# RBA

S. Fischer - Biosys - MAIAGE

November 28, 2016

# Contents

# 1 Principles underlying algorithm

Figure 1 shows how the fluxes of molecules are described by RBA and where the constraint are placed.
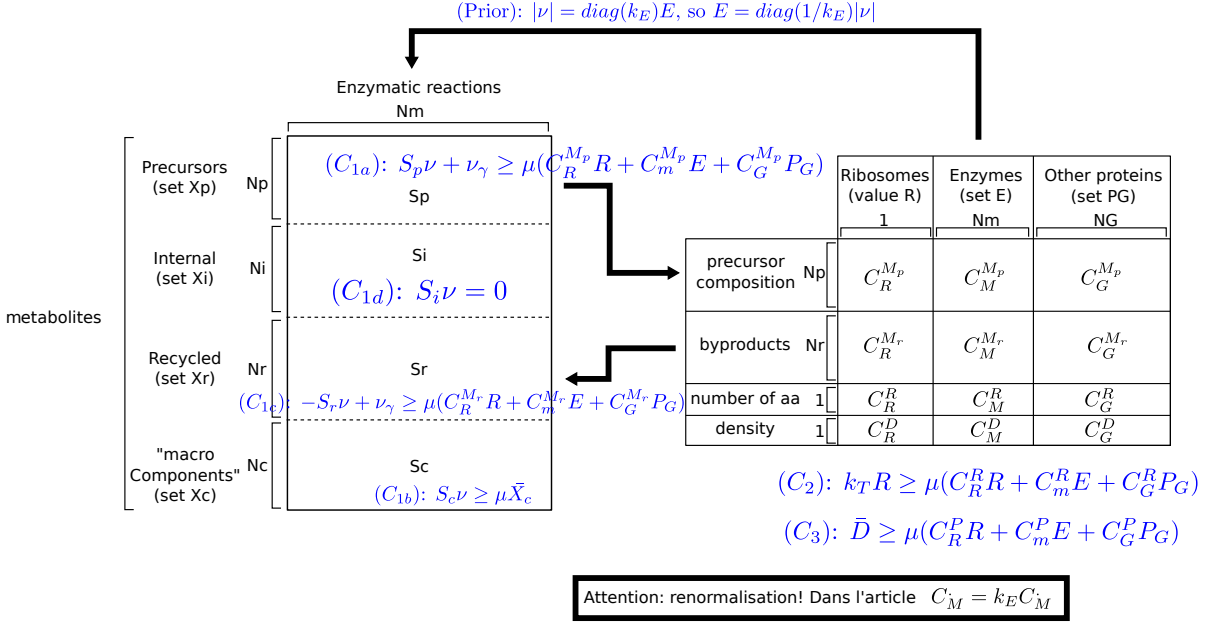


Figure 1: Description of cell by RBA.

# 2 Analysis of existing algorithm (RBAV01)

## 2.1 Structure of RBAV01

The original program is built around a multitude of functions and structures displayed on Figure 2. Figure 3 shows how matrices are built in the original algorithm.
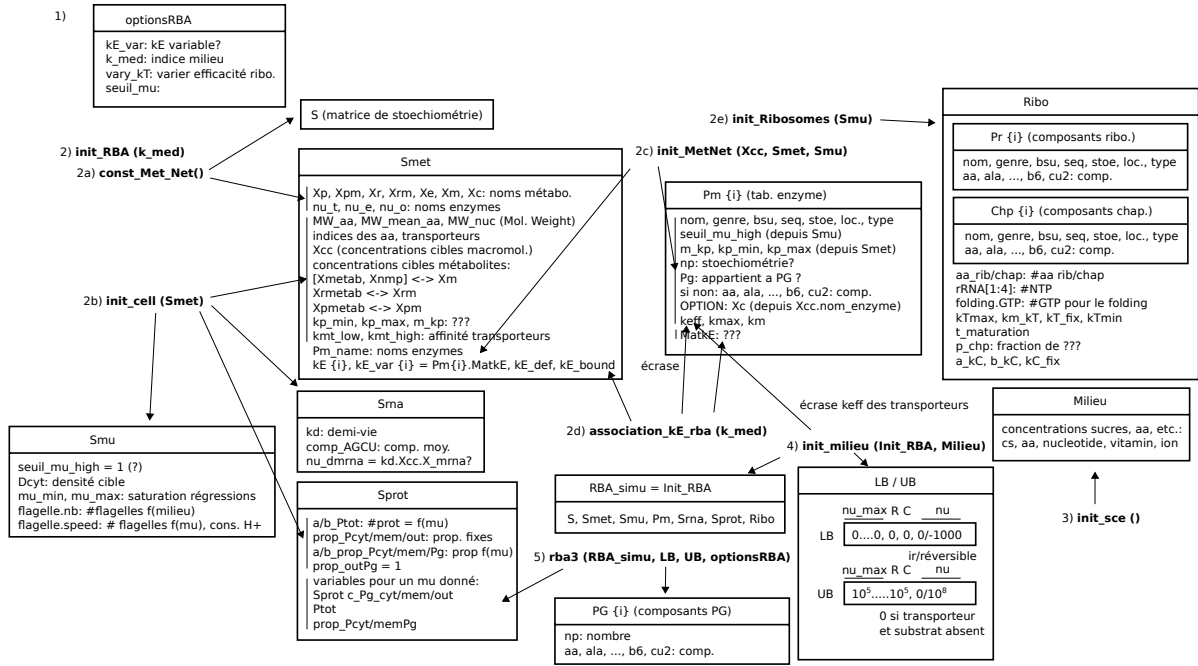
**1) optionsRBA**
kE_var: kE variable?
k_med: indice milieu
vary_kT: varier efficacité ribo.
seuil_mu:

S (matrice de stoechiométrie)

**2) init_RBA (k_med)**
**2a) const_Met_Net()**

**Smet**
Xp, Xpm, Xr, Xrm, Xe, Xm, Xc: noms métabo.
nu_t, nu_e, nu_o: noms enzymes
MW_aa, MW_mean_aa, MW_nuc (Mol. Weight)
indices des aa, transporteurs
Xcc (concentrations cibles macromol.)
concentrations cibles métabolites:
[Xmetab, Xnmp] <-> Xm
Xrmetab <-> Xrm
Xpmetab <-> Xpm
kp_min, kp_max, m_kp: ???
kmt_low, kmt_high: affinité transporteurs
Pm_name: noms enzymes
kE {i}, kE_var {i} = Pm{i}.MatkE, kE_def, kE_bound

**2b) init_cell (Smet)**

**Smu**
seuil_mu_high = 1 (?)
Dcyt: densité cible
mu_min, mu_max: saturation régressions
flagelle.nb: #flagelles f(milieu)
flagelle.speed: # flagelles f(mu), cons. H+

**Srna**
kd: demi-vie
comp_AGCU: comp. moy.
nu_dmrna = kd.Xcc.X_mrna?

**Sprot**
a/b_Ptot: #prot = f(mu)
prop_Pcyt/mem/out: prop. fixes
a/b_prop_Pcyt/mem/Pg: prop f(mu)
prop_outPg = 1
variables pour un mu donné:
Sprot c_Pg_cyt/mem/out
Ptot
prop_Pcyt/memPg

**2c) init_MetNet (Xcc, Smet, Smu)**

**Pm {i} (tab. enzyme)**
nom, genre, bsu, seq, stoe, loc., type
seuil_mu_high (depuis Smu)
m_kp, kp_min, kp_max (depuis Smet)
np: stoechiométrie?
Pg: appartient a PG ?
si non: aa, ala, ..., b6, cu2: comp.
OPTION: Xc (depuis Xcc.nom_enzyme)
keff, kmax, km
MatkE: ???
écrase

**2e) init_Ribosomes (Smu)**

**Ribo**
*Pr {i} (composants ribo.)*
nom, genre, bsu, seq, stoe, loc., type
aa, ala, ..., b6, cu2: comp.

*Chp {i} (composants chap.)*
nom, genre, bsu, seq, stoe, loc., type
aa, ala, ..., b6, cu2: comp.

aa_rib/chap: #aa rib/chap
rRNA[1:4]: #NTP
folding.GTP: #GTP pour le folding
kTmax, km_kT, kT_fix, kTmin
t_maturation
p_chp: fraction de ???
a_kC, b_kC, kC_fix

**2d) association_kE_rba (k_med)**
**4) init_milieu (Init_RBA, Milieu)**
écrase keff des transporteurs

**Milieu**
concentrations sucres, aa, etc.:
cs, aa, nucleotide, vitamin, ion

**3) init_sce ()**

RBA_simu = Init_RBA
S, Smet, Smu, Pm, Srna, Sprot, Ribo

**5) rba3 (RBA_simu, LB, UB, optionsRBA)**

**LB / UB**
nu_max R C    nu
LB   $0....0, 0, 0, 0/-1000$
ir/réversible
nu_max R C    nu
UB   $10^5.....10^5, 0/10^8$
0 si transporteur
et substrat absent

**PG {i} (composants PG)**
np: nombre
aa, ala, ..., b6, cu2: comp.

Figure 2: Algorithm used in RBAV01.

|  | nu_max | R | C | nu |
|---|---|---|---|---|
| **f (fonction objective)** | 1 | 0 | 0 | 0 |
| **D (densité cytosol)** | CD_M | CD_R | CD_C | |
| **M (densité membrane)** | CD_M | | | |
| **taille Nm (flux)** | -Pm.keff | | | 1 |
| **taille Nm (flux)** | -Pm.keff | | | -1 |
| **ATP maintenance** | | | | nu_ATP=-1 |

$= A \leq b =$

$$\begin{bmatrix} D - CD\_G.C\_PGcyt \\ Ptot.Pmem - CD\_G.C\_PGmem \\ \\ \\ -atp\_maintenance \end{bmatrix}$$

| | | | |
|---|---|---|---|
| **Ribosomes** | muCR_M | muCR_R -kT.p_act | muCR_C |
| **Chaperones** | muCC_M | muCC_R | muCC_C -kC |
| **Xe** | | | Se |
| **Xp** | muCMp_M | muCMp_R | muCMp_C  -Sp |
| **Xpm** | | | -Spm |
| **Xr** | muCMr_M | muCMr_R | muCMr_C  Sr |
| **Xrm** | | | Srm |
| **Xc** | | | -Sc |
| **Xm** | | | -Sm |
| **Coût flagelles** | | | nu_H+=-1 |

$= Aeq = beq =$

$$\begin{bmatrix} -CR\_G.mu.C\_PG \\ -CR\_G.pch.mu.C\_PG \\ 0 \\ -muCMp\_G + \begin{bmatrix}0 \\ -mu.Xp\_metab\end{bmatrix} \\ -muCMr\_G + \begin{bmatrix}0 \\ -mu.Xr\_m\end{bmatrix} \\ -muXc \\ -muX\_metab \\ -muX\_nmp \\ -h\_flagelle \end{bmatrix} + \begin{bmatrix} \\ \\ \\ -nu\_dmrna\ (X\_H2O) \\ \\ nu\_dmrna\ (X\_H+) \\ \\ -kd.Xc\ (X\_mrna) \\ comp\_AGCU.nu\_dmrna \end{bmatrix}$$

**mRNA**

Figure 3: Matrices used in RBAV01.

4

## 2.2 Orientations for improving RBA algorithm

It would be more elegant and flexible to not name any metabolite or flux sets. These sets should be defined by data, not the program. The existence of $X_{e/c/r/rm/p/pm/m}$ is unnecessary. Note that there are even more subsets than defined in theory, we probably need to avoid proliferation of such subsets.

We need to find clear abstraction that can represent RBA entities as generically as possible. All macromolecules should be representable with a composition matrix that would help build the final matrices rapidly. We also need to express the second member as data, in terms of processes.

Final objectives are

- Completely separate data from code: Create data files that are meaningful to users.

- Make code more compact and easier to read (and potentially quicker): by using a full matrix formalism (avoiding loops). This can only be achieved by finding the right abstraction to represent RBA elements.

# 3 Rewriting RBAs in full matrix form

**Composition vector for macromolecules**   In the following, each macromolocule is described by its **composition vector**. It is a column vector containing the metabolites necessary to build it, with a **minus** sign for metabolites consumed and a **plus** sign for byproducts generated.

**Metabolism constraint**   The metabolism constraint $(C_1)$ can be rewritten

$$\overbrace{\underbrace{S\nu}_{\substack{\text{metabolic flux}\\\text{generated by}\\\text{metabolism}}} + \underbrace{\mu[C_E, C_R, C_C][E; R; C]}_{\substack{\text{precursors}\\\text{used/byproducts}\\\text{generated by produc-}\\\text{ing new molecules}}}}^{\text{'Variable' terms}} = \overbrace{\underbrace{-\mu[C_G, C_{X_c}][P_G; X_c]}_{\substack{\text{precursors}\\\text{used/byproducts}\\\text{generated by produc-}\\\text{ing new molecules}}}}^{\text{Constant terms}}$$

$$-\operatorname{diag}(k_E) E \leq \nu \leq \operatorname{diag}(k_E) E$$

where $C_i$ are composition matrices as defined above.

Written this way, it actually does not matter to which group a metabolite belongs. Its group is only dictated by the structures of the composition matrices. Based on this constraint only it does not make sense to create metabolic groups *in the program*.

In full matrix form, the above equations become

$$\begin{cases} [S, \mu C_E, \mu C_R, \mu C_C][\nu; E; R; C] & = & -[\mu C_G, \mu C_{X_c}][P_G; X_c] \\ \begin{bmatrix} I, & -\operatorname{diag}(k_E), & 0, & 0 \\ -I, & -\operatorname{diag}(k_E), & 0, & 0 \end{bmatrix}[\nu; E; R; C] & \leq & 0 \end{cases}$$

The most important here is the composition matrix. In the data, the user gives the list of metabolites used for synthesis and the list of byproducts. The program then has to figure out how to reorder terms to build the matrix.

**Ribosome/chaperone constraints**   We can rewrite the ribosome constraint as follows:

$$k_T R + \mu[C_E^R, C_R^R, C_C^R][E; R; C] = -\mu C_G^R P_G$$

In full form:

$$[0, \mu C_E^R, k_T + \mu C_R^R, \mu C_C^R][\nu; E; R; C] = [-\mu C_G^R, 0][P_G; X_c]$$

The number of aas needed to build a protein can be included in the description of molecule sets. The same applies for folding by chaperones.

These additionnal constraints have the form

$$a[\nu; E; R; C] = b_0 + b_1[P_G; X_c]$$

this means they can be appended to $C_1$ by **simple line concatenation**.

**RNA degradation/replication**   In this type of constraints, we simply add new fluxes of metabolites. These fluxes are simply added up on the right hand side of $(C_1)$

$$
\begin{aligned}
... &= -\mu[C_G, C_{X_c}][P_G; X_c] - \mu^d C_{process} X_{process} \\
&= -[\mu C_G, \mu C_{X_c}, \mu^d C_{process}][P_G; X_c; X_{process}]
\end{aligned}
$$

where $d = 0$ when the target flux is an abolute flux and $d = 1$ when it compensates dilution.

**Density constraints**   The density constraint writes

$$
[0, C_E^D, C_R^D, C_C^D][\nu; E; R; C] \leq \overline{D} - [C_G^D, C_G^c][P_G; X_c]
$$

In general, these constraints can be expected to be of the form

$$
a[\nu; E; R; C] \leq b_0 + b_1[P_G; X_c]
$$

**Maintenance ATP constraint (metabolite specific constraints)**   It is defined as

$$
\mathbf{1}_{\nu_{constraint}}[\nu; E; R; C] = b
$$

where $\mathbf{1}_{\nu_{constraint}}$ is an indicator matrix selecting the reaction associated with the production of maintenance ATP/flagella fuel. Generally this means a specific reaction has to be added in $S$, as well as specific metabolites like ATP maintenance or protons dedicated to fuelling flagella.

**Summary**

$$
\begin{aligned}
[S, \mu C_E, \mu C_R, \mu C_C] \qquad\quad [\nu; E; R; C] &= -\sum \mu^{d_i} C_i X_i && \text{(Base metabolism} \\
&&& \text{and process production)} \\
a \qquad\quad [\nu; E; R; C] &= b_0 + b_1[\{X_i\}] && \text{(Process capacity)} \\
\mathbf{1}_{\nu_{constraint}} \qquad\quad [\nu; E; R; C] &= b && \text{(Metabolite constraints)} \\
\begin{bmatrix} I, & -\text{diag}\,(k_E), & 0, & 0 \\ -I, & -\text{diag}\,(k_E), & 0, & 0 \end{bmatrix} [\nu; E; R; C] &\leq 0 && \text{(Enzymatic fluxes)} \\
a \qquad\quad [\nu; E; R; C] &\leq b_0 + b_1[\{X_i\}] && \text{(Density constraints)}
\end{aligned}
$$

A more graphic representation of the matrix is provide in Figure 4.

**f (objective function)** $\begin{bmatrix} \begin{bmatrix} 0 \end{bmatrix} & \begin{bmatrix} 1 \end{bmatrix} & 0 & 0 \end{bmatrix}$

**Metabolism** — S, muC_E, muC_R, muC_C

**Ribosomes** — muCR_M, muCR_R -kT, muCR_C

**Chaperones** — muCC_M, muCC_R, muCC_C -kC

**= Aeq = beq = -**

muC_G  muC_Xc  muC_Xm  C_mrna

muCR_G
muCC_G

**x** [ nu  E  R  C ]    **x** [ C_PG  Xc  Xmetab  Xmrna ]

**Flux** — I, -Pm.keff

**Flux** — -I, -Pm.keff

**D (densité cytosol)** — CD_M, CD_R, CD_C

**M (densité membrane)** — CD_M

**ATP maintenance** — nu_ATP=-1

**= A ≤ b =**

D_c
D_m
-atp_m

**-**

muCD_G
muCD_G

**x** [ nu  E  R  C ]    **x** [ C_PG  Xc  Xmetab  Xmrna ]

Figure 4: Structure of matrices. Note that the left hand side and the right hand side have a similar structure (except for the $\nu$ submatrix). Each column represents a set of molecules: on the right hand side, the production flux of the set has to be determined, on the left hand side it is already given *a priori*.

# 4 New formalism

The new formalism must allow simple interactions between user readable data and efficient construction of the RBA problem. The data is going to be written is XML format, inspired by SBML standards.

## 4.1 Final differences between RBAv01 and new formalism

|  | RBAv01 mat12 | RBAv01 mat15 | RBAnew mat12 | RBAnew mat15 |
|---|---|---|---|---|
| Input files | Custom files + partially hard coded | | XML files (nothing hard coded) | |
| Code length (commented) | $\simeq 3500$ ($\simeq 1000$) | | $\simeq 1500$ ($\simeq 500$) | |
| Parsing | 15s | 8s | 12s | 12s |
| Solving (23 rounds) | 15s | 5s | 1.5-2s | 1.5-2s |
| 1 matrix update | 580ms | 100ms | 3-10ms | 3-10ms |
| 1 CPLEX round | 50ms | 58ms | 50ms | 58ms |

Table 1: Comparison between original algorithm (RBAv01) and new formalism (RBAnew). Because performance for the old algorithm depended on matlab version, we included performance for matlab R2012a (mat12) and R2015b (mat15).

Table 1 shows the main differences between RBAv01 and the new algorithm. The new algorithm uses generic XML files and is thus easier to modify for the end user. Parsing files is less efficient than the original algoirthm but this is not very important as it is only done once. More important is the solving time, which is significantly lower than the original algorithm.

The new algorithm could be improved further by handling block-allocation more efficiently for sparse matrices, but current performance seems pretty good.

## 4.2 Input data



Standard SBML file, see SBML specifications. Compartments are needed for density constraints. Species and reactions are needed to build the stoichiometry matrix.

Figure 5: Metabolism file is a standard SBML file.

**Metabolism file**  The metabolism file is a standard SBML file (Fig. 5). It contains information about cell compartments, metabolite species and metabolism reactions. Concentration of external metabolites are defined in the species section.



```
−<RBAParameters>
  −<listOfMaximalDensities>
     <maximalDensity compartment="c" value="4.8972"/>
    −<maximalDensity compartment="mp">
       <functionReference function="protein_concentration"/>
       <functionReference function="average_protein_weight"/>
       <functionReference function="fraction_membrane_protein"/>
     </maximalDensity>
   </listOfMaximalDensities>
  −<listOfFunctions>
    −<function id="protein_concentration" type="linear">
      −<listOfParameters>
         <parameter id="LINEAR_COEF" value="-0.0048302"/>
         <parameter id="LINEAR_CONSTANT" value="0.031256"/>
         <parameter id="X_MIN" value="0.25"/>
         <parameter id="X_MAX" value="1.6"/>
         <parameter id="Y_MIN" value="-Inf"/>
         <parameter id="Y_MAX" value="Inf"/>
       </listOfParameters>
     </function>
    +<function id="average_protein_weight" type="constant"></function>
    +<function id="fraction_cytosol_protein" type="constant"></function>
    +<function id="fraction_membrane_protein" type="constant"></function>
    +<function id="fraction_external_protein" type="constant"></function>
    +<function id="fraction_nonenzymatic_cytosol_protein" type="linear"></function>
    +<function id="fraction_nonenzymatic_membrane_protein" type="linear"></function>
    +<function id="fraction_nonenzymatic_external_protein" type="constant"></function>
    +<function id="number_flagella" type="linear"></function>
    +<function id="flagella_speed" type="constant"></function>
    +<function id="flagella_H_consumption" type="constant"></function>
    +<function id="ribosomeEfficiencyMM" type="michaelisMenten"></function>
    +<function id="ribosomeEfficiencyCM" type="constant"></function>
    +<function id="fractionActiveRibosomes" type="exponential"></function>
    +<function id="chaperoneEfficiencyLM" type="linear"></function>
    +<function id="chaperoneEfficiencyMedium1" type="constant"></function>
    +<function id="chaperoneEfficiencyMedium2" type="constant"></function>
    +<function id="chaperoneEfficiencyMedium3" type="constant"></function>
    +<function id="chaperoneEfficiencyMedium4" type="constant"></function>
    +<function id="chaperoneEfficiencyMedium5" type="constant"></function>
    +<function id="maintenanceATP" type="linear"></function>
   </listOfFunctions>
 </RBAParameters>
```

Density constraints: up to one constraint per compartment defined in the SBML file.

User defined functions: they describe how a number of cell parameters evolve with growth rate (the variable of these functions is mu). Any number of functions can be defined, but a restricted number of type of functions is avalaible (indicator, constant, linear, exponential, michaelisMenten). For every type of function a given set of parameters must be specified.
These functions can then be used to define other parameters, such as the maximal density per compartment or target values for processes.

Figure 6: Structure of the parameter file used by RBA.

**Parameter file**  The parameter file is an XML file composed of two subsections (Fig. 6). `listOfMaximalDensities` contains density constraints. `listOfFunctions` contains user-defined functions that can be used to set $\mu$-dependent maximal densities or $\mu$-dependent process targets (see below).

**Macromolecule files**  Currently, RBA defines three sets of macromolecules: proteins, RNAs and DNA. One XML file is created for each set (Fig. 7). All macromolecules have a `listOfComponents` describing the building blocks used to produce them (*e.g.* amino acids, vitamins and ions for proteins). Then a `listOfSpecies` contains all molecules of the set, including their description in terms of components. Additionally, proteins have a `listOfEfficiencyFunctions` that lists efficiency models for enzymes. Enzymes contain an `enzymaticActivity` structure that defines the reaction they catalyze and the parame-

ters for the efficiency models. Finally, transporters have a `transporterEfficiency` that modulates the enzymatic activity depending on substrate and cofactor availability.

**Process file** The process file is an XML file containing a `listOfProcesses` and a `listOfComponentMaps`. Each `process` can contain up to 3 subsections. The `capacityConstraint` defines a machinery used by the process that has a limited capacity. The `operatingCosts` defines which macromolecules the process produces/degrades/modifies and the cost associated with these operations. The `targets` are set fluxes that a process must maintain in order for the cell to work properly. Target fluxes can apply to metabolites (`targetValue`) and reactions (`targetReaction`). Target fluxes scale with $\mu$ if they contribute to `dilution_compensation` or if they are defined using a $\mu$-dependent user-function. Finally `componentMaps` are used to compute the costs in the `operatingCosts` section.

```xml
-<process cellDesign="yes" compartment="c" id="P_TA" name="Translation apparatus">——— Process definition requires id and compartment, rest is optional.
  -<capacityConstraint>
    -<machineryComposition>
      +<listOfReactants></listOfReactants>
      +<listOfProducts></listOfProducts>
    </machineryComposition>
    -<capacity>
      <functionReference function="ribosomeEfficiencyMM"/>
      <functionReference function="fractionActiveRibosomes"/>
    </capacity>
  </capacityConstraint>
  -<operatingCosts>
    <production componentMap="translation" set="protein"/>
  </operatingCosts>
  -<targets>
    -<targetValue species="Pg_c">
      <functionReference function="protein_concentration"/>
      <functionReference function="fraction_cytosol_protein"/>
      <functionReference function="fraction_nonenzymatic_cytosol_protein"/>
    </targetValue>
    +<targetValue species="Pg_e"></targetValue>
    +<targetValue species="Pg_mp"></targetValue>
  </targets>
</process>
```

(optional) A capacity constraint can be associated with the process. In this case, a machinery must be associated with the process as a list of reactants (metabolites or macromolecules) used to assemble a functional machine and a list of byproducts generated in the process. A capacity function, defined as the product of user-defined function, sets how many units a machine can process.

(optional) If the process produces, degrades or modifies a set of macromolecules. it is specified in this node. A componentMap is used to compute metabolites used and byproducts generated during the process.

(optional) Fluxes that must be produced by the process at a preset value. These can be metabolite or reaction fluxes. The flux value can be a preset value or the product of user-defined functions. It can be an absolute flux or a dilution_compensation flux. In latter case, it serves as a way to maintain the target at constant concentration, and the flux value actually scales with the growth rate.

```xml
  -<targets>
    -<targetReaction reaction="Th">
      <functionReference function="number_flagella"/>
      <functionReference function="flagella_speed"/>
      <functionReference function="flagella_H_consumption"/>
    </targetReaction>
  </targets>
  -<targets>
    <targetValue species="mrna" value="0.01"/>
    <targetValue species="trna" value="0.02625"/>
    <targetValue species="trna2" value="0.01125"/>
    <targetValue dilution_compensation="0" species="mrna" value="0.15996"/>
  </targets>
```

```xml
<RBAProcesses>
  +<listOfProcesses></listOfProcesses>
  +<listOfComponentMaps></listOfComponentMaps>
</RBAProcesses>
```

Global structure: a list of processes of the cell and a list of component maps used by processes to build/degrade polymers.

```xml
-<listOfComponentMaps>
  -<componentMap id="translation">
    -<constantCost>
      -<listOfReactants>
        <speciesReference species="m_tfmet" stoichiometry="1"/>
        <speciesReference species="m_gtp" stoichiometry="1"/>
        <speciesReference species="m_h2o" stoichiometry="1"/>
      </listOfReactants>
      -<listOfProducts>
        <speciesReference species="m_trna" stoichiometry="1"/>
        <speciesReference species="m_fmet" stoichiometry="1"/>
        <speciesReference species="m_gdp" stoichiometry="1"/>
        <speciesReference species="m_p" stoichiometry="1"/>
        <speciesReference species="m_h" stoichiometry="1"/>
      </listOfProducts>
    </constantCost>
    -<cost component="ala" processingCost="1">
      -<listOfReactants>
        <speciesReference species="m_tala" stoichiometry="1"/>
        <speciesReference species="m_gtp" stoichiometry="2"/>
        <speciesReference species="m_h2o" stoichiometry="2"/>
      </listOfReactants>
      -<listOfProducts>
        <speciesReference species="m_trna" stoichiometry="1"/>
        <speciesReference species="m_gdp" stoichiometry="2"/>
        <speciesReference species="m_p" stoichiometry="2"/>
        <speciesReference species="m_h" stoichiometry="2"/>
      </listOfProducts>
    </cost>
    +<cost component="arg" processingCost="1"></cost>
    +<cost component="asn" processingCost="1"></cost>
```
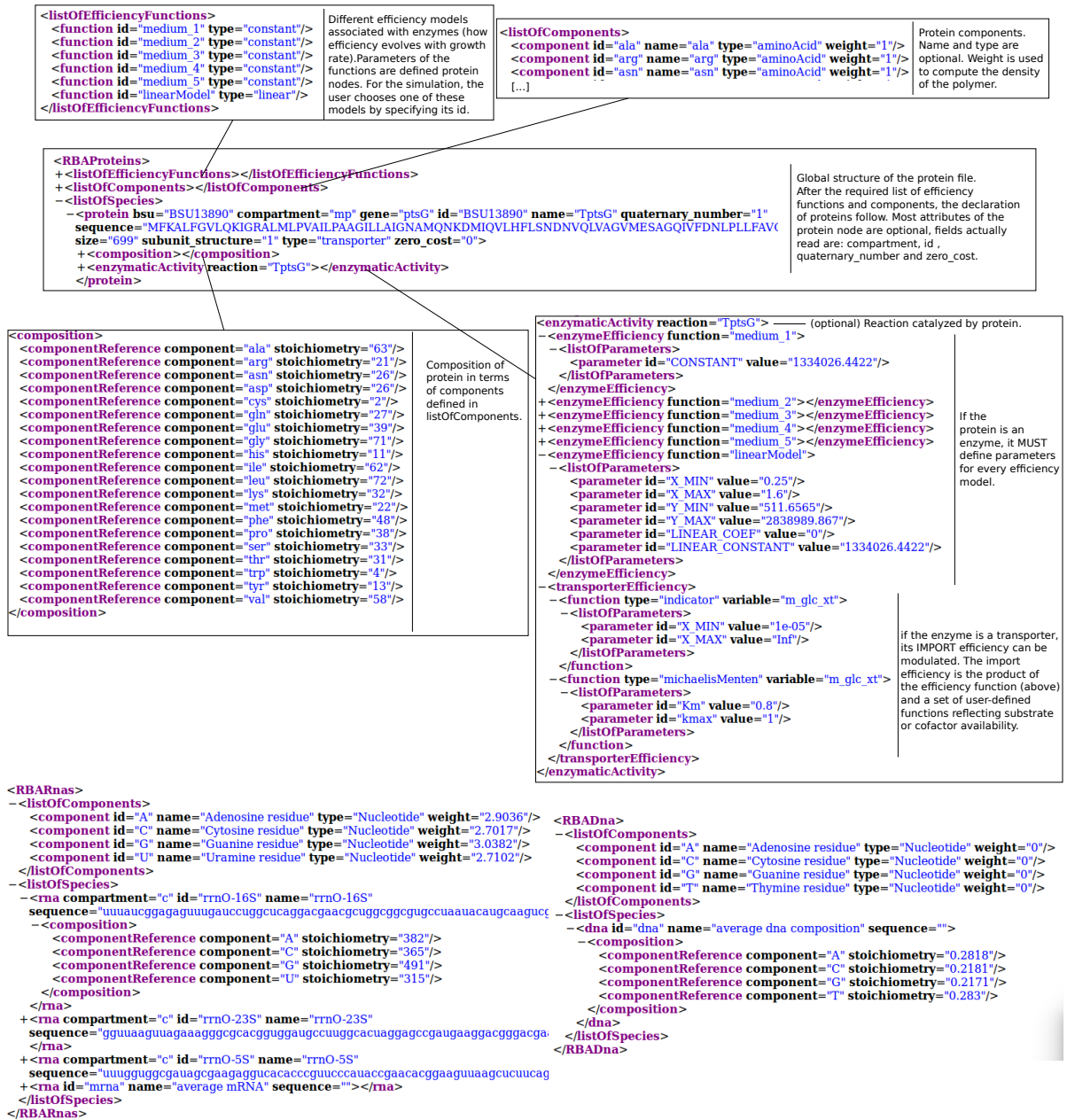
(optional) A component map may have a constant cost part that is independent of a polymer's composition. A processing cost may be associated with this node.

For every component of the macromolecule, the component map defines the metabolites consumed/generated during its processing (production, degradation or modification). If the process has a capacity constraint, a processingCost can also be associated to each component. The number of components the machinery can process is then limited by its capacity (explicitly: the sum of processing costs cannot exceed the machine's capacity).

Figure 8: Structure of process file.

## 4.3 From data to structures

**Algorithm overview**  Figure 9 shows how the XML files are transformed into matlab structures. Note that a lot of elements are already in matrix form.
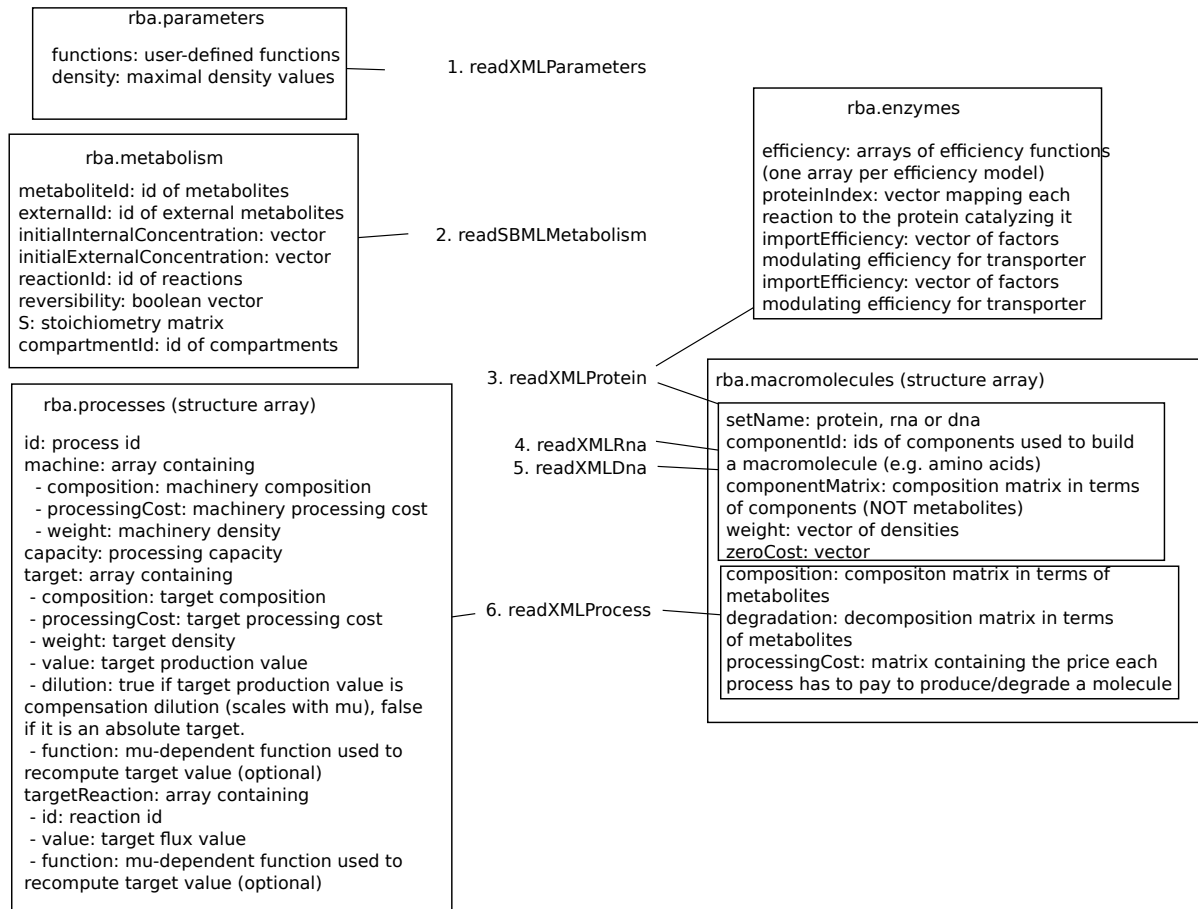


Figure 9: Algorithm used to read data in the new formalism.

**Composition matrices**  Figure 10 shows how the matrices stored in `rba.macromolecules` are built.

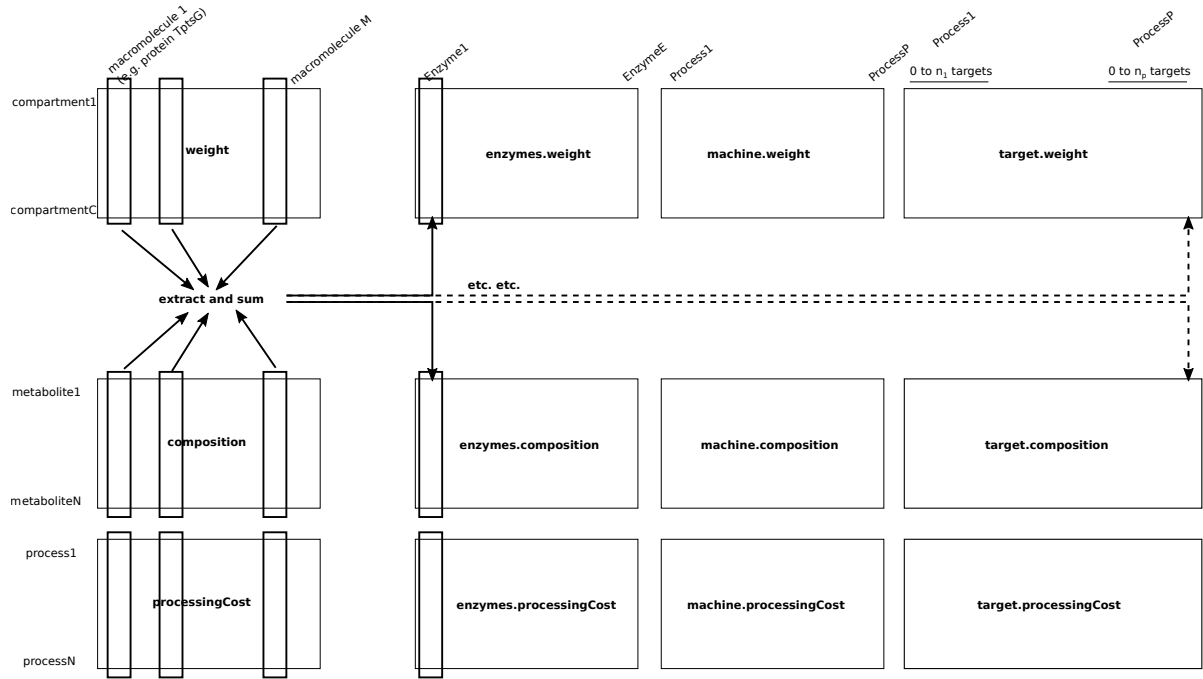**Process-related matrices**  Figure 11 shows how the matrices stored in `rba.processes` are built.

## 4.4 Building final matrices

Figure 12 shows how the optimization matrices stored are built.

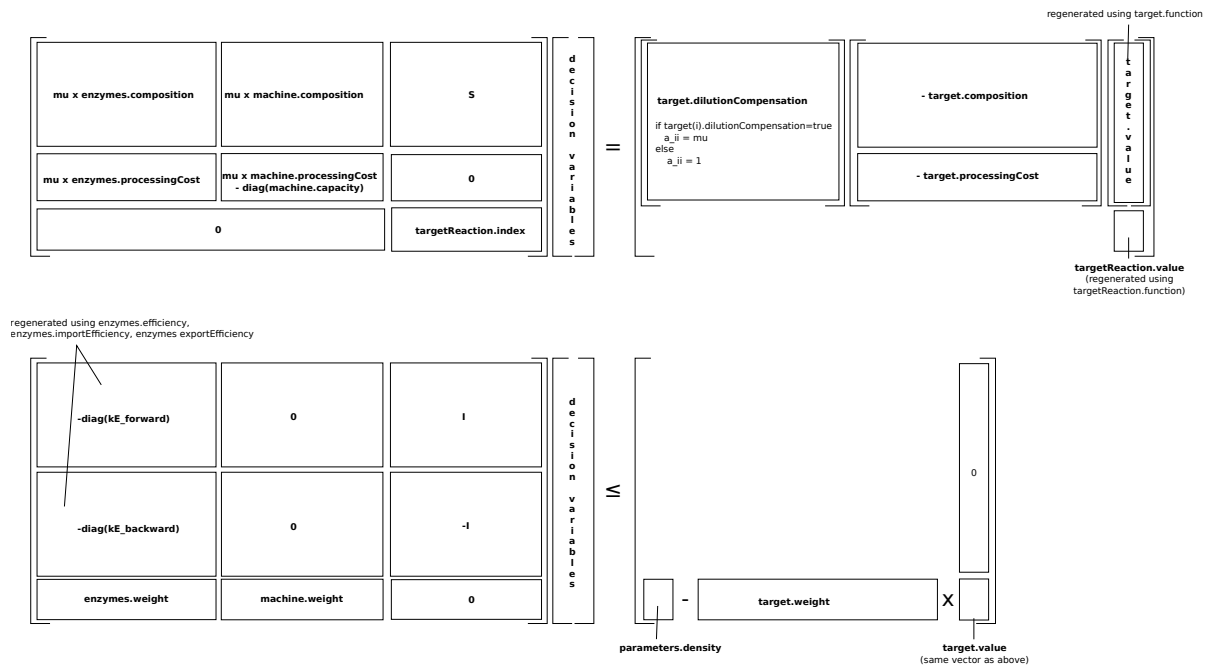Figure 10: Macromolecule matrices.

Figure 11: Process-related matrices.



Figure 12: Matrices used to solve the optimization problem.