



**HAL**  
open science

# Rewrite Rules and Strategies for Molecular Graphs

Oana Andrei, H el ene Kirchner

► **To cite this version:**

Oana Andrei, H el ene Kirchner. Rewrite Rules and Strategies for Molecular Graphs. 2007. inria-00146362v1

**HAL Id: inria-00146362**

**<https://inria.hal.science/inria-00146362v1>**

Submitted on 14 May 2007 (v1), last revised 12 Nov 2007 (v3)

**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 destin ee au d ep ot et  a la diffusion de documents scientifiques de niveau recherche, publi es ou non,  emanant des  tablissements d'enseignement et de recherche fran ais ou  trangers, des laboratoires publics ou priv es.

# Rewrite Rules and Strategies for Molecular Graphs

Oana Andrei<sup>1</sup> and H el ene Kirchner<sup>2</sup>

<sup>1</sup> INRIA - LORIA

<sup>2</sup> CNRS - LORIA

Campus scientifique BP 239

F-54506 Vandoeuvre-l es-Nancy Cedex, France

Firstname.Lastname@loria.fr

**Abstract.** In this paper, we present a rewriting framework for modeling molecular complexes, biochemical reaction rules and generation of biochemical networks based on the representation of molecular complexes as a particular type of multigraphs with ports called molecular graphs. The advantage of this approach is to obtain for free a rewriting calculus which allows defining at the same level transformation rules and strategies to control rule application and network generation.

## 1 Introduction

Many formal notations have already been used for modeling biological systems, and, in particular, protein-protein interactions concerned with the connectivity inside molecular complexes. These models can be classified as qualitative or quantitative, and the employed formal notations range from process calculi with different extensions or variations [9, 10, 16, 24], to term-rewriting [27] and graph rewriting [7].

In the biochemical model we consider, the behavior of a protein is given by its functional domains that determine which other protein it can bind to or interact with. These domains are usually abstracted as *sites* that can be bound or free, visible or hidden. A protein is characterized by the collection of interaction sites on its surface, also called *interface*. Proteins can bind to each other forming molecular complexes. Membranes can also form complexes, called tissues, due to the binding proteins on their surfaces. The structure of a complex is naturally described as a particular class of graphs, a mathematical structure that is easy to understand and use by computer scientists as well as by chemists and biologists.

We propose in this paper a calculus based on multigraph rewriting and rewrite strategies for modeling biochemical systems, in particular the structure and interactions of protein complexes. We use an extended version of multigraphs with ports [3] for describing the structure of complexes as molecular graphs and multigraph rewriting for modeling the interactions between them (Sect. 3). We encode the molecular graphs as terms, the reaction patterns as rewrite rules, and the transformation on molecular graphs as a rewriting relation (Sect. 4). We define a rewriting calculus for molecular graphs  $\rho_{bio}$ -calculus, obtained from the

rewrite calculus for labeled multigraphs with ports, the  $\rho_{mg}$ -calculus, introduced in [3], by adding state information on ports and imposing some conditions on edges. Strategic rewriting allows modeling the control mechanism in biochemical systems and the automatic generation of biochemical networks. This is illustrated in Sect. 5 on a fragment of the epidermal growth factor receptor (EGFR) signaling cascade.

## 2 Background

In this section we briefly review some basic definitions of sorted term algebra, term rewriting [4] and strategic rewriting, graph theory and labeled multigraphs with ports developed in [3].

**Term Algebra.** A *many-sorted signature* is a pair  $(S, \Sigma)$ , where  $S$  is called the *sort set* and  $\Sigma$  is an  $S^* \times S$ -sorted family  $\{\Sigma_{w,s} \mid w \in S^* \text{ and } s \in S\}$ . For  $\mathcal{X} = \{\mathcal{X}_s\}$  an  $S$ -sorted family of disjoint sets of variables,  $\mathcal{T}_\Sigma(\mathcal{X})$  is the smallest set of  $\Sigma$ -terms over  $\mathcal{X}$  built with operators from  $\Sigma$  and variables from  $\mathcal{X}$ . If  $\mathcal{X}$  is empty, then we write  $\mathcal{T}_\Sigma$  for  $\mathcal{T}_\Sigma(\emptyset)$  and call it the set of *ground  $\Sigma$ -terms*. The set of all terms of sort  $s$  is denoted by  $\mathcal{T}_{\Sigma,s}(\mathcal{X})$ . A *ground substitution* is a partial mapping from  $\mathcal{X}$  to ground  $\Sigma$ -terms and it uniquely extends to a  $\Sigma$ -homomorphism from  $\mathcal{T}_\Sigma(\mathcal{X})$  to  $\mathcal{T}_\Sigma$ . A finite ground substitution has the form  $\sigma = \{x_1 \mapsto t_1, \dots, x_n \mapsto t_n\}$ , with  $x_i \in \mathcal{X}$  and  $t \in \mathcal{T}_\Sigma$ , for all  $i = 1..n$ . Considering  $S$  as a partially ordered set leads to quite similar definitions for order-sorted signature, terms and substitutions [20].

**Term Rewriting.** A set  $\mathcal{R}$  of *rewrite rules* is a set of ordered pairs of terms of  $\mathcal{T}_\Sigma(\mathcal{X})$ , denoted  $l \rightarrow r$ , such that  $l$  and  $r$  belong to the same sort,  $l \notin \mathcal{X}$  and  $Var(r) \subseteq Var(l)$ . In this paper, we only consider finite sets of rewrite rules. The *rewriting relation* induced by  $\mathcal{R}$  is denoted by  $\rightarrow_{\mathcal{R}}$  ( $\rightarrow$  if there is no ambiguity on  $\mathcal{R}$ ), and defined by  $s \rightarrow t$  iff there exists a substitution  $\sigma$  and a position  $p$  in  $s$  such that  $s = s[\sigma l]_p$  for some rule  $l \rightarrow r$  of  $\mathcal{R}$ , and  $t = s[\sigma r]_p$ . This is written  $s \xrightarrow{\mathcal{R}, p, l \rightarrow r, \sigma} t$  where either  $\mathcal{R}$ ,  $p$ ,  $l \rightarrow r$ , or  $\sigma$  may be omitted. The reflexive transitive closure of the rewriting relation induced by  $\mathcal{R}$  is denoted by  $\xrightarrow{*}_{\mathcal{R}}$ , and by  $\xrightarrow{+}_{\mathcal{R}}$  its respective transitive closure. A rewriting derivation is a chain of terms  $t_1 \rightarrow t_2 \rightarrow \dots \rightarrow t_n$ . The *source* of a derivation  $t_1 \rightarrow t_2 \rightarrow \dots$  is  $t_1$ , and when the derivation is finite, its last term is called its *target*.

A *rewrite theory* over  $\mathcal{T}_\Sigma(\mathcal{X})$  is a triple  $(\Sigma, E, \mathcal{R})$  where  $E$  is a finite set of (sort-preserving) equalities and  $\mathcal{R}$  a finite set of rewrite rules. The relation  $\rightarrow_{\mathcal{R}/E}$  on  $\mathcal{T}_\Sigma(\mathcal{X})$  is  $=_E; \rightarrow_{\mathcal{R}}; =_E$  and it induces a relation  $\rightarrow_{\mathcal{R}/E}$  on the quotient algebra  $\mathcal{T}_\Sigma(\mathcal{X})/_E$  by  $[t]_E \rightarrow_{\mathcal{R}/E} [t']_E$  iff  $t \rightarrow_{\mathcal{R}, E} t'$ . The relation  $\rightarrow_{\mathcal{R}, E}$  (called rewriting modulo  $E$ ) on  $\mathcal{T}_\Sigma(\mathcal{X})$  is defined by  $s \rightarrow_{\mathcal{R}, E} t$  iff there exists a substitution  $\sigma$  and a position  $p$  in  $s$  such that  $s|_p =_E \sigma l$  for some rule  $l \rightarrow r$  of  $\mathcal{R}$ , and  $t = s[\sigma r]_p$ . In this paper, we will consider theories with  $E$  defined by specific axioms, namely associativity (A), associativity and commutativity (AC) and unit element (U) of some function symbols.

**Strategic Rewriting.** The notions of strategy and strategic rewriting have been introduced in order to control rewriting derivations. The notion of strategy

$[\mathbf{id}](t) = \{t\}$	$[\mathbf{onestep}(\zeta_1, \dots, \zeta_n)](t) = \bigcup_i [\zeta_i](t)$
$[\mathbf{fail}](t) = \emptyset$	$[\mathbf{seq}(\zeta_1, \zeta_2)](t) = [\zeta_2]([\zeta_1](t))$
$[l \rightarrow r](t) = \{s \mid t \xrightarrow{l \rightarrow r} s\}$	$[\mathbf{choice}(\zeta_1, \zeta_2)](t) = \begin{cases} [\zeta_1](t), & \text{if } [\zeta_1](t) \neq \emptyset \\ [\zeta_2](t), & \text{if } [\zeta_1](t) = \emptyset \end{cases}$

Fig. 1. Strategy constructors

can be defined in a general way: a *rewrite strategy*  $\zeta$  for the rewrite system  $\mathcal{R}$  is a subset of the set of all derivations of  $\mathcal{R}$ . A derivation of the strategy is called a *strategic rewriting derivation*. The *application of a strategy*  $\zeta$  on a term  $t$  is denoted  $[\zeta](t)$  and defined as the set of all targets  $t'$  of the derivations of source  $t$  in  $\zeta$ . When no derivation in  $\zeta$  has  $t$  as source, we say that the strategy application on  $t$  fails. The result of the application of a failing strategy on a term  $t$  is the empty set. We extend the application of a strategy to a set of terms as the union of the application of the strategy on each element of the set.

A strategy can be described by enumerating all its elements or more suitably by a *strategy language*. From elementary strategy expressions directly issued from a rewrite system  $\mathcal{R}$ , more elaborated strategy expressions are built like in ELAN [22], Stratego [28], TOM [5] or, more recently, Maude [17]. The semantics of such a language is naturally described in the rewriting calculus [12, 14]. We describe below some constructs of the strategy language of interest in this paper and available in TOM<sup>1</sup>.

Given a rewrite system  $\mathcal{R}$  over  $\mathcal{T}_\Sigma(\mathcal{X})$ , a strategy expression is either a rewrite rule in  $\mathcal{R}$  or an expression described below. A strategy operator  $\zeta$  may take arguments  $\zeta_1, \dots, \zeta_n$ , and the result is expressed functionally by  $\zeta(\zeta_1, \dots, \zeta_n)$ . We present in Fig. 1 some strategy constructors used in this paper. The strategy **onestep** computes all derivations issued from the application of a nonempty set of strategies. The strategy **seq**( $\zeta_1, \zeta_2$ ) first applies  $\zeta_1$  and then, if that succeeds, it applies  $\zeta_2$ ; it fails if either  $\zeta_1$  fails, or  $\zeta_2$  fails. **choice**( $\zeta_1, \zeta_2$ ) applies the strategy  $\zeta_1$ ; if the application of  $\zeta_1$  fails, it applies  $\zeta_2$ . Therefore **choice**( $\zeta_1, \zeta_2$ ) fails if both  $\zeta_1$  and  $\zeta_2$  fail. Both strategy operators **seq** and **choice** extend naturally to be applicable to a list of strategies.

Along with the elementary rewrite rules, the strategy constructors represent the key-components necessary for defining more complex strategies. For instance, **try** and **repeat** allow iterating the application of a strategy and never fail.

$$\begin{aligned} \mathbf{try}(\zeta) &= \mathbf{choice}(\zeta, \mathbf{id}) \\ \mathbf{repeat}(\zeta) &= \mathbf{try}(\mathbf{seq}(\zeta, \mathbf{repeat}(\zeta))) \end{aligned}$$

Other high-level strategies implement term traversal and normalization on terms and are well-known in the term rewriting literature (see [5] for instance).

**Labeled Graphs.** A *label alphabet*  $\mathcal{L} = (\mathcal{L}_V, \mathcal{L}_E)$  is a pair of sets of node labels and edge labels. A (finite) *graph* over  $\mathcal{L}$  is a triple  $G = (V, E, l)$  where  $V$  is a

<sup>1</sup> Web page <http://tom.loria.fr>

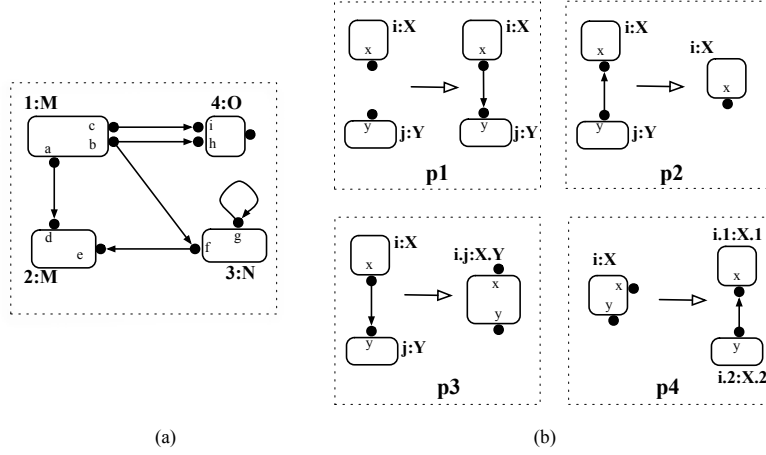


Fig. 2. A multigraph with ports (a) and some multigraph transformations (b)

set  $\{v_1, \dots, v_n\}$  of elements called *nodes*,  $E$  is a family  $(e_1, \dots, e_m)$  of elements of the Cartesian product  $V \times V$  called *edges*, and  $l = (l_V, l_E)$  is the *labeling function* for nodes ( $l_V : V \rightarrow \mathcal{L}_V$ ) and edges ( $l_E : E \rightarrow \mathcal{L}_E$ ). If  $G$  is a graph, we denote by  $V_G$  its node set, by  $E_G$  its edge set, and by  $l_G$  its labeling function. An edge of the form  $(v, v)$  is called a *loop*. For an edge  $(v_1, v_2)$ ,  $v_1$  and  $v_2$  are called *end nodes* (or *end points*) with  $v_1$  the *source* and  $v_2$  the *target*; moreover we say that  $v_1$  and  $v_2$  are *adjacent* or *neighbouring* nodes, with  $v_2$  neighbour of  $v_1$ . An edge is *incident* to a node if the node is one of its endnodes. An edge is *multiple* if there is another edge with the same source and target; otherwise it is *simple*. A graph is *simple* if it has no multiple edges or loops, otherwise it is a *multigraph*. An *adjacency list* for a node is given by a list of pairs consisting of a neighbour and the corresponding edge label.

**Labeled Multigraphs with Ports.** A *labeled multigraph with ports* is obtained from a labeled multigraph by associating to each node a set of *ports* such that each elementary edge is determined by two ports, one from each endnode. A node label consists of a node identifier and a set of ports, while a simple edge label is the pair of source and target ports. Hereinafter we use multigraph instead of labeled multigraph with ports if there is no risk of confusion. We illustrate in Fig. 2(a) such a multigraph.

**Transforming Multigraphs with Ports.** A *transformation rule on multigraphs with ports*  $L \rightsquigarrow R$  consists of two multigraphs  $L$  and  $R$  called the left- and right-hand side respectively, and a correspondence between nodes of  $L$  and nodes of  $R$ ,  $\xi : N_L \rightarrow \mathcal{P}(N_R)$ , which we call *node-substitution*. This correspondence is provided by some unique identifiers associated to nodes. The application of a multigraph transformation rule  $L \rightsquigarrow R$  to a multigraph  $G$  produces a new multigraph  $G'$  by the following steps: first, find and remove a subgraph of  $G$

isomorphic with  $L$  via a matching morphism  $m$ , resulting in a context graph  $G^-$ ; then add and reconnect  $m(R)$  to  $G^-$  with the help of  $m(\xi)$ .

We illustrate in Fig. 2 (b) some simple multigraph transformations: **p1** creates an edge between two ports, **p2** removes the source node of an edge, **p3** merges two adjacent nodes, and **p4** splits a node containing (at least) two ports in two adjacent nodes by distributing a port to each node.

**The Rewriting Calculus for Multigraphs with Ports.** The rewriting calculus for labeled multigraphs with ports (or  $\rho_{mg}$ -calculus), fully detailed in [3], is a variation of the  $\rho$ -calculus [13] with particular features related to multigraph transformation. It provides a clear operational semantics to multigraph transformation. We benefit from its capability of encoding rewrite strategies and conditional rewrite rules. Moreover the  $\rho$ -calculus integrates  $\lambda$ -calculus and rewriting in a powerful higher-order pattern calculus.

### 3 The Biochemical Model

A *molecular complex* consists of molecules connected through bonds between sites. A molecule is characterized by a pair of name and interface (set of sites). In order to make a distinction between molecules having the same name, we assign to each molecule a *unique identifier* consisting of a sequence of integers. Each site is characterized by a name and a state: **b** for bound, **v** and **h** for free site that are visible and hidden respectively.

We represent a molecular complex as a connected labeled multigraph with ports where molecules are nodes, sites are ports with states, and bonds are simple edges, together with the restriction that a port can be the end node of at most one edge. We call this kind of multigraphs, not necessarily connected, **molecular graphs**. In Fig. 3 we define their abstract syntax, with “,” the juxtaposition operator and  $\epsilon$  the empty set. An unbound molecule has also a molecular graph representation as a single node with an empty list of bonds.

Using the same graphical representation as for multigraphs, inspired by [7, 16], we represent a molecule as an empty box having the identifier placed at the exterior and the sites as small points on the surface of the box. We represent the state of a site as follows: a filled circle for bound sites, an empty one for visible sites, and a slashed circle for hidden sites.

<b>Id</b>	$i$	$::= Int^+$
<b>Name</b>	$m, s$	$::= String$
<b>State</b>	$t$	$::= \mathbf{b} \mid \mathbf{v} \mid \mathbf{h}$
<b>Sites</b>	$S$	$::= s \frown t \mid S, S \mid \epsilon$
<b>Molecules</b>	$M$	$::= \langle i : m \parallel S \rangle \mid M, M \mid \epsilon$
<b>Bonds</b>	$B$	$::= (i \frown s, i \frown s) \mid B, B \mid \epsilon$
<b>Molecular graph</b>	$G$	$::= (M, B)$

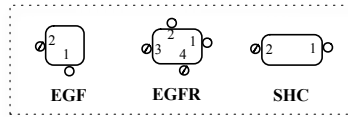
**Fig. 3.** The syntax of molecular graphs

A biochemical system is represented as a discrete system consisting of interacting components which give rise to structural and behavioral transformation of the components and of the system as a whole. Such a system is dynamic, has an emergent behavior, is highly concurrent and non-deterministic.

The interactions are abstracted as reaction patterns: each of the reaction patterns corresponds to a class of reactions characterized by common features. There are two types of protein interactions: (i) state changing by attaching or removing small phosphate groups at specific sites, and (ii) complexation by binding sites or decomplexation by breaking a bond between two proteins. An interaction may reveal new sites or hide some of the existing ones.

Given an initial set of biochemical species (as molecular graphs) and reactions described by a set of patterns (as molecular graph transformation rules), we automatically obtain a network by iteratively applying the reactions on the existing biochemical species until a termination condition is satisfied or until the exhaustive generation of all possible species is achieved.

*Example 1.* We illustrate a fragment of the EGFR signaling cascade already formalized in papers like [16, 24]. The protagonists of this model, depicted in Fig. 4, are: the epidermal growth factor *EGF* situated outside the cell acting as a ligand, the transmembrane protein *EGFR* with two extracellular sites and two intracellular sites as a receptor, and the adapter protein *SHC* inside the cell.



**Fig. 4.** The species occurring in the EGFR signaling cascade fragment

The signalling information is propagated from outside the cell to its interior following the reaction patterns depicted in Fig. 5. We note that, for the molecules participating in the reaction patterns, only the relevant sites are made precise, the rest of them remaining unchanged. The ligand *EGF* does not enter the cell, instead it transmits the information as a signal across the cell membrane. The function of information transfer across the membrane is performed by the receptor *EGFR*. One extracellular binding site of the receptor recognizes the signal protein in a dimeric form *EGF.EGF* produced by **(r1)** leading to the creation of a bond between the *EGF* dimer and the receptor, and to the activation by phosphorylation of a second extracellular site of the receptor **(r2)**. Two receptors activated by the same *EGF* dimer bind by **(r3)** creating an active dimer – a receptor tyrosine kinase (RTK). The RTK attaches phosphate groups to certain tyrosines, first on itself **(r4)**, and then on other proteins by activating them, as it is the case for the adapter *SHC* **(r5)**. Some of these latter proteins are also tyrosine kinases and, consequently, a cascade of phosphorylations occurs. We

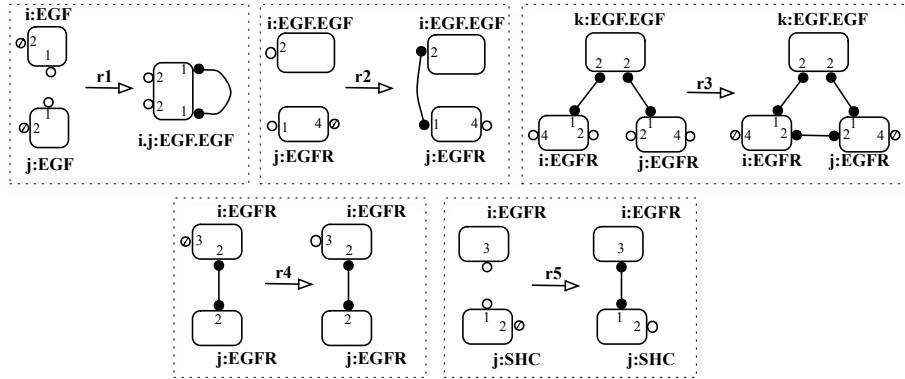


Fig. 5. The reaction patterns in the EGFR signaling cascade fragment

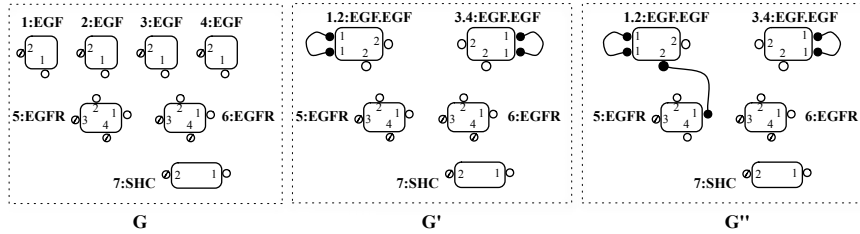


Fig. 6.  $\mathbf{G}$  is the initial configuration;  $\mathbf{G}'$  results from applying twice  $\mathbf{r1}$  on  $\mathbf{G}$ ;  $\mathbf{G}''$  results from an application of  $\mathbf{r2}$  on  $\mathbf{G}'$

illustrate in Fig. 6 an initial molecular graph  $\mathbf{G}$  we consider for this example, and two molecular graphs  $\mathbf{G}'$  and  $\mathbf{G}''$  resulting from it by applying some of the reactions above.

After the initiation of the signaling process, the signaling cascade must be terminated at a certain point since cells must be able to stop responding to a signal. If not, the signaling process may lead to uncontrolled cell growth as, for instance, in cancer. This process is ended by quickly engulfing and destroying the ligand-receptor complex by receptor-mediated endocytosis.

## 4 Term Rewriting Semantics for Molecular Graphs Transformations

We present in this section a signature for terms encoding molecular graphs, an encoding of reaction patterns as rewriting rules, and a rewriting relation encoding the application of reactions patterns on molecular graphs.



#### 4.1 Encoding the Molecular Graphs

We encode molecular graphs as terms using an order-sorted signature  $\Sigma = (\mathcal{S}, <, \mathcal{F})$  similarly to the term rewriting encoding of multigraphs with ports [3] up to sort renaming and extension of site information with a state. We give here only the operator set:

$$\begin{aligned}
\epsilon_X & : \longrightarrow XSet \\
-, - & : XSet \times XSet \longrightarrow XSet \ [ACU(\epsilon_X)] \\
- \hat{-} & : SiteName \times \{\mathbf{b}, \mathbf{v}, \mathbf{h}\} \rightarrow Site \\
\langle - : - \parallel - \rangle & : Id \times MoleculeName \times SiteSet \longrightarrow Molecule \\
(-, -) & : SiteName \times SiteName \longrightarrow Bond \\
- \hat{\wedge} - & : Id \times BondSet \longrightarrow Neighbour \\
- \simeq - & : Id \times NeighbourSet \longrightarrow AdjacencyEq \\
-(\downarrow) & : MoleculeSet \times AdjacencyEqSet \longrightarrow MGraph
\end{aligned}$$

where  $X$  takes sort values from the set  $\{Molecule, Bond, Neighbour, AdjacencyEq\}$ . The associative-commutative operator  $-, -$  (union) is overloaded on each of the set sorts, and  $\epsilon_X$  denotes the identity element (the empty set) for the operation  $-, -$ . We use  $\epsilon$  instead of  $\epsilon_X$  whenever the sort  $X$  can be easily deduced from the context. The subsorting relation  $X < XSet$  states that each term of sort  $X$  can be seen as a set with a single element.

Since the molecular graphs represent a particular case of multigraphs with ports and we encoded them in the oriented version, we impose an orientation of bonds by choosing an arbitrary order  $>$  on molecule names in the model and by considering that:

- a bond between two molecules of different name corresponds to a bond with the source labeled by the greater molecule name and the target labeled by the smaller molecule name, and
- a bond between two molecules of the same name is not oriented (or we can represent it equivalently as a bidirectional edge).

Sometimes biochemical reactions may occur without knowing all intermediate states or molecule names. This motivates us in using “unknown” names for molecules and sites as in [26] for beta-binders. We then choose to represent a bond between two molecules where at least one of them has an unknown name by an unoriented bond (or, equivalently, a bidirectional bond), and call this assumption *supra-estimation*. When a new molecule name is produced by a reaction, the total order on names is extended for including it. One simple way to handle this extension is to consider any new name to be smaller than all existing names.

Let  $\mathcal{X}$  be an  $(\mathcal{S}, <)$ -sorted family of variables. The axiomatizing involves theories with associativity, commutativity and unit element axioms (ACU) for the constructors  $-, -$  of each set-like sort with  $\epsilon$  the unit.

**Definition 2.** *We encode a molecular graph  $G = (M, B)$  as an algebraic term  $t_1(t_2)$  of sort  $MGraph$  where:*

- $t_1 \in \mathcal{T}_{\Sigma, MoleculeSet}(\mathcal{X})$  is the set of all molecules in  $G$ , and

- $t_2 \in \mathcal{T}_{\Sigma, \text{AdjacencyEqSet}}(\mathcal{X})$  is the set of adjacency equations providing the neighbours for each molecule in  $M$  (if any) together with the labels of the incident bonds.

Additionally, algebraic terms encoding molecular graphs must satisfy some structural properties in order to be considered *well-formed*:

- each molecule identifier occurs at most once: in the molecule set, in the adjacency equation set as left-hand side of an adjacency equation, in the neighbour set of a molecule identifier;
- each molecule identifier or site occurring in the adjacency equation set must also occur in the molecule set;
- each site has at most one incident bond and its state is correct w.r.t. the incidence information provided in the adjacency equation set;

and some *canonical form* requirements:

- right-hand sides of adjacency equations are non-empty sets of neighbours;
- only non-empty set of bonds occur in neighbour terms.

*Example 3.* The molecular graph occurring in the left-hand side of the reaction pattern **r3** illustrated in Fig.5 is encoded as the following term, assuming the total order on molecule names  $EGFR \succ EGF \succ SHC \succ EGF.EGF$ :

$$\begin{aligned}
 & \langle i : EGFR \parallel 1 \hat{\wedge} b, 2 \hat{\wedge} v, 4 \hat{\wedge} v \rangle, \langle j : EGFR \parallel 1 \hat{\wedge} b, 2 \hat{\wedge} v, 4 \hat{\wedge} v \rangle, \langle k : EGF.EGF \parallel 2 \hat{\wedge} b, 2 \hat{\wedge} b \rangle \\
 & (i \simeq k \frown (1, 2), j \simeq k \frown (1, 2))
 \end{aligned}$$

## 4.2 Encoding the Reaction Patterns

For all rewrite rules over  $\mathcal{T}_{\Sigma, MGraph}(\mathcal{X})$  we impose molecule identifiers occurring in the left-hand side to be variables. We say that a rewrite rule  $t_1 \rightarrow t_2$  over  $\mathcal{T}_{\Sigma, MGraph}(\mathcal{X})$  is well-formed (in canonical form) if both  $t_1$  and  $t_2$  are well-formed (in canonical form respectively). We call *bio-rewrite rule* a well-formed rewrite rule in canonical form.

**Definition 4.** We encode a molecular graph rewrite rule  $G_1 \rightsquigarrow G_2$  as a term rewrite rule  $t_1 \rightarrow t_2$  where  $t_i$  encodes  $G_i$ , for  $i = 1, 2$ .

The encoding of a molecular graph rewrite rule is a *bio-rewrite rule* since by definition the term encoding of a molecular graph is well-formed and in canonical form.

The *node-substitution* corresponding to a *bio-rewrite rule*, like for the encoding of transformations on multigraphs with ports, can be extracted automatically by means of a procedure `GetMap` analysing the identifier occurrences of both sides of a rule. We usually omit the identity mappings. We make explicit the deletion of a molecule (resp. site) in the node-substitution by mapping it to a new molecule (resp. site)  $\bullet$  called *black hole*.

*Example 5.* The molecular graph transformation **r1** illustrated in Fig. 5 has the following encoding as a rewrite rule with the corresponding node-substitution:

$$r_1 = (\langle i : EGF \parallel 1 \hat{v}, 2 \hat{h} \rangle, \langle j : EGF \parallel 1 \hat{v}, 2 \hat{h} \rangle)(\epsilon) \rightarrow \\ \langle i.j : EGF.EGF \parallel 1 \hat{b}, 1 \hat{b}, 2 \hat{v}, 2 \hat{v} \rangle (\langle i.j \simeq i.j \wedge (1, 1) \rangle) \\ \text{with } \xi_1 = \langle i : EGF \parallel 1 \hat{v}, 2 \hat{h} \rangle \mapsto \langle i.j : EGF.EGF \parallel 1 \hat{b}, 1 \hat{b}, 2 \hat{v}, 2 \hat{v} \rangle, \\ \langle j : EGF \parallel 1 \hat{v}, 2 \hat{h} \rangle \mapsto \langle i.j : EGF.EGF \parallel 1 \hat{b}, 1 \hat{b}, 2 \hat{v}, 2 \hat{v} \rangle.$$

### 4.3 Encoding the Reaction Pattern Application

The rewriting relation is the encoding of the multigraph transformation as presented in Sect. 2 with some particularities given by the instantiation of multigraphs with ports as molecular graphs.

The main ideas behind a rewriting step  $t \xrightarrow{\bar{r}} t'$  with  $r : t_1 \rightarrow t_2$  are the following:

- $\bar{r} : \overline{t_1} \rightarrow \overline{t_2}$  is the extended rule obtained from  $r$  by appending extension variables to each set-sorted subterm in the left-hand side and, accordingly, in the right-hand side.
- **GetMap** computes for a rewrite rule on *MGraph*-sorted terms the associated node-substitution using an analysis on the identifier occurrences.
- We use an ACU matching algorithm together with a matching rule which allows us to solve match equations where the left hand-side is a *NeighbourSet*- or *BondSet*-sorted term  $t_1, \dots, t_k$  (with  $k \geq 2$ ) not containing extension variables, and the right-hand side is a singleton,  $t$ . Such cases arise due to supra-estimated bonds. The problem is solved by matching only one term  $t_i$  from the left-hand side against the term  $t$  from the right-hand side, discarding the other terms  $t_j$ , with  $j \neq i$ , and matching the extension variables against the set unit  $\epsilon$ .
- We apply the substitution  $\sigma$  as usual. However, the result of the application of a substitution  $\{x \mapsto t\}$  with  $x, t : \text{MoleculeId}$  on adjacency equations requires a special treatment. For example, let us assume by supra-estimation that the nodes represented by  $x$  and  $y$  are connected by a bidirectional bond. If, after the application of a substitution  $\sigma$ , the name of  $\sigma(y)$  is greater than the name of  $\sigma(x)$ , then the bond becomes unidirectional with  $\sigma(x)$  the target and  $\sigma(y)$  the source of the bond.
- We apply the instantiated node-substitution  $\xi^\sigma$  on  $\overline{t_2}$  by sequentially applying each of its elementary node-substitution on the term. The application of a node-substitution on an adjacency equation (or a neighbour) transforms it in  $n$  adjacency equations (resp. neighbours), one for each corresponding node in the right-hand side of the mapping, and propagates the mapping application on the set of neighbours. After applying a node-substitution, the *MGraph*-terms may be neither well-formed, nor in canonical form.
- $\mathcal{S}$  is the rewrite system defining transformations on terms into well-formed *MGraph*-sorted terms in canonical form. Such rules (i) delete the adjacency equations for black holes and the black hole neighbours taking care of updating the states of incident sites, (ii) delete the extra-bonds (bonds whose

endpoints do not appear among the sites of the connected nodes), (iii) merge elements having the same identifier, and (iv) delete neighbours with empty bond sets and adjacency equations with an empty neighbour sets.

**Definition 6.** A term  $t$  of sort  $MGraph$  rewrites to a term  $t'$  using a *bio-rewrite rule*  $r : t_1 \rightarrow t_2$  where  $t_1, t_2 \in \mathcal{T}_{\Sigma, MGraph}(\mathcal{X})$  with  $\bar{r} : \overline{t_1} \rightarrow \overline{t_2}$  and  $\xi = \text{GetMap}(\bar{r})$ , if there exists a substitution  $\sigma$  solution of the ACU-matching problem  $\overline{t_1} \ll t$  such that  $t' = \xi^\sigma(\sigma(\overline{t_2})) \downarrow_{\mathcal{S}}$ . We call this relation *bio-rewriting*, we say that  $t$  *bio-rewrites to  $t'$  by  $r$*  and denote it by  $t \xrightarrow{\bar{r}} t'$ .

We do not illustrate here an example of a *bio-rewriting* step since it would be too technical (see [3] for such an example). However, it is important to emphasize that, according to the correspondence theorem in [3], the rewriting relation defined here on  $MGraph$ -terms encodes the transformation relation on molecular graphs.

#### 4.4 $\rho_{bio}$ -Calculus

Following closely the same steps for encoding the more general multigraphs with ports, we obtain for free a high-level calculus for molecular complexes which we call  $\rho_{bio}$ -calculus, similar to the  $\rho_{mg}$ -calculus. By extending the  $\lambda$ -calculus and term rewriting, the  $\rho_{bio}$ -calculus provides a way to combine rewriting and high-order functions. In addition, it fits in the direction given by Fontana and Buss [19] of using the  $\lambda$ -calculus as a formalism for studying biochemical systems.

## 5 Automatic Biochemical Network Generation

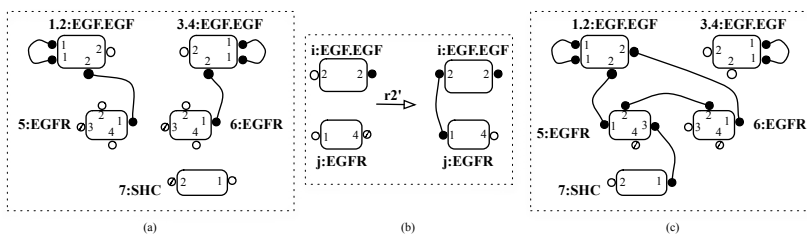
A biochemical system is not completely described by its components and the way they interact by means of reactions, but also by the behavior of the system as a whole. Modeling the automatic generation of a biochemical network amounts to defining how a set of reaction patterns is applied on a collection (set or multiset) of molecules. *Strategic rewriting* provides a formal model for defining the control on the reaction rule application.

One rough way of describing the behavior of the EGFR signalling cascade is to try to apply repeatedly all reaction patterns on the initial collection of molecules:  $\zeta_1 = \text{repeat}(\text{onestep}(r_1, r_2, r_3, r_4, r_5))$ . We can easily prove that this strategy terminates since the number of free binding sites decreases or remains constant with every rule application.

The execution can be separated in two stages: the *first* one,  $\{r_1, r_2, r_3\}$ , is concerned with the extracellular interactions between the signals and the receptors, while the *second* one,  $\{r_4, r_5\}$ , with the RTK cascade. Consequently, the control is specified as a repeated sequential composition of two strategies, one for each stage:

$$\zeta_2 = \text{seq}(\text{repeat}(\text{onestep}(r_1, r_2, r_3)), \text{repeat}(\text{onestep}(r_4, r_5))).$$

We require that all *EGF* dimers are created in a first step, and only after that



**Fig. 7.** (a) A possible graph obtained from applying  $r_2$  on  $G''$  from Fig. 6. (b) An additional rule which urges the binding of a dimeric form of *EGF* to a second receptor. (c) A molecular graph resulting from applying the strategy  $\zeta_5$  on the initial molecular graph from Fig. 6.

they bind to receptors by means of the strategy:

$$\zeta_3 = \text{seq}(\text{repeat}(r_1), \text{repeat}(r_2), \text{repeat}(r_3), \text{repeat}(\text{onestep}(r_4, r_5))),$$

or we require that an *EGF* dimeric form binds a receptor as soon as it is created by giving  $r_2$  priority over  $r_1$  using the **choice** strategy:

$$\zeta_4 = \text{seq}(\text{repeat}(\text{choice}(r_2, r_1)), \text{repeat}(r_3), \text{repeat}(\text{onestep}(r_4, r_5))).$$

In order to disallow the occurrence of molecular graphs like the one in Fig. 7(a) where each of the *EGF* dimeric forms binds one of the receptor and the binding of receptors is no longer possible, we introduce a new reaction rule  $r_2'$  depicted in Fig. 7(b). We impose a dimerization of the receptors as soon as possible by providing  $r_2'$  with a higher priority than  $r_2$ . We specify this using the strategy  $\text{choice}(r_2', r_2)$  which tries to apply  $r_2'$  first, and only if this does not succeed, it applies the rule  $r_2$ . Hence the *EGF* dimer binds two receptors as soon as it is created by means of the strategy:

$$\zeta_5 = \text{seq}(\text{repeat}(\text{choice}(r_2', r_2, r_1)), \text{repeat}(r_3), \text{repeat}(\text{onestep}(r_4, r_5))).$$

If we apply this last strategy on the initial molecular graph  $G$  illustrated in Fig. 6, the result of the *EGFR* signaling cascade fragment is represented by the molecular graph from Fig. 7(c) modulo a relabeling of the dimeric forms of *EGF*. Usually we obtain an exhaustive generation of the network by repeating the application of the control strategies, for instance  $\text{repeat}(\zeta_i)$ , for  $i \in \{2, \dots, 5\}$ .

The use of strategies at this stage provides our calculus with a declarative formalism for modeling a flexible control of rule applications. Different control strategies may be easily prototyped and experimented by selecting specific rules, grouping rules and applying each group with different strategies, and then linking them up. The advantages of using strategic rewriting for the automatic network generation have already been experimented in [8, 2] for a chemical application.

## 6 Related Work and Perspectives

Using graph rewriting and rule-based formalisms for modeling biological systems are already developed by several authors. A review of such rule-based formalisms can be found in [21] with emphasis on the capability of representing the topology

of complexes.  $\kappa$ -calculus [16] is a language of formal proteins which models complexes as graphs-with-sites and their interactions as a particular graph-rewriting operation. It is derived from process calculus, and bonds are represented in complexes by shared names. A step forward,  $\text{bio}\kappa$ -calculus [24] combines the  $\kappa$ -calculus with the brane calculi [9] providing a more expressive formalism by modeling the effects of protein interactions on the interaction capabilities of membranes. In this paper, we initiate the development of a similar approach based on a different calculus.

Pathway Logic [27] is a rewriting system formalism, where proteins and cells are modeled by algebraic terms, and reactions by rewrite rules. It is designed to work on two levels of abstraction, one concerning the protein states, and the other concerning the protein-protein interactions handled by means of graph rewriting. The *Maude* [15] system is used for implementing Pathway Logic, providing executability of the specifications and analytic tools. We use as well algebraic terms and rewrite rules for modeling molecules and reactions respectively as in Pathway Logic, with the differences that we represent graphs using adjacency lists, and the control of rewrite rule applications is expressed in the calculus.

An inspiring starting point for our work was the graphical formalism presented in [7] for modeling biochemical networks where the protein complexes are represented by typed attributed graphs, and classes of reactions are modeled by graph transformation rules. This model considers also quantitative information on reactions. Our approach, developed on a rule-based modeling framework and extensively using expressive graphical representations (like the one of [7]), focuses on providing in addition a strategic rewriting dimension. This new component allows a flexible modeling for automatic generation of biochemical networks.

The  $\rho_{\text{bio}}$ -calculus introduced in this paper for modeling the interactions between molecular complexes is an extension of the  $\rho_{\text{mg}}$ -calculus [3]. The connectivity information in molecular complexes is essential in the definition of some reaction patterns since it can directly affects their reactivity; this kind of information is easily modeled by graph transformation rules. It is also possible to express in the term encoding the absence of certain molecules or forbidden bonds in reaction patterns by means of recently developed anti-patterns [23], a concept already implemented in TOM. This allows us to model halting conditions as forbidden occurrences of some particular molecular graphs.

Based on our previous experience on implementing molecular chemistry [8, 2], we are currently implementing this calculus using TOM, a tool providing matching, normalization, and strategic rewriting in high-level programming languages like Java. An aspect we cannot neglect in modeling protein-protein interactions is the combinatorial complexity [21]. In order to obtain an efficient implementation, we make extensive use of the maximal term sharing provided by TOM.

Several further extensions of this work are possible.

We have conceived a qualitative model; by imposing kinetic information on reactions and on the initial set of biochemical species we should obtain biochemical networks allowing us to derive quantitative models as ODE-based models,

stoichiometric models, time depending. An inspiring rule-based formalism for a qualitative modeling of biochemical process is found in BIOCHAM [11].

We focus in this paper on interactions between proteins at the level of functional domains. However, the model we propose can be easily tuned to represent other types of biomolecular interactions, such as protein-DNA or protein-lipid interactions [18], or even membrane interactions [9, 10].

Strategic rewriting is a suitable formalism for modeling the highly concurrent and non-deterministic behavior of complex systems in general, and biological ones in particular. A first link between strategies and computation models inspired by biology is presented in [1] where control mechanisms in membrane systems are expressed by means of rewrite strategies.

Strategic rewriting opens many possibilities for modeling different aspects of biological systems at different abstraction levels. In particular, we plan to explore the capabilities of this formalism to model and reason about the adaptability and flexibility of the cell behavior. We do not limit our aim to modeling some well-known biological systems, but to help understand their behaviour and deduce new organization and behavioral principles. In the same vein as Păun in [25], we consider that reasoning at the level of strategies of computing (rewrite strategies), rather than at the tactic level (rewrite rules), is an incentive direction of formally studying biological systems.

## References

1. O. Andrei, G. Ciobanu, and D. Lucanu. Expressing Control Mechanisms in P systems by Rewriting Strategies. In H. Hoogeboom, G. Paun, G. Rozenberg, and A. Salomaa, eds., *Proc. of WMC7*, LNCS, vol. 4361, pag. 154–169. Springer, 2006.
2. O. Andrei, L. Ibanescu, and H. Kirchner. Non-intrusive Formal Methods and Strategic Rewriting for a Chemical Application. In K. Futatsugi, J.-P. Jouannaud, and J. Meseguer, eds., *Algebra, Meaning, and Computation, Essays Dedicated to Joseph A. Goguen*, LNCS, vol. 4060, pag. 194–215. Springer, 2006.
3. O. Andrei, H. Kirchner. A Rewriting Calculus for Multigraphs with Ports. In *Proc. of RULE'07 (to appear)*. (2007)
4. F. Baader and T. Nipkow. *Term Rewriting and All That*. Cambridge University Press, 1998.
5. E. Baland, P. Brauner, R. Kopetz, P.E. Moreau, A. Reilles. Tom: Piggybacking Rewriting on Java. In *Proc. of RTA'07 (to appear)*. (2007)
6. M. Bezem, J. Klop, and R. de Vrijer, eds. *Term Rewriting Systems.*, volume I of *Cambridge Tracts in Theoretical Computer Science*. Cambridge Univ. Press, 2003.
7. M. L. Blinov, J. Yang, J. R. Faeder, and W. S. Hlavacek. Graph theory for rule-based modeling of biochemical networks. In *Proc. of BioCONCUR 2005*.
8. O. Bournez, G.-M. Côme, V. Conraud, H. Kirchner, and L. Ibanescu. A Rule-Based Approach for Automated Generation of Kinetic Chemical Mechanisms. In R. Nieuwenhuis, ed., *Proc. of RTA'03*, LNCS, vol. 2706, pag. 30–45. Springer, 2003.
9. L. Cardelli. Brane Calculi. In V. Danos and V. Schächter, eds., *Proc. of CMSB 2004*, LNCS, vol. 3082, pag. 257–278. Springer, 2005.
10. L. Cardelli and S. Pradalier. Where Membranes Meet Complexes. In *Proc. of BioCONCUR 2005*.

11. N. Chabrier-Rivier, F. Fages, and S. Soliman. The Biochemical Abstract Machine BIOCHAM. In V. Danos and V. Schächter, eds., *Proc. of CMSB 2004*, LNCS, vol. 3082, pag. 172–191. Springer, 2005.
12. H. Cirstea and C. Kirchner. The rewriting calculus - Part I and II. *Logic Journal of the IGPL*, 9(3):427–498, 2001.
13. H. Cirstea, C. Kirchner, and L. Liquori. Matching Power. In A. Middeldorp, ed., *Proc. of RTA*, LNCS, vol. 2051, pag. 77–92. Springer, 2001.
14. H. Cirstea, C. Kirchner, L. Liquori, and B. Wack. Rewrite strategies in the rewriting calculus. *ENTCS*, 86(4):593–624, June 2003.
15. M. Clavel, F. Durán, S. Eker, P. Lincoln, N. Martí-Oliet, J. Meseguer, and J. F. Quesada. Maude: specification and programming in rewriting logic. *TCS*, 285(2):187–243, 2002.
16. V. Danos and C. Laneve. Formal molecular biology. *TCS*, 325(1):69–110, 2004.
17. S. Eker, N. Martí-Oliet, J. Meseguer, and A. Verdejo. Deduction, Strategies, and Rewriting. In M. Archer, T. Boy de la Tour, and C. A. Muñoz, eds., *Proc. of STRATEGIES'06*, ENTCS, 2007.
18. J. R. Faeder, M. L. Blinov, and W. S. Hlavacek. Graphical rule-based representation of signal-transduction networks. In H. Haddad, L. M. Liebrock, A. Omicini, and R. L. Wainwright, editors, *Proc. of SAC 2005*, pag. 133–140. ACM, 2005.
19. W. Fontana and L. W. Buss. The barrier of objects: From dynamical systems to bounded organizations. In J. Casti and A. Karlqvist, eds., *Boundaries and Barriers.*, pag. 56–116. Addison-Wesley, 1996.
20. J. A. Goguen and J. Meseguer. Order-Sorted Algebra I: Equational Deduction for Multiple Inheritance, Overloading, Exceptions and Partial Operations. *TCS*, 105(2):217–273, 1992.
21. W. S. Hlavacek, J. R. Faeder, M. L. Blinov, R. G. Posner, M. Hucka, and W. Fontana. Rules for Modeling Signal-Transduction Systems. *Science STKE* 334, 2006.
22. C. Kirchner, H. Kirchner, and M. Vittek. Designing Constraint Logic Programming Languages using Computational Systems. In P. Van Hentenryck and V. Saraswat, eds., *Principles and Practice of Constraint Programming. The Newport Papers.*, ch. 8, pag. 131–158. MIT Press, 1995.
23. C. Kirchner, R. Kopetz, and P.-E. Moreau. Anti-Pattern Matching. In *Proc. of ESOP 2007*, LNCS, vol. 4421, pag. 110–124. Springer, 2007.
24. C. Laneve and F. Tarissan. A simple calculus for proteins and cells. In *Proc. MeCBIC*, 2006.
25. Păun, G.: Twenty Six Research Topics About Spiking Neural P Systems. Available at the P Systems Web Page: <http://psystems.disco.unimib.it>. (2006)
26. C. Priami and P. Quaglia. Beta Binders for Biological Interactions. In V. Danos and V. Schächter, eds., *Proc. of CMSB'04*, LNCS, vol. 3028, pag. 20–33. Springer, 2004.
27. C. L. Talcott, S. Eker, M. Knapp, P. Lincoln, and K. Laderoute. Pathway Logic Modeling of Protein Functional Domains in Signal Transduction. In R. B. Altman, A. K. Dunker, L. Hunter, T. A. Jung, and T. E. Klein, eds., *Proc. of Biocomputing 2004*, pages 568–580. World Scientific, 2004.
28. E. Visser. Stratego: A Language for Program Transformation based on Rewriting Strategies. System Description of Stratego 0.5. In A. Middeldorp, editor, *Proc. of RTA '01*, LNCS, vol. 2051, pag. 357–361. Springer, 2001.