*Systems Biology*

# PARROT: Prediction of enzyme abundances using protein-constrained metabolic models

Maurício Alexander de Moura Ferreira[1], Wendel Batista da Silveira[1] and Zoran Nikoloski [2,3,*]

[1]Department of Microbiology, Federal University of Viçosa, Viçosa, Minas Gerais, Brazil, [2]Bioinformatics, Institute of Biochemistry and Biology, University of Potsdam, Potsdam, Germany, [3]Systems Biology and Mathematical Modelling, Max Planck Institute of Molecular Plant Physiology, Potsdam, Germany.

*To whom correspondence should be addressed.

## Abstract

**Motivation:** Protein allocation determines activity of cellular pathways and affects growth across all organisms. Therefore, a variety of experimental and machine learning approaches has been developed to quantify and predict protein abundances, respectively. Yet, despite advances in protein quantification, it remains challenging to predict condition-specific allocation of enzymes in metabolic networks.

**Results:** Here we propose a family of constrained-based approaches, termed PARROT, to predict enzyme allocations based on the principle of minimizing the enzyme allocation adjustment using protein-constrained metabolic models. To this end, PARROT variants model the minimization of enzyme reallocation using four different (combinations of) distance functions. We demonstrate that the PARROT variant that minimizes the Manhattan distance of enzyme allocations outperforms existing approaches based on the parsimonious distribution of fluxes or enzymes for both *Escherichia coli* and *Saccharomyces cerevisiae*. Further, we show that the combined minimization of flux and enzyme allocation adjustment leads to poor and inconsistent predictions. Together, our findings indicate that minimization of resource rather than flux redistribution is a governing principle determining steady-state pathway activity for microorganism grown in suboptimal conditions.

**Availability and implementation:** The implementation of PARROT can be found in the GitHub repository: https://github.com/mauricioamf/PARROT

**Contact:** zoran.nikoloski@uni-potsdam.de

**Supplementary information:** Supplementary data are available online.

## 1   Introduction

Constraint-based approaches have been employed to simulate and predict phenotypes based on genome-scale metabolic models (GEMs) (Price *et al.*, 2003). While already useful for predicting a wide range of phenotypes, the predictive performance of GEMs has been further improved by integrating protein constraints, such as enzyme catalytic rates and the allocation of enzyme abundances across reactions (Adadi *et al.*, 2012; Sánchez *et al.*, 2017). These protein-constrained GEMs (pcGEMs) have been used to predict complex phenotypes, such as the overflow metabolism, in which fermentation predominates over respiration when microorganisms grow in high sugar concentrations (Basan *et al.*, 2015; Sánchez *et al.*, 2017), and diauxic growth, when multiple carbon sources are available and the microbial growth presents two or more growth phases (Beg *et al.*, 2007). The models also allow for the incorporation of proteomics data, and thus provide a framework for multi-omics data analysis and integration (Bekiaris and Klamt, 2020; Sánchez *et al.*, 2017).

The parameters included in pcGEMs are: (i) the enzyme turnover numbers, $k_{cat}$, a first-order rate constant with the unit of s⁻¹ that describes the

limiting rate of reactions catalysed by enzymes when these are fully occupied at their saturation point; and (ii) enzyme abundances (in mmol/gDW), obtained from quantitative proteomics experiments. Values of $k_{cat}$ can be measured from biochemical assays or be estimated from computational methods based on constraint-based and data-driven approaches (Ferreira *et al.*, 2022), while enzyme abundances are obtained from absolute proteomics measurements. More specifically, they are obtained from peptide intensity-based quantification or spectral counting (Lindemann *et al.*, 2017). However, proteomics experiments for absolute quantification are still difficult to perform, given the challenges put forward by the diversity of physicochemical properties of protein (Otto *et al.*, 2014), lack of standards and problems in reproducibility (Calderón-Celis *et al.*, 2018), and overall inaccessibility given the high costs of equipment and supplies (Swiatly *et al.*, 2018).

Computational methods have also been developed to predict protein abundance, mostly based on data-driven models. These models often explore the central dogma of molecular biology, by assessing the relationship between transcription and protein biosynthesis. Notable approaches to estimate protein abundance include the joint learning approach devised by Li *et al.* (2019), where an ensemble model was constructed by combining different supervised learning algorithms, outperforming competing

approaches in the NCI-CPTAC DREAM Proteogenomics Challenge. Another approach, developed by Terai and Asai (2020), uses features such as the accessibility around the Shine-Dalgarno sequence, minimum free energy of the mRNA molecule, Viterbi score, and inside-outside score. Further, Ferreira *et al.* (2021) explored codon usage bias information to train an AdaBoost regression model, achieving higher correlations than previous approaches without the usage of transcriptomics data.

Aside from machine learning models, constraint-based approaches have also been used to predict protein abundance. Using approaches such as MOMENT (Adadi *et al.*, 2012) or GECKO (Sánchez *et al.*, 2017), it is possible to calculate the optimal concentration of enzymes necessary to carry the provided flux with the provided catalytic rate, given the relationship:

$$v_j \leq k_{cat}^{ij} \cdot [E_i] \tag{1}$$

where $v_j$ is the metabolic flux of reaction $j$, $[E_i]$ is the concentration of an enzyme $i$, and $k_{cat}^{ij}$ is the catalytic rate of an enzyme i catalyzing a reaction $j$. This allows for deriving $k_{cat}^{ij}$ values given the other two are available. This relationship was explored by Heckmann *et al.* (2018) by using pcGEMs to predict enzyme concentrations given catalytic rates predicted computationally, achieving a 43% lower root mean squared error.

Assuming that pcGEMs that integrate proteomics data predict flux distributions that reflect the corresponding metabolic state, we ask whether the reverse operation could be employed to predict proteomics data that match a given physiological state. Moreover, as cells are exposed to stresses or changing environmental conditions, the optimal growth state is disturbed, leading to a suboptimal growth state in which gene expression, regulatory pathways and metabolic flux are changed to aid the cell in adjusting to this new physiological condition (Lahtvee *et al.*, 2016). Despite the aforementioned advances in predicting protein abundances, the problem of predicting enzyme allocation under suboptimal growth conditions remains largely unexplored. Here we propose PARROT (Figure 1), for **P**rotein allocation **A**djustment fo**R** st**R**ess c**O**ndi**T**ions, a family of constraint-based approaches for prediction of protein abundances for suboptimal conditions using protein abundances measured in a reference, optimal state. Our proposed approach is inspired by Minimization of Metabolic Adjustment (MOMA) (Segrè *et al.*, 2002), which minimizes the distance between a reference state and a gene knock-out state while ensuring cell survival in the later. We show that PARROT predicted enzyme concentrations in very good agreement with experimental data and outperformed competing methods for minimizing flux distributions. Therefore, PARROT can be used to parameterize pcGEMs for unseen, suboptimal conditions from which metabolic phenotypes can further be analysed.

## 2 Methods

### 2.1 The principle of minimizing the change in enzyme usage between a suboptimal and reference state

To find the enzyme distribution vector that matches the enzyme usage of a cell growing in suboptimal growth conditions, we propose PARROT that minimizes the distance between a reference enzyme allocation $\mathbf{E_{ref}}$ and a suboptimal growth enzyme allocation $\mathbf{E_s}$ (Figure 1). This is consistent with observations that micro-organisms minimize expenditures to perform a growth and associated flux state (Goelzer *et al.*, 2015). We define and compare four different objectives to model the distance between enzyme allocations in suboptimal and reference states: (i) the Manhattan distance; (ii) the Euclidean distance; (iii) the weighted sum of the Manhattan distance between enzyme allocations and the Manhattan distance between

flux distributions; (iv) the weighted sum of the Euclidean distance between enzyme allocations and the Euclidean distance between flux distributions. The first can be formulated as a linear optimization problem (LP1), specified as follows:

$$\min \left\| \frac{\mathbf{E_{ref}}}{E_{ref}^{tot}} - \frac{\mathbf{E_s}}{E_s^{tot}} \right\|_1 \tag{2}$$

$$\text{s.t. } \mathbf{Nv} = \mathbf{0} \tag{3}$$

$$v_{s,min} \leq v_s \leq v_{s,max} \tag{4}$$

$$v_s \leq k_{cat} \cdot [E_s] \tag{5}$$

$$\sum E_s = E_s^{tot} \tag{6}$$

$$v_{bio} = \mu \tag{7}$$

where $E_{ref}^{tot}$ and $E_s^{tot}$ represent the total enzyme usage in the model for the reference and suboptimal states, respectively; $\mathbf{N}$ is the stoichiometric matrix; $\mathbf{v}$ is the flux distribution vector; $v_{bio}$ is the flux through the biomass pseudo-reaction; and $\mu$ is the specific growth rate, determined from measurements. The other objectives are captured by the following:

$$\text{QP1: } \left\| \frac{\mathbf{E_{ref}}}{E_{ref}^{tot}} - \frac{\mathbf{E_s}}{E_s^{tot}} \right\|_2 \tag{8}$$

$$\text{LP2: } \left\| \frac{\mathbf{E_{ref}}}{E_{ref}^{tot}} - \frac{\mathbf{E_s}}{E_s^{tot}} \right\|_1 + \lambda \|\mathbf{v_{ref}} - \mathbf{v_s}\|_1 \tag{9}$$
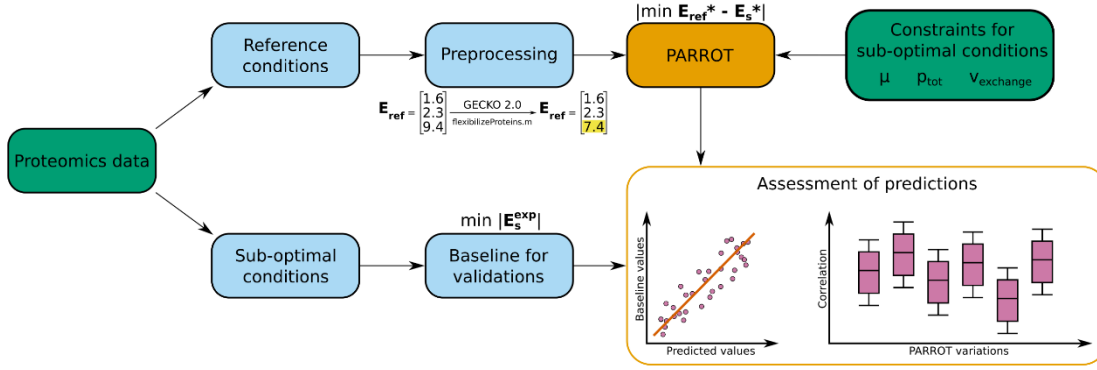
$$\text{QP2: } \left\| \frac{\mathbf{E_{ref}}}{E_{ref}^{tot}} - \frac{\mathbf{E_s}}{E_s^{tot}} \right\|_2 + \lambda \|\mathbf{v_{ref}} - \mathbf{v_s}\|_2 \tag{10}$$

where the parameter $\lambda$ is a weighting factor chosen by inspecting the difference between the norms of enzyme allocation and the flux distributions. We solved the corresponding problems under the same constraints as in Eq. 2. We implemented and solved the problems in MATLAB (The Math-Works Inc., Natick, Massachusetts) using the COBRA Toolbox (Heirendt *et al.*, 2019) and the Gurobi solver v9.1.1 (Gurobi Optimization, 2020). The implementation of PARROT can be found in the GitHub repository: https://github.com/mauricioamf/PARROT

### 2.2 Experimental data and simulation constraints

To test the variants of the proposed approach, PARROT, we used the pcGEMs of *Saccharomyces cerevisiae*, ecYeast8 (Lu *et al.*, 2019), and *Escherichia coli*, eciML1515 (Domenzain *et al.*, 2022). We employed quantitative proteomics measurements for both species performed in a number of growth conditions, ranging from optimal growth in standard physiological conditions to stress conditions, alternative nutrient usage and chemostat cultivation.

For *S. cerevisiae*, we used the protein measurements from Chen and Nielsen (2021) for 19 different growth conditions, which were collected from four studies (Lahtvee *et al.*, 2017; Yu *et al.*, 2020; Di Bartolomeo *et al.*, 2020; Yu *et al.*, 2021). These included proteomics measurements in yeast growing in ethanol, osmolarity, and high temperature stresses (Lahtvee *et al.*, 2017); yeast growing in chemostats with reducing nitrogen availability (Yu *et al.*, 2020); and yeast growing in chemostats limited by the nitrogen source in increasing dilution rates and in chemostats with alternative nitrogen sources (Yu *et al.*, 2021). We also made use the measurements of nutrient uptake rates, growth rates and protein content from these studies to constrain the batch model, which does not consider protein measurements and rely on the protein pool constraint.

**Fig. 1. Workflow of PARROT to predict enzyme usage for suboptimal growth conditions.** PARROT uses experimental proteomics data from an optimal growth condition as a reference point, and experimental physiological parameters from a suboptimal growth condition in a protein-constrained model. The proteomics data from the reference state is preprocessed by integrating the data in a pcGEM using the GECKO Toolbox 2 and flexibilizing its values. The proteomics data from the suboptimal state is used to generate a baseline, which is in turn used for comparison with predictions from the PARROT variants.

For *E. coli*, we used the proteomics data for 20 different growth conditions collected in Davidi *et al.* (2016) from three different studies (Valgepea *et al.*, 2013; Peebo *et al.*, 2015; Schmidt *et al.*, 2016). These include batch cultivations of *E. coli* growing with different carbon sources and a glucose-limited chemostat culture, with dilution rates ranging from $0.12\ h^{-1}$ to $0.5\ h^{-1}$ performed by Schmidt *et al.* (2016), a second chemostat limited by glucose at dilution rates ranging from $0.11\ h^{-1}$ to $0.49\ h^{-1}$ (Valgepea *et al.*, 2013), and a third chemostat limited by glucose at dilutions rates ranging from $0.21\ h^{-1}$ to $0.51\ h^{-1}$ (Peebo *et al.*, 2015). Similar to *S. cerevisiae*, the batch model was constrained with the nutrient uptake rates, growth rates and protein content measured in the studies where the protein measurements were taken. For both species, we excluded the conditions that did not have measured uptake rates, growth rates, or protein content. In addition, we excluded the temperature stress conditions from Lahtvee *et al.* (2017), as temperature can severely impact the function of enzymes (Li *et al.*, 2021), and temperature stress responses entail changes beyond metabolic flux redistribution (Lahtvee *et al.*, 2016).

## 2.3 Pre-processing of protein measurements for the reference state

From the protein measurements obtained from Davidi *et al.* (2016) and Chen and Nielsen (2021) we separated the measurements according to each experiment performed in the original studies. From each experiment, we selected the control sample to represent the reference state in our approach PARROT. We corrected the protein measurements for the reference state measurements by integrating the values into the pcGEMs ecYeast8 and eciML1515 for *S. cerevisiae* and *E. coli*, respectively, using the GECKO Toolbox 2 (Domenzain *et al.*, 2022). The GECKO Toolbox 2 identifies the enzyme usage values that most limit growth and flexibilizes the values to prevent over-constraining the model. We then used for the $\mathbf{E_{ref}}$ vector of each experiment the values for flexibilized proteins along with values for proteins that were unchanged.

## 2.4 Assessment of predicted enzyme usage distributions

The protein measurements, $\mathbf{E_s^{exp}}$, for the suboptimal growth conditions obtained from Davidi *et al.* (2016) and Chen and Nielsen (2021) were not used directly in simulations. These experimental measurements were instead employed to calculate a baseline to which predictions of $\mathbf{E_s}$ were compared. Assuming that simulations performed with pcGEMs use only the optimal concentration of enzymes necessary to carry a given metabolic flux, the model-allocated protein usage should underestimate the *in vivo* enzyme concentrations. To allow for a fair comparison, we devised a baseline by integrating the experimental proteomics measurements of each experiment into the pcGEMs using the GECKO Toolbox 2. Then, we minimized the total enzyme allocation given the following optimization problem:

$$\min \left\| \mathbf{E_s^{exp}} \right\|_1 \tag{11}$$

$$\text{s.t. } \mathbf{Nv} = \mathbf{0} \tag{12}$$

$$\mathbf{v_{s,min}} \le \mathbf{v_s} \le \mathbf{v_{s,max}} \tag{13}$$

$$v_{s,j} \le k_{cat}^{ij} \cdot \left[ E_s^{exp,i} \right] \tag{14}$$

$$\sum E_s^{exp} = E_s^{exp,tot} \tag{15}$$

$$v_{bio} = \mu \tag{16}$$

The resulting enzyme usage distribution, $\mathbf{E_s^{exp}}$, was then defined as the baseline for each sample of each proteomics experiment. We compared the predicted $\mathbf{E_s}$ values from the four variants of PARROT to $\mathbf{E_s^{exp}}$ by calculating the Pearson and Spearman correlations of each sample. Further, we calculated the root-median square error (RMdSE) to measure the difference between predicted and baseline values. For assessing both correlations and the RMdSE, we log10-transformed the values for the predictions and the baseline.

We also performed a robustness analysis to check the effect of using the minimization of the second norm in constructing a baseline. In addition, we compared the predictions of our approaches to those obtained using an extension of parsimonious enzyme usage FBA (pFBA) (Lewis *et al.*, 2010) to consider enzyme constraints. To this end, for each sample of each experiment, we defined the optimization problem as:

$$\min \sum_{j=1}^{m} v_{j,s,irrev} \tag{17}$$

$$\text{s.t. } \mathbf{N_{s,irrev}} \cdot \mathbf{v_{s,irrev}} = \mathbf{0} \tag{18}$$

$$0 \leq \mathbf{v_{s,irrev}} \leq \mathbf{v_{s,irrev,max}} \qquad (19)$$

$$v_{s,irrev,j} \leq k_{cat}^{ij} \cdot \left[E_{s,i}\right] \qquad (20)$$

$$\sum E_s = E_s^{tot} \qquad (21)$$

$$v_{bio} = \mu \qquad (22)$$

where $v_{j,s,irrev}$ corresponds to the flux distribution of an irreversible model in a non-optimal growth condition. We also assessed a modified version of pFBA with enzyme constraints with the following objective:

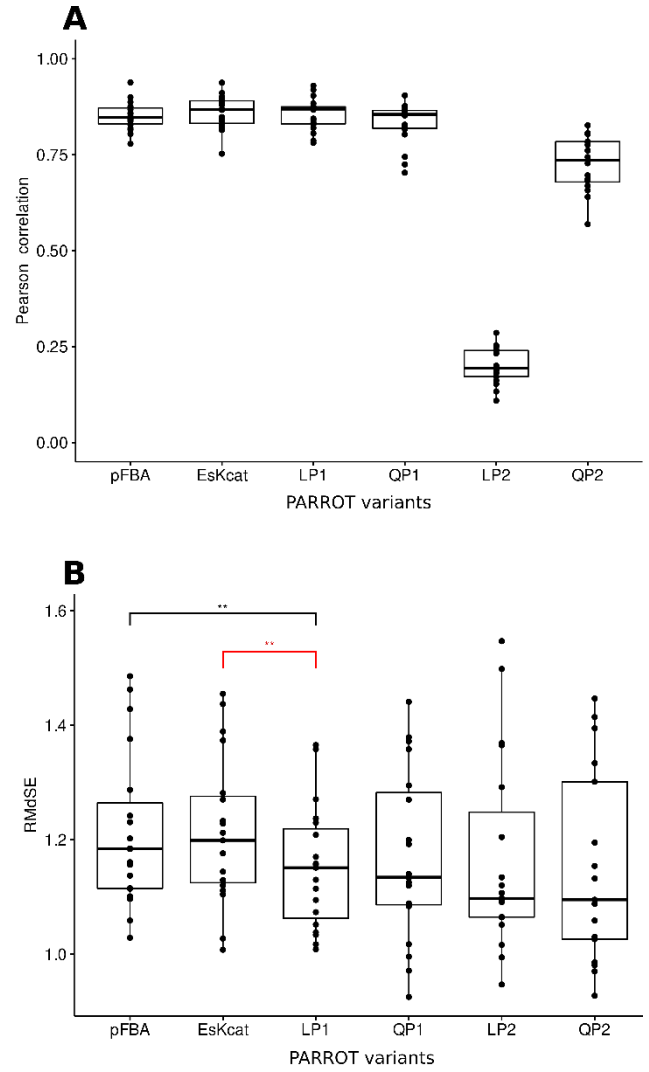$$\min \sum E_s \cdot k_{cat} . \qquad (23)$$

For pFBA and the modified implementation, we applied the same constraints on nutrient uptake rates and growth rates as for the four approaches assessed previously, and calculated the Pearson and Spearman correlations, and the RMdSE. We calculated the correlation values and RMdSE for all assessed optimization problems and compared them to the predictions of pFBA using a Pairwise Wilcoxon rank sum test with Bonferroni correction.

## 3    Results

### 3.1 The minimization of the Manhattan distance captures protein allocation changes in yeast

We used PARROT to predict the enzyme usage distribution for 19 growth conditions under constraints provided by experimental data. First, we built a baseline for comparison with predictions from PARROT (Figure 1). To this end, we integrated the experimental proteomics measurements obtained from Lahtvee *et al.* (2017), Yu *et al.* (2020), Di Bartolomeo *et al.* (2020), and Yu *et al.* (2021) (Table S1) in the ecYeast8 model and minimized the enzyme allocation (Methods). The resulting allocation of enzymes $\mathbf{E_s^{exp}}$ included 286 to 336 enzymes with abundance in all considered conditions. For the reference condition, we used the experimental proteomics measurements from optimal (control) growth conditions in the respective four groups of experiments, after flexibilization following GECKO 2.0 (see Methods) (Table S1). The number of enzymes contained in $\mathbf{E_{ref}}$ ranged from 533 to 744.

With the resulting enzyme allocation at the reference and the baseline of the suboptimal condition, $\mathbf{E_{ref}}$ and $\mathbf{E_s^{exp}}$, we used the four variants of PARROT to predict the enzyme allocation, $\mathbf{E_s}$, for the suboptimal condition. The number of enzymes contained in the predicted $\mathbf{E_s}$ ranged from 114 to 267 over the considered experiments. When comparing the median of the calculated Pearson correlations between the baseline and predicted enzyme allocation correlations, we found the minimization of the Manhattan distance in PARROT achieved the highest median correlations (Figure 2a, Figure S1). We also evaluated the RMdSE between predictions and the baselines, and observed that the minimization of the Manhattan distance also resulted in the smallest median error compared to the other optimization problems. Further, PARROT exhibited significantly smaller RMdSE than pFBA (p-value = 0.002, pairwise Wilcoxon rank sum test) and its modified implementation (p-value = 0.000965, pairwise Wilcoxon rank sum test) (Figure 2b) Taken together, the results demonstrated that the minimization of the Manhattan distance in PARROT showed the best prediction performance based on the data from *S. cerevisiae*.
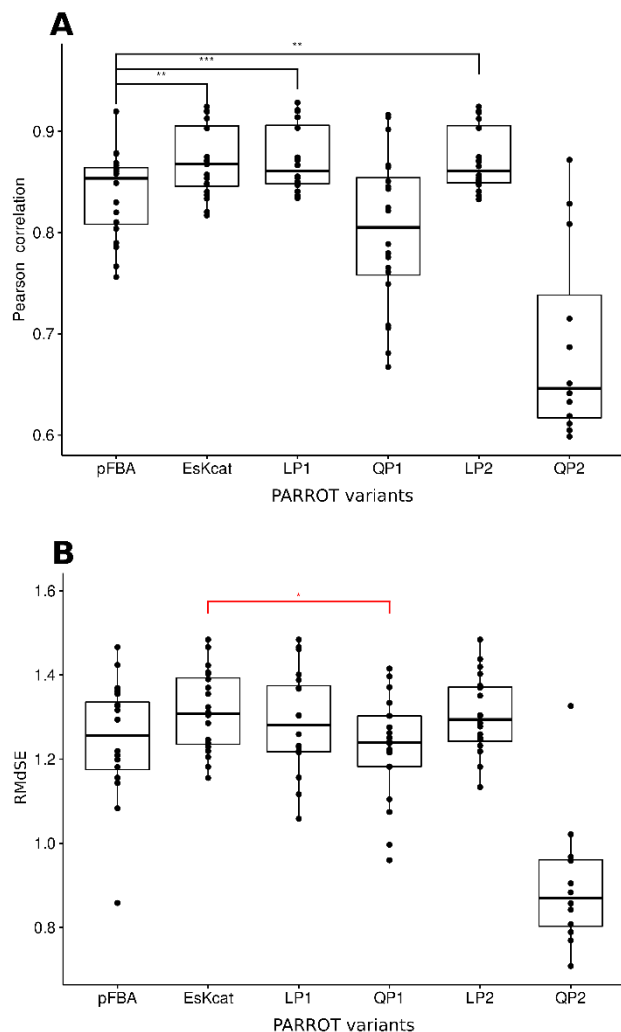
**Fig. 2. Comparative performance analysis of PARROT with proteomics data from *S. cerevisiae*.** All values were log10-transformed prior to comparisons. **a**. Pearson correlation calculated between predicted enzyme distribution and the baseline obtained from minimizing the first norm of the experimental enzyme usage distribution. The four variants of PARROT are denoted as LP1 (Manhattan distance of enzyme distributions), LP2 (weighted Manhattan distance, considering flux and enzyme distributions), QP1 (Euclidean distance of enzyme distributions), and QP2 (weighted Euclidean distance of flux and enzyme distributions). The performance of PARROT was compared to pFBA and its modified version EsKcat (first norm of enzyme usage), see Methods. **b**. Assessment of model performance based on the root median squared error (RMdSE). A pairwise Wilcoxon rank sum assesses the statistical significance: ** p-value < 0.002. Black significance bar indicates comparisons to pFBA. Red significance bar indicates comparison to EsKcat.

### 3.2 Different variants of PARROT outperformed contending methods for *E. coli*

To verify if the conclusions from PARROT hold in another unicellular model organism, we applied it to predict enzyme allocation $\mathbf{E_s}$ in suboptimal conditions for *E. coli* given constraints provided by growth experiments. As in the case of *S. cerevisiae*, we built a baseline for comparison

with the predictions obtained from PARROT by integrating the experimental proteomics measurements from Valgepea *et al.* (2013), Peebo *et al.* (2015) and (Schmidt *et al.*, 2016) (Table S2) in the eciML1515 model, and minimized the total enzyme allocation (see Methods). The resulting



**Fig. 3. Comparative performance analysis of PARROT with proteomics data from E. coli.** All values were log10-transformed prior to comparisons. **a.** Pearson correlation calculated between predicted enzyme usage distribution and the baseline obtained from minimizing the first norm of the experimental enzyme usage distribution. A pairwise Wilcoxon rank sum assesses the statistical significance: *** p-value < 1e-4, ** p-value < 0.002. **b.** Assessment of model performance based on the RMdSE in E. coli. All values were log10-transformed prior to comparisons. A pairwise Wilcoxon rank sum assesses the statistical significance: * p-value < 0.05. Black significance bar indicates comparisons to pFBA. Red significance bar indicates comparison to EsKcat.

$E_s^{exp}$ included protein allocation for 164 to 176 enzymes. Further, as reference condition we considered the control samples or the chemostat measurements with the smallest dilution rate (Table S2). The number of enzymes contained in $E_{ref}$ ranged from 152 to 188.

The prediction of $E_s$ distributions and their assessment were similar to *S. cerevisiae*, with the number of predicted values ranging from 42 to 106.

After performing a comparison of Pearson and Spearman correlations between variants of PARROT, pFBA and its modified implementation, we observed that different variants outperformed pFBA. The minimization of the Manhattan distance had the highest median correlations, showing significant difference to pFBA (p-value = 1.34e-4 and 0.004, for Pearson and Spearman correlations respectively, pairwise Wilcoxon rank sum test). Another variant with a significant difference to pFBA was the minimization of the weighted sum of Manhattan distance of enzyme usage and Manhattan distance of flux distributions (Figure 3a, Figure S2). Regarding the RMdSE, the minimization of the Euclidean distance achieved the lowest errors, outperforming both pFBA, its modified implementation and the minimization of the Manhattan distance (Figure 3b), with a significant difference when compared to the modified version of pFBA (p-value = 0.009, pairwise Wilcoxon rank sum test). These findings demonstrated that PARROT is applicable with data from different microorganisms without decreasing its performance.

### 3.3 Robustness analysis shows the consistency of prediction from PARROT

To further evaluate the predictions made by PARROT, we investigated how the usage of a baseline constructed by minimizing the second norm of the vector $E_s^{exp}$ impacts the comparisons. To this end, we repeated all comparisons as performed for a baseline constructed by minimizing the first norm, using the predicted $E_s$ obtained by the PARROT variants. Importantly, the results were consistent between the two baseline approaches. For *S. cerevisiae*, the minimization of the Manhattan distance maintained its performance by achieving higher mean correlations than pFBA and its modified implementation (Figures S3-S4). When considering the Pearson correlation, the minimization of the Euclidean distance also achieved higher correlation values than the two contending methods. For the RMdSE, the minimization of the Manhattan distance was still the variant with the lowest errors, with a significant difference to pFBA (p-value < 0.0002, pairwise Wilcoxon rank sum test) (Figure S5).

The comparisons performed using predictions obtained for *E. coli* were also consistent with different variants of PARROT that outperformed pFBA. Considering the Pearson and Spearman correlations, the minimization of the Manhattan distance also had the highest median correlations and were significantly different to pFBA (p-value = 8.36e-4 and 1.91e-6, for Pearson and Spearman correlations respectively, pairwise Wilcoxon rank sum test). Likewise, the modified implementation of pFBA, and the minimization of the weighted sum of Manhattan distance of enzyme usage and Manhattan distance of flux distributions also had a significant difference to pFBA (Figures S6-S7). The comparison of RMdSE values were also consistent, as the minimization of the Euclidean distance achieved the lowest errors among variants and pFBA (Figure S8). Altogether, these results highlight the robustness of estimations of $E_s$ obtained from PARROT.

## 4    Discussion

Here we proposed a family of constraint-based approaches, termed PARROT, that address the problem of predicting reallocation of protein abundance from an optimal condition to a suboptimal condition. PARROT is based on the principle that organisms tend to minimally adjust cellular physiology between growth conditions to make effective use of resources (Goelzer *et al.*, 2015). The predictions of enzyme allocation generated by PARROT rely on quantitative proteomics data for a reference condition. The resulting optimization problems constructed are thus similar to

MOMA, which depends on a model representing a wild-type strain to predict a minimally adjusted flux distribution for a mutant strain. By comparing the predictions to a baseline constructed with experimental proteomics measurements for suboptimal conditions, we found that PARROT predicted protein abundances with very good agreement with the baseline. In addition, we demonstrated that these predictions were consistent and robust to how the baseline is constructed. The performance of PARROT also holds for two model organisms, *S. cerevisiae* and *E. coli*, highlighting the general application of the principle of minimal protein adjustment on which the predictions are based.

From the different variants of PARROT, the minimization of the Manhattan distance (LP1) was the most promising, with QP1 as the second-best contender. However, LP2 and QP2 – that combine fluxes and enzyme predictions for both organisms – resulted in highly inconsistent performance between yeast and *E. coli*. This suggests that the joint minimization of fluxes and enzymes is not a principle of flux redistribution and the principle is guided by minimization of resource redistribution, as captured by LP1 and QP1. Altogether, we demonstrated that minimizing the readjustment of enzyme resource allocation is one principle underpinning microbial adjustment to a suboptimal condition. Thus, PARROT may allow for study and engineering of microbial cell factories, as these are often under suboptimal growth conditions in industrial settings (Deparis *et al.*, 2017).

The baseline approach devised to assess the predictions allows for a fair comparison between the predicted enzyme usage distribution and the experimental protein abundance values. In constraining the pcGEMs using the proteomics measurements, the experimental values are first readjusted to match the enzyme levels that actually carry flux in the model, since more protein is produced than actually needed by the cell (O'Brien *et al.*, 2016). This, however, implies that the predicted values are not directly comparable to experimental proteomics values, which affect the determined measures of performance. By adjusting the experimental values to levels that are compatible with what is actually employed to carry metabolic flux, we could more adequately assess the correlation with enzyme allocation predicted from the pcGEMs, albeit losing the direct correspondence to experimental data.

Despite the advantages of using a baseline, predictions of enzyme levels using Equation 1 still underestimates protein abundance, leading to a disparity between predictions and *in vivo* concentrations. This remaining portion of proteins, termed the "proteome reserve", is useful for the cell to quickly adapt to unstable environments, being an evolutionary conserved strategy (Mori *et al.*, 2017). This falls in line with the evolutionary conservation of protein stoichiometries at the pathway level as demonstrated by Lalanne *et al.* (2018). Although it is still not understood how preferred enzyme stoichiometry is determined, it was observed that the preferred range of enzyme stoichiometry follows a narrow distribution among pathways in Gram-positive and -negative bacteria, likely a result of evolutionary conservation or convergence. As suggested in the study, protein biosynthesis and consequently its usage is bound to a cost-benefit trade-off, where the optimal level of enzymes is balanced with the need for a buffer zone in case of changing environments. Similar to our approach, the works of Mori *et al.* (2017) and Lalanne *et al.* (2018) deals with proteome reallocation in a suboptimal growth condition. However, the first deals with proteome sectors, while the latter concerns with pathway-centric stoichiometries. Our approach thus differs as we consider protein reallocation for each enzyme individually.

Nevertheless, other approaches for estimating *in vivo* protein concentrations would still need to overcome the underestimating capacity of pcGEMs, especially by considering the proteome reserve. These ap-

proaches could include features such as cellular machinery beyond enzymes that participate in metabolism, or by integrating constraint-based approaches with data-driven approaches.

## References

Adadi,R. *et al.* (2012) Prediction of Microbial Growth Rate versus Biomass Yield by a Metabolic Network with Kinetic Parameters. *PLoS Comput. Biol.*, **8**, e1002575.

Di Bartolomeo,F. *et al.* (2020) Absolute yeast mitochondrial proteome quantification reveals trade-off between biosynthesis and energy generation during diauxic shift. *Proc. Natl. Acad. Sci. U. S. A.*, **117**, 7524–7535.

Basan,M. *et al.* (2015) Overflow metabolism in Escherichia coli results from efficient proteome allocation. *Nature*, **528**, 99–104.

Beg,Q.K. *et al.* (2007) Intracellular crowding defines the mode and sequence of substrate uptake by Escherichia coli and constrains its metabolic activity. *Proc. Natl. Acad. Sci. U. S. A.*, **104**, 12663–12668.

Bekiaris,P.S. and Klamt,S. (2020) Automatic construction of metabolic models with enzyme constraints. *BMC Bioinformatics*, **21**, 1–13.

Calderón-Celis,F. *et al.* (2018) Standardization approaches in absolute quantitative proteomics with mass spectrometry. *Mass Spectrom. Rev.*, **37**, 715–737.

Chen,Y. and Nielsen,J. (2021) In vitro turnover numbers do not reflect in vivo activities of yeast enzymes. *Proc. Natl. Acad. Sci. U. S. A.*, **118**, e2108391118.

Davidi,D. *et al.* (2016) Global characterization of in vivo enzyme catalytic rates and their correspondence to in vitro kcat measurements. *Proc. Natl. Acad. Sci. U. S. A.*, **113**, 3401–3406.

Deparis,Q. *et al.* (2017) Engineering tolerance to industrially relevant stress factors in yeast cell factories. *FEMS Yeast Res.*, **17**, 1–35.

Domenzain,I. *et al.* (2022) Reconstruction of a catalogue of genome-scale metabolic models with enzymatic constraints using GECKO 2.0. *Nat. Commun.*, **13**, 1–13.

Ferreira,M. *et al.* (2021) Protein Abundance Prediction Through Machine Learning Methods. *J. Mol. Biol.*, **433**, 167267.

Ferreira,M.A. de M. *et al.* (2022) Protein constraints in genome-scale metabolic models: data integration, parameter estimation, and prediction of metabolic phenotypes. *Authorea Prepr.*

Goelzer,A. *et al.* (2015) Quantitative prediction of genome-wide resource allocation in bacteria. *Metab. Eng.*, **32**, 232–243.

Gurobi Optimization,L. (2020) Gurobi Optimizer Reference Manual.

Heckmann,D. *et al.* (2018) Machine learning applied to enzyme turnover numbers reveals protein structural correlates and improves metabolic models. *Nat. Commun.*, **9**, 1–10.

Heirendt,L. *et al.* (2019) Creation and analysis of biochemical constraint-based

models using the COBRA Toolbox v.3.0. *Nat. Protoc.*, **14**, 639–702.

Lahtvee,P.-J. *et al.* (2016) Adaptation to different types of stress converge on mitochondrial metabolism. *Mol. Biol. Cell*, **27**, 2505–2514.

Lahtvee,P.J. *et al.* (2017) Absolute Quantification of Protein and mRNA Abundances Demonstrate Variability in Gene-Specific Translation Efficiency in Yeast. *Cell Syst.*, **4**, 495-504.e5.

Lalanne,J.B. *et al.* (2018) Evolutionary Convergence of Pathway-Specific Enzyme Expression Stoichiometry. *Cell*, **173**, 749-761.e38.

Lewis,N.E. *et al.* (2010) Omic data from evolved E. coli are consistent with computed optimal growth from genome-scale models. *Mol. Syst. Biol.*, **6**.

Li,G. *et al.* (2021) Bayesian genome scale modelling identifies thermal determinants of yeast metabolism. *Nat. Commun.*, **12**, 1–12.

Li,H. *et al.* (2019) Joint learning improves protein abundance prediction in cancers. *BMC Biol.*, **17**, 1–14.

Lindemann,C. *et al.* (2017) Strategies in relative and absolute quantitative mass spectrometry based proteomics. *Biol. Chem.*, **398**, 687–699.

Lu,H. *et al.* (2019) A consensus S. cerevisiae metabolic model Yeast8 and its ecosystem for comprehensively probing cellular metabolism. *Nat. Commun.*, **10**.

Mori,M. *et al.* (2017) Quantifying the benefit of a proteome reserve in fluctuating environments. *Nat. Commun.*, **8**, 1–8.

O'Brien,E.J. *et al.* (2016) Quantification and Classification of E. coli Proteome Utilization and Unused Protein Costs across Environments. *PLOS Comput. Biol.*, **12**, e1004998.

Otto,A. *et al.* (2014) Quantitative proteomics in the field of microbiology. *Proteomics*, **14**, 547–565.

Peebo,K. *et al.* (2015) Proteome reallocation in Escherichia coli with increasing specific growth rate. *Mol. Biosyst.*, **11**, 1184–1193.

Price,N.D. *et al.* (2003) Genome-scale microbial in silico models: The constraints-based approach. *Trends Biotechnol.*, **21**, 162–169.

Sánchez,B.J. *et al.* (2017) Improving the phenotype predictions of a yeast genome-scale metabolic model by incorporating enzymatic constraints. *Mol. Syst. Biol.*, **13**, 935.

Schmidt,A. *et al.* (2016) The quantitative and condition-dependent Escherichia coli proteome. *Nat. Biotechnol.*, **34**, 104–110.

Segrè,D. *et al.* (2002) Analysis of optimality in natural and perturbed metabolic networks. *Proc. Natl. Acad. Sci. U. S. A.*, **99**, 15112–15117.

Swiatly,A. *et al.* (2018) Mass spectrometry-based proteomics techniques and their application in ovarian cancer research. *J. Ovarian Res.*, **11**, 1–13.

Terai,G. and Asai,K. (2020) Improving the prediction accuracy of protein abundance in Escherichia coli using mRNA accessibility. *Nucleic Acids Res.*, **48**.

Valgepea,K. *et al.* (2013) Escherichia coli achieves faster growth by increasing catalytic and translation rates of proteins. *Mol. Biosyst.*, **9**, 2344–2358.

Yu,R. *et al.* (2020) Nitrogen limitation reveals large reserves in metabolic and translational capacities of yeast. *Nat. Commun.*, **11**, 1–12.

Yu,R. *et al.* (2021) Quantifying absolute gene expression profiles reveals distinct regulation of central carbon metabolism genes in yeast. *Elife*, **10**.