

# Metabolic Network Modeling of *Clostridium thermocellum* for Systems Biology and Metabolic Engineering



R. ADAM THOMPSON<sup>2,3</sup>, SERGIO GARCIA<sup>1,2</sup>, SANJEEV DAHAL<sup>2,4</sup>, INTAWAT NOOKAEW<sup>2,4</sup>, DONOVAN S. LAYTON<sup>1,2</sup>, ADAM M. GUSS<sup>2,3,5</sup>, DAN G. OLSON<sup>2,6</sup>, LEE R. LYND<sup>2,6</sup>, and CONG T. TRINH<sup>1,2,3</sup>

<sup>1</sup>Department of Chemical and Biomolecular Engineering, University of Tennessee; Knoxville, Tennessee; <sup>2</sup>BioEnergy Science Center, Oak Ridge, Tennessee; <sup>3</sup>Bredesen Center for Interdisciplinary Research and Education; <sup>4</sup>Comparative Genomics Group, Biosciences Division, Oak Ridge National Laboratory; <sup>5</sup>Biosciences Division, Oak Ridge National Laboratory, Oak Ridge, TN, USA; <sup>6</sup>Thayer School of Engineering, Dartmouth College, Hanover, NH, USA

**Abstract:** *Clostridium thermocellum* is a gram-positive thermophile that can directly convert lignocellulosic material into commercially relevant chemicals such as biofuels. Its metabolism contains many branches and redundancies, which limit the production of biofuels at industrially relevant yields and titers. In order to guide the experimental efforts required to overcome these barriers, we built two models of *C. thermocellum* metabolism. Through an extensive literature review, we first constructed a model of the core metabolism of *C. thermocellum*. This model was experimentally validated and served to investigate the range of phenotypes of *C. thermocellum* in response to significant perturbation of energy and redox pathways. The results revealed a complex, robust redox metabolism of *C. thermocellum*. By incorporating experimental data into this core model, we identified redox bottlenecks hindering high-yield ethanol production in *C. thermocellum*. With the recently published sequence of a genetically-tractable strain *C. thermocellum* DSM 1313, the KEGG database as a scaffold, and further literature review, we expanded the core model into a genome scale model (iAT601). This model constitutes a knowledge base for the organism, including detailed metabolic information, as well as gene-protein-reaction association. These features allow us to conduct studies on the impact of secondary metabolisms, isozymes, media composition, and provide a more solid basis for computational strain design. We used several sets of experimental data to train the model, e.g., estimation of the ATP requirement for growth-associated maintenance (13.5 mmol ATP/g DCW/hr) and cellulose synthesis (57 mmol ATP/g cellulose/hr). Using our tuned model, we predicted the experimentally observed differences in cell biomass yield based on which cellodextrin species is assimilated. We further employed our tuned model to analyze the experimentally quantified differences in fermentation profiles (i.e., the ethanol to acetate ratio) between cellbiose- and cellulose-grown cultures, for which we inferred potential regulatory mechanisms to explain the phenotypic differences. Finally, we used the model to design over 250 genetic modification strategies with the potential to optimize ethanol production, 6,155 for hydrogen production, and 28 for isobutanol production. Our developed genome-scale model iAT601 is capable of accurately predicting complex cellular phenotypes under a variety of conditions by integration of low- and high-throughput data, and serves as a high-quality platform for model-guided strain design to produce industrial biofuels and chemicals of interest.

### Theory

**Stoichiometric Metabolic Modeling**

$$S \cdot r = 0$$

$$r_{irrev} \geq 0$$

**Alternative strategies**

- Elementary Mode Analysis finds basis pathways to study network properties.
- Estimate flux distributions:
  - Metabolic Flux Analysis (Best experimental data fit)
  - Flux Balance Analysis (Growth rate maximization)

Admissible Flux Space

- Defined by network

### Modeling Methods

**Flux Balance Analysis**

$$\max r_{BIO}$$

$$\sum_j S_{ij} r_j = 0 \text{ metabolites } i$$

$$r_j^{min} \leq r_j \leq r_j^{max} \forall \text{ reactions } j$$

- Where reactions are bound by  $r_j^{min}$  and  $r_j^{max}$

**Metabolic Flux Analysis**

$$r_{rxn} = EM_{rxn} \cdot w_{k \times 1}$$

- $r$  is the flux distribution vector,  $EM$  is the elementary mode matrix, and  $w$  ( $R^{*}$ ) is the vector of weighting factors
- $r_m = EM_m \cdot w$
- $w = pinv(EM_m) \cdot r_m$

- Weighting factors can be calculated from the measured fluxes ( $r_m$ ) and the Moore-Penrose pseudoinverse of  $EM_m$

**Metabolic Flux Ratio Analysis**

$$MFR_{i-} = \frac{r_{i-}}{\sum_k r_{k-i}}$$

- $r$  is the incoming flux(es) of node  $i$  through reaction  $j$ . can be defined similarly where outgoing flux(es) of node  $i$  is represented with the opposite arrow.

**Constrained Minimal Cut Sets**

- $T \cdot r \leq b$  Target flux polyhedron. Fluxes below specified production formation rate.
- $D \cdot r \leq d$  Desired flux polyhedron. Fluxes with growth rate above desired minimum.
- Enumerate all reaction knockouts satisfying constraints.

### Questions We Can Address:

- What media/culturing conditions are optimal for production of target metabolites?
- What genotypes can constrain the phenotypic space of *C. thermocellum* for high-yield biofuel production?
- What understanding can we gain of metabolic and regulatory phenotypes from OMICs data?
- How can we balance thermodynamics, kinetics, and enzyme levels to optimize biofuel production?

### Primary References

- Thompson et al. Biot. Biofuels, under review.
- Thompson et al. Metab. Eng., 2015.

### Core Metabolic Model Capabilities

**Experimental Yields**

Experimental yields between parent and  $\Delta hydG$  are not significantly altered.

- In  $\Delta hydG \Delta pta$ , reducing acetate yield forces carbon from acetyl-CoA to ethanol, pulling electrons away from hydrogen production. The lack of ATP generation via acetyl-kinase manifests in a lower DCW yield.
- The significant drop in hydrogen yield in  $\Delta hydG \Delta ech$  is correlated with a dramatic increase in ethanol yield, as well as an unexpected increase in formate yield.

**Increase in ethanol yield is more dependent on constrained electron flow than constrained carbon flow.**

**METAFor Analysis**

- METAFor analysis reiterates the robustness of hydrogen production, as ECH can compensate for the lack of hydG related activity.
- METAFor analysis of  $\Delta hydG \Delta pta$  shows that the push of carbon to ethanol instead of acetate pulls electrons with it. This is seen in the increase in flux through PFOR, RNF, and NFN.
- In  $\Delta hydG \Delta ech$ , RNF and NFN fluxes increase dramatically to recycle reduced ferredoxin. However, there is still an apparent redox imbalance that favors diverting pyruvate flux through the redox neutral PFL reaction instead of PFOR.

**Alterations of Flux Distribution**

Major outliers are associated with redox metabolism with  $\Delta hydG \Delta ech$  being the most effective perturbation presented here:

### iAT601

The refined genome scale model consists of 601 reactions encoded by 601 genes, spanning multiple KEGG ontology categories.

Reconstruction → Curation → Validation → Prediction

Genome sequence → Reaction database → Literature database

KBbase DOE Systems Biology Knowledgebase

### Genome Scale Metabolic Model Capabilities

**Finding ATP Requirements: Cellbiose vs Cellulose**

**Cellbiose**

Initially, cellbiose was used as the sole constraint for flux balance analysis, which lead to poor predictions (highlighted in red).

Setting fermentation parameters, we sequentially varied the growth associated maintenance term to fit the experimental data.

Combining experimental constraints with the fitted GAM coefficient gives good agreement to Flux Balance Analysis values.

**Cellulose**

For simulation of cellulosic growth, we set the cellulose to be 20% of DCW, while setting flux constraints to experimental values. This leads to poor growth prediction (highlighted in red).

In a similar manner to above, we calibrated the ATP requirement for cellulosic production. This value had to be calibrated because secretion and turnover of the cellulose costs more ATP than normal cellular proteins.

Combining experimental fluxes, our tuned GAM coefficient, and the tuned cellulose coefficient leads to accurate phenotype prediction from the simulations.

**Model Driven Elucidation of Redox Bottlenecks**

Left, the core model predicts possible phenotypes reachable by the strains presented (Lines) and experimental data falls within these spaces (symbols). The Quad mutant contains  $\Delta hydG \Delta pta \Delta pfl \Delta dh$ . Right, the model predicts no growth with a  $\Delta hydG \Delta ech \Delta pfl$  genotype.

**Predicting Optimal Genotype**

Minimal Metabolic Functionality algorithm predicts that *C. therm* can be constrained to a high ethanol yielding phenotype by eliminating hydrogen, acetate, and lactate production and valine secretion.

**Strain Design via Constrained Minimal Cut Sets**

Minimal cut sets are metabolic engineering strategies which seek to maximize product yields while maintaining a minimum growth rate. We set two constraints for all of them, minimum growth rate, and minimum product yield.

Products	Target cut set sizes	# Strain Designs
Ethanol	6	67
Hydrogen	7	185
Isobutanol	4	12
	5	221
	6	1198
	7	4816
	8	28

With our genome scale model, we are able to find over 200 strategies for optimizing ethanol production by eliminating 6 or 7 genes.

Hydrogen production can be optimized with as little as 4 cuts, while isobutanol needs 7 cuts.

These strategies are excellent starting points for strain design and require more in depth analysis.

### Core Metabolic Model

**Cell Composition**

**iAT\_core**

**Most Relevant Reaction Abbreviations**

- pyruvate phosphate dikinase (PPDK)
- oxaloacetate decarboxylase (ODC)
- malic enzyme (MAE)
- pyruvate formate lyase (PFL)
- pyruvate ferredoxin oxidoreductase (PFOR)
- lactate dehydrogenase (LDH)
- alcohol dehydrogenase (AdhE)
- phosphotransacetylase (PTA)
- citrate synthase (TCA1)
- Ni-Fe energy conserving hydrogenase (ECH)
- Ni-Fe NADH-Fd oxidoreductase (RNF)
- NADH-Fd:NADP+ oxidoreductase (NFN)
- NADH-Fd:ribose hydrogenase (BIF)
- Fe-Fe NADPH-dependent hydrogenase (Fe-H2)

**Metabolic nodes of interest:** PYR, PFL, NFN, Fd, FdH2, LDH, ECH, H2, ETOH, ACE, VAL, HPP.

**Determination of Cellulose Composition**

To determine cellulose composition, proteomics data from multiple substrate conditions were compared.

Abundance values and protein sequences allowed for the calculation of amino acid requirements, and the median value across conditions was used for our *in silico* cellulose composition.

PR: GEM Dry cell weight protein composition, CB: Experimental data from growth on cellbiose, C: Experimental growth on cellulose, CX: Cellulose / Xylan, CP: Cellulose / pectin, CPX: Cellulose / Pectin / Xylan, SWG: Switchgrass, ZT: Z-Tim dietary fiber, Cell: median values for application in cellulose component.

**Effect of Cellodextrin Length**

*C. thermocellum* processes cellodextrins phosphorylatically, yielding more ATP per glucose unit as the oligomers increase in length.

Our results simulating protein yield correlate well with literature values.

**Strain Design Via Constrained Minimal Cut Sets**

Top 15 knockout set strategies

Products	Target cut set sizes	# Strain Designs
Ethanol	6	67
Hydrogen	7	185
Isobutanol	4	12
	5	221
	6	1198
	7	4816
	8	28

These strategies are excellent starting points for strain design and require more in depth analysis.