# Phenomics validation of the Escherichia coli underground metabolic reconstruction

**Claudio Angione**

UNIVERSITY OF
CAMBRIDGE

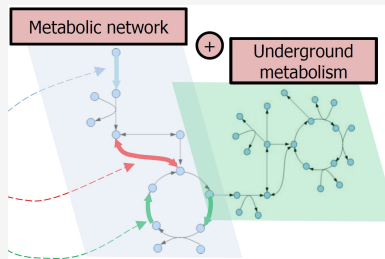## Studying a physiological system in silico [Mo and Palsson, 2009]



- Genome-scale models
- Huge variable space and objective space
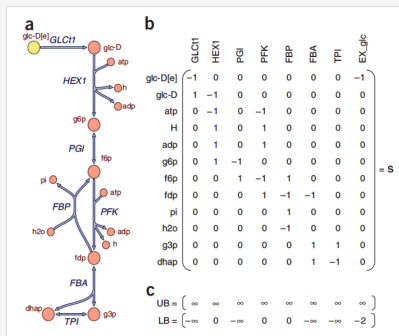- Computationally expensive to explore

# Underground metabolism

- Metabolism = set of chemical reactions taking place in living cells, with the aim of maintaining cellular functions

- Once considered only a passive result of the state of a cell

- Now widely recognized as a main contributor to cell behavior

Underground metabolism = reactions catalyzed with less catalytic efficiency, result of the weak side activity of preexisting enzymes

- 1380 genes

- 3027 reactions

- 2151 metabolites



Metabolic network (+) Underground metabolism

# Mathematical interpretation of FBA [Orth et al., 2010]



For every metabolite $X_i$, $i = 1, \ldots, m$ a material balance is $\dfrac{dX_i}{dt} = \sum_{j=1}^{n} S_{ij} v_j$.

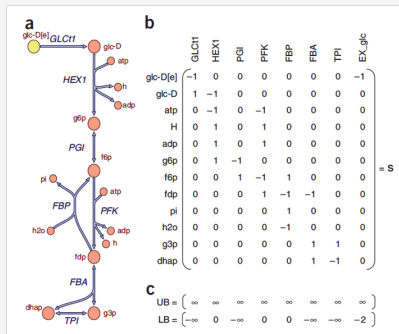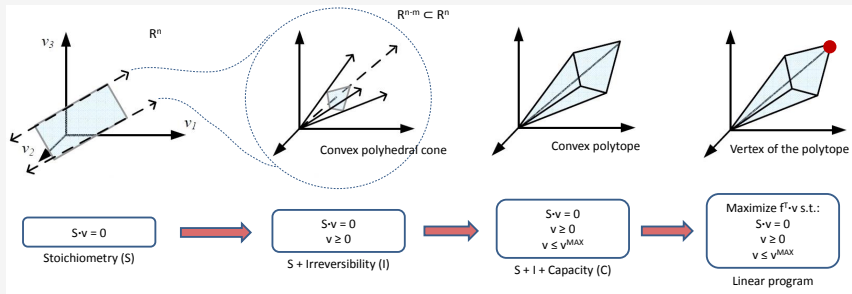FBA: Steady state conditions $dX/dt = 0 \Rightarrow Sv = 0$

maximise (or minimise) $\quad f^\mathsf{T} v$

such that $Sv = 0$ $\hspace{3cm}$ (1)

$V_j^{min} \leq v_j \leq V_j^{max}, \quad j = 1, \ldots, n$

# Mathematical interpretation of FBA [Orth et al., 2010]



For every metabolite $X_i$, $i = 1, \ldots, m$ a material balance is $\dfrac{dX_i}{dt} = \displaystyle\sum_{j=1}^{n} S_{ij} v_j$.

FBA: Steady state conditions $dX/dt = 0 \Rightarrow Sv = 0$

maximise (or minimise) $\quad f^\mathsf{T} v$

such that $Sv = 0$ (1)

$V_j^{min} \leq v_j \leq V_j^{max}, \quad j = 1, \ldots, n$

# Geometric interpretation of FBA [Orth et al., 2010]



**Mathematical formulation of flux balance analysis.** The stoichiometric matrix $S$ of $n$ reactions and $m$ metabolites restricts the search of possible flux distributions to the hyperplane $\mathbb{R}^{n-m}$. Thermodynamic constraints (irreversibility of reactions) limit the space of feasible solutions, which becomes a polyhedral cone. Capacity constraints (enzyme or transport capacities) constitute an upperbound for the flux rates; if this is available for every flux in the network (more specifically, it is sufficient that the upper bound is available for the edges of the cone), the feasible space reduces to a convex polytope. Once an objective function has been defined, the final linear program finds the final flux distribution as a vertex of the convex polytope, and then reconstructs the solution in the initial space $\mathbb{R}^n$.
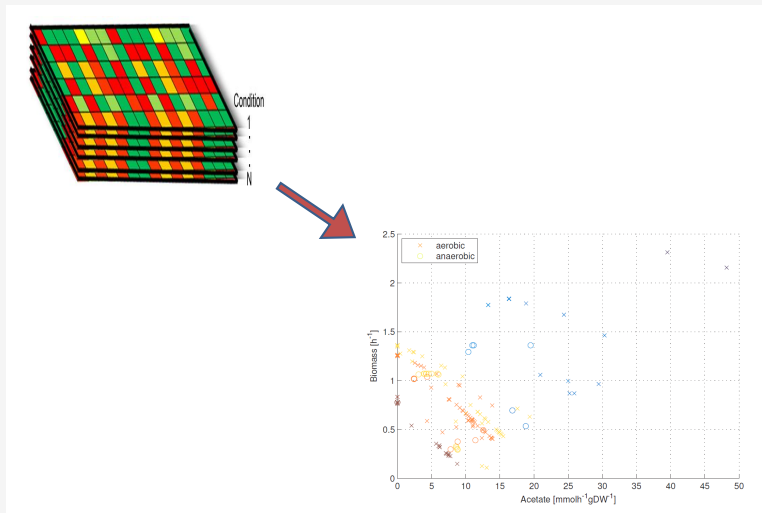
# FBA: pool in a waterfall



Input fluxes

Volume of pool of water = metabolite concentration
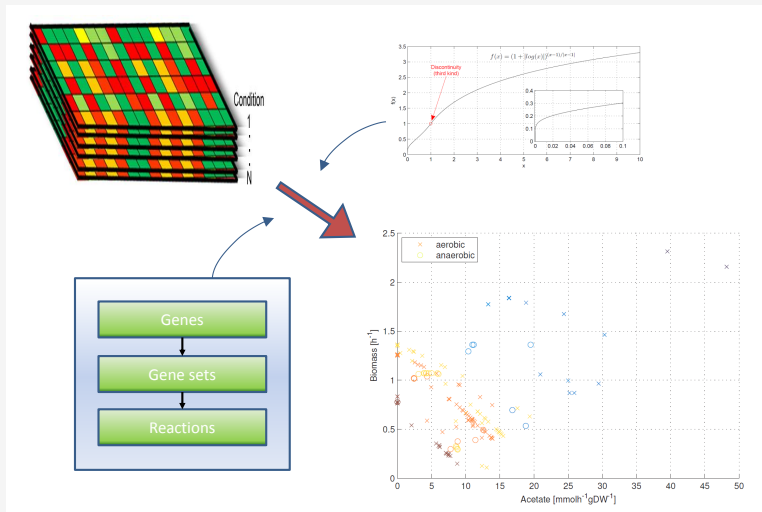
Output fluxes

**Slide Credit: Jeremy Zucker**

# Gene expression profiles to Pareto optimality

# Gene expression profiles to Pareto optimality



Condition = gene expression profile

## Redefining Constraints for Linear Programming

The bounds of the fluxes depend on the gene expression:

$$
\begin{aligned}
\max \quad & g^{\mathsf{T}} v \\
\text{such that} \quad \max \quad & f^{\mathsf{T}} v \\
\text{such that} \quad & Sv = 0 \\
& v_i \geq V_i^{min} h(y_i) \\
& v_i \leq V_i^{max} h(y_i)
\end{aligned}
\tag{2}
$$

$$
h(y_i) = \begin{cases} (1 + |log(y_i)|)^{\text{sgn}(y_i - 1)} & \text{if } y_i \in \mathbb{R}^+ \setminus \{1\} \\ 1 & \text{if } y_i = 1 \end{cases}
\tag{3}
$$

where $\text{sgn}(y_i - 1) = (y_i - 1)/|y_i - 1|$. $y_i$ is the gene set expression of the $i$th gene set, responsible for the $i$th reaction of the model.

The bilevel problem can be converted to a single-level problem using linear programming duality theory: for every linear programming problem (primal) there exists a unique optimization problem (dual) whose optimal objective value is equal to that of the primal problem.

## Redefining Constraints for Linear Programming

The bounds of the fluxes depend on the gene expression:

$$
\begin{aligned}
\max \quad & g^{\mathsf{T}} v \\
\text{such that} \quad \max \quad & f^{\mathsf{T}} v \\
\text{such that} \quad & Sv = 0 \\
& v_i \geq V_i^{min} h(y_i) \\
& v_i \leq V_i^{max} h(y_i)
\end{aligned}
\tag{2}
$$

$$
h(y_i) = \begin{cases} (1 + |log(y_i)|)^{\mathsf{sgn}(y_i - 1)} & \text{if } y_i \in \mathbb{R}^+ \setminus \{1\} \\ 1 & \text{if } y_i = 1 \end{cases}
\tag{3}
$$

where $\mathrm{sgn}(y_i - 1) = (y_i - 1)/|y_i - 1|$. $y_i$ is the gene set expression of the $i$th gene set, responsible for the $i$th reaction of the model.

The bilevel problem can be converted to a single-level problem using linear programming duality theory: for every linear programming problem (primal) there exists a unique optimization problem (dual) whose optimal objective value is equal to that of the primal problem.

## Redefining Constraints for Linear Programming

The bounds of the fluxes depend on the gene expression:

$$
\begin{aligned}
\max \quad & g^{\mathsf{T}} v \\
\text{such that} \quad & \max \quad f^{\mathsf{T}} v \\
& \text{such that} \quad Sv = 0 \\
& \qquad\qquad v_i \geq V_i^{min} h(y_i) \\
& \qquad\qquad v_i \leq V_i^{max} h(y_i)
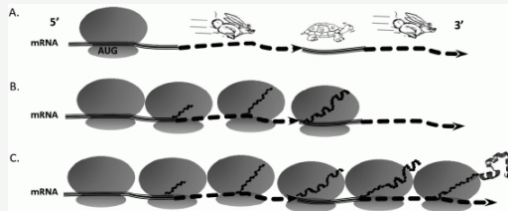\end{aligned} \tag{2}
$$

$$
h(y_i) = \begin{cases} (1 + |log(y_i)|)^{\mathsf{sgn}(y_i - 1)} & \text{if } y_i \in \mathbb{R}^+ \setminus \{1\} \\ 1 & \text{if } y_i = 1 \end{cases} \tag{3}
$$

where $\mathrm{sgn}(y_i - 1) = (y_i - 1)/|y_i - 1|$. $y_i$ is the gene set expression of the $i$th gene set, responsible for the $i$th reaction of the model.

The bilevel problem can be converted to a single-level problem using linear programming duality theory: for every linear programming problem (primal) there exists a unique optimization problem (dual) whose optimal objective value is equal to that of the primal problem.

# Codon usage bias affects protein production

- Codon = mRNA nucleotide triplet
- In a protein, each amino acid is encoded by up to six synonymous codons
- $4^3$ codons, 61 of which actually encode for amino acids plus 3 stop codons
- but only 20 different translated amino acids
- some codons are slow, some are fast in the translation process
- therefore, the choice between these codons affects the translation rate and therefore the final amount of protein produced
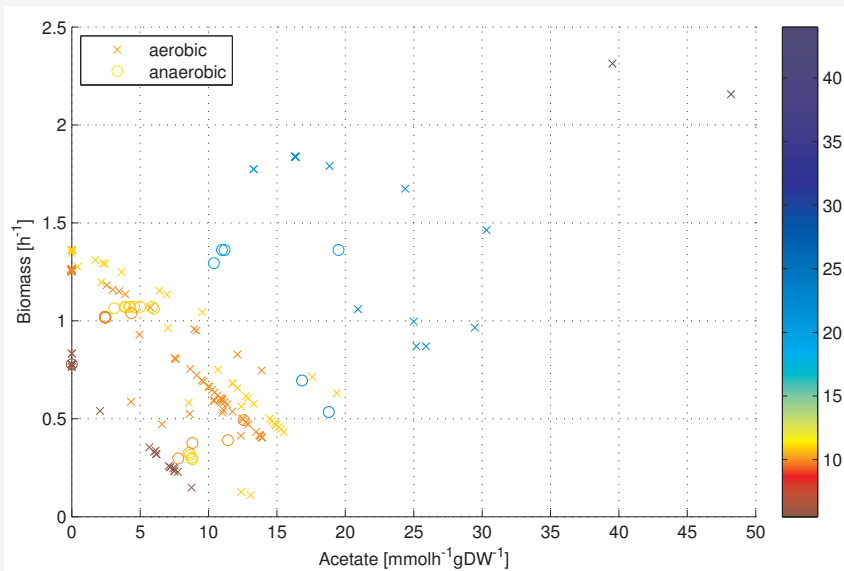
*Applications of the multi-omic model*

## Why mapping gene expression profiles to a metabolic model?

Goal: analyzing (e.g. clustering) gene expression profiles using metabolism, rather than performing the analysis directly on gene expression data.
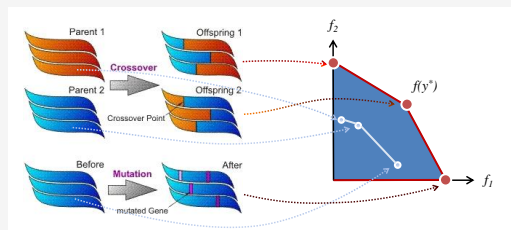
1 Multi-omic models cluster genes in their relative pathway; pathways can then be clustered and ranked through an effect-based approach (i.e., looking at the output outcomes in the phenotypic space)

2 The model acts as a ranking and noise-reduction tool: effect of low-importance genes is filtered out even if their expression is highly variable across conditions; without the multi-omic model, these genes would be incorrectly regarded as key genes to differentiate conditions

3 Performing inference directly on gene expression values may lead to incorrect prediction of the centrality of a gene whose level seems to be highly correlated with many other genes, but with only a marginal role in the metabolism (e.g., no impact on the biomass and on key metabolites)

# 466 *E. coli* experimental conditions

# Multi-objective optimization

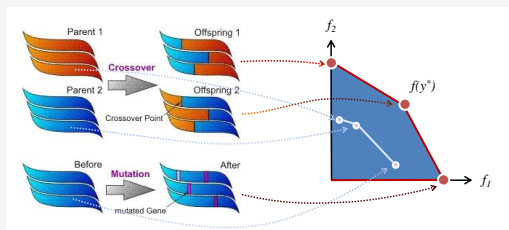Let $f$ be the vector of $r$ objective functions to optimize in the objective space



- Solution of a multi-objective problem: set of points called Pareto front

- Represents the best trade-off between two or more requirements

- A point $y^*$ in the solution space is Pareto optimal if there does not exist a point $y$ such that $f(y)$ dominates $f(y^*)$, i.e.
$\nexists\, y$ s.t. $f_i(y) > f_i(y^*), \forall i = 1, ..., r$

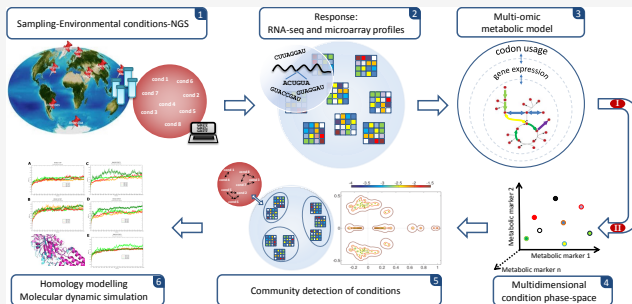[Angione et al. - Theoretical Computer Science, 2015]

# Multi-objective optimization

Let $f$ be the vector of $r$ objective functions to optimize in the objective space



- Solution of a multi-objective problem: set of points called Pareto front

- Represents the best trade-off between two or more requirements

- A point $y^*$ in the solution space is Pareto optimal if there does not exist a point $y$ such that $f(y)$ dominates $f(y^*)$, i.e.
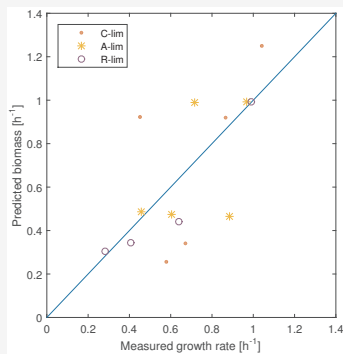$\nexists\ y$ s.t. $f_i(y) > f_i(y^*), \forall i = 1, ..., r$

[Angione et al. - Theoretical Computer Science, 2015]

# METRADE: the common framework



- Multi-omic model to associate conditions to a phenotypic outcome in a set of objective spaces
- Microarray expression profiles (environmental conditions) are mapped to a multi-omic model of metabolism
- A spectral method for community detection infers condition similarities according to the metabolic response in the multi-objective space.
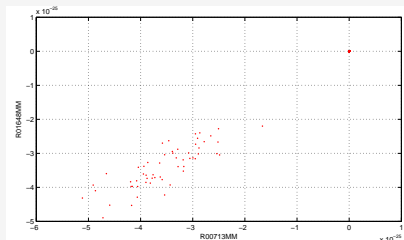
[Angione and Lió - Nature Scientific Reports, 2015 (in press)]
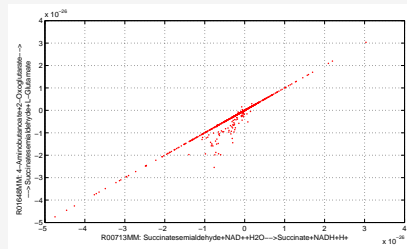
## Validation of METRADE in *E. coli* underground metabolism



- Pearson's $r = 0.680$ (*p*-value $= 0.007$)

- Spearman's $\rho = 0.678$ (*p*-value $= 0.008$)

*Predictions of biological processes through multi-omic models*

# Mitochondrial diseases (identifiability analysis)
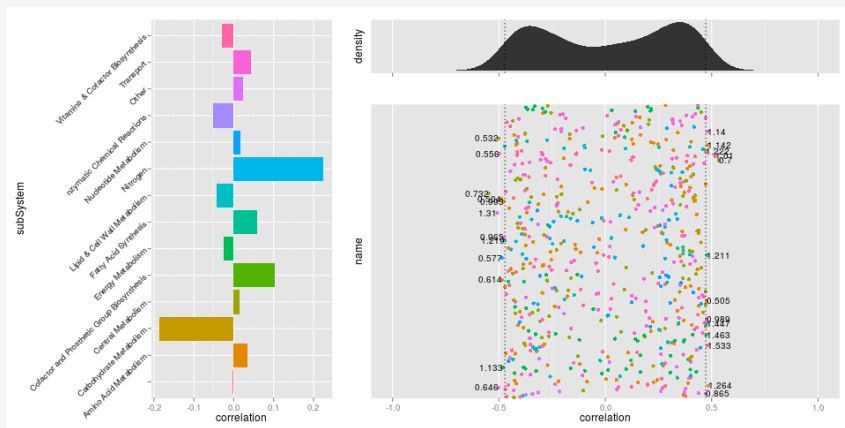


(a) Fumarase deficiency



(b) Succinate dehydrogenase deficiency

- Monogenic diseases
- Inferring functional relations between flux rates
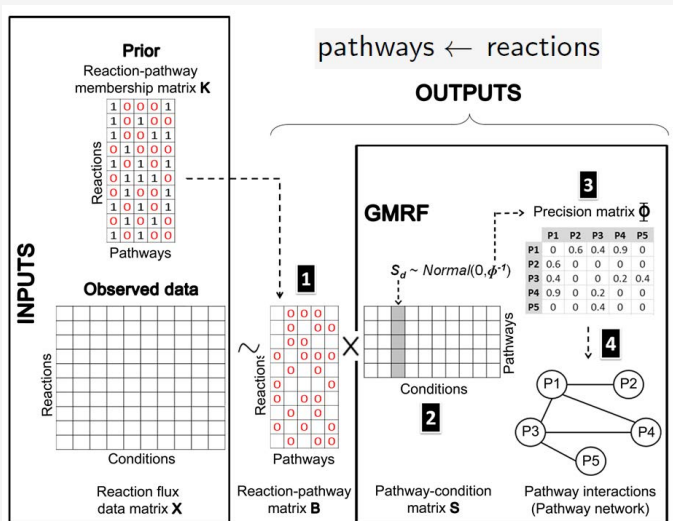- Predicting the type of monogenic disorder through the shape of the functional relation

[Angione et al. - PLoS One, 2015 (in press)]

# Geobacter sulfurreducens - correlation gene expression-position on the Pareto front



[Conway, Angione, Lió - Current Bioinformatics, 2015 (in press)]

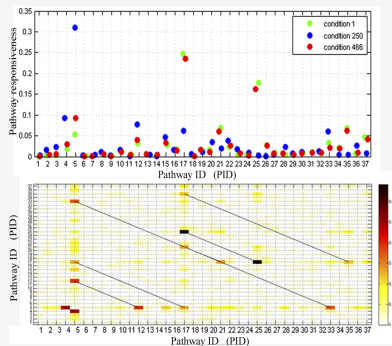# Probabilistic Matrix Factorization with Gaussian Markov Random Fields



[Angione, Pratanwanich, Lió - ACS Synthetic Biology, 2015]
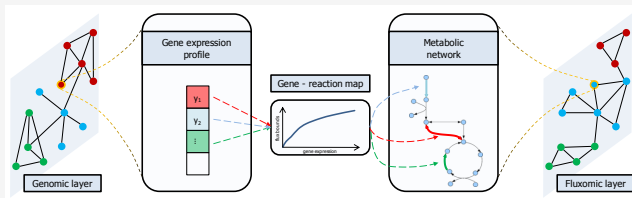
## Bayesian pathway analysis in *E. coli*

- Comparison of pathway responsiveness to different conditions (top)

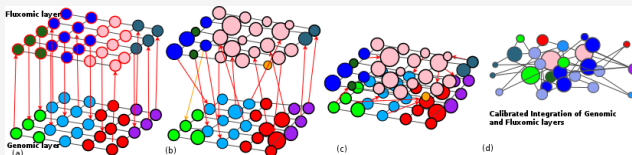- Crosstalks between pathways can be viewed as a correlation matrix (bottom).



| High glucose | | | Low glucose | |
|---|---|---|---|---|
| PID | Pathway | Aerobic | Anaerobic | Aerobic | Anaerobic |
| 5 | Nucleotide salvage | 0.0865 | 0.0965 | 0.1366 | 0.1714 |
| 17 | Valine, leucine, and isoleucine metabolism | 0.2219 | 0.2147 | 0.1974 | 0.1590 |
| 25 | Alanine and aspartate metabolism | 0.1544 | 0.1487 | 0.1285 | 0.1076 |

- Average responsiveness of the most responsive pathways across aerobic and anaerobic conditions of high and low glucose in *E. coli*.

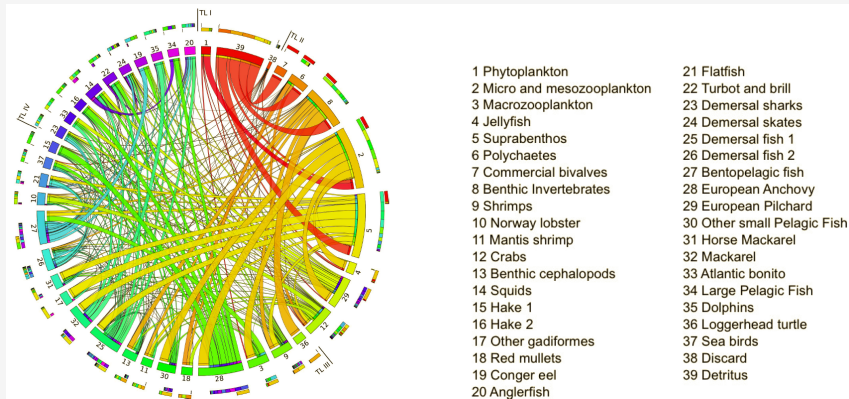# Weighted integration of multi-omic layers of conditions



Integration of omic layers (genomics, metabolomics, fluxomics) using multiplex/multilayer network theory
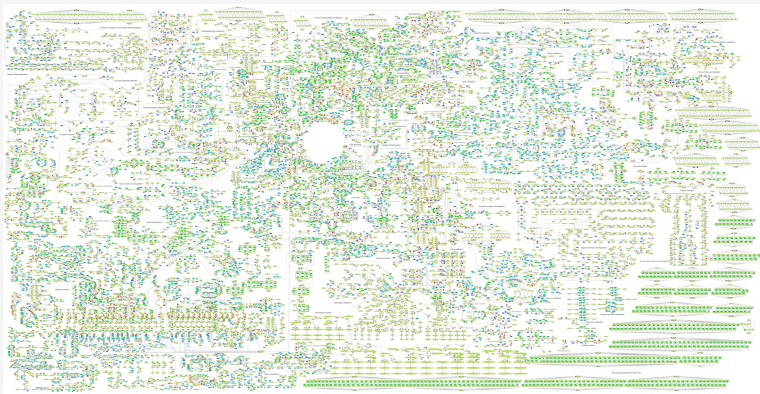


[Angione, Conway, Lió - BMC Bioinformatics, under review]
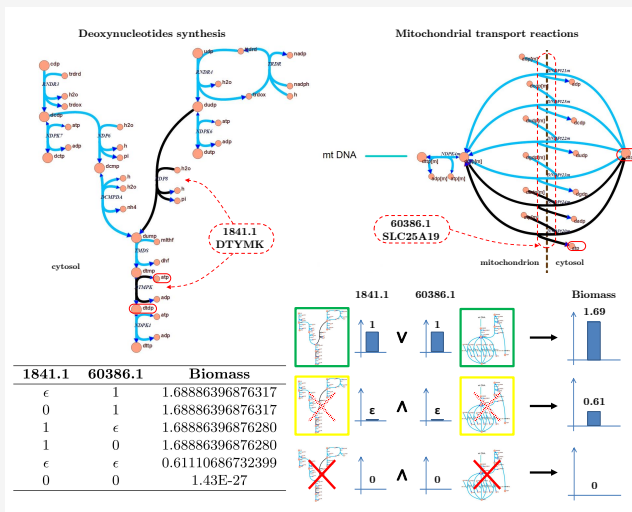
# Adriatic Sea food web: PCB bioaccumulation network



1 Phytoplankton
2 Micro and mesozooplankton
3 Macrozooplankton
4 Jellyfish
5 Suprabenthos
6 Polychaetes
7 Commercial bivalves
8 Benthic Invertebrates
9 Shrimps
10 Norway lobster
11 Mantis shrimp
12 Crabs
13 Benthic cephalopods
14 Squids
15 Hake 1
16 Hake 2
17 Other gadiformes
18 Red mullets
19 Conger eel
20 Anglerfish

21 Flatfish
22 Turbot and brill
23 Demersal sharks
24 Demersal skates
25 Demersal fish 1
26 Demersal fish 2
27 Bentopelagic fish
28 European Anchovy
29 European Pilchard
30 Other small Pelagic Fish
31 Horse Mackarel
32 Mackarel
33 Atlantic bonito
34 Large Pelagic Fish
35 Dolphins
36 Loggerhead turtle
37 Sea birds
38 Discard
39 Detritus

[Taffi et al. - Frontiers in Genetics, 2014]
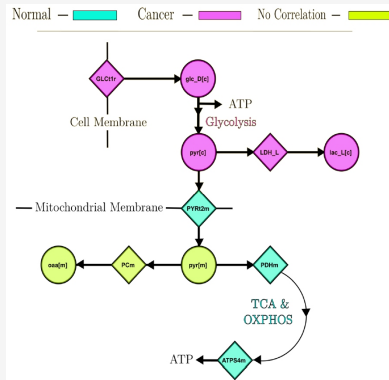
## Human metabolism



- 8 cellular compartments
- $7,440$ reactions
- $1,789$ enzyme-encoding genes
- $2,626$ unique metabolites

# Genome-scale exploration of simultaneous gene effects

## Human cancer metabolism

- We map METABRIC to the human metabolism using METRADE

- We independently find the Warburg effect

- 8 new biomarkers (reactions) significantly different in cancer/normal conditions

- All these new biomarkers have been previously associated with breast cancer

## Contacts



claudio.angione@cl.cam.ac.uk

`https://www.cl.cam.ac.uk/users/ca394/`