

A genome-scale metabolic network model and machine learning predict amino acid concentrations in Chinese Hamster Ovary cell cultures

Song-Min Schinn¹, Carly Morrison², Wei Wei², Lin Zhang², Nathan E. Lewis^{1,3,4}

1. Department of Pediatrics, University of California, San Diego

2. Pfizer, Biotherapeutics Pharmaceutical Sciences, Andover, MA

3. Department of Bioengineering, University of California, San Diego

4. Novo Nordisk Foundation Center for Biosustainability at UC San Diego

Corresponding author: Dr. Nathan E. Lewis, nlewisres@ucsd.edu

Abstract

The control of nutrient availability is critical to large-scale manufacturing of biotherapeutics. However, the quantification of proteinogenic amino acids is time-consuming and thus is difficult to implement for real-time *in situ* bioprocess control. Genome-scale metabolic models describe the metabolic conversion from media nutrients to proliferation and recombinant protein production, and therefore are a promising platform for *in silico* monitoring and prediction of amino acid concentrations. This potential has not been realized due to unresolved challenges: (1) the models assume an optimal and highly efficient metabolism, and therefore tend to underestimate amino acid consumption, and (2) the models assume a steady state, and therefore have a short forecast range. We address these challenges by integrating machine learning with the metabolic models. Through this we demonstrate accurate and time-course dependent prediction of individual amino acid concentration in culture medium throughout the production process. Thus, these models can be deployed to control nutrient feeding to avoid premature nutrient depletion or provide early predictions of failed bioreactor runs.

Keywords: bioprocess, Chinese Hamster Ovary, metabolism, Systems Biology, Metabolic Network Modeling

Short Communication

Chinese Hamster Ovary (CHO) cells are widely used to manufacture complex biotherapeutic molecules at large scales. Industrial bioprocesses ensure high product yield and quality by maintaining favorable growth conditions in cell culture environments, which requires careful monitoring and control of nutrient availability. Chemically-defined serum-free media can contain dozens or >100 components (Ritacco et al., 2018), but key nutrients include proteinogenic amino acids, which are direct substrates and regulators (Duarte et al., 2014; Fomina & Yadlin et al., 2014) of proliferation and protein synthesis. Unfortunately, conventional methods for amino acid quantification based on liquid chromatography and mass spectrometry are time-consuming and difficult to use for decision making and control of cell culture. Alternate spectroscopic approaches have been sensitive to a limited number of amino acid species (Bhatia et al., 2018). Here we present a computational method to forecast time-course amino acid concentrations from routine bioprocess measurements, facilitating a timely and anticipatory control of the bioprocess (Fig. 1).

At the foundation of our method is a genome-scale metabolic network model, which accounts for the complex conversion from media nutrients to biomass and recombinant protein production. Such models have been increasingly utilized for CHO cells (Hefzi et al., 2016; Calmels et al., 2019; Huang & Yoon, 2020) and bioprocess applications (Sommeregger et al., 2017), such as predicting clonal performances (Popp et al., 2016), identifying metabolic bottlenecks (Zhuangrong & Seongkyu, 2020), and optimizing media formulation (Fouladiha et al., 2020; Traustason et al., 2019). Metabolic network models can also estimate amino acid

uptake rates necessary to experimentally support observed proliferation and productivity (Chen et al., 2019). However, several challenges have limited their practical application.

First, metabolic network models are typically highly complex but under-constrained, and therefore are easy to overfit. This is mitigated by training the model on a variety of bioprocess conditions and metabolic phenotypes. Second, metabolic network models assume that cells operate at some metabolic optimum, and thus tend to describe an idealized metabolism specifically fit to the assumed objective, e.g., biomass production (Feist & Palsson, 2010; Szeliöva et al., 2020), minimization of redox (Savinell & Palsson, 1992). Third, for the present purpose, these models need to predict amino acid consumption fluxes, typically on the order of $10^{-3} \text{ mmol} \cdot \text{g}_{\text{DW}}^{-1} \cdot \text{hr}^{-1}$ (see Methods), from input data that are multiple magnitudes larger, such as growth rate and glucose consumption (10^{-1} to $10^{-2} \text{ mmol} \cdot \text{g}_{\text{DW}}^{-1} \cdot \text{hr}^{-1}$). The preceding two challenges increase prediction error. Lastly, metabolic network models assume a steady state, which reduces the range of forecast. Typically, input data from one day are used to make predictions for the same day. However, such predictions cannot be extended to multiple days or subsequent culture phases, as cross-temporal shifts in metabolism would violate the steady state assumption. In summary, model predictions of amino acid concentrations can be overfit, ideal, and near-sighted – all of which dilutes their practicality for industrial bioprocess control. Here we demonstrate that these weaknesses can be addressed in a data-driven manner by coupling a metabolic network model with machine learning.

We developed this hybrid approach on a diverse set of 10 CHO clones with different growth and productivity profiles from two different fed-batch production processes. These CHO clones were subject to different bioprocess conditions and recombinant antibody identities (see Methods), resulting in a variety of phenotypes and productivity performances (Fig. S1). For

example, several high-performing clones were exceptionally proliferative or productive, suggesting an efficient conversion from nutrients to biomass or recombinant protein product. Other clones performed these conversions at lower rates, suggesting attenuated metabolic activity or inefficient resource utilization. The CHO cells adjusted their nutrient uptake according to these various metabolic phenotypes, leading to diverse amino acid consumption patterns (Fig, S2). For example, the consumption of glucose and serine differed by several fold across conditions and time. Furthermore, different clones varied in their consumption or secretion of key metabolites such as lactate, alanine, glycine, and glutamine.

We sought to predict these diverse consumption behaviors using a tailored model of CHO metabolism (Table S1, S2). As input information, we utilized the following routinely measured industrial bioprocess data: (1) viable cell density and titer measurements, from which growth rate and specific productivity are calculated (Methods, equation 1), and (2) bioreactor concentrations of glucose, lactate, glutamate and glutamine, from which their respective consumption rates are calculated. These measurements were used as boundary conditions by constraining the fluxes of biomass production, recombinant protein synthesis and consumption of the four metabolites to observed values. Subsequently, we used Markov chain Monte Carlo sampling of metabolic fluxes (Schellenberger et al., 2011) to sample the range and magnitude of all reaction fluxes to calculate the likely uptake fluxes of the remaining 18 proteinogenic amino acids (see Methods). These predictions were applied to the CHO clones across 8 days of a 12-day production run (days 4 to 11), resulting in a total of 80 individual predictions.

Predictions from the metabolic model agreed well with experimental measurements. Prediction errors were small compared to the scale of input data (Fig. 2A), suggesting that metabolic models can describe the conversion from nutrients to biomass and recombinant

proteins. However, the model also underestimated consumption rates for almost all amino acids (Fig. 2B, x-axis), on average by about half a fold. This is likely because the model doesn't consider certain metabolic inefficiencies – e.g. futile cycles or cytotoxic byproduct synthesis (Mulukutla et al., 2017).

Notably, the predicted consumption rates correlated well with measurements for many amino acids (Fig. 2B, y-axis). Therefore, we constructed a series of linear regression models to 'correct' the metabolic model predictions, using the predicted values and growth rate as explanatory variables (Methods, equation 2). This substantially improved predictions for most amino acids (Fig. S4). As exceptions, predictions for alanine and glycine did not sufficiently improve due to their high fold change error and low correlation to experimental measurements (Fig. 2B). These amino acids are non-essential and can be synthesized from glucose cost-efficiently. Therefore, their consumption may be regulated distinctly and more independent from growth requirements, as observed previously in other organisms (M. Zampieri et al., 2019). Indeed, alanine and glycine were the only two amino acid species that were variously consumed and secreted in significant amounts (Fig. S2). In short, the investigated CHO cells seem to consume them in a 'less ideal' manner than other amino acids.

Overall, our hybrid modeling approach estimated most amino acid consumptions well at a small timescale (1 day), when the steady state assumption holds true. This assumption is not valid at larger timescales of multiple days, where nutrient consumption declines asymptotically as cellular metabolism shifts from exponential growth phase to stationary phase. We addressed this limitation by modeling the multi-phase consumption profile with an exponential decay function (Methods, equation 3; Figure S5). Specifically, we first predicted amino acid consumption rates of several early culture days as heretofore described (Fig. 1, red datapoints).

Then, these datapoints were used to fit an exponential decay function that describes the entire consumption profile, including later culture days (Fig. 1, orange line).

Our approach accurately predicted daily consumption rates for each amino acid excluding alanine and glycine (Fig. 3A). This included amino acids that are highly abundant in recombinant antibodies (e.g. serine, valine, and leucine) (Fan et al., 2015), or that complicate media formulation due to low solubility (e.g. tyrosine). We also estimated the total amounts of amino acid consumed over the 8 culture days to within 86% of experimental values (Fig. 3B). These results highlight the method's value in monitoring and forecasting the bioreactor environment.

In summary, the presented modeling workflow forecasted the entire amino acid consumption profile from early bioprocess measurements, facilitating anticipatory and *in situ* control of bioreactor nutrient availability. This was realized by a novel combination of metabolic and statistical models. A metabolic network model estimated amino acid uptake rates necessary for observed proliferation and productivity, assuming an ideally efficient metabolism and steady state conditions. Two subsequent statistical models refined these predictions by offsetting prediction errors empirically and by describing the time-course relationship of individual predictions. These statistical models can easily be adjusted and re-trained for changes in cell-lines or bioprocesses. Our efforts are part of a growing trend of synergizing metabolic network models with machine learning methods (G. Zampieri et al., 2019), and demonstrates the power of hybrid modeling for on-line control of bioprocesses.

Methods

Cell culture experiments

Two production fed batch processes were used, Fed batch 1 and Fed batch 2. Both fed batch processes used chemically defined media and feeds over the 12-day cell culture. Fed batch 1 used a glucose restricted fed batch process called HiPDOG (Gagnon et al., 2011). Glucose concentration is kept low during the initial phase of the process, Day 2-7, through intermittent addition of feed medium containing glucose at the high end of pH dead-band and then glucose was maintained above 1.5 g/L thereafter, restricting lactate production without compromising the proliferative capability of cells. In Fed batch 2 a conventional cell culture process was used where glucose was maintained above 1.5 g/L throughout the process.

For both process conditions, bioreactor vessels were inoculated at 2×10^6 viable cells/mL. The following bioprocess characteristics were quantified daily using a NOVA Flex BioProfile Analyzer (Nova Biomedical, Waltham, MA): viable cell density, average live cell diameter and concentrations of glucose, lactate, glutamate, and glutamine. Viable cell density data were converted to growth rates by following equation to be compared to model-predicted growth rates.

$$(1) \quad \text{Growth rate} = \frac{1}{vcd} \cdot \frac{\Delta vcd}{\Delta \text{time}} = \frac{1}{vcd_0} \cdot \frac{vcd_{+1} - vcd_{-1}}{\text{time}_{+1} - \text{time}_{-1}}$$

Flash-frozen cell pellets (10^6 cells) and supernatant (1 mL) were collected from bioreactor runs for each sampling day. Collected samples were sent to Metabolon (Metabolon Inc, Morrisville, NC) for metabolomics analyses. Metabolomics measurements were used as input data to the model by converting their units to model units of mmol per gram of dry weight of cell per hour.

Metabolic network modeling

We used a previously described metabolic network model that is tailored to the investigated CHO clones (Schinn et al., 2020). Experimental measurements for clone and culture day were used to constrain model reactions for biomass production, monoclonal antibody secretion and consumption of glucose, lactate, glutamate, and glutamine. Then, we computed distributions of likely amino acid consumption rates by stochastically sampling 5000 points within the model's solution space via a Markov chain Monte Carlo sampling algorithm, as described previously (Nam et al., 2012), using *optGpSampler* (Megchelenbrink et al., 2014) and COBRApy (Ebrahim et al., 2013).

Statistical methods

For each amino acid, the mean of the sample distribution was interpreted as likely consumption rates predicted by the metabolic model. These predictions were refined and extended by statistical models, as explained below. The modeling workflow is visualized in a detailed diagram (Fig. S6) and demonstrated by sample code (Supplementary Data). Specifically, the statistical models were trained and validated by randomly dividing the 80 observations into two sets, consisting of 48 and 32 observations, respectively. Quantified snapshots of the validation data throughout the analysis workflow are detailed in supplementary tables, from experimental measurement to final model prediction (Table S3, S4, S6, S10); priors and inferences derived from the training data set are also provided (Table S5, S7, S8, S9).

The first statistical model refined metabolic model predictions (equation 2). Growth rate was included as an explanatory variable as it also correlated well with consumption rates of many amino acid species (Fig. S3).

$$(2) \text{ Corrected prediction} = \gamma_0 + \gamma_1 \cdot \text{prediction} + \gamma_2 \cdot \text{growth rate}$$

The second statistical model described time-course amino acid consumption by an exponential decay function (equation 3). The coefficient β_0 represents the minimum consumption rate which the cells asymptotically approach during later stationary phase. First, regression coefficients were calculated from the training dataset to be used as priors and constraints for nonlinear optimization. Specifically, the mean, minimum and maximum values of these training coefficients were used as initial guess values, lower bounds, and upper bounds, respectively. Then, regression coefficients were fitted to minimize two values: (1) the difference between outputs of equation 2 and equation 3 for early culture days, (2) the difference between fitted β_0 and previously observed asymptotic values.

$$(3) \text{ Consumption rate} = \beta_1 \cdot \exp\left(\frac{-\text{time}}{\beta_2}\right) \cdot (\text{time} + \beta_3)^3 + \beta_0$$

These analyses were carried out and visualized using COBRA Toolbox 2.0 (Schellenberger et al., 2011) in MATLAB R2018b (MathWorks; Natick, Massachusetts, USA)

References

- Bhatia, H., Mehdizadeh, H., Drapeau, D., & Yoon, S. (2018). In-line monitoring of amino acids in mammalian cell cultures using raman spectroscopy and multivariate chemometrics models. *Engineering in Life Sciences*, 18(1), 55–61.
<https://doi.org/10.1002/elsc.201700084>
- Calmels, C., McCann, A., Malphettes, L., & Andersen, M. R. (2019). Application of a curated genome-scale metabolic model of CHO DG44 to an industrial fed-batch process. *Metabolic Engineering*, 51, 9–19. <https://doi.org/10.1016/j.ymben.2018.09.009>
- Chen, Y., McConnell, B. O., Gayatri Dhara, V., Mukesh Naik, H., Li, C.-T., Antoniewicz, M. R., & Betenbaugh, M. J. (2019). An unconventional uptake rate objective function approach enhances applicability of genome-scale models for mammalian cells. *NPJ Systems Biology and Applications*, 5. <https://doi.org/10.1038/s41540-019-0103-6>
- Duarte, T. M., Carinhas, N., Barreiro, L. C., Carrondo, M. J. T., Alves, P. M., & Teixeira, A. P. (2014). Metabolic responses of CHO cells to limitation of key amino acids. *Biotechnology and Bioengineering*, 111(10), 2095–2106.
<https://doi.org/10.1002/bit.25266>
- Ebrahim, A., Lerman, J. A., Palsson, B. O., & Hyduke, D. R. (2013). COBRApy: CONstraints-Based Reconstruction and Analysis for Python. *BMC Systems Biology*, 7(1), 74.
<https://doi.org/10.1186/1752-0509-7-74>
- Fan, Y., Val, I. J. D., Müller, C., Sen, J. W., Rasmussen, S. K., Kontoravdi, C., Weilguny, D., & Andersen, M. R. (2015). Amino acid and glucose metabolism in fed-batch CHO cell culture affects antibody production and glycosylation. *Biotechnology and Bioengineering*, 112(3), 521–535. <https://doi.org/10.1002/bit.25450>

- 1 Feist, A. M., & Palsson, B. O. (2010). The biomass objective function. *Current Opinion in*
2 *Microbiology*, 13(3), 344–349. <https://doi.org/10.1016/j.mib.2010.03.003>
- 3 Fomina□Yadlin, D., Gosink, J. J., McCoy, R., Follstad, B., Morris, A., Russell, C. B., &
4 McGrew, J. T. (2014). Cellular responses to individual amino-acid depletion in antibody-
5 expressing and parental CHO cell lines. *Biotechnology and Bioengineering*, 111(5), 965–
6 979. <https://doi.org/10.1002/bit.25155>
- 7 Fouladiha, H., Marashi, S.-A., Torkashvand, F., Mahboudi, F., Lewis, N. E., & Vaziri, B. (2020).
8 A metabolic network-based approach for developing feeding strategies for CHO cells to
9 increase monoclonal antibody production. *BioRxiv*, 751347.
10 <https://doi.org/10.1101/751347>
- 11 Gagnon, M., Hiller, G., Luan, Y.-T., Kittredge, A., DeFelice, J., & Drapeau, D. (2011). High-
12 End pH-controlled delivery of glucose effectively suppresses lactate accumulation in
13 CHO Fed-batch cultures. *Biotechnology and Bioengineering*, 108(6), 1328–1337.
14 <https://doi.org/10.1002/bit.23072>
- 15 Hefzi, H., Ang, K. S., Hanscho, M., Bordbar, A., Ruckerbauer, D., Lakshmanan, M., Orellana, C.
16 A., Baycin-Hizal, D., Huang, Y., Ley, D., Martinez, V. S., Kyriakopoulos, S., Jiménez, N.
17 E., Zielinski, D. C., Quek, L.-E., Wulff, T., Arnsdorf, J., Li, S., Lee, J. S., ... Lewis, N. E.
18 (2016). A Consensus Genome-scale Reconstruction of Chinese Hamster Ovary Cell
19 Metabolism. *Cell Systems*, 3(5), 434–443.e8. <https://doi.org/10.1016/j.cels.2016.10.020>
- 20 Huang, Z., & Yoon, S. (2020). Integration of Time-Series Transcriptomic Data with Genome-
21 Scale CHO Metabolic Models for mAb Engineering. *Processes*, 8(3), 331.
22 <https://doi.org/10.3390/pr8030331>

Megchelenbrink, W., Huynen, M., & Marchiori, E. (2014). optGpSampler: An Improved Tool for Uniformly Sampling the Solution-Space of Genome-Scale Metabolic Networks. *PLOS ONE*, 9(2), e86587. <https://doi.org/10.1371/journal.pone.0086587>

Mulukutla, B. C., Kale, J., Kalomeris, T., Jacobs, M., & Hiller, G. W. (2017). Identification and control of novel growth inhibitors in fed-batch cultures of Chinese hamster ovary cells. *Biotechnology and Bioengineering*, 114(8), 1779–1790. <https://doi.org/10.1002/bit.26313>

Nam, H., Lewis, N. E., Lerman, J. A., Lee, D.-H., Chang, R. L., Kim, D., & Palsson, B. O. (2012). Network Context and Selection in the Evolution to Enzyme Specificity. *Science*, 337(6098), 1101–1104. <https://doi.org/10.1126/science.1216861>

Popp, O., Müller, D., Didzus, K., Paul, W., Lipsmeier, F., Kirchner, F., Niklas, J., Mauch, K., & Beaucamp, N. (2016). A hybrid approach identifies metabolic signatures of high-producers for chinese hamster ovary clone selection and process optimization. *Biotechnology and Bioengineering*, 113(9), 2005–2019. <https://doi.org/10.1002/bit.25958>

Ritacco, F. V., Wu, Y., & Khetan, A. (2018). Cell culture media for recombinant protein expression in Chinese hamster ovary (CHO) cells: History, key components, and optimization strategies. *Biotechnology Progress*, 34(6), 1407–1426. <https://doi.org/10.1002/btpr.2706>

Savinell, J. M., & Palsson, B. O. (1992). Network analysis of intermediary metabolism using linear optimization. I. Development of mathematical formalism. *Journal of Theoretical Biology*, 154(4), 421–454. [https://doi.org/10.1016/s0022-5193\(05\)80161-4](https://doi.org/10.1016/s0022-5193(05)80161-4)

Schellenberger, J., Que, R., Fleming, R. M. T., Thiele, I., Orth, J. D., Feist, A. M., Zielinski, D. C., Bordbar, A., Lewis, N. E., Rahmanian, S., Kang, J., Hyduke, D. R., & Palsson, B. Ø. (2011). Quantitative prediction of cellular metabolism with constraint-based models: The

COBRA Toolbox v2.0. *Nature Protocols*, 6(9), 1290–1307.

<https://doi.org/10.1038/nprot.2011.308>

Schinn, S.-M., Morrison, C., Wei, W., Zhang, L., & Lewis, N. (2020). *Systematic evaluation of parameterization for genome-scale metabolic models of cultured mammalian cells*.

Sommeregger, W., Sissolak, B., Kandra, K., von Stosch, M., Mayer, M., & Striedner, G. (2017). Quality by control: Towards model predictive control of mammalian cell culture bioprocesses. *Biotechnology Journal*, 12(7). <https://doi.org/10.1002/biot.201600546>

Szeliova, D., Ruckerbauer, D., Galleguillos, S., Petersen, Hanscho, M., Troyer, Causon, Schoeny, Christensen, Lee, D. Y., Lewis, N. E., Koellensperger, Hann, Nielsen, L. K., Borth, N., & Zanghellini, J. (2020). What CHO is made of: Variations in the biomass composition of Chinese hamster ovary cell lines. *Metabolic Engineering*.

Traustason, B., Cheeks, M., & Dikicioglu, D. (2019). Computer-Aided Strategies for Determining the Amino Acid Composition of Medium for Chinese Hamster Ovary Cell-Based Biomanufacturing Platforms. *International Journal of Molecular Sciences*, 20(21), 5464. <https://doi.org/10.3390/ijms20215464>

Zampieri, G., Vijayakumar, S., Yaneske, E., & Angione, C. (2019). Machine and deep learning meet genome-scale metabolic modeling. *PLoS Computational Biology*, 15(7). <https://doi.org/10.1371/journal.pcbi.1007084>

Zampieri, M., Hörl, M., Hotz, F., Müller, N. F., & Sauer, U. (2019). Regulatory mechanisms underlying coordination of amino acid and glucose catabolism in Escherichia coli. *Nature Communications*, 10(1), 3354. <https://doi.org/10.1038/s41467-019-11331-5>

Zhuangrong, H., & Seongkyu, Y. (2020). Identifying metabolic features and engineering targets for productivity improvement in CHO cells by integrated transcriptomics and genome-

1 scale metabolic model. *Biochemical Engineering Journal*, 107624.

2 <https://doi.org/10.1016/j.bej.2020.107624>

3

4

1 Figures

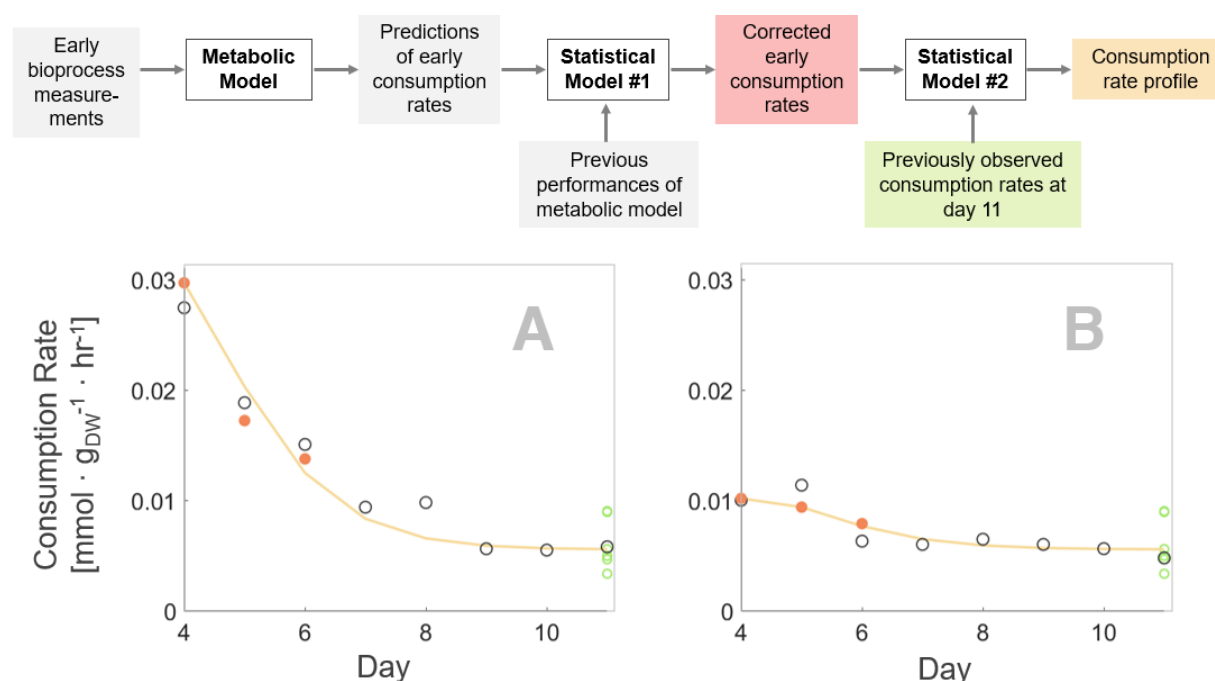


Figure 1: Overview of method. A novel combination of a metabolic and statistical models forecast the time-course amino acid consumption profiles in CHO cell cultures. A workflow of the prediction procedure is provided; key data are color-coded and visualized in plots A & B by the same color. First, a metabolic model predicts amino acid consumption rates for days 4-6 based on routine bioprocess measurements such as viable cell density and glucose uptake rate. Then, a statistical model refines these predictions (red) by considering the metabolic model's previous performances. Based on these predictions from early culture days, a second statistical model predicts the complete consumption profile (orange). The model references asymptotic behavior of previous consumption profiles as priors (green). The predicted consumption profiles agreed well with experimental data (black empty markers). The two plots show distinct leucine consumption profiles from CHO clones C2 and Z3 (Fig. S1) with disparate early consumption patterns. A more detailed workflow can be found in Fig. S6.

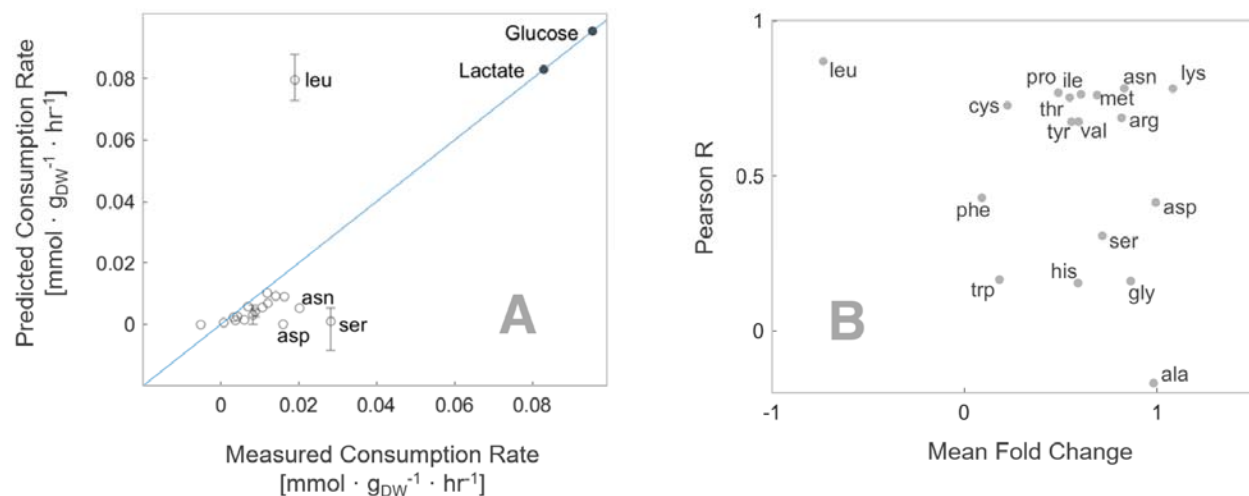


Figure 2: Metabolic network model estimates amino acid consumption rates. (a) Model predictions compared well to experimental observations, given the scale of input data such as the consumption rates of glucose and lactate (upper right, filled circles). (b) However, the fold change between model predictions and experimental measurements could be significant for several amino acids (x-axis). Fortunately, the relatively high linear correlation between predictions and measurements (y-axis) suggests that predictions could be improved empirically.

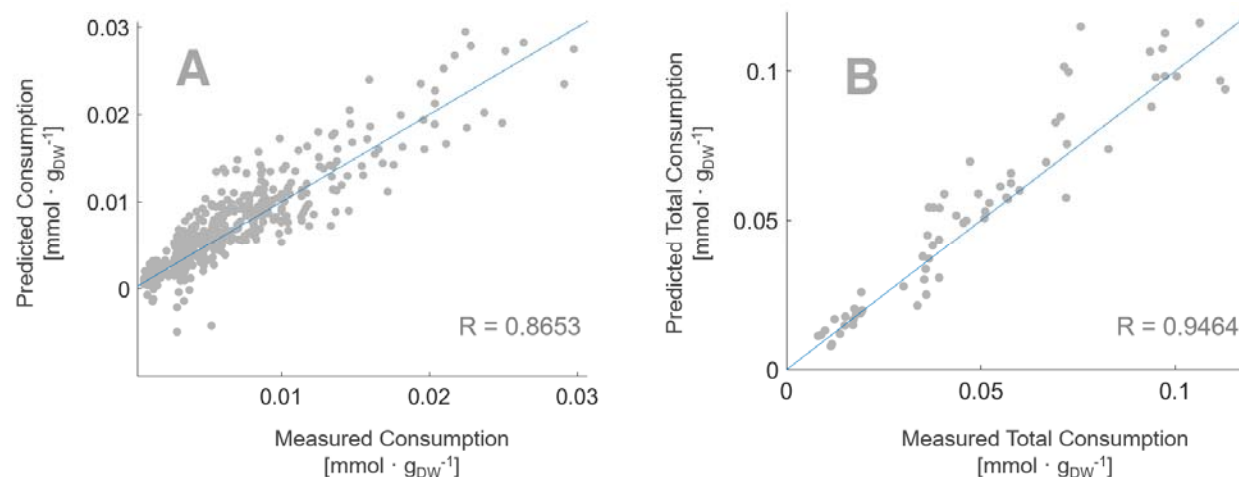


Figure 3: Statistical model forecasts consumption profiles. (A) For validation, the daily consumption rates were calculated for 16 amino acids (excluding alanine and glycine). On average, the predicted values agreed with experimental measurements to within 83%. (B) Then, the total amount of amino acid consumed across the investigated culture period were calculated by summing the daily consumption rates, which agreed with experimental measurements to within 86% on average. This suggests that the method can track and forecast the bioreactor environment.