

Supplementary Material: A principal components method constrained by elementary flux modes: Analysis of flux data sets

Moritz von Stosch (1), Cristiana Rodrigues de Azevedo (1), Mauro Luis (1), Sebastiao Feyeo de Azevedo (2), Rui Oliveira* (1)

1 – REQUIMTE/DQ, Faculty of Science and Technology, University Nova de Lisboa, Campus de Caparica, 2829-516 Caparica, Portugal

2 – DEQ, Faculty of Engineering, University do Porto, Rua Dr. Roberto Frias s/n, 4200-465 Porto, Portugal

*Corresponding author

A: The iterative calculation of the p's and its relation to the Gram Schmidt orthonormalization

For two normalized EMs, $e_{j,n}$ and $e_{k,n}$, the procedure described in the article for calculating the respective weightings from the flux matrix V_{mes} yields:

$$p_{j,n} = e_{j,n}^T \cdot V_{mes} , \quad (A 1)$$

such that:

$$V_{j,n} = e_{j,n} \cdot p_{j,n} \quad (A 2)$$

and

$$V_{\Delta mes} = V_{mes} - V_{j,n} \quad (A 3)$$

wherefore:

$$p_{k,n} = e_{k,n}^T \cdot V_{\Delta mes} . \quad (A 4)$$

Inserting the introduced relations in this equation yields:

$$p_{k,n} = e_{k,n}^T \cdot (V_{mes} - V_{j,n}) = e_{k,n}^T \cdot (V_{mes} - e_{j,n} \cdot p_{j,n}) \quad (A 5)$$

and further:

$$p_{k,n} = e_{k,n}^T \cdot (V_{mes} - e_{j,n} \cdot (e_{j,n}^T \cdot V_{mes})) . \quad (A 6)$$

Reformulating:

$$p_{k,n} = (e_{k,n}^T - e_{k,n}^T \cdot e_{j,n} \cdot e_{j,n}^T) \cdot V_{mes} \quad (A 7)$$

Introducing the relation

$$p_{k,n} = e_{k,c,n}^T \cdot V_{mes} \quad (A 8)$$

Where $e_{k,c,n}$ could be interpreted as the interference corrected $e_{k,n}$, yields:

$$e_{k,c,n}^T \cdot V_{mes} = (e_{k,n}^T - e_{k,n}^T \cdot e_{j,n} \cdot e_{j,n}^T) \cdot V_{mes} \quad (A 9)$$

and rearrangement

$$e_{k,c,n}^T = e_{k,n}^T - e_{k,n}^T \cdot e_{j,n} \cdot e_{j,n}^T , \quad (A 10)$$

which is equal to what the Gram Schmidt orthonormalization for the normalized EMs would yield in case that the weight values of $p_{j,n}$ and $p_{k,n}$ are greater than zero. The same can be shown for more factors. The orthonormalization of the elementary modes in the proposed method is implicitly achieved by the subtraction of the captured variance from the data, Eq. (A3), and subsequent analysis of the variance remaining in the data, Eq. (A4,) as in principle component analysis. This is an important feature, as this means that:

- i) the calculated contributions can simply be summed since they are independent;
- ii) the iteratively calculated variances for the EMs can be summed (see also Eq. (B6)),

wherefore:

$$(\sum_l p_{l,n})^2 = \sum_l (p_{l,n})^2, \quad (\text{A } 11)$$

which allows us to simplify equations (7) resulting in equation (8), see paper;

- iii) allows us to assess how much information two EMs share, see following.

B: Analysis of the interference between two successive EMs

The flux patterns which are captured by the combination of the two normalized EMs, $e_{j,n}$ and $e_{k,n}$, is described by the sum:

$$V_{est} = V_{j,n} + V_{k,n} \quad (\text{B } 1)$$

When inserting the relations presented before

$$V_{est} = V_{j,n} + V_{k,n} = e_{j,n} \cdot p_{j,n} + e_{k,n} \cdot p_{k,n} \quad , \quad (\text{B } 2)$$

$$V_{est} = e_{j,n} \cdot (e_{j,n}^T \cdot V_{mes}) + e_{k,n} \cdot (e_{k,n}^T \cdot V_{\Delta mes}) \quad , \quad (\text{B } 3)$$

$$V_{est} = e_{j,n} \cdot (e_{j,n}^T \cdot V_{mes}) + e_{k,n} \cdot (e_{k,n}^T \cdot (V_{mes} - V_{j,n})) \quad , \quad (B 4)$$

$$V_{est} = e_{j,n} \cdot (e_{j,n}^T \cdot V_{mes}) + e_{k,n} \cdot \left(e_{k,n}^T \cdot \left(V_{mes} - e_{j,n} \cdot (e_{j,n}^T \cdot V_{mes}) \right) \right), \quad (B 5)$$

and rearranging

$$V_{est} = e_{j,n} \cdot (e_{j,n}^T \cdot V_{mes}) + e_{k,n} \cdot (e_{k,n}^T \cdot V_{mes}) - \left(e_{k,n}^T \cdot e_{j,n} \cdot (e_{j,n}^T \cdot V_{mes}) \right). \quad (B 6)$$

If the inner vector product of $e_{k,n}^T \cdot e_{j,n}$ would be zero then the contributions would be independent, i.e.:

$$e_{k,n}^T \cdot e_{j,n} = |e_{k,n}| \cdot |e_{j,n}| \cdot \cos(\angle e_{k,n}, e_{j,n}) = 0 \quad .$$

Wherefore the last term in equation (B 6) would vanish. Otherwise the subtracted term is the information, which both EMs share and, if it would not be subtracted, would contribute two times to the summation of flux pattern. Note that it is exactly the contribution of this term, which is eliminated by the Gram-Schmidt approach through orthonormalization and in the proposed algorithm is implicitly accounted for by the subtraction of the captured variance from the data, Eq. (A3), and subsequent analysis of the variance remaining in the data, Eq. (A4). The decompositions of V_{mes} using e_j and then e_k are thus independent!

C: Pichia pastoris Simulation Case

This simulation case study is based on the metabolic network of *Pichia pastoris* proposed by Tortajada et al [1] and the provided 98 EMs are herein used for flux data generation. Twelve different experimental conditions were simulated, see Table 1. Considering all of the different conditions the EMs were filtered employing the following rational. If a compound is not present and cannot be produced then its contribution in the EM must be zero. If a compound is present

and is consumed then it cannot be secreted at the same time. The reduced set of EMs possibly active considering these scenarios comprised 76 EMs. The rank of the reduced set of EMs was analyzed (rank=20) and from the reduced set a subset of EMs was chosen that spans the space of (provides a basis to) the reduced set (20 EMs). For each scenario only those EMs (out of the subset of the reduced set) that can be active were assumed to contribute randomly to the flux pattern, but obeying to the minimum and maximum values of the fluxes given in Tortajada et al[1]. In total only 16 EM were active, namely EMs=[1, 3, 7, 12, 13, 14, 16, 19, 20, 22, 23, 24, 28, 32, 33, 37].

Table 1: Biomass growth in different simulated experimental scenarios

| Exp | Experimental Scenario | Condition for fluxes |
|-----|---|--|
| 1 | Growth on Glucose | $v_{Gly}=0, v_{Pyr}=0, v_{Met}=0$ |
| 2 | Growth on Methanol | $v_{Glc}=0, v_{Gly}=0, v_{Pyr}=0$ |
| 3 | Growth on Glycerol | $v_{Glc}=0, v_{Pyr}=0, v_{Met}=0$ |
| 4 | Growth on Glucose and Pyruvate | $v_{Gly}=0, v_{Met}=0$ |
| 5 | Growth on Methanol and Pyruvate | $v_{Glc}=0, v_{Gly}=0$ |
| 6 | Growth on Glycerol and Pyruvate | $v_{Glc}=0, v_{Met}=0$ |
| 7 | Growth on Glucose without Ethanol formation | $v_{Eth}=0, v_{Gly}=0, v_{Pyr}=0, v_{Met}=0$ |
| 8 | Growth on Methanol without Ethanol formation | $v_{Eth}=0, v_{Glc}=0, v_{Gly}=0, v_{Pyr}=0$ |
| 9 | Growth on Glycerol without Ethanol formation | $v_{Eth}=0, v_{Glc}=0, v_{Pyr}=0, v_{Met}=0$ |
| 10 | Growth on Glucose and Pyruvate without Ethanol formation | $v_{Eth}=0, v_{Gly}=0, v_{Met}=0$ |
| 11 | Growth on Methanol and Pyruvate without Ethanol formation | $v_{Eth}=0, v_{Glc}=0, v_{Gly}=0$ |
| 12 | Growth on Glycerol and Pyruvate without Ethanol formation | $v_{Eth}=0, v_{Glc}=0, v_{Met}=0$ |

D: Results Simulation Case – Impact of Noise on the EM identification

Table 2: Selected EMs and the respective captured variance (ϑ) values for one to 10 number of n_{Fac} s obtained for the simulated data with 2% Gaussian noise. The set of truly active EMs for data generation was EMs=[1, 3, 7, 12, 13, 14, 16, 19, 20, 22, 23, 24, 28, 32, 33, 37]. n_{lv}^{**} : Number of latent variables for PCA.

| n_{Fac} | EM/ ϑ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|-----------|-----------------|-------|---|---|---|---|---|---|---|---|----|
| 1 | EM | 70 | | | | | | | | | |
| 1 | ϑ | 29.39 | | | | | | | | | |
| 2 | EMs | 70 | 7 | | | | | | | | |

| | | | | | | | | | | | |
|-----|---------------|-------|-------|-------|-------|-------|-------|--------|--------|--------|-------|
| 2 | ϑ | 29.39 | 53.80 | | | | | | | | |
| 3 | EMs | 70 | 7 | 40 | | | | | | | |
| 3 | ϑ | 29.39 | 53.80 | 71.13 | | | | | | | |
| 4 | EMs | 7 | 69 | 13 | 40 | | | | | | |
| 4 | ϑ | 24.53 | 46.87 | 65.83 | 83.53 | | | | | | |
| 5 | EMs | 7 | 71 | 13 | 33 | 37 | | | | | |
| 5 | ϑ | 24.53 | 45.98 | 64.40 | 82.47 | 93.95 | | | | | |
| 6 | EMs | 7 | 33 | 13 | 3 | 37 | 23 | | | | |
| 6 | ϑ | 24.53 | 44.08 | 62.34 | 78.92 | 91.37 | 97.07 | | | | |
| 7 | EMs | 7 | 33 | 13 | 3 | 37 | 23 | 12 | | | |
| 7 | ϑ | 24.53 | 44.08 | 62.34 | 78.92 | 91.37 | 97.07 | 97.21 | | | |
| 8 | EMs | 7 | 33 | 13 | 3 | 37 | 23 | 12 | 19 | | |
| 8 | ϑ | 24.53 | 44.08 | 62.34 | 78.92 | 91.37 | 97.07 | 97.21 | 97.25 | | |
| 9 | EMs | 7 | 33 | 13 | 3 | 37 | 23 | 12 | 19 | 16 | |
| 9 | ϑ | 24.53 | 44.08 | 62.34 | 78.92 | 91.37 | 97.07 | 97.21 | 97.25 | 97.26 | |
| 10 | EMs | 7 | 33 | 13 | 3 | 37 | 23 | 12 | 19 | 16 | 17 |
| 10 | ϑ | 24.53 | 44.08 | 62.34 | 78.92 | 91.37 | 97.07 | 97.21 | 97.25 | 97.26 | 97.26 |
| PCA | n_{iv}^{**} | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | |
| PCA | ϑ | 50.13 | 82.14 | 91.80 | 96.98 | 99.27 | 99.96 | 100.00 | 100.00 | 100.00 | |

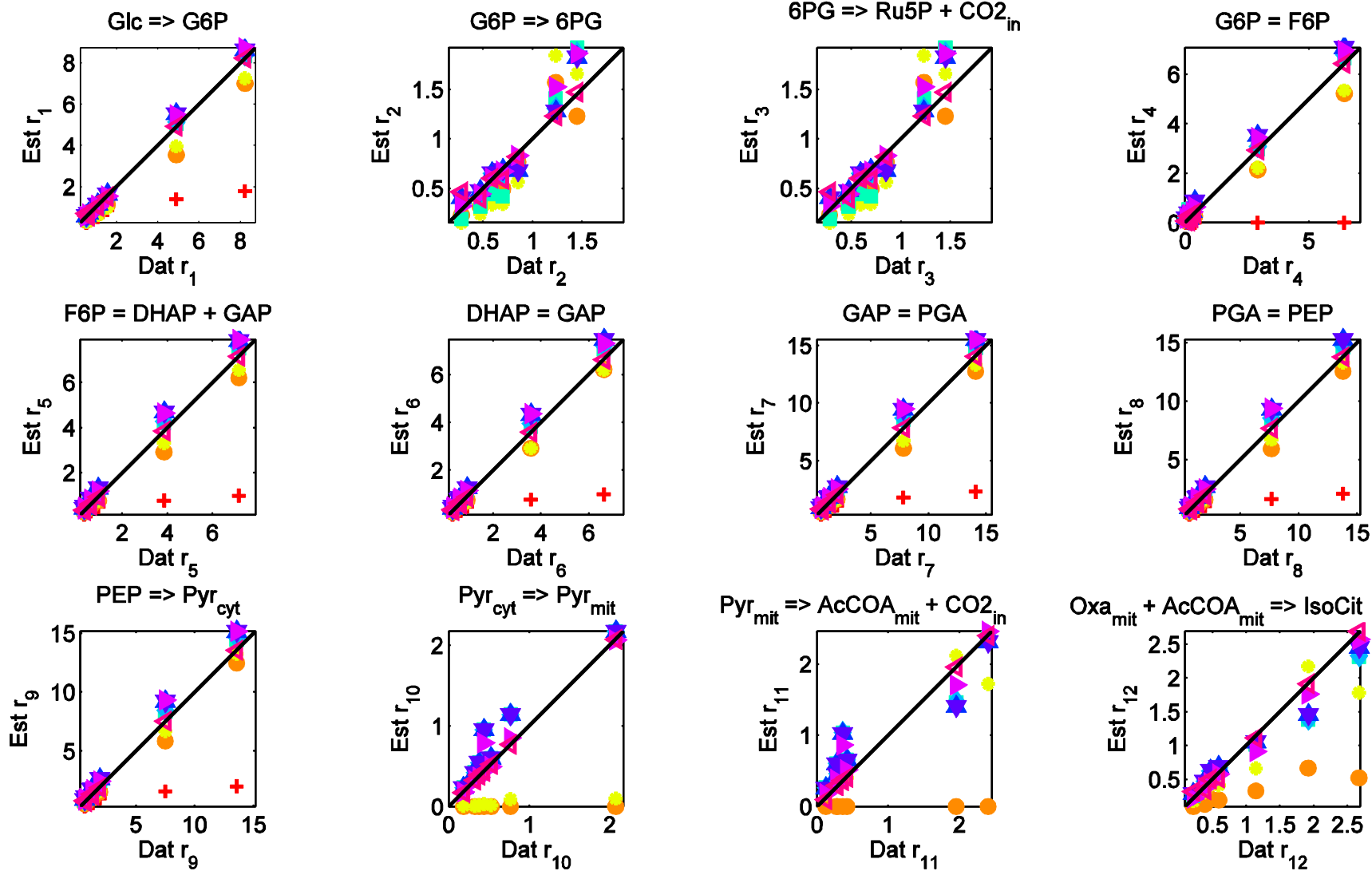
E: Saccharomyces cerevisiae experimental case study – Reaction Network

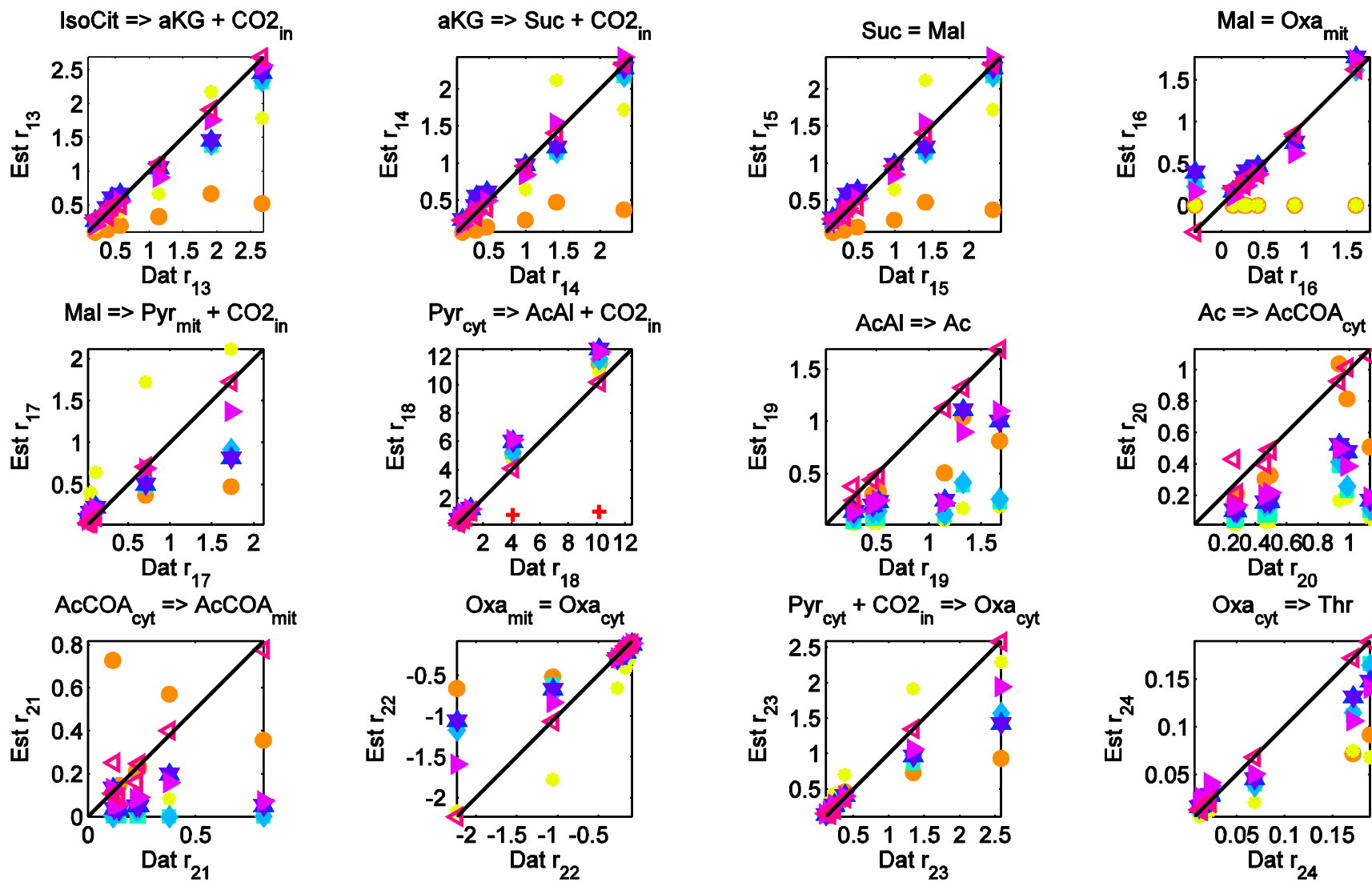
The metabolic network proposed by Hayakawa et al [2] was adopted in this study. To improve the performance of the algorithm the biomass synthesis equation of Gianchandani et al [3] was used instead of the variable biomass synthesis rates used in [2]. Consistency of the metabolic model and data was verified using Metabolic Flux Analysis. The system was observed to be redundant and consistent. The data and the network reactions can be found in an excel-file in the Supplementary Material.

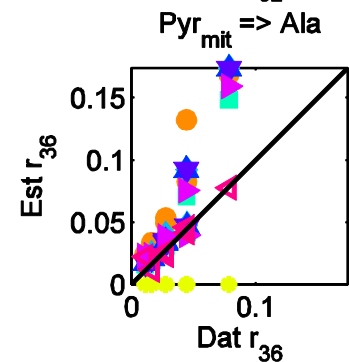
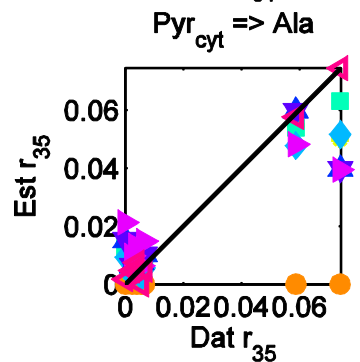
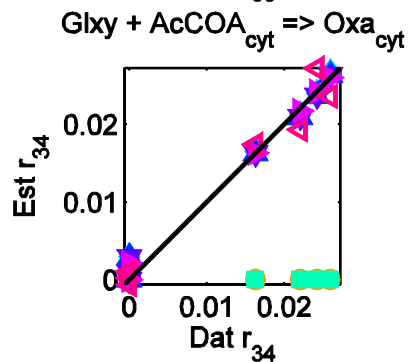
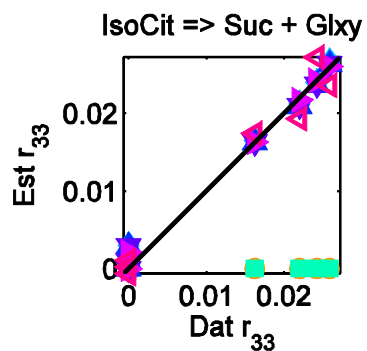
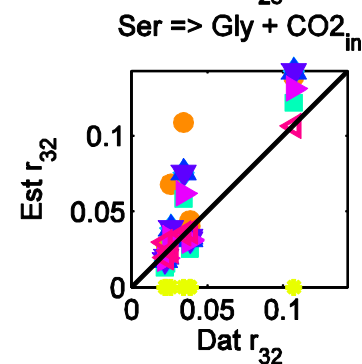
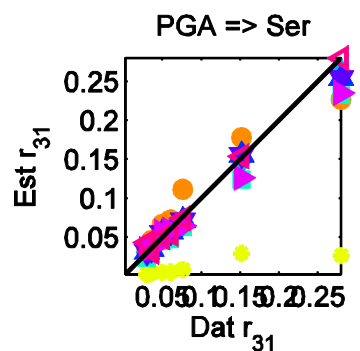
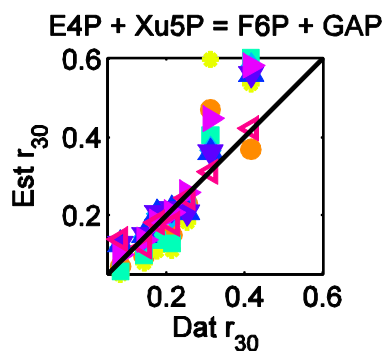
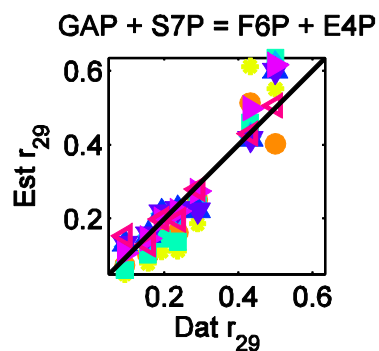
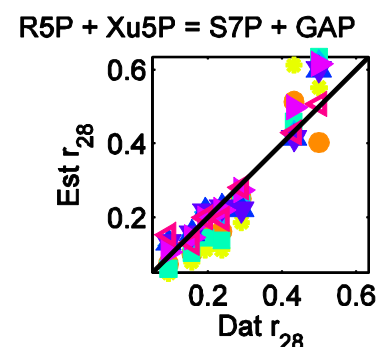
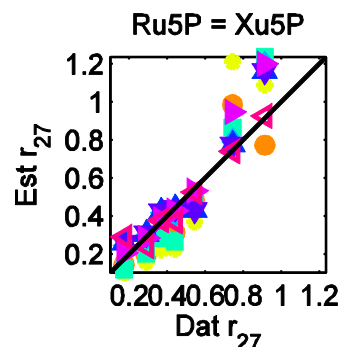
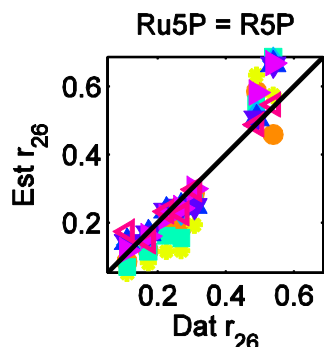
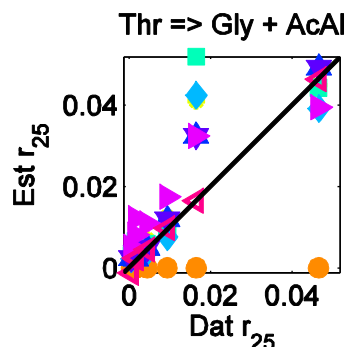
F: Saccharomyces cerevisiae experimental case study – Additional Results

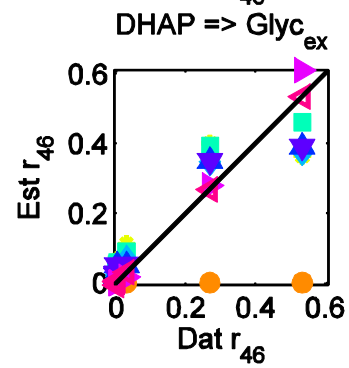
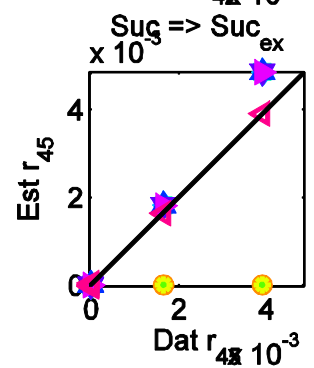
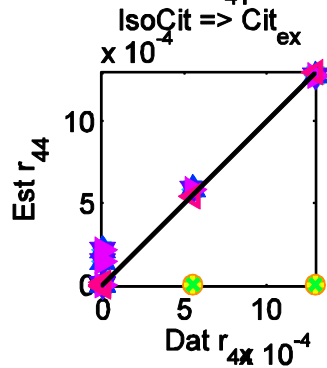
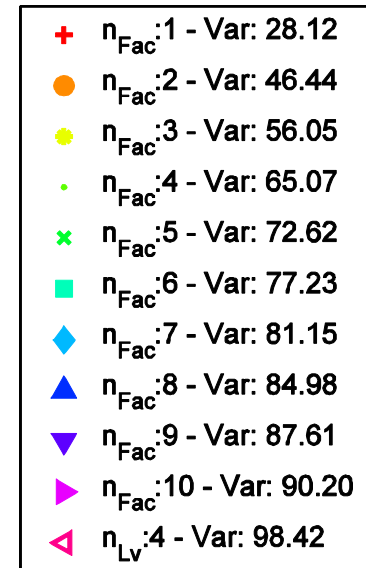
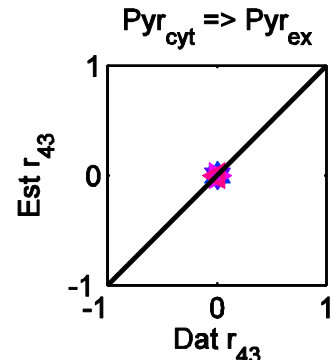
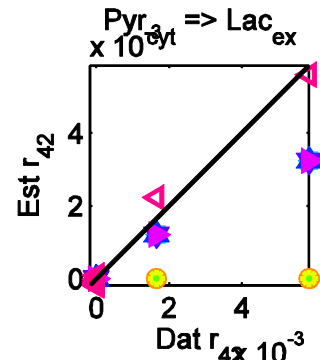
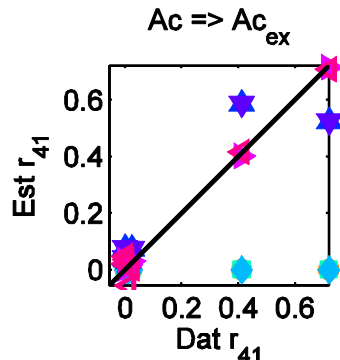
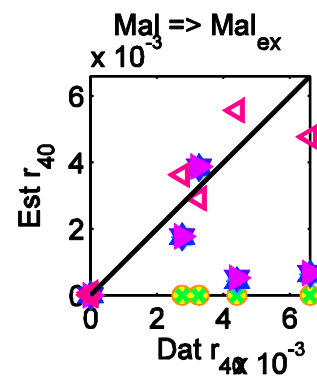
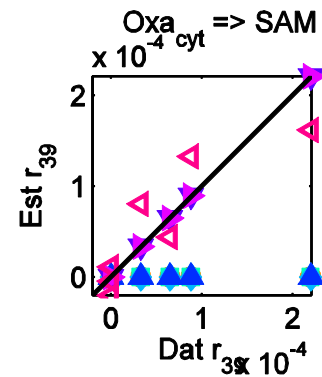
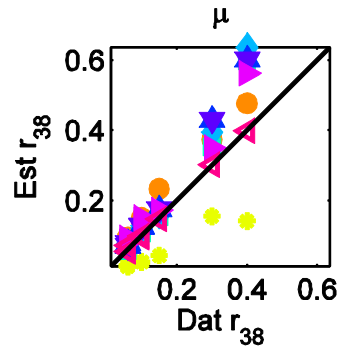
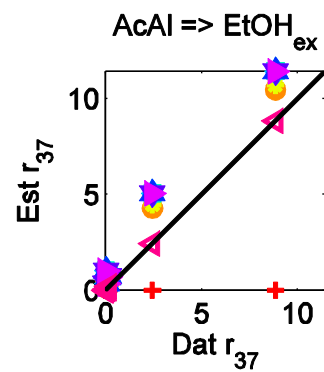
Estimations of reactions for PEMA with different numbers of factors and PCA

The plots of the estimated reaction over the measured values for different numbers of factors is in the following shown for all reactions. The legend in the last plot contains the information of the correspondences between color/symbols and the number of factors.



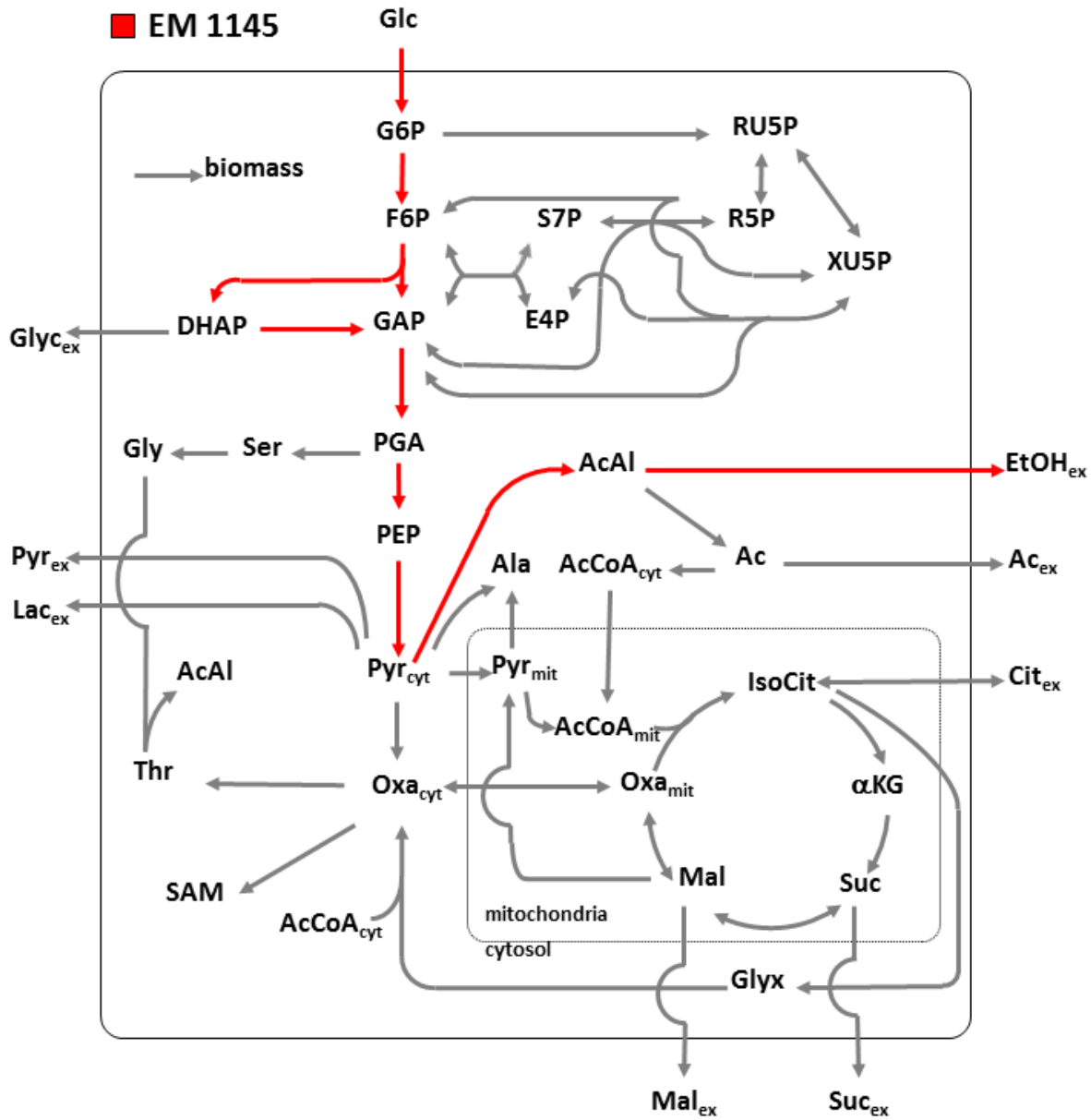


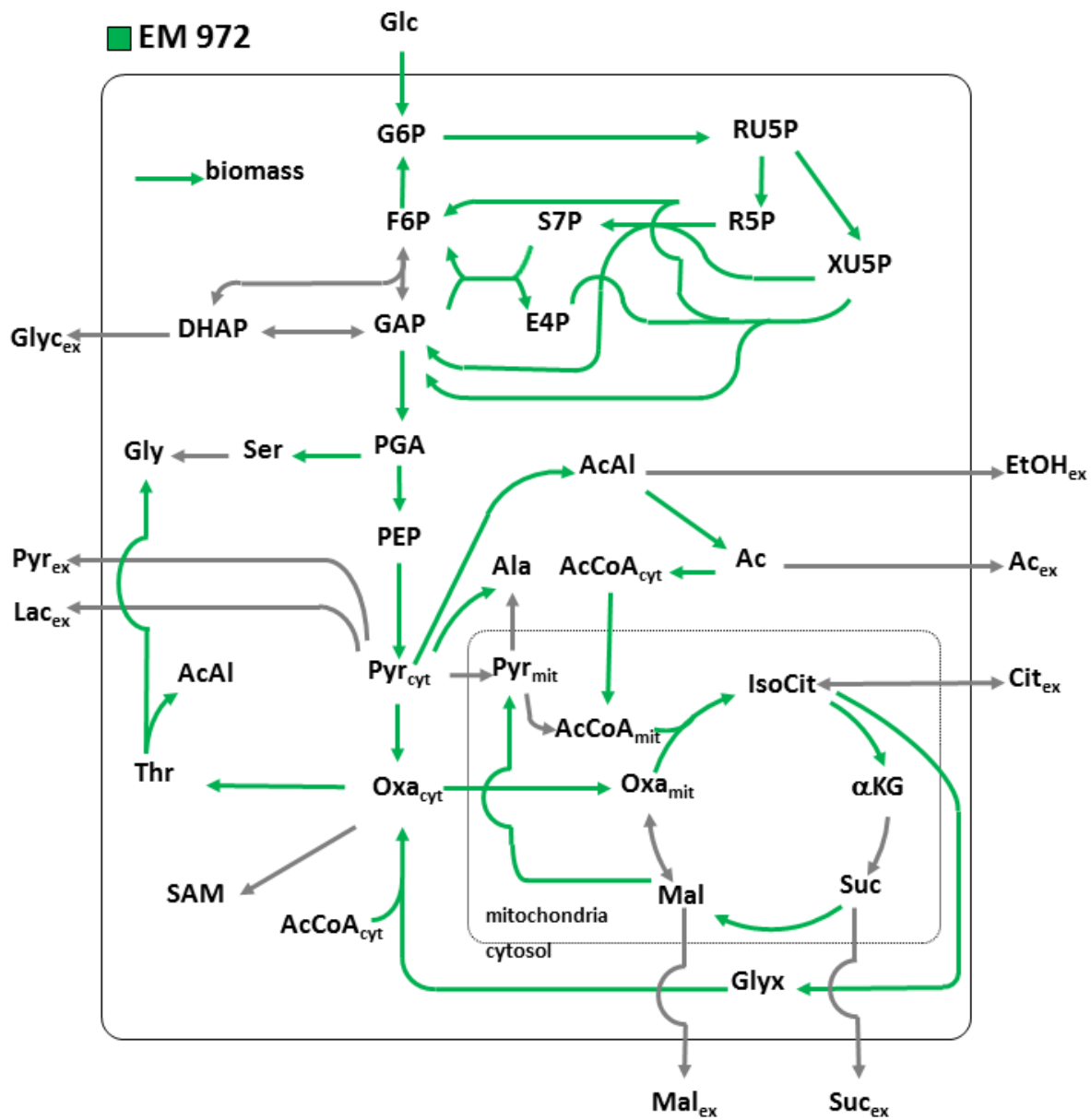


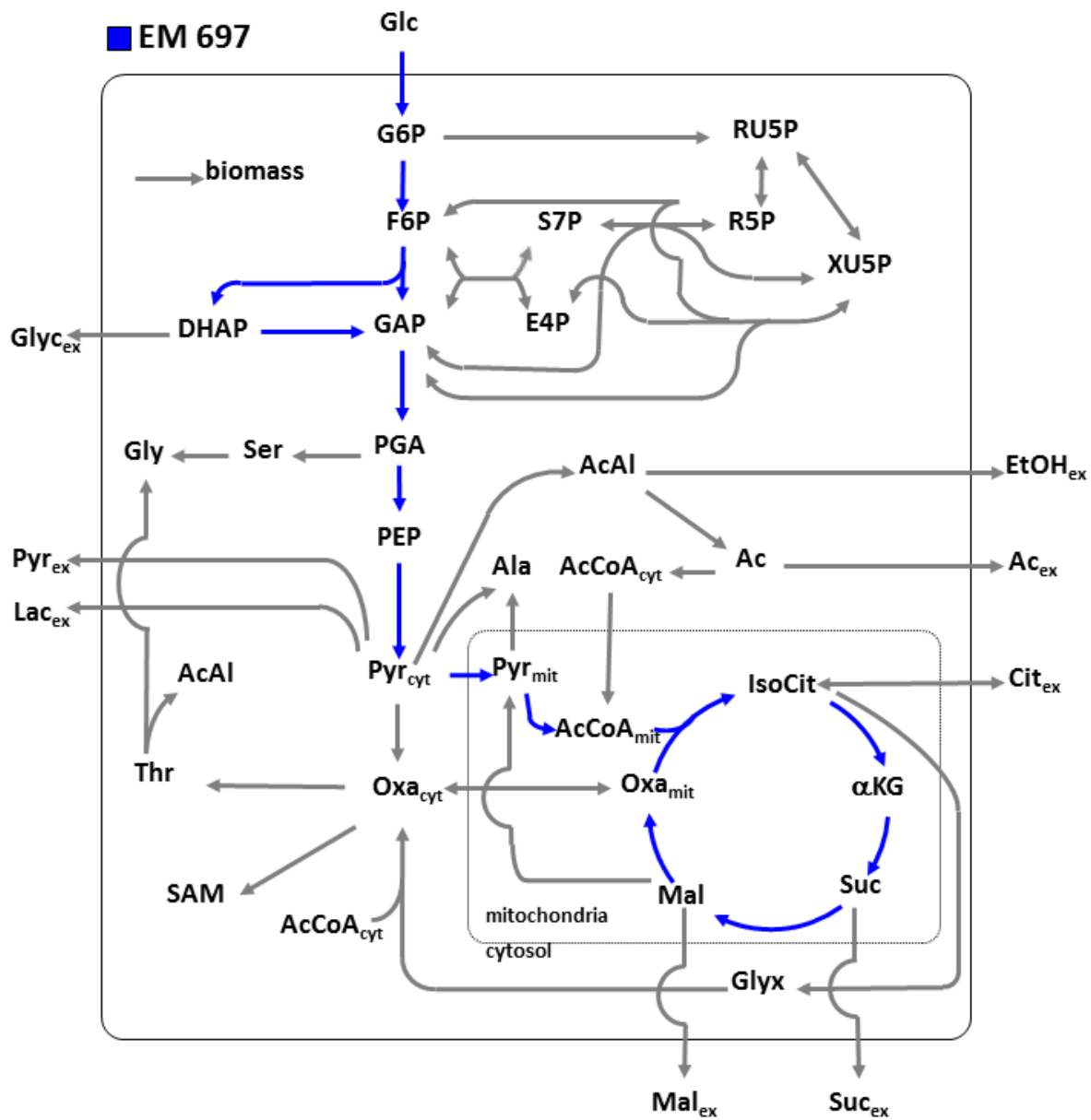


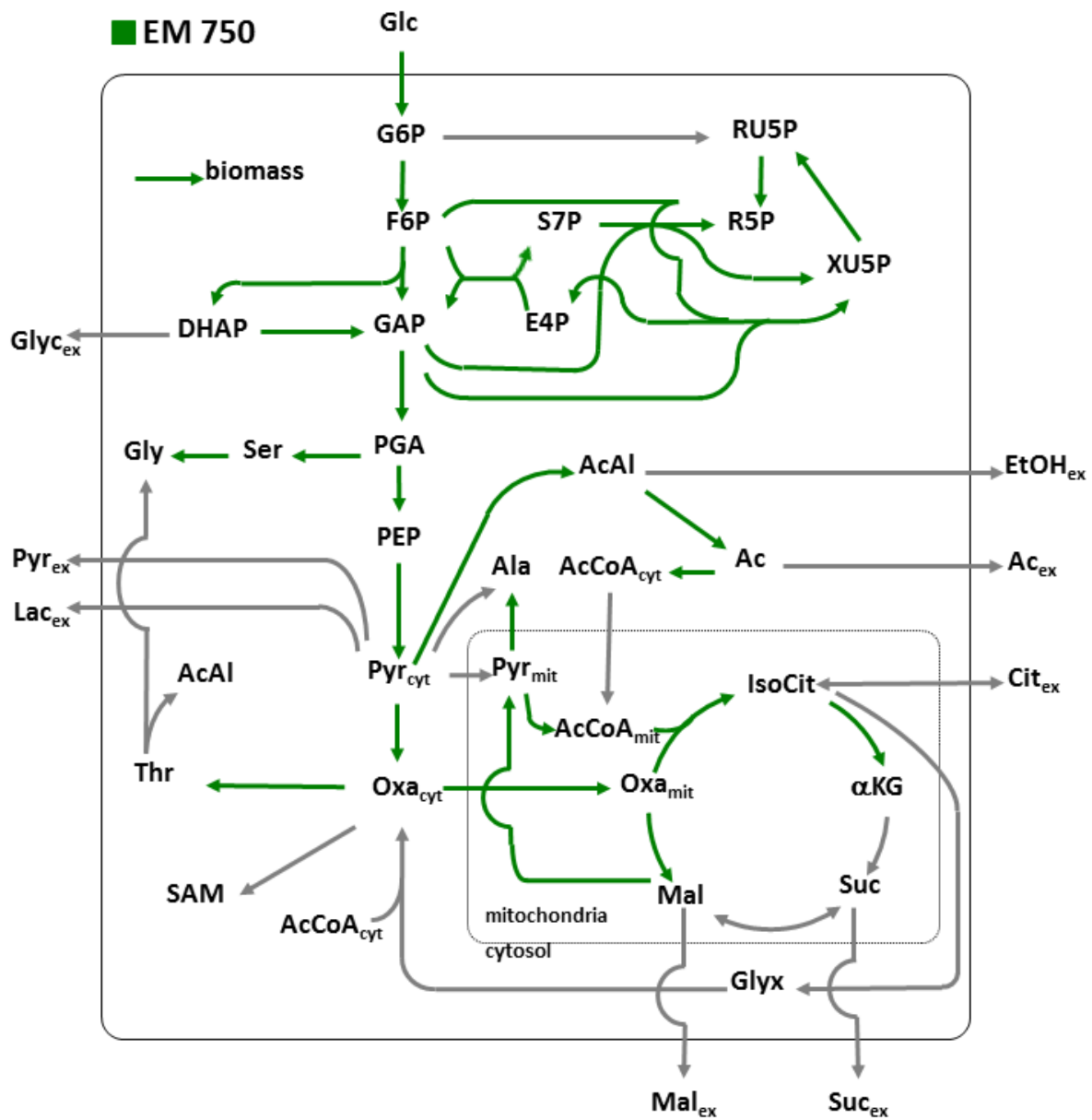
Active Elementary Modes

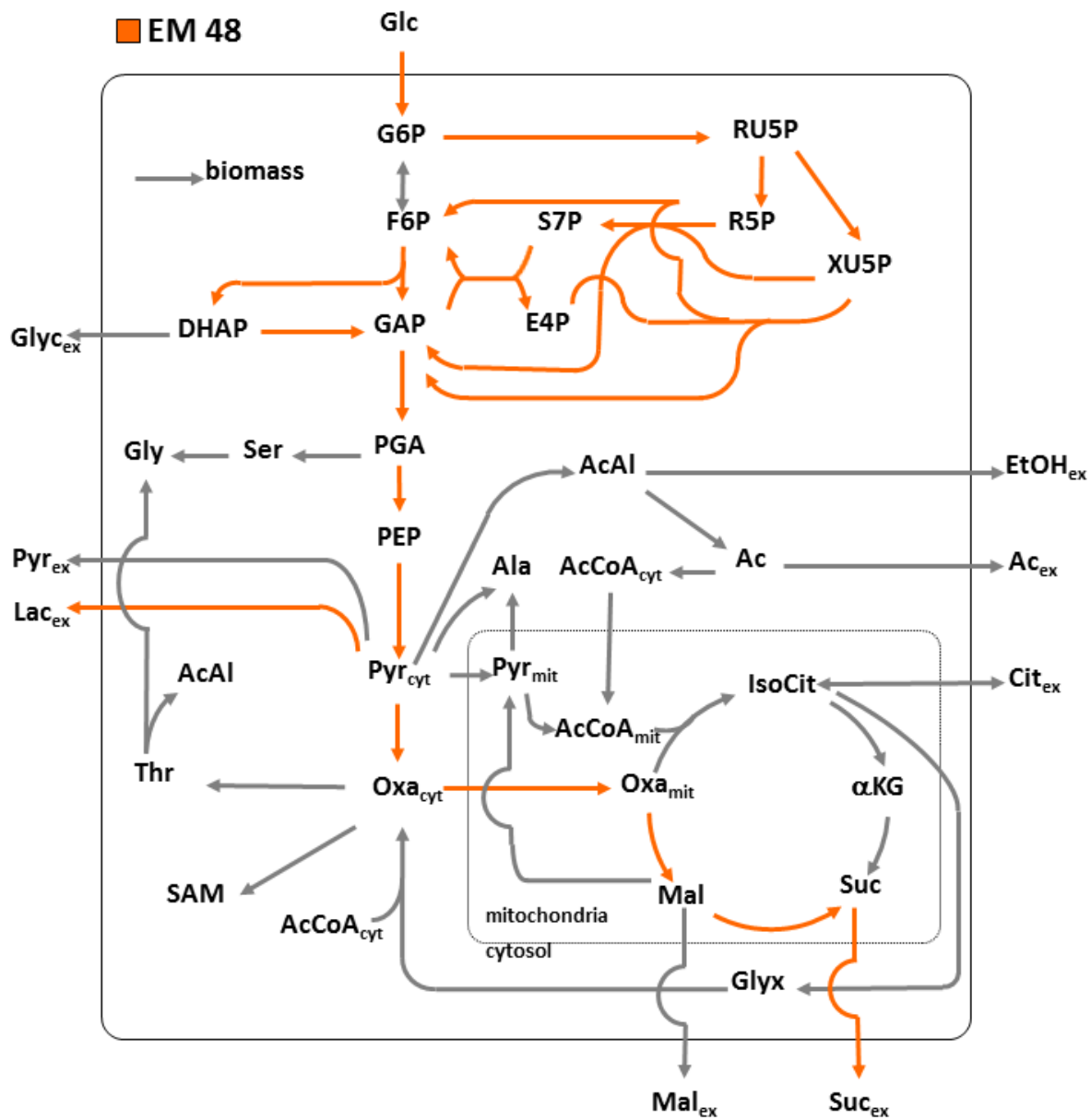
In the following ten pages, the active Elementary Flux Modes selected by PEMA with ten factors are highlighted in different colors in the metabolic network, which was adapted from Hayakawa et al [2].

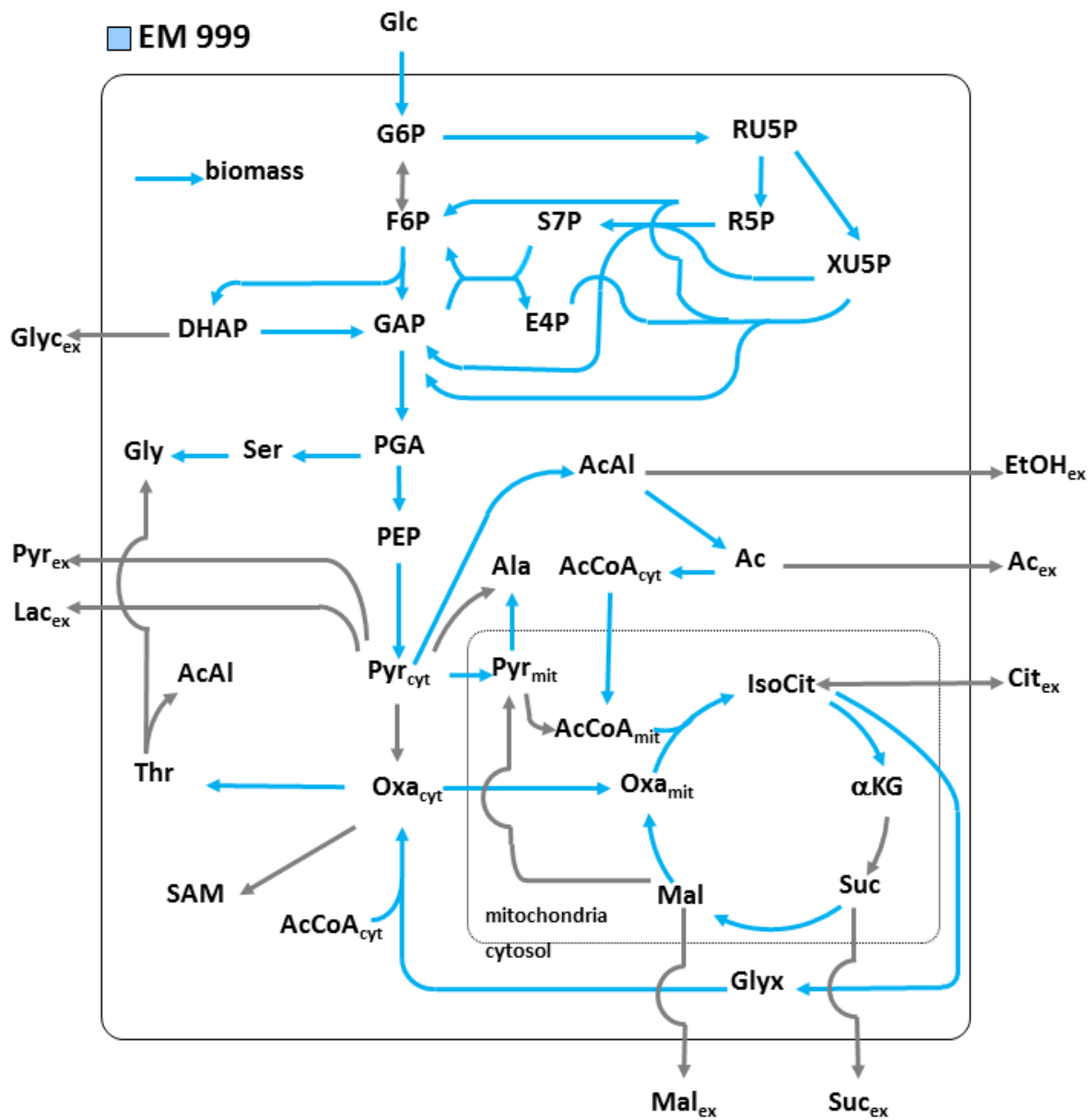


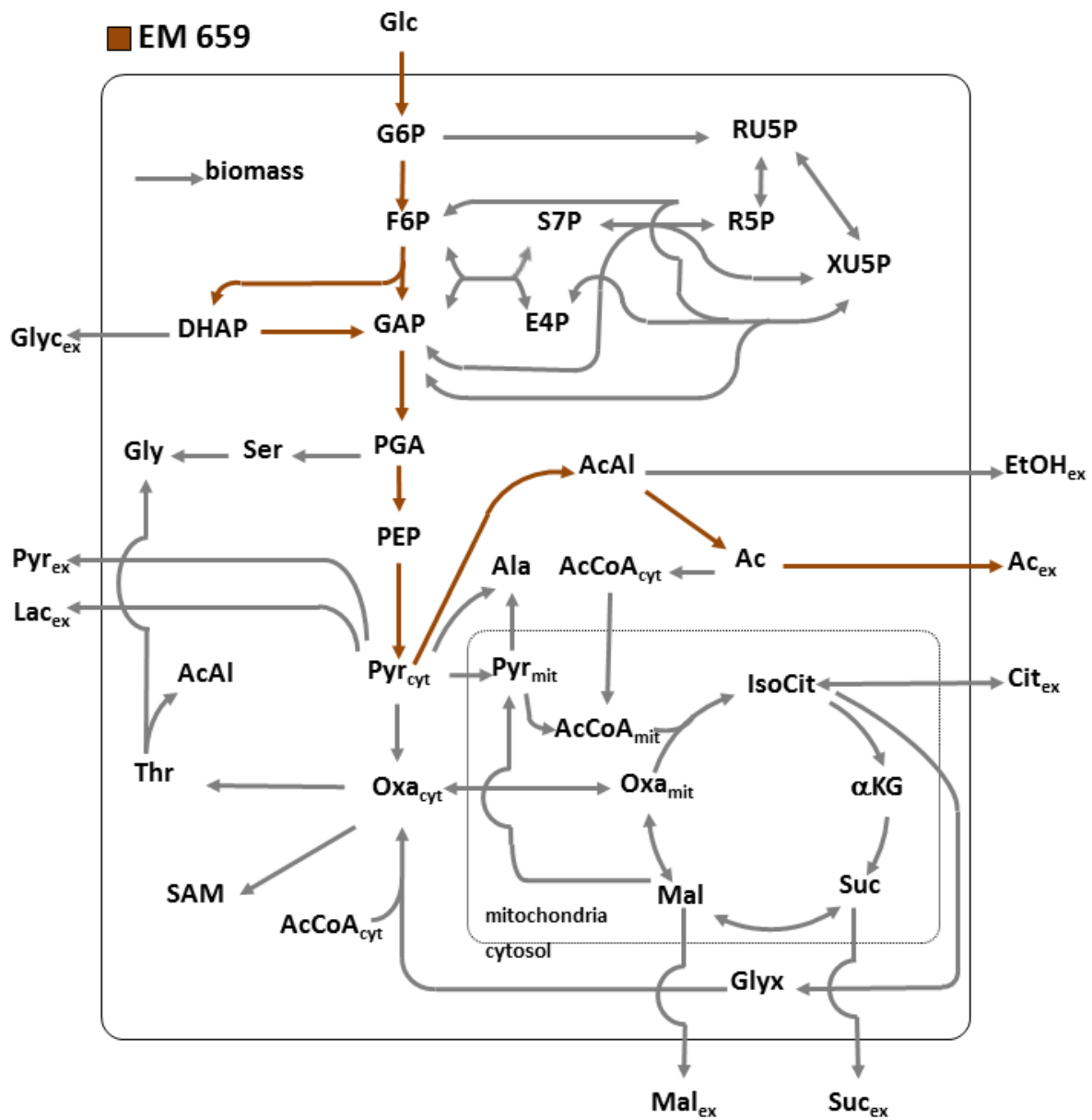


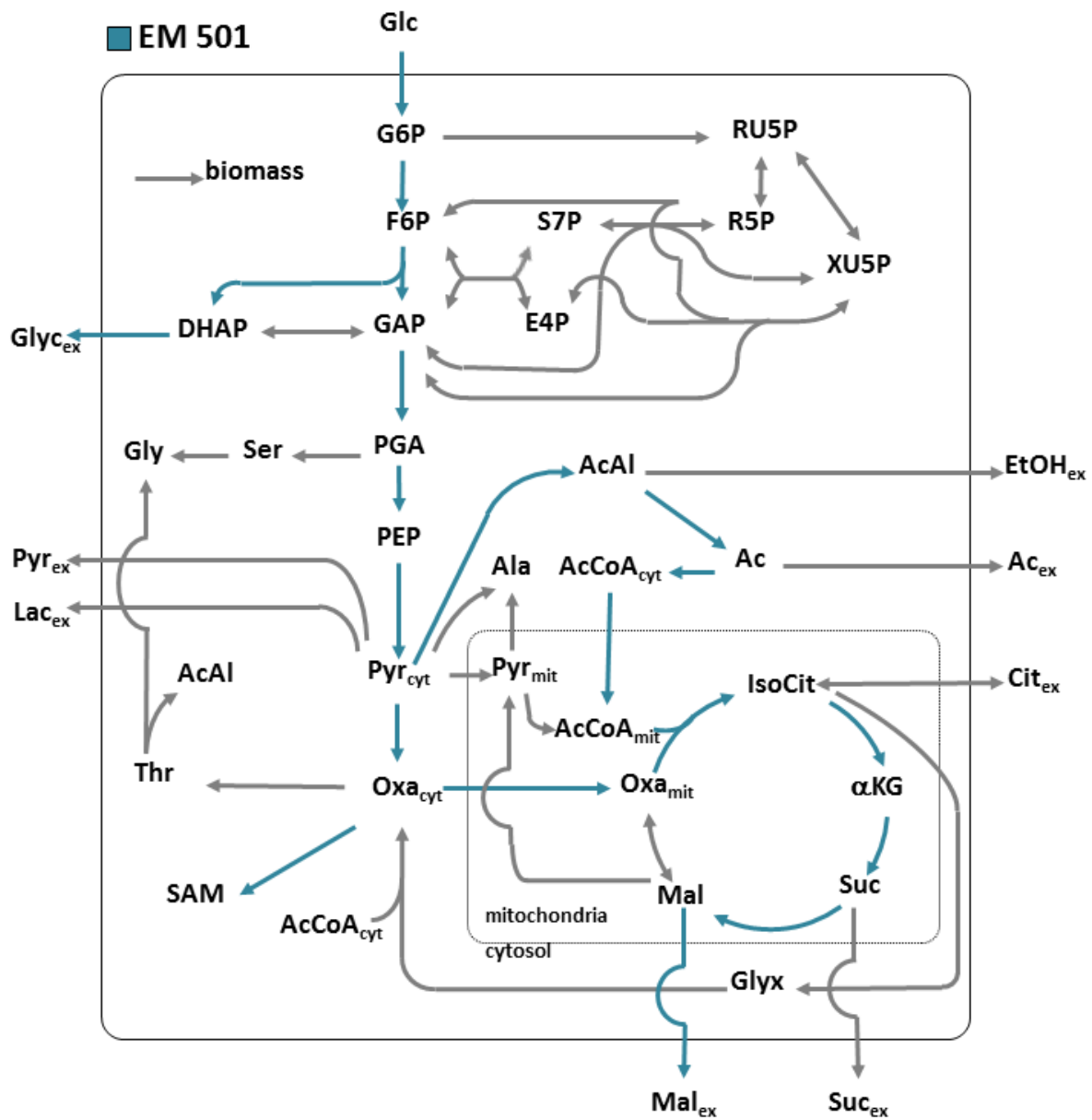












References

1. Tortajada M, Llaneras F, Pico J (2010) Validation of a constraint-based model of *Pichia pastoris* metabolism under data scarcity. *BMC Systems Biology* 4: 115.
2. Hayakawa K, Kajihata S, Matsuda F, Shimizu H ¹³C-metabolic flux analysis in S-adenosyl-L-methionine production by *Saccharomyces cerevisiae*. *Journal of Bioscience and Bioengineering*.
3. Gianchandani EP, Oberhardt MA, Burgard AP, Maranas CD, Papin JA (2008) Predicting biological system objectives de novo from internal state measurements. *BMC Bioinformatics* 9: 43-43.