# The Robot Scientist Genesis: Abduction for Metabolic Modelling

Alexander H. Gower, Konstantin Korovin, Daniel Brunnsåker, Filip Kronstrom, Gabriel K. Reder, Ievgeniia A. Tiukova, Ronald S. Reiserer, John Wikswo, Ross D. King

The original discovery problem

How we formulated the problem in computational terms

What data and knowledge we provided to our system

How we represented the system's inputs and outputs

The space of candidate models that the system searched

What criteria it used to evaluate candidate models

How we interpreted results that the system generated

# Metabolic modelling *Saccharomyces cerevisiae*

- Yeast is the model eukaryote

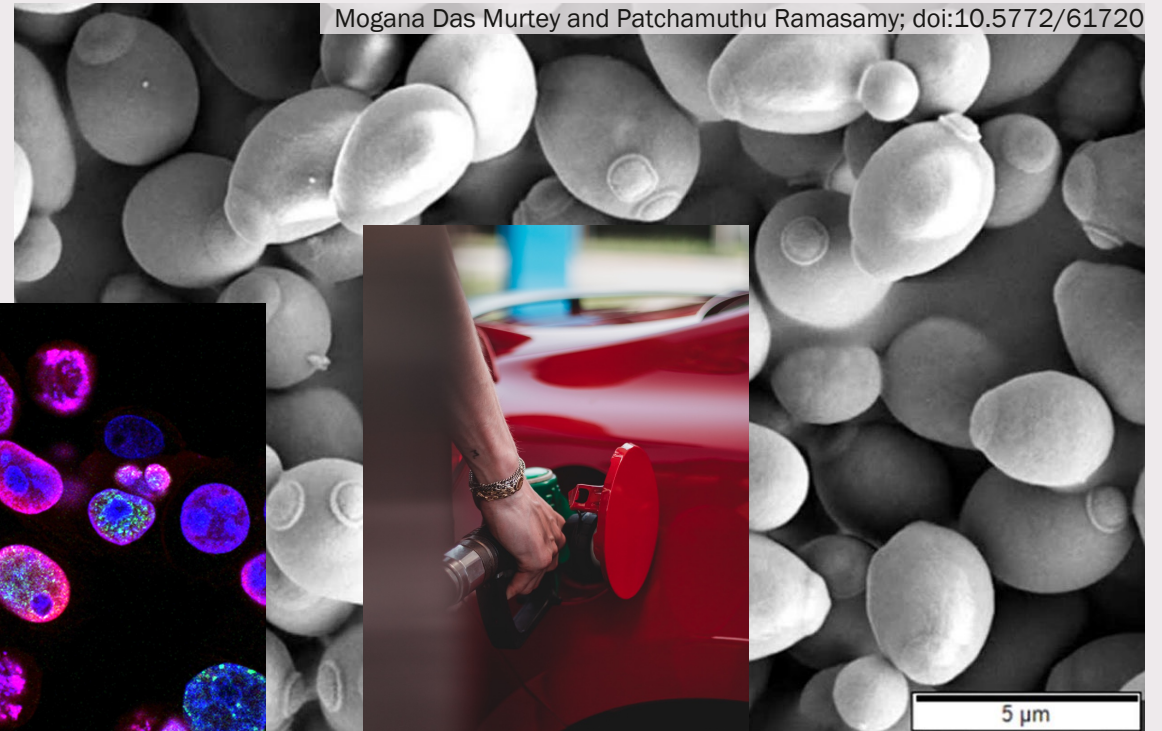- Exist tools to conduct experiments (e.g. CRISPR/Cas9)

- Cell factory

Mogana Das Murtey and Patchamuthu Ramasamy; doi:10.5772/61720

5 µm

Photo by Jeff Siepman on Unsplash

Photo by Jude Infantini on Unsplash

Photo by National Cancer Institute on Unsplash

Photo by Wassim Chouak on Unsplash

# Metabolic modelling *Saccharomyces cerevisiae*

D-glucopyranose 6-phosphate

fructose degradation
mannose degradation

pentose phosphate pathway

β-D-fructofuranose 6-phosphate

β-D-fructose 1,6-bisphosphate

glycerone phosphate

D-glyceraldehyde 3-phosphate

superpathway of phosphatidate biosynthesis
glycerol biosynthesis

3-phospho-D-glycerate

2-phospho-D-glycerate

gluconeogenesis I

phospho*enol*pyruvate

chorismate biosynthesis

pyruvate

TCA cycle, aerobic respiration
superpathway of acetoin and butanediol biosynthesis
Amino Acid Biosynthesis

Mogana Das Murtey and Patchamuthu Ramasamy; doi:10.5772/61720

"The ultimate goal of **genome-scale metabolic network** reconstruction in the future is to have a well-annotated network including all parts of the metabolism without any missing reactions or gaps; however it is not yet possible due to incomplete knowledge of the yeast metabolism."
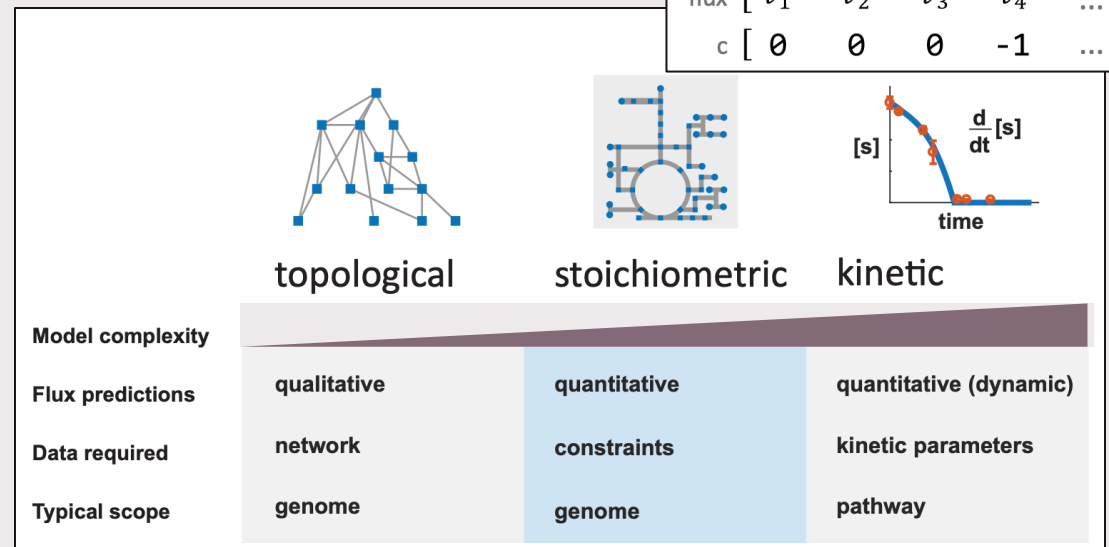
— *Österlund et. al (2012)*

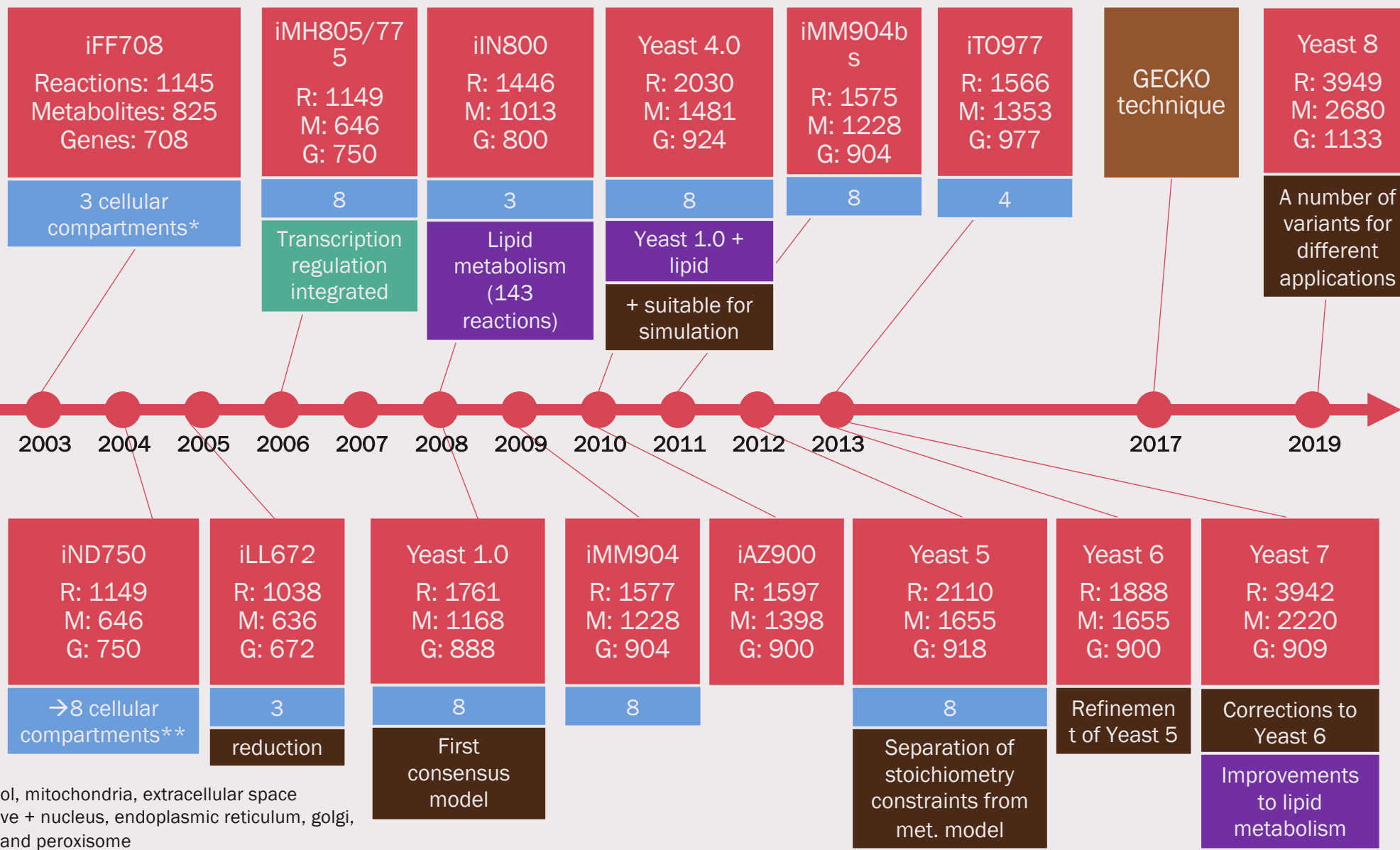# Metabolic modelling *Saccharomyces cerevisiae*

## Genome scale metabolic model

- Some quantities can be measure directly

- Others are abstractions (e.g. metabolic fluxes)

- Common approach is to encode biological knowledge as constraints
  - Either from observed experiments or
  - Biophysical knowledge



|  | R1 | R2 | R3 | R4 | ... |
|------|------|------|------|------|------|
| m1 | -1 | 0 | 0 | 0 | ... |
| m2 | +1 | -1 | 0 | 0 | ... |
| m3 | 0 | +1 | -1 | 0 | ... |
| m4 | 0 | +1 | +1 | -1 | ... |
| m5 | 0 | 0 | 0 | +1 | ... |
| ... | ... | ... | ... | ... | ... |

| flux | $v_1$ | $v_2$ | $v_3$ | $v_4$ | ... |
| c | 0 | 0 | 0 | -1 | ... |

|  | topological | stoichiometric | kinetic |
|---|---|---|---|
| **Model complexity** | | | |
| **Flux predictions** | qualitative | quantitative | quantitative (dynamic) |
| **Data required** | network | constraints | kinetic parameters |
| **Typical scope** | genome | genome | pathway |

Figures modified from Avlant Nilsson PhD thesis (2019)

**iFF708**
Reactions: 1145
Metabolites: 825
Genes: 708

3 cellular compartments*

**iMH805/775**
R: 1149
M: 646
G: 750

8

Transcription regulation integrated

**iIN800**
R: 1446
M: 1013
G: 800

3

Lipid metabolism (143 reactions)

**Yeast 4.0**
R: 2030
M: 1481
G: 924

8

Yeast 1.0 + lipid

+ suitable for simulation

**iMM904bs**
R: 1575
M: 1228
G: 904

8

**iTO977**
R: 1566
M: 1353
G: 977

4

**GECKO technique**

**Yeast 8**
R: 3949
M: 2680
G: 1133

A number of variants for different applications

2003  2004  2005  2006  2007  2008  2009  2010  2011  2012  2013  2017  2019

**iND750**
R: 1149
M: 646
G: 750

→8 cellular compartments**

**iLL672**
R: 1038
M: 636
G: 672

3

reduction

**Yeast 1.0**
R: 1761
M: 1168
G: 888

8

First consensus model

**iMM904**
R: 1577
M: 1228
G: 904

8

**iAZ900**
R: 1597
M: 1398
G: 900

**Yeast 5**
R: 2110
M: 1655
G: 918

8

Separation of stoichiometry constraints from met. model

**Yeast 6**
R: 1888
M: 1655
G: 900

Refinement of Yeast 5

**Yeast 7**
R: 3942
M: 2220
G: 909

Corrections to Yeast 6
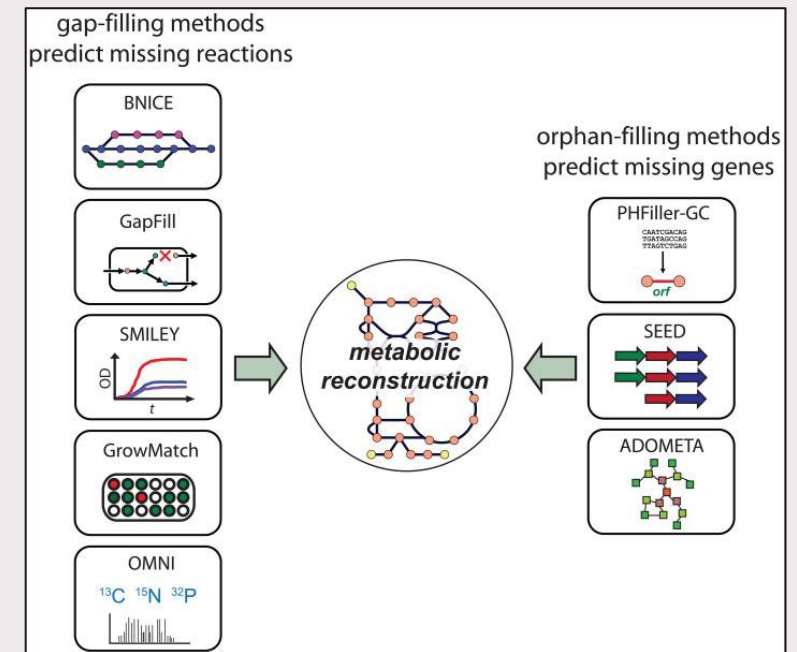
Improvements to lipid metabolism

* - cytosol, mitochondria, extracellular space
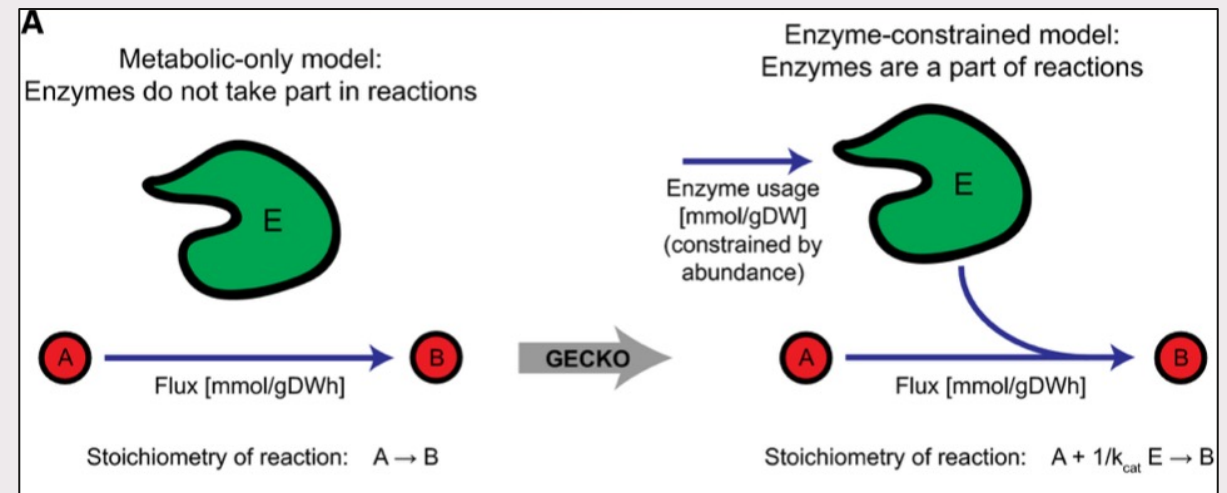** - above + nucleus, endoplasmic reticulum, golgi, vacuole and peroxisome

# Model improvement

- Model reduction
- Model expansion (new annotation)
  - Regulatory interaction
  - Assign gene to known reaction
  - Predict missing reaction
- Condition-specific effects
- Compartmentalisation (split model over parts of cell)
- New mathematical rules (enzymatic rate equation

$$v_i = \overline{k_{cat,i}} \cdot \sum E_i \cdot \rho_i$$

- New constraint mechanism (e.g. introducing enzymes explicitly into equations, GECKO)



Orth & Palsson (2010)



Sanchez, Zhang et. al. (2017)

# Humans and machines working together

- Model reduction
- Model expansion (new annotation)
  - Regulatory interaction
  - Assign gene to known reaction
  - Predict missing reaction
- Condition-specific effects

Algorithms exist or are being developed for these methods

- Compartmentalisation (split model over parts of cell)
- New mathematical rules (enzymatic rate equation

$$v_i = \overline{k_{cat,i}} \cdot \sum E_i \cdot \rho_i$$

- New constraint mechanism (e.g. introducing enzymes explicitly into equations, GECKO)

Currently proposed by human scientists – require a level of abstraction

# How to compare models

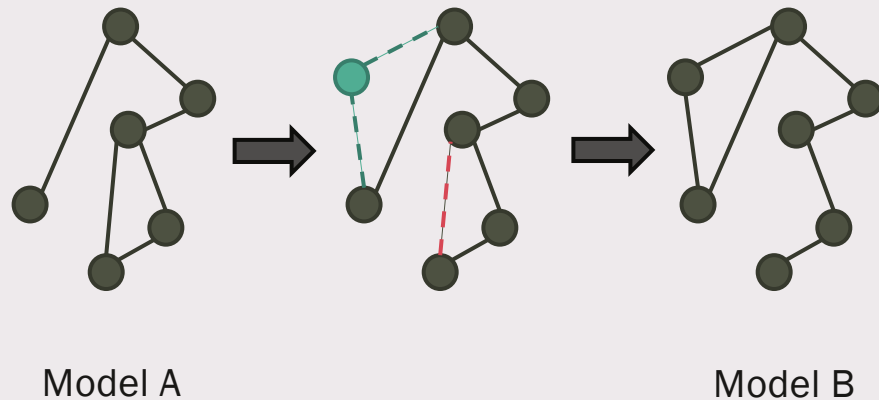Good models have:

| explanatory power | predictive power |
|---|---|
| consistency across contexts | consistency with other scientific models |

Many possible metrics one could use:

- genomic coverage;
- overlap of annotated metabolites;
- predictive ability for single gene essentiality;
- biomass production prediction;
- ...



Model A                    Model B

- It is difficult to find a single metric that can summarise a model's quality
- Among *S. cerevisiae* models there is evidence of tradeoffs between predictive accuracy and gene network coverage (Heavner and Price, 2015)
- Models are often developed for specific applications

The original discovery problem

**How we formulated the problem in computational terms**

What data and knowledge we provided to our system

How we represented the system's inputs and outputs

The space of candidate models that the system searched

What criteria it used to evaluate candidate models

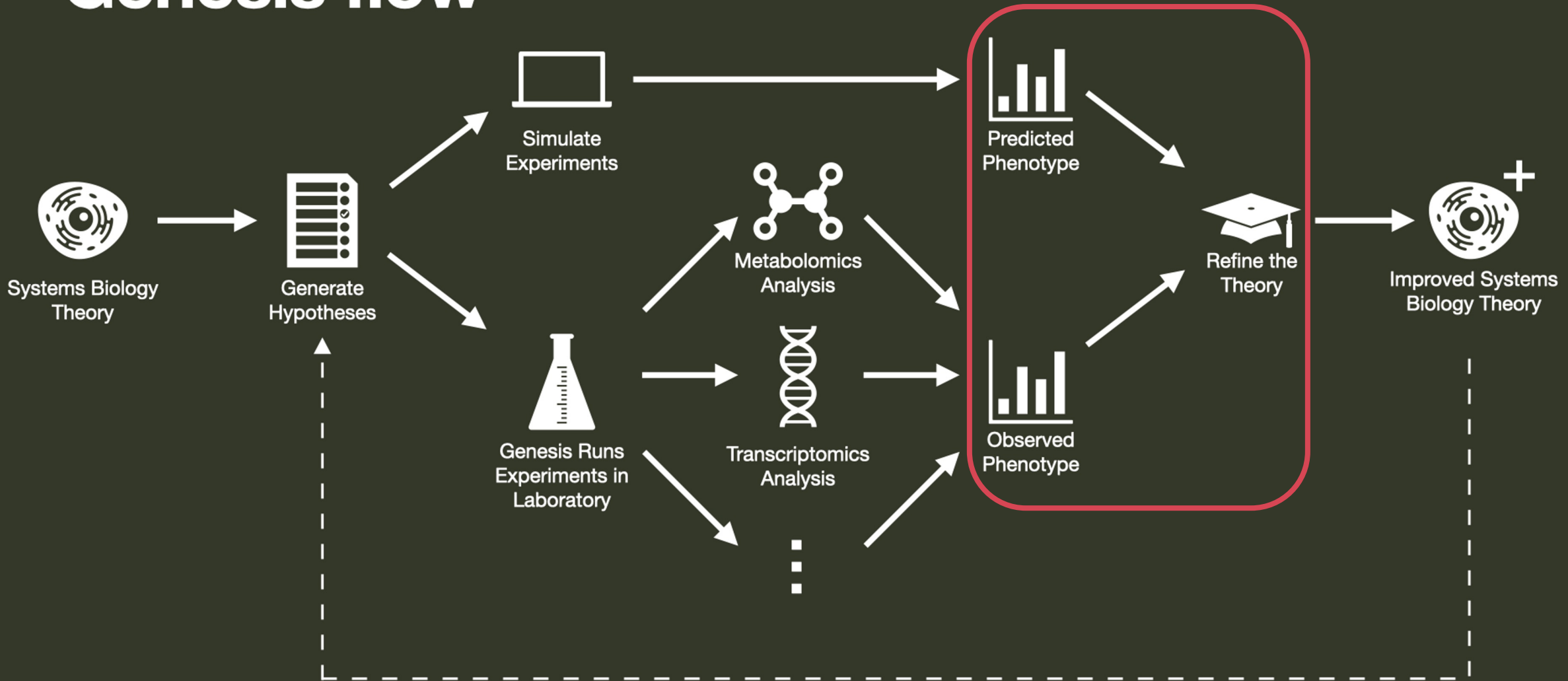How we interpreted results that the system generated

# Logical inference

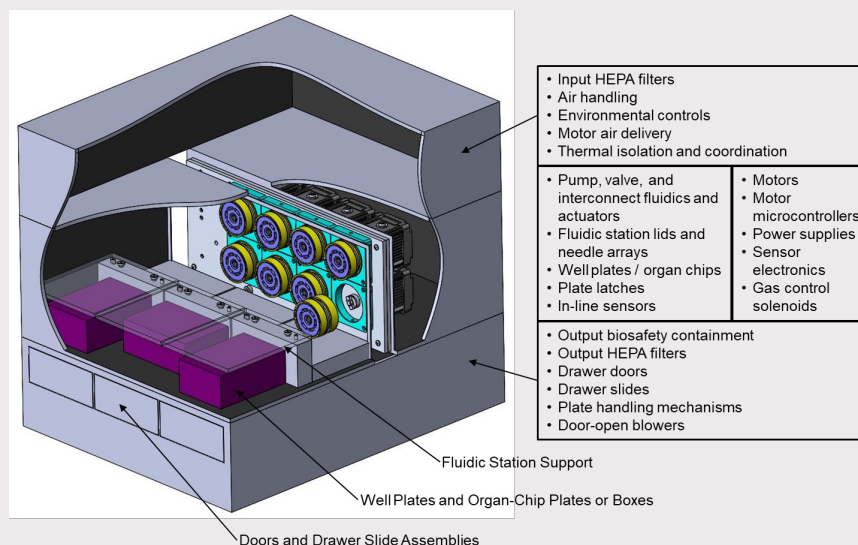| | |
|---|---|
| induction | allows us to generalise models from data |
| deduction | given a theory what conclusions can we draw |
| abduction | how can we "fix" the theory to be consistent with empirical data? |

# Active learning

Machine learning paradigm where the learning agent has agency over the selection of the next data to learn from — analagous to the scientific method

# Genesis flow

# Background: Genesis

- Scalable automated biological experimentation

- Small volume chemostat cultivations – vision is for thousands of parallel experiments

- AI-driven laboratory machine

- Measurements via high-throughput metabolomics and transcriptomics

John Wikswo group VIIBRE

The original discovery problem

How we formulated the problem in computational terms

What data and knowledge we provided to our system

How we represented the system's inputs and outputs

The space of candidate models that the system searched

What criteria it used to evaluate candidate models

How we interpreted results that the system generated

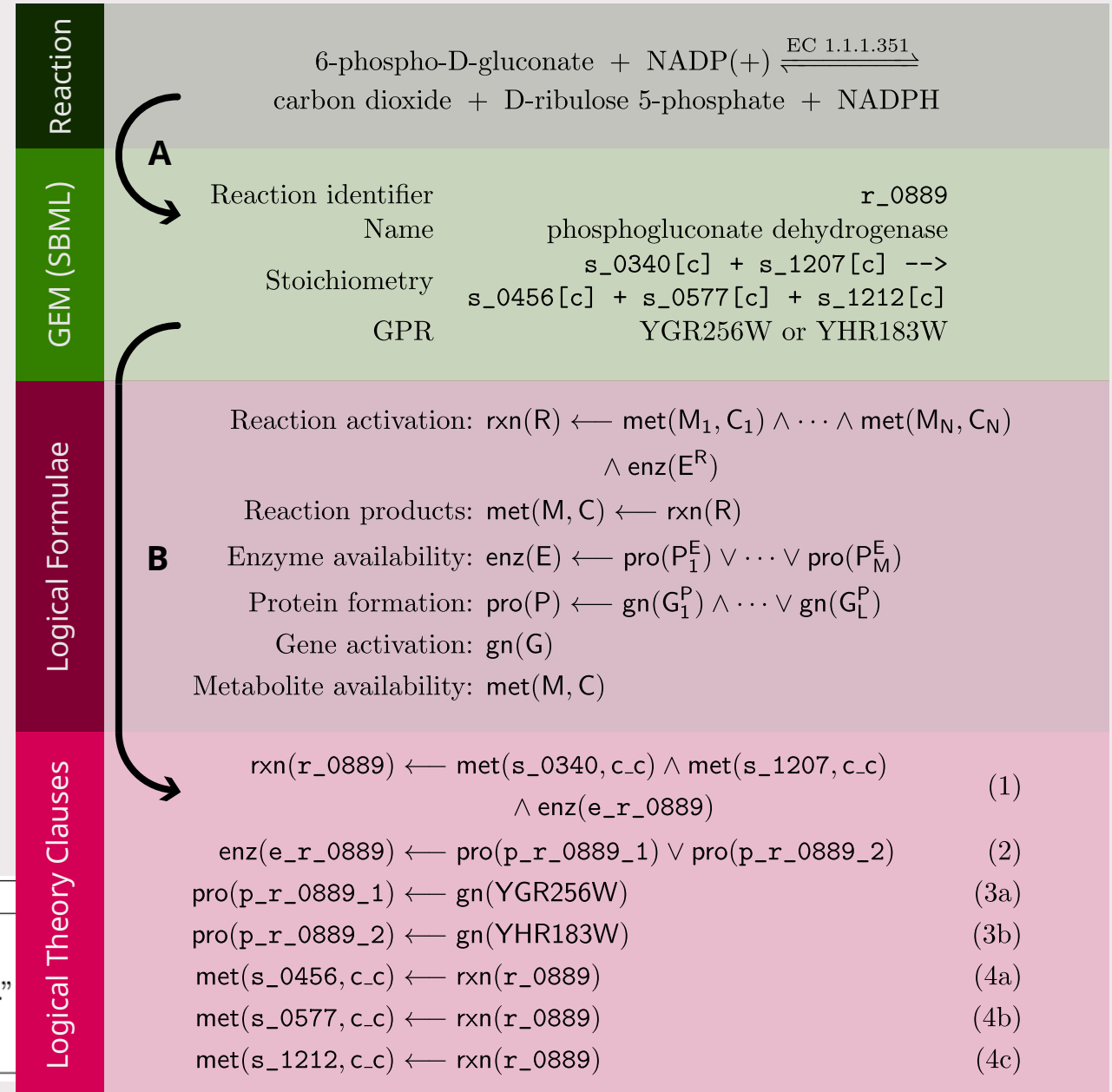# Constructing a logical theory of metabolism

- Background knowledge is encoded in the curated genome-scale metabolic models (GEMs)

- Each reaction in the GEM is translated to a set first-order logic clauses

- Clauses are written in conjugate normal form to produce theory

- Inference is performed using iProver: a theorem prover for first-order logic with support for arithmetical reasoning (Korovin, 2008)

- iProver has the efficiency required for a logical theory on this scale

| Predicate | Natural language interpretation |
|---|---|
| met\2 | "Metabolite X is present in cellular compartment Y." |
| gn\1 | "Gene X is expressed." |
| pro\1 | "Protein complex X is available (in every cellular compartment)." |
| enz\1 | "Isoenzyme category X is available." |
| rxn\1 | "There is positive flux through reaction X." |

**Reaction**

$$\text{6-phospho-D-gluconate} + \text{NADP(+)} \xrightleftharpoons{\text{EC 1.1.1.351}}$$
$$\text{carbon dioxide} + \text{D-ribulose 5-phosphate} + \text{NADPH}$$

**A**

**GEM (SBML)**

| | |
|---|---|
| Reaction identifier | r_0889 |
| Name | phosphogluconate dehydrogenase |
| Stoichiometry | s_0340[c] + s_1207[c] --> s_0456[c] + s_0577[c] + s_1212[c] |
| GPR | YGR256W or YHR183W |

**Logical Formulae**

**B**

Reaction activation: $\text{rxn}(R) \longleftarrow \text{met}(M_1, C_1) \wedge \cdots \wedge \text{met}(M_N, C_N) \wedge \text{enz}(E^R)$

Reaction products: $\text{met}(M, C) \longleftarrow \text{rxn}(R)$

Enzyme availability: $\text{enz}(E) \longleftarrow \text{pro}(P_1^E) \vee \cdots \vee \text{pro}(P_M^E)$

Protein formation: $\text{pro}(P) \longleftarrow \text{gn}(G_1^P) \wedge \cdots \vee \text{gn}(G_L^P)$

Gene activation: $\text{gn}(G)$

Metabolite availability: $\text{met}(M, C)$

**Logical Theory Clauses**

$\text{rxn}(\text{r\_0889}) \longleftarrow \text{met}(\text{s\_0340}, \text{c\_c}) \wedge \text{met}(\text{s\_1207}, \text{c\_c}) \wedge \text{enz}(\text{e\_r\_0889})$    (1)

$\text{enz}(\text{e\_r\_0889}) \longleftarrow \text{pro}(\text{p\_r\_0889\_1}) \vee \text{pro}(\text{p\_r\_0889\_2})$    (2)

$\text{pro}(\text{p\_r\_0889\_1}) \longleftarrow \text{gn}(\text{YGR256W})$    (3a)

$\text{pro}(\text{p\_r\_0889\_2}) \longleftarrow \text{gn}(\text{YHR183W})$    (3b)

$\text{met}(\text{s\_0456}, \text{c\_c}) \longleftarrow \text{rxn}(\text{r\_0889})$    (4a)

$\text{met}(\text{s\_0577}, \text{c\_c}) \longleftarrow \text{rxn}(\text{r\_0889})$    (4b)

$\text{met}(\text{s\_1212}, \text{c\_c}) \longleftarrow \text{rxn}(\text{r\_0889})$    (4c)

The original discovery problem

How we formulated the problem in computational terms

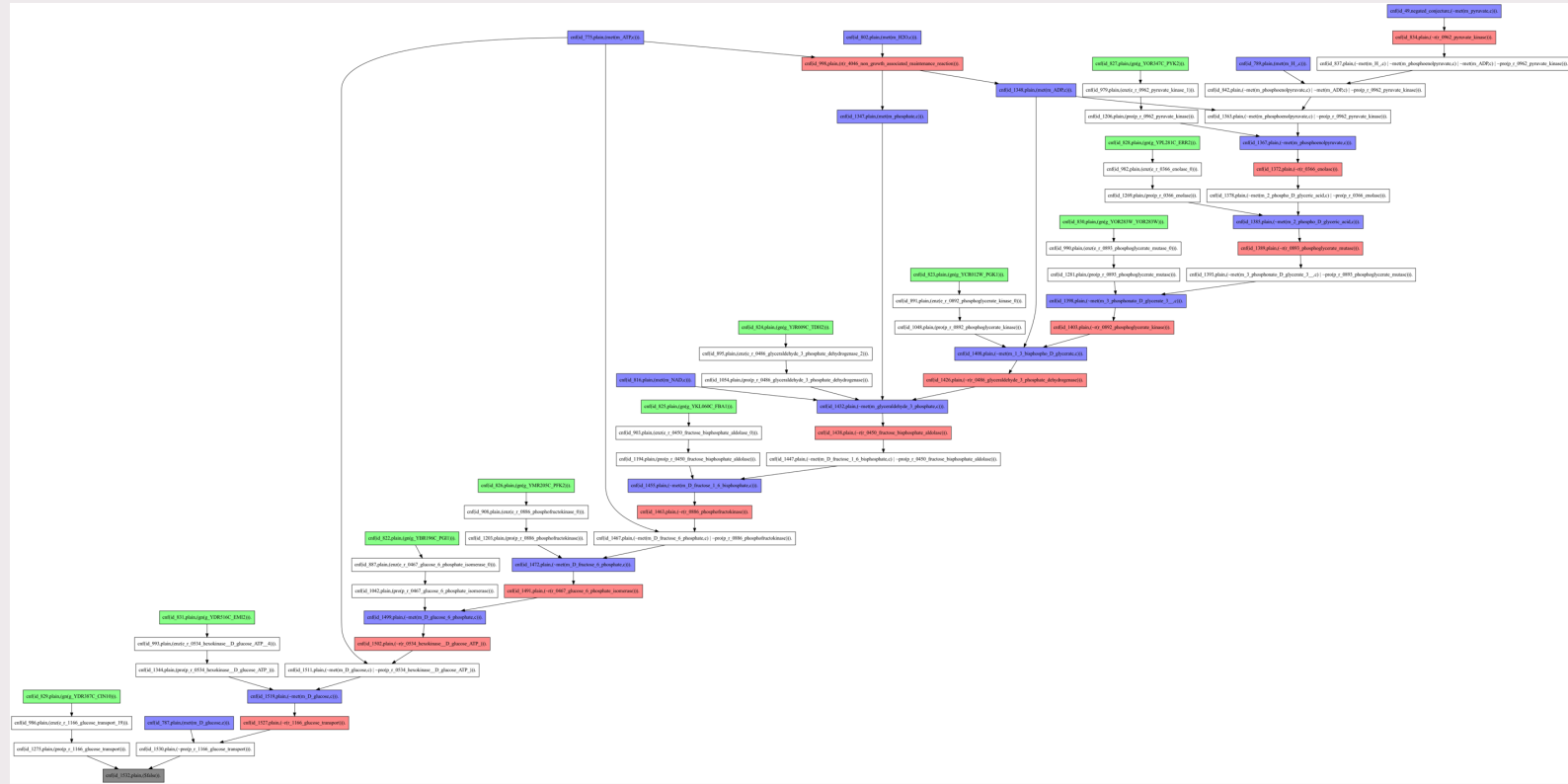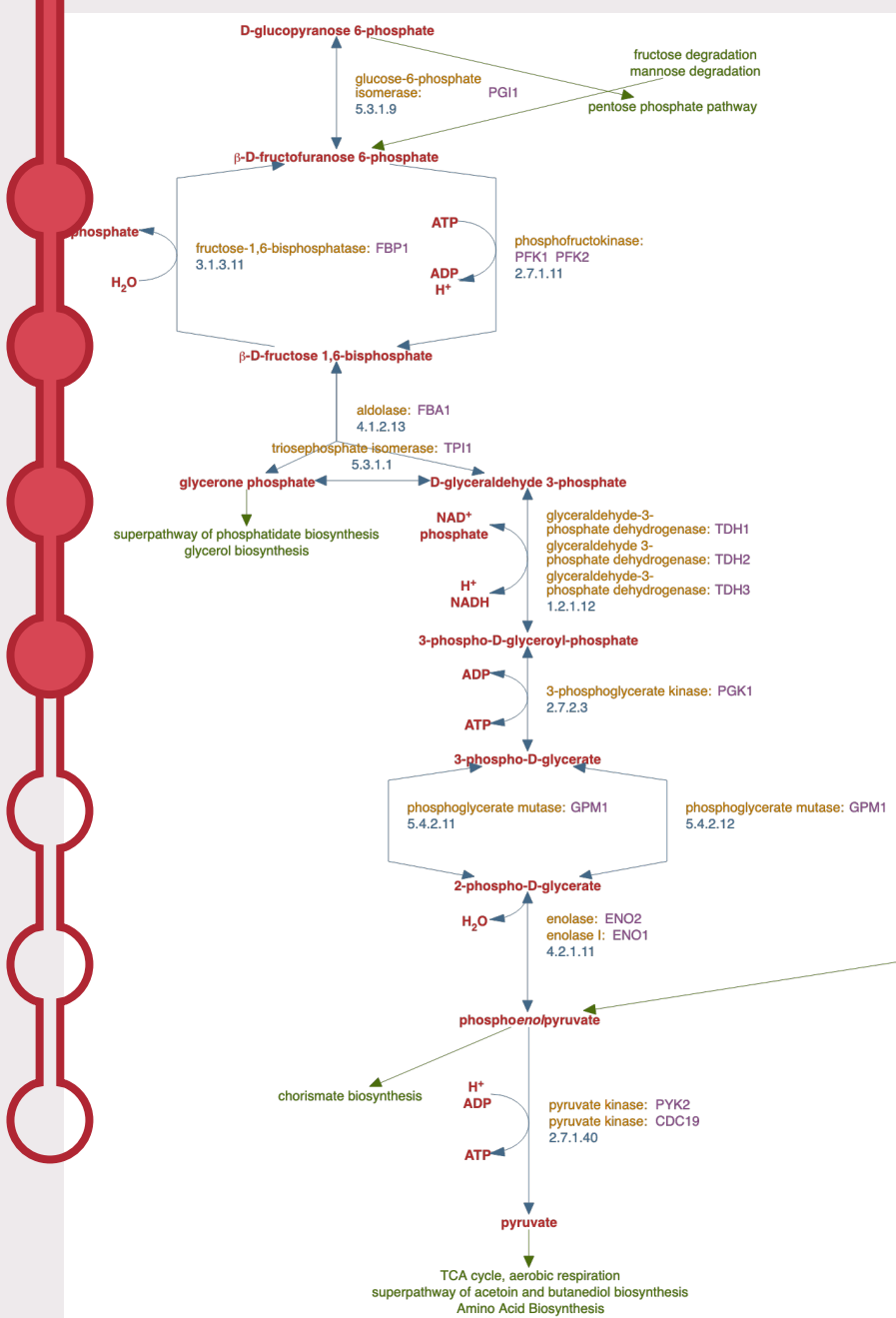What data and knowledge we provided to our system

How we represented the system's inputs and outputs

The space of candidate models that the system searched

What criteria it used to evaluate candidate models

How we interpreted results that the system generated

Glycolysis pathway (from https://pathway.yeastgenome.org/)

Proof graph from iProver (simplified)

The original discovery problem

How we formulated the problem in computational terms

What data and knowledge we provided to our system

How we represented the system's inputs and outputs

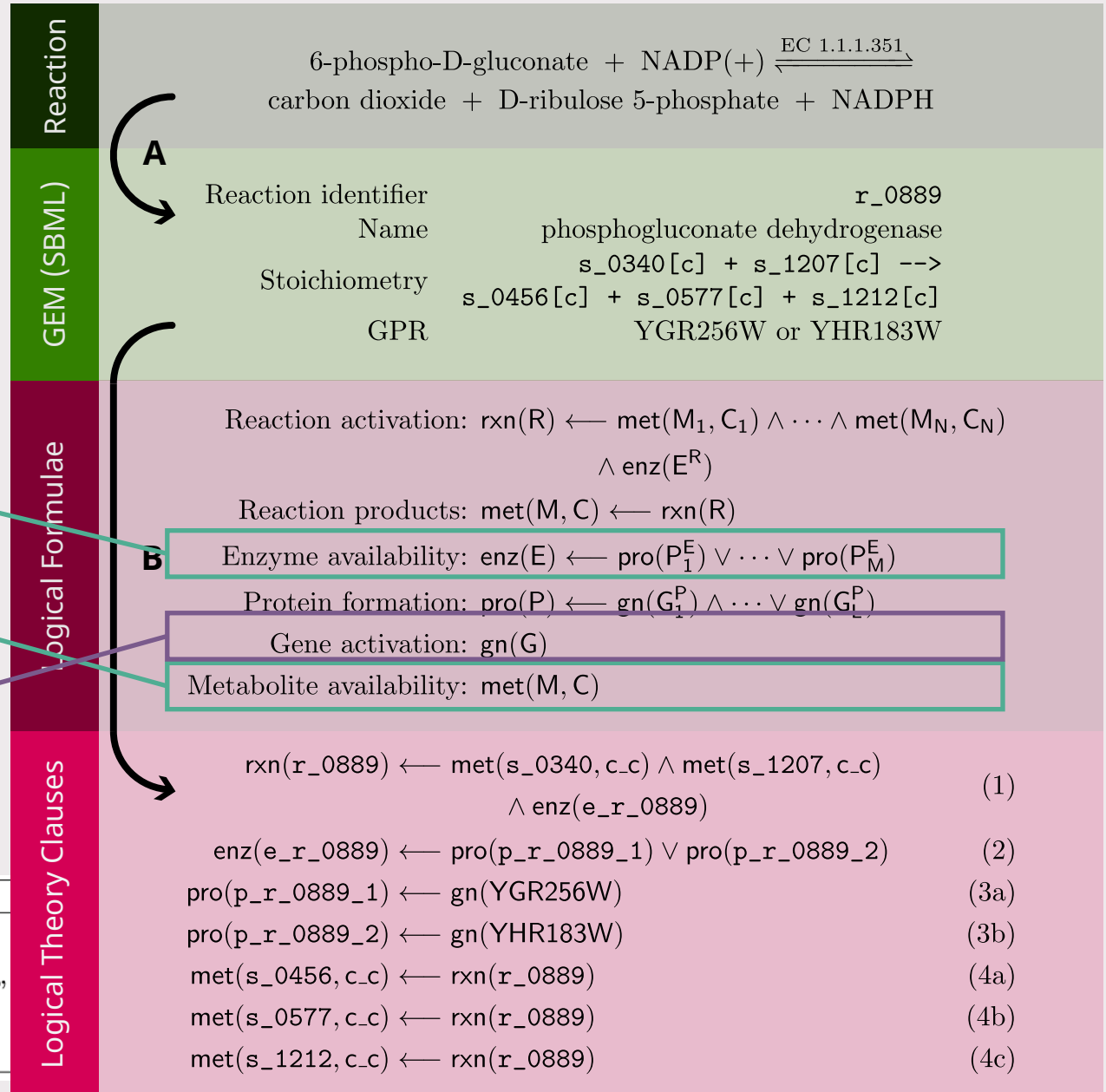The space of candidate models that the system searched

What criteria it used to evaluate candidate models

How we interpreted results that the system generated

# Approach: abduction opportunities

- Some parts of the model are well-known due to the chemistry (reaction stoichiometry, protein formation)

- We seek to:
  a) learn rules about which enzymes catalyse which reactions
  b) identify possible missing reactions by finding compound presence that will repair a broken pathway

- Future work will be to learn rules for gene expression and activation

| Predicate | Natural language interpretation |
|---|---|
| met\2 | "Metabolite X is present in cellular compartment Y." |
| gn\1 | "Gene X is expressed." |
| pro\1 | "Protein complex X is available (in every cellular compartment)." |
| enz\1 | "Isoenzyme category X is available." |
| rxn\1 | "There is positive flux through reaction X." |

**Reaction**

$$\text{6-phospho-D-gluconate} + \text{NADP(+)} \xrightleftharpoons{\text{EC 1.1.1.351}}$$
$$\text{carbon dioxide} + \text{D-ribulose 5-phosphate} + \text{NADPH}$$

**A**

**GEM (SBML)**

Reaction identifier      `r_0889`
Name      phosphogluconate dehydrogenase
Stoichiometry      `s_0340[c] + s_1207[c] -->`
`s_0456[c] + s_0577[c] + s_1212[c]`
GPR      YGR256W or YHR183W

**Logical Formulae**

Reaction activation: $\text{rxn}(R) \longleftarrow \text{met}(M_1, C_1) \wedge \cdots \wedge \text{met}(M_N, C_N)$
$\wedge \text{enz}(E^R)$

Reaction products: $\text{met}(M, C) \longleftarrow \text{rxn}(R)$

**B**    Enzyme availability: $\text{enz}(E) \longleftarrow \text{pro}(P_1^E) \vee \cdots \vee \text{pro}(P_M^E)$

Protein formation: $\text{pro}(P) \longleftarrow \text{gn}(G_1^P) \wedge \cdots \vee \text{gn}(G_L^P)$

Gene activation: $\text{gn}(G)$

Metabolite availability: $\text{met}(M, C)$

**Logical Theory Clauses**

$\text{rxn}(\texttt{r\_0889}) \longleftarrow \text{met}(\texttt{s\_0340}, \texttt{c\_c}) \wedge \text{met}(\texttt{s\_1207}, \texttt{c\_c})$
$\wedge \text{enz}(\texttt{e\_r\_0889})$    (1)

$\text{enz}(\texttt{e\_r\_0889}) \longleftarrow \text{pro}(\texttt{p\_r\_0889\_1}) \vee \text{pro}(\texttt{p\_r\_0889\_2})$    (2)

$\text{pro}(\texttt{p\_r\_0889\_1}) \longleftarrow \text{gn}(\text{YGR256W})$    (3a)

$\text{pro}(\texttt{p\_r\_0889\_2}) \longleftarrow \text{gn}(\text{YHR183W})$    (3b)

$\text{met}(\texttt{s\_0456}, \texttt{c\_c}) \longleftarrow \text{rxn}(\texttt{r\_0889})$    (4a)

$\text{met}(\texttt{s\_0577}, \texttt{c\_c}) \longleftarrow \text{rxn}(\texttt{r\_0889})$    (4b)

$\text{met}(\texttt{s\_1212}, \texttt{c\_c}) \longleftarrow \text{rxn}(\texttt{r\_0889})$    (4c)

The original discovery problem

How we formulated the problem in computational terms

What data and knowledge we provided to our system

How we represented the system's inputs and outputs

The space of candidate models that the system searched

What criteria it used to evaluate candidate models

How we interpreted results that the system generated

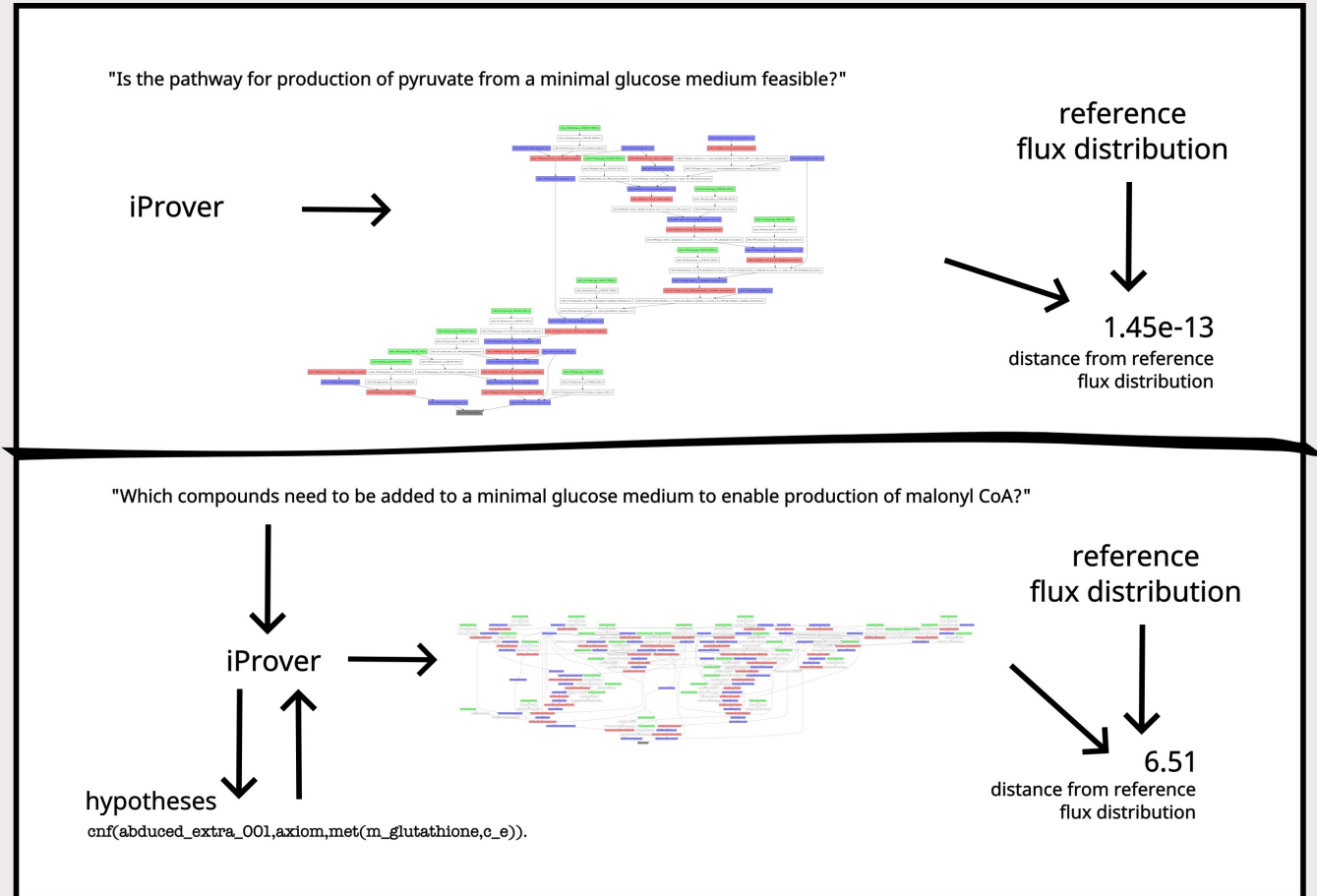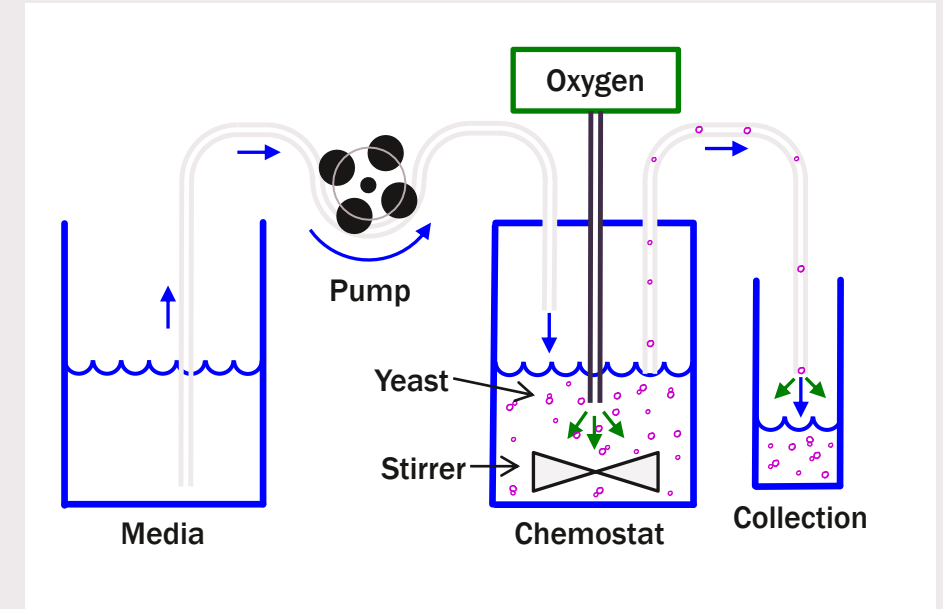# Evaluating models: predicting single-gene essentiality

- Which of the ~6000 genes are necessary for healthy growth?
- We assume the presence of extracellular metabolites that correspond to YNB (yeast nitrogen base) plus glucose
- Systematically remove a gene from the model by negating the relevant clause

    **g(g_to_knock_out) becomes ~g(g_to_knock_out)**

- Then query for the essential metabolites

| Base GEM | Yeast 8 |
|---|---|
| Number of predictions (number of genes in model) | 1068 (1150) |
| NG Recall (NGNG / NG) | 0.16 (25/161) |
| NG Precision (NGNG / (NGNG + NGG)) | 0.31 (25/81) |
| GNG Rate (GNG / NG) | 0.845 (136/161) |
| NGG Rate (NGG / G) | 0.062 (56/907) |
| F1 score | 0.207 |

# Evaluating models: constraining flux balance simulations

Stoichiometric matrix – $S$

Fluxes – $\boldsymbol{\nu}$

$$\underset{\boldsymbol{\nu} \in \mathbb{R}^n}{\text{maximize}} \quad f(\nu_1, \ldots, \nu_n)$$

$$\text{subject to} \quad S\boldsymbol{\nu} = \mathbf{0},$$

$$\nu_i^{\text{LB}} \leq \nu_i \leq \nu_i^{\text{UB}}, \quad i = 1, \ldots, n.$$



"Is the pathway for production of pyruvate from a minimal glucose medium feasible?"

iProver

reference flux distribution

1.45e-13
distance from reference flux distribution

"Which compounds need to be added to a minimal glucose medium to enable production of malonyl CoA?"

iProver

hypotheses

cnf(abduced_extra_001,axiom,met(m_glutathione,c_e)).

reference flux distribution

6.51
distance from reference flux distribution

The original discovery problem

How we formulated the problem in computational terms

What data and knowledge we provided to our system

How we represented the system's inputs and outputs

The space of candidate models that the system searched

What criteria it used to evaluate candidate models

How we interpreted results that the system generated

# Interpreting results

- Can use automated theorem prover (iProver) for deduction but also abduction—has the efficiency required for models of this size

- Finding models is easy, finding good models is hard

- Using multiple deduction techniques can help check model consistency—bridge the divide

- Conflicting evidence in literature—e.g. pathway for L-arginine production

- Limits to what one can do with others' data

# Future work

- Hypothesis testing with Genesis platform
- Integration of metabolomics and transciptomics measurements
- Learn rules for gene expression and regulation
- Integrate with signalling pathways

# Acknowledgements

Konstantin Korovin

King Lab

Erik Bjurström
Daniel Brunnsåker
Filip Kronström
Praphapan Lasin
Gabriel Reder
Ievgeniia Tiukova
Ross D. King

Vanderbilt University

John Wikswo
Ron Reiserer
Clayton Britt
Greg Gerken
Kyle Hawkins
Dmitry Markov
Philip Samson
David Schaffer
Eric Spivey
Erik Werner

Larisa N. Soldatova

Thoughtworks

Rushikesh Halle
Vinay Mahamuni
Amit Patel
Harshal Hayatnagarkar

…and more!

# Areas of improvement

| Challenge - *from Chen Y, Li F, Nielsen J (2022)\** | Hypotheses | Techniques to test hypotheses |
|---|---|---|
| Annotation of the model | Correct formulae, charge etc. for metabolites; enzyme numbers; reaction directionality | Mass balance analysis; ?? |
| Noise from low-confidence components | Removal of certain elements from model; additional dimension of model "explains" noise | ; condition-specific evaluation techniques |
| "Dead-end" metabolites | Add reactions | Predictive accuracy of metabolic activity with/without reactions; thermodynamic balance analysis |
| Un- or poorly-annotated reactions (in particular transport reactions) | New or changed gene-protein rule | Comparing transcriptomics or proteomics data with metabolic activity; mutant strain cultivation |
| Changes to biomass equation itself | Coefficient change; variable addition or removal; condition-specific biomass equations | Comparative prediction analysis |
| Enzyme turnover rate estimation | Values for enzyme turnover rates | Experimentally measure enzyme levels; comparitive prediction analysis |
| Integration of **subcellular constraints** | Reaction constraints e.g. within mitochondrion | Quantification of sub-cellular proteomes; ?? |
| Integration of regulation mechanisms | More precise mathematical formulations of regulatory mechanisms; better models of currently understood mechanisms of regulation | Comparative prediction analysis; multi-omics analysis; mutant strain cultivation |

2023-03-28

# Role of automation

Automated laboratory processes and AI techniques bring:

1. Efficiency

2. Broader reasoning (as opposed to deep reasoning right now)

3. Precision and repeatability

I would summarise by saying the strength of automation (in scientific discovery for S. *cerevisiae* models) will be in embracing complexity by **avoiding excessive simplification** during reasoning and exploiting big datasets to **extract small signals from large noise**, particularly when it comes to testing condition-specific models.