

# Mathematical Programming for Identification of Functional Differences in Genome-Scale Models

**Joshua J. Hamilton** and Jennifer L. Reed

Department of Chemical and Biological Engineering  
University of Wisconsin-Madison

243<sup>rd</sup> ACS National Meeting  
March 26, 2012



**WISCONSIN**  
UNIVERSITY OF WISCONSIN-MADISON



# Outline

- 1 Introduction and Background
- 2 CONGA Algorithm
- 3 Comparison of *E. coli* Metabolic Models

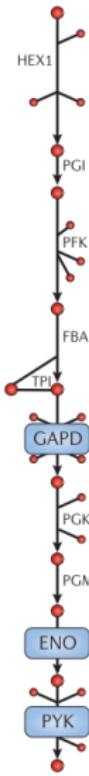
## Why models?

- Understanding cellular behavior
- Evaluate and design metabolic engineering strategies
- Contextualize high-throughput data

# Genome-Scale Models & Constraint-Based Modeling

## Why models?

- Understanding cellular behavior
- Evaluate and design metabolic engineering strategies
- Contextualize high-throughput data



# Genome-Scale Models & Constraint-Based Modeling

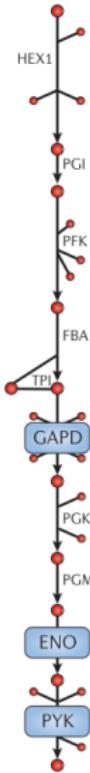
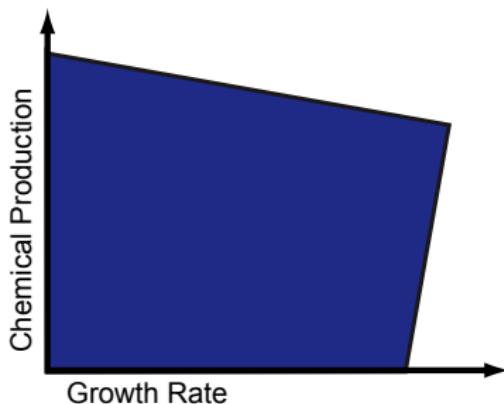
## Why models?

- Understanding cellular behavior
- Evaluate and design metabolic engineering strategies
- Contextualize high-throughput data

Flux distribution subject to

Solution Space

- Steady-state mass balance constraints
- Limits on fluxes



# Genome-Scale Models & Constraint-Based Modeling

## Why models?

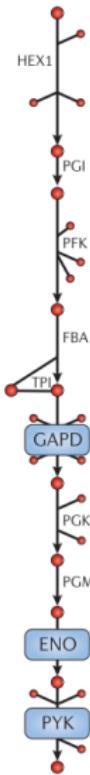
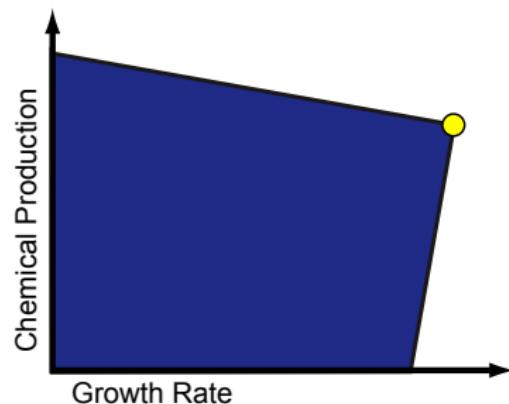
- Understanding cellular behavior
- Evaluate and design metabolic engineering strategies
- Contextualize high-throughput data

Flux distribution subject to

- Steady-state mass balance constraints
- Limits on fluxes

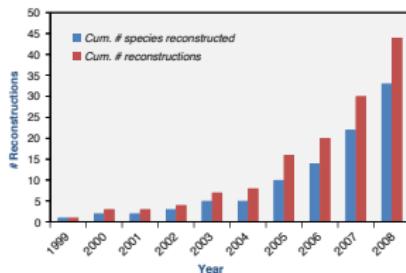
And want to

- Maximize growth



# Advances in Model Development

- Number of available models growing exponentially <sup>1</sup>

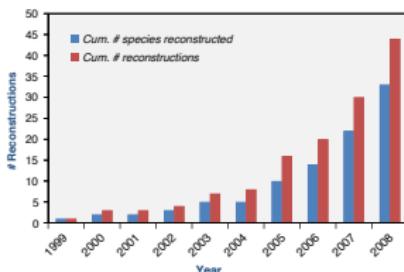


<sup>1</sup>Oberhardt et al, Mol Sys Bio 2009

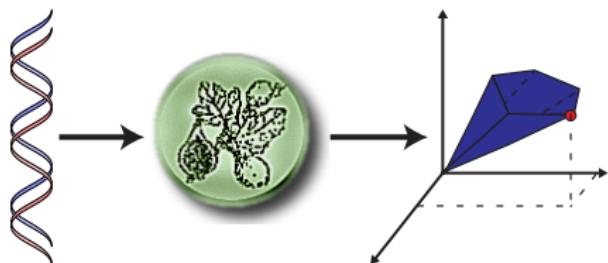
<sup>2</sup>Henry et al, Nature Biotech 2010

# Advances in Model Development

- Number of available models growing exponentially <sup>1</sup>



- New tools enable rapid conversion from genome to model <sup>2</sup>



<sup>1</sup>Oberhardt et al, Mol Sys Bio 2009

<sup>2</sup>Henry et al, Nature Biotech 2010



# Existing Approaches to Model Comparison

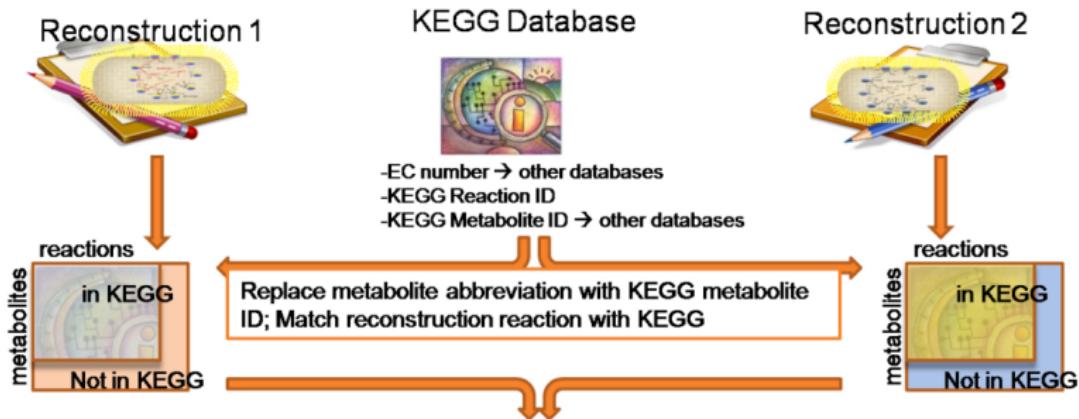
- Jamborees and network reconciliation align models at the reaction level and identify unique and shared network components

---

<sup>3</sup>Thiele et al, BMC Sys Bio 2011

# Existing Approaches to Model Comparison

- Jamborees and network reconciliation align models at the reaction level and identify unique and shared network components
- Require mapping of compounds and reactions across networks <sup>3</sup>



<sup>3</sup>Thiele et al, BMC Sys Bio 2011



# Existing Approaches to Model Comparison

- Jamborees and network reconciliation align models at the reaction level and identify unique and shared network components
- Require mapping of compounds and reactions across networks <sup>3</sup>
- Primarily descriptive, identifying only structural network differences <sup>4</sup>

---

<sup>3</sup>Thiele et al, BMC Sys Bio 2011

<sup>4</sup>Oberhardt et al, PLoS Comp Bio 2011

# CONGA: Comparison of Networks by Gene Alignment

## Advantages to Our Approach:

- Aligning models at the gene level can be done automatically
- Serves as a proxy for reaction-level alignments
- Mathematical programming can identify conditions under which structural network differences have a *functional* impact



# Outline

- 1 Introduction and Background
- 2 CONGA Algorithm
- 3 Comparison of *E. coli* Metabolic Models

# CONGA Formulation

- Bilevel Program

## CONGA Formulation <sup>5</sup>

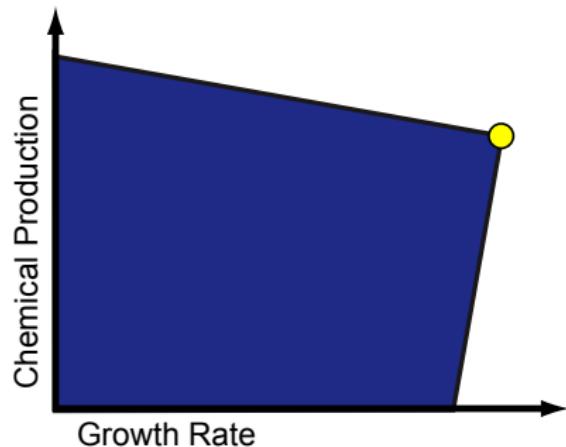
**max**

s.t. **max** cellular objective

s.t. cellular constraints

**max** cellular objective

s.t. cellular constraints



<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

# CONGA Formulation

- Bilevel Program
  - Identify genetic deletions

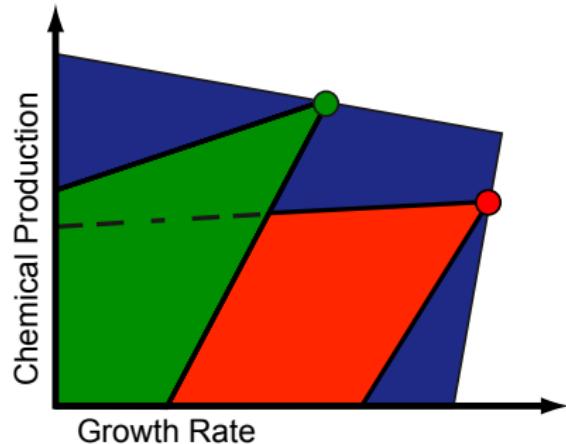
## CONGA Formulation <sup>5</sup>

**max**

s.t. **max** cellular objective  
s.t. cellular constraints  
gene deletions

**max** cellular objective

s.t. cellular constraints  
gene deletions



<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

# CONGA Formulation

- Bilevel Program
  - Identify genetic deletions
  - Giving rise to a difference in phenotype

## CONGA Formulation <sup>5</sup>

**max** difference in cellular phenotype

s.t. **max** cellular objective

s.t. cellular constraints

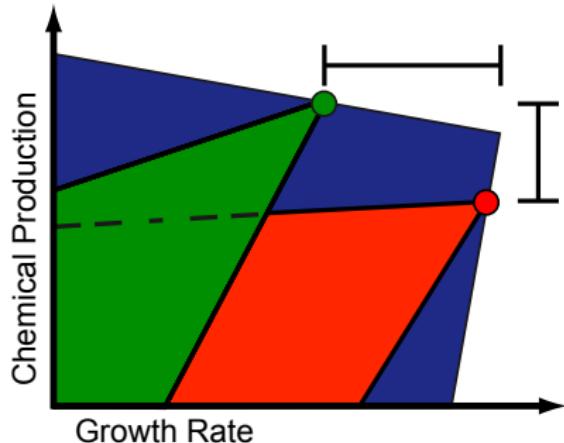
gene deletions

**max** cellular objective

s.t. cellular constraints

gene deletions

number of deletions  $\leq$  limit



<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

# CONGA Formulation

- Bilevel Program
  - Identify genetic deletions
  - Giving rise to a difference in phenotype
  - Which point to genetic differences w.r.t. that phenotype

## CONGA Formulation <sup>5</sup>

**max** difference in cellular phenotype

s.t. **max** cellular objective

s.t. cellular constraints

gene deletions

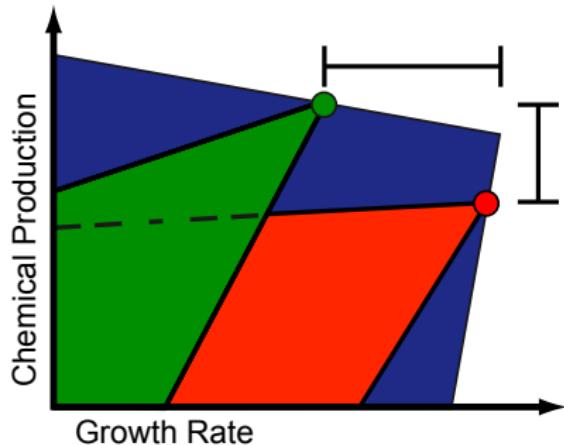
**max** cellular objective

s.t. cellular constraints

gene deletions

number of deletions  $\leq$  limit

*orthologs deleted in common*



<sup>5</sup>Hamilton and Reed, PLoS ONE 2012



# Outline

- 1 Introduction and Background
- 2 CONGA Algorithm
- 3 Comparison of *E. coli* Metabolic Models



# Comparison of *E. coli* Metabolic Models

- Genome-scale models of *E. coli* are frequently used in metabolic engineering studies
- Updated version released every few years: *iJR904* (2003)<sup>6</sup>, *iAF1260* (2007)<sup>7</sup>
- How does the updated model affect computational predictions?

---

<sup>6</sup>Reed et al, Genome Bio 2003

<sup>7</sup>Feist et al, Mol Sys Bio 2007



# Comparison of *E. coli* Metabolic Models

- Genome-scale models of *E. coli* are frequently used in metabolic engineering studies
- Updated version released every few years: *iJR904* (2003)<sup>6</sup>, *iAF1260* (2007)<sup>7</sup>
- How does the updated model affect computational predictions?
  - Can we identify conditions where one model predicts higher yields than the other?
  - Are these strategies desirable metabolic engineering strategies?
  - What structural network differences are responsible?

---

<sup>6</sup>Reed et al, Genome Bio 2003

<sup>7</sup>Feist et al, Mol Sys Bio 2007

# Identifying Model-Dominant Strategies

## Formulation

**max** difference in biochemical production

s.t. **max** cellular growth

s.t. cellular constraints

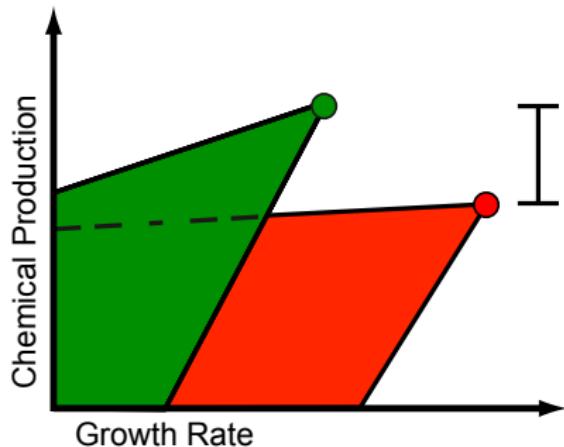
gene deletions

**max** cellular growth

s.t. cellular constraints

gene deletions

number of deletions  $\leq$  limit  
orthologs deleted in common



# Identifying Model-Dominant Strategies

## Formulation

**max** difference in biochemical production

s.t. **max** cellular growth

s.t. cellular constraints

gene deletions

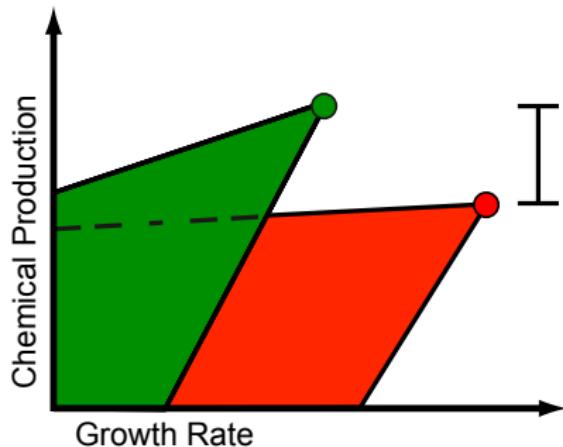
**max** cellular growth

s.t. cellular constraints

gene deletions

number of deletions  $\leq$  limit

orthologs deleted in common

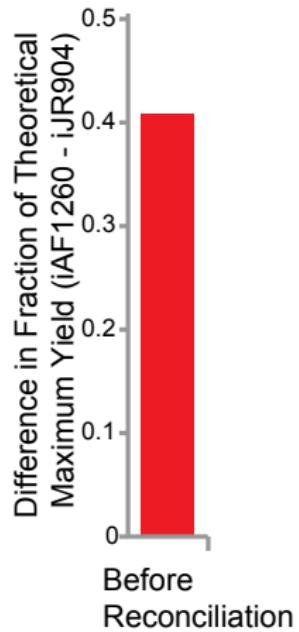


## CONGA identified model-dominant strategies for:

- Three fermentation products: ethanol, lactate, succinate
- Each *E. coli* model (*iAF1260* and *iJR904*)
- Across multiple numbers of gene deletions

# Production of Ethanol: Deletion of *mhpF* and *adhC*<sup>5</sup>

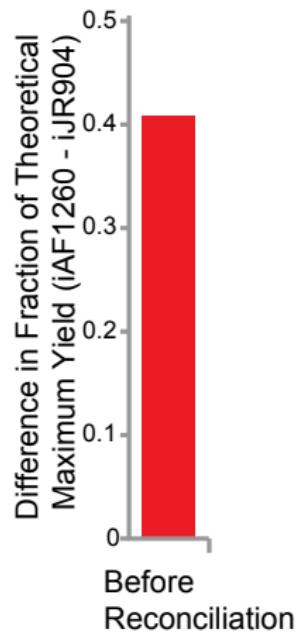
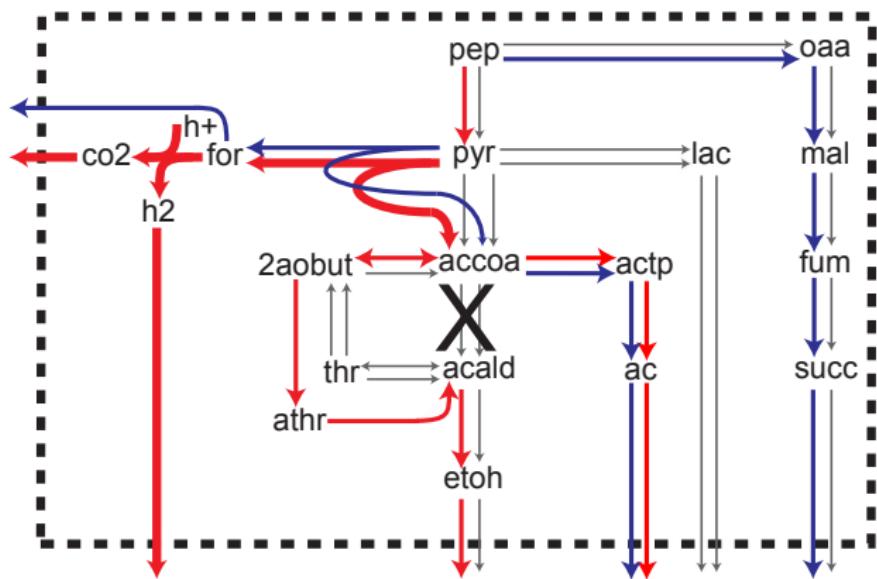
- Deletion of *mhpF* and *adhC* (acetaldehyde dehydrogenase) is model-dominant in the *iAF1260* model



<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

## Production of Ethanol: Deletion of *mhpF* and *adhC*<sup>5</sup>

- Deletion of *mhpF* and *adhC* (acetaldehyde dehydrogenase) is model-dominant in the *iAF1260* model

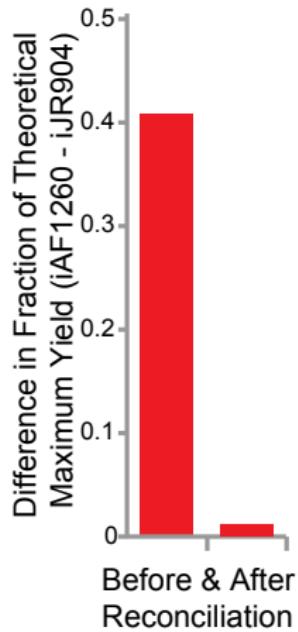


<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

# Production of Ethanol: Deletion of *mhpF* and *adhC*<sup>5</sup>

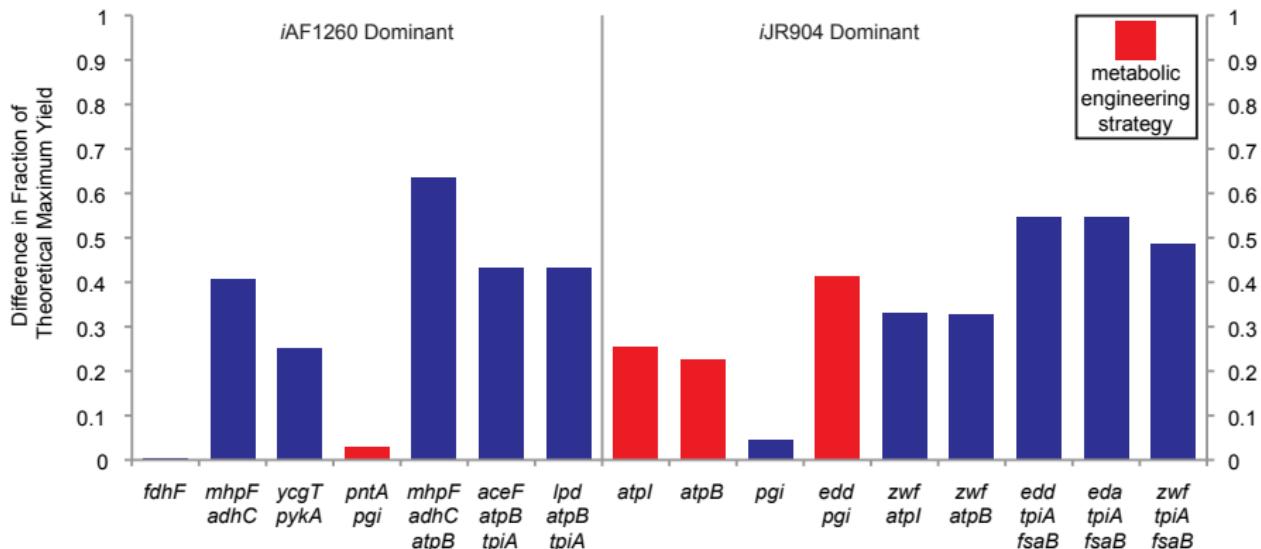
- Reconciliation of the *iJR904* model eliminates the phenotypic difference

Metabolite	% of Theoretical Yield	
	<i>iAF1260</i>	<i>iJR904</i>
Acetate	30.7	30.8
<b>Ethanol</b>	40.8	39.7
Succinate	0.5	0.5



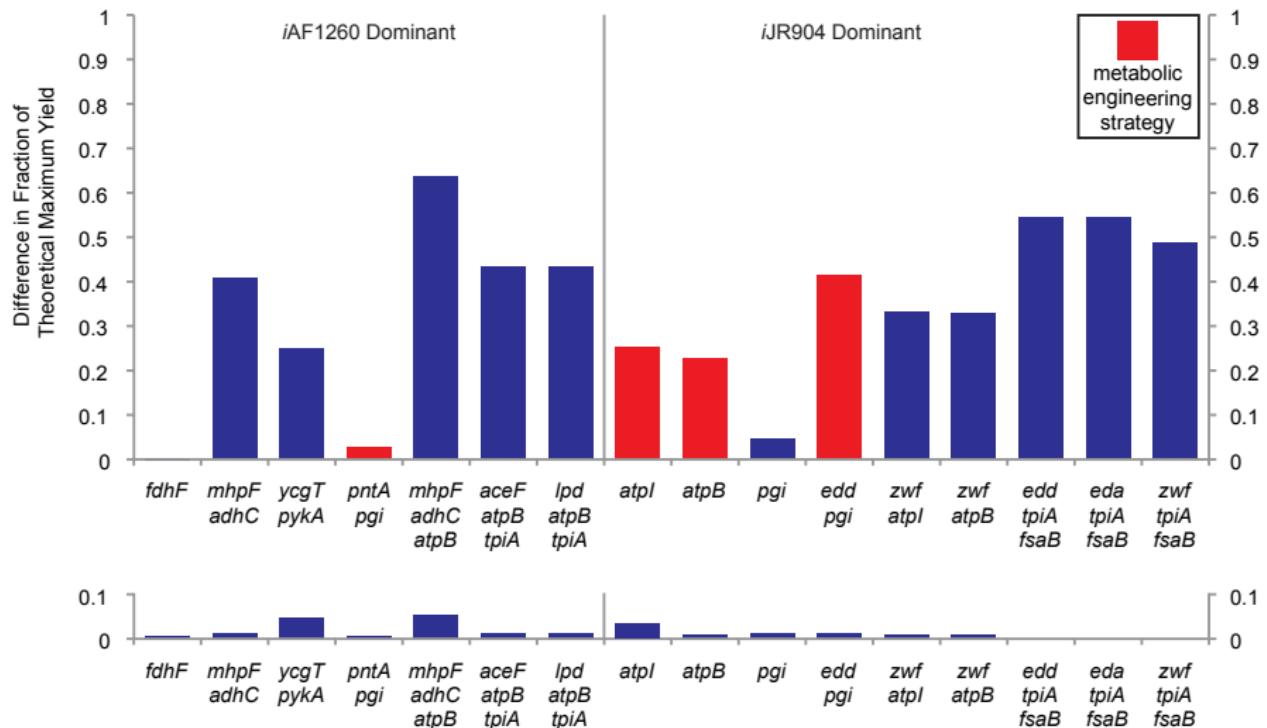
<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

# Results for Ethanol<sup>5</sup>



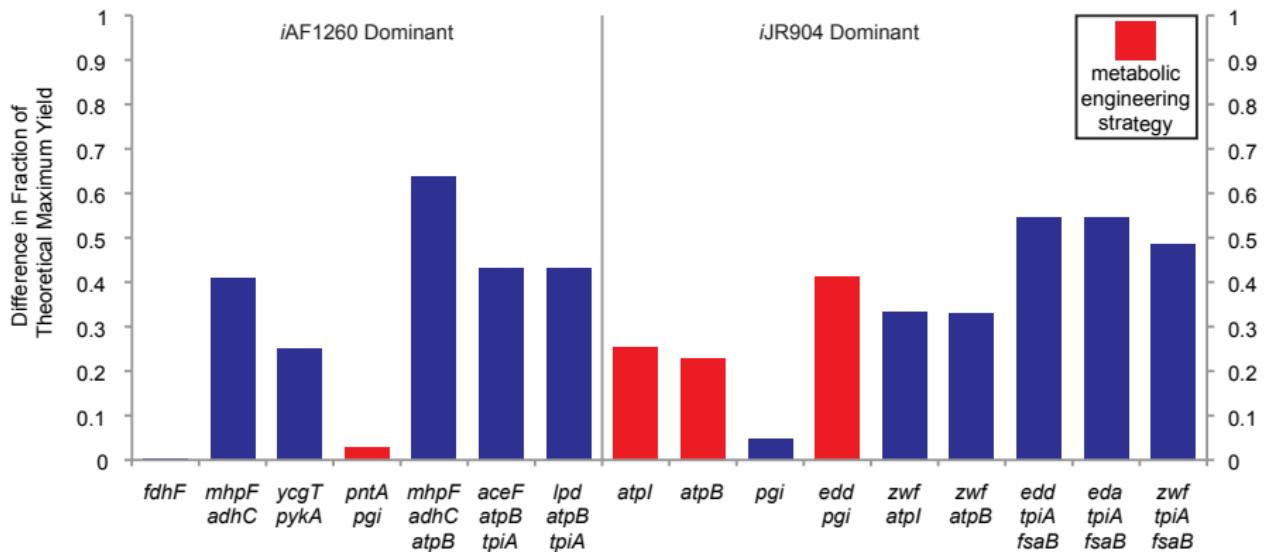
<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

# Results for Ethanol<sup>5</sup>



<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

# Results for Ethanol<sup>5</sup>



- Metabolic engineering strategies: among the top three strategies identified by OptORF<sup>8</sup>

<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

<sup>8</sup>Kim and Reed, PLoS Comp Bio 2010



# Results Across All Fermentation Products<sup>5</sup>

- Five metabolic differences accounted for 75% of all strategies we identified!

## **Metabolic Difference**

1,2-Propanediol Synthesis

Ethanol Synthesis

Hexokinase

Hydrogen Transport

Succinate Transport

- Model-dominant strategies were also metabolic engineering strategies in 40% of cases

---

<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

# Results Across All Fermentation Products<sup>5</sup>

- Five metabolic differences accounted for 75% of all strategies we identified!

## Metabolic Difference

1,2-Propanediol Synthesis  
Ethanol Synthesis  
Hexokinase  
Hydrogen Transport  
Succinate Transport

- Model-dominant strategies were also metabolic engineering strategies in 40% of cases
- Good models require an accurate representation of the pathway of interest

---

<sup>5</sup>Hamilton and Reed, PLoS ONE 2012



# Conclusions

- Developed a mathematical programming method for comparison of genome-scale network models



# Conclusions

- Developed a mathematical programming method for comparison of genome-scale network models
  - Identifies structural network differences relevant in a particular context
  - Facilitates network reconciliation - multiple models exist for *C. acetobutylicum*, *S. cerevisiae*, and more!
- Selection of background strains for metabolic engineering
- Can be applied to other perturbation strategies



# Conclusions

- Developed a mathematical programming method for comparison of genome-scale network models
  - Identifies structural network differences relevant in a particular context
  - Facilitates network reconciliation - multiple models exist for *C. acetobutylicum*, *S. cerevisiae*, and more!
  - Selection of background strains for metabolic engineering
  - Can be applied to other perturbation strategies
- Motivate a shift from *identifying* to *understanding* impact of differences between organisms

# Acknowledgements

## Prof. Jennifer Reed

- Dr. Dave Baumler
- Camo Cotten
- Joonhoon Kim
- Wai Kit Ong
- Chris Tervo
- Trang Vu
- Xiaolin Zhang

## Funding

- National Science Foundation: Graduate Research Fellowship Program



- Department of Energy: Genomic Science Program





# Questions?

- Developed a mathematical programming method for comparison of genome-scale network models
  - Identifies structural network differences relevant in a particular context
  - Facilitates network reconciliation and model development
  - Selection of background strains for metabolic engineering
  - Can be applied to other deletion strategies
- Motivate a shift from *identifying* to *understanding* impact of differences between organisms



# Formulation: Textual

maximize difference in flux  
subject to

*maximize cellular growth*

*subject to mass balance constraints*

*enzyme capacity constraints*

*thermodynamic constraints*

*reaction deletions*

GPR constraints

*maximize cellular growth*

*subject to mass balance constraints*

*enzyme capacity constraints*

*thermodynamic constraints*

*reaction deletions*

GPR constraints

orthologs deleted from both models

unique genes deleted from each model

limited number of deletions

*Model A*

*Model B*

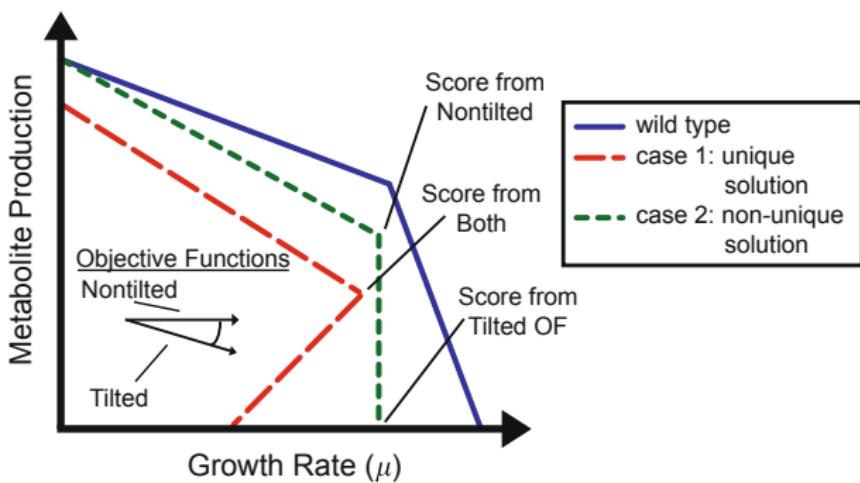


# Formulation: Mathematical

$$\begin{array}{ll} \max & v_{chem_A} - v_{chem_B} \\ \text{s.t.} & \max \quad v_{BM_A} - \gamma v_{chem_A} \\ & \text{s.t.} \quad \sum_j S_{ij} v_j = 0 \quad \forall i \in I \quad \forall \text{ Species A} \\ & \quad \alpha_j \leq v_j \leq \beta_j \quad \forall j \in J \quad \forall \text{ Species A} \\ & \quad v_j = 0 \quad \forall j \in J \mid y_j = 0 \quad \forall \text{ Species A} \\ \text{s.t.} & \max \quad v_{BM_B} - \gamma v_{chem_B} \\ & \text{s.t.} \quad \sum_j S_{ij} v_j = 0 \quad \forall i \in I \quad \forall \text{ Species B} \\ & \quad \alpha_j \leq v_j \leq \beta_j \quad \forall j \in J \quad \forall \text{ Species B} \\ & \quad v_j = 0 \quad \forall j \in J \mid y_j = 0 \quad \forall \text{ Species B} \\ & y_j = f(z_{\hat{g}}, w_{\hat{p}}) \quad \forall \text{GPR}(j, \hat{p}, \hat{g}) \in J, P, G \quad \forall \text{ Species A and B} \\ & \sum_g (1 - z_g) \leq K \quad \forall \text{ Species A and B} \\ & z_{g_A} = z_{g_B} \quad \forall (z_{g_A}, z_{g_B}) \in O \end{array}$$

# Nonuniqueness: Tilting the Inner Objective

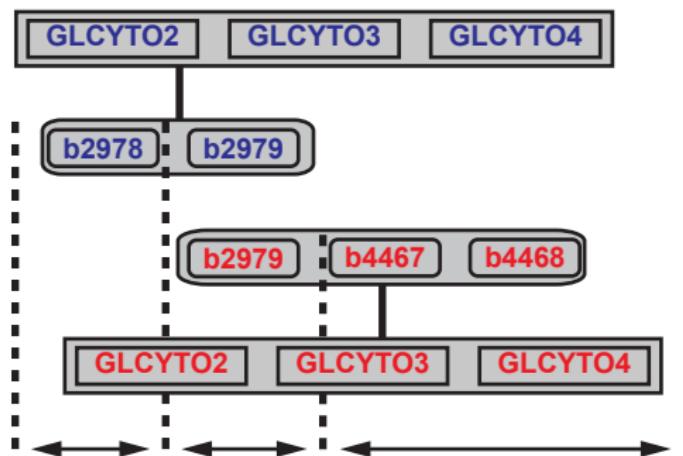
- In some instances a production flux may be nonunique
- May result in CONGA identifying incorrect solutions
- By imposing a penalty on the inner objective, we ensure the lowest production level is selected <sup>9</sup>



<sup>9</sup>Feist et al, Metab Eng 2010

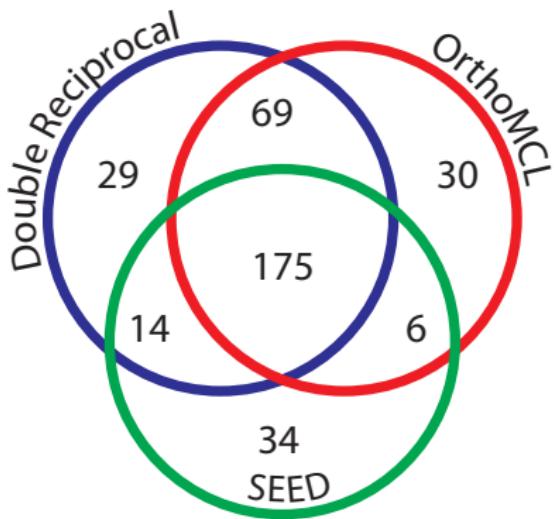
# Variable Reduction Procedure

- Computational difficulty of integer programs scales as the number of binary variables (genes)
- Consolidating groups of genes can ease the computational burden

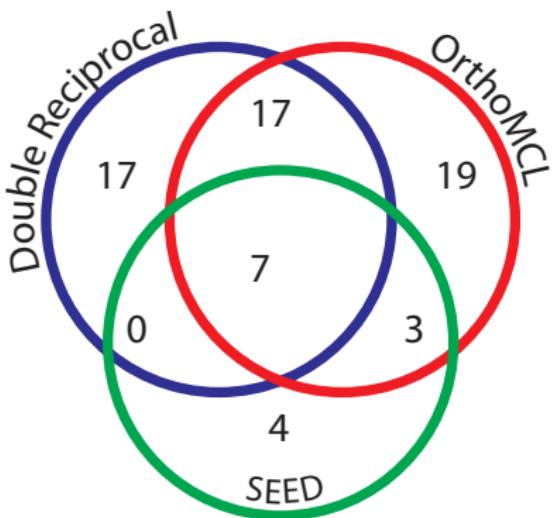


# Comparison of Methods for Orthology Prediction

Total Ortholog Calls



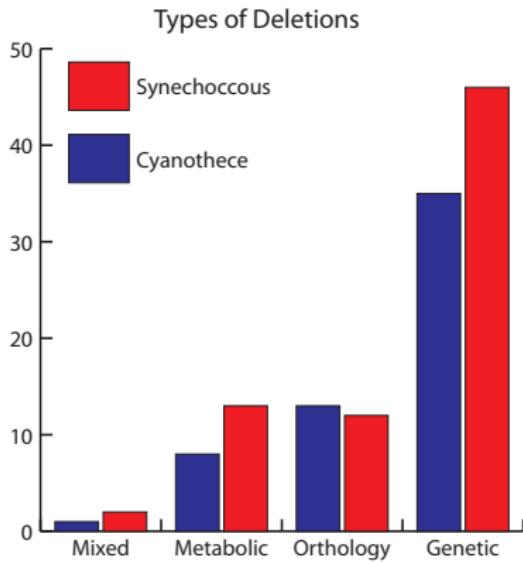
Incorrect Ortholog Calls



- BLAST and OrthoMCL call the most orthologs

# Model Development

- Developed a metabolic model for *Synechococcus* sp. PCC 7002<sup>5</sup> from a model of *Cyanothece* sp. ATCC 51142<sup>10</sup>



	Before	After
Genes	542	611
Reactions	491	552

<sup>5</sup>Hamilton and Reed, PLoS ONE 2012

<sup>10</sup>Vu et al, PLoS Comp Bio 2012