Model SEED Tutorial Part 1: Technical Primer

Christopher Henry, Scott Devoid, Matt DeJongh, Aaron Best, Ross Overbeek, and Rick Stevens

Presented by: Christopher Henry

August 31-September 2



Fellowship for Interpretation of Genomes







Predicting Phenotypes from Genotypes — the prediction of system level behavior from collections of functional components



Metabolic Modeling is One Key to Predicting Phenotype from Genotype

- What is a metabolic model?
- 1.) A list of all reactions involved in the metabolic pathways
- 2.) A list of rules associating reaction activity to gene activity
- 3.) A biomass reaction listing essential building blocks needed for growth and Gene B division Gene A **Function** Function Nutrients Enzyme OLIC PATHWAYS Amino acids Nucleotides ►Lipids Cofactors **Biomass** Cell walls Energy

Metabolic Modeling is One Key to Predicting Phenotype from Genotype

- What can a metabolic model do?
- 1.) Predict culture conditions and possible responses to environment changes.
- 2.) Predict metabolic capabilities from genotype.
- 3.) Predict impact of genetic perturbations



Metabolic Modeling is One Key to Predicting Phenotype from Genotype

- What can a metabolic model do?
- 1.) Predict culture conditions and possible responses to environment changes.
- 2.) Predict metabolic capabilities from genotype.
- 3.) Predict impact of genetic perturbations
- 4.) Linking annotations to observed organism behavior enabling validation and correction of annotations



Model reconstruction lags behind genome sequencing



- •≈1000 completely sequenced prokaryotes vs ≈30 published genome-scale models
- •Models are often constructed one-at-a-time by individuals working independently
- •Model building typically begins by identifying bidirectional best hits with *E. coli*
- •Current process results in replication of work, propagation of errors, and extensive manual curation
- •Bottom line: it previously required approximately one year to produce a complete model



Biochemistry Database in the SEED

•A biochemistry database was constructed combining content from the **KEGG** and **13** published genome-scale models into a non-redundant set of compounds and reactions



•Reactions were then mapped to the functional roles in the SEED based on EC number, substrate names, and enzyme names:



Biomass Objective Function

•To test growth of the model, we build a biomass objective function template





•Each biomass component may be rejected from the biomass reaction of a model based on the following criteria:

Subsystem representation

•Functional role presence

Taxonomy

•Cell wall types



Flux Balance Analysis



Flux Balance Analysis



Genome Annotations Contain Knowledge Gaps



Flux Balance Analysis



Model Auto-completion Optimization



Genome Annotation: the Subsystems Approach





Accuracy Before Optimization

Biolog phenotype data

•SEED models were used to predict the output of 14 biolog phenotyping arrays

•Average accuracy: 60%

Essentiality data

•SEED models were used to predict essential genes for 14 experimental gene essentiality datasets

•Average accuracy: 72%

Overall accuracy: 66%





Biolog Consistency Analysis

Biolog phenotype data

•Add transporters for Biolog nutrients if missing from models

•69 transporters added to each model on average

•Average accuracy: 70%

Essentiality data

•Accuracy unchanged: 72%

Overall accuracy: 71%





Annotation Consistency Analysis





Model Optimization: Gap Filling





•Fix false negative predictions by adding reactions to models

Biolog accuracy

•Average accuracy: 83%

Essentiality accuracy

•Average accuracy: 81%

Overall accuracy: 82%





Model Optimization: Gap Generation



•Fix false positive predictions by removing reactions from models

Biolog accuracy

•Average accuracy: 88%

Essentiality accuracy

•Average accuracy: 85%

Overall accuracy: 87%





Seed Models vs Published Models

•Single-genome Seed models compare favorably with published single genome models

Organism name	Published model	Published reactions	SEED Reactions	Published genes	SEED genes
Acinetobacter	iAbaylyiv4	868	1196	775	785
B. subtilis	iYO844	1020	1463	844	1041
C. acetobutylicum	iJL432	502	989	432	721
E. coli	iAF1260	2013	1529	1261	1083
G. sulfurreducens	iRM588	523	721	588	468
H. influenzae	iCS400	461	969	400	575
H. pylori	iIT341	476	731	341	421
L. plantarum	iBT721	643	908	721	699
L. lactis	iAO358	621	965	358	646
M. succiniciproducens	iTK425	686	1048	425	659
M. tuberculosis	iNJ661	939	1021	661	728
M. genitalium	iPS189	264	294	189	214
N. meningitidis	iGB555	496	903	555	560
P. gingivalis	iVM679	679	744	0*	399
P. aeruginosa	iMO1056	883	1386	1056	1094
P. putida	iNJ746	950	1261	746	1053
R. etli	iOR363	387	1264	363	1242
S. aureus	iSB619	641	1115	619	770
S. coelicolor	iIB700	700	1159	700	987

Assessing Subsystem Annotations From Auto-completion

•We identify how *complete* the annotations are for each of the Seed subsystems by calculating the following ratio:

auto-completion reactions in subsystem	_ Fraction of subsystem reactions with
total reactions in subsystem	= missing genes

•Highest scoring subsystems:

•Cell Wall and Capsule Biosynthesis (15%)

•21 reactions per model added during auto-completion

- •LOS Core Oligosaccharide Biosynthesis (Gram negative)
- Teichoic and Lipoteichoic Acids Biosynthesis (Gram positive)
- •KDO2-Lipid A Biosynthesis

•Cofactors, Vitamins, and Prosthetic Group Biosynthesis (5%)

- •10 reaction per model added during auto-completion
- •Ubiquinone Biosynthesis
- •Menaquinone and Phylloquinone Biosynthesis
- •Thiamin Biosynthesis

•Six subsystems account for 31/56 reactions added to each model during the autocompletion process

Model statistics across the phylogenetic tree



Reaction Activity Across All Models



Essential Genes Across All Models



Essential Nutrients Across All Models



Acknowledgements

ANL/U. Chicago Team

- Robert Olson
- Terry Disz
- Daniela Bartels
- Tobias Paczian
- Daniel Paarmann
- Scott Devoid
- Andreas Wilke
- Bill Mihalo
- Elizabeth Glass
- Folker Meyer
- Jared Wilkening
- Rick Stevens
- Alex Rodriguez
- Mark D'Souza
- Rob Edwards
- Christopher Henry

FIG Team

- Ross Overbeek
- Gordon Pusch
- Bruce Parello
- Veronika Vonstein
- Andrei Ostermann
- Olga Vassieva
- Olga Zagnitzko
- Svetlana Gerdes

Hope College Team

- Aaron Best
- Matt DeJongh
- Nathan Tintle
- Hope college students



National Institute of Allergy and Infectious Diseases National Institutes of Health

www.theseed.org



Fellowship for Interpretation of Genomes