

# The Calculus of Looping Sequences for Modeling Biological Membranes

Roberto Barbuti<sup>1</sup>, Andrea Maggiolo-Schettini<sup>1</sup>,  
Paolo Milazzo<sup>1</sup>, and Angelo Troina<sup>2,3</sup>

<sup>1</sup> Dip. di Informatica, Università di Pisa, Largo B. Pontecorvo 3, 56127 - Pisa, Italy

<sup>2</sup> LIX - École Polytechnique Rue de Saclay, 91128 - Palaiseau, France

<sup>3</sup> LSV - ENS Cachan 61 Avenue du Président Wilson, 94235 - Cachan, France

**Abstract.** We survey the formalism Calculus of Looping Sequences (CLS) and a number of its variants from the point of view of their use for describing biological membranes. The formalism CLS is based on term rewriting and allows describing biomolecular systems. A first variant of CLS, called Stochastic CLS, extends the formalism with stochastic time, another variant, called LCLS (CLS with links), allows describing proteins interaction at the domain level. A third variant is introduced for easier description of biological membranes. This extension can be encoded into CLS as well as other formalisms capable of membrane description such as Brane Calculi and P Systems. Such encodings allow verifying and simulating descriptions in Brane Calculi and P Systems by means of verifiers and simulators developed for CLS.

## 1 Introduction

Cell biology, the study of the morphological and functional organization of cells, is now an established field in biochemical research. Computer Science can help the research in cell biology in several ways. For instance, it can provide biologists with models and formalisms capable of describing and analyzing complex systems such as cells. In the last few years many formalisms originally developed by computer scientists to model systems of interacting components have been applied to Biology. Among these, there are Petri Nets [16], Hybrid Systems [1], and the  $\pi$ -calculus [9, 25]. Moreover, new formalisms have been defined for describing biomolecular and membrane interactions [2, 7, 8, 11, 21, 23]. Others, such as P Systems [17, 18], have been proposed as biologically inspired computational models and have been later applied to the description of biological systems.

The  $\pi$ -calculus and new calculi based on it [21, 23] have been particularly successful in the description of biological systems, as they allow describing systems in a compositional manner. Interactions of biological components are modeled as communications on channels whose names can be passed; sharing names of private channels allows describing biological compartments. However, these calculi offer very low-level interaction primitives, and this causes models to become very large and difficult to read. Calculi such as those proposed in [7, 8, 11] give a more abstract description of systems and offer special biologically motivated

operators. However, they are often specialized to the description of some particular kinds of phenomena such as membrane interactions or protein interactions. Finally, P Systems have a simple notation and are not specialized to the description of a particular class of systems, but they are still not completely general. For instance, it is possible to describe biological membranes and the movement of molecules across membranes, and there are some variants able to describe also more complex membrane activities. However, the formalism is not so flexible to allow describing easily new activities observed on membranes without extending the formalism to model such activities.

Therefore, we conclude that there is a need for a formalism having a simple notation, having the ability of describing biological systems at different levels of abstraction, having some notions of compositionality and being flexible enough to allow describing new kinds of phenomena as they are discovered, without being specialized to the description of a particular class of systems. For this reason in [3] we have introduced the Calculus of Looping Sequences (CLS).

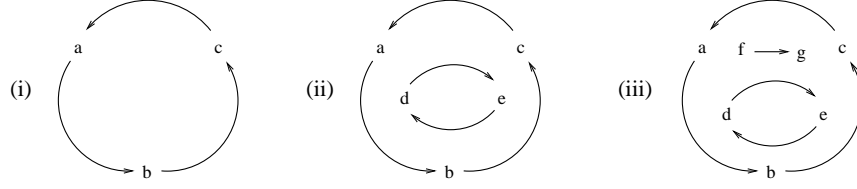
CLS is a formalism based on term rewriting with some features, such as a commutative parallel composition operator, and some semantic means, such as bisimulations, which are common in process calculi. This permits to combine the simplicity of notation of rewriting systems with the advantage of a form of compositionality. Actually, in [4] we have defined bisimilarity relations on CLS terms which are congruences with respect to the operators. The bisimilarity relation may be used to verify a property of a system by assessing its bisimilarity with a system one knows to enjoy that property. The fact that bisimilarity is a congruence is very important for a compositional account of behavioural equivalence.

In [5, 6], we have defined two extensions of CLS. The first, Stochastic CLS, allows describing quantitative aspects of the modeled systems such as the time spent by occurrences of chemical reactions. The second, CLS with links, allows describing protein interaction more precisely at a lower level of abstraction, namely at the domain level.

In this paper, after recalling CLS and the two mentioned extensions, we focus on the modeling of biological membranes by means of CLS. Now, CLS does not offer an easy representation for membranes whose nature is fluid and for proteins which consequently move freely on membrane surfaces. For this reason, in [15] we have defined a CLS variant, called CLS+, which introduces a new operator allowing commutativity on membrane surfaces. We show how CLS+ can be encoded into CLS.

In [3, 15] we have shown how Brane Calculi [7] and P Systems [18] can be translated into CLS. Here we recall the ideas on which the translations are based.

CLS appears to allow description and manipulation of biological membranes and, moreover, offers, via translations, verification and simulation tools to other formalisms for membrane description.



**Fig. 1.** (i) represents  $(a \cdot b \cdot c)^L$ ; (ii) represents  $(a \cdot b \cdot c)^L \mid (d \cdot e)^L$ ; (iii) represents  $(a \cdot b \cdot c)^L \mid ((d \cdot e)^L \mid f \cdot g)$ .

## 2 The Calculus of Looping Sequences (CLS)

In this section we recall the Calculus of Looping Sequences (CLS) and we give some guidelines for the modeling of biological systems. CLS is essentially based on term rewriting, hence a CLS model consists of a term and a set of rewrite rules. The term is intended to represent the structure of the modeled system, and the rewrite rules to represent the events that may cause the system to evolve.

### 2.1 Formal Definition

We start with defining the syntax of terms. We assume a possibly infinite alphabet  $\mathcal{E}$  of symbols ranged over by  $a, b, c, \dots$

**Definition 1 (Terms).** Terms  $T$  and sequences  $S$  of CLS are given by the following grammar:

$$\begin{aligned} T &::= S \mid (S)^L \mid T \mid T \mid T \\ S &::= \epsilon \mid a \mid S \cdot S \end{aligned}$$

where  $a$  is a generic element of  $\mathcal{E}$ , and  $\epsilon$  represents the empty sequence. We denote with  $\mathcal{T}$  the infinite set of terms, and with  $\mathcal{S}$  the infinite set of sequences.

In CLS we have a sequencing operator  $\cdot$ , a looping operator  $(-)^L$ , a parallel composition operator  $\mid$  and a containment operator  $\mid$ . Sequencing can be used to concatenate elements of the alphabet  $\mathcal{E}$ . The empty sequence  $\epsilon$  denotes the concatenation of zero symbols. A term can be either a sequence or a looping sequence (that is the application of the looping operator to a sequence) containing another term, or the parallel composition of two terms. By definition, looping and containment are always applied together, hence we can consider them as a single binary operator  $(-)^L \mid$  which applies to one sequence and one term.

Brackets can be used to indicate the order of application of the operators, and we assume  $(-)^L \mid$  to have precedence over  $\mid$ . In Figure 1 we show some examples of CLS terms and their visual representation.

In CLS we may have syntactically different terms representing the same structure. We introduce a structural congruence relation to identify such terms.

**Definition 2 (Structural Congruence).** *The structural congruence relations  $\equiv_S$  and  $\equiv_T$  are the least congruence relations on sequences and on terms, respectively, satisfying the following rules:*

$$\begin{aligned}
S_1 \cdot (S_2 \cdot S_3) &\equiv_S (S_1 \cdot S_2) \cdot S_3 & S \cdot \epsilon &\equiv_S \epsilon \cdot S \equiv_S S \\
S_1 &\equiv_S S_2 \text{ implies } S_1 \equiv_T S_2 \text{ and } (S_1)^L \rfloor T &\equiv_T (S_2)^L \rfloor T \\
T_1 \mid T_2 &\equiv_T T_2 \mid T_1 & T_1 \mid (T_2 \mid T_3) &\equiv_T (T_1 \mid T_2) \mid T_3 & T \mid \epsilon &\equiv_T T \\
(\epsilon)^L \rfloor \epsilon &\equiv_T \epsilon & (S_1 \cdot S_2)^L \rfloor T &\equiv_T (S_2 \cdot S_1)^L \rfloor T
\end{aligned}$$

Rules of the structural congruence state the associativity of  $\cdot$  and  $\mid$ , the commutativity of the latter and the neutral role of  $\epsilon$ . Moreover, axiom  $(S_1 \cdot S_2)^L \rfloor T \equiv_T (S_2 \cdot S_1)^L \rfloor T$  says that looping sequences can rotate. In the following, for simplicity, we will use  $\equiv$  in place of  $\equiv_T$ .

Rewrite rules will be defined essentially as pairs of terms, with the first term describing the portion of the system in which the event modeled by the rule may occur, and the second term describing how that portion of the system changes when the event occurs. In the terms of a rewrite rule we allow the use of variables. As a consequence, a rule will be applicable to all terms which can be obtained by properly instantiating its variables. Variables can be of three kinds: two of these are associated with the two different syntactic categories of terms and sequences, and one is associated with single alphabet elements. We assume a set of term variables  $TV$  ranged over by  $X, Y, Z, \dots$ , a set of sequence variables  $SV$  ranged over by  $\tilde{x}, \tilde{y}, \tilde{z}, \dots$ , and a set of element variables  $\mathcal{X}$  ranged over by  $x, y, z, \dots$ . All these sets are possibly infinite and pairwise disjoint. We denote by  $\mathcal{V}$  the set of all variables,  $\mathcal{V} = TV \cup SV \cup \mathcal{X}$ , and with  $\rho$  a generic variable of  $\mathcal{V}$ . Hence, a pattern is a term that may include variables.

**Definition 3 (Patterns).** *Patterns  $P$  and sequence patterns  $SP$  of CLS are given by the following grammar:*

$$\begin{aligned}
P &::= SP \mid (SP)^L \rfloor P \mid P \mid P \mid X \\
SP &::= \epsilon \mid a \mid SP \cdot SP \mid \tilde{x} \mid x
\end{aligned}$$

where  $a$  is a generic element of  $\mathcal{E}$ , and  $X, \tilde{x}$  and  $x$  are generic elements of  $TV, SV$  and  $\mathcal{X}$ , respectively. We denote with  $\mathcal{P}$  the infinite set of patterns.

We assume the structural congruence relation to be trivially extended to patterns. An *instantiation* is a partial function  $\sigma : \mathcal{V} \rightarrow \mathcal{T}$ . An instantiation must preserve the type of variables, thus for  $X \in TV, \tilde{x} \in SV$  and  $x \in \mathcal{X}$  we have  $\sigma(X) \in \mathcal{T}, \sigma(\tilde{x}) \in \mathcal{S}$  and  $\sigma(x) \in \mathcal{E}$ , respectively. Given  $P \in \mathcal{P}$ , with  $P\sigma$  we denote the term obtained by replacing each occurrence of each variable  $\rho \in \mathcal{V}$  appearing in  $P$  with the corresponding term  $\sigma(\rho)$ . With  $\Sigma$  we denote the set of all the possible instantiations and, given  $P \in \mathcal{P}$ , with  $Var(P)$  we denote the set of variables appearing in  $P$ . Now we define rewrite rules.

**Definition 4 (Rewrite Rules).** *A rewrite rule is a pair of patterns  $(P_1, P_2)$ , denoted with  $P_1 \mapsto P_2$ , where  $P_1, P_2 \in \mathcal{P}$ ,  $P_1 \neq \epsilon$  and such that  $Var(P_2) \subseteq Var(P_1)$ . We denote with  $\mathfrak{R}$  the infinite set of all the possible rewrite rules.*

Biomolecular Entity	CLS Term
Elementary object (genes, domains, other molecules, etc...)	Alphabet symbol
DNA strand	Sequence of elements repr. genes
RNA strand	Sequence of elements repr. transcribed genes
Protein	Sequence of elements repr. domains or single alphabet symbol
Molecular population	Parallel composition of molecules
Membrane	Looping sequence

**Table 1.** Guidelines for the abstraction of biomolecular entities into CLS.

A rewrite rule  $P_1 \mapsto P_2$  states that a term  $P_1\sigma$ , obtained by instantiating variables in  $P_1$  by some instantiation function  $\sigma$ , can be transformed into the term  $P_2\sigma$ . We define the semantics of CLS as a transition system, in which states correspond to terms, and transitions correspond to rule applications.

**Definition 5 (Semantics).** *Given a set of rewrite rules  $\mathcal{R} \subseteq \mathfrak{R}$ , the semantics of CLS is the least transition relation  $\rightarrow$  on terms closed under  $\equiv$ , and satisfying the following inference rules:*

$$\frac{P_1 \mapsto P_2 \in \mathcal{R} \quad P_1\sigma \not\equiv \epsilon \quad \sigma \in \Sigma}{P_1\sigma \rightarrow P_2\sigma} \quad \frac{T_1 \rightarrow T_2}{T \mid T_1 \rightarrow T \mid T_2} \quad \frac{T_1 \rightarrow T_2}{(S)^L \mid T_1 \rightarrow (S)^L \mid T_2}$$

where the symmetric rule for the parallel composition is omitted.

A *model* in CLS is given by a term describing the initial state of the system and by a set of rewrite rules describing all the events that may occur.

## 2.2 Modeling Guidelines

We describe how CLS can be used to model biomolecular systems analogously to what done by Regev and Shapiro in [24] for the  $\pi$ -calculus. An abstraction is a mapping from a real-world domain to a mathematical domain, which may allow highlighting some essential properties of a system while ignoring other, complicating, ones. In [24], Regev and Shapiro show how to abstract biomolecular systems as concurrent computations by identifying the biomolecular entities and events of interest and by associating them with concepts of concurrent computations such as concurrent processes and communications. In particular, they give some guidelines for the abstraction of biomolecular systems to the  $\pi$ -calculus, and give some simple examples.

The use of rewrite systems, such as CLS, to describe biological systems is founded on a different abstraction. Usually, entities (and their structures) are abstracted by terms of the rewrite system, and events by rewriting rules. We have already introduced the biological interpretation of CLS operators in the previous section. Here we want to give more general guidelines.

Biomolecular Event	Examples of CLS Rewrite Rule
State change	$a \mapsto b$ $\tilde{x} \cdot a \cdot \tilde{y} \mapsto \tilde{x} \cdot b \cdot \tilde{y}$
Complexation	$a   b \mapsto c$ $\tilde{x} \cdot a \cdot \tilde{y}   b \mapsto \tilde{x} \cdot c \cdot \tilde{y}$
Decomplexation	$c \mapsto a   b$ $\tilde{x} \cdot c \cdot \tilde{y} \mapsto \tilde{x} \cdot a \cdot \tilde{y}   b$
Catalysis	$c   P_1 \mapsto c   P_2$ where $P_1 \mapsto P_2$ is the catalyzed event
State change on membrane	$(a \cdot \tilde{x})^L \rfloor X \mapsto (b \cdot \tilde{x})^L \rfloor X$
Complexation on membrane	$(a \cdot \tilde{x} \cdot b \cdot \tilde{y})^L \rfloor X \mapsto (c \cdot \tilde{x} \cdot \tilde{y})^L \rfloor X$ $a   (b \cdot \tilde{x})^L \rfloor X \mapsto (c \cdot \tilde{x})^L \rfloor X$ $(b \cdot \tilde{x})^L \rfloor (a   X) \mapsto (c \cdot \tilde{x})^L \rfloor X$
Decomplexation on membrane	$(c \cdot \tilde{x})^L \rfloor X \mapsto (a \cdot b \cdot \tilde{x})^L \rfloor X$ $(c \cdot \tilde{x})^L \rfloor X \mapsto a   (b \cdot \tilde{x})^L \rfloor X$ $(c \cdot \tilde{x})^L \rfloor X \mapsto (b \cdot \tilde{x})^L \rfloor (a   X)$
Catalysis on membrane	$(c \cdot \tilde{x} \cdot SP_1 \cdot \tilde{y})^L \rfloor X \mapsto (c \cdot \tilde{x} \cdot SP_2 \cdot \tilde{y})^L \rfloor X$ where $SP_1 \mapsto SP_2$ is the catalyzed event
Membrane crossing	$a   (\tilde{x})^L \rfloor X \mapsto (\tilde{x})^L \rfloor (a   X)$ $(\tilde{x})^L \rfloor (a   X) \mapsto a   (\tilde{x})^L \rfloor X$ $\tilde{x} \cdot a \cdot \tilde{y}   (\tilde{z})^L \rfloor X \mapsto (\tilde{z})^L \rfloor (\tilde{x} \cdot a \cdot \tilde{y}   X)$ $(\tilde{z})^L \rfloor (\tilde{x} \cdot a \cdot \tilde{y}   X) \mapsto \tilde{x} \cdot a \cdot \tilde{y}   (\tilde{z})^L \rfloor X$
Catalyzed membrane crossing	$a   (b \cdot \tilde{x})^L \rfloor X \mapsto (b \cdot \tilde{x})^L \rfloor (a   X)$ $(b \cdot \tilde{x})^L \rfloor (a   X) \mapsto a   (b \cdot \tilde{x})^L \rfloor X$ $\tilde{x} \cdot a \cdot \tilde{y}   (b \cdot \tilde{z})^L \rfloor X \mapsto (b \cdot \tilde{z})^L \rfloor (\tilde{x} \cdot a \cdot \tilde{y}   X)$ $(b \cdot \tilde{z})^L \rfloor (\tilde{x} \cdot a \cdot \tilde{y}   X) \mapsto \tilde{x} \cdot a \cdot \tilde{y}   (b \cdot \tilde{z})^L \rfloor X$
Membrane joining	$(\tilde{x})^L \rfloor (a   X) \mapsto (a \cdot \tilde{x})^L \rfloor X$ $(\tilde{x})^L \rfloor (\tilde{y} \cdot a \cdot \tilde{z}   X) \mapsto (\tilde{y} \cdot a \cdot \tilde{z} \cdot \tilde{x})^L \rfloor X$
Catalyzed membrane joining	$(b \cdot \tilde{x})^L \rfloor (a   X) \mapsto (a \cdot b \cdot \tilde{x})^L \rfloor X$ $(\tilde{x})^L \rfloor (a   b   X) \mapsto (a \cdot \tilde{x})^L \rfloor (b   X)$ $(b \cdot \tilde{x})^L \rfloor (\tilde{y} \cdot a \cdot \tilde{z}   X) \mapsto (\tilde{y} \cdot a \cdot \tilde{z} \cdot \tilde{x})^L \rfloor X$ $(\tilde{x})^L \rfloor (\tilde{y} \cdot a \cdot \tilde{z}   b   X) \mapsto (\tilde{y} \cdot a \cdot \tilde{z} \cdot \tilde{x})^L \rfloor (b   X)$
Membrane fusion	$(\tilde{x})^L \rfloor (X)   (\tilde{y})^L \rfloor (Y) \mapsto (\tilde{x} \cdot \tilde{y})^L \rfloor (X   Y)$
Catalyzed membrane fusion	$(a \cdot \tilde{x})^L \rfloor (X)   (b \cdot \tilde{y})^L \rfloor (Y) \mapsto$ $(a \cdot \tilde{x} \cdot b \cdot \tilde{y})^L \rfloor (X   Y)$
Membrane division	$(\tilde{x} \cdot \tilde{y})^L \rfloor (X   Y) \mapsto (\tilde{x})^L \rfloor (X)   (\tilde{y})^L \rfloor (Y)$
Catalyzed membrane division	$(a \cdot \tilde{x} \cdot b \cdot \tilde{y})^L \rfloor (X   Y) \mapsto$ $(a \cdot \tilde{x})^L \rfloor (X)   (b \cdot \tilde{y})^L \rfloor (Y)$

**Table 2.** Guidelines for the abstraction of biomolecular events into CLS.

First of all, we select the biomolecular entities of interest. Since we want to describe cells, we consider molecular populations and membranes. Molecular populations are groups of molecules that are in the same compartment of the cell. Molecules can be of many types: we classify them as DNA and RNA strands, proteins, and other molecules. DNA and RNA strands and proteins can be seen as non-elementary objects. DNA strands are composed by genes, RNA strands are composed by parts corresponding to the transcription of individual genes, and proteins are composed by parts having the role of interaction sites (or domains). Other molecules are considered as elementary objects, even if they are complexes. Membranes are considered as elementary objects, in the sense that we do not describe them at the level of the lipids they are made of. The only interesting properties of a membrane are that it may contain something (hence, create a compartment) and that it may have molecules on its surface.

Now, we select the biomolecular events of interest. The simplest kind of event is the change of state of an elementary object. Then, we may have interactions between molecules: in particular complexation, decomplexation and catalysis. These interactions may involve single elements of non-elementary molecules (DNA and RNA strands, and proteins). Moreover, we may have interactions between membranes and molecules: in particular a molecule may cross or join a membrane. Finally, we may have interactions between membranes: in this case there may be many kinds of interactions (fusion, division, etc. . .).

The guidelines for the abstraction of biomolecular entities and events into CLS are given in Table 1 and Table 2, respectively. Entities are associated with CLS terms: elementary objects are modeled as alphabet symbols, non-elementary objects as CLS sequences and membranes as looping sequences. Biomolecular events are associated with CLS rewrite rules. In the figure we give some examples of rewrite rules for each type of event. The list of examples is not complete: one could imagine also rewrite rules for the description of complexation/decomplexation events involving more than two molecules, or catalysis events in which the catalyzing molecule is on a membrane and the catalyzed event occurs in its content, or more complex interactions between membranes. We remark that in the second example of rewrite rule associated with the complexation event we have that one of the two molecules which are involved should be either an elementary object or a protein modeled as a single alphabet symbol. As before, this is caused by the problem of modeling protein interaction at the domain level. This problem is solved by the extension of CLS with links, called LCLS, we shall describe in the following.

### 2.3 Examples

A well-known example of biomolecular system is the epidermal growth factor (EGF) signal transduction pathway[26, 19]. If EGF proteins are present in the environment of a cell, they should be interpreted as a proliferation signal from the environment, and hence the cell should react by synthesizing proteins which stimulate its proliferation. A cell recognizes the EGF signal because it has on its

membrane some EGF receptor proteins (EGFR), which are transmembrane proteins (they have some intra-cellular and some extra-cellular domains). One of the extra-cellular domains binds to one EGF protein in the environment, forming a signal-receptor complex on the membrane. This causes a conformational change on the receptor protein that enables it to bind to another one signal-receptor complex. The formation of the binding of the two signal-receptor complexes (called dimerization) causes the phosphorylation of some intra-cellular domains of the dimer. This, in turn, causes the internal domains of the dimer to be recognized by a protein that is inside the cell (in the cytoplasm), called SHC. The protein SHC binds to the dimer, enabling a long chain of protein-protein interactions, which finally activate some proteins, such as one called ERK, which bind to the DNA and stimulate synthesis of proteins for cell proliferation.

Now, we use CLS to build a model of the first steps of the EGF signaling pathway up to the binding of the signal-receptor dimer to the SHC protein. In the following we shall refine the model by using the LCLS extension to describe interactions at the domain level.

We model the EGFR, EGF and SHC proteins as the alphabet symbols  $EGFR$ ,  $EGF$  and  $SHC$ , respectively. The cell is modeled as a looping sequence (representing its external membrane), initially composed only by  $EGFR$  symbols, containing  $SHC$  symbols and surrounded by  $EGF$  symbols. The rewrite rules modeling the first steps of the pathway are the following:

$$EGF \mid (EGFR \cdot \tilde{x})^L \mid X \mapsto (CMPLX \cdot \tilde{x})^L \mid X \quad (R1)$$

$$(CMPLX \cdot \tilde{x} \cdot CMPLX \cdot \tilde{y})^L \mid X \mapsto (DIM \cdot \tilde{x} \cdot \tilde{y})^L \mid X \quad (R2)$$

$$(DIM \cdot \tilde{x})^L \mid X \mapsto (DIMp \cdot \tilde{x})^L \mid X \quad (R3)$$

$$(DIMp \cdot \tilde{x})^L \mid (SHC \mid X) \mapsto (DIMpSHC \cdot \tilde{x})^L \mid X \quad (R4)$$

Rule R1 describes the binding of a EGF protein to a EGFR receptor protein on the membrane surface. The result of the binding is a signal-receptor complex denoted  $CMPLX$ . Rule R2 describes the dimerization of two signal-receptor complex, the result is denoted  $DIM$ . Rule R3 describes the phosphorylation (and activation) of a signal-receptor dimer, that is the replacement of a  $DIM$  symbol with a  $DIMp$  symbol. Finally, rule R4 describes the binding of an active dimer  $DIMp$  with a SHC protein contained in the cytoplasm. The result is a  $DIMpSHC$  symbol placed on the membrane surface.

A possible initial term for the model in this example is given by a looping sequence composed by some  $EGFR$  symbols, containing some  $SHC$  symbols and with some  $EGF$  symbols outside. A possible evolution of such a term by means of application of the given rewrite rules is the following (we write on each transition the name of the rewrite rule applied):

$$\begin{aligned} & EGF \mid EGF \mid (EGFR \cdot EGFR \cdot EGFR \cdot EGFR)^L \mid (SHC \mid SHC) \\ \xrightarrow{(R1)} & EGF \mid (EGFR \cdot CMPLX \cdot EGFR \cdot EGFR)^L \mid (SHC \mid SHC) \\ \xrightarrow{(R1)} & (EGFR \cdot CMPLX \cdot EGFR \cdot CMPLX)^L \mid (SHC \mid SHC) \end{aligned}$$



$$\xrightarrow{(R2)} (EGFR \cdot DIM \cdot EGFR)^L \mid (SHC \mid SHC)$$

$$\xrightarrow{(R3)} (EGFR \cdot DIM_p \cdot EGFR)^L \mid (SHC \mid SHC)$$

$$\xrightarrow{(R4)} (EGFR \cdot DIM_p SHC \cdot EGFR)^L \mid SHC$$

We show another example of modeling of a biomolecular system with CLS, that is the modeling of a simple gene regulation process. This kind of processes are essential for cell life as they allow a cell to regulate the production of proteins that may have important roles for instance in metabolism, growth, proliferation and differentiation.

The example we consider is as follows: we have a simple DNA fragment consisting of a sequence of three genes. The first, denoted  $p$ , is called *promoter* and is the place where a *RNA polymerase* enzyme (responsible for translation of DNA into RNA) binds to the DNA. The second, denoted  $o$ , is called *operator* and it is the place where a *repressor* protein (responsible for regulating the activity of the RNA polymerase) binds to the DNA. The third, denoted as  $g$ , is the gene that encodes for the protein whose production is regulated by this process.

When the repressor is not bound to the DNA, the RNA polymerase can scan the sequence of genes and transcribe gene  $g$  into a piece of RNA that will be later translated into the protein encoded by  $g$ . When the repressor is bound to the DNA, it becomes an obstacle for the RNA polymerase that cannot scan any more the sequence of genes.

The CLS model of this simple regulation process is as follows. The sequence of genes is represented as the CLS sequence  $p \cdot o \cdot g$ , the RNA polymerase enzyme as *polym*, the repressor protein as *repr*, and the piece of RNA obtained by the translation of gene  $g$  as *rna*. The rewrite rules describing the process are the following:

$$polym \mid p \cdot \tilde{x} \mapsto pp \cdot \tilde{x} \tag{R1}$$

$$repr \mid \tilde{x} \cdot o \cdot \tilde{y} \mapsto \tilde{x} \cdot ro \cdot \tilde{y} \tag{R2}$$

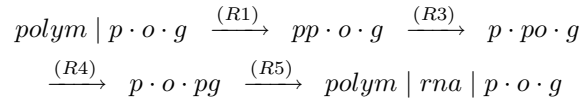
$$pp \cdot o \cdot \tilde{x} \mapsto p \cdot po \cdot \tilde{x} \tag{R3}$$

$$\tilde{x} \cdot po \cdot g \mapsto \tilde{x} \cdot o \cdot pg \tag{R4}$$

$$\tilde{x} \cdot pg \mapsto polym \mid rna \mid \tilde{x} \cdot g \tag{R5}$$

Rules R1 and R2 describe the binding of the RNA polymerase and of the repressor to the corresponding genes in the DNA sequences. The results of these bindings are that the symbols representing the two genes are replaced by  $pp$  and  $ro$ , respectively. Rules R3, R4 and R5 describe the activity of the RNA polymerase enzyme in the absence of the repressor: it moves from gene  $p$  to gene  $o$  in rule R3, then it moves from gene  $o$  to gene  $g$  in rule R4, and finally it produces the RNA fragment and leaves the DNA in rule R5. Note that, in order to apply rule R3, the repressor must be not bound to the DNA.

The only possible evolution of a term representing an initial situation in which no repressors are present is



that represent the case in which the RNA polymerase enzyme can scan the DNA sequence and transcribe gene  $g$  into a piece of RNA. When the repressor is present, instead, a possible evolution is



and it corresponds to a situation in which the repressor stops the transcription of the gene by hampering the activity of the RNA polymerase.

### 3 Two Extensions of CLS

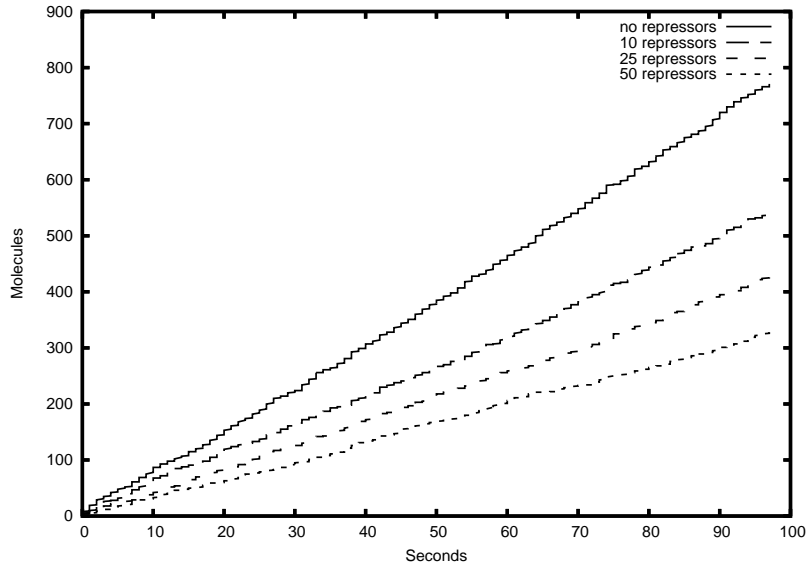
In this section we describe two extensions of CLS. The first, Stochastic CLS, allows describing quantitative aspects of the modeled systems, such as the time spent by occurrences of chemical reactions. The second, CLS with links, allows describing protein interaction more precisely at a lower level of abstraction, namely at the domain level.

#### 3.1 Stochastic CLS

In CLS only qualitative aspects of biological systems are considered, such as their structure and the presence (or the absence) of certain molecules. As a consequence, on CLS models it is only possible to verify properties such as the reachability of particular states or causality relationships between events. It would be interesting to verify also properties such as the time spent to reach a particular state, or the probability of reaching it. To face this problem, in [6] we have developed a stochastic extension of CLS, called *Stochastic CLS*, in which quantitative aspects, such as time and probability are taken into account.

The standard way of extending a formalism to model quantitative aspects of biological systems is by incorporating the stochastic framework developed by Gillespie with its simulation algorithm for chemical reactions [12] in the semantics of the formalism. This has been done, for instance, for the  $\pi$ -calculus [20, 22]. The idea of Gillespie's algorithm is that a rate constant is associated with each chemical reaction that may occur in the system. Such a constant is obtained by multiplying the kinetic constant of the reaction by the number of possible combinations of reactants that may occur in the system. The resulting rate constant is then used as the parameter of an exponential distribution modeling the time spent between two occurrences of the considered chemical reaction.

The use of exponential distributions to represent the (stochastic) time spent between two occurrences of chemical reactions allows describing the system as a

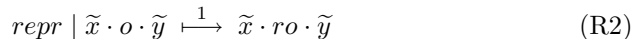
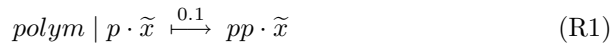


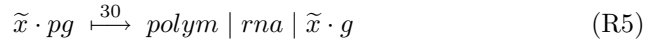
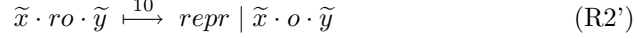
**Fig. 2.** Simulation result of the regulation process: number of RNA molecules over time.

Continuous Time Markov Chain (CTMC), and consequently it allows verifying properties of the described system by means of analytic means and by means of stochastic model checkers.

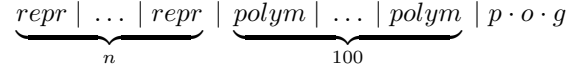
In Stochastic CLS, incorporating Gillespie’s stochastic framework is not a simple exercise. The main difficulty is counting the number of possible reactant combinations of the chemical reaction described by a rewrite rule. This means counting the number of different positions where the rewrite rule can be applied, by taking into account that rules may contain variables. We have defined the Stochastic CLS in [6], and showed how to derive a CTMC from the semantics of a system modeled in Stochastic CLS. This allows performing simulation and verification of properties of the described systems, for instance by using stochastic model checkers, such as PRISM [13].

Let us consider the simple regulation process we modeled with CLS in Section 2.3. We now extend the CLS model by including a kinetic constant in each rewrite rule. The result is a Stochastic CLS model. In order to make the model a little more realistic we add two rewrite rules describing the unbinding of the RNA polymerase and of the repressor from the DNA. Hence, the rewrite rules of the Stochastic CLS model are the following:





We developed a simulator based on Stochastic CLS, and we used it to study the behaviour of the regulation process. In particular, we performed simulations by varying the quantity of repressors and we observed the production of RNA fragments in each case. The initial configuration of the system is given by the following term



and we performed simulations with  $n = 0, 10, 25$  and  $50$ . The results of the simulations are shown in Figure 2. By varying the number of repressors from 0 to 50 the rate of transcription of the DNA into RNA molecules decreases.

### 3.2 CLS with Links (LCLS)

A formalism for modeling proteins interactions at the domain level was developed in the seminal paper by Danos and Laneve [11], and extended in [14]. This formalism allows expressing proteins by a node with a fixed number of domains; binding between domains allow complexating proteins. In this section we extend CLS to represent proteins interaction at the domain level. Such an extension, called Calculus of Linked Looping Sequences (LCLS), is obtained by labelling elementary components of sequences. Two elements with the same label are considered to be linked.

To model a protein at the domain level in CLS it would be natural to use a sequence with one symbol for each domain. However, the binding between two domains of two different proteins, that is the linking between two elements of two different sequences, cannot be expressed in CLS. To represent this, we extend CLS by labels on basic symbols. If in a term two symbols appear having the same label, we intend that they represent domains which are bound to each other. If in a term there is a symbol with a label and no other symbol with the same label, we intend that the term represents only a part of a system we model, and that the symbol will be linked to some other symbol in another part of the term representing the full model.

As membranes create compartments, elements inside a looping sequence cannot be linked to elements outside. Elements inside a membrane can be linked either to other elements inside the membrane or to elements of the membrane itself. An element can be linked at most to another element.

As an example, we model in LCLS the first steps of the EGF pathway described before. We model the EGFR protein as the sequence  $R_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2}$  where  $R_{E1}$  and  $R_{E2}$  are two extra-cellular domains and  $R_{I1}$  and  $R_{I2}$  are two

intra-cellular domains. The membrane of the cell is modeled as a looping sequence which could contain EGFR proteins. Outside the looping sequence (i.e. in the environment) there could be EGF proteins, and inside (i.e. in the cytoplasm) there could be SHC proteins. Rewrite rules modeling the pathway are the following:

$$EGF \mid (R_{E1} \cdot \tilde{x})^L \mid X \mapsto (sR_{E1} \cdot \tilde{x})^L \mid X \quad (R1)$$

$$\begin{aligned} & (sR_{E1} \cdot R_{E2} \cdot x \cdot y \cdot \tilde{x} \cdot sR_{E1} \cdot R_{E2} \cdot z \cdot w \cdot \tilde{y})^L \mid X \mapsto \\ & (sR_{E1} \cdot R_{E2}^1 \cdot x \cdot y \cdot sR_{E1} \cdot R_{E2}^1 \cdot z \cdot w \cdot \tilde{x} \cdot \tilde{y})^L \mid X \end{aligned} \quad (R2)$$

$$\begin{aligned} & (R_{E2}^1 \cdot R_{I1} \cdot \tilde{x} \cdot R_{E2}^1 \cdot R_{I1} \cdot \tilde{y})^L \mid X \mapsto \\ & (R_{E2}^1 \cdot PR_{I1} \cdot \tilde{x} \cdot R_{E2}^1 \cdot R_{I1} \cdot \tilde{y})^L \mid X \end{aligned} \quad (R3)$$

$$\begin{aligned} & (R_{E2}^1 \cdot PR_{I1} \cdot \tilde{x} \cdot R_{E2}^1 \cdot R_{I1} \cdot \tilde{y})^L \mid X \mapsto \\ & (R_{E2}^1 \cdot PR_{I1} \cdot \tilde{x} \cdot R_{E2}^1 \cdot PR_{I1} \cdot \tilde{y})^L \mid X \end{aligned} \quad (R4)$$

$$\begin{aligned} & (R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot \tilde{x} \cdot R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot \tilde{y})^L \mid (SHC \mid X) \mapsto \\ & (R_{E2}^1 \cdot PR_{I1} \cdot R_{I2}^2 \cdot \tilde{x} \cdot R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot \tilde{y})^L \mid (SHC^2 \mid X) \end{aligned} \quad (R5)$$

Rule R1 represents the binding of the EGF protein to the receptor domain  $R_{E1}$  with  $sR_{E1}$  as a result. Rule R2 represents that when two EGFR proteins activated by proteins EGF occur on the membrane, they may bind to each other to form a dimer (shown by the link 1). Rule R3 represents the phosphorylation of one of the internal domains  $R_{I1}$  of the dimer, and rule R4 represents the phosphorylation of the other internal domain  $R_{I1}$  of the dimer. The result of each phosphorylation is  $pR_{I1}$ . Rule R5 represents the binding of the protein SHC in the cytoplasm to an internal domain  $R_{I2}$  of the dimer. Remark that the binding of SHC to the dimer is represented by the link 2, allowing the protein SHC to continue the interactions to stimulate cell proliferation.

Let us denote the sequence  $R_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2}$  by EGFR. By starting from a cell with some EGFR proteins on its membrane, some SHC proteins in the cytoplasm and some EGF proteins in the environment, a possible evolution is the following:

$$\begin{aligned} & EGF \mid EGF \mid (EGFR \cdot EGFR \cdot EGFR)^L \mid (SHC \mid SHC) \\ & \xrightarrow{(R1)} EGF \mid (sR_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2} \cdot EGFR \cdot EGFR)^L \mid (SHC \mid SHC) \\ & \xrightarrow{(R1)} (sR_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2} \cdot EGFR \cdot sR_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2})^L \mid (SHC \mid SHC) \\ & \xrightarrow{(R2)} (sR_{E1} \cdot R_{E2}^1 \cdot R_{I1} \cdot R_{I2} \cdot sR_{E1} \cdot R_{E2}^1 \cdot R_{I1} \cdot R_{I2} \cdot EGFR)^L \mid (SHC \mid SHC) \\ & \xrightarrow{(R3)} (sR_{E1} \cdot R_{E2}^1 \cdot pR_{I1} \cdot R_{I2} \cdot sR_{E1} \cdot R_{E2}^1 \cdot R_{I1} \cdot R_{I2} \cdot EGFR)^L \mid (SHC \mid SHC) \end{aligned}$$

$$\begin{aligned} & \xrightarrow{(R4)} (sR_{E1} \cdot R_{E2}^1 \cdot pR_{I1} \cdot R_{I2} \cdot sR_{E1} \cdot R_{E2}^1 \cdot pR_{I1} \cdot R_{I2} \cdot EGFR)^L \mid (SHC \mid SHC) \\ & \xrightarrow{(R5)} (sR_{E1} \cdot R_{E2}^1 \cdot pR_{I1} \cdot R_{I2}^2 \cdot sR_{E1} \cdot R_{E2}^1 \cdot pR_{I1} \cdot R_{I2} \cdot EGFR)^L \mid (SHC^2 \mid SHC) \end{aligned}$$

## 4 CLS and Membranes

What could seem strange in CLS is the use of looping sequences for the description of membranes, as sequencing is not a commutative operation and this do not correspond to the usual fluid representation of membranes in which objects can move freely. What one would expect is to have a multiset or a parallel composition of objects on a membrane. In the case of CLS, what could be used is a parallel composition of sequences. To address this problem, we define an extension of CLS, called CLS+, in which the looping operator can be applied to a parallel composition of sequences, and we show that we can translate quite easily CLS+ models into CLS ones.

### 4.1 Definition of CLS+

Terms in CLS+ are defined as follows.

**Definition 6 (Terms).** Terms  $T$ , branes  $B$ , and sequences  $S$  of CLS+ are given by the following grammar:

$$\begin{aligned} T & ::= S \mid (B)^L \mid T \mid T \\ B & ::= S \mid B \mid B \\ S & ::= \epsilon \mid a \mid S \cdot S \end{aligned}$$

where  $a$  is a generic element of  $\mathcal{E}$ . We denote with  $\mathcal{T}$  the infinite set of terms, with  $\mathcal{B}$  the infinite set of branes, and with  $\mathcal{S}$  the infinite set of sequences.

The structural congruence relation of CLS+ is a trivial extension of the one of CLS. The only difference is that commutativity of branes replaces rotation of looping sequences.

**Definition 7 (Structural Congruence).** The structural congruence relations  $\equiv_S$ ,  $\equiv_B$  and  $\equiv_T$  are the least congruence relations on sequences, on branes and on terms, respectively, satisfying the following rules:

$$\begin{aligned} S_1 \cdot (S_2 \cdot S_3) & \equiv_S (S_1 \cdot S_2) \cdot S_3 & S \cdot \epsilon & \equiv_S \epsilon \cdot S \equiv_S S \\ S_1 & \equiv_S S_2 \text{ implies } S_1 & \equiv_B S_2 \\ B_1 \mid B_2 & \equiv_B B_2 \mid B_1 & B_1 \mid (B_2 \mid B_3) & \equiv_B (B_1 \mid B_2) \mid B_3 & B \mid \epsilon & \equiv_B B \\ S_1 & \equiv_S S_2 \text{ implies } S_1 & \equiv_T S_2 \\ B_1 & \equiv_B B_2 \text{ implies } (B_1)^L \mid T & \equiv_T (B_2)^L \mid T \\ T_1 \mid T_2 & \equiv_T T_2 \mid T_1 & T_1 \mid (T_2 \mid T_3) & \equiv_T (T_1 \mid T_2) \mid T_3 & T \mid \epsilon & \equiv_T T & (\epsilon)^L \mid \epsilon & \equiv \epsilon \end{aligned}$$

Now, to define patterns in CLS+ we consider an additional type of variables with respect of CLS, namely brane variables. We assume a set of brane variables  $BV$  ranged over by  $\bar{x}, \bar{y}, \bar{z}, \dots$

**Definition 8 (Patterns).** Patterns  $P$ , brane patterns  $BP$  and sequence patterns  $SP$  of CLS+ are given by the following grammar:

$$\begin{aligned} P & ::= SP \mid (BP)^L \mid P \mid P \mid X \\ BP & ::= SP \mid BP \mid BP \mid \bar{x} \\ SP & ::= \epsilon \mid a \mid SP \cdot SP \mid \tilde{x} \mid x \end{aligned}$$

where  $a$  is a generic element of  $\mathcal{E}$ , and  $X, \bar{x}, \tilde{x}$  and  $x$  are generic elements of  $TV, BV, SV$  and  $\mathcal{X}$ , respectively. We denote with  $\mathcal{P}$  the infinite set of patterns.

As usual, rewrite rules are pairs of patterns.

**Definition 9 (Rewrite Rules).** A rewrite rule is a pair of patterns  $(P_1, P_2)$ , denoted with  $P_1 \mapsto P_2$ , where  $P_1, P_2 \in PP$ ,  $P_1 \neq \epsilon$  and such that  $\text{Var}(P_2) \subseteq \text{Var}(P_1)$ . We denote with  $\mathfrak{R}$  the infinite set of all the possible rewrite rules.

Now, differently from CLS, we have that a rule such as  $a \mid b \mapsto c$  could be applied to elements of a looping sequence. For instance,  $a \mid b \mapsto c$  can be applied to the term  $(a \mid b)^L \mid d$  so to obtain the term  $(c)^L \mid d$ . However, a rule such as  $(a)^L \mid b \mapsto c$  still cannot be applied to elements of a looping sequences, as  $((a)^L \mid b)^L \mid c$  is not a CLS+ term.

The rules that can be applied to elements of a looping sequence are those having the form  $(B_1, B_2)$  with  $B_1, B_2 \in \mathcal{B}$ . We call these rules *brane rules* and we denote as  $\mathfrak{R}_{\mathcal{B}} \subset \mathfrak{R}$  their infinite set. Now, in the semantics of CLS+ we have to take into account brane rules and allow them to be applied also to elements of looping sequences. Hence, we define the semantics as follows.

**Definition 10 (Semantics).** Given a set of rewrite rules  $\mathcal{R} \subseteq \mathfrak{R}$ , and a set of brane rules  $\mathcal{R}_{\mathcal{B}} \subseteq \mathfrak{R}$ , such that  $(\mathcal{R} \setminus \mathcal{R}_{\mathcal{B}}) \cap \mathfrak{R}_{\mathcal{B}} = \emptyset$ , the semantics of CLS is the least transition relation  $\rightarrow$  on terms closed under  $\equiv$ , and satisfying the following inference rules:

$$\begin{aligned} & \frac{(P_1, P_2) \in \mathcal{R} \quad P_1 \sigma \neq \epsilon \quad \sigma \in \Sigma}{P_1 \sigma \rightarrow P_2 \sigma} \quad \frac{T_1 \rightarrow T_2}{T \mid T_1 \rightarrow T \mid T_2} \quad \frac{T_1 \rightarrow T_2}{(B)^L \mid T_1 \rightarrow (B)^L \mid T_2} \\ & \frac{(BP_1, BP_2) \in \mathcal{R}_{\mathcal{B}} \quad BP_1 \sigma \neq \epsilon \quad \sigma \in \Sigma}{BP_1 \sigma \rightarrow_{\mathcal{B}} BP_2 \sigma} \\ & \frac{B_1 \rightarrow_{\mathcal{B}} B_2}{B \mid B_1 \rightarrow_{\mathcal{B}} B \mid B_2} \quad \frac{B_1 \rightarrow_{\mathcal{B}} B_2}{(B_1)^L \mid T \rightarrow (B_2)^L \mid T} \end{aligned}$$

where  $\rightarrow_{\mathcal{B}}$  is a transition relation on branes, and where the symmetric rules for the parallel composition of terms and of branes are omitted.

In the definition of the semantics of CLS+ we use an additional transition relation  $\rightarrow_{\mathcal{B}}$  on branes. This relation is used to describe the application of a brane rule to elements of a looping sequence. As usual, a CLS+ model is composed by a term, representing the initial state of the modeled system, and a set of rewrite rules.

In the following section we show that CLS+ models can be translated into CLS models. The translation into CLS is compositional and preserves the semantics of the model.

## 4.2 Translating CLS+ into CLS

The first step of the translation of a CLS+ model into CLS is a preprocessing procedure. For each brane rule  $(BP_1, BP_2)$  in the CLS+ model, we add to the set of rules of the model a new rule, namely  $((BP_1 \mid \bar{x})^L \mid X, (BP_2 \mid \bar{x})^L \mid X)$ . This new rule is redundant in the model, as every time it can be applied to a CLS+ term, also the original one can be applied with the same result. However, the translation we are going to define will translate the original rule into a CLS rule that will be applicable only inside looping sequences, or at the top level of the term, and will translate the new rule into a CLS rule applicable only to elements that compose a looping sequence.

Now, the translation of CLS+ into CLS consists mainly of an encoding function, denoted  $\{\{\cdot\}\}$ , which maps CLS+ patterns into CLS patterns. This encoding function will be used to translate each rewrite rule of the CLS+ model into a rewrite rule for the corresponding CLS model, and to translate the term representing the initial state of the system in the CLS+ model into a CLS term for the corresponding CLS model.

The encoding function for CLS+ patterns is defined as follows. We assume a total and injective function from brane variables into a subset of term variables that are never used in CLS models. More easily, we assume brane variables to be a subset of the term variables of CLS. Moreover, we assume *in* and *out* to be symbols of the alphabet  $\mathcal{E}$  never used in CLS models.

The encoding follows the “ball-bearing” technique described by Cardelli in [7]. Intuitively, every CLS+ looping sequence is translated into a couple of CLS looping sequences, one contained in the other, with the brane patterns of the CLS+ looping sequence between the two corresponding CLS looping sequences.

**Definition 11 (Encoding Function).** *The encoding function  $\{\{\cdot\}\}$  maps CLS+ patterns into CLS patterns, and is given by the following recursive definition:*

$$\begin{aligned} \{\{SP\}\} &= SP \\ \{\{X\}\} &= X \\ \{\{(BP)^L \mid P\}\} &= (out)^L \mid (BP \mid (in)^L \mid \{\{P\}\}) \\ \{\{P_1 \mid P_2\}\} &= \{\{P_1\}\} \mid \{\{P_2\}\} \end{aligned}$$



A CLS rewrite rule is obtained from each CLS+ rewrite rule of the translated model by applying the encoding function to the two patterns of the rule. More precisely, given a CLS+ rule  $P_1 \mapsto P_2$ , the corresponding CLS rule is  $(in)^L \rfloor (\{\{P_1\}\} \mid X) \mapsto (in)^L \rfloor (\{\{P_2\}\} \mid X)$  where  $X$  is a term variable that does not occur in  $P_1$  and  $P_2$ . For example, by applying the encoding to the two patterns of the CLS+ rewrite rule

$$R = b \cdot x \mid c \mapsto b \cdot x$$

we obtain

$$R_{\{\cdot\}} = (in)^L \rfloor (b \cdot x \mid c \mid X) \mapsto (in)^L \rfloor (b \cdot x \mid X).$$

The encoding of a CLS+ term into a CLS term is as follows: given a CLS+ term  $T$  the corresponding CLS term is  $(in)^L \rfloor \{\{T\}\}$ . In this case we have that the encoding function never encounters variables. Consider, as an example, the following CLS+ term:

$$T = a \mid (c \mid d \mid b \cdot b \mid d)^L \rfloor d$$

the corresponding CLS term is as follows:

$$T_{\{\cdot\}} = (in)^L \rfloor (a \mid (out)^L \rfloor (c \mid d \mid b \cdot b \mid d \mid (in)^L \rfloor d))$$

Now, it is easy to see that  $R$  can be applied to  $T$ , because parallel components in the looping sequence can be commuted, and the result of the application is

$$T' = a \mid (b \cdot b \mid d \mid d)^L \rfloor d$$

but the corresponding CLS rewrite rule  $R_{\{\cdot\}}$  cannot be applied to  $T_{\{\cdot\}}$ . However, we have that  $R \in \mathcal{R}_{\mathcal{B}}$ , and hence, by the preprocessing phase described above, we have that also

$$R' = (b \cdot x \mid c \mid \bar{x})^L \rfloor X \mapsto (b \cdot x \mid \bar{x})^L \rfloor X$$

is a rule of the CLS+ model. By translating rule  $R'$  we obtain

$$R'_{\{\cdot\}} = (in)^L \rfloor ((out)^L \rfloor (b \cdot x \mid c \mid \bar{x} \mid (in)^L \rfloor X) \mid Y) \mapsto (in)^L \rfloor ((out)^L \rfloor (b \cdot x \mid \bar{x} \mid (in)^L \rfloor X) \mid Y)$$

that can be applied to  $T_{\{\cdot\}}$ . The result of the application is

$$(in)^L \rfloor (a \mid (out)^L \rfloor (b \cdot b \mid d \mid d \mid (in)^L \rfloor d))$$

that corresponds exactly to the encoding of  $T'$ .

### 4.3 CLS, Brane Calculi and P Systems

Brane Calculi are a family of process calculi specialized in the description of membrane activity, and they allow associating processes with membranes of a membrane structure. Each process is composed by actions whose execution has an effect on the membrane structure. Some examples of actions are phagocytosis (a membrane engulfs another one), exocytosis (a membrane expels another one), and pinocytosis (a new membrane is created inside another one). These three actions are enough to define the simplest of Brane Calculi, namely the PEP calculus. Other actions, such as fusion of membranes and mitosis can be used to define different calculi of the family. Moreover, extensions of Brane Calculi allow describing interactions with molecules and complexes, such as letting them enter and exit membranes.

We have given a sound and complete encoding of the PEP calculus into CLS in [3, 15]. Here, to recall shortly the encoding technique, we give a very simple example of PEP system and we show its translation into CLS. The PEP system we consider is the following

$$\phi(| \diamond |) \circ \phi^\perp(\mathbf{0})(| \diamond |)$$

representing two adjacent membranes  $\phi_n(| \diamond |)$  and  $\phi_n^\perp(\mathbf{0})(| \diamond |)$  ( $\circ$  denotes juxtaposition) both containing nothing of relevant (what is between brackets  $(| \diamond |)$  is the content of the membrane and  $\diamond$  is the null system). The processes associated with the two membranes are  $\phi$  and  $\phi^\perp(\mathbf{0})$ , respectively, representing two complementary phagocytosis actions: the first says that the membrane it is associated with can be engulfed by another membrane, and the second that the membrane it is associated with can engulf another membrane, that will be surrounded by another new membrane whose associated process is the parameter of the action (in this case it is the idle process  $\mathbf{0}$ ). Hence, in accordance with the semantics of the PEP calculus, we have that the only transition that can be performed by the system is the following, leading to a system that is equivalent to the null system  $\diamond$ :

$$\phi(| \diamond |) \circ \phi^\perp(\mathbf{0})(| \diamond |) \rightarrow \mathbf{0}(| \mathbf{0}(| \mathbf{0}(| \diamond |) |) |) \equiv \diamond$$

By applying the encoding to the system we obtain the following CLS term  $T$ :

$$act \cdot circ \cdot e \cdot brane \cdot b \cdot \phi \cdot a \cdot \mathbf{0} \cdot a \cdot b \cdot \mathbf{0} \cdot e \cdot brane \cdot d \cdot \phi^\perp \cdot c \cdot \mathbf{0} \cdot c \cdot \mathbf{0} \cdot c \cdot d \cdot \mathbf{0}$$

where *act* is a sort of program counter that precedes the symbol representing the next action to be executed, symbol *circ* represents  $\circ$ , symbol *brane* represents a membrane  $(| \diamond |)$ , symbols  $\phi$  and  $\phi^\perp$  represent the corresponding actions, symbol  $\mathbf{0}$  represents the idle process and symbols *a*, *b*, *c*, *d* and *e* are used as separators of actions and parameters. The translation consists also of a set of CLS rewrite rules to be applied to terms obtained by the encoding of PEP systems. Such a set of rewrite rules does not depend on the encoded PEP system, hence it is always the same. By applying rewrite rules, the long sequence obtained from the encoding is transformed into a hierarchy of looping sequences corresponding to the membrane hierarchy in the original PEP system, then rewrite rules are

applied that correspond to the semantics of the actions occurring in the processes associated with membranes.

Hence, by means of application of rewrite rules, the result of the encoding of the PEP system may evolve as follows (where  $\rightarrow^*$  represents a sequence of rewrite rule applications):

$$\begin{aligned}
T &\rightarrow act \cdot brane \cdot b \cdot \phi \cdot a \cdot \mathbf{0} \cdot a \cdot b \cdot \mathbf{0} \mid act \cdot brane \cdot d \cdot \phi^\perp \cdot c \cdot \mathbf{0} \cdot c \cdot \mathbf{0} \cdot c \cdot d \\
&\rightarrow^* (act \cdot \phi \cdot a \cdot \mathbf{0} \cdot a)^L \mid act \cdot \mathbf{0} \mid (act \cdot \phi^\perp \cdot c \cdot \mathbf{0} \cdot c \cdot \mathbf{0} \cdot c)^L \mid act \cdot \mathbf{0} \\
&\rightarrow (act \cdot \mathbf{0})^L \mid (act \cdot \mathbf{0} \mid (act \cdot \mathbf{0})^L \mid (act \cdot \mathbf{0})^L \mid act \cdot \mathbf{0}) \\
&\rightarrow^* act \cdot \mathbf{0}
\end{aligned}$$

Differently from Brane Calculi, P Systems (in their most common formulation) do not allow describing complex membrane activities such as phagocytosis and exocytosis. However, they are specialized in the description of reactions between molecules which are placed in a compartment of a complex membrane structure.

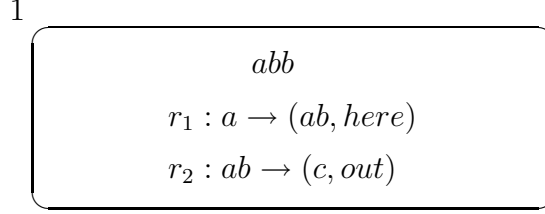
A P System is a membrane structure (a nesting of membranes) in which there could be multisets of objects representing molecules. A set of multiset rewrite rules is associated with each membrane, and describe the reactions that may occur between the molecules contained in the membrane. The result of the application of a rewrite rule can either remain in the same membrane, or exit the membrane, or enter an inner membrane. Priorities can be imposed on rewrite rules, meaning that some rules can be applied only if some others cannot, and it is possible for a membrane to dissolve and release its content to in the environment.

A peculiarity of P Systems is that rewrite rules are applied in a fully-parallel manner, namely in one step of evolution of the system all rules are applied as many times as possible (to different molecules), and this is one of the main differences with respect to CLS in which at each step one only rewrite rule is applied. We show that P Systems can be translated into CLS, and that the execution of a (fully parallel) step of a P System is simulated by a sequence of steps in CLS. A variant of P Systems, called Sequential P Systems, in which rules are applied sequentially is described in [10]. We do not consider the translation of this variant into CLS as it would be quite trivial and of little interest.

To recall the encoding technique, we give a simple example of P System and we show its translation into CLS. We focus on the translation of multiple parallelism, hence we consider a P System (depicted in Figure 3) consisting of a single membrane with only two rules, without priorities and without membrane dissolutions. We give a simplified translation: more details can be found in [15].

The alphabet of objects in the considered P System is  $\{a, b, c\}$ . A multiset of objects from this alphabet is represented by a CLS term as follows: let  $n_a, n_b$  and  $n_c$  be the number of occurrences of  $a, b$  and  $c$  in the multiset, respectively, then

$$a \cdot \overbrace{1 \cdot \dots \cdot 1}^{n_a} \mid b \cdot \overbrace{1 \cdot \dots \cdot 1}^{n_b} \mid c \cdot \overbrace{1 \cdot \dots \cdot 1}^{n_c}$$



**Fig. 3.** A simple example of P System.

is the term representing the multiset. We choose this representation as it allows us checking whether an object is absent, by checking whether the corresponding symbol if followed by zero 1s. An empty multiset is represented as  $a \mid b \mid c$ .

The CLS term obtained by the translation of the considered P System is the following:

$$(1)^L \mid (Check \mid a \cdot 1 \mid b \cdot 1 \cdot 1 \mid c \mid r_1 \mid r_2 \mid (next)^L \mid (a \mid b \mid c))$$

where the membrane of the P System is represented by a looping sequence composed by the membrane label (in this case 1). Inside the looping sequence there is a *Check* symbol representing the current state of the system, the translation of the multiset of objects of the membrane, two symbols  $r_1$  and  $r_2$  corresponding to the evolutionary rules of the membrane, and an empty multiset surrounded by a looping sequence *next*. This empty multiset is used to store temporary information on the result of the application of evolutionary rules.

The CLS rewrite rules obtained by the encoding of the considered P System are the following:

$$(1)^L \mid (X \mid Check \mid a \cdot 1 \cdot \tilde{x} \mid r_1) \mapsto (1)^L \mid (X \mid Check \mid a \cdot 1 \cdot \tilde{x} \mid r_1 \cdot 1) \quad (C1)$$

$$(1)^L \mid (X \mid Check \mid a \mid r_1) \mapsto (1)^L \mid (X \mid Check \mid a \mid r_1 \cdot 0) \quad (C2)$$

$$(1)^L \mid (X \mid Check \mid a \cdot 1 \cdot \tilde{x} \mid b \cdot 1 \cdot \tilde{y} \mid r_1 \cdot z \mid r_2) \mapsto \\ (1)^L \mid (X \mid Check \mid a \cdot 1 \cdot \tilde{x} \mid b \cdot 1 \cdot \tilde{y} \mid r_1 \cdot z \mid r_2 \cdot 1) \quad (C3)$$

$$(1)^L \mid (X \mid Check \mid a \mid r_1 \cdot z \mid r_2) \mapsto (1)^L \mid (X \mid Run \mid a \mid r_1 \cdot z \mid r_2 \cdot 0) \quad (C4)$$

$$(1)^L \mid (X \mid Check \mid b \mid r_1 \cdot z \mid r_2) \mapsto (1)^L \mid (X \mid Run \mid b \mid r_1 \cdot z \mid r_2 \cdot 0) \quad (C5)$$

$$(1)^L \mid (X \mid Run \mid a \cdot 1 \cdot \tilde{x} \mid r_1 \cdot 1 \mid (next)^L \mid (Y \mid a \cdot \tilde{y} \mid b \cdot \tilde{z})) \mapsto \\ (1)^L \mid (X \mid Run \mid a \cdot \tilde{x} \mid r_1 \cdot 1 \mid (next)^L \mid (Y \mid a \cdot 1 \cdot \tilde{y} \mid b \cdot 1 \cdot \tilde{z})) \quad (R1)$$

$$\begin{aligned}
(1)^L \rfloor (X \mid Run \mid a \cdot 1 \cdot \tilde{x} \mid b \cdot 1 \cdot \tilde{y} \mid r_2 \cdot 1 \mid (next)^L \rfloor (Y \mid c \cdot \tilde{z})) &\mapsto \\
(1)^L \rfloor (X \mid Run \mid a \cdot \tilde{x} \mid b \cdot \tilde{y} \mid r_2 \cdot 1 \mid (next)^L \rfloor (Y \mid c \cdot 1 \cdot \tilde{z})) &\quad (R2)
\end{aligned}$$

$$(1)^L \rfloor (X \mid Run \mid a \mid r_1 \cdot 1) \mapsto (1)^L \rfloor (X \mid Run \mid a \mid r_1 \cdot 0) \quad (R3)$$

$$(1)^L \rfloor (X \mid Run \mid a \mid r_2 \cdot 1) \mapsto (1)^L \rfloor (X \mid Run \mid a \mid r_2 \cdot 0) \quad (R4)$$

$$(1)^L \rfloor (X \mid Run \mid b \mid r_2 \cdot 1) \mapsto (1)^L \rfloor (X \mid Run \mid b \mid r_2 \cdot 0) \quad (R5)$$

$$(1)^L \rfloor (X \mid Run \mid r_1 \cdot 0 \mid r_2 \cdot 0) \mapsto (1)^L \rfloor (X \mid Update \mid r_1 \cdot 0 \mid r_2 \cdot 0) \quad (R6)$$

$$\begin{aligned}
(1)^L \rfloor (X \mid Update \mid x \cdot \tilde{x} \mid (next)^L \rfloor (Y \mid x \cdot 1 \cdot \tilde{y})) &\mapsto \\
(1)^L \rfloor (X \mid Update \mid x \cdot 1 \cdot \tilde{y} \cdot \tilde{x} \mid (next)^L \rfloor (Y \mid x)) &\quad (U1)
\end{aligned}$$

$$\begin{aligned}
(1)^L \rfloor (X \mid Update \mid (next)^L \rfloor (a \mid b \mid c)) &\mapsto \\
(1)^L \rfloor (X \mid Check \mid (next)^L \rfloor (a \mid b \mid c)) &\quad (U2)
\end{aligned}$$

Rules (C1)–(C5) describe the steps performed by the system while it is in *Check* state: the objective of this phase is to test whether each evolutionary rule is applicable or not. When all rules have been tested, the system moves into a state called *Run*, whose steps are given by the application of rules (R1)–(R6). In this second phase, evolutionary rules previously identified as applicable are actually applied, and the result of the application is stored inside the looping sequence *next*. Finally, when no evolutionary rule is further applicable, the system moves into a state called *Update*, in which the content of the looping sequence *next* is used to reset the multiset of objects of the membrane by applying rule (U1)–(U2). When this update operation has been performed, the system moves back to the *Check* state.

## 5 Conclusions

We have surveyed the formalism CLS and a number of its variants from the point of view of its use for describing biological membranes. Verification and simulation tools have been developed for CLS and its variants and can be used to study properties of membrane systems. Via translations, these tools can be used to study systems described by other formalisms such as Brane Calculi and P Systems, capable of describing biological membranes.

## References

1. R. Alur, C. Belta, F. Ivancic, V. Kumar, M. Mintz, G.J. Pappas, H. Rubin and J. Schug. “Hybrid Modeling and Simulation of Biomolecular Networks”. Hybrid Systems: Computation and Control, LNCS 2034, pages 19–32, Springer, 2001.

2. R. Barbuti, S. Cataudella, A. Maggiolo-Schettini, P. Milazzo and A. Troina. “A Probabilistic Model for Molecular Systems”. *Fundamenta Informaticae*, volume 67, pages 13–27, 2005.
3. R. Barbuti, A. Maggiolo-Schettini, P. Milazzo and A. Troina. “A Calculus of Looping Sequences for Modelling Microbiological Systems”. *Fundamenta Informaticae*, volume 72, pages 21–35, 2006.
4. R. Barbuti, A. Maggiolo-Schettini, P. Milazzo, and A. Troina. “Bisimulation Congruences in the Calculus of Looping Sequences”. *Int. Colloquium on Theoretical Aspects of Computing (ICTAC’06)*, LNCS 4281, pages 93–107, Springer, 2006.
5. R. Barbuti, A. Maggiolo-Schettini, and P. Milazzo. “Extending the Calculus of Looping Sequences to Model Protein Interaction at the Domain Level”. *Int. Symposium on Bioinformatics Research and Applications (ISBRA’07)*, LNBI 4463, pages 638–649, Springer, 2006.
6. R. Barbuti, A. Maggiolo-Schettini, P. Milazzo, P. Tiberi and A. Troina. “Stochastic CLS for the Modeling and Simulation of Biological Systems”. Submitted for publication. Draft available at: <http://www.di.unipi.it/~milazzo/>.
7. L. Cardelli. “Brane Calculi. Interactions of Biological Membranes”. *CMSB’04*, LNCS 3082, pages 257–280, Springer, 2005.
8. N. Chabrier-Rivier, M. Chiaverini, V. Danos, F. Fages and V. Schachter. “Modeling and Querying Biomolecular Interaction Networks”. *Theoretical Computer Science*, volume 325, number 1, pages 25–44, 2004.
9. M. Curti, P. Degano, C. Priami and C.T. Baldari. “Modelling Biochemical Pathways through Enhanced pi-calculus”. *Theoretical Computer Science*, volume 325, number 1, pages 111–140, 2004.
10. Z. Dang and O.H. Ibarra. “On P Systems Operating in Sequential and Limited Parallel Modes”, *Workshop on Descriptive Complexity of Formal Systems*, pages 164–177, 2004.
11. V. Danos and C. Laneve. “Formal Molecular Biology”. *Theoretical Computer Science*, volume 325, number 1, pages 69–110, 2004.
12. D. Gillespie. “Exact Stochastic Simulation of Coupled Chemical Reactions”. *Journal of Physical Chemistry*, volume 81, pages 2340–2361, 1977.
13. M. Kwiatkowska, G. Norman, and D. Parker. “Probabilistic Symbolic Model Checking with PRISM: a Hybrid Approach”. *Int. Journal on Software Tools for Technology Transfer*, volume 6, number 2, pages 128–142, 2004.
14. C. Laneve and F. Tarissan. “A Simple Calculus for Proteins and Cells”. *Workshop on Membrane Computing and Biological Inspired Process Calculi (MeCBIC’06)*, to appear in *ENTCS*.
15. P. Milazzo. “Qualitative and Quantitative Formal Modeling of Biological Systems”. PhD Thesis, Università di Pisa, 2007.
16. H. Matsuno, A. Doi, M. Nagasaki and S. Miyano. “Hybrid Petri Net Representation of Gene Regulatory Network”. *Pacific Symposium on Biocomputing*, World Scientific Press, pages 341–352, 2000.
17. G. Păun. “Computing with Membranes”. *Journal of Computer and System Sciences*, volume 61, number 1, pages 108–143, 2000.
18. G. Păun. “Membrane Computing. An Introduction”. Springer, 2002.
19. M.J. Pérez-Jiménez and F.J. Romero-Campero. “A Study of the Robustness of the EGFR Signalling Cascade Using Continuous Membrane Systems”. *IWINAC’05*, LNCS 3561, pages 268–278, Springer, 2005.
20. C. Priami. “Stochastic  $\pi$ -Calculus”. *The Computer Journal*, volume 38, number 7, pages 578–589, 1995.

21. C. Priami and P. Quaglia “Beta Binders for Biological Interactions”. CMSB’04, LNCS 3082, pages 20–33, Springer, 2005.
22. C. Priami, A. Regev, W. Silvermann, and E. Shapiro. “Application of a Stochastic Name–Passing Calculus to Representation and Simulation of Molecular Processes”. *Information Processing Letters*, volume 80, pages 25–31, 2001.
23. A. Regev, E.M. Panina, W. Silverman, L. Cardelli and E. Shapiro. “BioAmbients: An Abstraction for Biological Compartments”. *Theoretical Computer Science*, volume 325, number 1, pages 141–167, 2004.
24. A. Regev and E. Shapiro. “The  $\pi$ –Calculus as an Abstraction for Biomolecular Systems”. *Modelling in Molecular Biology*, pages 219–266, Natural Computing Series, Springer, 2004.
25. A. Regev, W. Silverman and E.Y. Shapiro. “Representation and Simulation of Biochemical Processes Using the pi-calculus Process Algebra”. *Pacific Symposium on Biocomputing*, World Scientific Press, pages 459–470, 2001.
26. H.S. Wiley, S.Y. Shvartsman and D.A. Lauffenburger. “Computational Modeling of the EGF–Receptor System: a Paradigm for Systems Biology”. *Trends in Cell Biology*, volume 13, number 1, pages 43–50, Elsevier, 2003.