



# Deriving Gene Associations for Metabolic Models

Razanne Issa, David James Sherman

## ► To cite this version:

Razanne Issa, David James Sherman. Deriving Gene Associations for Metabolic Models. [Research Report] RR-9246, Inria Bordeaux Sud-Ouest. 2018. hal-01976744

**HAL Id: hal-01976744**

**<https://hal.inria.fr/hal-01976744>**

Submitted on 10 Jan 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# Deriving Gene Associations for Metabolic Models

Razanne Issa, David James Sherman

**RESEARCH  
REPORT**

**N° 9246**

December 2018

Project-Team Pleiade

ISRN INRIA/RR--9246--FR+ENG

ISSN 0249-6399





## Deriving Gene Associations for Metabolic Models

Razanne Issa, David James Sherman

Project-Team Pleiade

Research Report n° 9246 — December 2018 — 25 pages

**Abstract:** We define a formal procedure for inferring gene-protein-reaction (GPR) relations from complete metabolic models, using a logical representation of knowledge and a small set of inference rules. We show that different use cases of metabolic models requires difference GPR relations. Three examples from the yeast *Saccharomyces cerevisiae* illustrate the procedure and demonstrate its usefulness.

**Key-words:** metabolic networks, comparative genomics, inference

**RESEARCH CENTRE  
BORDEAUX – SUD-OUEST**

200 avenue de la Vieille Tour  
33405 Talence Cedex

## Dériver des associations de gène pour des modèles métaboliques

**Résumé :** Nous définissons une procédure formelle pour l'inférence de relations gène-protéine-réaction (GPR) à partir de modèles métaboliques complets, par moyen d'une représentation logique des connaissances et un petit ensemble de règles d'inférence. Nous montrons que les différents cas d'utilisation requièrent des différentes relations GPR. Trois exemples tirés de la levure *Saccharomyces cerevisiae* illustrent la procédure et démontrent son utilité.

**Mots-clés :** réseaux métaboliques, génomique comparée, inférence

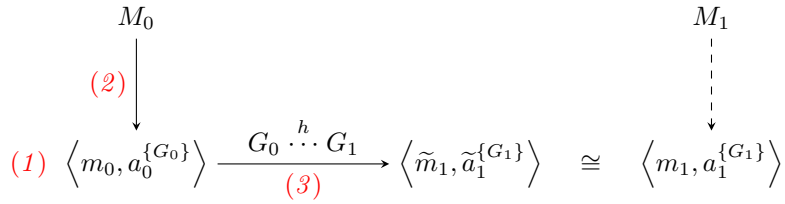


Figure 1: Incremental modeling: we (1) define a tuple  $\langle m, a^{\{G_0\}} \rangle$  that annotates a minimal enzymatic model  $m$  with annotations  $a_0$  defined over a set of variables  $G$ , and (2) project the tuple from a complete model  $M$ . The objective is that (3) the inferred model  $\langle \tilde{m}_1, \tilde{a}_1^{\{G_1\}} \rangle$  is congruent to the ideal model  $\langle m_1, a_1^{\{G_1\}} \rangle$  that we would have obtained had  $M_1$  been available.

## 1 Introduction

*Gene Associations* or Gene-Protein-Reaction relations (GPR) [RVSP03, NvEF<sup>+</sup>06] in metabolic models record knowledge about the direct and indirect relations between genes and biochemical reactions. Whether written as boolean formulas or as collections of  $\langle \text{gene, protein, reaction} \rangle$  triplets, gene associations define boolean functions that assert whether a reaction may take place. One widespread use of gene associations is in metabolic model inference for new target organisms. For such organisms there is typically a genome but very little systematic experimental data. The genes in the genome are used to evaluate the gene associations for every reaction in a reaction database; a reaction is included in the inferred model if its gene association is satisfied.

Gene associations play an essential role in incremental model inference (figure 1) as defined in [Iss15]. The incremental method is designed for cases when the target organism is closely related to an organism for which a minimal model  $\langle m, a^{\{G_0\}} \rangle$ , derived (or derivable) from a complete model  $M_0$ , already exists. The gene associations  $a^{\{G_0\}}$  annotate the metabolic model  $m$ , and are essential in inferring the target model  $\langle \tilde{m}_1, \tilde{a}_1^{\{G_1\}} \rangle$ .

In the general case, however, gene associations are an annotation determined by a human curator on the basis of existing knowledge about the biochemistry and gene regulation of the target organism. While expert creation is of course the ultimate criterion for correctness, it would be useful to also have a formal criterion defined in terms of the knowledge base.

In this report we describe how a generalized form of gene associations can be used to annotate a minimal enzymatic model, in order to infer specific metabolic models for related organisms based on genomic comparisons. The gene associations will implicitly represent knowledge derived from a hypothetical complete model. We will moreover consider how to rewrite the gene associations during model inference.

Different kinds of models, with different degrees of sophistication, can be defined for the same organism. For this work we have selected four complementary use cases for model inference using gene associations. We will show that each requires a different kind of assertion on the reaction, and consequently a different kind of gene association. Consequently, the step (1) of figure 1, which projects a complete model to a minimal model annotated by gene associations, depends on the type of models we wish to infer.

- *Model inference* uses an annotated reference model and homology between the genes in that reference and the genes in a target organism, to infer a model for the target organism [cite]. The assertion that must be tested for each reaction is whether the genome of the target contains the required genes, for example those coding for the enzymes that catalyze

the reaction.

- *Flux balance analysis* (FBA) calculates the fluxes through the reactions that maximize an objective function under a steady-state assumption [cite]. FBA can be used to simulate the effect of gene knockouts on metabolic fluxes. The assertion that must be tested is whether a reaction should be removed from the system because its requisite genes are knocked out or deactivated.
- *Hierarchical modeling* decomposes the complex cellular system into a hierarchy of independently-modeled stochastic transition systems, that communicate through flows and can synchronously perform state transitions through firing of global events [cite]. The assertion that must be tested is whether the *guard* on a transition, defined in terms of state variables such as protein expression or metabolite concentration, permits it to be executed.
- *Dynamic simulation* converts the metabolic model into a system of differential equations that define the evolution of molecule concentrations over time. Gene associations in dynamic simulation assert dependencies between variables that ultimately might be any molecular species involved in the model.

For each of these uses we define the appropriate gene associations and show how to derive them from a complete metabolic model. This derivation will provide a different gene association correctness criterion for each use case.

To illustrate the differences between the four use cases and the corresponding gene associations, we use three examples from *Saccharomyces cerevisiae*:

- Ribonuclease reductase, reaction *RNR*;
- Ferrocytochrome-C oxidoreductase, reaction *CYC*;
- 3-methyl-2-oxopentanoate decarboxylase, reaction *PDC*.

Each of these examples presents a different challenge: the RNR system must take the cell cycle into account, the CYC system must take compartments into account, and the PDR system must take gene regulation into account.

By showing how the appropriate gene associations can be derived from a complete model, we provide the means to define a gold standard for annotation of a metabolic model. However, in practical applications we rarely have a complete model for which gene associations could be derived; rather, gene association annotations are defined at the same time as the definition or the inference of the biochemical reactions in the model. The interesting question in these applications is how to improve the quality of the reactions and their annotations. We show how the use-specific definition of gene associations provides a consistency check on the model being defined.

## 2 Definitions and Notation

Gene associations or Gene-Protein-Reaction (GPR) relations are boolean formulas constructed using the operators “AND” and “OR” over two levels: genes to proteins and proteins to reactions. In the first level, the formulas express how proteins are coded by their genes. In the second level, the formulas express the dependencies of reactions on proteins (Figure 2). Identifying GPR relations is essential for the reconstruction of metabolic networks[FP08, FST05, FHT<sup>+</sup>09] and a critical step in that reconstruction.

The gene associations in a metabolic network can be identified through manual curation, a task that requires extensive experimental data [FGAT09] (e.g UniProtDB [C<sup>+</sup>11]), or by automatic methods. The latter are founded on knowledge engineering, and are principally based on sequence similarity [KDWV13, WLT05, AJW<sup>+</sup>08]. Because extensive experimental data are rarely available, efficient automatic methods are necessary to decipher correct annotations [RCO<sup>+</sup>13].

Gene associations are stored in many public data banks such as KEGG [KG00] and MetaCyc [CFF<sup>+</sup>06].

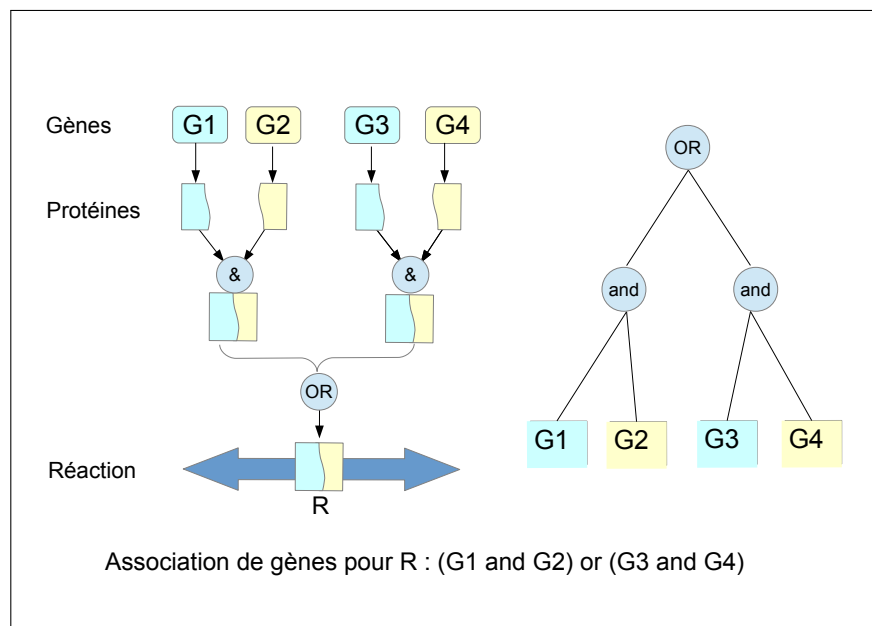


Figure 2: (left) The gene association for a reaction **R** catalyzed by a protein complex comprised either of the aggregation of products of genes **G1** and **G2**, or of the aggregation of the products of **G3** and **G4**. (right) The formula for the gene association, written as a binary tree.

The knowledge represented in the model by gene associations is indispensable for predicting cell phenotypes under genetic modifications such as gene knockouts in flux balance analysis (FBA) [JLP11], and for taking into account the effects of gene regulation in dynamic simulation [KSM12].

Let us commence by defining the logical representation we will use to document a state of knowledge about a set of reactions in a metabolic model and a set of gene-protein-reaction relations associated with those reactions. In this representation, genes and proteins are atomic formulas (or *atoms*), and the relations between these objects and the state of the system are *clauses*. For example, the clause "**A:B**  $\vdash$  **A**, **B**" means that, to the complex **A:B** can only be formed if the proteins **A** and **B** are already present. We define inference rules that combine and rewrite clauses in order to derive the formulas that are required for each use.

The following notation is used for the assertions employed by the logical representation:

- (**A**, **c**) - the molecule **A** is present in compartment **c**
- **A:B** - a macromolecular complex is formed from **A** and **B**
- **DEGRAD(A)** - the molecule **A** is degraded by the cell



- **DIS(A:B)** - the complex **A:B** is dissociated by the cell
- **OX(A)** - the molecule **A** is oxidized
- **RED(A)** - the molecule **A** is reduced
- **OX-RED(A)** - the redox cycle of the molecule **A** is present

This representation permits one to write clauses of the kind  $\mathbf{C}, \mathbf{D} \vdash \mathbf{A}, \mathbf{B}$ , where **A** and **B** imply **C** and **D**, for example as a result of a reaction  $\mathbf{A} + \mathbf{B} \rightarrow \mathbf{C} + \mathbf{D}$ . We will see, in what follows, that the formula required for each use can be obtained by rewriting clauses, up to a use-specific halting criterion.

BioRica [AGS<sup>+</sup>11] is a high-level modeling framework that integrates discrete, continuous, stochastic, nondeterministic, and timed behaviors. It provides a means to hierarchically compose models in a nonambiguous way. BioRica extends the ideas of AltaRica [ABCC94] [APGR99] to biological systems, allowing specifications that use of differential equations as well as automata with stochastic and nondeterministic behaviors. BioRica is designed for modeling complex systems through the combination of existing models of subsystems, in particular those imported from SBML models.

The basic building block of BioRica is the *node*, which has discrete or continuous behavior. The specification of a node has eight fields: *state*, *flow*, *sub*, *event*, *trans*, *assert*, *init* and *extern*. The field *state* declares state variables and their value domains. Hierarchical connections to other BioRica nodea are specified using the field *flow*, which defines input-output connections, and the field *sub*, which lists nodes that the current node uses hierarchically. The field *event* declares the name of the discrete events that are possible, and the field *trans* defines the state transitions that are caused by these events. The field *assert* imposes constraints on input-output flows and on states. The field *init* declares the initial values of state variables. Finally, the field *extern* is used to specify event delay distributions and priorities between events.

### 3 Strategies for Deriving Gene Associations

In what follows we propose strategies for deriving gene associations for each of the use cases presented in 1. All of the strategies are based on a common method: substitution and simplification rules are applied top-down until a halting criterion is reached. The choice of halting criterion to use depends on the nature of the use case being studied.

In the case of **model inference**, to determine whether a reaction can occur in an organism, it suffices to know whether the enzymes that catalyze it are present or not. Consequently, gene association rewriting must halt as soon as the genes that code for the enzymes catalyzing the reaction are observed.

Steady-state **flux balance analysis** supposes that a reaction can carry any flux that is physiologically possible, as long as the necessary genes are expressed. Assuming steady-state, gene association rewriting must halt before any substitution, such as gene expression, for example, that would correspond to a change of state.

In the case of **hierarchical modeling**, the nodes in the BioRica specification correspond to a model that defines a natural limit of knowledge. Consequently, gene association rewriting is limited to the confines of each node. These limited can be identified syntactically in the model's specification.

Finally, in **dynamic simulation**, since any reaction may contribute to the simulation, gene association rewriting must continue as long as clauses can be rewritten, and can in principle concern the entire model. In the following, we will not illustrate this strategy.

The three examples that follow have been chosen to illustrate the process of deriving gene associations in metabolic networks. Each example starts from one central reaction of interest, but may include up to fifteen associated reactions, required for understanding the prerequisite conditions of the central reaction.

- RNR, the ribonucleotide reductase system modifying **RNR** by **Trx1p**, specifically the reaction catalyzed by the bisulfite **Trx1p** ( $r_8$  dans Yeastnet).
- CYC, the complex **Cyc7p:Ccp1p** that catalyzes the transformation of ferricytochrome-C to ferrocycytochrome-C ( $r_{0469}$  dans Yeastnet).
- PDC, the transcription factor complex **Thi3p:Pdc2p** or **Thi3p:Thi2p** that induces the four paralogs of **Pcd1p**, who catalyze the transformation of 3-methyl-2-oxopentanoate en 2-methyl butanol ( $r_{0067}$  dans Yeastnet)

## 4 Ribonucleotide reductase system (RNR)

Ribonucleotide reductase (RNR) converts ribonucleotides to deoxyribonucleotides during the **S** phase of the cell cycle. Recall that the cell cycle corresponds to the set of steps constituting and delimiting the life of a cell, and is composed of several growth stages: the **Mitose/Méiose** phase, where the cell divides, the **Gap1** phase, where the cell grows until it reaches a critical size that signals transition to the **S** phase, where DNA is replicated, before transitioning to the **Gap2** phase.

During the **Gap1**, **Gap2** and **Mitose/Méiose** phases, when ribonucleotide reduction is not necessary, the heterodimer complex **Rnr2p:Rnr4p** is imported into the nucleus through the action first of the addressing protein **Dif1p**, which brings the complex **Rnr2p :Rnr4p** in proximity to the nucleus, and second of the protein **Kap122p** that imports the complex into the nucleus. The **Rnr2p:Rnr4p** complex is retained in the nucleus by the anchoring protein **Wtm1p**.

During the **S** phase of the cell cycle, the **Rnr2p:Rnr4p** complex is moved from the nucleus to the cytoplasm by two cooperating mechanisms:

- The degradation **Dif1p**, which prevents entry of new **Rnr2p:Rnr4p** complexes into the nucleus.
- The action of the export protein **Crm1p**, which transports the complex **Rnr2p:Rnr4p** from the nucleus to the cytoplasm.

The **Rnr2p:Rnr4p** complex associates with the homodimère complex **Rnr1p:Rnr1p** to form the new complex **Rnr1p:Rnr1p:Rnr2p:Rnr4p**, which is then reduced by one of the other of the isozymes of thioredoxine **Trx1p** or **Trx2p**. These proteins are coded respectively by the **TRX1** and **TRX2**. The reduced complex **Rnr1p:Rnr1p:Rnr2p:Rnr4p** in turn catalyzes the reduction of diphosphate ribonucleoside (**NDP**) into diphosphate deoxyribonucleoside (**dNDP**) and is thus oxidized.

A graphic representation of the model used here is presented in figure 3. Starting from this representation, we determined a logical representation and a BioRica specification corresponding to the reaction  $\text{NDP} \wedge \text{O}_2 \wedge 2\text{H}_2 \rightarrow \text{dNDP} \wedge 2\text{H}_2\text{O}$ .

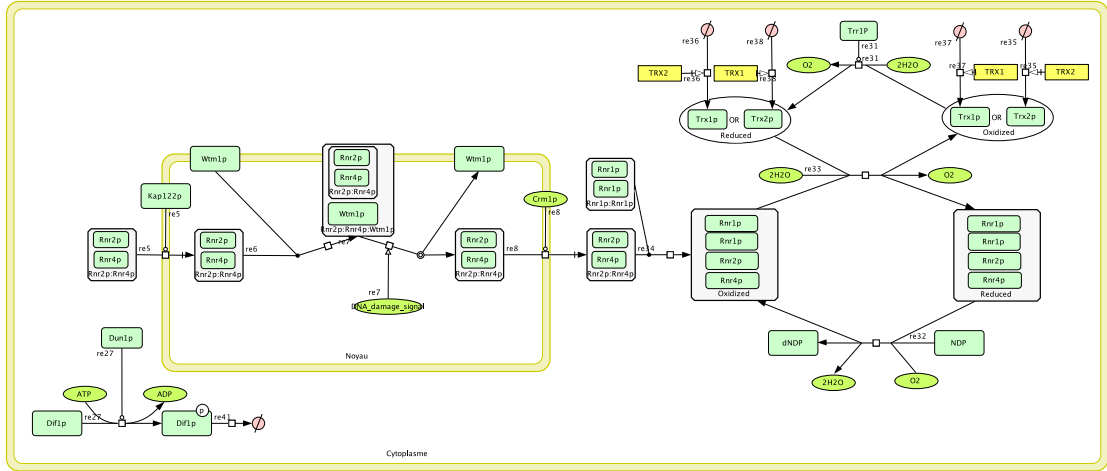


Figure 3: Explicit model for the ribonucleotide reductase reaction (RNR)

#### 4.1 Logical Description

$$\text{NDP} \wedge \text{O}_2 \wedge 2\text{H}_2 \rightarrow \text{dNDP} \wedge 2\text{H}_2\text{O} \vdash \text{OX-RED}(\text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p}) \quad (1)$$

$$\text{dNDP} \vdash \text{NDP} \wedge \text{OX-RED}(\text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p}) \quad (2)$$

$$\begin{aligned} \text{OX-RED}(\text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p}) \vdash & \text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p} \\ & \wedge (\text{OX-RED}(\text{Trx1p}) \vee \text{OX-RED}(\text{Trx2p})) \end{aligned} \quad (3)$$

$$\text{OX-RED}(\text{Trx1p}) \vee \text{OX-RED}(\text{Trx2p}) \vdash (\text{Trx1p} \vee \text{Trx2p}) \wedge \text{Trr1p} \quad (4)$$

$$\text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p} \vdash \text{Rnr1p}:\text{Rnr1p} \wedge \text{Rnr2p}:\text{Rnr4p} \quad (5)$$

$$\text{Rnr1p}:\text{Rnr1p} \vdash \text{Rnr1p} \wedge \text{Rnr1p} \quad (6)$$

$$\text{Rnr2p}:\text{Rnr4p} \vdash \text{Rnr2p} \wedge \text{Rnr4p} \quad (7)$$

$$(\text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p}, c) \vdash (\text{Rnr1p}:\text{Rnr1p}, c) \wedge (\text{Rnr2p}:\text{Rnr4p}, c) \quad (8)$$

$$(\text{Rnr1p}:\text{Rnr1p}, c) \vdash (\text{Rnr1p}:\text{Rnr1p}, n) \wedge \text{Crm1p} \wedge \text{Phase}(S) \quad (9)$$

$$\begin{aligned} (\text{Rnr1p}:\text{Rnr1p}, n) \vdash & (\text{Rnr1p}:\text{Rnr1p}, c) \wedge \text{Kap122p} \wedge \text{Dif1p} \\ & \wedge \neg\text{Phase}(S) \end{aligned} \quad (10)$$

$$\text{Dif1p} \vdash \neg\text{Dun1p} \wedge \neg\text{Phase}(S) \quad (11)$$

$$\text{Rnr2p}:\text{Rnr4p}:\text{Wrm1p} \vdash \text{Rnr2p} \wedge \text{Rnr4p} \wedge \text{Wrm1p} \quad (12)$$

$$\text{Crm1p} \vdash \text{CRM1} \quad (13)$$

$$(14)$$

$$\text{Dif1p} \vdash \text{DIF1} \quad (15)$$

$$\text{Dun1p} \vdash \text{DUN1} \quad (16)$$

$$\text{Kap122p} \vdash \text{KAP122} \quad (17)$$

$$\text{Rnr1p} \vdash \text{RNR1} \quad (18)$$

$$\text{Rnr2p} \vdash \text{RNR2} \quad (19)$$

$$\text{Rnr4p} \vdash \text{RNR4} \quad (20)$$

$$\text{Trr1p} \vdash \text{TRR1} \quad (21)$$

$$\text{Trx1p} \vdash \text{TRX1} \quad (22)$$

$$\text{Trx2p} \vdash \text{TRX2} \quad (23)$$

$$\text{Wrm1p} \vdash \text{WRM1} \quad (24)$$

## 4.2 BioRica Description

The hierarchical model of the RNR system is comprised of three BioRica nodes:

- **Subcomplex**, which describes the formation of the **Rnr2p:Rnr4p** heterodimer and its importation into the nucleus. We suppose that the formation of the complex is the limiting step in the reaction; the kinetics of the addressing and importing reactions are assumed to be much more rapid, and are not explicitly modeled.
- **Complex**, which describes the formation of the **Rnr1p:Rnr1p:Rnr2p:Rnr4p** heterotetramer in the cytoplasm during **S** phase. We suppose that this formation is only limited by the presence of **Rnr2p:Rnr4p** in the cytoplasm, which is in turn controlled by the transition to **S** phase. The transporter **Crml** is not explicitly modeled.
- **Reaction**, which described the system of reactions neighboring ribonucleotide reductase. This includes the cycles of reduction-oxidation of the reductase heterotetramer and both of the thioredoxines **Trx1** and **Trx2**.

The **Main** node links these three sub-nodes by boolean variables that indicate the presence in the cytoplasm of the heterodimer **Rnr2p:Rnr4p** and the heterotetramer **Rnr1p:Rnr1p:Rnr2p:Rnr4p**. It also synchronizes the transitions to the **S** phase by the sub-nodes **Subcomplex** and **Complex**. The complete BioRica specification is as follows.

### node Subcomplex

#### state

$$\begin{aligned} T_c, T_n & : \{0, 1\}; \\ \varepsilon & : \{0, 1\}; \end{aligned}$$

#### flow

$$(\text{Rnr2p:Rnr4p}, c) : \text{FLOAT};$$

#### init

$$(\text{Rnr2p:Rnr4p}, c_i) ;$$

#### diff

$$\begin{aligned} d(\text{Rnr2p:Rnr4p}, n) = & T_n \cdot (\text{cell.V1} \cdot (1 + 1 \cdot \ln(\frac{(\text{Rnr2p:Rnr4p}, c)}{ic(\text{Rnr2p:Rnr4p}, c)})) - \ln(\frac{(\text{Rnr2p:Rnr4p}, n)}{ic((\text{Rnr2p:Rnr4p}, n)}))) \\ & + \varepsilon(\text{Rnr2p:Rnr4p}, c_i); \end{aligned}$$

#### event

$$\begin{aligned} & \text{phaseS}; \\ & \text{sortieS}; \end{aligned}$$

```

trans
  true ⊢ phaseS → Tn := 0, Tc := 1, ε := 1;
  true ⊢ sortieS → Tn := 1, Tc := 0, ε := 0;
edon

node Complex

state
  (Rnr1p:Rnr1p, c) = BOOL;
flow
  (Rnr2p:Rnr4p, c) = BOOL;
  (Rnr1p:Rnr1p:Rnr2p:Rnr4p, c) = BOOL;
event
  phaseS;
  associate;
trans
  (Rnr2p:Rnr4p, c), (Rnr1p:Rnr1p, c) ⊢ true ⊢ phaseS → (Rnr2p:Rnr4p, c) := true;
  (Rnr1p:Rnr1p:Rnr2p:Rnr4p, c) ⊢ associate → (Rnr1p:Rnr1p:Rnr2p:Rnr4p, c) := true;
edon

node Reaction

state
  TRX1 : BOOL;
  TRX2 : BOOL;
  Trx1p : BOOL;
  Trx2p : BOOL;
  Trr1p : BOOL;
  ox(Trx1p or Trx2p) = BOOL ;
  D : {0, 1} ;
  NDP : FLOAT;
  dNDP : FLOAT;
  dNDP : FLOAT;
  O2 : FLOAT;
  H2O : FLOAT;
  ox(Rnr1p:Rnr1p:Rnr2p:Rnr4p, c) : BOOL;
flow
  (Rnr1p:Rnr1p:Rnr2p:Rnr4p, c) : BOOL;
  ox(Rnr1p:Rnr1p:Rnr2p:Rnr4p, c) : BOOL;
event
  ox(TRX);
  ox(RNR);
  G-expr(TRX1);
  G-expr(TRX2);
  active;
constants
  iC(Ndp);
  iC(dNdp);
  iC(O2) ;
  iC(H2O);
diff
  d(dNDP) = D.(cell .V0.(1 + 1. ln( $\frac{C(NDP)}{ic(NDP)}$ ) + ln( $\frac{C(O_2)}{ic(O_2)}$ ) - 2 ln( $\frac{C(H_2O)}{ic(H_2O)}$ )
    - ln( $\frac{C(dNDP)}{ic(dNDP)}$ )));
trans

```

$$\begin{array}{l}
 \text{TRX1} \vdash \text{G-expr}(\text{TRX1}) \rightarrow \text{Trx1p} := \text{true}; \\
 \text{TRX2} \vdash \text{G-expr}(\text{TRX2}) \rightarrow \text{Trx2p} := \text{true}; \\
 \text{Trr1p} \vdash \text{ox}(\text{TRX}) \rightarrow \text{ox}(\text{Trx1p} \mid \text{Trx2p}) := \text{true}; \\
 (\text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p}, c), \text{ox}(\text{Trx1p} \mid \text{Trx2p}) \vdash \text{ox}(\text{RNR}) \rightarrow \\
 \text{ox}(\text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p}, c) := \text{true}; \\
 \text{ox}(\text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p}, c) \vdash \text{active} \rightarrow \text{D} = 1 ; \\
 \text{edon} \\
 \text{node Main} \\
 \text{sub} \\
 \text{N} : \text{Subcomplex}; \\
 \text{C} : \text{Complex}; \\
 \text{R} : \text{Reaction}; \\
 \text{assert} \\
 \text{C}.\langle \text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p}, c \rangle = \text{R}.\langle \text{Rnr1p}:\text{Rnr1p}:\text{Rnr2p}:\text{Rnr4p}, c \rangle; \\
 \text{N}.\langle \text{Rnr2p}:\text{Rnr4p}, c \rangle = \text{C}.\langle \text{Rnr2p}:\text{Rnr4p}, c \rangle; \\
 \text{sync} \\
 \langle \text{N.phaseS}, \text{C.phaseS} \rangle; \\
 \text{edon}
 \end{array}$$

### 4.3 Gene Associations for Model Inference

We can conclude that the ribonucleotide reductase reaction can occur in a target organism if there is sufficient evidence to rewrite its gene association. The reaction is catalyzed by the oxidation cycle of the **Rnr1p:Rnr1p:Rnr2p:Rnr4p** complex, so in order to instantiate the reaction we must be assured that the genes coding for the constituent proteins of the complex are present in the target organism (figure 4).

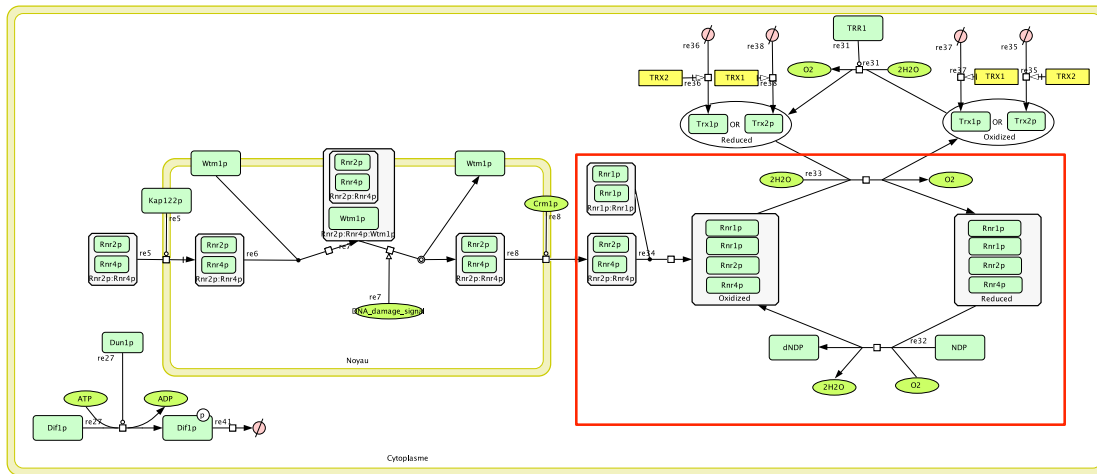


Figure 4: Explicit model of the ribonucleotide reductase reaction, outlining in red the part of the model considered when performing model inference

The gene association of the ribonucleotide reductase reaction is given by  $\text{RNR1} \wedge \text{RNR2} \wedge \text{RNR4}$ . This is derived formally from the presence of the oxidation cycle of the **Rnr1p:Rnr1p:Rnr2p:Rnr4p** complex:



thioredoxine reductase. Thus, to describe these dependencies, we require the genes coding the **Rnr1p:Rnr1p:Rnr2p:Rnr4p** heterotetramer reductase complex, the thioredoxine that performs the reduction, as well as the thioredoxine reductase (Figure 6).

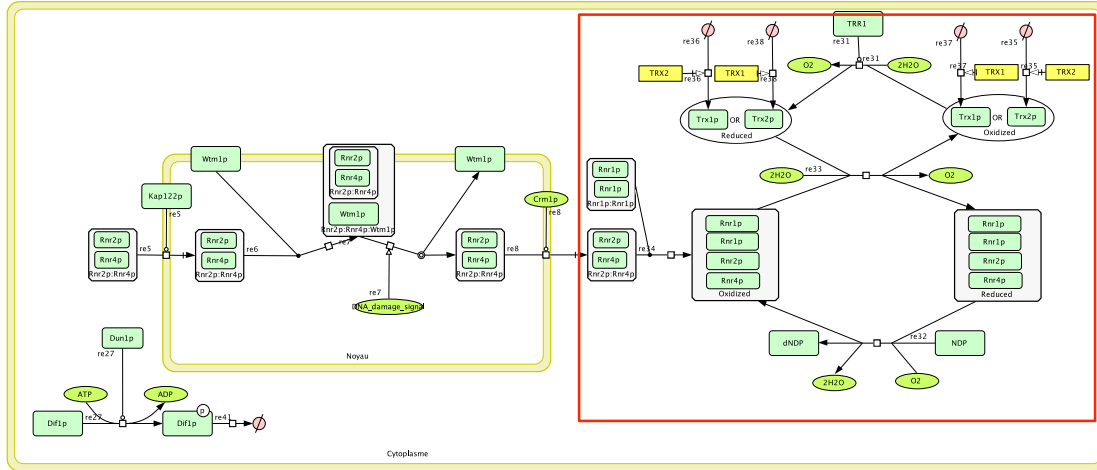


Figure 6: Explicit model of the ribonucleotide reductase reaction, outlining in red the part of the model considered when performing hierarchical modeling

$$\begin{aligned}
 & \text{Rnr1p:Rnr1p:Rnr2p:Rnr4p} \wedge (\text{Trx1p} \vee \text{Trx2p}) \wedge \text{Trr1p} \\
 \vdash & \text{Rnr1p:Rnr1p} \wedge \text{Rnr2p:Rnr4p} \wedge (\text{Trx1p} \vee \text{Trx2p}) \wedge \text{Trr1p} && \text{eqs(4)(5)} \\
 \vdash & \text{Rnr1p} \wedge \text{Rnr2p} \wedge \text{Rnr4p} \wedge (\text{Trx1p} \vee \text{Trx2p}) \wedge \text{Trr1p} && \text{eqs(6)(7)} \\
 \vdash & \text{RNR1} \wedge \text{RNR2} \wedge \text{RNR4} \wedge (\text{TRX1} \vee \text{TRX2}) \wedge \text{TRR1} && \text{eqs(18) (19)} \\
 & && \text{(20) (22)} \\
 & && \text{(23) (21)}
 \end{aligned}$$



## 5 Ferrocyclochrome c System (CYC)

The **Cyc1p** and **Cyc7p** isoforms of cytochrome c are electron transport proteins taking part in the electron transport chain in mitochondria. In this system, the **Cyc1p** or **Cyc7p** proteins form a complex with cytochrome c peroxidase **Ccp1p**. In this complex, **Ccp1p** catalyzes the oxidation of 2-ferrocyclochrome c, which is the reduced form of **Cyc1p** or **Cyc7p** (figure 7).

Les cytochrome c isoformes **Cyc1p** et **Cyc7p** sont des protéines de transport d'électrons dans la chaîne de transport d'électrons mitochondriale. Dans ce système, la protéine **Cyc1p** ou **Cyc7p** forme d'abord un complexe avec cytochrome c peroxidase **Ccp1p**, et dans ce complexe, **Ccp1p** catalyse l'oxydation du 2-ferrocyclochrome c, qui est la forme réduite de **Cyc1p** or **Cyc7p** (Figure 7).

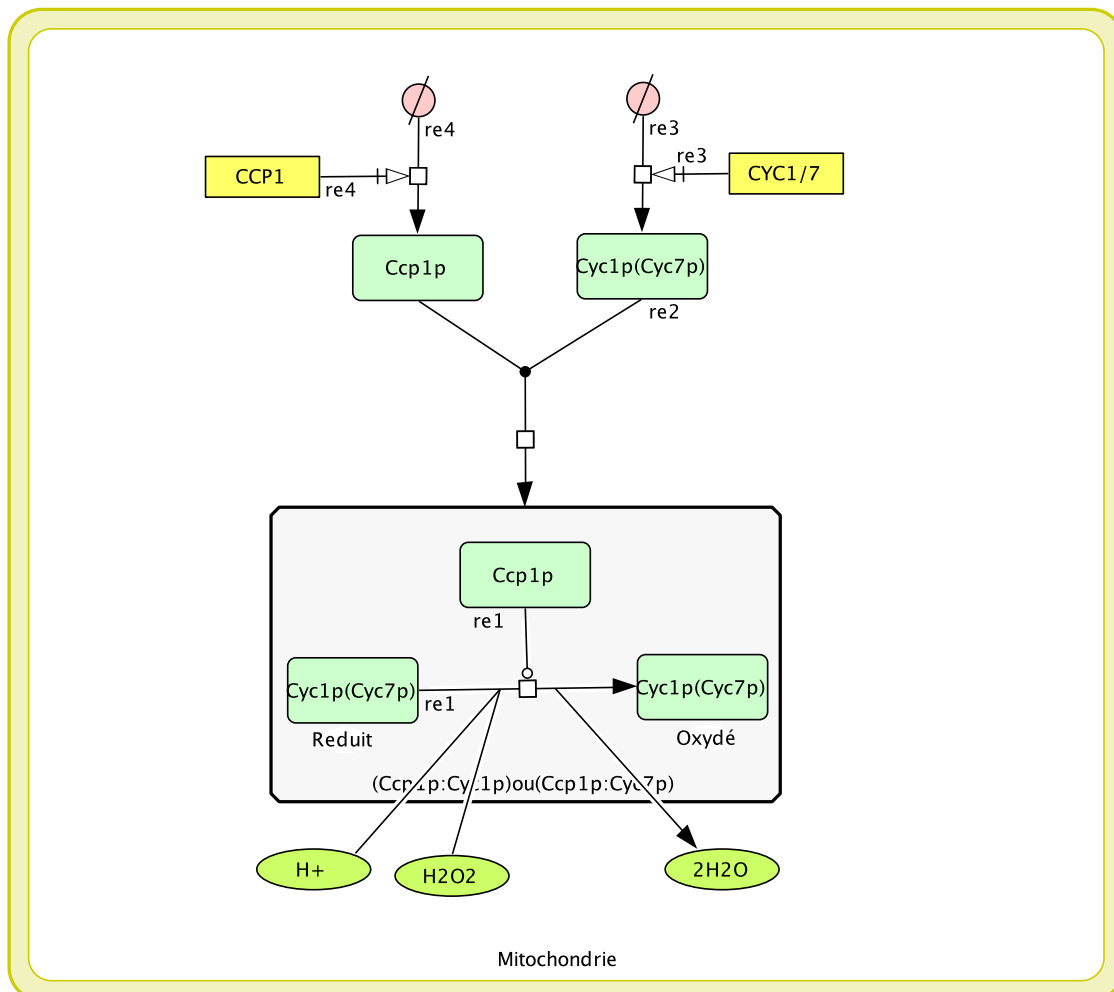


Figure 7: Explicit model of the ferrocyclochrome c reaction

## 5.1 Logical Description

$$R = \text{RED}(\text{Cyc1p}) \wedge \text{H}^+ \wedge \text{H}_2\text{O}_2 \rightarrow \text{ox}(\text{Cyc1p}) \wedge 2\text{H}_2\text{O}$$

$$R \vdash \text{Ccp1p} \wedge \text{Cyc1p:Ccp1p} \quad (25)$$

$$\text{Cyc1p:Ccp1p} \vdash \text{Cyc1p} \wedge \text{Ccp1p} \quad (26)$$

$$\text{Cyc1p} \vdash \text{CYC1} \quad (27)$$

$$\text{Ccp1p} \vdash \text{CCP1} \quad (28)$$

$$R' = \text{RED}(\text{Cyc7p}) \wedge \text{H}^+ \wedge \text{H}_2\text{O}_2 \rightarrow \text{ox}(\text{Cyc7p}) \wedge 2\text{H}_2\text{O}$$

$$R' \vdash \text{Ccp7p} \wedge \text{Cyc7p:Ccp1p} \quad (29)$$

$$\text{Cyc7p:Ccp1p} \vdash \text{Cyc7p} \wedge \text{Ccp1p} \quad (30)$$

$$\text{Cyc7p} \vdash \text{CYC7} \quad (31)$$

$$\text{Ccp1p} \vdash \text{CCP1} \quad (32)$$

## 5.2 BioRica Description

The hierarchical model of the **CYC** system contains a single hybrid node **Main**, composed of discrete transitions that represent gene expression and the formation of the **Ccp1p:cytochrome c** complex, and one continuous equation describing the oxidation of cytochrome c by **Ccp1**. The oxidation of **Ccp1** may only occur once the **Ccp1p:cytochrome c** complex is formed.

node Main

state

CCP1	:	BOOL;
CYC1	:	BOOL;
CYC7	:	BOOL;
Ccp1p	:	BOOL;
Cyc1p	:	BOOL;
Cyc7p	:	BOOL;
Ccp1p:(Cyc1p Cyc7p)	:	BOOL;
$R_0$	=	{0, 1};
(Cyc1p Cyc7p) <sub>O</sub>	:	FLOAT;
(Cyc1p Cyc7p) <sub>R</sub>	:	FLOAT;

diff

$$d((\text{Cyc1p}|\text{Cyc7p})_O) = R_0(\text{cell.V}_0.(1 + 1.\ln(\frac{(\text{Cyc1p}|\text{Cyc7p})_R}{ic(\text{Cyc1p}|\text{Cyc7p})_R}) + \ln(\frac{H^+}{ic(H^+)}) + \ln(\frac{H_2O_2}{ic(H_2O_2)}) - \ln(\frac{(\text{Cyc1p}|\text{Cyc7p})_O}{ic(\text{Cyc1p}|\text{Cyc7p})_O}) - \ln(\frac{H_2O}{ic(H_2O)}))));$$

event

G-expr(CCP1); G-expr(CYC1); G-expr(CYC7);  
associateCOM; Oxydate;

trans

CCP1	⊢	G-expr(CCP1)	→	Ccp1p := true;
CYC1	⊢	G-expr(CYC1)	→	Cyc1p := true;
CYC7	⊢	G-expr(CYC7)	→	Cyc7p := true;
Ccp1p, (Cyc1p   Cyc7p)	⊢	associateCOM	→	Ccp1p:(Cyc1p   Cyc7p) := true;
Ccp1:(Cyc1p   Cyc7p)	⊢	Oxydate	→	$R_0 := 1;$

edon

### 5.3 Gene Associations for Model Inference

The genes needed so that this reaction may occur in the target organism are the genes coding for the catalyst **Ccp1p** and the complex **Cyc1p|Cyc7p:Ccp1p** (figure 8).

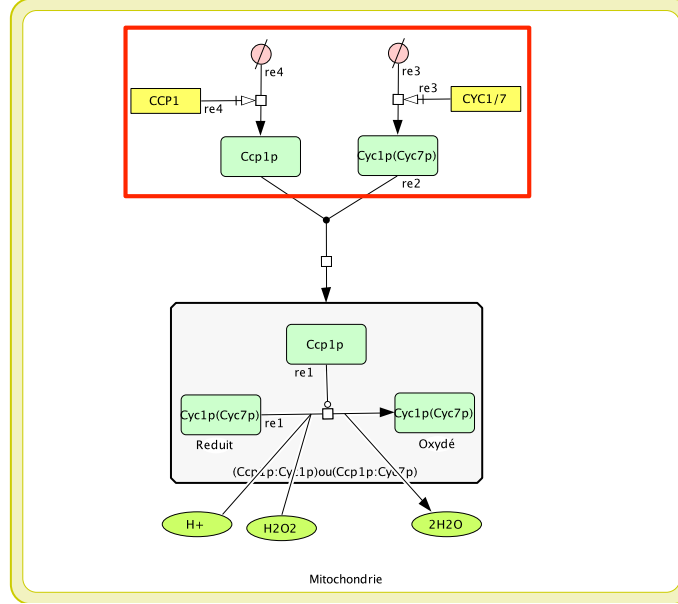
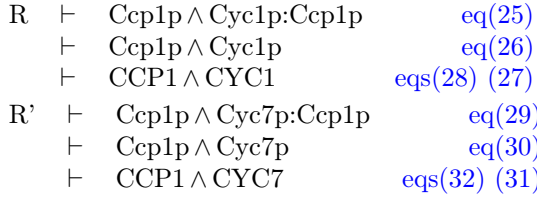
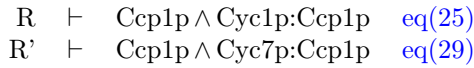


Figure 8: Explicit model for the Ferrocycytochrome c system, outlining in red the part of the model considered when performing model inference



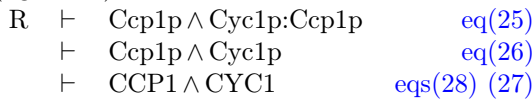
### 5.4 Gene Associations for Flux Balance Analysis

In this example no state changes occur, so the reaction may be instantiated if the enzyme and the complex are already present. (figure 9).



### 5.5 Gene Associations for Hierarchical Modeling

The oxidation of cytochrome c (**Cyc1p** ou **Cyc7p**) by **Ccp1p** takes place within the **Cyc1p:Ccp1p** or the **Cyc7p:Ccp1p** complex, which associates cytochrome c with cytochrome c peroxidase (figure 10).



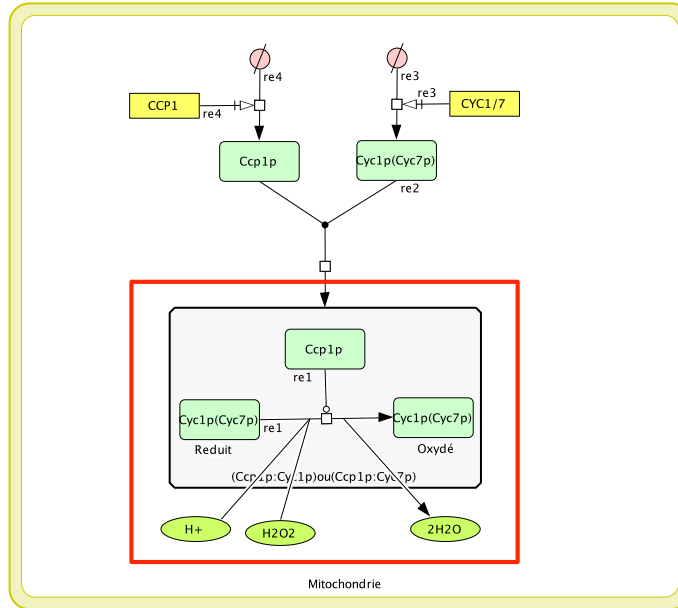


Figure 9: Explicit model of the Ferrocycytochrome c system, outlining in red the part of the model considered when performing flux balance analysis

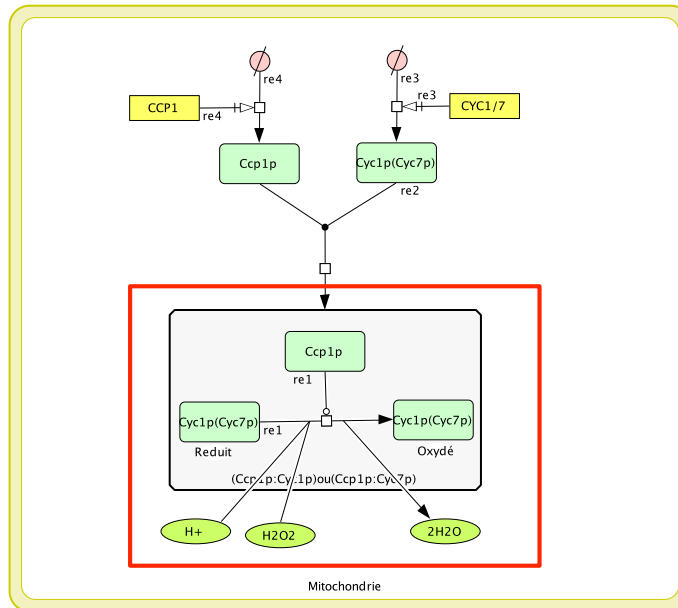


Figure 10: Explicit model of the Ferrocycytochrome c system, outlining in red the part of the model considered when performing hierarchical modeling

$$\begin{aligned}
 R' &\vdash \text{Ccp1p} \wedge \text{Cyc7p:Ccp1p} && \text{eq(29)} \\
 &\vdash \text{Ccp1p} \wedge \text{Cyc1p} && \text{eq(30)} \\
 &\vdash \text{CCP1} \wedge \text{CYC7} && \text{eqs(32) (31)}
 \end{aligned}$$

## 6 3-methyl-2-oxopentanoate Decarboxylase System (PDC)

The enzymes of 3-methyl-2-oxopentanoate decarboxylase, **Pdc1p**, **Pdc5p**, **Pdc6p**, et **ARO10**, react to 2-oxo-acids like 3-methyl-2-oxopentanoate, which is an intermediate in the synthesis of isoleucine, (figure 11).

**Thi3p** is a regulatory protein of this reaction. It binds to the transcription factors **Thi2p** and **Pdc2p**. In the presence of thiamine, the **Pdc2p:Thi3p** and **Thi2p:Thi3p** complexes dissociate, and thus deactivate the transcription factors. The genes **PDC1** and **PDC5**, which code respectively for the proteins **Pdc1p** and **Pdc5p**, are regulated by the **Pdc2p:Thi3p** complex or by the **Thi2p:Thi3p** complex, and by the absence of the protein **Pdc1p**. However, the genes **PDC6** and **ARO10**, which code respectively for the proteins **Pdc6p** and **Aro10p**, are not regulated by these mechanisms.

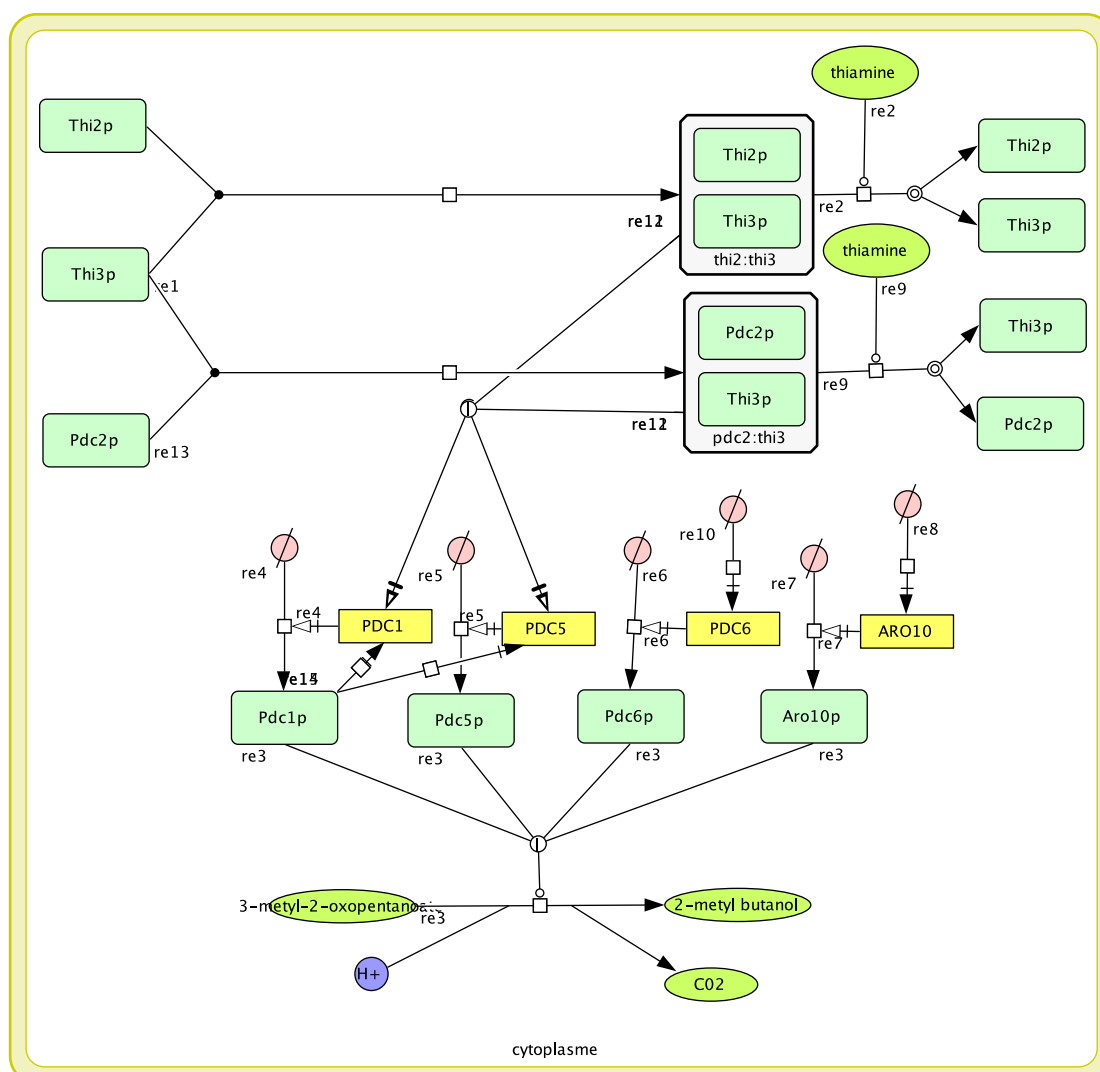


Figure 11: Explicit model of the 3-méthyl-2-oxopentanoate décarboxylase reaction

## 6.1 Logical Description

$$R(\text{PDC}) = 3\text{-methyl-2-oxepentanal} \wedge \text{H}^+ \rightarrow 2\text{-methylbutanal} \wedge \text{CO}_2$$

$$R(\text{PDC}) \vdash \text{Pdc1p} \vee \text{Pdc5p} \vee \text{Pdc6p} \vee \text{Aro10} \quad (33)$$

$$\text{Pdc1p} \vdash \text{PDC1} \wedge \text{Thi3p} : (\text{Thi2p} \vee \text{Pdc2p}) \wedge \neg \text{Pdc1p} \quad (34)$$

$$\text{Pdc5p} \vdash \text{PDC5} \wedge \text{Thi3p} : (\text{Thi2p} \vee \text{Pdc2p}) \wedge \neg \text{Pdc1p} \quad (35)$$

$$\text{Thi3p:Pdc2p} \vdash \text{Thi3p} \wedge \text{Pdc2p} \wedge \neg \text{Thiamine} \quad (36)$$

$$\text{Thi3p:Thi2p} \vdash \text{Thi3p} \wedge \text{Thi2p} \wedge \neg \text{Thiamine} \quad (37)$$

$$\text{DIS}(\text{Thi3p:Pdc2p}) \vdash \text{Thi3p:Pdc2p} \wedge \text{Thiamine} \quad (38)$$

$$\text{DIS}(\text{Thi3p:Thi2p}) \vdash \text{Thi3p:Thi2p} \wedge \text{Thiamine} \quad (39)$$

$$\text{Pdc6p} \vdash \text{PDC6} \quad (40)$$

$$\text{Aro10p} \vdash \text{ARO10} \quad (41)$$

## 6.2 BioRica Description

The hierarchical model of the **PDC** system is composed of two sub=nodes:

- **Complex**, which describes the association of **Thi3p** with either **Thi2p** or with **Pdc2p**, to form the transcription factor that regulates the expression of the genes **PDC1** and **PDC5**.
- **Factor**, which is a hybrid node composed of discrete transitions that represent the expression of the genes **PDC1**, **PDC5**, **PDC6**, or **ARO10**, and a continuous equation representing the decarboxylase reaction. The expression of genes **PDC1** and **PDC5** depends on the transcription factor.

The **Main** node is connected to **Complex** and **Factor** by boolean variables that indicate the presence of one or the other transcription factor complexes.

### node Complex

#### state

Thi3p : BOOL;  
 Thi2p : BOOL;  
 Pdc2p : BOOL;  
 Thiamin : BOOL;

#### flow

Thi2p:Thi3p : BOOL;  
 Pdc2p:Thi3p : BOOL;

#### event

associate(Thi2p:Thi3p); associate(Pdc2p:Thi3p);

#### trans

Thi3p , Thi2p , ¬Thiamin  $\vdash$  associate(Thi2p:Thi3p)  $\rightarrow$  Thi2p:Thi3p := true;  
 Thi3p , Pdc2p , ¬Thiamin  $\vdash$  associate(Pdc2p:Thi3p)  $\rightarrow$  Pdc2p:Thi3p := true; **edon**

### node Factor

#### state

```

PDC1 : BOOL;
PDC5 : BOOL;
PDC6 : BOOL;
ARO10 : BOOL;
Pcd1p : BOOL;
Pcd5p : BOOL;
Pcd6p : BOOL;
Aro10p : BOOL;
3-metyl-2-oxepentanal : FLOAT;
R0 : {0,1};

flow
  Thi2p:Thi3p : BOOL;
  Pdc2p:Thi3p : BOOL;

const
  ic(3-metyl-2-exepontanal)

diff
  d( 3-metyl-2-oxepentanal ) = R0(cell.V0.(1 + 1. ln(  $\frac{3\text{-metyl-2-oxepentanal}}{ic(3\text{-metyl-2-exepontanal})}$  )
    + ln(  $\frac{H^+}{ic(H^+)}$  ) - ln(  $\frac{2\text{-metylbutanal}}{ic(2\text{-metylbutanal})}$  ) - ln(  $\frac{CO2}{ic(CO2)}$  )))

event
  G-expr(PDC1); G-expr(PDC5);
  G-expr(PDC6); G-expr(ARO10);
  React;

trans
  PDC1, (Thi2p:Thi3p | Pdc2p:Thi3p),  $\neg$  Pdc1p  $\vdash$  G-expr(PDC1)  $\rightarrow$  Pdc1p := true ;
  PDC5, (Thi2p:Thi3p | Pdc2p:Thi3p),  $\neg$  Pdc5p  $\vdash$  G-expr(PDC5)  $\rightarrow$  Pdc5p := true ;
  PDC6  $\vdash$  G-expr(PDC6)  $\rightarrow$  Pdc6p := true ;
  ARO10  $\vdash$  G-expr(ARO10)  $\rightarrow$  Aro10p := true ;
  Pdc1p | Pdc5p | Pdc6p | Aro10p  $\vdash$  React  $\rightarrow$  R0 :=1 ;

edon

node Main

sub
  F : Factor;
  C: complex ;

assert
  C.Thi2p:Thi3p = F.Thi2p:Thi3p;
  C.Pdc2p:Thi3p = F.Pdc2p:Thi3p;

edon

```

### 6.3 Gene Associations for Model Inference

The reaction is catalyzed by one of the enzymes **Pdc1p**, **Pdc5p**, **Pdc6p**, and **Aro10p**, which are products of the genes **PDC1**, **PDC5**, **PDC6**, and **ARO10**. Consequently, the reaction may be inferred if at least one of these genes is present (figure 12).

$$\begin{aligned}
 R(\text{PDC}) \vdash & \text{Pdc1p} \vee \text{Pdc5p} \vee \text{Pdc6p} \vee \text{Aro10p} && \text{eq(33)} \\
 \vdash & \text{PDC1} \vee \text{PDC5} \vee \text{PDC6} \vee \text{ARO10} && \text{eqs(34) (35) (40) (41)}
 \end{aligned}$$

### 6.4 Gene Associations for Flux Balance Analysis

In this specific case, we suppose that the enzymes and enzymatic complexes are present (or not), and we only need to model the cycles of reduction reactions.





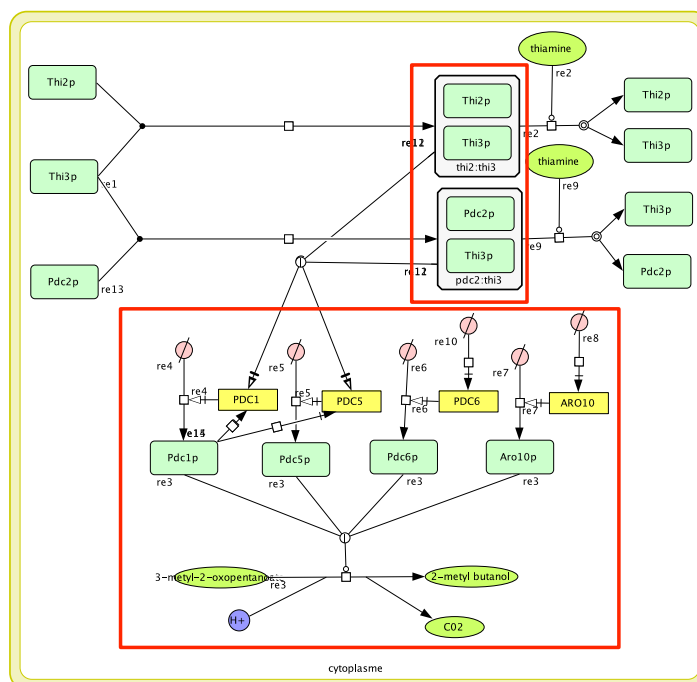


Figure 13: Explicit model of the 3-methyl-2-oxopentanoate decarboxylase system, outlining in red the part of the model considered when performing flux balance analysis

## References

- [ABCC94] André Arnold, Didier Bégay, Paul Crubillé, and P Crubille. *Construction and analysis of transition systems with MEC*. Number 3 in AMAST series in computing. World Scientific, 1994.
- [AGS<sup>+</sup>11] Rodrigo Assar, Alice Garcia, David James Sherman, et al. Modeling stochastic switched systems with biorica. In *Journées Ouvertes en Biologie, Informatique et Mathématiques JOBIM 2011*, pages 297–304, 2011.
- [AJW<sup>+</sup>08] Michiel E Adriaens, Magali Jaillard, Andra Waagmeester, Susan LM Coort, Alex R Pico, and Chris TA Evelo. The public road to high-quality curated biological pathways. *Drug discovery today*, 13(19):856–862, 2008.
- [APGR99] André Arnold, Gérard Point, Alain Griffault, and Antoine Rauzy. The altarica formalism for describing concurrent systems. *Fundam. Inf.*, 40(2,3):109–124, August 1999.
- [C<sup>+</sup>11] UniProt Consortium et al. Ongoing and future developments at the universal protein resource. *Nucleic acids research*, 39(suppl 1):D214–D219, 2011.
- [CFF<sup>+</sup>06] Ron Caspi, Hartmut Foerster, Carol A Fulcher, Rebecca Hopkinson, John Ingraham, Pallavi Kaipa, Markus Krummenacker, Suzanne Paley, John Pick, Seung Y Rhee, et al. Metacyc: a multiorganism database of metabolic pathways and enzymes. *Nucleic acids research*, 34(suppl 1):D511–D516, 2006.



- [KG00] Minoru Kanehisa and Susumu Goto. Kegg: kyoto encyclopedia of genes and genomes. *Nucleic acids research*, 28(1):27–30, 2000.
- [KSM12] Akhil Kumar, Patrick F Suthers, and Costas D Maranas. Metrxn: a knowledge-base of metabolites and reactions spanning metabolic models and databases. *BMC bioinformatics*, 13(1):6, 2012.
- [NvEF<sup>+</sup>06] R. A. Notabaart, F. H. van Enckevort, C. Francke, R. J. Siezen, and B. Teusink. Accelerating the reconstruction of genome-scale metabolic networks. *BMC Bioinformatics*, 7:296, Jun 2006.
- [RCO<sup>+</sup>13] Predrag Radivojac, Wyatt T Clark, Tal Ronnen Oron, Alexandra M Schnoes, Tobias Wittkop, Artem Sokolov, Kiley Graim, Christopher Funk, Karin Verspoor, Asa Ben-Hur, et al. A large-scale evaluation of computational protein function prediction. *Nature methods*, 10(3):221–227, 2013.
- [RVSP03] J. L. Reed, T. D. Vo, C. H. Schilling, and B. O. Palsson. An expanded genome-scale model of Escherichia coli K-12 (iJR904 GSM/GPR). *Genome Biol.*, 4(9):R54, 2003.
- [WLT05] James D Watson, Roman A Laskowski, and Janet M Thornton. Predicting protein function from sequence and structural data. *Current opinion in structural biology*, 15(3):275–284, 2005.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Definitions and Notation</b>	<b>4</b>
<b>3</b>	<b>Strategies for Deriving Gene Associations</b>	<b>6</b>
<b>4</b>	<b>Ribonucleotide reductase system (RNR)</b>	<b>7</b>
4.1	Logical Description . . . . .	8
4.2	BioRica Description . . . . .	9
4.3	Gene Associations for Model Inference . . . . .	11
4.4	Gene associations for Flux Balance Analysis . . . . .	12
4.5	Gene Associations for Hierarchical Modeling . . . . .	12
<b>5</b>	<b>Ferrocycytochrome c System (CYC)</b>	<b>14</b>
5.1	Logical Description . . . . .	15
5.2	BioRica Description . . . . .	15
5.3	Gene Associations for Model Inference . . . . .	16
5.4	Gene Associations for Flux Balance Analysis . . . . .	16
5.5	Gene Associations for Hierarchical Modeling . . . . .	16
<b>6</b>	<b>3-methyl-2-oxopentanoate Decarboxylase System (PDC)</b>	<b>18</b>
6.1	Logical Description . . . . .	19
6.2	BioRica Description . . . . .	19
6.3	Gene Associations for Model Inference . . . . .	20
6.4	Gene Associations for Flux Balance Analysis . . . . .	20
6.5	Gene Associations for Hierarchical Modeling . . . . .	21

**7 Discussion**

**21**



**RESEARCH CENTRE  
BORDEAUX – SUD-OUEST**

200 avenue de la Vieille Tour  
33405 Talence Cedex

Publisher  
Inria  
Domaine de Voluceau - Rocquencourt  
BP 105 - 78153 Le Chesnay Cedex  
[inria.fr](http://inria.fr)

ISSN 0249-6399