

Soft sensor for monitoring dynamic changes in cell composition

Sebastián Espinel-Ríos^{*,**} Bruno Morabito^{*}
Katja Bettenbrock^{**} Steffen Klamt^{**} Rolf Findeisen^{****}

^{*} *Laboratory for System Theory and Automatic Control,
Otto von Guericke University Magdeburg, Germany*

^{**} *Analysis and Redesign of Biological Networks, Max Planck Institute
for Dynamics of Complex Technical Systems, Germany*

^{***} *Bioprocess Engineering, Max Planck Institute for Dynamics of
Complex Technical Systems, Germany*

^{****} *Control and Cyber-Physical Systems Laboratory, TU Darmstadt,
Germany (e-mail: rolf.findeisen@iat.tu-darmstadt.de)*

Abstract: Precise and accurate online monitoring methods are needed to enable smart biomanufacturing and automation. Most of the sensors available focus on process parameters such as metabolite and dissolved gas concentrations, cell density or viability, among other variables like pH, temperature, etc. In this work, we develop a soft sensor algorithm to estimate the cell composition online, a very important aspect often overlooked in the bioprocess monitoring literature. Our strategy is based on full information estimation, an optimization-based estimator that takes into account the dynamics of the cell metabolism and considers all the available measurements from the beginning of the process, thus it has a *memory* effect. Being able to track dynamic changes in cell composition can open the door to promising applications, e.g., predictive control and automation of biosystems. As a case study, we consider the *Escherichia coli*'s metabolism growing on glycerol under different levels of oxygen supply. We compare the performance of our soft sensor method against resource balance analysis, a previously proposed estimator based on steady-state assumptions. Overall, the presented full information estimator was able to track the dynamic changes in cell composition significantly more accurately. We also discuss how our estimation strategy can be transformed into a moving horizon estimation, where only the available measurements in a fixed and moving window are considered, thereby reducing possible computational burdens.

Copyright © 2022 The Authors. This is an open access article under the CC BY-NC-ND license (<https://creativecommons.org/licenses/by-nc-nd/4.0/>)

Keywords: Cell composition, soft sensor, full information estimation, resource balance analysis, dynamic enzyme-cost flux balance analysis.

1. INTRODUCTION

Concepts from Industry 4.0 and smart manufacturing are driving the development of novel bioprocess monitoring methods to enable optimization, control and automation of biological systems while satisfying industry quality guidelines (Reyes et al., 2022). The state-of-the-art literature in the field deals mainly with online sensors for monitoring metabolite concentrations such as substrates, intermediates and products, cell density/viability, dissolved gases, as well as some operational parameters like pH, temperature, etc. (Reardon, 2021; Reyes et al., 2022; Fung Shek and Betenbaugh, 2021) In this work, we focus on monitoring dynamic changes in cell composition online, a process parameter that is often overlooked but that can open the door to promising biotechnological applications.

Information on the cell composition, i.e., the intracellular proportion of enzymes, non-catalytic proteins, storage

elements, RNA, DNA, cell wall/ membrane components, lipids and other small molecules, can provide immense insight into the overall state of the cell. For instance, the intracellular levels of proteins such as enzymes can be linked to the activity of metabolic pathways (Noor et al., 2016). Therefore, knowing the cellular concentration of certain proteins can facilitate the modeling of biological systems and the understanding of many cellular processes.

Monitoring the cell composition online can be very valuable in bioprocess control. Recently in our group, we have proposed the use of model predictive control (MPC), an advanced feedback control scheme, to maximize the production efficiency of microbial cell factories via temporal manipulations of key metabolic fluxes for several case studies (Jabarivelisdeh et al., 2020; Espinel-Ríos et al., 2022a,b). In general, MPC is based on a repeated solution of an optimization problem which considers the dynamics of the system and predicts its behaviour over a prediction horizon. The MPC algorithm is solved at every sampling time with the measurements of the current states of the plant thereby generating an optimal input policy. In the referred case studies, we have used dynamic enzyme-cost

^{*} This work was supported by the International Max Planck Research School for Advanced Methods in Process and Systems Engineering (IMPRS ProEng).

flux balance analysis (deFBA) (Waldherr et al., 2015) (cf. Section 2) to accurately model the dynamics of the cell metabolism. The deFBA model contains two time-varying states, one for the extracellular metabolites and another for the intracellular components. Although there are several options for measuring extracellular metabolite concentrations online (Reardon, 2021; Reyes et al., 2022; Fung Shek and Betenbaugh, 2021), this is not the case for the intracellular composition which imposes a big technical limitation to deFBA-based control.

To tackle this technical challenge, Jabarivelisdeh et al. (2020) proposed the use of resource balance analysis (RBA) to estimate the cell composition from cell dry weight measurements and the known values of manipulated metabolic fluxes. It is based on resource allocation theory, i.e., it assumes that the cell allocates resources optimally to maximize the cell growth. Despite its relative simplicity, RBA does not always provide a reliable quantitative estimation of the cell composition. This is in big part explained by the fact that RBA considers quasi-steady-state conditions, which is a very optimistic assumption given the dynamic nature of metabolism. Furthermore, it is limited to the current state of the cell and does not use past measurements for the estimation.

In this work, we outline the use of a full information estimation (FIE) algorithm (Rawlings et al., 2017) as a better alternative to RBA for inferring dynamic changes in cell composition. Our estimation strategy offers in principle better theoretical properties in terms of optimality since it is based on the dynamic model of the metabolism and can consider all process measurements starting from the initial time, hence it has a *memory* effect. As an application example, we consider the *E. coli*'s glycerol metabolism under different oxygen limitation levels. Note that we focus our analysis to state estimation and do not consider feedback control in the current study. The remainder of this paper is as follows: in Section 2 we outline the deFBA modeling framework, in Section 3 we summarize the RBA and FIE algorithms and highlight their differences, and in Section 4 we present our application example.

2. CONSTRAINT-BASED DYNAMIC MODEL

We provide a summarized version of the deFBA framework, for more detailed information about the model derivation and assumptions refer to Waldherr et al. (2015). The dynamics of the cell metabolism is described as the following optimization problem with constraints

$$\max_{V(\cdot)} \int_{t_0}^{t_{\text{deFBA}}} B(p) dt \quad (1a)$$

$$\text{s.t.} \quad \left[\frac{dz(t)}{dt} \quad \frac{dp(t)}{dt} \right]^T = S_{zp}V(t), \quad (1b)$$

$$[z(t_0) \quad p(t_0)]^T = [z_0 \quad p_0]^T, \quad (1c)$$

$$0 = \frac{dm(t)}{dt} = S_m V(t), \quad (1d)$$

$$B = b^T p(t), \quad (1e)$$

$$\sum_{j \in \text{cat}_i} \left| \frac{V_j(t)}{k_{\text{cat},j}} \right| \leq p_i, \quad \forall i \in [1, n_{p_i}] \quad (1f)$$

$$\varphi_Q b^T p(t) \leq p_Q(t), \quad (1g)$$

$$V_{\min}(t) \leq V(t) \leq V_{\max}(t). \quad (1h)$$

The objective function is the maximization of the cell dry weight $B \in \mathbb{R}$ integral from t_0 to t_{deFBA} . The decision variable is the vector of reaction fluxes $V \in \mathbb{R}^{n_V}$. The optimization is then subject to equality and inequality constraints that reduce the solution space. The change of extracellular metabolites $z \in \mathbb{R}^{n_z}$ and cell components $p \in \mathbb{R}^{n_p}$ (e.g., catalytic enzymes, ribosomes and non-catalytic/quota components) with their corresponding initial conditions are described by Eqs. (1b)-(1c). The model assumes quasi-steady-state conditions of the intracellular metabolites $m \in \mathbb{R}^{n_m}$, Eq. (1d). $S_{zp} \in \mathbb{R}^{n_z+n_p, n_V}$ is the stoichiometric matrix of z and p , and $S_m \in \mathbb{R}^{n_m, n_V}$ is the stoichiometric matrix of m . The cell dry weight is expressed as $b^T p$, where $b \in \mathbb{R}^{n_b}$ contains the molecular weights of p , Eq. (1e). The upper bound of V is limited by the product of the catalytic enzyme concentrations and the catalytic constants $k_{\text{cat}} \in \mathbb{R}^{k_{\text{cat}}}$, Eq. (1f). A minimal fraction $\varphi_Q \in [0, 1]$ of the cell dry weight corresponds to a lumped quota compound p_Q , Eq. (1g). Finally, the fluxes are further narrowed down to feasible or biologically sound bounds, Eq. (1h). For example, the upper and lower bounds of externally-regulated or manipulated fluxes $V_{\text{reg}} \in V$ can be set to be known and equal.

3. ESTIMATION OF CELL COMPOSITION

Our goal is to estimate the cell composition at every sampling time t_k . Let $\hat{p}(t_k)$ be this estimate.

3.1 Resource balance analysis

The RBA estimation strategy is summarized in Fig. 1-A. The algorithm is formulated as (Jabarivelisdeh et al., 2020)

$$\max_{\mu, V, p} \quad \mu \quad (2a)$$

$$\text{s.t.} \quad S_p V - \mu p = 0, \quad (2b)$$

$$\tilde{B} = b^T p, \quad (2c)$$

$$0 = S_m V(t), \quad (2d)$$

$$\sum_{j \in \text{cat}_i} \left| \frac{V_j}{k_{\text{cat},j}} \right| \leq p_i, \quad \forall i \in [1, n_{p_i}] \quad (2e)$$

$$\varphi_Q b^T p \leq p_Q, \quad (2f)$$

$$V_{\min} \leq V \leq V_{\max}. \quad (2g)$$

This static optimization maximizes the growth rate μ for each cell dry weight measurement \tilde{B} by allocating p to render an optimal V distribution. $S_p \in \mathbb{R}^{n_p, n_V}$ is the stoichiometric matrix of p . Note that \tilde{B} constrains the solution space of p in Eq. (2c). In contrast to deFBA, RBA assumes that the rate of production of p is equally diluted by cell growth, hence there is no accumulation of p . The steady-state assumption turns the optimization problem overall simpler and easier to solve. Let, μ^* , p^* and V^* be the optimal values of μ , p and V , respectively; thus, $\hat{p}(t_k) := p^*$.

3.2 Full information estimation

The FIE strategy is summarized in Fig. 1-B. Let us denote a general optimization variable calculated at time t_i with

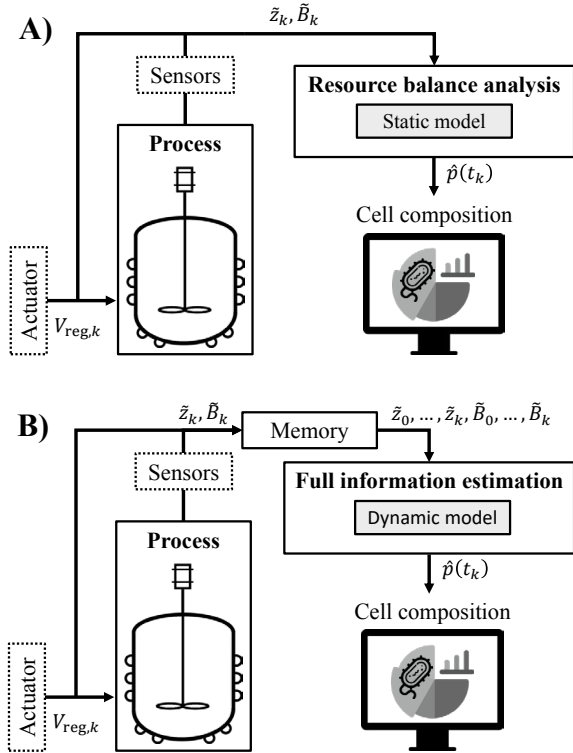


Fig. 1. Estimation strategy of the cell composition with the A) RBA and B) FIE algorithms.

(\cdot) _{i} . For simplicity, we assume that the measurements are taken at equidistant sampling times, with sampling interval Δt , although non-equidistant sampling times are also conceivable. Let us define $x(t) := [z(t)^T, p(t)^T]^T$. The FIE algorithm (Rawlings et al., 2017) solves at every sampling time the following optimization problem

$$\min_{x_0, w, k_{\text{cat}}} \left\| \begin{bmatrix} x_0 \\ k_{\text{cat}} \end{bmatrix} - \begin{bmatrix} \bar{x}_0 \\ \bar{k}_{\text{cat}} \end{bmatrix} \right\|_P^2 + \sum_{i=0}^k \|y(t_i) - \tilde{y}_i\|_R^2 +$$

$$+ \|w_i\|_Q^2$$

$$\text{s.t.} \quad \max_{V(\cdot)} \int_{t_i}^{t_i+\Delta t} B(p(t)) dt \quad (3b)$$

$$\text{s.t.} \quad x(t_i + \Delta t) = x(t_i)^T +$$

$$+ \int_{t_i}^{t_i+\Delta t} S_{zp} V(t) dt + w_i, \quad (3c)$$

$$x(t_0) = x_0, \quad (3d)$$

$$y(t_i) = [z(t_i) \ B(p(t_i))], \quad (3e)$$

$$\text{Eqs. (1d)-(1h)}, \quad (3f)$$

$$i \in [0, \dots, k],$$

where $P \in \mathbb{R}^{n_z+n_p+n_{k_{\text{cat}}}, n_z+n_p+n_{k_{\text{cat}}}}$, $R \in \mathbb{R}^{n_z+1, n_z+1}$ and $Q \in \mathbb{R}^{n_z+1, n_z+1}$ are weighting matrices and $k+1$ is the number of samples collected up to time t_k ¹. The notation ($\bar{\cdot}$) indicates the prior information of that variable. For example, this could be the best measurement or best guess available of that specific variable. For example, \bar{k}_{cat} is the prior information of the k_{cat} parameters, while \bar{z}_0 of

the extracellular states at time t_0 . The objective function (3a) has three terms. The first term weights the difference between the prior information and the estimated states at the initial time and k_{cat} . The second term, weights the difference between the predicted measurements and the real measurements \tilde{y}_i for every $i \in [0, \dots, k]$. Finally, the last term weights the effect of the *state noise*. To reflect possible model uncertainty, the state noise w_i is added to the model at every sampling time (cf. (3c)). The role of the state noise is to help reconcile the model with the measurements by modifying the right-hand side of the ordinary differential equation. Using state noise is useful in the presence of model mismatch due to, e.g., parameter or structural model uncertainties. The weights P, Q , and R should be chosen accordingly to the importance of the different terms. For example, in the presence of low measurements noise and high model uncertainty, the weight R should be larger than Q to reflect that we “trust” more the measurements than our model. Similarly, if we trust our prior values more than the measurements, the matrix P should be relatively larger than R .

Note that the optimization variables are the extracellular concentration z_0 and the cellular components p_0 , the k_{cat} and w which is the tuple collecting the state noise for every sampling time, i.e., $w = (w_0, w_1, \dots, w_{N(t_i)-1})$ where $N(t_i) \in \mathbb{R}$ is the number of measurements taken at time t_i . Once the optimal values z_0^* , p_0^* , w^* and k_{cat}^* are found, by using (3c) the estimated cell composition at the current time t_k is $\hat{p}(t_k) := p^*(t_k)$. The FIE can be seen as a constrained least square estimation where the dynamics enter the constraints of the optimization problem. In this case, the resulting optimization problem is bilevel since the deFBA model is part of the constraints of the FIE.

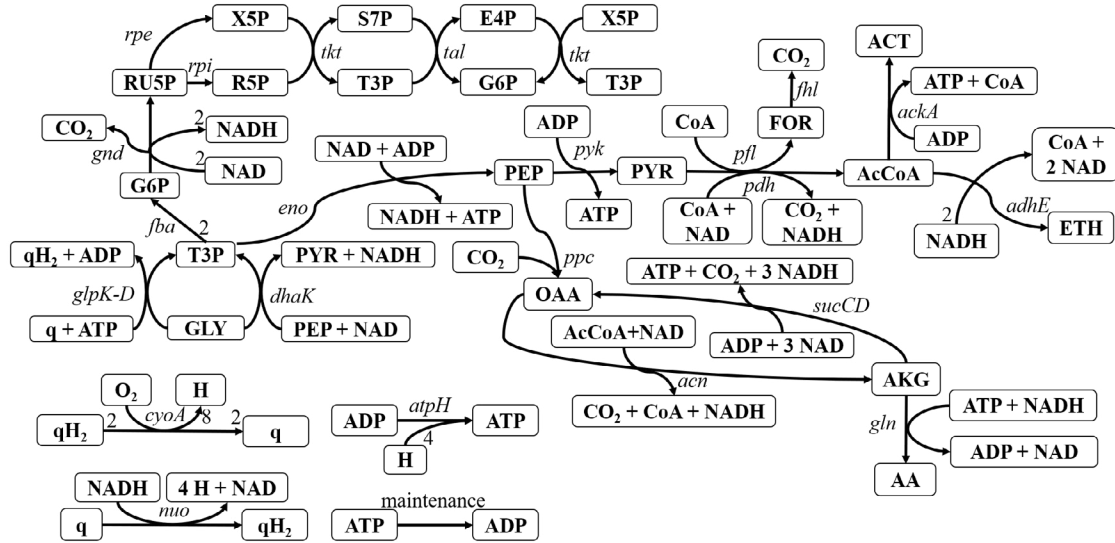
4. CASE STUDY: GLYCEROL METABOLISM

We consider *E. coli* cells growing on glycerol with an oxygen uptake rate (OUR) input policy such that the cells experience several metabolic states throughout the process. Different levels of oxygen supply impose the cells with redox and energy constraints that lead to temporal metabolic adaptations. Therefore, the cell composition is expected to change dynamically to cope with the fluctuating environment, providing us with an *ad hoc* setup to test our soft sensor method.

To model the dynamics of this system, we used an adapted resource allocation model from Jabarivelisdeh et al. (2020). The network covers relevant metabolic reactions related to glycerol catabolism, glycolysis, the pentose phosphate pathway, anaerobic fermentation and respiration (see Fig. 2). Six extracellular metabolites are included in z : glycerol, acetate, ethanol, formate, O_2 and CO_2 . Although not shown in Fig. 2, the model also comprises production reactions for all p -elements such as catalytic enzymes, ribosomes, and a lumped quota compound.

We seek to elucidate whether using FIE, which is based on the dynamics of the metabolism and takes into account current and previous measurements, can result in a better estimation of the cell composition compared to RBA, which is based on quasi-steady-state assumptions and only considers the current measurements. Therefore, for

¹ Our measurements start from t_0 , hence at t_k we have $k+1$ measurements.



AA: amino acid; AcCoA: acetyl coenzyme A; ACT: acetate; ADP: Adenosine diphosphate; AKG: α -ketoglutarate; ATP: adenosine triphosphate; CO₂: carbon dioxide; CoA: coenzyme A; E4P: erythrose-6-phosphate; ETH: ethanol; FOR: formate; G6P: glucose-6-phosphate; GLY: glycerol; NAD: nicotinamide adenine dinucleotide; NADH: nicotinamide adenine dinucleotide reduced; O₂: oxygen; OAA: oxaloacetate; PEP: phosphoenolpyruvate; PYR: pyruvate; q: ubiquinol; qH₂: ubiquinol; R5P: ribose 5-phosphate; RU5P: ribulose 5-phosphate; S7P: sedoheptulose 7-phosphate; T3P: glyceraldehydes-3-phosphate; X5P: xylulose 5-phosphate.

Fig. 2. Metabolic pathway of *E. coli* growing on glycerol. Genes of catalytic enzymes are written in italics. For example, *cyoA* is the gene of the enzyme p_{cyoA} . Production reactions of cell components are not depicted.

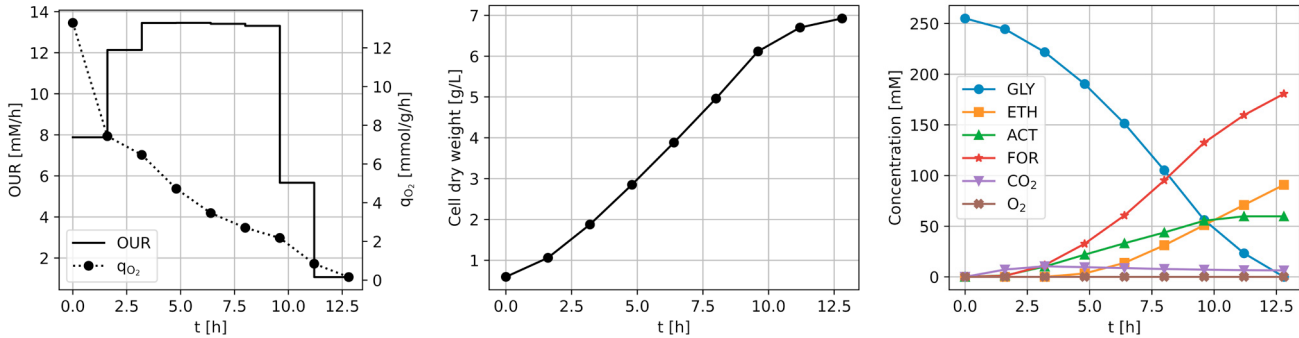


Fig. 3. Dynamic changes of z and B with varying oxygen uptake rates at different sampling times.

the sake of simplicity, we consider that measurements of z and B can be obtained online with negligible noise. Furthermore, we assume that the applied OUR is equal to the flux catalyzed by the enzyme p_{cyoA} . We also consider no accumulation of oxygen, meaning that the provided oxygen is immediately depleted by the cells. In other words, the OUR is equal to the oxygen transfer rate from the gas to the liquid, which is a common assumption when modeling oxygen dynamics in processes without substrate starvation (Humbird and Fei, 2016).

The required measurements for the estimation algorithms were obtained from deFBA simulations (see Fig. 3). To facilitate the discussions, we also plotted the specific oxygen uptake rate q_{O_2} -normalized per cell mass-. At $q_{O_2} \gtrsim 5$, acetate and formate were the main products, with little-to-none ethanol formation. Afterwards, at lower q_{O_2} values, the cells started to produce significant amounts of ethanol, while the acetate production rate began to slow-down and finally stagnated. Similarly, the cell growth rate decreased with decreasing q_{O_2} . Overall, this was the

expected behaviour of *E. coli* cells growing on glycerol (Durnin et al., 2009).

The estimation of the cell composition at the different sampling times is presented in Fig. 4. The optimizations were solved using CasADi (Andersson et al., 2019). The FIE bilevel optimization was converted into a single-level problem by applying the Karush–Kuhn–Tucker conditions to the inner problem, following an optimistic approach (Dempe and Franke, 2019). We penalized the complementarity slackness constraint in the objective function to facilitate the numerics. Note that we only considered the second term of the FIE objective function (3a) for our case study, and k_{cat} was assumed constant and not part of the optimization variables.

To compare the quality of the estimation by RBA and FIE, we computed the standard error of the estimate (SE) for each cell component p_i as

$$SE = \sqrt{\frac{\sum_{j=0}^{t_N} (p_{i,j} - \hat{p}_{i,j})^2}{E_i}}, \quad \forall i \in [1, n_{p_i}], \quad (4)$$

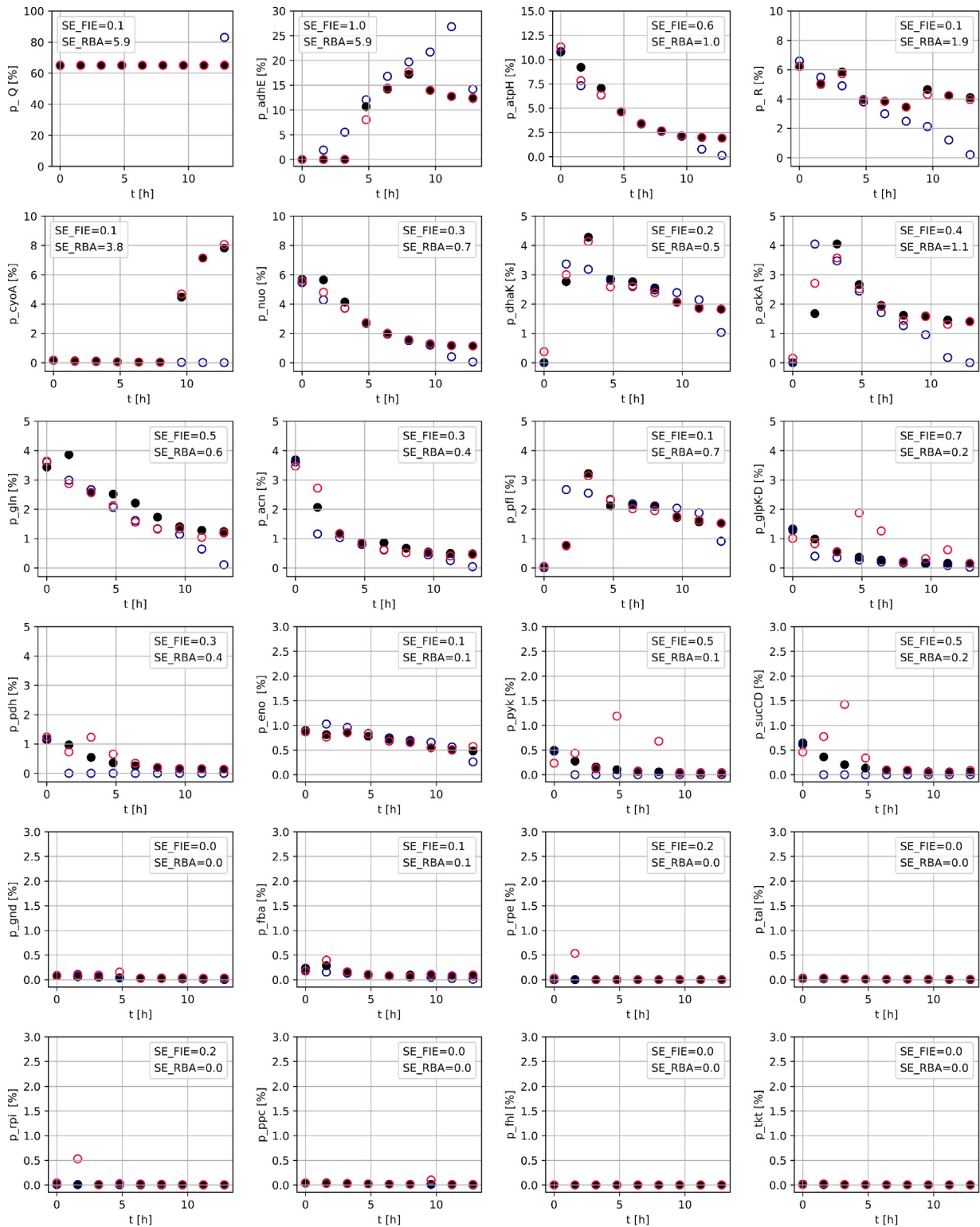


Fig. 4. Estimation of the cell composition in percentage of cell dry weight. ● Exact value, ○ FIE estimate, ○ RBA estimate. The standard errors of the estimates for the FIE (SE_FIE) and the RBA (SE_RBA) algorithms are presented. Q: lumped quota element, R: ribosome; p_i : catalytic enzyme encoded by the gene i . Note that, in principle, Q could be lumped to a greater or lesser extent depending on the user needs.

where E_i is the total number of estimates. The SE is a measure of the variability or dispersion of the predictions, and it has the same units of the measurements. The smaller the SE, the closer the predictions are to the exact values. The *real* values of p were extracted from the deFBA simulations.

Before we discuss the estimation quality, it is worth noting that, as expected, the cell components experienced dynamic changes throughout the cultivation as a result of the metabolic adaptations and resource allocation. See e.g., the trends for enzymes p_{adhE} , p_{atpH} , p_{cyoA} , p_{nuo} , p_{dhaK} , p_{ackA} , p_{gln} , and p_{pff} , where the dynamic changes are more noticeable. The observed p dynamics are also correlated to the z dynamics in Fig. 3. For instance, p_{ackA} increases during the first 5 h (high oxygen supply), while p_{adhE} remains zero. Afterwards, with increasing oxygen limitation, p_{adhE} starts to accumulate while p_{ackA} decreases. This matches the rise in ethanol production and drop in acetate concentration previously described.

Based on the SE, the FIE strategy performed better than RBA for 12 out of 24 cell components. This is already quite an improvement considering that these were basically the most representative elements in terms of percentage cell weight. For other 7 cell components, there was no difference between RBA and FIE, having the same SE values. Only for the cases of p_{glpk-D} , p_{pyk} , p_{sucCD} , p_{rpe} and p_{rpi} , RBA performed slightly better. However, each of the latter components was never higher than 1.5 % of the cell dry weight. Moreover, the FIE predictions were still within acceptable ranges. Based on these results, we conclude that FIE is a very good candidate to substitute previously proposed methods like RBA as it can estimate dynamic changes in cell composition more accurately.

Remark that the FIE algorithm considers all measurements starting from the initial time, which in principle enhances the estimation performance as the process proceeds due to its *memory* effect. Although this is at first glance an advantage, for long processes with a high number of past measurements to consider, this can turn out to be a burden because of the increasing computation effort. If that happens, we suggest to use a moving horizon estimator Rawlings et al. (2017) to reduce the computational complexity. This estimator is similar to an FIE, but instead of considering all the measurements starting from t_0 , only the measurements within a fixed and moving time window are used. In this way, the computational burden is limited and can be adjusted by changing the length of the window. Nevertheless, in our application example this was not concerning since the process time was short and the number of measurements manageable.

5. CONCLUSION AND OUTLOOK

We have developed an FIE algorithm to monitor dynamic changes in cell composition during the operation of bioprocesses. Our method can be regarded as an online soft sensor that is based on the dynamics of the metabolism and uses present and past information of available measurements to infer the intracellular composition. We applied our FIE strategy to estimate the cell composition of

E. coli growing on glycerol with varying levels of OUR. We compared the predictions against previously proposed observers of cell composition such as RBA. In general, FIE showed better estimation performance than RBA using the standard error as a criterion. Future work focuses on using FIE to facilitate the implementation of predictive control strategies applied to metabolic cybergenetic systems and other bioprocesses.

REFERENCES

- Andersson, J.A.E., Gillis, J., Horn, G., Rawlings, J.B., and Diehl, M. (2019). CasADi: a software framework for nonlinear optimization and optimal control. *Math. Program. Comput.*, 11(1), 1–36.
- Dempe, S. and Franke, S. (2019). Solution of bilevel optimization problems using the KKT approach. *Optimization*, 68(8), 1471–1489.
- Durnin, G., Clomburg, J., Yeates, Z., Alvarez, P.J., Zygorakis, K., Campbell, P., and Gonzalez, R. (2009). Understanding and harnessing the microaerobic metabolism of glycerol in *Escherichia coli*. *Biotechnol. Bioeng.*, 103(1), 148–161.
- Espinel-Ríos, S., Bettenbrock, K., Klamt, S., and Findeisen, R. (2022a). Maximizing batch fermentation efficiency by constrained model-based optimization and predictive control of adenosine triphosphate turnover. *AIChE J.*, 68(4), e17555.
- Espinel-Ríos, S., Morabito, B., Pohlodek, J., Bettenbrock, K., Klamt, S., and Findeisen, R. (2022b). Optimal control and dynamic modulation of the ATPase gene expression for enforced ATP wasting in batch fermentations. *IFAC-PapersOnLine (to appear)*.
- Fung Shek, C. and Betenbaugh, M. (2021). Taking the pulse of bioprocesses: at-line and in-line monitoring of mammalian cell cultures. *Curr. Opin. Biotechnol.*, 71, 191–197.
- Humbird, D. and Fei, Q. (2016). Scale-up considerations for biofuels. In *Biotechnology for Biofuel Production and Optimization*, 513–537. Elsevier.
- Jabarivelisdeh, B., Carius, L., Findeisen, R., and Waldherr, S. (2020). Adaptive predictive control of bioprocesses with constraint-based modeling and estimation. *Comput. Chem. Eng.*, 135, 106744.
- Noor, E., Flamholz, A., Bar-Even, A., Davidi, D., Milo, R., and Liebermeister, W. (2016). The protein cost of metabolic fluxes: prediction from enzymatic rate laws and cost minimization. *PLoS Comput. Biol.*, 12(11), e1005167.
- Rawlings, J.B., Mayne, D.Q., and Diehl, M. (2017). *Model predictive control: theory, computation, and design*, volume 2. Nob Hill Publishing Madison.
- Reardon, K.F. (2021). Practical monitoring technologies for cells and substrates in biomanufacturing. *Curr. Opin. Biotechnol.*, 71, 225–230.
- Reyes, S.J., Durocher, Y., Pham, P.L., and Henry, O. (2022). Modern sensor tools and techniques for monitoring, controlling, and improving cell culture processes. *Processes*, 10(2), 189.
- Waldherr, S., Oyarzún, D.A., and Bockmayr, A. (2015). Dynamic optimization of metabolic networks coupled with gene expression. *J. Theor. Biol.*, 365, 469–485.