



Method article

Guided extraction of genome-scale metabolic models for the integration and analysis of omics data



Andrew Walakira^a, Damjana Rozman^a, Tadeja Režen^a, Miha Mraz^b, Miha Moškon^{b,*}

^a Centre for Functional Genomics and Bio-Chips, Institute for Biochemistry and Molecular Genetics, Faculty of Medicine, University of Ljubljana, Ljubljana, Slovenia

^b Faculty of Computer and Information Science, University of Ljubljana, Ljubljana, Slovenia

ARTICLE INFO

Article history:

Received 8 April 2021

Received in revised form 4 June 2021

Accepted 4 June 2021

Available online 8 June 2021

Keywords:

Genome-scale metabolic model

Model extraction methods

Context-specific metabolic model

Omics data integration

Subsystem enrichment analysis

Model interpretability

ABSTRACT

Omics data can be integrated into a reference model using various model extraction methods (MEMs) to yield context-specific genome-scale metabolic models (GEMs). How to choose the appropriate MEM, thresholding rule and threshold remains a challenge. We integrated mouse transcriptomic data from a *Cyp51* knockout mice diet experiment (GSE58271) using five MEMs (GIMME, iMAT, FASTCORE, INIT and tINIT) in a combination with a recently published mouse GEM iMM1865. Except for INIT and tINIT, the size of extracted models varied with the MEM used (t-test: p-value <0.001). The Jaccard index of iMAT models ranged from 0.27 to 1.0. Out of the three factors under study in the experiment (diet, gender and genotype), gender explained most of the variability (>90%) in PC1 for FASTCORE. In iMAT, each of the three factors explained less than 40% of the variability within PC1, PC2 and PC3. Among all the MEMs, FASTCORE captured the most of the true variability in the data by clustering samples by gender. Our results show that for the efficient use of MEMs in the context of omics data integration and analysis, one should apply various MEMs, thresholding rules, and thresholding values to select the MEM and its configuration that best captures the true variability in the data. This selection can be guided by the methodology as proposed and used in this paper. Moreover, we describe certain approaches that can be used to analyse the results obtained with the selected MEM and to put these results in a biological context.

© 2021 The Author(s). Published by Elsevier B.V. on behalf of Research Network of Computational and Structural Biotechnology. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

1. Introduction

The advent of high-throughput technologies generate large volumes of omics data hence making it possible to study organisms at the cellular level. These studies have enabled a better understanding of the underlying biological processes of many diseases such as diabetes [1,2] and asthma [3], yielded useful products such as insulin [4,5] and antibiotics [6,7], and also improved production of commercial products such as wines [8]. However, there is still a big knowledge gap in the finer details of how organisms, unicellular or multicellular, are able to maintain life and how disruptions at the molecular level affect their phenotypes.

The phenotypic characteristics of an organism are determined by intricately connected reactions and consequently pathways that generate energy and other forms of biological products necessary to sustain life and organise organismal development. These pathways are connected to form and function as biological systems

and thus it is vital to study organisms at system level. The degree of complexity of biological systems differs greatly among organisms, for example, humans are far more complex than *Drosophila melanogaster* despite the latter's importance in modelling human diseases [9–11]. It is extremely difficult, or even impossible for higher level organisms, to study the entirety of their pathways either *in vitro* or *in vivo*. However, mathematical tools such as genome-scale models can be used to gain insight into how these biological systems function [12].

Genome-scale metabolic models (GEMs) are increasingly becoming a popular tool for studying biological processes *in silico* [13–15]. GEMs are formulated to contain all known biochemical reactions involved in maintaining the living state of a cell or an organism (metabolism). Context-specific metabolic models (also known as tissue-specific models) are GEMs in which inactive reactions for a given condition (context) are removed [16] and thus represent the context better since not all possible reactions are active in different cell types and/or in different contexts. Moreover, GEMs can be used to perform *in silico* studies and observe the dynamical response of the system in a given condition using

* Corresponding author.

E-mail address: Miha.Moskon@fri.uni-lj.si (M. Moškon).

computer simulations. As such, they provide better understanding of the organism as a functional system. One of the main problems of using GEMs in a combination with omics data lies in the complexity of the obtained models and the number of models produced. The obtained results are thus difficult to interpret in the context of biology. Moreover, it is hard to select an appropriate model extraction method (MEM) for a specific dataset. When the appropriate MEM is selected, it also needs to be configured, e.g., gene activity thresholds and thresholding rules need to be defined [17,18] to yield accurate and biologically relevant results.

We aim to address these problems by suggesting a methodology for performing analyses of omics data using GEMs in combination with different MEMs. This methodology will serve as an essential step towards the development of a pipeline that will automatically select a suitable MEM for a specific dataset, perform its configuration, and extract the models. Moreover, such a pipeline could provide its users with a set of publication-ready figures and tables, describing the results of simulations performed upon the derived models as well as the results of different statistical and enrichment analyses ready for a straightforward interpretation.

2. Materials and methods

In this section, we describe how a GEM is reconstructed from available biological information to produce a reference model. This is followed by a summary of the data integration algorithms that were applied in our analysis. Finally, we describe our analysis of the dataset used in our case study, namely, mouse diet experiment data.

2.1. Genome-scale metabolic models: from reconstruction to simulation

The construction of GEMs can be summarised in four main steps as described by Feist and colleagues [19]. First, a draft reconstruction of the biological network of an organism is extracted using information about reactions, enzymes, and pathways from databases such as KEGG, BRENDA, etc. The second step is the manual curation of the reconstructed draft model. This involves checking and filling the gaps and correcting misplaced reactions. Here, organism-specific databases and literature are used. Computational algorithms such as GAUGE [20], FastGapFill [21], and FBA-Gap [22] can also be applied. The third step is the conversion of the reconstructed model into a mathematical representation that can be used for subsequent simulations. The final step is the refining, validation, and application of the model to inform decisions. Here, the model outputs are compared with known information to confirm consistency. Model validation can be done by testing how well a model performs a set of tasks, comparing the simulation results with experimental data for a particular objective such as growth and use of gene essentiality analysis i.e. identification of genes that are required for survival of an organism [23].

The most efficient approach for prediction of phenotype from genotype using GEMs is constraint-based modelling [24,25]. Here, it is common to assume that the system is in a steady state, i.e., the concentrations of all metabolites involved are constant (see Eq. 1). Constraint-based modelling is focused on predicting flux distributions while optimising a selected cellular function (or a set of functions). Flux Balance Analysis (FBA) is the most widely used technique to predict flux distributions in GEMs [26,27]. It requires a mathematical representation of the model in the form of a stoichiometric matrix S with rows representing metabolites and columns representing reactions. We can evaluate a vector v of fluxes through the observed metabolic reactions constrained by the upper and lower flux bounds by using the equation

$$\sum_{j=1}^N S_{ij} \cdot v_j = 0, \quad (1)$$

$$\forall i \in \{1, 2, \dots, M\},$$

$$\forall j \in \{1, 2, \dots, N\},$$

where M is the number of observed metabolites and N the number of observed reactions. FBA reduces the problem to a linear program, hence lowering the computational requirements involved [26]. However, there are certain limitations. First, FBA does not yield a unique solution and is highly dependent on the choice of the objective function, i.e., a description of the phenotype relevant to the problem being studied [28]. While in some cases biomass optimisation is a plausible biological objective (e.g., cancer cells, cell lines, single cell organisms), different optimisation criteria need to be applied depending on the question of interest. How to select an appropriate objective function is an unanswered question. Algorithms that can select the objective function automatically have been proposed [29]. There is also increasing evidence that multiple objective functions are required to allow metabolic flexibility and improve accuracy of the model [30,31]. The other challenges are that FBA is dependent on the choice of the solver (used to solve the set of linear equations) and the quality of reconstruction [27,32]. Variations of FBA have been proposed to ease these limitations for example, parsimonious FBA (pFBA) [33] and Flux Variability Analysis (FVA) [34], among others [27,35]. Another method for predicting metabolic fluxes is flux sampling, which can be used to estimate probability distributions of reaction fluxes without assuming any particular cellular objective [36,37].

2.2. Computational approaches for experimental data integration

Model extraction methods can be classified as members of a “GIMME-like” family which minimises flux through reactions associated with low gene expression, an “iMAT-like” family which finds an optimal trade-off between keeping reactions whose genes are highly expressed and removing reactions associated with low gene expression, and an “MBA-like” family which retains a selected set of core reactions [16,38]. Below is a summary of the data integration algorithms that were used in this analysis.

2.2.1. GIMME

Gene Inactivity Moderated by Metabolism and Expression (GIMME) [16,39,40] uses gene expression data and one or more objective functions to produce a context-specific model. The GIMME algorithm takes three inputs, i.e., expression data, genome scale reconstruction, and one or more required metabolic functionalities (RMFs) that the cell is expected to achieve. The algorithm first finds a flux distribution that optimises the given objective(s) and then minimises the use of inactive reactions (reactions whose expression is below a predefined threshold). The expression data is used directly as weights in the objective function. A threshold is used to determine if a weight in the objective is positive or negative. A weight of zero is assigned to reactions without expression data. The method yields an inconsistency score i.e. $(flux * (threshold - data))$, a score that shows the disagreement between the expression data and the metabolic objective function. The normalised version of this score shows how well each gene in the expression data agrees with a particular metabolic function. GIMME has been successfully used in studies such as [41] aimed at understanding the impact of drought stress on *Arabidopsis thaliana*.

2.2.2. iMAT

The Integrative Metabolic Analysis Tool (iMAT) [42] takes three-valued expression data as inputs. The data are categorised as lowly, moderately, or highly expressed and coded as -1 , 0 , and 1 , respectively. A Boolean gene-to-reaction mapping is used to identify the state of a reaction, i.e., if the genes encoding enzymes of the reaction are low, moderate or highly expressed. This leads to classifying the reactions in the model as either highly or lowly expressed. This is followed by finding the steady state flux distribution that satisfies the stoichiometric and dynamic constraints and maximises the number of reactions whose activity is consistent with their expression state. A reaction is considered active if it carries a significant positive flux (or negative flux for reversible reactions) that is greater than a threshold. A reaction is inactive if it carries a flux of zero (0). The algorithm returns a vector showing the predicted activity state (fluxes) of each reaction. iMAT performs a pathway enrichment analysis and identifies up- and down-regulated genes thus shedding light on the active pathways in the conditions under study.

2.2.3. FASTCORE

FASTCORE [43] is a data integration algorithm that accepts a core set of reactions that are known to be active in regard to the context under study. The core set of reactions can be determined by considering reactions in which highly expressed genes (genes whose expression level is above a predefined threshold) are involved. This is followed by a search for a flux consistent subnetwork, i.e., a network in which each reaction has a non zero flux in at least one feasible flux distribution. Such subnetwork presents a context-specific model which contains no blocked reactions.

2.2.4. INIT and tINIT

Integrative Network Inference for Tissues (INIT) [44] formulates a mixed integer-linear problem (MILP) designed to use data from the Human Protein Atlas (HPA) and other omics data as inputs. INIT does not apply a strict steady state assumption for all internal metabolites. Instead it allows a small positive net production of metabolites which are given positive weights in the optimization. Consequently, all reactions in the resulting model are able to carry flux. The algorithm produces networks that are snapshots of active metabolism [45]. tINIT is the task-driven version of INIT. Here, a set of tasks that must be carried out by the resulting model are defined first and then followed by the INIT algorithm.

2.3. Case study using gene expression data from *Cyp51* knockout mice

The mouse gene expression dataset was downloaded from the GEO database (accession number GSE58271) and processed to obtain normalised gene expression values. Briefly, this dataset was generated from a study in which the mice were divided into three groups and fed on three diets i.e. low fat without cholesterol (LFnC), high fat without cholesterol (HFnC) and high fat with cholesterol (HFC). Each diet group contained both wild type and the *Cyp51* knockout genotype in female and male mice. The detailed description of the dataset can be found in [46].

Extraction of context-specific models was performed in Matlab R2019b (MathWorks Inc., Natick, Massachusetts, USA) using normalised gene expression data and each of the model extraction methods (MEMs) above as described by their respective authors. A recently published mouse model, iMM1865 was used as a reference model. This model has 10612 reactions, 5839 metabolites and 93 subsystems and has no dead-end metabolites or blocked reactions [47]. Highly expressed genes for model extraction were determined by setting the threshold values at the 50th, 70th, 75th,

and 80th percentiles of the normalised gene expression data per sample. The rationale of setting this threshold per sample is that all individuals are biologically distinct even when under similar conditions. They yield unique expression patterns implying that their set of highly expressed genes may differ because of these intricate biological differences. Additionally, we assessed the impact of using thresholds set per each gene separately, i.e., thresholds defined within a gene. Here, we considered the 80th, 90th, 95th percentile and the mean for each gene within the observed population. The thresholds were chosen based on the trend of the variance explained by the first principal component (PC1) (Fig. 1c and d). The extracted models from each algorithm were compared with one another (pairwise comparison) using the Jaccard index (distance) metric [48] to identify the distance between models extracted under the same conditions.

Principal component analysis (PCA) was performed as described in Opdam et al. 2017 [16]. For every MEM, a matrix showing if a reaction was present (1) or absent (0) was generated. Reactions that were present or absent in all observed models were removed, and row means of the matrix were zero-centered. Principal component analysis on reactions was then performed. Furthermore, the variance explained by each factor (MEM, diet, gender, genotype) within a principal component was calculated by taking the square of the maximum Pearson correlation coefficient (R^2) of the component scores across all possible orderings of the factors as described in [16]. This was reported as a percentage. The validity of model separation observed in PCA was confirmed using t-distributed stochastic neighbour embedding (t-SNE), which is particularly suitable for high-dimensional data [49].

To assess the dynamical response of the models, we performed flux sampling using the artificial centering hit-and-run (ACHR) algorithm [50] on the extracted models. 1000 flux samples were generated for each of the models, and the mean flux of each reaction was used to identify the reactions that are either down- or up-regulated in pairwise comparisons of specific factors (diet, gender, and genotype) according to Spearman's rank correlation. The reactions identified to be significantly changed were then used to perform the enrichment analysis of metabolic subsystems using the hypergeometric test. The obtained p-values were adjusted for multiple testing using the Benjamini-Hochberg procedure. The implementation of the described analysis is available as a set of IPython (IPYNB) notebooks at https://github.com/CompBioLj/GEMS_and_MEMS.

3. Results

3.1. Thresholds affect the extracted models

Our results indicate that the models extracted with the iMAT methodology were able to explain the highest amount of variance in comparison to other model extraction methods (see Section 3.2 and Fig. 2d). We thus opted to use the iMAT GEMs to assess the impact of thresholding on the extracted models. Two types of thresholds were used. First, thresholds were considered for each sample. This was done by taking a certain percentile of all the data within each sample to get the cutoff for highly expressed genes in that sample. This was achieved by taking a certain percentile of data to get the cutoff for highly expressed genes in that sample. The 50th, 70th, 75th and 80th percentile were considered. With the exception of models extracted at the 50th percentile threshold, the range and distribution of the Jaccard index and the size of the models were generally similar (see Figs. 1a and b). The percentage of variance explained by the PC1 was the smallest at the 50th

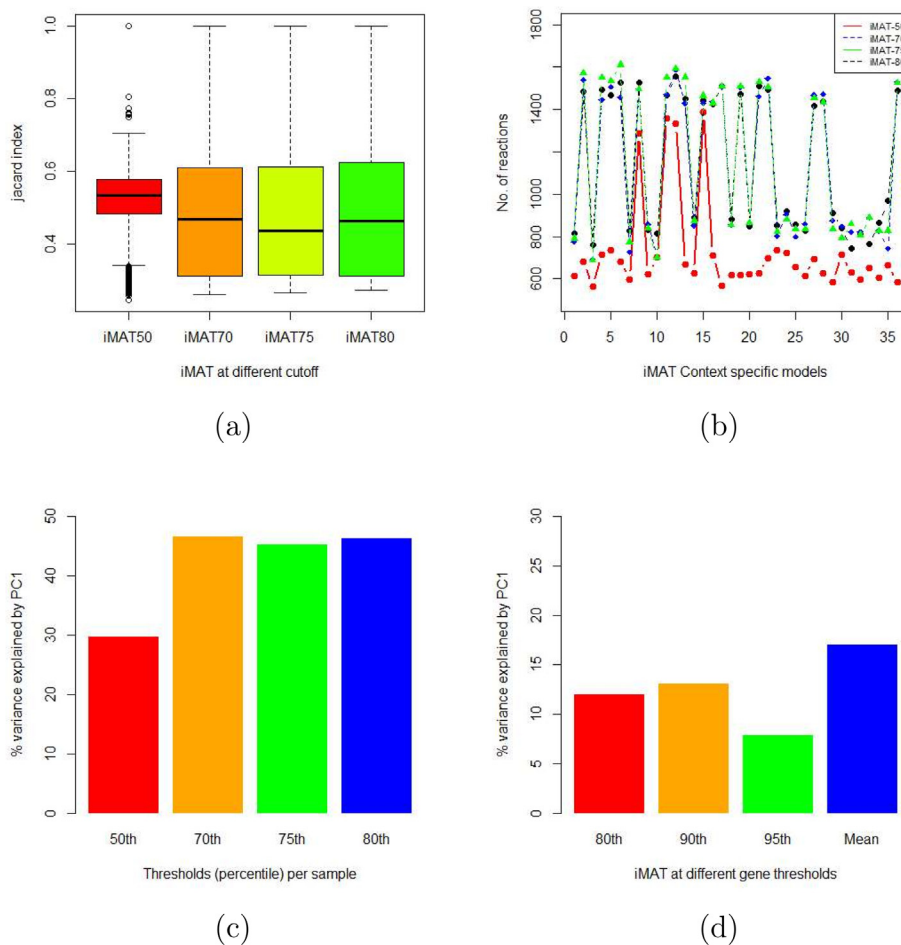


Fig. 1. Analysis of iMAT derived models extracted using different thresholding rules and values. Fig. 1 (a), (b), and (c) present the analysis of the models extracted by thresholding per sample, where figure (a) shows the Jaccard index from the pairwise comparisons of all models for each threshold, figure (b) how these models are varied in size, and figure (c) the percentage of variance explained by PC1 for each threshold. Figure (d) shows the variance explained by PC1 in the models extracted when thresholding is performed per gene. PC1: first principal component.

percentile threshold but similar for other threshold values (see Fig. 1c).

The second type of thresholds was set per gene, i.e., by taking a percentile per gene for all samples. In this case the same threshold values were used for all samples, but were different among different genes. We considered the 80th, 90th and 95th percentile gene expression for each gene as a cutoff to classify if a gene is highly expressed or not. We also considered taking the mean of each gene as a cutoff for the same purpose. From the perspective of principal component analysis, the highest variance explained by PC1 was 17% (see Fig. 1d). Considering thresholds at 80th, 90th and 95th percentile, the highest variance explained by PC1 was achieved at the 90th percentile threshold.

PC1s from models where the threshold was taken per sample explained more variance than from models where the threshold was taken per gene. In general, the choice of thresholding strongly affected the content of the extracted models and their ability to capture variance in the data. Thresholding per sample with an 80th percentile as a cutoff value was considered as the most appropriate and was used in further analyses.

3.2. Extracted models vary with algorithm

Context-specific models were extracted using the GIMME, iMAT, FASTCORE, INIT, and tINIT model extraction methods

(MEMs). We performed the extraction process using the COBRA [51,52] and RAVEN Toolboxes [53] in Matlab R2019b (MathWorks Inc., Natick, Massachusetts, USA) using the Gurobi solver [54].

The number of remaining reactions in an extracted model was considered to represent the size of a model. Different MEMs generated models of different sizes (see Fig. 2a). iMAT-produced models were significantly different in size from other MEMs (t-test: p-value < 0.001). For each model within each MEM, we identified reactions that were present (1) or absent (0) and calculated the Jaccard index between all possible pairwise combinations to compare the similarities between models. For GIMME, FASTCORE, tINIT, and INIT, the Jaccard index was between 0.8 and 1.0, implying that these models are very similar (see Fig. 2b). For iMAT, the Jaccard index ranged from 0.27 to 1.0, indicating that GEMs extracted with iMAT varied substantially (see Figs. 2 b,c). Furthermore, we performed principal component analysis (PCA) on the matrix of reactions for each MEM. The PC1 explained the highest variance in iMAT derived models in comparison to PC1s of other MEMs (see Fig. 2d).

In addition, we analysed how well do the clusters observed in the feature space described by the first two principal components (PC1 and PC2) comply with the groups defined by the experiment, namely diet, gender, and genotype. The PCA plot (PC1 versus PC2) of models extracted with iMAT did not show consistent clustering of samples to predefined groups (see Supplementary Fig. 1). This was similar for GIMME (see Supplementary Fig. 2). Separation by

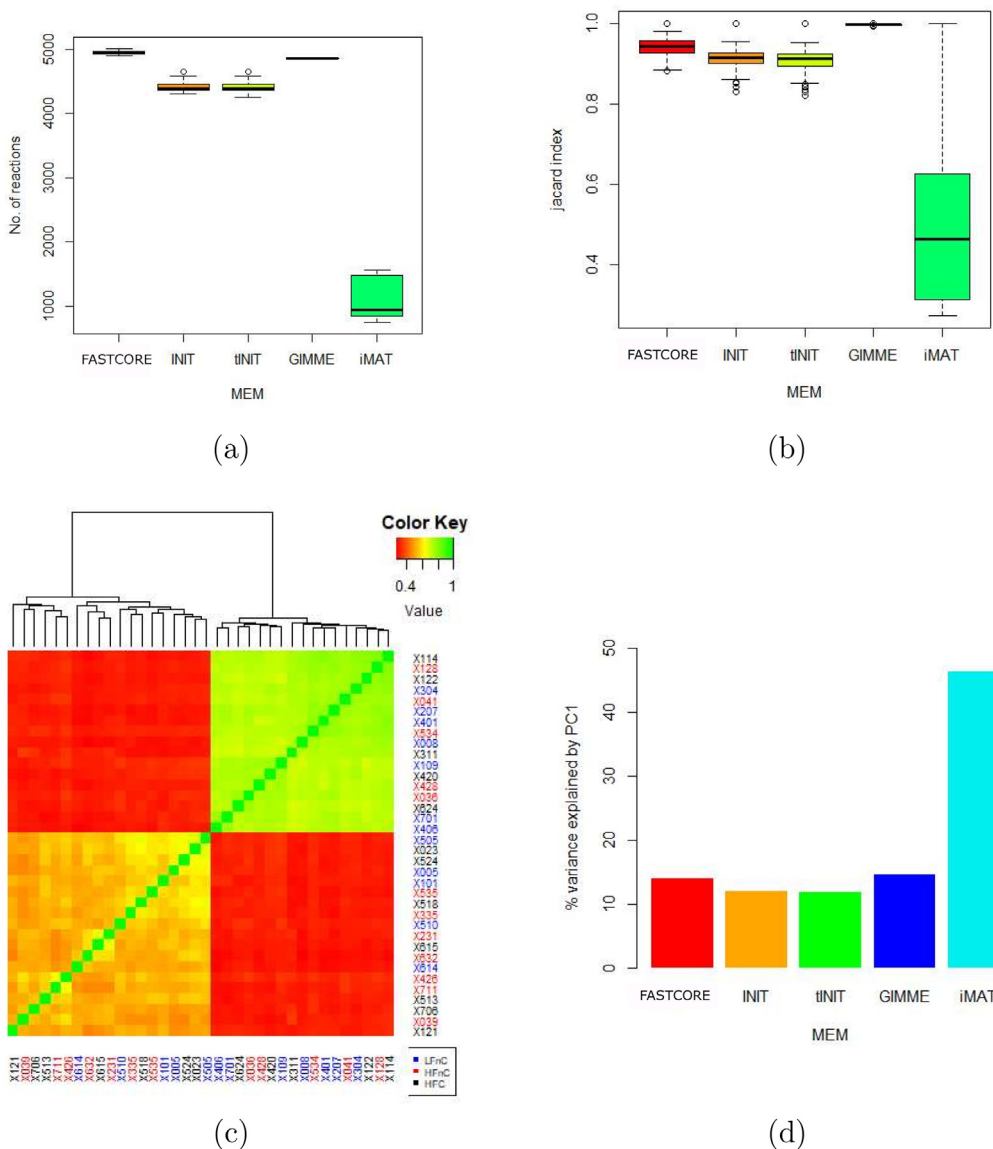


Fig. 2. Comparison of models produced by different MEMs. Fig. 2 (a) summarizes the size of the models, (b) summarizes the Jaccard indices, (c) shows the Jaccard indices of models produced by iMAT, and (d) shows the variance explained by PC1 of each MEM. PC1: first principal component; MEM: model extraction method.

gender was consistent with the clustering by PC2 of INIT explaining approximately 8% of variance, and at least partially consistent with the clustering by PC2 of tINIT explaining approximately 7.4% of variance in the data (see Supplementary Figs. 3 and 4). GEMs extracted using FASTCORE were appropriately clustered by gender by the PC1 explaining approximately 13% of variance in the data (see Fig. 3 and Supplementary Fig. 5). All PCA plots performed on the FASTCORE extracted models are available as Supplementary Fig. 5.

To verify the validity of the clustering obtained using PCA, we additionally performed t-distributed stochastic neighbour embedding (t-SNE) using different values of perplexity parameter [49]. Consistent clustering was obtained when the perplexity value was at least 15. The observed separation of models complied with the PCA results (see Supplementary Figs. 6 and 7). Separation of models was consistent with gender in FASTCORE, INIT and tINIT, whereas using the PCA separation by gender was observed in FASTCORE and INIT and partially also in tINIT. Clustering by gender and only by gender was also observed in the original study [46], which confirms the adequacy of our analysis.

We further assessed how much each of the factors (MEM, diet, gender and genotype) contributes to the variation explained by the first three principal components (PCs). This was evaluated by taking each factor and calculating the squares of Pearson correlation coefficient (R^2) of component scores across all possible orderings of the factors [16]. Considering all factors, the selection of a MEM explains the most variability in the PCs (Fig. 4a). Furthermore, in FASTCORE (Fig. 4b) and GIMME (Fig. 4c) extracted models, gender explained the most variability in the PCs as was also observed in the original study [46].

In the models extracted with iMAT (Fig. 4d), the observed three factors, namely, gender, genotype, and diet, generally contributed equally with genotype having a slight edge in comparison to diet and gender. Gender explained the most variability in models extracted with INIT and tINIT (Fig. 4e,f).

3.3. Enrichment analysis of metabolic subsystems

The models extracted with FASTCORE were able to capture the largest amount of true variability of the observed data (see Fig. 3).

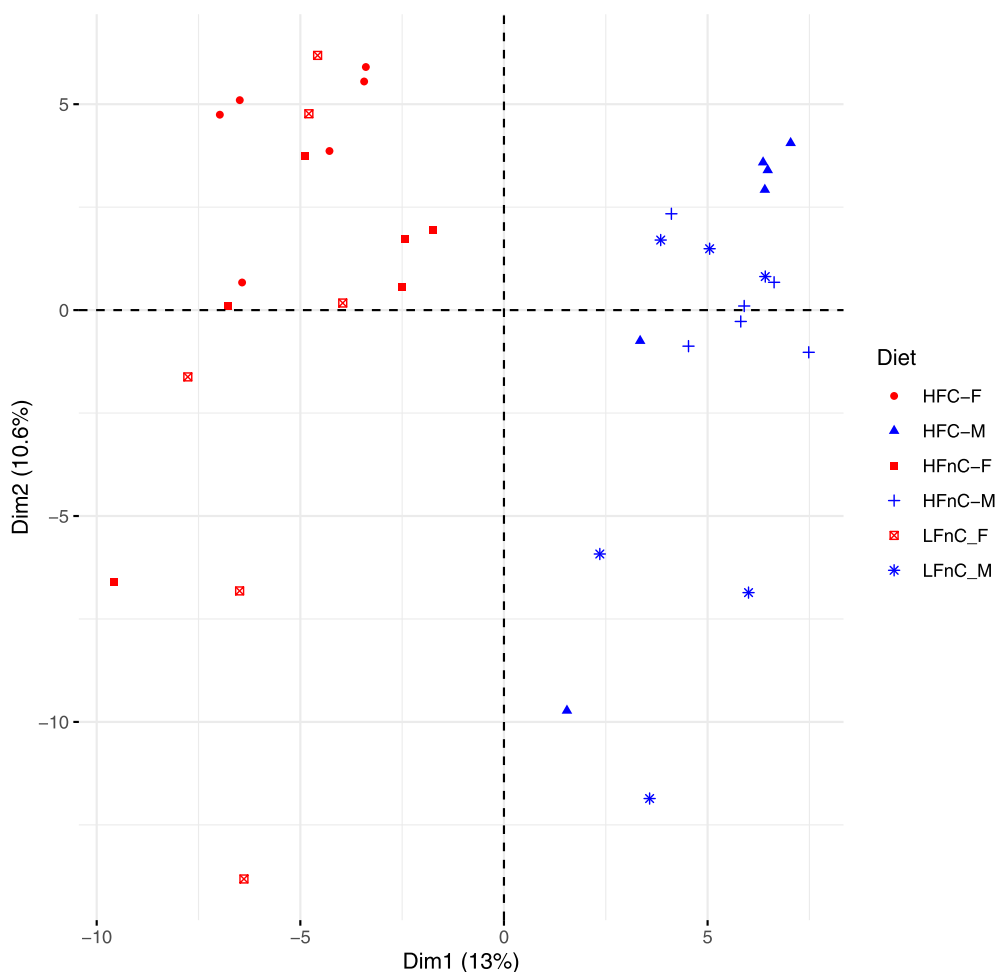


Fig. 3. PCA plot showing separation by gender performed on the FASTCORE extracted models. Blue colour indicates male and red colour female samples, respectively. PCA: principal component analysis; F: female; M: male; LFnC: low fat without cholesterol; HFnC: high fat without cholesterol; HFC: high fat with cholesterol.

We thus opted to perform further analysis only on these models. Since the basic version of flux balance analysis (FBA) is unable to provide unique solutions, we opted to analyse the activity of the observed metabolic reactions using flux sampling using the artificial centering hit-and-run (ACHR) algorithm [50]. It was performed on the models extracted with the FASTCORE algorithm to obtain the mean values of reaction fluxes in each of the models. We performed the flux sampling using COBRA [51,52] and the Gurobi solver [54] in Matlab R2019b.

We identified up- and down-regulated reactions between different groups of extracted models (e.g., wild type versus knockout groups) and their combinations using Spearman's rank correlation. Genome-scale metabolic models are usually composed of several subsystems containing metabolic reactions with a specific function (e.g. cholesterol synthesis subsystem) [52,53]. Based on the identification of differentially regulated reactions, we further analysed the enriched subsystems between different groups using the hypergeometric test. We compared the models on the basis of diet (see Fig. 5).

To check if the obtained models comply with our expectations, we further observed how different diets affect the cholesterol synthesis and metabolism subsystems. Fig. 5 shows that cholesterol metabolism was enriched in mice on LFnC or HFnC diet compared to HFC diet. However, cholesterol metabolism was down-regulated in mice on LFnC diet compared to HFnC diet. Furthermore, Cholesterol synthesis was up-regulated in female wild type mice on LFnC diet compared to the HFC diet. In female knockout mice,

cholesterol synthesis was down-regulated in mice on the LFnC diet in comparison to the mice on the HFnC diet. In wild type mice, cholesterol synthesis was up-regulated in mice on the HFnC diet compared to the HFC diet. There was no significant subsystem enrichment due to genotype or gender alone. A 5% level of significance after the adjustment for multiple testing (p -value: < 0.05) was considered in all analyses.

4. Discussion

The aim of this work was to highlight the factors that have the strongest influence on the context-specific extraction of genome-scale metabolic models, especially in relation to the model-based analysis of omics data. Additionally, we proposed a methodology that could be followed in such analyses (see Fig. 6). This is composed of the following steps: (1) identification of the most suitable reference model; (2) extraction of context-specific GEMs using different MEMs with different configurations; (3) identification of the MEM and its configuration that is able to capture the variance of the observed data as well as the groups defined in the experiment; (4) analysis of obtained models using different approaches, such as PCA, t-SNE, and metabolic subsystems enrichment analysis.

We extracted GEMs using five model extraction methods i.e. GIMME, iMAT, FASTCORE, INIT and tINIT. The results show that the choice of a MEM affects the structure and contents of the output models. Additionally, the models varied greatly in size between MEMs. INIT and tINIT are very similar algorithms and so are their

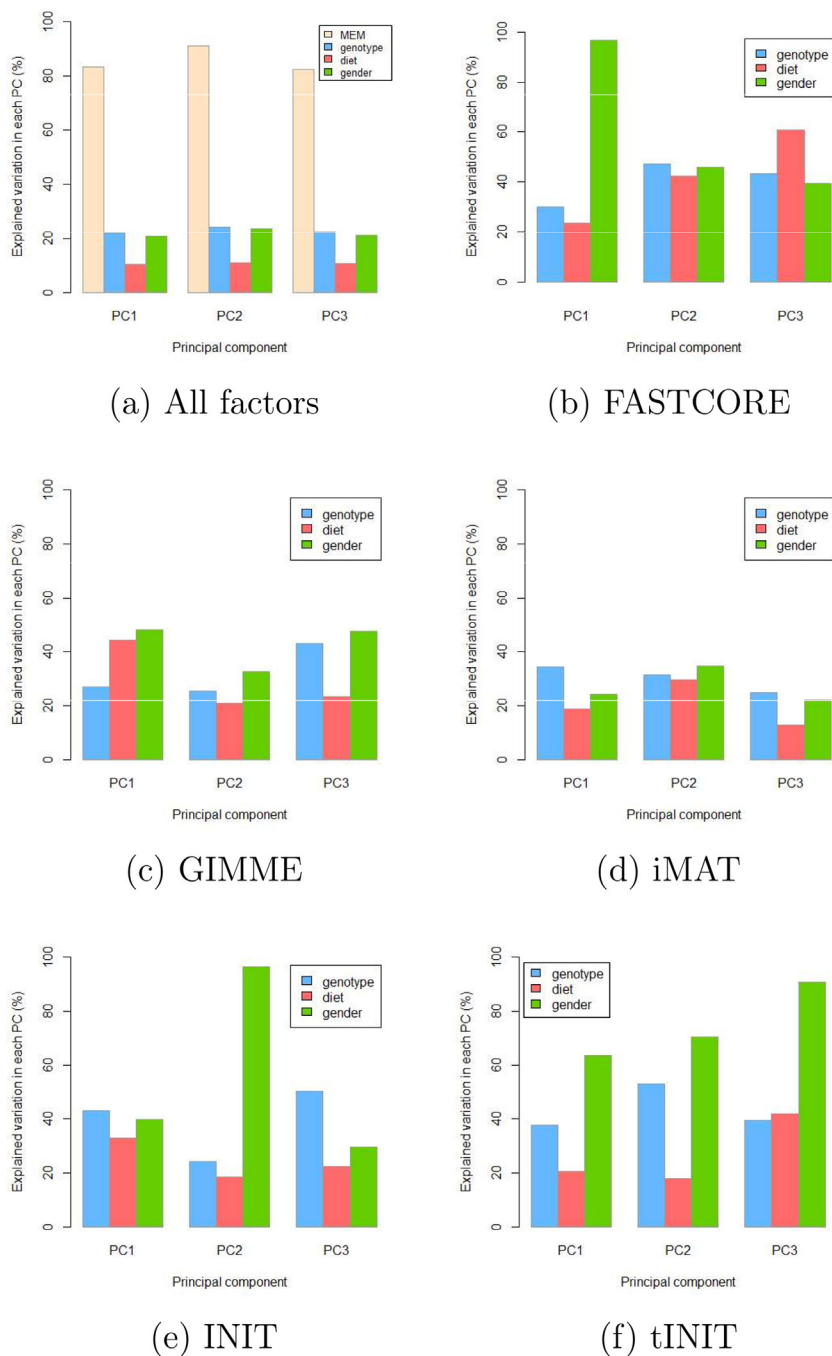


Fig. 4. Contribution of different factors, namely, MEM (yellow), diet (red), genotype (blue) and gender (green) to PC1, PC2 and PC3. Fig. 4 (a) shows the contribution of all factors. Remaining figures show the contribution of diet, gender and genotype in models extracted with FASTCORE (b), GIMME (c), iMAT (d), INIT (e) and tINIT (f). MEM: model extraction method; PC: principal component.

output models. These results are consistent with published reports [16,17,18]. We used Jaccard index to assess how much the models differed within each MEM and PCA to assess how much variance in the data each MEM was able to capture. GIMME produced very similar models and captured the least variability in the data. iMAT captured the most variability in the data (46% by PC1) but was unable to capture the true variability in the experiment. However, FASTCORE with 13% of variance explained by PC1 in this example was able to capture the largest amount of variability due to any of the observed factors, namely gender. Similar clustering was also observed in the original study [46]. This implies that it is important to select a MEM that not only captures the largest variance in the

data, but also captures the groups that are aligned with predefined experimental groups. These groups could reflect different metabolic signatures and thus have a potential to guide downstream analyses. Our findings assert that the choice of a MEM greatly affects the output models.

We further assessed the impact of using different thresholds on the output models. Thresholds were set on gene expression data to identify highly expressed genes. For each model extracted with iMAT, we applied two kinds of thresholding, i.e., within a sample and within a gene. Thresholding within a sample captured unique individual differences in expression even if individuals were from the same experimental group. We see that the type of thresholding

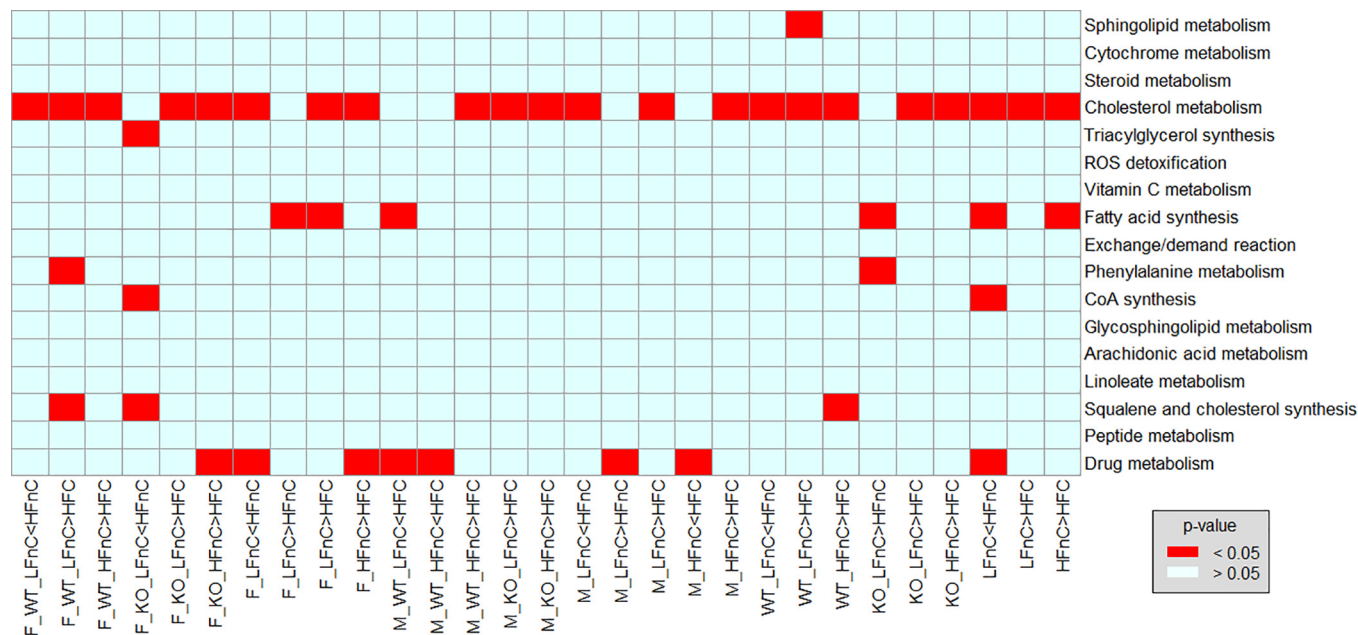


Fig. 5. Metabolic subsystems enrichment analysis between different diets. Subsystems that were differentially enriched between two groups and those associated with cholesterol are shown. Cholesterol metabolism was enriched in mice on HFnC diet compared to HFC diet. The less than symbol (<) between the groups corresponds to down-regulation and the more than symbol (>) corresponds to up-regulation. F: female; M: male; WT: wild type; KO: knockout; LFnC: low fat without cholesterol; HFnC: high fat without cholesterol; HFC: high fat with cholesterol.

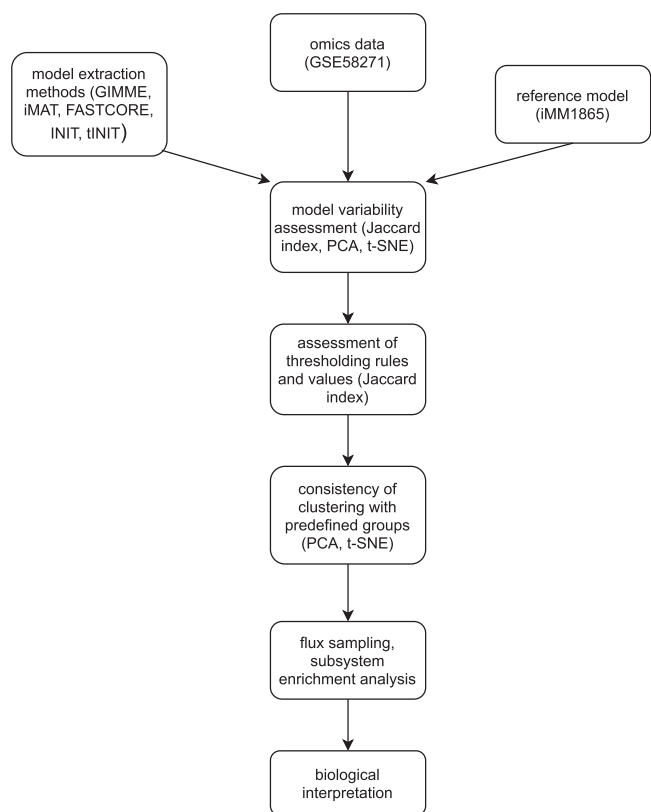


Fig. 6. Overview of the proposed methodology. PCA: principal component analysis; t-SNE: t-distributed stochastic neighbour embedding.

applied and the threshold values used strongly affect the output models.

We performed subsystem enrichment analysis to assess how cholesterol metabolism and synthesis varied between the groups

separated by different diets. This is because knocking out the *Cyp51* gene interrupts cholesterol synthesis and metabolism [46,55–57]. In our analysis, models extracted with FASTCORE showed that cholesterol metabolism significantly changed between groups with different diets. Cholesterol metabolism was enriched in mice on LFnC or HFnC diet compared to HFC diet. However, cholesterol metabolism was down-regulated in mice on LFnC diet compared to HFnC diet.

The proposed methodology, like other analyses using GEMs, requires selection of a suitable reference model. The choice of a reference model strongly affects the output models. The model should optimally describe the same tissue or at least the same organism as the samples used in the experiment. Moreover, the model should contain accurate gene-product-reaction (GPR) rules, which are used by MEMs to connect the omics data with metabolic reactions. In our case, we used the most recently published mouse model, namely iMM1865 [47].

Different MEMs implement different algorithms for data integration. As such, the output models differ greatly in their contents. In this analysis, we opted to use a selection of the state-of-the-art MEMs, namely, GIMME, iMAT, FASTCORE, INIT, and tINIT. However, other algorithms, such as CORDA [58], mCADRE [59] and MBA [60] could as well be applied. Selected MEMs also need to be configured as required by a particular algorithm, e.g., with the identification of the most suitable thresholding rules and values. The choice of values for a particular configuration needs to be guided by known knowledge and further analyses performed on selected MEMs. Identification of the most suitable MEM(s) needs to be based on the results from the analysis of the obtained models using different approaches, such as PCA, t-SNE, various distance metrics, and metabolic subsystem enrichment analysis. Ideally, a MEM of choice will explain the most variance in the observed dataset and will be able to separate the models in compliance with the groups defined in the experiment. Moreover, the obtained results of further analyses of the models extracted with the selected MEM will reflect the biological relevance and will provide a platform to generate novel knowledge and hypotheses. Such an ideal

situation may not be obtained easily. When compromises are required, the MEM that best captures the groups defined by the experiment should be selected. In this context, the MEM that has the largest percentage of explained variance within the PC1 for a selected group should be selected (see Fig. 4).

Even though the proposed methodology was demonstrated using only transcriptomics data, different omics data could as well be used in the process. For example, metabolomics and lipidomics data can be mapped into the context of metabolic reactions [61]. These can be used to identify a set of tasks or core reactions that need to be executed by a model. Moreover, specific model extraction methods, such as tINIT, allow a direct integration of non-transcriptomics data. Namely, tINIT algorithm accepts a list of metabolites the model should be able to produce [45].

We propose that the analysis of omics data using GEMs should be initiated with an extraction of different context-specific models using different MEMs. This should be followed with the application of various thresholding rules and different threshold values in the process of reconstruction. The resulting models should then be analysed to select the most appropriate MEM(s), threshold values and thresholding rules that best capture the variability in the data while capturing the known experimental groups.

CRediT authorship contribution statement

Andrew Walakira: Methodology, Software, Formal analysis, Investigation, Writing - original draft, Writing - review & editing, Visualization. **Damjana Rozman:** Conceptualization, Writing - review & editing, Funding acquisition. **Tadeja Režen:** Conceptualization, Writing - review & editing. **Miha Mraz:** Writing - review & editing, Funding acquisition, Supervision. **Miha Moškon:** Conceptualization, Methodology, Software, Formal analysis, Investigation, Writing - review & editing, Supervision.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

This work was supported by the Horizon 2020 TransSYS Marie Curie Initial Training Network Grant Agreement No. 860895, by the scientific research programs Pervasive Computing (P2-0359) and Functional Genomics and Biotechnology for Health (P1-0390) both financed by the Slovenian Research Agency, and by the basic research project CholesteROR in metabolic liver diseases (J1-9176) financed by the Slovenian Research Agency. We would like to acknowledge John Hancock for proofreading the manuscript.

Appendix A. Supplementary data

Supplementary data associated with this article can be found, in the online version, at <https://doi.org/10.1016/j.csbj.2021.06.009>.

References

- [1] Wang R, Li B, Lam SM, Shui G. Integration of lipidomics and metabolomics for in-depth understanding of cellular mechanism and disease progression. *J Genet Genomics* 2020;47(2):69–83.
- [2] Faulkner A, Dang Z, Avolio E, Thomas AC, Batstone T, Lloyd GR, Weber RJ, Najdekr L, Jankevics A, Dunn WB, et al. Multi-omics analysis of diabetic heart disease in the db/db model reveals potential targets for treatment by a longevity-associated gene. *Cells* 2020;9(5):1283.
- [3] Ivanova O, Richards LB, Vijverberg SJ, Neerincx AH, Sinha A, Sterk PJ, Maitland-van der Zee AH. What did we learn from multiple omics studies in asthma?. *Allergy* 2019;74(11):2129–45.
- [4] Hwang H-G, Kim K-J, Lee S-H, Kim C-K, Min C-K, Yun J-M, Lee SU, Son Y-J, et al. Recombinant glargine insulin production process using escherichia coli. *J Microbiol Biotechnol* 2016;26(10):1781–9.
- [5] Govender K, Naicker T, Lin J, Baijnath S, Chaturgoon AA, Abdul NS, Docrat T, Kruger HG, Govender T. A novel and more efficient biosynthesis approach for human insulin production in escherichia coli (e. coli). *AMB Express* 2020;10(1):1–9.
- [6] Pulido MR, García-Quintanilla M, Gil-Marqués ML, McConnell MJ. Identifying targets for antibiotic development using omics technologies. *Drug Discov Today* 2016;21(3):465–72.
- [7] Mohite OS, Weber T, Kim HU, Lee SY. Genome-scale metabolic reconstruction of actinomycetes for antibiotics production. *Biotechnol J* 2019;14(1):1800377.
- [8] Siren K. Multi-omics methods to unravel microbial diversity in fermentation of riesling wines [Ph.D. thesis]. TU Kaiserslautern: Germany; 2019..
- [9] Mirzoyan Z, Sollazzo M, Allocca M, Valenza AM, Grifoni D, Bellostà P. *Drosophila melanogaster*: a model organism to study cancer. *Front Genet* 2019;10:51.
- [10] Xiong Y, Yu J. Modeling parkinson's disease in drosophila: what have we learned for dominant traits?. *Front Neurol* 2018;9:228.
- [11] Ugur B, Chen K, Bellen HJ. *Drosophila* tools and assays for the study of human diseases. *Disease Models Mech* 2016;9(3):235–44.
- [12] Gu C, Kim GB, Kim WJ, Kim HU, Lee SY. Current status and applications of genome-scale metabolic models. *Genome Biol* 2019;20(1):1–18.
- [13] Blagotinšek Cokan K, Urlep Ž, Moškon M, Mraz M, Kong XY, Eskild W, Rozman D, Juvan P, Režen T. Common transcriptional program of liver fibrosis in mouse genetic models and humans. *Int J Mol Sci* 2021;22(2):832.
- [14] Schneider WM, Luna JM, Hoffmann H-H, Sánchez-Rivera FJ, Leal AA, Ashbrook AW, Le Pen J, Ricardo-Lax I, Michailidis E, Peace A, et al. Genome-scale identification of SARS-CoV-2 and pan-coronavirus host factor networks. *Cell* 2021;184(1):120–32.
- [15] Hozhabri H, Lashkari A, Razavi S-M, Mohammadian A. Integration of gene expression data identifies key genes and pathways in colorectal cancer. *Med Oncol* 2021;38(1):1–14.
- [16] Opdam S, Richelle A, Kellman B, Li S, Zielinski DC, Lewis NE. A systematic evaluation of methods for tailoring genome-scale metabolic models. *Cell Syst* 2017;4(3):318–29.
- [17] Richelle A, Joshi C, Lewis NE. Assessing key decisions for transcriptomic data integration in biochemical networks. *PLoS Comput Biol* 2019;15(7):e1007185.
- [18] Richelle A, Chiang AW, Kuo C-C, Lewis NE. Increasing consensus of context-specific metabolic models by integrating data-inferred cell functions. *PLoS Comput Biol* 2019;15(4):e1006867.
- [19] Feist AM, Herrgård MJ, Thiele I, Reed JL, Palsson BØ. Reconstruction of biochemical networks in microorganisms. *Nat Rev Microbiol* 2009;7(2):129–43.
- [20] Hosseini Z, Marashi S-A. Discovering missing reactions of metabolic networks by using gene co-expression data. *Scientific Rep* 2017;7(1):1–12.
- [21] Thiele I, Vlassis N, Fleming RM. fastgapfill: efficient gap filling in metabolic networks. *Bioinformatics* 2014;30(17):2529–31.
- [22] Brooks JP, Burns WP, Fong SS, Gowen CM, Roberts SB. Gap detection for genome-scale constraint-based models. *Adv Bioinf* 2012;2012.
- [23] Jensen CS, Norsigian CJ, Fang X, Nielsen XC, Christensen JJ, Palsson BO, Monk JM. Reconstruction and validation of a genome-scale metabolic model of streptococcus oralis (iCJ415), a human commensal and opportunistic pathogen. *Front Genet* 2020;11:116.
- [24] de Arroyo Garcia L, Jones PR. In silico co-factor balance estimation using constraint-based modelling informs metabolic engineering in escherichia coli. *PLoS Comput Biol* 2020;16(8):e1008125.
- [25] Reed JL. Genome-scale metabolic modeling and its application to microbial communities. In *The Chemistry of Microbiomes: Proceedings of a Seminar Series*, National Academies Press (US); 2017. pp. 85–92..
- [26] Orth JD, Thiele I, Palsson BØ. What is flux balance analysis?. *Nat Biotechnol* 2010;28(3):245–8.
- [27] Volkova S, Matos MR, Mattanovich M, Marin de Mas I. Metabolic modelling as a framework for metabolomics data integration and analysis. *Metabolites* 2020;10(8):303.
- [28] Lularevic M, Racher AJ, Jaques C, Kiparissides A. Improving the accuracy of flux balance analysis through the implementation of carbon availability constraints for intracellular reactions. *Biotechnol Bioeng* 2019;116(9):2339–52.
- [29] Lachance J-C, Lloyd CJ, Monk JM, Yang L, Sastry AV, Seif Y, Palsson BO, Rodrigue S, Feist AM, King ZA, et al. Bofdat: Generating biomass objective functions for genome-scale metabolic models from experimental data. *PLoS Comput Biol* 2019;15(4):e1006971.
- [30] Budinich M, Bourdon J, Larhlmi A, Eveillard D. A multi-objective constraint-based approach for modeling genome-scale microbial ecosystems. *PLoS One* 2017;12(2):e0171744.
- [31] Patané A, Jansen G, Conca P, Carapezza G, Costanza J, Nicosia G. Multi-objective optimization of genome-scale metabolic models: the case of ethanol production. *Ann Oper Res* 2019;276(1–2):211–27.
- [32] Edwards JS, Palsson BO. Metabolic flux balance analysis and the in silico analysis of Escherichia coli K-12 gene deletions. *BMC Bioinf* 2000;1(1):1–10.
- [33] Lewis NE, Hixson KK, Conrad TM, Lerman JA, Charusanti P, Polpitiya AD, Adkins JN, Schramm G, Purvine SO, Lopez-Ferrer D, et al. Omic data from evolved E. coli are consistent with computed optimal growth from genome-scale models. *Mol Syst Biol* 2010;6(1):390.

- [34] Gudmundsson S, Thiele I. Computationally efficient flux variability analysis. *BMC Bioinf* 2010;11(1):1–3.
- [35] Lewis NE, Nagarajan H, Palsson BO. Constraining the metabolic genotype–phenotype relationship using a phylogeny of in silico methods. *Nat Rev Microbiol* 2012;10(4):291–305.
- [36] Megchelenbrink W, Huynen M, Marchiori E. optgpsampler: an improved tool for uniformly sampling the solution-space of genome-scale metabolic networks. *PLoS one* 2014;9(2):e86587.
- [37] Herrmann HA, Dyson BC, Vass L, Johnson GN, Schwartz J-M. Flux sampling is a powerful tool to study metabolism under changing environmental conditions. *NPJ Syst Biol Appl* 2019;5(1):1–8.
- [38] Robaina Estévez S, Nikoloski Z. Generalized framework for context-specific metabolic model extraction methods. *Front Plant Sci* 2014;5:491.
- [39] Becker SA, Palsson BO. Context-specific metabolic networks are consistent with experiments. *PLoS Comput Biol* 2008;4(5):e1000082.
- [40] Kim MK, Lun DS. Methods for integration of transcriptomic data in genome-scale metabolic models. *Comput Struct Biotechnol J* 2014;11(18):59–65.
- [41] Siritwach R, Matsuda F, Yano K, Hirai MY. Drought stress responses in context-specific genome-scale metabolic models of *Arabidopsis thaliana*. *Metabolites* 2020;10(4):159.
- [42] Zur H, Ruppín E, Shlomi T. iMAT: an integrative metabolic analysis tool. *Bioinformatics* 2010;26(24):3140–2.
- [43] Vlassis N, Pacheco MP, Sauter T. Fast reconstruction of compact context-specific metabolic network models. *PLoS Comput Biol* 2014;10(1):e1003424.
- [44] Agren R, Bordel S, Mardinoglu A, Pornputtapong N, Nookaew I, Nielsen J. Reconstruction of genome-scale active metabolic networks for 69 human cell types and 16 cancer types using INIT. *PLoS Comput Biol* 2012;8(5):e1002518.
- [45] Agren R, Mardinoglu A, Asplund A, Kampf C, Uhlen M, Nielsen J. Identification of anticancer drugs for hepatocellular carcinoma through personalized genome-scale metabolic modeling. *Mol Syst Biol* 2014;10(3):721.
- [46] Lorbek G, Perše M, Jeruc J, Juvan P, Gutierrez-Mariscal FM, Lewinska M, Gebhardt R, Horvat S, Björkhem I, Rozman D, et al. Lessons from hepatocyte-specific Cyp51 knockout mice: impaired cholesterol synthesis leads to oval cell-driven liver injury. *Scientific Rep* 2015;5(1):1–11.
- [47] Khodae S, Asgari Y, Totonchi M, Karimi-Jafari MH. imm1865: a new reconstruction of mouse genome-scale metabolic model. *Scientific Rep* 2020;10(1):1–13.
- [48] Chung NC, Miasojedow B, Startek M, Gambin A. Jaccard/Tanimoto similarity test and estimation methods for biological presence-absence data. *BMC Bioinf* 2019;20(15):1–11.
- [49] Van der Maaten L, Hinton G. Visualizing data using t-SNE. *J Mach Learn Res* 2008;9(11).
- [50] Fallahi S, Skaug HJ, Alendal G. A comparison of monte carlo sampling methods for metabolic network models. *PLoS one* 2020;15(7):e0235393.
- [51] Becker SA, Feist AM, Mo ML, Hannum G, Palsson BØ, Herrgard MJ. Quantitative prediction of cellular metabolism with constraint-based models: the cobra toolbox. *Nat Protocols* 2007;2(3):727–38.
- [52] Heirendt L, Arreckx S, Pfau T, Mendoza SN, Richelle A, Heinken A, Haraldsdóttir HS, Wachowiak J, Keating SM, Vlasov V, et al. Creation and analysis of biochemical constraint-based models using the cobra toolbox v. 3.0. *Nat Protocols* 2019;14(3):639–702.
- [53] Wang H, Marčišauskas S, Sánchez BJ, Domenzain I, Hermansson D, Agren R, Nielsen J, Kerkhoven EJ. Raven 2.0: A versatile toolbox for metabolic network reconstruction and a case study on *Streptomyces coelicolor*. *PLoS Comput Biol* 2018;14(10):e1006541.
- [54] Gurobi Optimization, LLC. Gurobi optimizer reference manual; 2021. <http://www.gurobi.com>.
- [55] Urlep Z, Lorbek G, Perše M, Jeruc J, Juvan P, Matz-Soja M, Gebhardt R, Björkhem I, Hall JA, Bonneau R, et al. Disrupting hepatocyte Cyp51 from cholesterol synthesis leads to progressive liver injury in the developing mouse and decreases rosc signalling. *Scientific Rep* 2017;7(1):1–13.
- [56] Cokan KB, Urlep Z, Lorbek G, Matz-Soja M, Skubic C, Perše M, Jeruc J, Juvan P, Režen T, Rozman D. Chronic disruption of the late cholesterol synthesis leads to female-prevalent liver cancer. *Cancers* 2020;12(11):3302.
- [57] Skubic C, Rozman D. Sterols from the post-lanosterol part of cholesterol synthesis: novel signaling players. In *Mammalian Sterols*. Springer; 2020. pp. 1–22..
- [58] Schultz A, Qutub AA. Reconstruction of tissue-specific metabolic networks using CORDA. *PLoS Comput Biol* 2016;12(3):e1004808.
- [59] Wang Y, Eddy JA, Price ND. Reconstruction of genome-scale metabolic models for 126 human tissues using mCADRE. *BMC Syst Biol* 2012;6(1):153.
- [60] Jerby L, Shlomi T, Ruppín E. Computational reconstruction of tissue-specific metabolic models: application to human liver metabolism. *Mol Syst Biol* 2010;6(1):401.
- [61] Poupin N, Vinson F, Moreau A, Batut A, Chazalviel M, Colsch B, Fouillen L, Guez S, Khoury S, Dalloux-Chioccioli J, et al. Improving lipid mapping in genome scale metabolic networks using ontologies. *Metabolomics* 2020;16(4):1–11.