5 was further increased from 0.5791 to
2 0.6397 and 0.6534, when the uptake of ethanol and of well-known five metabolic by-products of yeast (ethanol,
3 carbon dioxide, succinate, glycerol, and acetate) was blocked respectively. The different sensitivity of the
4 performance of each model to changes in the number of carbon sources (see the trend lines in the figure) may
5 indicate a different degree of interconnectivity among intracellular reactions and exchange reactions inherent to the
6 model. Considering that only a small number of preferred carbon sources are consumed by most microorganisms,
7 which is a well-known phenomenon called Carbon Catabolite Repression (CCR) [2, 3], reducing the number of
8 possible carbon source uptake reactions by analyzing growth media composition is a sensible way to improve the
9 predictive accuracy of AC+SPOT.

1 Supplementary Methods

2 A mathematical justification for dropping the $\|v\|$ term in the objective function of SPOT

3 Removing the $\|v\|$ term in the objective function of SPOT can be justified if a solution of problem 2 (see below, on
4 the right side) also optimizes problem 1 (on the left side):

$$\begin{array}{ccc} \text{Problem 1} & & \text{Problem 2} \\ \max \frac{v \cdot g}{\|v\|} & \rightarrow & \max v \cdot g \\ \text{subject to } \begin{cases} Av = 0 \\ 0 \leq v \end{cases} & & \text{subject to } \begin{cases} Av = 0 \\ 0 \leq v \\ \|v\| \leq 1 \end{cases} \end{array}$$

5 where v and g has dimension n , and A is a stoichiometric matrix.

6 **Lemma.** If v^* optimizes problem 2, then v^* optimizes problem 1.

7 This lemma can be proved by contradiction.

8 Let us assume that this statement is false. Then its negation, i.e. if v^* optimizes (i.e. maximizes) problem 2, then v^*
9 does not optimize problem 1, should be true. In other words, we assume that another vector v exists such that

$$\frac{v \cdot g}{\|v\|} > \frac{v^* \cdot g}{\|v^*\|}$$

10 where $v \neq v^*$.

11 Let $v' = \frac{v}{\|v\|}$, then $\|v'\| = \frac{1}{\|v\|} \times \|v\| = 1$.

12 Substituting v' into the term on the left side of the inequality above gives $v' \cdot g$.

1 In addition, since v^* optimizes problem 2 by assumption, $\|v^*\| \leq 1$. Note that no optimal solution to problem 2 has
2 $\|v\| < 1$ because g is a non-negative vector and we assume that $g \neq 0$ (for otherwise, the problem is trivial). Thus,
3 if $\|v\| < 1$, the objective of problem 2 can always be increased by scaling up v until $\|v\| = 1$. Hence $\|v^*\| = 1$.

4 The right side of the inequality, thus, is equal to $v^* \cdot g$. Therefore,

$$v' \cdot g > v^* \cdot g$$

5 where $\|v'\| = 1$.

6 This contradicts our assumption that v^* is a vector that optimizes problem 2. Since the negation is impossible (false),
7 the original statement is true. Note that the lemma is still valid if the number used to limit $\|v\|$ in problem 2 is any
8 constant $c > 0$.

9

1 **A mathematical proof of the uniqueness of SPOT solutions**

2 SPOT (equation (8) in the main manuscript) is defined as the following optimization problem:

$$\begin{aligned} & \max f'v \\ & \text{subject to } \begin{cases} Av = 0 \\ 0 \leq v \\ \|v\| \leq 1 \end{cases} \end{aligned}$$

3 where v has dimension n , $f \geq 0$ and A is a stoichiometric matrix.

4 **Lemma.** If the optimal value of the SPOT problem is strictly positive, then its solution is unique.

5 **Proof.** *Step 1.* First we prove that $\|v\| = 1$. If $\|v\| < 1$ then for ϵ sufficiently small $\omega = v + \epsilon \frac{v}{\|v\|}$ satisfies
6 $\|v\| < 1$ (and the other constraints) and

$$f'\omega = f'v \left(1 + \frac{\epsilon}{\|v\|}\right) > f'v$$

7 because $f'v > 0$, contradicting the maximality of $f'v$.

8 *Step 2.* Because of the maximality of $f'v$, the affine plane $\{\omega: f'\omega = f'v\}$ is a supporting plane for the convex set
9 $C = \{\omega: \omega \in \ker(A), \omega \geq 0, \|\omega\| \leq 1\}$. Moreover, by Step 1, such a plane cannot be orthogonal to any linear
10 space containing the vector v . We conclude that the affine plane intersects C only at v , thus v is the only point of the
11 maximum for the linear functional $\omega \rightarrow f'\omega$.

12 **Remark.** Notice that the condition of strict positivity of the maximum if necessary. Indeed if $\ker(A)$ is contained on
13 a coordinate plane, say the first: $\ker(A) \subset \{\omega: \omega_1 = 0\}$, then the vector $f = e_1$ satisfies $f'\omega = 0$ for every vector
14 in $\ker(A)$ and there is no uniqueness.

15
16

1 **Calculation of the possible range of correlation between the measured fluxes and the**
 2 **predicted fluxes**

3 For standard FBA and E-Flux, the methods that do not give a unique metabolic flux distribution, the possible range
 4 of correlation between the measured fluxes and the predicted fluxes were calculated. Given information on the
 5 measured fluxes, the minimum and the maximum correlations of standard FBA can be calculated using the
 6 following two steps of optimization:

<p>7 <u>Step 1. standard FBA</u></p> <p>8 $z^* = \max f'v$</p> <p>9 subject to $\begin{cases} Sv = 0 \\ a_j \leq v_j \leq b_j \end{cases}$</p>	<p>→</p>	<p>7 <u>Step 2. calculation of the possible range of correlation</u></p> <p>8 $\min/\max \frac{v_p \cdot v_m}{\ v_p\ \ v_m\ }$</p> <p>9 subject to $\begin{cases} Sv = 0 \\ a_j \leq v_j \leq b_j \\ f'v = z^* \end{cases}$</p>
--	----------	--

10 where v is a flux vector representing the reaction rates of the n reactions in the network, f is a coefficient vector
 11 defining the organism's objective function, S is the stoichiometric matrix, and a_j and b_j are the minimum and
 12 maximum reaction rates through reaction j . The vectors v_p and v_m shown in the objective function of Step 2 are the
 13 predicted and measured vectors of intracellular fluxes, respectively, and $\|\cdot\|$ denotes the l^2 norm.

14 The average correlation of standard FBA in Table 3 was calculated using solutions obtained in Step 1 under our
 15 computational settings (see Methods in the main manuscript for the detailed settings). Step 2 allows us to find a
 16 metabolic flux distribution which achieves theoretically maximal or minimal correlation with the measured fluxes
 17 while maintaining the optimal biomass flux, denoted as z^* here. The nonlinear optimization problem in Step 2 was
 18 solved using the sequential quadratic programming (SQP) algorithm provided by the MATLAB function `fmincon`.

19 Importantly, the maximum possible correlation can be calculated only when we already have the known measured
 20 flux datasets. There is no way to force each method to produce a metabolic flux distribution which achieves the best
 21 correlation with the measured fluxes. Our methods were developed during the process of finding that way and of
 22 rigorously testing various strategies.

23 In the same way, the lower and upper bound of correlations of E-Flux can be calculated as follows:

1 **References**

- 2 1. Nevoigt E: **Progress in metabolic engineering of *Saccharomyces cerevisiae***. *Microbiol Mol Biol Rev* 2008,
3 **72**:379–412.
- 4 2. Görke B, Stülke J: **Carbon catabolite repression in bacteria: many ways to make the most out of nutrients**.
5 *Nat Rev Microbiol* 2008, **6**:613–624.
- 6 3. Vinuselvi P, Kim MK, Lee SK, Ghim CM: **Rewiring carbon catabolite repression for microbial cell factory**.
7 *BMB Rep* 2012:59–70.

8