

# Learning metabolic regulatory rules from time series data

***MERRIN:** MEtabolic Regulation Rules INference from time series data<sup>\*</sup>*



Kerian Thuillier<sup>1</sup>, Caroline Baroukh<sup>2</sup>, Alexander Bockmayr<sup>3</sup>,  
Ludovic Cottret<sup>2</sup>, Loïc Paulevé<sup>4</sup>, Anne Siegel<sup>1</sup>

<sup>1</sup> Univ Rennes, Inria, CNRS, IRISA, Rennes, France

<sup>2</sup> LIPME, INRAe, CNRS, Université de Toulouse, Castanet-Tolosan, France

<sup>3</sup> Freie Universität Berlin, Institute of Mathematics, D-14195 Berlin, Germany

<sup>4</sup> Univ. Bordeaux, Bordeaux INP, CNRS, LaBRI, UMR5800, F-33400 Talence, France

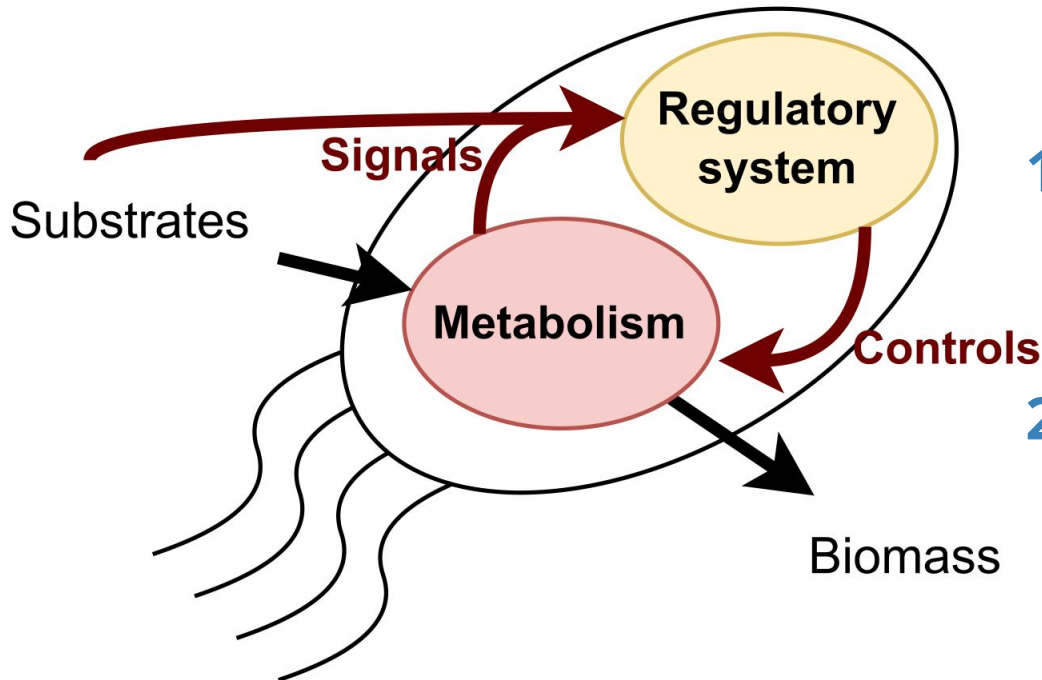
26th March 2024

---

<sup>\*</sup> K. Thuillier et al., **Oxford Bioinformatics**, 2022

# Cells: hybrid multi-layered structures

Model as two interconnected systems



## 1. Metabolic system

*Chemicals reactions converting substrates to energy and biomass*

## 2. Regulatory system

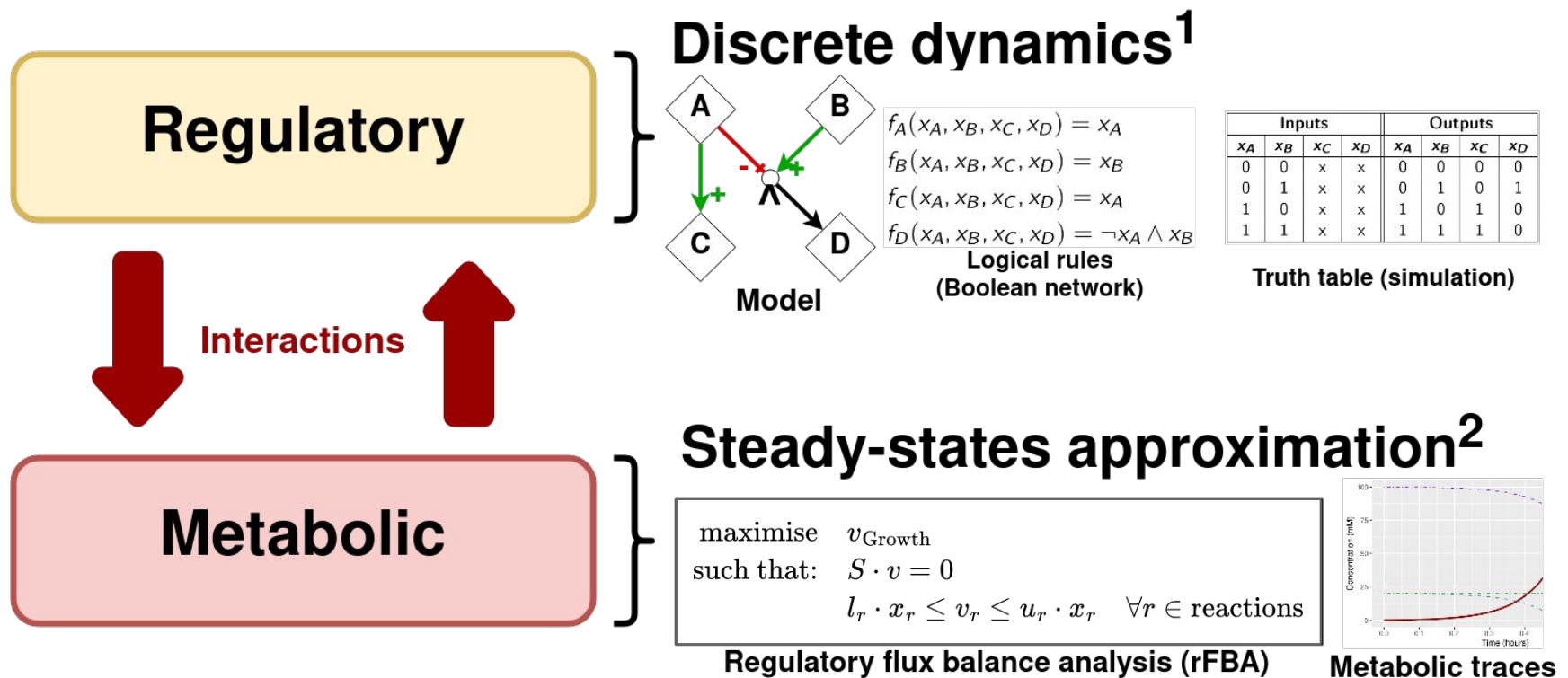
*Rules constraining the metabolism to adapt itself to its environment*

### Objective:

Inferring the **regulatory systems** from time series observations of the cells  
(*metabolism and regulation*)

# Multiplicity of modelling formalisms

Two systems models with different dynamics

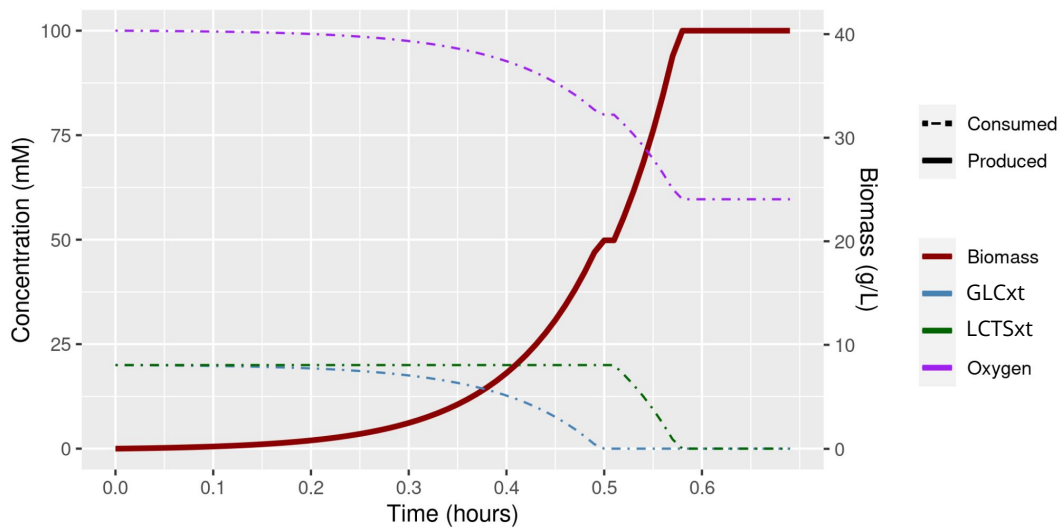


<sup>1</sup> S. Videla et al., **Bioinformatics**, 2016

<sup>2</sup> M. W. Covert et al., **Journal of theoretical biology**, 2001

# Example: diauxic shift (*Monod et al., 1953*)

Evolution of cell biomass



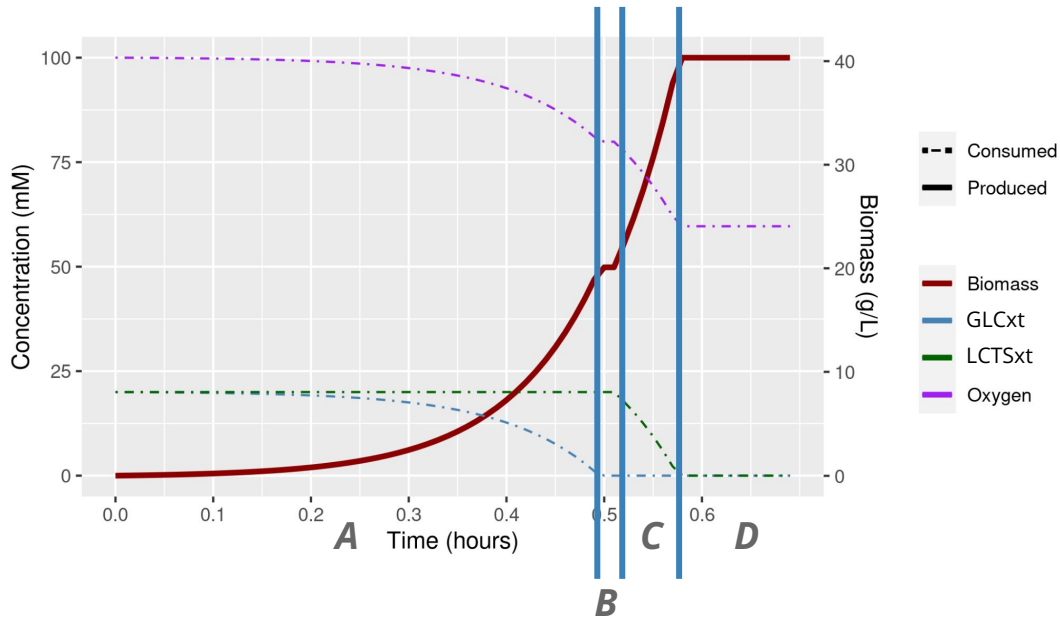
## Diauxic shift

→ Successive growth phases on different mediums

→ Control by regulations

# Example: diauxic shift *(Monod et al., 1953)*

Evolution of cell biomass



## Diauxic shift

→ Successive growth phases on different mediums

→ Control by regulations

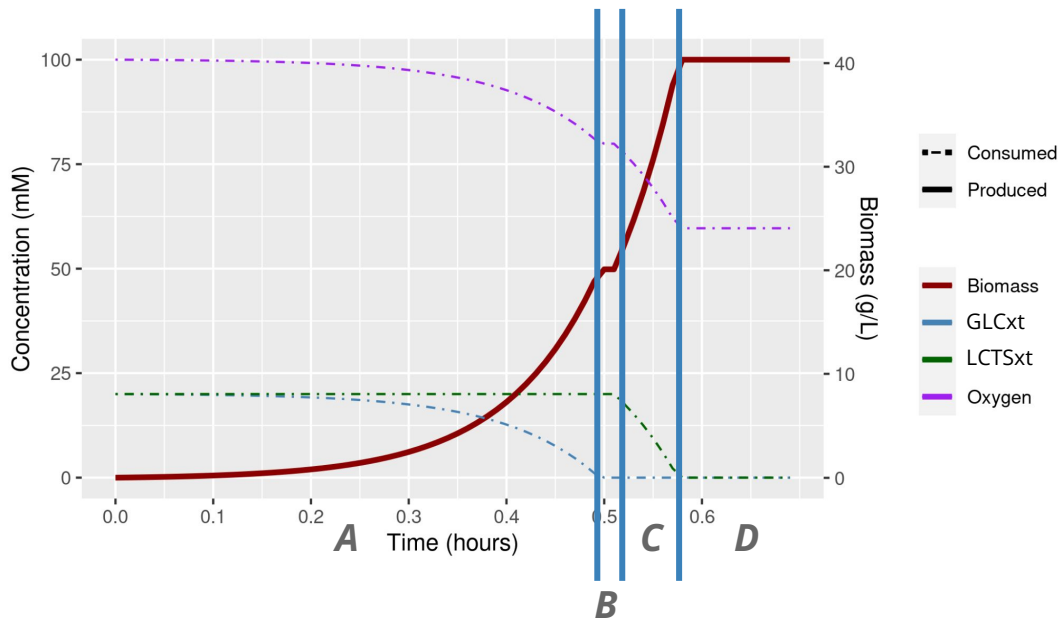
## Divide in 4 phases

Characterize by different qualitative behaviours (e.g. growth medium)

$A \rightarrow \text{Growth}$	$B \rightarrow \text{No growth}$
$C \rightarrow \text{Growth}$	$D \rightarrow \text{No growth}$

# Example: diauxic shift *(Monod et al., 1953)*

Evolution of cell biomass



## Diauxic shift

→ Successive growth phases on different mediums

→ Control by regulations

## Divide in 4 phases

Characterize by different qualitative behaviours (e.g. growth medium)

**A** → Growth

**B** → No growth

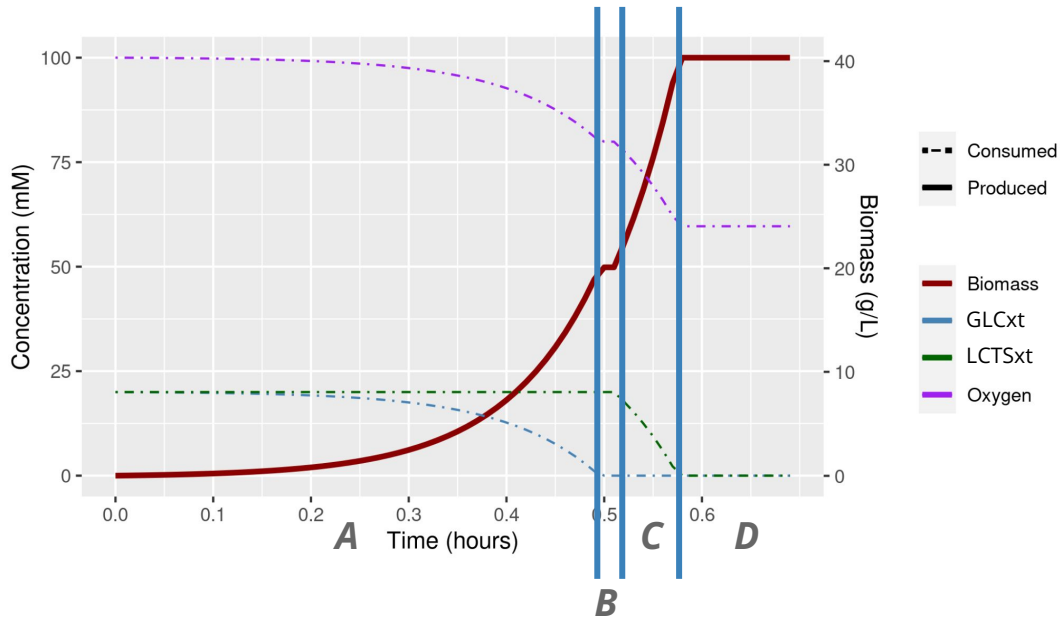
**C** → Growth

**D** → No growth

Both **regulatory system** and **metabolic system** dynamics must be considered to reproduce experimental observations

# Example: diauxic shift (*Monod et al., 1953*)

Evolution of cell biomass



How can we learn the regulatory rules from such observations?

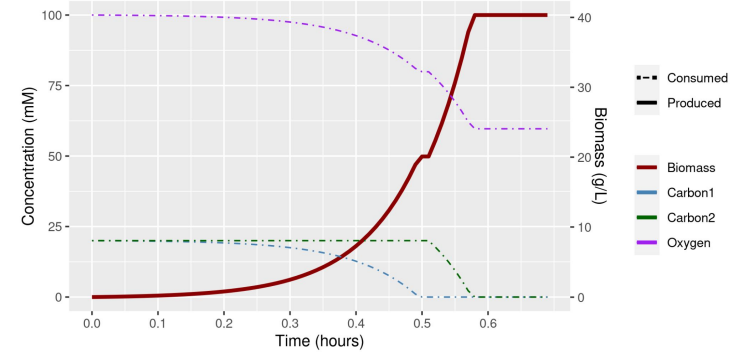
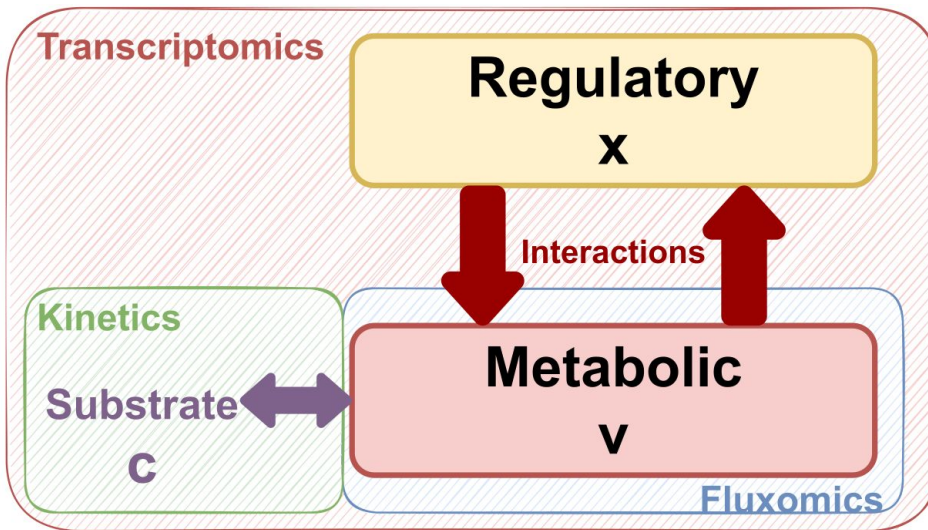
$$\begin{aligned}f_{\text{LacI}}(x) &= \neg x_{\text{LCTSxt}} \\f_{\text{GalR}}(x) &= \neg x_{\text{LCTSxt}} \\f_{\text{lacZ}}(x) &= \neg x_{\text{GLCxt}} \wedge \neg x_{\text{LacI}} \\f_{\text{galKTEU}}(x) &= \neg x_{\text{GLCxt}} \wedge \neg x_{\text{GalR}}\end{aligned}$$

Both **regulatory system** and **metabolic system** dynamics must be considered to reproduce experimental observations

# Time series data

## Observations of the *regulatory* and *metabolic* system activities

- *Quantitative* and *qualitative* measurements
- *Compatible* with observation from different mutant strains



GLCxt	LCTSxt	lacZ	galP	galKTEU	LacI	GalR
0	0	0	0	0	1	1

## 3 data types:

- **Transcriptomics** (qualitative)  
*Analysis of the RNA transcripts*
- **Fluxomics** (quantitative)  
*Rates of metabolic reactions*
- **Kinetics** (quantitative)  
*Substrate concentrations*

Compatible with any combination of those datatypes

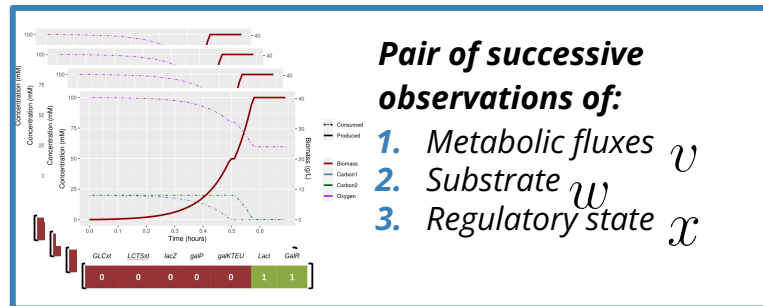


# Problems tackled by MERRIN

## Inputs:

**Metabolic**

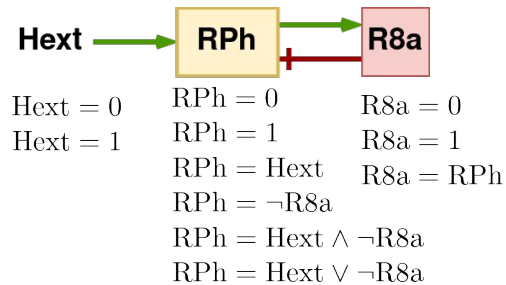
*Metabolic network  $\mathcal{N}$*



*Set of time series  $\{T_i\}_i$*

*(kinetics, fluxomics, transcriptomics)*

**Set of authorised interactions: activation and inhibition effects**



**36 compatible regulatory networks**

**$O(2^{2^n})$  in the number  $n$  of interactions**

**Prior Knowledge Network (PKN)**

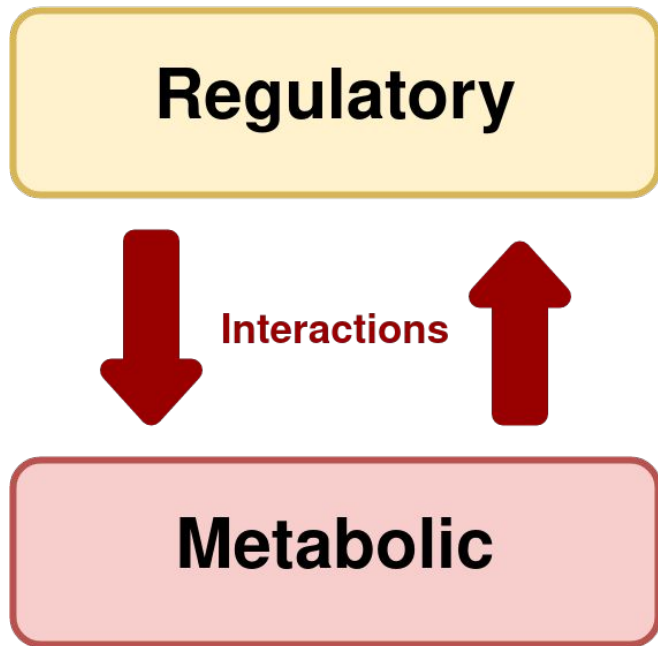
*Define a search space  $\mathcal{F}$*

## Outputs:

All the regulatory networks  $f \in \mathcal{F}$  compatible with the PKN and matching with the time series  $\{T_i\}_i$

# Underlying formalism

## Regulatory Flux Balance Analysis<sup>1</sup> – rFBA



rFBA timestep:

1. Update the **regulatory system**

*1 synchronous update  
of the Boolean network*

$$\begin{aligned} f_A(x_A, x_B, x_C, x_D) &= x_A \\ f_B(x_A, x_B, x_C, x_D) &= x_B \\ f_C(x_A, x_B, x_C, x_D) &= x_A \\ f_D(x_A, x_B, x_C, x_D) &= \neg x_A \wedge x_B \end{aligned}$$

2. Update the **metabolic system**

*Solve a FBA — LP problem*

$$\begin{aligned} &\text{maximise} && v_{\text{Growth}} \\ &\text{such that:} && S \cdot v = 0 \\ & && l_r \cdot x_r \leq v_r \leq u_r \cdot x_r \quad \forall r \in \text{reactions} \end{aligned}$$

3. Update the cell environment

<sup>1</sup> M. W. Covert et al., **Journal of theoretical biology**, 2001

# Inferring problem — *formal definition*

**Input:** metabolic network  $\mathcal{N}$ , PKN  $\mathcal{F}$ , set of time series  $\{T_i\}_i$

**Output:** all regulatory networks  $f \in \mathcal{F}$  such that:

$$\bigwedge_{T_i} \bigwedge_{(s,s') \in T_i} \left( f(x) = x' \right. \\ \left. \wedge \exists \hat{v} \in \mathbb{R}^{\text{Reactions}}, \left( S \cdot \hat{v} = 0 \wedge \bigwedge_{r \in \text{Reactions}} l_r \cdot x'_r \leq \hat{v}_r \leq u_r \cdot x'_r \wedge \hat{v}_{\text{growth}} \geq v'_{\text{growth}} - \epsilon \right) \right. \\ \left. \wedge \forall \hat{v} \in \mathbb{R}^{\text{Reactions}}, \left( S \cdot \hat{v} = 0 \wedge \bigwedge_{r \in \text{Reactions}} l_r \cdot x'_r \leq \hat{v}_r \leq u_r \cdot x'_r \right) \implies \hat{v}_{\text{growth}} \leq v'_{\text{growth}} + \epsilon \right)$$

**Hybrid problem: combinatorial + quantified linear constraints**

# Inferring problem — *formal definition*

**Input:** metabolic network  $\mathcal{N}$ , PKN  $\mathcal{F}$ , set of time series  $\{T_i\}_i$

**Output:** all regulatory networks  $f \in \boxed{\mathcal{F}}$  such that:

$$\bigwedge_{T_i} \bigwedge_{(s,s') \in T_i} \left( \boxed{f(x) = x'} \right) \quad \text{Boolean constraints}$$

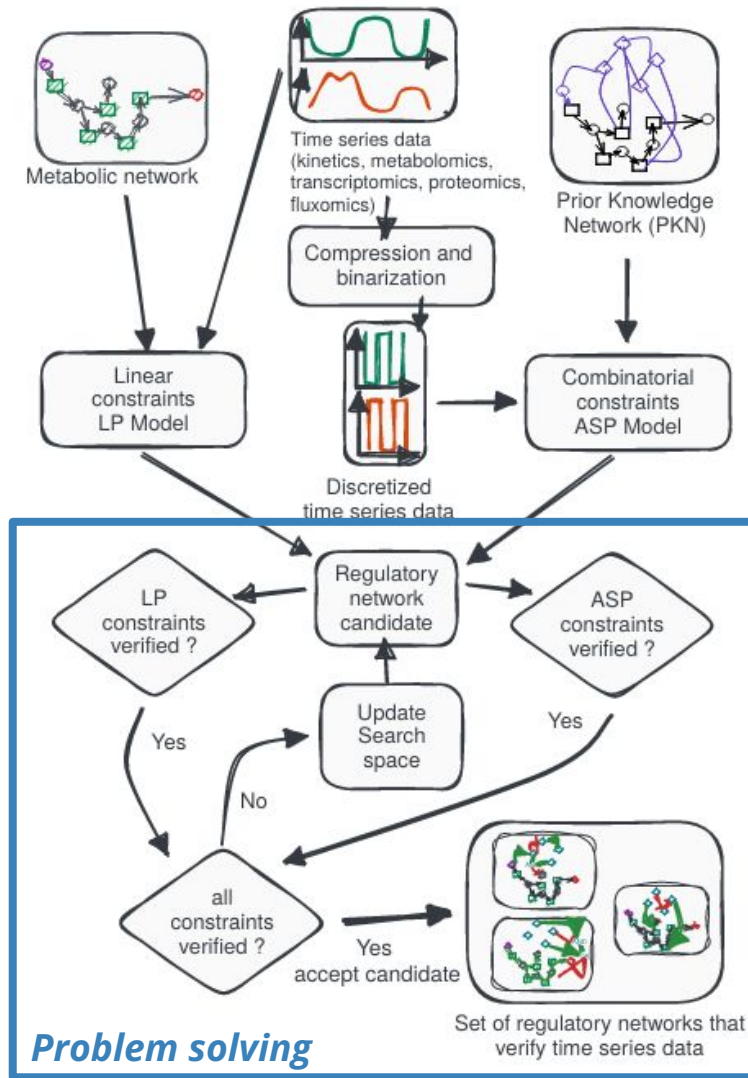
$$\bigwedge \exists \hat{v} \in \mathbb{R}^{\text{Reactions}}, \left( S \cdot \hat{v} = 0 \wedge \bigwedge_{r \in \text{Reactions}} l_r \cdot x'_r \leq \hat{v}_r \leq u_r \cdot x'_r \wedge \hat{v}_{\text{growth}} \geq v'_{\text{growth}} - \epsilon \right)$$

$$\bigwedge \forall \hat{v} \in \mathbb{R}^{\text{Reactions}}, \left( S \cdot \hat{v} = 0 \wedge \bigwedge_{r \in \text{Reactions}} l_r \cdot x'_r \leq \hat{v}_r \leq u_r \cdot x'_r \right) \implies \hat{v}_{\text{growth}} \leq v'_{\text{growth}} + \epsilon$$

Quantified linear constraints

**Hybrid problem: combinatorial + quantified linear constraints**

# Contribution: MERRIN's workflow

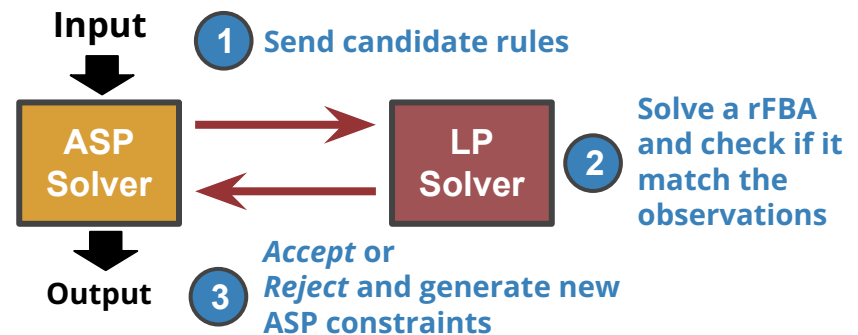


## Problem

**Input:** *Metabolic network, Prior Knowledge Network (PKN), Time series data*  
+ several solving parameters

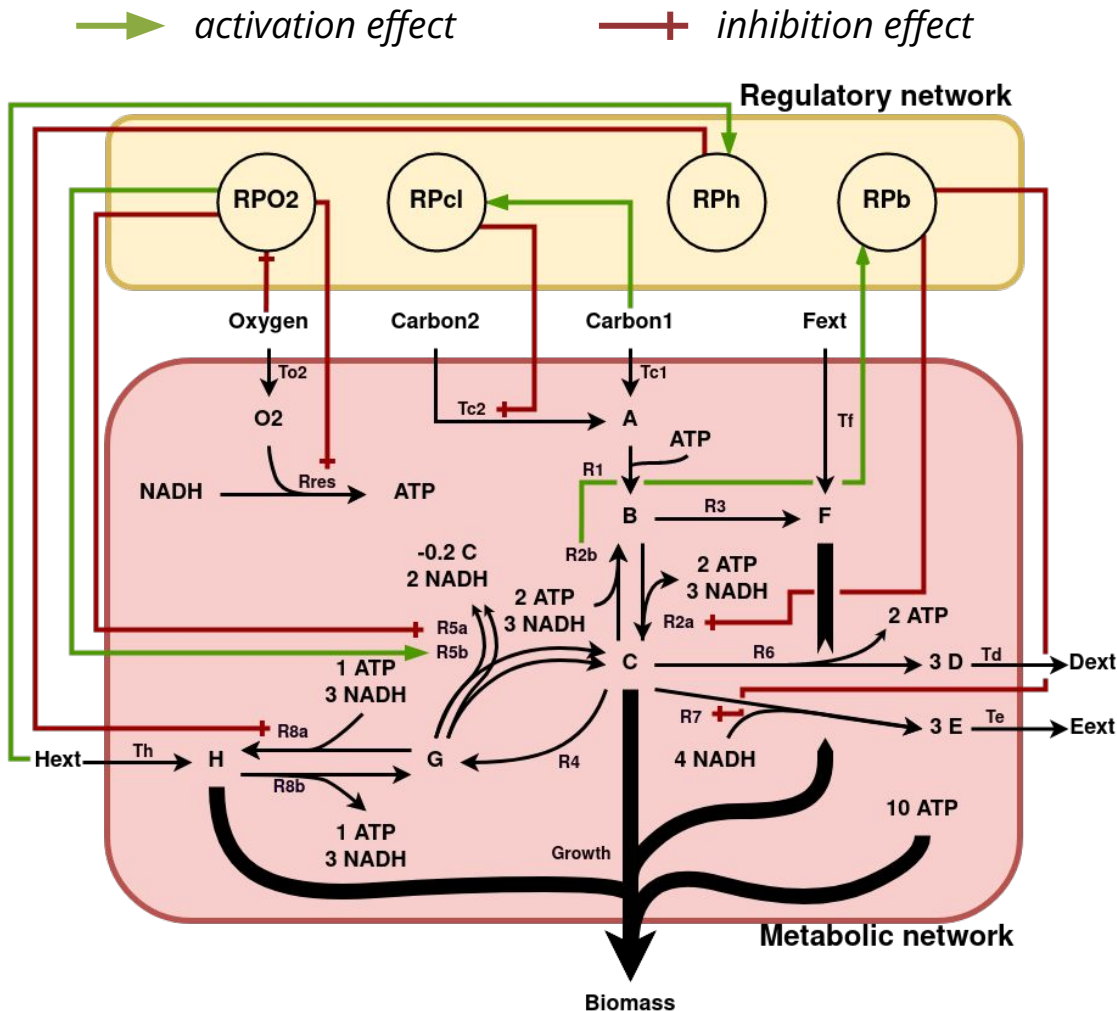
**Output:** *All the subset minimal regulatory metabolic network **satisfying the PKN** and **matching time series data***

## Rely on a hybrid resolution framework<sup>1</sup>



<sup>1</sup> K. Thuillier et al., **Proceedings of the AAAI Conference**, 2024

# Gold standard instance (Covert et al, 2001)

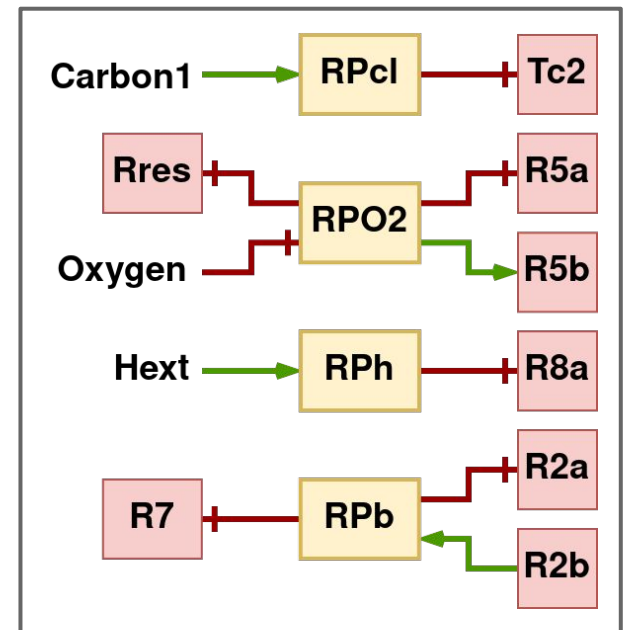


## Toy model based on *E.coli*

20 reactions, 4 regulatory protein, 11 regulations

## Model complex behaviours

Diauxic shift, aerobic/anaerobic growth, etc.



## Influence graph

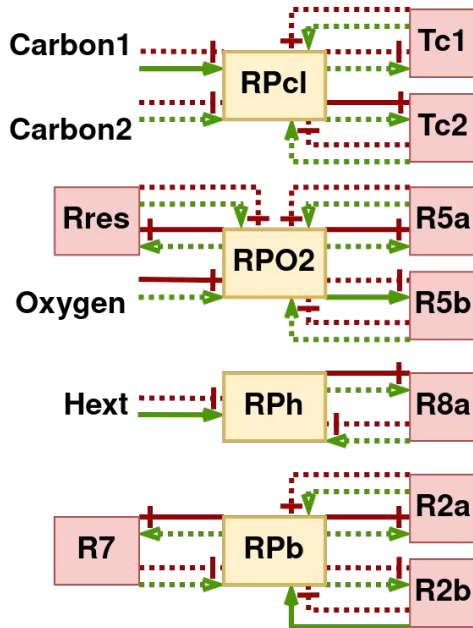
<sup>1</sup> M. W. Covert et al., **Journal of theoretical biology**, 2001

# Instance generation

## MERRIN inputs

### Prior Knowledge Network

Add hypothetical regulations (e.g. *RPcl* and *Tc1*)  
Remove sign + direction of interactions

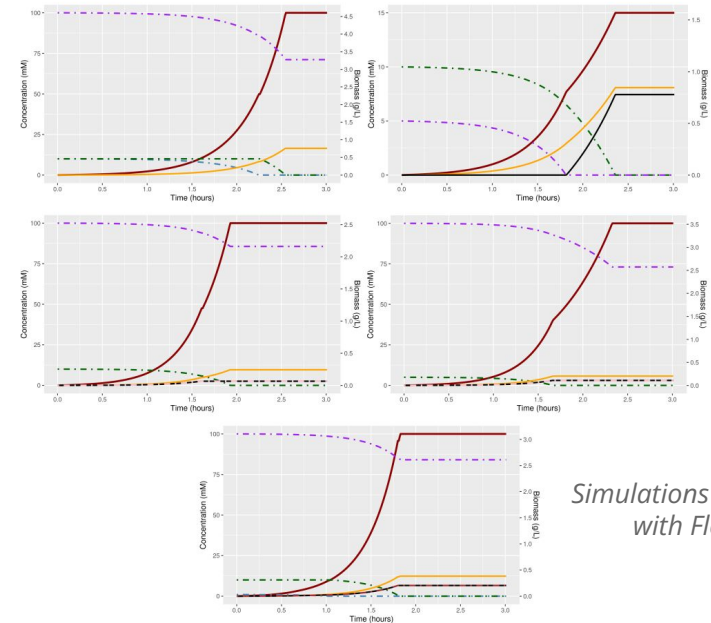


~2.9x10<sup>12</sup> BNs compatibles



### 5 simulations<sup>1</sup>

Kinetics, fluxomics and transcriptomics  
Perfect observations (no noise)

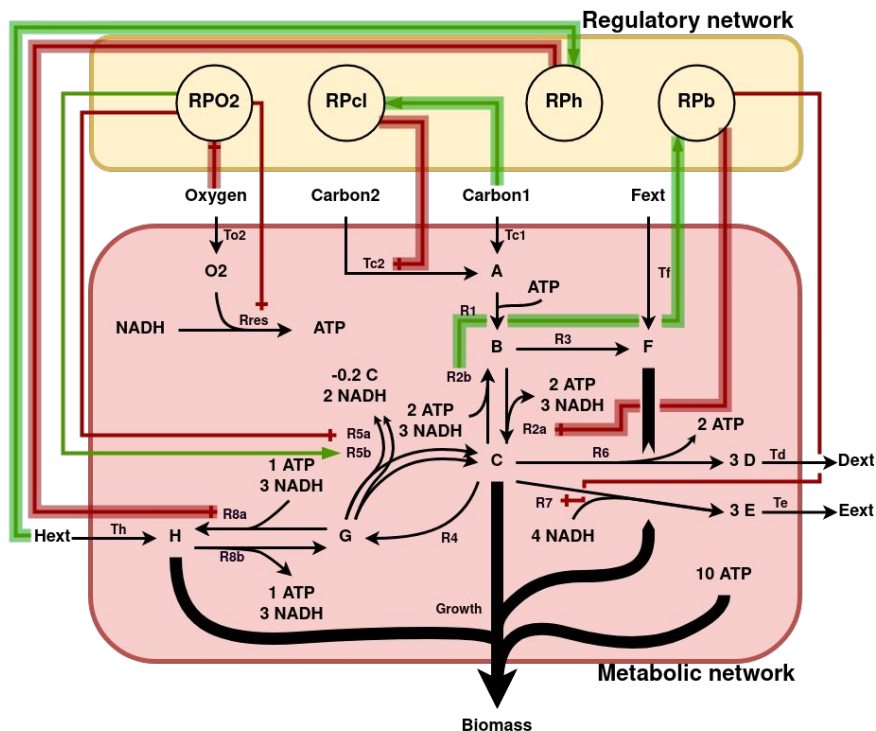


Simulations made with FlexFlux

240 instances with 6 degradation levels

<sup>1</sup> M. W. Covert et al., **Journal of theoretical biology**, 2001

# Validation and robustness testing



## Learn more parsimonious model than ground truth

- Reproduce exactly the input time series
- Unrecovered regulations can be explained

## Validation on a benchmark of 240 instances (*in silico*)

- 4 data types
- 6 level of degradations (0% to 50%)

## Perfectly reproduce the time series with:

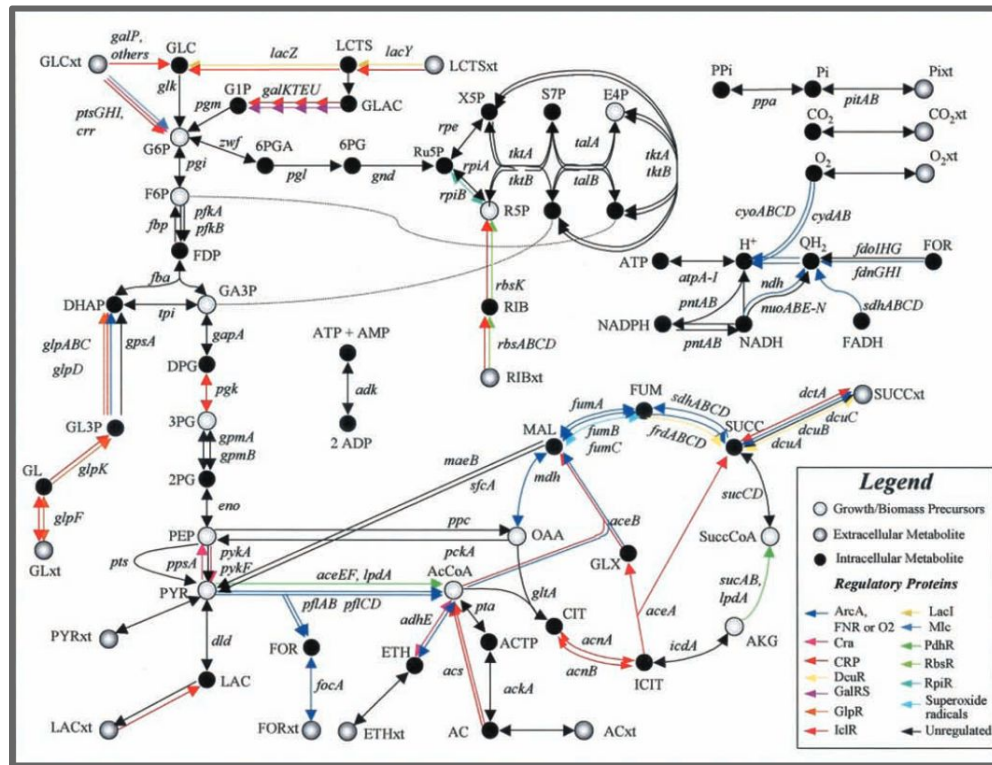
- *kinetics* and *transcriptomics* data
- < 20% of degradation

<sup>1</sup> M. W. Covert et al., **Journal of theoretical biology**, 2001



# Scalability of MERRIN

*E.coli* medium scale instance<sup>1</sup>



Metabolic network<sup>1</sup>

## Description:

- 60 regulatory rules
  - 19 regulatory proteins
  - 41 (regulated) genes
- 113 reactions
  - 66 are reversible

## 3 experimental conditions<sup>1</sup>

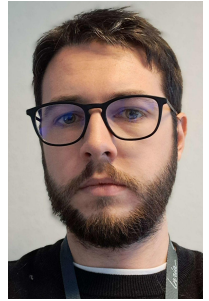
- rFBA time series *in silico*
- Mutant strains**

**Computation time: ~15 minutes**

**MERRIN scales on bigger models**

<sup>1</sup> M. W. Covert and B. Ø. Palsson, **Journal of biological chemistry**, 2002

# Conclusion



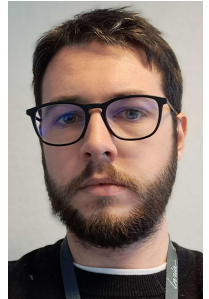
- **MERRIN<sup>1</sup>: inferring regulatory rules from metabolic traces**
  - **Hybrid (ASP + LP) resolution** based on SMT approaches
  - Compatible with **kinetics, fluxomics and/or transcriptomics data**
  - Compatible with **mutant strains**
- **Validation on simulated benchmark**
  - Find smaller RN than gold standard — Consistent with state of the art
  - Study the **impact of noise and data type** on the inferring — 240 instances
- **Scalability**
  - *E.coli* medium scale instance

---

<sup>1</sup> Implementation available on <https://github.com/bioasp/merrin/>

# Perspective — Work In Progress —

## *Model correction with MERRIN*



### Description:

- 1.473 regulatory rules
  - 600 genes and regulatory proteins
  - 873 regulated reactions
- 1075 reactions

| *Escherichia coli str. K-12 substr. MG1655*

### New input datatype: *Biolog data*<sup>1 2</sup>

- 111 mutant strains
- 124 mediums

| *Work in progress*

| **13.764 observations**  
| *(in silico and in vivo)*

**Existing models can be incompatible with new experimental results**  
— *How to update them?* —

<sup>1</sup> M. W. Covert and B. Ø. Palsson, **Nature**, 2004

<sup>2</sup> J. D. Glasner et al. **Nucleic Acids Res.**, 2003