

modelBorgifier Manual

J. T. Sauls and J. M. Buescher

December 2, 2013

BRAIN Aktiengesellschaft

Microbial Production Technologies

Platform for Quantitative Biology and Sequencing

Darmstaeter Str. 34-36, 64673 Zwingenberg, Germany

www.brain-biotech.de

jrb@brain-biotech.de

“We will add your biological and technological distinctiveness to our own. Resistance is futile.”
 – Borg hail.

Contents

1	Introduction	3
2	Set-up	3
2.1	Installation	3
2.2	Dependencies	3
2.3	Test script	3
3	Procedure	4
3.1	Load the comparison model.	4
3.1.1	Verify format	4
3.1.2	Add information from the SEED database	5
3.1.3	Verify model holds flux	6
3.2	Load the template model	6
3.2.1	Load Tmodel from SBML	6
3.2.2	Load Tmodel from .mat file	6
3.3	Compare models	6
3.4	Match reactions and metabolites	7
3.4.1	Optimizing scoring parameters	8
3.4.2	Auto-matching based on scores	9
3.4.3	Matching flow	9
3.5	Merge model with database	10
3.6	Extracting models	11
4	Appendix	11
4.1	Description of scoring parameters	11
4.2	Description of scripts	12
4.3	Important terms	12

1 Introduction

When using this tool, please reference

J.T. Sauls and J.M. Buescher, Assimilating genome-scale metabolic reconstructions with modelBorgifier
submitted

modelBorgifier is published under Creative Commons BY-NC-SA

2 Set-up

2.1 Installation

Assuming that the openCOBRA Toolbox is installed and functional on your system, installation only requires that you add the modelBorgifier directory to the Matlab path.

2.2 Dependencies

Before you start, please download and install these dependencies. We suggest you locate them in the directory */modelBorgifier/dependencies* together with their respective accompanying licenses. The linear and exponential optimization techniques use functions available through the Matlab Optimization Toolbox (non-free).

- modelBorgifier relies on the openCOBRA Toolbox (<http://opencobra.sourceforge.net/opencOBRA/Welcome.html>) [6] for reading and writing functions (themselves dependent on the Matlab sbmltoolbox, included with the standard openCOBRA install), and any FBA analysis.
- Reading CSV files utilizes the `csv2cell.m` by David Groppe and `csvimport.m` by Ashish Sadanandan functions from MatlabCentral (<http://www.mathworks.com/matlabcentral/fileexchange/20836-csv2cell> and <http://www.mathworks.com/matlabcentral/fileexchange/23573-csvimport>).
- For machine learning, the LIBSVM library[2] (<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>) is needed for SVM optimization.
- The script *findjob.m* by Yair Altman (<http://www.mathworks.com/matlabcentral/fileexchange/14317-findjob-find-java-handles-of-matlab-graphic-objects>) is used for GUI presentation.

2.3 Test script

The script *driveModelBorgifier.m* demonstrates the order of use and implementation of the functions in the modelBorgifier suite. It is designed to be a guide when adding a new model and uses the two popular models iJO1366 and iBSU1103 as an for example comparison and merging. It also utilizes the open databases from the Model SEED to add additional annotation information based on SEED IDs. To use the test suite as is, download both the above models and place the SBML files in the */modelBorgifier/test* directory.

- iJO1366[4] is a popular *E. coli* model which can be downloaded in SBML format from <http://www.nature.com/msb/journal/v7/n1/extref/msb201165-s3.xml> (paper: <http://www.nature.com/msb/journal/v7/n1/full/msb201165.html>; license <http://creativecommons.org/licenses/by-nc-sa/3.0/> note this is a non-commercial license).
- iBSU1103[3] is a *B. subtilis* model which uses IDs consistent with the Model SEED (model in SBML: <http://genomebiology.com/content/supplementary/gb-2009-10-6-r69-s4.xml>; paper: <http://genomebiology.com/2009/10/6/R69>; license: <http://creativecommons.org/licenses/by/2.0/>).
- The Model SEED[5] (<http://seed-viewer.theseed.org/>) is a collection of stoichiometric models, and the way to access models created by the RAST[1] annotation server. Models in the Model SEED draw from two databases, for reactions (seed-viewer.theseed.org/ModelSEEDdownload.cgi?biochemistry=1) and metabolites (seed-viewer.theseed.org/ModelSEEDdownload.cgi?biochemCompounds=1). These databases are open and freely available. For use in modelBorgifier, first convert them to CSV files, semicolon delimited. We suggest you place them in */modelBorgifier/test/SEED_db* with the appropriate license.

3 Procedure

To quickly demo the modelBorgifier, simply follow *driveModelBorgifier.m*. Additionally, all scripts have documentation located in the .m file that can be accessed via Matlab's help command. The following is a description of the comparison and modeling process.

3.1 Load the comparison model.

Load the model you wish to compare (Cmodel) and merge with the template model (Tmodel), from its current format.

If the model is already in openCOBRA compatible SBML format, then use the openCOBRA function *readCbmodel.m*.

```
>> Cmodel = readCbModel(fileName);
```

If the model is in a less readily machine readable format, such as .xls, then use an appropriate custom written script (converting .xls to .csv may be necessary if you are not on a Windows machine). Some model may be provided in SBML format, but an additional file contains useful information. If so, consider using using a script to extract information from this additional file(s).

3.1.1 Verify format

Use *verifyModel.m* to ensure that Cmodel has all relevant information and can be successfully compared against Tmodel for already existing reactions and metabolites.

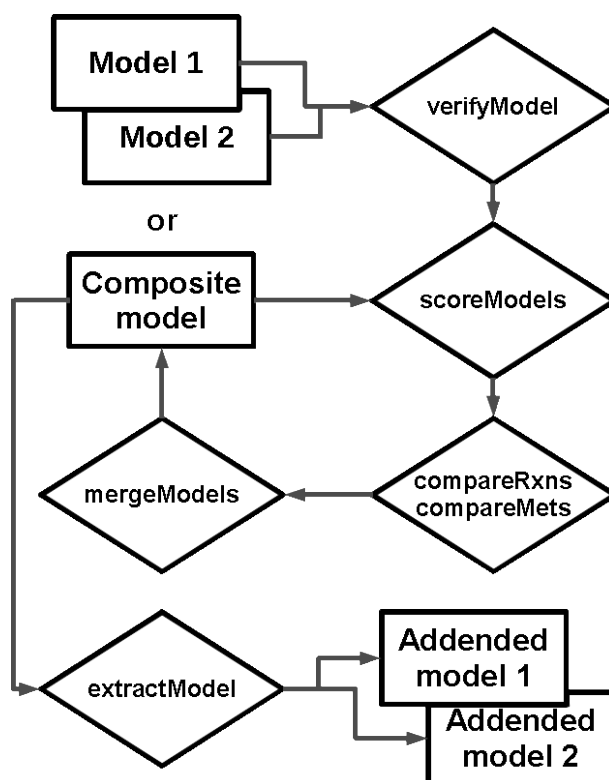


Figure 1: Program flow.

```
>> Cmodel = verifyModel(Cmodel);
```

Verify model ensures all required information arrays are present, that all reactions are forward or reversible (not reverse), and that all reaction and metabolite names are unique. If reaction names are not unique, *verifyModel.m* will prompt you to put in new names, or you can use *checkCobraNamesUnique.m* from the openCOBRA toolbox to systematically rename all duplicate reactions. Finally, *verifyModel.m* formats reaction and metabolite names so they are lower case and do not contain obtuse characters, gives the fields a conserved order, and makes sure the model has a name.

3.1.2 Add information from the SEED database

Optional. If your model either came from the Model SEED or uses SEED identifiers (reaction names in the form rxnXXXXXX and metabolites in the form cpdXXXXXX), then the script *addSEEDInfo.m* uses those identifiers too add information to the model via the databases mentioned above.

```
>> Cmodel = addSEEDInfo(Cmodel, rxnDbFileName, cpdDbFileName);
```

This greatly aids comparison and matching, especially when the template model does not reference the SEED IDs otherwise. It also replaces the metabolite names with the appropriate

abbreviations in the reaction equations. The databases should be in semicolon delimited CSV format.

3.1.3 Verify model holds flux

Optional. Run a simple FBA simulation to see if the model holds flux. Make sure there is an objective function set, and that the bounds are reasonable. If the model doesn't hold flux before it is added to merged, than it is more difficult to determine if it was added correctly afterwards.

```
>> solverOkay = changeCobraSolver('glpk','LP');  
>> CmodelSol = optimizeCbModel(Cmodel);
```

3.2 Load the template model

For the very first model assimilation, use a regular genome scale metabolic reconstruction of your choice as template model. For subsequent assimilations, we suggest you use the composite model generated in previous model assimilations as template model.

3.2.1 Load Tmodel from SBML

In this case, load and verify the model as above, optionally adding information from the SEED database and checking if the model holds flux. Then, convert the model to an appropriate template for comparison using *buildTmodel.m*.

```
>> Tmodel = buildTmodel(Tmodel);
```

3.2.2 Load Tmodel from .mat file

If previous models have been merged, then this composite model can be used as Tmodel. Simply load the Matlab workspace (.mat file) containing the composite model.

3.3 Compare models

Compare Cmodel to Tmodel by running the script *compareCbModels.m*.

```
>> [Cmodel,Tmodel,score,Stats] = compareCbModels(Cmodel,Tmodel);
```

This script first assembles additional information arrays in both models for the purpose of comparison, returning the information with the original models. It then compares Cmodel to Tmodel, pairwise by reaction, allocating a normalized score based on similarity. This may take ~4 hours for two models of around 2000 reactions each.

These scores are held in the 3D array *score*, which is of size [number of reactions in Cmodel x number of reactions in Tmodel x number of scoring parameters]. Stats contains information on best matches.

3.4 Match reactions and metabolites

Pair reactions and metabolites from Cmodel to Tmodel using *reactionCompare.m*. Matching is initiated with the following command.

```
>> [rxnList, metList, Stats] = reactionCompare(Cmodel, Tmodel, score);
```

reactionCompare.m calls a GUI that allows the user to choose a match for a given reaction in Cmodel, and also to match the metabolites for that reaction. Note that the best matches are presented (default 3, more can be viewed).

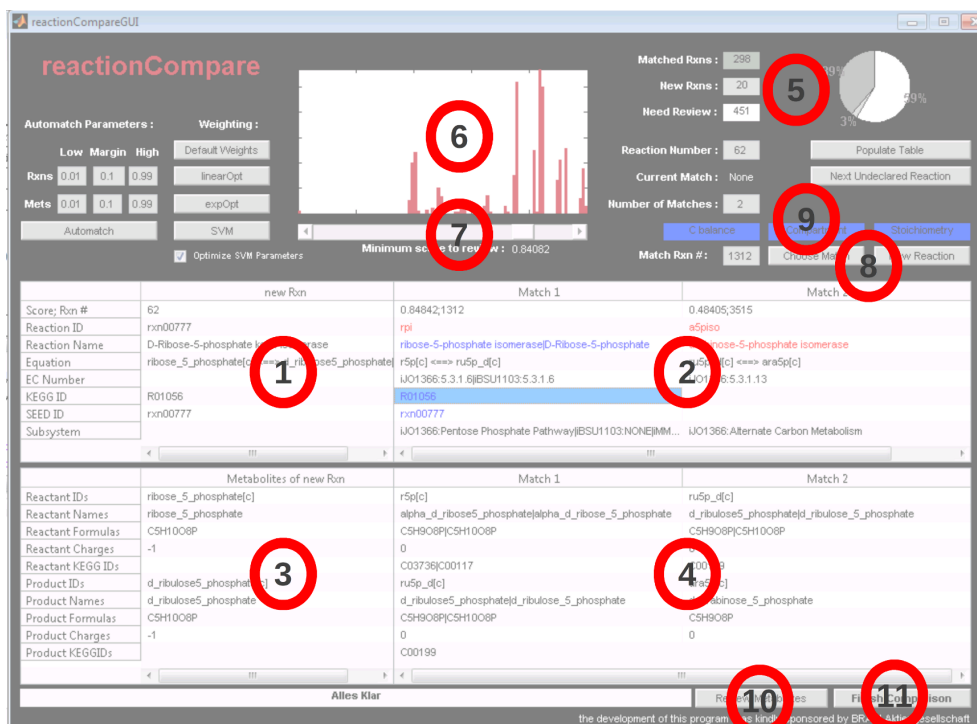


Figure 2: reactionCompare GUI displays details for the next reaction from Cmodel to be matched (1) and details for the highest scoring reactions in Tmodel (2). Likewise, details for the metabolites of this reaction in Cmodel are displayed (3) next to the metabolites of the corresponding reactions in Tmodel (4). The current status of the progress of the matching is given in numbers and in a pie chart (5). A histogram of the scores of the best matching is displayed (6) above a slider (7) that lets the user select the minimum score of the reactions to be reviewed. The matching reaction from Tmodel is selected by clicking on the reaction in (2) and clicking the button “Choose Match” in (8). If no match is available, the reaction can be declared new with the button “New Reaction”. To facilitate the review process, common mistakes are indicated by red fields in (9), while passed checks are blue. If the check cannot be performed, the fields remain white. If metabolites remain to be reviewed, this can be started by clicking “Review Metabolites” (10). A click on the button “Finish Comparison” (11) ends the matching (e.g. to save the workspace).

reactionCompare.m gives auto-matching options described below. *reactionCompare.m* outputs rxnList and metList, which designate reactions and metabolites in Cmodel to the index of their match in Tmodel, or indicate them as new or need review. Stats contains weighting information in addition to best matches. Matching can be suspended at anytime, so subsequent calls to *reactionCompare.m* should include these output arguments.

```
>> [rxnList, metList, Stats] = reactionCompare(Cmodel, Tmodel, score, ...
rxnList, metList, Stats);
```

3.4.1 Optimizing scoring parameters

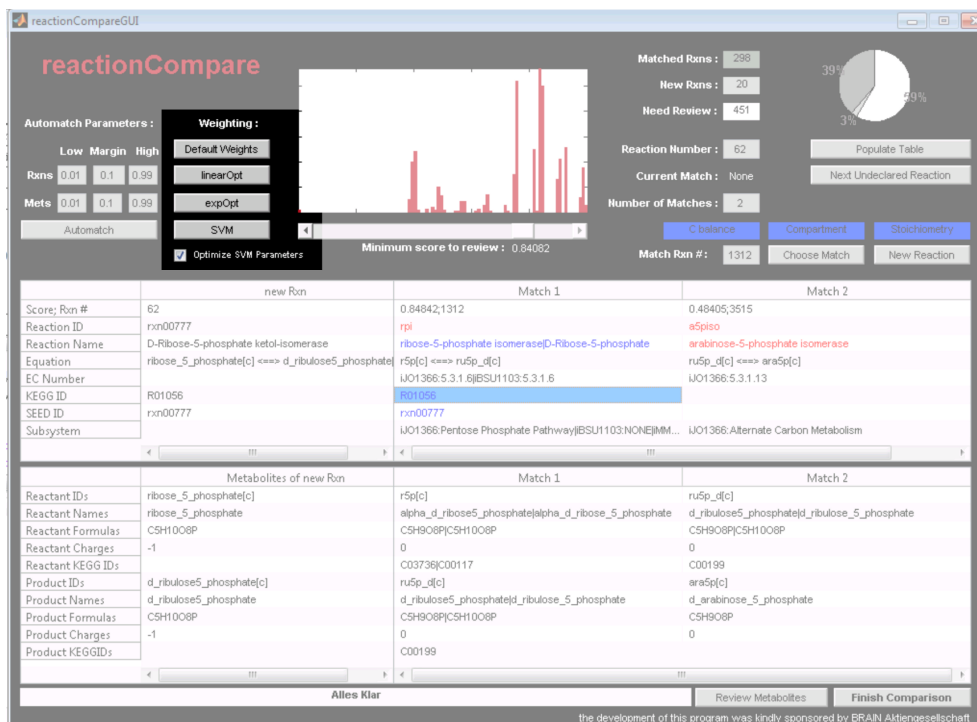


Figure 3: Select optimization of scoring parameters in highlighted area

The initial comparison between models is based on forty parameters which are arbitrarily weighted. The parameters cover a large swath of information available in genome-scale reconstructions, but do not necessarily emphasize the relevant information in the models you have at hand. For this reason it is important to optimize the weighting of the scoring parameters. This is accomplished through the *reactionCompare.m* GUI via one of three methods; support vector machines (SVM), linear optimization, or exponential optimization. It is recommended to check the 'optimize parameters' box in the GUI when using SVM.

All three methods use previously declared reactions as a training set to determine which parameters are pertinent, and weighs them accordingly. For example, if KEGG IDs are not available in one of the models, then that corresponding parameter will be down weighted (most likely to zero), as it holds no bearing on the matching process. Likewise, if reaction EC numbers prove to be particularly good at predicting matches, then it will be more heavily weighted.

3.4.2 Auto-matching based on scores

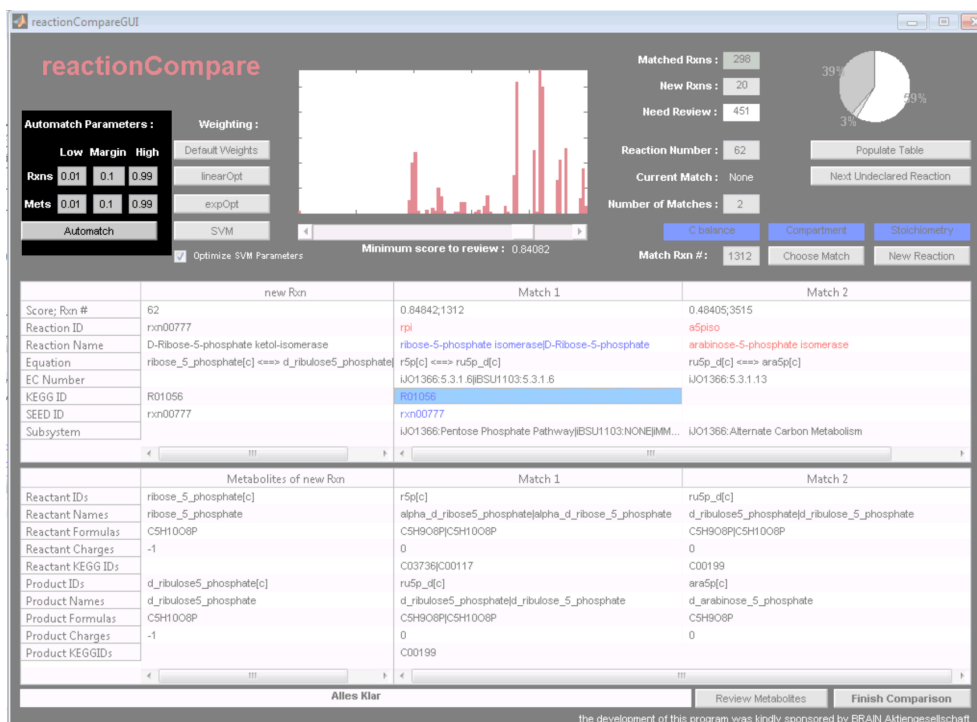


Figure 4: Set cutoffs for automated matching and start automatching in highlighted area

Optimizing the weight of the scoring parameters helps widen the scoring distance between probable matches and probable new reactions and metabolites. Thus, auto-matching based on cutoff scores can be done more confidently after optimization. Simply set a low cutoff to automatically declare new reactions when no match exists with a score above a certain value. The same applies for a high cutoff for automatic pairing. Additionally, a margin can be set, such that a reaction will only be automatically paired with a match, when the score for the match is better than the second best match by a certain margin. Auto-matching in the same manner exist for metabolites.

As metabolites become declared, the status of the reactions which contain those metabolites can be elucidated. For example, if all of a reaction's metabolites have matches in the template model, and the template model also contains a reaction with those exact metabolites, then the original reaction is automatically paired with the one from the template. Conversely, if a reaction contains a metabolite that does not exist in the template model, then it is declared as new.

3.4.3 Matching flow

The following is a general guide to effective and efficient matching.

1. Consider a small subset of reactions, ie 20-50. You can choose to consider reactions with at least one match with a score above a certain cutoff by adjusting the slider under the score frequency histogram in the GUI. Note that comparing high scoring reactions (likely matches) is more helpful towards the generation of a training data set.

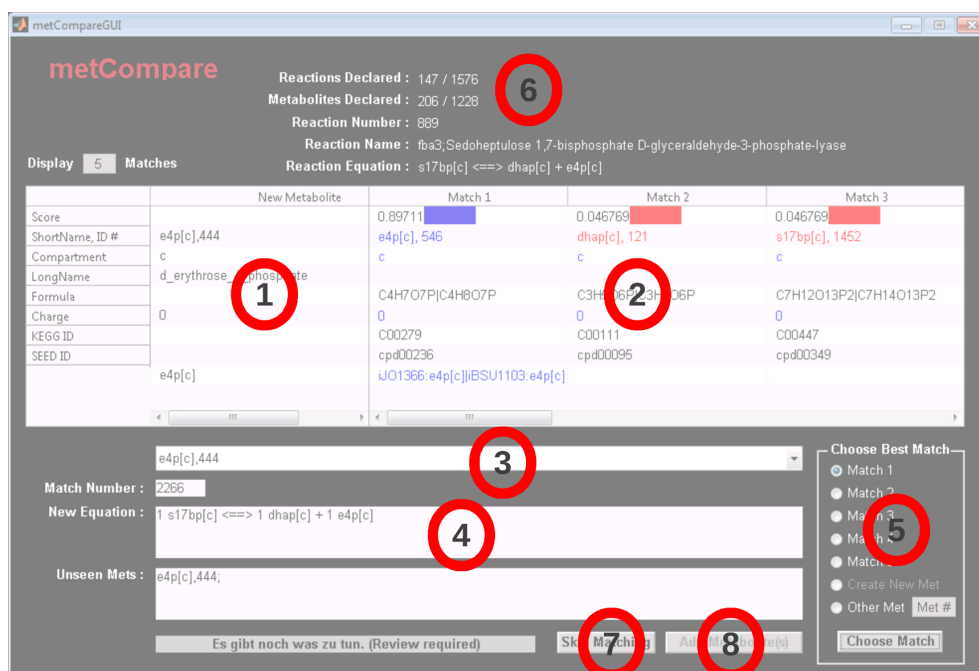


Figure 5: The metCompare GUI is called viay the reactionCompare GUI. It displays the available details on the metabolite from Cmodel (1) and details on the highest scoring metabolites from Tmodel. The other metabolites from the same reaction can also be chosen (3) and the corresponding reaction equation is displayed (4). The matching metabolite can be selected by clicking it in (2) or by selecting it in (5). Clicking the button “Choose Match” finalizes the selection. The current progress of the matching is also displayed (6). The matching of metabolites can be skipped at any time (7) to return to the reactionCompare GUI. If all metabolites of the reaction have been matched or declared new, the GUI can be closed by clicking “Add Metabolite(s)” (8).

2. Declare these reactions as either having a match or new.
3. After each reaction is declared, a sub-GUI that allows the user to compare metabolites will start. Similarly, declare metabolites as having a match or new. Reactions that are given a match must also have matches for all of their metabolites.
4. Use one of the optimizing functions to partition scores, then use the auto-match function to declare for-sure matches and for-sure new reactions. Adjust the auto-match parameters based on the scores you observe. The metabolite comparison GUI will most likely re-launch for these newly declared reactions.
5. Repeat steps 1-4 until all reactions and metabolites have matches in Tmodel or are declared as new.

3.5 Merge model with database

Use *mergeModels.m* to merge Cmodel with Tmodel. It returns the composite model, TmodelC, and Cmodel as extracted from the composite model. Stats now additionally contains comparison information between models in the composite model.

```
>> [TmodelC, Cspawn, Stats] = mergeModels(Cmodel, Tmodel, rxnList, metList, ...  
Stats);
```

For reactions and metabolites for which matches have been found, the script adds information to the current entries in Tmodel. New entries are made for new reactions and metabolites from Cmodel. *mergeModels.m* calls *cleanToModel.m*, which may prompt the user to clarify duplicate reaction and metabolite names to ensure the composite model only contains unique entries.

Merging is also done in a manner that ensures the original models can be retrieved in their original mathematical form. This may require that some reactions and metabolites be re-reviewed. Note that reaction bounds, reversibility of reactions, and the index of the objective function are stored individually for all models. Furthermore, protons and water are ignored when checking for stoichiometric fidelity, because these two metabolites are notoriously not balanced in some published metabolic models.

3.6 Extracting models

Models merged together or into an existing composite model can be later retrieved with *readCbTmodel.m*. This reproduces the initial model with additional annotation information.

```
>> Cspawn = readCbTmodel('modelName', Tmodel);
```

The composite model can also be flattened and treated as a stand alone model with *squishTmodel.m*.

4 Appendix

4.1 Description of scoring parameters

The reaction scoring parameters are defined in *compareCbModels.m*. Most parameters come in a pair, one allotting points for a match, the other removing points for a miss between any to given reactions. String comparison supports multiple pieces of information contained within one field as separated by a pipe ('|').

rxnName	Reaction names match.
rxnNameL	Reaction long name match. Also cross checks to ID.
ec	Reaction EC number match.
rxnKEGG	Reaction KEGG ID match.
rev	Reaction reversibility match.
sub	Reaction subsystem match.
gr	Genes match.
metNum	Same number of metabolites.
metSto	Same stoichiometry.
metName	Metabolites have same name. Points allotted per metabolite.

metForm	Metabolites have same formula. Points allotted per metabolite.
metKEGG	Metabolites have same formula. Points allotted per metabolite.
metSEED	Metabolites have same SEED ID. Points allotted per metabolite.
reacNum	Same number of reactants.
reacSto	Same reactant stoichiometry.
prodNum	Same number of products.
prodSto	Same product stoichiometry.
rxnComp	Reactions involve the same compartments.
rxnSEED	Reaction SEED ID match.
metBonus	Bonus allotted if all the metabolite names, formulas, KEGG or SEED IDs match.
rxnNet	Network similarity around reaction.

Metabolites are also scored individually on the following parameters: ID, compartment, name, formula, charge, KEGGID, SEEDID, ChEBIID, PubChemID, and InChIString.

4.2 Description of scripts

All scripts have usage descriptions located within the .m file itself. Scripts in */modelBorgifier/src* are core to the comparison process, while those in */modelBorgifier/tools* are mostly concerned with formatting data to aid comparison.

4.3 Important terms

Cmodel	Compare Model. The model that is to be compared and merged to Tmodel.
Tmodel	Template Model. The model to which Cmodel is compared and merged. This can simply be another model, or a composite model of many previously combined models.
rxnList	The array which pairs reactions from Cmodel to their match in Tmodel, or declares them as new.
metList	Analogous array for metabolites.

References

- [1] Aziz, R. K. et al. The RAST Server: rapid annotations using subsystems technology. BMC genomics 9, 75 (2008).
- [2] C.-C. Chang and C.-J. Lin. LIBSVM : a library for support vector machines. ACM Transactions on Intelligent Systems and Technology, 2:27:1–27:27 (2011).
- [3] Henry, C. S., Zinner, J. F., Cohoon, M. P. & Stevens, R. L. iBsu1103: a new genome-scale metabolic model of Bacillus subtilis based on SEED annotations. Genome biology 10, R69 (2009).

- [4] Orth, J. D. et al. A comprehensive genome-scale reconstruction of *Escherichia coli* metabolism—2011. *Molecular Systems Biology* 7, (2011).
- [5] Overbeek, R. et al. The subsystems approach to genome annotation and its use in the project to annotate 1000 genomes. *Nucleic acids research* 33, 5691–702 (2005).
- [6] Schellenberger, J. et al. Quantitative prediction of cellular metabolism with constraint-based models: the COBRA Toolbox v2.0. *Nature Protocols* 6, 1290–1307 (2011).