# Qualitative and Quantitative Formal Modeling of Biological Systems

Paolo Milazzo

Dipartimento di Informatica, Università di Pisa, Italy

Camerino – April, 2007

# Outline of the talk

# Cells: complex systems of interactive components



- Two classifications of cell:
  - procaryotic
  - eucaryotic
- Main actors:
  - membranes
  - proteins
  - DNA/RNA strands
- Interaction networks:
  - metabolic pathways
  - signaling pathways
  - gene regulatory networks

Computer Science can provide biologists with formalisms for the description of interactive systems and tools for their analysis.

# Examples of interaction networks: the EGF pathway

# Examples of interaction networks: the *lac* operon

# Outline of the talk

# The Calculus of Looping Sequences (CLS)

We assume an alphabet $\mathcal{E}$. **Terms** $T$ and **Sequences** $S$ of CLS are given by the following grammar:

$$T \ ::= \ S \ \mid \ (S)^L \rfloor \ T \ \mid \ T \mid T$$
$$S \ ::= \ \epsilon \ \mid \ a \ \mid \ S \cdot S$$

where $a$ is a generic element of $\mathcal{E}$, and $\epsilon$ is the empty sequence.

The operators are:

$\quad S \cdot S \quad$ : Sequencing

$\quad (S)^L \quad$ : Looping ($S$ is closed and it can rotate)

$\quad T_1 \rfloor T_2 \quad$ : Containment ($T_1$ contains $T_2$)

$\quad T \mid T \quad$ : Parallel composition (juxtaposition)

Actually, looping and containment form a single binary operator $(S)^L \rfloor T$.

# Example of Terms



$$(i) \quad (a \cdot b \cdot c)^L \rfloor \epsilon$$

$$(ii) \quad (a \cdot b \cdot c)^L \rfloor (d \cdot e)^L \rfloor \epsilon$$

$$(iii) \quad (a \cdot b \cdot c)^L \rfloor (f \cdot g \mid (d \cdot e)^L \rfloor \epsilon)$$

# 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:

$$S_1 \cdot (S_2 \cdot S_3) \equiv_S (S_1 \cdot S_2) \cdot S_3 \qquad S \cdot \epsilon \equiv_S \epsilon \cdot S \equiv_S S$$

$$T_1 \mid T_2 \equiv_T T_2 \mid T_1 \qquad T_1 \mid (T_2 \mid T_3) \equiv_T (T_1 \mid T_2) \mid T_3$$

$$T \mid \epsilon \equiv_T T \quad \left(\epsilon\right)^L \rfloor \epsilon \equiv_T \epsilon \quad \left(S_1 \cdot S_2\right)^L \rfloor T \equiv_T \left(S_2 \cdot S_1\right)^L \rfloor T$$

We write $\equiv$ for $\equiv_T$.

# CLS Patterns

Let us consider variables of three kinds:

- term variables $(X, Y, Z, \ldots)$
- sequence variables $(\widetilde{x}, \widetilde{y}, \widetilde{z}, \ldots)$
- element variables $(x, y, z, \ldots)$

**Patterns** $P$ and **Sequence Patterns** $SP$ of CLS extend CLS terms and sequences with variables:

$$P \ ::= \ SP \ \mid \ (SP)^L \rfloor P \ \mid \ P \mid P \ \mid \ X$$
$$SP \ ::= \ \epsilon \ \mid \ a \ \mid \ SP \cdot SP \ \mid \ x \ \mid \ \widetilde{x}$$

where $a$ is a generic element of $\mathcal{E}$, $\epsilon$ is the empty sequence, and $x, \widetilde{x}$ and $X$ are generic element, sequence and term variables

The structural congruence relation $\equiv$ extends trivially to patterns

# Rewrite Rules

$P\sigma$ denotes the term obtained by replacing any variable in $T$ with the corresponding term, sequence or element.

$\Sigma$ is the set of all possible instantiations $\sigma$

A **Rewrite Rule** is a pair $(P, P')$, denoted $P \mapsto P'$, where:

- $P, P'$ are patterns
- variables in $P'$ are a subset of those in $P$

A rule $P \mapsto P'$ can be applied to all terms $P\sigma$.

Example: $a \cdot x \cdot a \mapsto b \cdot x \cdot b$

- can be applied to $a \cdot c \cdot a$ (producing $b \cdot c \cdot b$)
- cannot be applied to $a \cdot c \cdot c \cdot a$

## Formal Semantics

Given a set of rewrite rules $\mathcal{R}$, evolution of terms is described by the transition system given by the least relation $\rightarrow$ satisfying

$$\frac{P \mapsto P' \in \mathcal{R} \qquad P\sigma \not\equiv \epsilon}{P\sigma \rightarrow P'\sigma}$$

$$\frac{T \rightarrow T'}{T \mid T'' \rightarrow T' \mid T''} \qquad \frac{T \rightarrow T'}{(S)^L \rfloor T \rightarrow (S)^L \rfloor T'}$$

and closed under structural congruence $\equiv$.

# Biomolecular entities as CLS terms

| Biomolecular Entity | CLS Term |
|---|---|
| Gene, protein domain, macro molecule, . . . | Alphabet symbol |
| DNA strand | Sequence of elements representing genes |
| RNA strand | Sequence of elements representing transcribed genes |
| Protein | Single alphabet symbols or sequence of elements representing domains |
| Molecular population | Parallel composition of molecules |
| Membrane | Looping sequence |

# Biomolecular events as CLS rewrite rules

| Biomolecular Event | Examples of CLS Rewrite Rule |
|---|---|
| Complexation | $a \mid b \mapsto c \qquad \widetilde{x} \cdot a \cdot \widetilde{y} \mid b \mapsto \widetilde{x} \cdot c \cdot \widetilde{y}$ |
| Catalysis | $c \mid P_1 \mapsto c \mid P_2$ <br> where $P_1 \mapsto P_2$ is the catalyzed event |
| Complexation <br> on membrane | $\left( a \cdot \widetilde{x} \cdot b \cdot \widetilde{y} \right)^L \rfloor X \mapsto \left( c \cdot \widetilde{x} \cdot \widetilde{y} \right)^L \rfloor X$ <br> $a \mid \left( b \cdot \widetilde{x} \right)^L \rfloor X \mapsto \left( c \cdot \widetilde{x} \right)^L \rfloor X$ |
| Membrane crossing | $a \mid \left( \widetilde{x} \right)^L \rfloor X \mapsto \left( \widetilde{x} \right)^L \rfloor (a \mid X)$ <br> $\widetilde{x} \cdot a \cdot \widetilde{y} \mid \left( \widetilde{z} \right)^L \rfloor X \mapsto \left( \widetilde{z} \right)^L \rfloor \left( \widetilde{x} \cdot a \cdot \widetilde{y} \mid X \right)$ |
| Membrane joining | $\left( \widetilde{x} \right)^L \rfloor (a \mid X) \mapsto \left( a \cdot \widetilde{x} \right)^L \rfloor X$ <br> $\left( \widetilde{x} \right)^L \rfloor \left( \widetilde{y} \cdot a \cdot \widetilde{z} \mid X \right) \mapsto \left( \widetilde{y} \cdot a \cdot \widetilde{z} \cdot \widetilde{x} \right)^L \rfloor X$ |
| Catalyzed <br> membrane fusion | $\left( a \cdot \widetilde{x} \right)^L \rfloor (X) \mid \left( b \cdot \widetilde{y} \right)^L \rfloor (Y) \mapsto$ <br> $\qquad \left( a \cdot \widetilde{x} \cdot b \cdot \widetilde{y} \right)^L \rfloor (X \mid Y)$ |
| Catalyzed <br> membrane division | $\left( a \cdot \widetilde{x} \cdot b \cdot \widetilde{y} \right)^L \rfloor (X \mid Y) \mapsto$ <br> $\qquad \left( a \cdot \widetilde{x} \right)^L \rfloor (X) \mid \left( b \cdot \widetilde{y} \right)^L \rfloor (Y)$ |

# CLS modeling examples: the EGF pathway (1)

# CLS modeling examples: the EGF pathway (2)

First steps of the EGF signaling pathway up to the binding of the signal-receptor dimer to the SHC protein

- The EGFR, EGF and SHC proteins are modeled as the alphabet symbols *EGFR*, *EGF* and *SHC*, respectively
- The cell is modeled as a looping sequence (representing its external membrane):

$$EGF \mid EGF \mid \left(EGFR \cdot EGFR \cdot EGFR \cdot EGFR\right)^{L} \rfloor (SHC \mid SHC)$$

Rewrite rules modeling the first steps of the pathway:

$$EGF \mid \left(EGFR \cdot \widetilde{x}\right)^{L} \rfloor X \;\mapsto\; \left(CMPLX \cdot \widetilde{x}\right)^{L} \rfloor X \qquad (R1)$$

$$\left(CMPLX \cdot \widetilde{x} \cdot CMPLX \cdot \widetilde{y}\right)^{L} \rfloor X \;\mapsto\; \left(DIM \cdot \widetilde{x} \cdot \widetilde{y}\right)^{L} \rfloor X \qquad (R2)$$

$$\left(DIM \cdot \widetilde{x}\right)^{L} \rfloor X \;\mapsto\; \left(DIMp \cdot \widetilde{x}\right)^{L} \rfloor X \qquad (R3)$$

$$\left(DIMp \cdot \widetilde{x}\right)^{L} \rfloor (SHC \mid X) \;\mapsto\; \left(DIMpSHC \cdot \widetilde{x}\right)^{L} \rfloor X \qquad (R4)$$

# CLS modeling examples: the EGFR pathway (2)

A possible evolution of the system:

$$EGF \mid EGF \mid (EGFR \cdot EGFR \cdot EGFR \cdot EGFR)^L \rfloor (SHC \mid SHC)$$

$$\xrightarrow{(R1)} EGF \mid (EGFR \cdot CMPLX \cdot EGFR \cdot EGFR)^L \rfloor (SHC \mid SHC)$$

$$\xrightarrow{(R1)} (EGFR \cdot CMPLX \cdot EGFR \cdot CMPLX)^L \rfloor (SHC \mid SHC)$$

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

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

$$\xrightarrow{(R4)} (EGFR \cdot DIMpSHC \cdot EGFR)^L \rfloor SHC$$

# CLS modeling examples: the *lac* operon (2)

$$Ecoli ::= (m)^L \rfloor (lacI \cdot lacP \cdot lacO \cdot lacZ \cdot lacY \cdot lacA \mid polym)$$

Rules for DNA transcription/translation:

$$lacI \cdot \widetilde{x} \mapsto lacI' \cdot \widetilde{x} \mid repr \qquad (R1)$$

$$polym \mid \widetilde{x} \cdot lacP \cdot \widetilde{y} \mapsto \widetilde{x} \cdot PP \cdot \widetilde{y} \qquad (R2)$$

$$\widetilde{x} \cdot PP \cdot lacO \cdot \widetilde{y} \mapsto \widetilde{x} \cdot lacP \cdot PO \cdot \widetilde{y} \qquad (R3)$$

$$\widetilde{x} \cdot PO \cdot lacZ \cdot \widetilde{y} \mapsto \widetilde{x} \cdot lacO \cdot PZ \cdot \widetilde{y} \qquad (R4)$$

$$\widetilde{x} \cdot PZ \cdot lacY \cdot \widetilde{y} \mapsto \widetilde{x} \cdot lacZ \cdot PY \cdot \widetilde{y} \mid betagal \qquad (R5)$$

$$\widetilde{x} \cdot PY \cdot lacA \mapsto \widetilde{x} \cdot lacY \cdot PA \mid perm \qquad (R6)$$

$$\widetilde{x} \cdot PA \mapsto \widetilde{x} \cdot lacA \mid transac \mid polym \qquad (R7)$$

# CLS modeling examples: the *lac* operon (3)

$$Ecoli ::= (m)^L \rfloor (lacI \cdot lacP \cdot lacO \cdot lacZ \cdot lacY \cdot lacA \mid polym)$$

Rules to describe the binding of the lac Repressor to gene o, and what happens when lactose is present in the environment of the bacterium:

$$repr \mid \widetilde{x} \cdot lacO \cdot \widetilde{y} \mapsto \widetilde{x} \cdot RO \cdot \widetilde{y} \tag{R8}$$

$$LACT \mid (m \cdot \widetilde{x})^L \rfloor X \mapsto (m \cdot \widetilde{x})^L \rfloor (X \mid LACT) \tag{R9}$$

$$\widetilde{x} \cdot RO \cdot \widetilde{y} \mid LACT \mapsto \widetilde{x} \cdot lacO \cdot \widetilde{y} \mid RLACT \tag{R10}$$

$$(\widetilde{x})^L \rfloor (perm \mid X) \mapsto (perm \cdot \widetilde{x})^L \rfloor X \tag{R11}$$

$$LACT \mid (perm \cdot \widetilde{x})^L \rfloor X \mapsto (perm \cdot \widetilde{x})^L \rfloor (LACT \mid X) \tag{R12}$$

$$betagal \mid LACT \mapsto betagal \mid GLU \mid GAL \tag{R13}$$

# CLS modeling examples: the *lac* operon (4)

$$Ecoli ::= (m)^L \rfloor (lacI \cdot lacP \cdot lacO \cdot lacZ \cdot lacY \cdot lacA \mid polym)$$

Example:

$Ecoli|LACT|LACT$

$\rightarrow^* (m)^L \rfloor (lacI' \cdot lacP \cdot lacO \cdot lacZ \cdot lacY \cdot lacA \mid polym \mid repr)|LACT|LACT$

$\rightarrow^* (m)^L \rfloor (lacI' \cdot lacP \cdot RO \cdot lacZ \cdot lacY \cdot lacA \mid polym)|LACT|LACT$

$\rightarrow^* (m)^L \rfloor (lacI' \cdot lacP \cdot lacO \cdot lacZ \cdot lacY \cdot lacA|polym|RLACT)|LACT$

$\rightarrow^* (perm \cdot m)^L \rfloor (lacI'{-}A|betagal|transac|polym|RLACT)|LACT$

$\rightarrow^* (perm \cdot m)^L \rfloor (lacI'{-}A|betagal|transac|polym|RLACT|GLU|GAL)$

# Outline of the talk

# Bisimulations

Bisimilarity is widely accepted as the finest extensional behavioral equivalence one may impose on systems.

- Two systems are bisimilar if they can perform step by step the same interactions with the environment.
- Properties of a system can be verified by assessing the bisimilarity with a system known to enjoy them.

Bisimilarities need semantics based on labeled transition relations capturing the potential interactions with the environment.

- In process calculi, transitions are usually labeled with actions.
- In CLS labels are contexts in which rules can be applied.

## Labeled semantics (1)

**Contexts** $\mathcal{C}$ are given by the following grammar:

$$\mathcal{C} ::= \square \quad | \quad \mathcal{C} \mid T \quad | \quad T \mid \mathcal{C} \quad | \quad (S)^L \rfloor \mathcal{C}$$

where $T \in \mathcal{T}$ and $S \in \mathcal{S}$. Context $\square$ is called the *empty context*.

**Parallel Contexts** $\mathcal{C}_P$ are given by the following grammar:

$$\mathcal{C}_P ::= \square \quad | \quad \mathcal{C}_P \mid T \quad | \quad T \mid \mathcal{C}_P.$$

where $T \in \mathcal{T}$.

$C[T]$ is context application and $C[C']$ is context composition.

## Labeled semantics (2)

Given a set of rewrite rules $\mathcal{R} \subseteq \Re$, the **labeled semantics** of CLS is the labeled transition system given by the following inference rules:

$$(\text{rule\_appl}) \ \frac{P \mapsto P' \in \mathcal{R} \quad C[T''] \equiv P\sigma \quad T'' \not\equiv \epsilon \quad \sigma \in \Sigma \quad C \in \mathcal{C}}{T'' \xrightarrow{C} T'\sigma}$$

$$(\text{cont}) \ \frac{T \xrightarrow{\square} T'}{(S)^L \rfloor T \xrightarrow{\square} (S)^L \rfloor T'} \qquad (\text{par}) \ \frac{T \xrightarrow{C} T' \quad C \in \mathcal{C}_P}{T \mid T'' \xrightarrow{C} T' \mid T''}$$

where the dual version of the *(par)* rule is omitted.

Rule (rule_appl) describes the (potential) application of a rule.

- $T'' \not\equiv \epsilon$ in the premise implies that $C$ cannot provide completely the left hand side of the rewrite rule.
- Example: let $R = a \mid b \mapsto c$, we have $a \xrightarrow{\square \mid b} c$, but $\epsilon \xrightarrow{a \mid b} \not\rightarrow$.

## Labeled semantics (3)

Given a set of rewrite rules $\mathcal{R} \subseteq \Re$, the **labeled semantics** of CLS is the labeled transition system given by the following inference rules:

$$\text{(rule\_appl)} \ \frac{P \mapsto P' \in \mathcal{R} \quad C[T''] \equiv T\sigma \quad T'' \not\equiv \epsilon \quad \sigma \in \Sigma \quad C \in \mathcal{C}}{T'' \xrightarrow{C} T'\sigma}$$

$$\text{(cont)} \ \frac{T \xrightarrow{\square} T'}{(S)^L \rfloor T \xrightarrow{\square} (S)^L \rfloor T'} \qquad \text{(par)} \ \frac{T \xrightarrow{C} T' \quad C \in \mathcal{C}_P}{T \mid T'' \xrightarrow{C} T' \mid T''}$$

where the dual version of the *(par)* rule is omitted.

Rule (cont) propagates $\square$–labeled transitions from the inside to the outside of a looping sequence.

- Transition labeled with a non–empty context cannot be propagated.
- Example: let $R = a \mid b \mapsto c$, we have $a \xrightarrow{\square \mid b} c$, but $(d)^L \rfloor a \xrightarrow{\square \mid b} \!\!\!\!\!/$.

## Labeled semantics (4)

Given a set of rewrite rules $\mathcal{R} \subseteq \Re$, the **labeled semantics** of CLS is the labeled transition system given by the following inference rules:

$$(\text{rule\_appl}) \ \frac{P \mapsto P' \in \mathcal{R} \quad C[T''] \equiv T\sigma \quad T'' \not\equiv \epsilon \quad \sigma \in \Sigma \quad C \in \mathcal{C}}{T'' \xrightarrow{C} T'\sigma}$$

$$(\text{cont}) \ \frac{T \xrightarrow{\square} T'}{(S)^L \rfloor T \xrightarrow{\square} (S)^L \rfloor T'} \qquad (\text{par}) \ \frac{T \xrightarrow{C} T' \quad C \in \mathcal{C}_P}{T \mid T'' \xrightarrow{C} T' \mid T''}$$

where the dual version of the *(par)* rule is omitted.

Rule (par) propagates transitions labeled with parallel contexts in parallel components.

- Example: let $R = (a)^L \rfloor b \mapsto c$, we have $b \xrightarrow{(a)^L \rfloor \square} c$, but $b \mid d \xrightarrow{(a)^L \rfloor \square} \!\!\!\!\!/ \ $ because $R$ cannot be applied $(a)^L \rfloor (b \mid d)$

# Bisimulations in CLS (1)

A binary relation $R$ on terms is a **strong bisimulation** if, given $T_1, T_2$ such that $T_1 R T_2$, the two following conditions hold:

- $T_1 \xrightarrow{C} T_1' \implies \exists T_2'$ s.t. $T_2 \xrightarrow{C} T_2'$ and $T_1' R T_2'$
- $T_2 \xrightarrow{C} T_2' \implies \exists T_1'$ s.t. $T_1 \xrightarrow{C} T_1'$ and $T_2' R T_1'$.

The *strong bisimilarity* $\sim$ is the largest of such relations.

A binary relation $R$ on terms is a **weak bisimulation** if, given $T_1, T_2$ such that $T_1 R T_2$, the two following conditions hold:

- $T_1 \xrightarrow{C} T_1' \implies \exists T_2'$ s.t. $T_2 \xRightarrow{C} T_2'$ and $T_1' R T_2'$
- $T_2 \xrightarrow{C} T_2' \implies \exists T_1'$ s.t. $T_1 \xRightarrow{C} T_1'$ and $T_2' R T_1'$.

The *weak bisimilarity* $\approx$ is the largest of such relations.

**Theorem:** Strong and weak bisimilarities are congruences.

# Bisimulations in CLS (2)

Consider the following set of rewrite rules:

$$\mathcal{R} = \{ \quad a \mid b \mapsto c \quad , \quad d \mid b \mapsto e \quad , \quad e \mapsto e \quad , \quad c \mapsto e \quad , \quad f \mapsto a \quad \}$$

We have that $a \sim d$, because

$$a \xrightarrow{\square \mid b} c \xrightarrow{\square} e \xrightarrow{\square} e \xrightarrow{\square} \ldots$$

$$d \xrightarrow{\square \mid b} e \xrightarrow{\square} e \xrightarrow{\square} \ldots$$

and $f \approx d$, because

$$f \xrightarrow{\square} a \xrightarrow{\square \mid b} c \xrightarrow{\square} e \xrightarrow{\square} e \xrightarrow{\square} \ldots$$

On the other hand, $f \not\sim e$ and $f \not\approx e$.

$$e \xrightarrow{\square} e \xrightarrow{\square} e \xrightarrow{\square} \ldots$$

# Bisimulations in CLS (3)

Let us consider systems $(T, \mathcal{R})$...

A binary relation $R$ is a **strong bisimulation on systems** if, given $(T_1, \mathcal{R}_1)$ and $(T_2, \mathcal{R}_2)$ such that $(T_1, \mathcal{R}_1)R(T_2, \mathcal{R}_2)$:

- $\mathcal{R}_1 : T_1 \xrightarrow{C} T_1' \implies \exists T_2'$ s.t. $\mathcal{R}_2 : T_2 \xrightarrow{C} T_2'$ and $(T_1', \mathcal{R}_1)R(T_2', \mathcal{R}_2)$
- $\mathcal{R}_2 : T_2 \xrightarrow{C} T_2' \implies \exists T_1'$ s.t. $\mathcal{R}_1 : T_1 \xrightarrow{C} T_1'$ and $(R_2, T_2')R(\mathcal{R}_1, T_1')$.

The *strong bisimilarity on systems* $\sim$ is the largest of such relations.

A binary relation $R$ is a **weak bisimulation on systems** if, given $(T_1, \mathcal{R}_1)$ and $(T_2, \mathcal{R}_2)$ such that $(T_1, \mathcal{R}_1)R(T_2, \mathcal{R}_2)$:

- $\mathcal{R}_1 : T_1 \xrightarrow{C} T_1' \implies \exists T_2'$ s.t. $\mathcal{R}_2 : T_2 \xRightarrow{C} T_2'$ and $(T_1', \mathcal{R}_1)R(T_2', \mathcal{R}_2)$
- $\mathcal{R}_2 : T_2 \xrightarrow{C} T_2' \implies \exists T_1'$ s.t. $\mathcal{R}_1 : T_1 \xRightarrow{C} T_1'$ and $(T_2', \mathcal{R}_2)R(T_1', \mathcal{R}_1)$

The *weak bisimilarity on systems* $\approx$ is the largest of such relations.

Strong and weak bisimilarities on systems are NOT congruences.

# Bisimulations in CLS (4)

Consider the following sets of rewrite rules

$$\mathcal{R}_1 = \{a \mid b \mapsto c\} \qquad \mathcal{R}_2 = \{a \mid d \mapsto c \,, \ b \mid e \mapsto c\}$$

We have that $\langle a, \mathcal{R}_1 \rangle \approx \langle e, \mathcal{R}_2 \rangle$ because

$$\mathcal{R}_1 : a \xrightarrow{\Box \mid b} c \qquad \mathcal{R}_2 : e \xrightarrow{\Box \mid b} c$$

and $\langle b, \mathcal{R}_1 \rangle \approx \langle d, \mathcal{R}_2 \rangle$, because

$$\mathcal{R}_1 : b \xrightarrow{\Box \mid a} c \qquad \mathcal{R}_2 : d \xrightarrow{\Box \mid a} c$$

but $\langle a \mid b, \mathcal{R}_1 \rangle \not\approx \langle e \mid d, \mathcal{R}_2 \rangle$, because

$$\mathcal{R}_1 : a \mid b \xrightarrow{\Box} c \qquad \mathcal{R}_2 : c \mid d \not\rightarrow$$

# Applying bisimulations to the *lac* operon (1)

$$Ecoli ::= (m)^L \rfloor (lacI \cdot lacP \cdot lacO \cdot lacZ \cdot lacY \cdot lacA \mid polym)$$

It can be easily proved that

$$lacI \cdot lacP \cdot lacO \cdot lacZ \cdot lacY \cdot lacA$$
$$\approx$$
$$lacP \cdot lacO \cdot lacZ \cdot lacY \cdot lacA \mid repr$$

and since weak bisimularity is a congruence the former can be replaced by the latter in the model.

## Applying bisimulations to the *lac* operon (2)

By using the weak bisimilarity on systems we can prove that from the state in which the repressor is bound to the DNA we can reach a state in which the enzymes are synthesized only if lactose appears in the environment.

We replace rule

$$\widetilde{x} \cdot RO \cdot \widetilde{y} \mid LACT \;\mapsto\; \widetilde{x} \cdot lacO \cdot \widetilde{y} \mid RLACT \qquad\qquad (R10)$$

with

$$\left(\widetilde{w}\right)^L \rfloor \left(\widetilde{x} \cdot RO \cdot \widetilde{y} \mid LACT \mid X\right) \mid START \;\mapsto$$
$$\left(\widetilde{w}\right)^L \rfloor \left(\widetilde{x} \cdot lacO \cdot \widetilde{y} \mid RLACT \mid X\right) \qquad (R10bis)$$

The obtained model is bisimilar to $(T_1, \mathcal{R})$ where $\mathcal{R}$ is

| | | | |
|---|---|---|---|
| $T_1 \mid LACT \;\mapsto\; T_2$ | (R1') | $T_2 \mid START \;\mapsto\; T_3$ | (R3') |
| $T_2 \mid LACT \;\mapsto\; T_2$ | (R2') | $T_3 \mid LACT \;\mapsto\; T_3$ | (R4') |

that is a system satisfying the property.

# Some theoretical results

CLS is Turing complete

- A Turing machine encoded into a CLS term and a single rewrite rule

Formalisms capable of describing membranes can be encoded into CLS

- Brane Calculi
- P Systems

Bisimilarities of Brane Calculi are preserved after translation into CLS

# Some variants of CLS

- Full–CLS
  - The looping operator can be applied to any term
  - Rule $a \mid b \mapsto c$ can be applied to $b \mid (a \cdot a \cdot a \cdot a)^L \rfloor d$

- CLS+
  - More realistic representation of the fluid nature of membranes: the looping operator can be applied to parallel compositions of sequences
  - Can be encoded into CLS

- Stochastic CLS
  - The application of a rule consumes a stochastic quantity of time

- LCLS (CLS with Links)
  - Description of protein–protein interactions at the domain level

# Outline of the talk

## Background: the kinetics of chemical reactions

Usual notation for chemical reactions:

$$\ell_1 S_1 + \ldots + \ell_\rho S_\rho \overset{k}{\underset{k_{-1}}{\rightleftharpoons}} \ell_1' P_1 + \ldots + \ell_\gamma' P_\gamma$$

where:

- $S_i, P_i$ are molecules (reactants)
- $\ell_i, \ell_i'$ are stoichiometric coefficients
- $k, k_{-1}$ are the kinetic constants
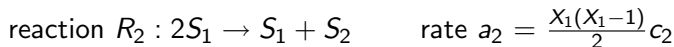
The kinetics is described by the *law of mass action*:

$$\frac{d[P_i]}{dt} = \ell_i' \underbrace{k[S_1]^{\ell_1} \cdots [S_\rho]^{\ell_\rho}}_{reaction\ rate} \qquad \frac{d[S_i]}{dt} = \ell_i \underbrace{k_{-1}[P_1]^{\ell_1'} \cdots [P_\gamma]^{\ell_\gamma'}}_{reaction\ rate}$$

# Background: Gillespie's simulation algorithm

- represents a chemical solution as a multiset of molecules
- computes the reaction rate $a_\mu$ by multiplying the kinetic constant by the number of possible combinations of reactants

Example: chemical solution with $X_1$ molecules $S_1$ and $X_2$ molecules $S_2$

reaction $R_1 : S_1 + S_2 \rightarrow 2S_1$      rate $a_1 = X_1 X_2 c_1$

reaction $R_2 : 2S_1 \rightarrow S_1 + S_2$      rate $a_2 = \frac{X_1(X_1-1)}{2} c_2$

Given a set of reactions $\{R_1, \ldots R_M\}$ and a current time $t$

- The time $t + \tau$ at which the next reaction will occur is randomly chosen with $\tau$ exponentially distributed with parameter $\sum_{\nu=1}^{M} a_\nu$;
- The reaction $R_\mu$ that has to occur at time $t + \tau$ is randomly chosen with probability $\frac{a_\mu}{\sum_{\nu=1}^{M} a_\nu}$.

At each step $t$ is incremented by $\tau$ and the chemical solution is updated.

# Stochastic CLS (1)

Stochastic CLS incorporates Gillespie's stochastic framework into the semantics of CLS

What is a reactant in Stochastic CLS?

- A *subterm* of a term $T$ is a term $T' \not\equiv \epsilon$ such that $T \equiv C[T']$ for some context $C$
- A *reactant* is an occurence of a subterm

Example: given $T = a \mid a \mid b \mid b$

- the set of subterms of $T$ is

$$\{ a ,\, b ,\, a \mid a ,\, a \mid b ,\, b \mid b ,\, a \mid a \mid b ,\, a \mid b \mid b ,\, T \}$$

- the multiset of reactants in $T$ is

$$\{ a ,\, a ,\, b ,\, b ,\, a \mid a ,\, a \mid b ,\, a \mid b ,\, a \mid b ,\, a \mid b ,\, b \mid b ,$$
$$a \mid a \mid b ,\, a \mid a \mid b ,\, a \mid b \mid b ,\, a \mid b \mid b ,\, T \}$$

# Stochastic CLS (2)

- Given $T = a \mid a \mid b \mid b$ the multiset of reactants in $T$ is

$$\{a\,,\ a\,,\ b\,,\ b\,,\ a \mid a\,,\ a \mid b\,,\ a \mid b\,,\ a \mid b\,,\ a \mid b\,,\ b \mid b\,,$$
$$a \mid a \mid b\,,\ a \mid a \mid b\,,\ a \mid b \mid b\,,\ a \mid b \mid b\,,\ T\}$$

Defining the stochastic semantics would be easy for rules without variables:

- Rewrite rules could be extended with kinetic constants (e.g. $a \mid b \overset{k}{\mapsto} c$)
- The number of possible combinations of molecules involved in the reactions corresponds to the number of reactants equivalent to the left-hand side of the rule (e.g. $2 \times 2 = 4$ corresponds to 4 occurrences of $a \mid b$)

# Stochastic CLS (3)

We consider rewrite rules containing variables as *rewrite rule schemata*

- at step we compute the set of ground rules that can be applied among those obtained by instantiating variables of the rewrite rule schema
- we reduce the problem of defining the semantics with rule schemata to the simpler problem of defining the semantics with ground rules only

Example: given $T = a \cdot b \mid a \cdot c$

From rule schema $a \cdot \widetilde{x} \mid a \cdot \widetilde{y} \stackrel{k}{\mapsto} d$ we can derive only $a \cdot b \mid a \cdot c \stackrel{k}{\mapsto} d$

From rule schema $a \cdot \widetilde{x} \mid a \cdot \widetilde{y} \stackrel{k}{\mapsto} \widetilde{y}$ we can derive both $a \cdot b \mid a \cdot c \stackrel{k}{\mapsto} b$ and $a \cdot b \mid a \cdot c \stackrel{k}{\mapsto} c$

Problem: the kinetic constant could be different for different instantiations

- We enrich rewrite rules with rate functions $f : \Sigma \to \mathbb{R}$ rather than $k$

## Stochastic CLS (4)

Given a finite set of rewrite rule schemata $\mathcal{R}$, the semantics of Stochastic CLS is given by the following inference rule
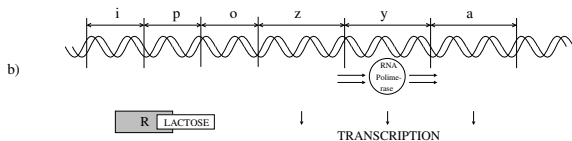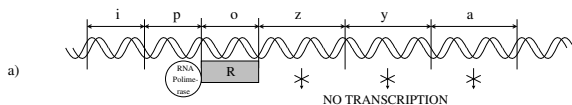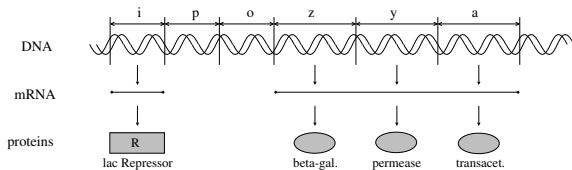
$$\frac{R = T_1 \overset{k}{\mapsto} T_2 \in AR(\mathcal{R}, T) \qquad T \equiv C[T_1]}{T \xrightarrow{R, k \cdot AC(R, T, C[T_2])} C[T_2]}$$

where:

- $AR(\mathcal{R}, T)$ is the set of ground rewrite rules obtained by schemata in $\mathcal{R}$ and applicable to $T$
- $AC(R, T, T')$ is the number of reactants in $T$ equivalent to the left–hand side of the ground rule $R$ and that allows obtaining term $T'$ after the application of $R$
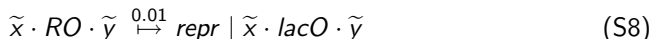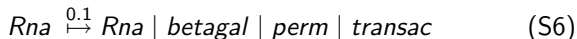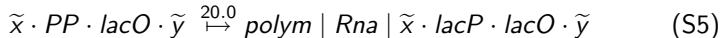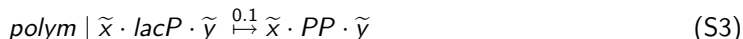
The transition system obtained can be easily transformed into a *Continuous Time Markov Chain*
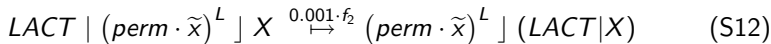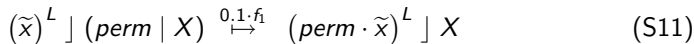
# A Stochastic CLS model of the *lac* operon (1)

# A Stochastic CLS model of the *lac* operon (2)

Transcription of DNA, binding of lac Repressor to gene o, and interaction between lactose and lac Repressor:

$$lacI \cdot \widetilde{x} \overset{0.02}{\mapsto} lacI \cdot \widetilde{x} \mid lrna \tag{S1}$$

$$lrna \overset{0.1}{\mapsto} lrna \mid repr \tag{S2}$$

$$polym \mid \widetilde{x} \cdot lacP \cdot \widetilde{y} \overset{0.1}{\mapsto} \widetilde{x} \cdot PP \cdot \widetilde{y} \tag{S3}$$

$$\widetilde{x} \cdot PP \cdot \widetilde{y} \overset{0.01}{\mapsto} polym \mid \widetilde{x} \cdot lacP \cdot \widetilde{y} \tag{S4}$$

$$\widetilde{x} \cdot PP \cdot lacO \cdot \widetilde{y} \overset{20.0}{\mapsto} polym \mid Rna \mid \widetilde{x} \cdot lacP \cdot lacO \cdot \widetilde{y} \tag{S5}$$

$$Rna \overset{0.1}{\mapsto} Rna \mid betagal \mid perm \mid transac \tag{S6}$$

$$repr \mid \widetilde{x} \cdot lacO \cdot \widetilde{y} \overset{1.0}{\mapsto} \widetilde{x} \cdot RO \cdot \widetilde{y} \tag{S7}$$

$$\widetilde{x} \cdot RO \cdot \widetilde{y} \overset{0.01}{\mapsto} repr \mid \widetilde{x} \cdot lacO \cdot \widetilde{y} \tag{S8}$$

$$repr \mid LACT \overset{0.005}{\mapsto} RLACT \tag{S9}$$

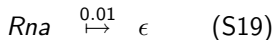$$RLACT \overset{0.1}{\mapsto} repr \mid LACT \tag{S10}$$

## A Stochastic CLS model of the *lac* operon (3)

The behaviour of the three enzymes for lactose degradation:

$$(\widetilde{x})^L \rfloor (perm \mid X) \overset{0.1 \cdot f_1}{\mapsto} (perm \cdot \widetilde{x})^L \rfloor X \qquad (S11)$$

$$LACT \mid (perm \cdot \widetilde{x})^L \rfloor X \overset{0.001 \cdot f_2}{\mapsto} (perm \cdot \widetilde{x})^L \rfloor (LACT \mid X) \qquad (S12)$$

$$betagal \mid LACT \overset{0.001}{\mapsto} betagal \mid GLU \mid GAL \qquad (S13)$$
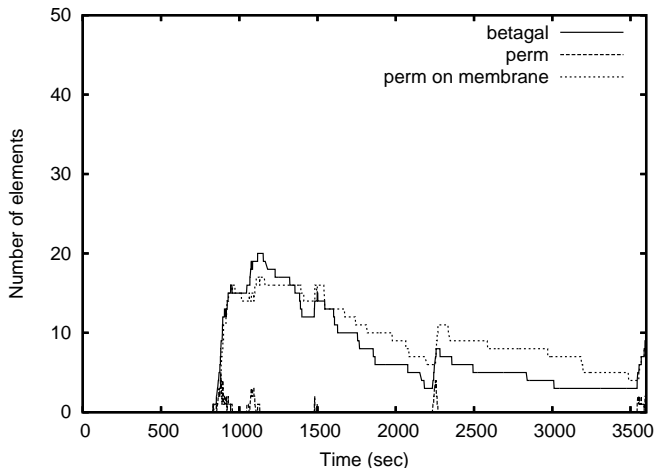
where $f_1(\sigma) = occ(perm, \sigma(X)) + 1$, $f_2(\sigma) = occ(perm, \sigma(\widetilde{x})) + 1$.

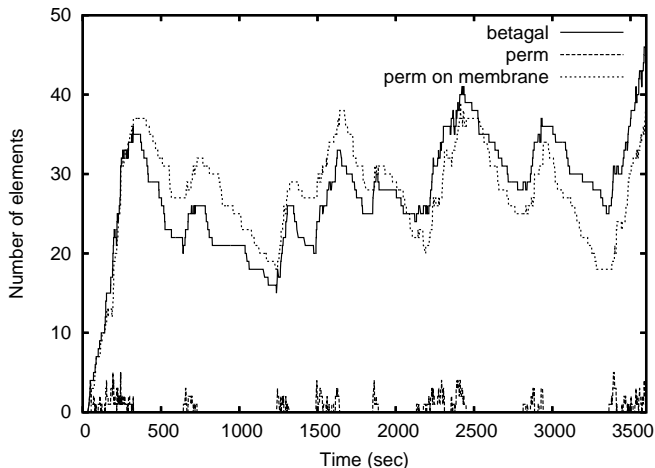Degradation of all the proteins and mRNA involved in the process:

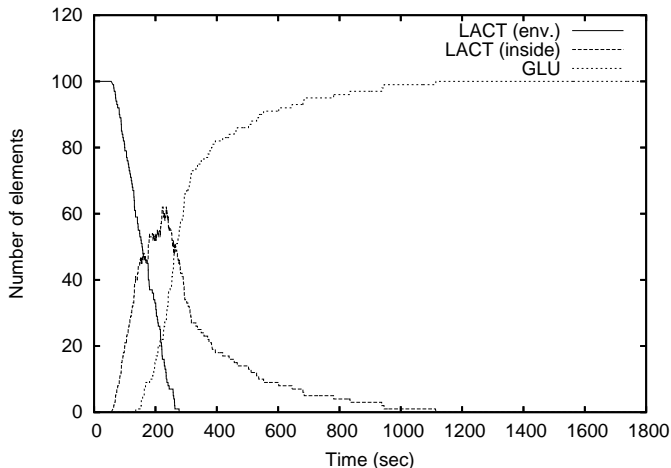| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $perm$ | $\overset{0.001}{\mapsto}$ | $\epsilon$ | (S14) | $betagal$ | $\overset{0.001}{\mapsto}$ | $\epsilon$ | (S15) |
| $transac$ | $\overset{0.001}{\mapsto}$ | $\epsilon$ | (S16) | $repr$ | $\overset{0.002}{\mapsto}$ | $\epsilon$ | (S17) |
| $lrna$ | $\overset{0.01}{\mapsto}$ | $\epsilon$ | (S18) | $Rna$ | $\overset{0.01}{\mapsto}$ | $\epsilon$ | (S19) |
| $RLACT$ | $\overset{0.002}{\mapsto}$ | $LACT$ | (S20) | | | | |

# Simulation results (1)



Production of enzymes in the absence of lactose
$$\left(m\right)^{L} \rfloor (lacI - A \mid 30 \times polym \mid 100 \times repr)$$
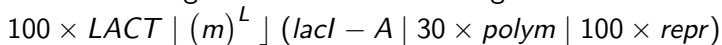
# Simulation results (2)



Production of enzymes in the presence of lactose

$$100 \times LACT \mid \big(m\big)^L \, \rfloor \, (lacI - A \mid 30 \times polym \mid 100 \times repr)$$

# Simulation results (3)



Degradation of lactose into glucose

$$100 \times LACT \mid (m)^L \rfloor (lacI - A \mid 30 \times polym \mid 100 \times repr)$$

# Outline of the talk

# Modeling proteins at the domain level

To model a protein at the domain level in CLS it would be natural to use a sequence with one symbol for each domain

The binding between two elements of two different sequences, cannot be expressed in CLS

LCLS extends CLS with labels on basic symbols

- two symbols with the same label represent domains that are bound to each other
- example: $a \cdot b^1 \cdot c \mid d \cdot e^1 \cdot f$

# Syntax of LCLS

**Terms** $T$ and **Sequences** $S$ of LCLS are given by the following grammar:

$$T ::= S \mid (S)^L \rfloor T \mid T \mid T$$
$$S ::= \epsilon \mid a \mid a^n \mid S \cdot S$$

where $a$ is a generic element of $\mathcal{E}$, and $n$ is a natural number.

**Patterns** $P$ and **sequence patterns** $SP$ of LCLS are given by the following grammar:

$$P ::= SP \mid (SP)^L \rfloor P \mid P \mid P \mid X$$
$$SP ::= \epsilon \mid a \mid a^n \mid SP \cdot SP \mid \widetilde{x} \mid x \mid x^n$$

where $a$ is an element of $\mathcal{E}$, $n$ is a natural number and $X, \widetilde{x}$ and $x$ are elements of $TV, SV$ and $\mathcal{X}$, respectively.

# Well–formedness of LCLS terms and patterns

An LCLS term (or pattern) is well–formed if and only if a label occurs no more than twice, and two occurrences of a label are always in the same compartment

Type system for well–formedness:

$$1. \ (\varnothing, \varnothing) \models \epsilon \qquad 2. \ (\varnothing, \varnothing) \models a \qquad 3. \ (\varnothing, \{n\}) \models a^n$$

$$4. \ (\varnothing, \varnothing) \models x \qquad 5. \ (\varnothing, \{n\}) \models x^n \qquad 6. \ (\varnothing, \varnothing) \models \widetilde{x} \qquad 7. \ (\varnothing, \varnothing) \models X$$

$$8. \ \frac{(N_1, N_1') \models SP_1 \quad (N_2, N_2') \models SP_2 \quad N_1 \cap N_2 = N_1' \cap N_2 = N_1 \cap N_2' = \varnothing}{(N_1 \cup N_2 \cup (N_1' \cap N_2'), (N_1' \cup N_2') \setminus (N_1' \cap N_2')) \models SP_1 \cdot SP_2}$$

$$9. \ \frac{(N_1, N_1') \models P_1 \quad (N_2, N_2') \models P_2 \quad N_1 \cap N_2 = N_1' \cap N_2 = N_1 \cap N_2' = \varnothing}{(N_1 \cup N_2 \cup (N_1' \cap N_2'), (N_1' \cup N_2') \setminus (N_1' \cap N_2')) \models P_1 \mid P_2}$$

$$10. \ \frac{(N_1, N_1') \models SP \quad (N_2, N_2') \models P \quad N_1 \cap N_2 = N_1' \cap N_2 = N_1 \cap N_2' = \varnothing \quad N_2' \subseteq N_1'}{(N_1 \cup N_2', N_1' \setminus N_2') \models (SP)^L \rfloor P}$$

# Instantiation of LCLS patterns

An instantiation function $\sigma$ is well–formed if it maps variables into well–formed CLOSED terms and sequences

- otherwise the w.f. pattern $(a)^L \rfloor X$ could be instantiated to the non–w.f. term $(a)^L \rfloor b^1$

The definition of pattern instantiation AVOIDS this kind situations:

- $P = a \cdot \widetilde{x} \mid X$ , $\sigma(\widetilde{x}) = b^1 \cdot c^1$ , $\sigma(X) = d^1 \cdot e^1$ are all w.f.
- $P\sigma = a \cdot b^1 \cdot c^1 \mid d^1 \cdot e^1$ is non–w.f.

Clashing labels are renamed during pattern instatiation

# Compartment safe rewrite rules

By applying a rewrite rule composed by w.f. patterns to a w.f. term by using a w.f. instantiation function we obtain a w.f. term

W.f. instantiations are closed

- Rule $(a)^L \rfloor (\widetilde{x} \mid \widetilde{y}) \mapsto \widetilde{x} \mid (a)^L \rfloor \widetilde{y}$ cannot be applied to $(a)^L \rfloor (b^1 \mid c^1)$ (so to obtain $b^1 \mid (a)^L \rfloor c^1$)

BUT

- Rule $\widetilde{x} \cdot a \mapsto \widetilde{x} \cdot b$ cannot be applied to $c^1 \mid d^1 \cdot a$ (so to obtain $c^1 \mid d^1 \cdot b$)

To allow application of the second kind of rules

- we relax the constraint on instantiations
- we add a constraint on rewrite rules

A *compartment safe* rewrite rule is such that

- it does not add/remove occurrences of variables
- it does not moves variables from one compartment (content of a looping sequence) to another one

# The semantics of LCLS

Given a set of compartment safe rewrite rules $\mathcal{R}^{CS}$ and a set of compartemnt unsafe rewrite rules $\mathcal{R}^{CU}$, the semantics of LCLS is given by the following rules

$$(\text{appCS}) \quad \frac{P_1 \mapsto P_2 \in \mathcal{R}^{CS} \quad P_1\sigma \neq \epsilon \quad \sigma \in \Sigma \quad \alpha \in \mathcal{A}}{P_1\alpha\sigma \to P_2\alpha\sigma}$$

$$(\text{appCU}) \quad \frac{P_1 \mapsto P_2 \in \mathcal{R}^{CU} \quad P_1\sigma \neq \epsilon \quad \sigma \in \Sigma_{wf} \quad \alpha \in \mathcal{A}}{P_1\alpha\sigma \to P_2\alpha\sigma}$$

$$(\text{par}) \quad \frac{T_1 \to T_1' \quad L(T_1) \cap L(T_2) = \{n_1, \ldots, n_M\} \quad n_1', \ldots, n_M' \text{ fresh}}{T_1 \mid T_2 \to T_1'\{n_1', \ldots, n_M'/n_1, \ldots, n_M\} \mid T_2}$$

$$(\text{cont}) \quad \frac{T \to T' \quad L(S) \cap L(T') = \{n_1, \ldots, n_M\} \quad n_1', \ldots, n_M' \text{ fresh}}{(S)^L \rfloor T \to (S)^L \rfloor T'\{n_1', \ldots, n_M'/n_1, \ldots, n_M\}}$$

where $\alpha$ is link renaming, $L(T)$ the set of links occurring twice in the top level compartment of $T$
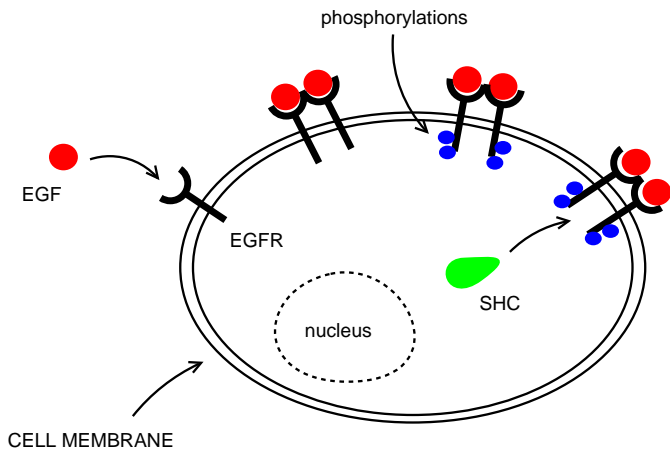
# Main theoretical result

**Theorem (Subject Reduction)**

Given a set of well–formed rewrite rules $\mathcal{R}$ and a well–formed term $T$

$$T \to T' \quad \implies \quad T' \text{ well–formed}$$

# An LCLS model of the EGF pathway (1)

# An LCLS model of the EGF pathway (2)

We model the EGFR protein as the sequence $R_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2}$

- $R_{E1}$ and $R_{E2}$ are two extra–cellular domains
- $R_{I1}$ and $R_{I2}$ are two intra–cellular domains

The rewrite rules of the model are

$$EGF \mid \left( R_{E1} \cdot \widetilde{x} \right)^L \rfloor X \;\mapsto\; \left( SR_{E1} \cdot \widetilde{x} \right)^L \rfloor X \tag{R1}$$

$$\left( SR_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2} \cdot \widetilde{x} \cdot SR_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2} \cdot \widetilde{y} \right)^L \rfloor X \;\mapsto\;$$
$$\left( SR_{E1} \cdot R_{E2}^1 \cdot R_{I1} \cdot R_{I2} \cdot SR_{E1} \cdot R_{E2}^1 \cdot R_{I1} \cdot R_{I2} \cdot \widetilde{x} \cdot \widetilde{y} \right)^L \rfloor X \tag{R2}$$

$$\left( R_{E2}^1 \cdot R_{I1} \cdot \widetilde{x} \cdot R_{E2}^1 \cdot R_{I1} \cdot \widetilde{y} \right)^L \rfloor X \;\mapsto\; \left( R_{E2}^1 \cdot PR_{I1} \cdot \widetilde{x} \cdot R_{E2}^1 \cdot R_{I1} \cdot \widetilde{y} \right)^L \rfloor X \tag{R3}$$

$$\left( R_{E2}^1 \cdot PR_{I1} \cdot \widetilde{x} \cdot R_{E2}^1 \cdot R_{I1} \cdot \widetilde{y} \right)^L \rfloor X \;\mapsto\; \left( R_{E2}^1 \cdot PR_{I1} \cdot \widetilde{x} \cdot R_{E2}^1 \cdot PR_{I1} \cdot \widetilde{y} \right)^L \rfloor X \tag{R4}$$

$$\left( R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot \widetilde{x} \cdot R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot \widetilde{y} \right)^L \rfloor \left( SHC \mid X \right) \;\mapsto\;$$
$$\left( R_{E2}^1 \cdot PR_{I1} \cdot R_{I2}^2 \cdot \widetilde{x} \cdot R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot \widetilde{y} \right)^L \rfloor \left( SHC^2 \mid X \right) \tag{R5}$$

# An LCLS model of the EGF pathway (3)

Let us write $EGFR$ for $R_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2}$

A possible evolution of the system is

$$EGF \mid EGF \mid \left( EGFR \cdot EGFR \cdot EGFR \cdot EGFR \right)^L \rfloor (SHC \mid SHC)$$

$$\xrightarrow{(R1)} EGF \mid \left( SR_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2} \cdot EGFR \cdot EGFR \cdot EGFR \right)^L \rfloor (SHC \mid SHC)$$

$$\xrightarrow{(R1)} \left( SR_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2} \cdot EGFR \cdot SR_{E1} \cdot R_{E2} \cdot R_{I1} \cdot R_{I2} \cdot EGFR \right)^L \rfloor (SHC \mid SHC)$$

$$\xrightarrow{(R2)} \left( SR_{E1} \cdot R_{E2}^1 \cdot R_{I1} \cdot R_{I2} \cdot SR_{E1} \cdot R_{E2}^1 \cdot R_{I1} \cdot R_{I2} \cdot EGFR \cdot EGFR \right)^L \rfloor (SHC \mid SHC)$$

$$\xrightarrow{(R3)} \left( SR_{E1} \cdot R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot SR_{E1} \cdot R_{E2}^1 \cdot R_{I1} \cdot R_{I2} \cdot EGFR \cdot EGFR \right)^L \rfloor (SHC \mid SHC)$$

$$\xrightarrow{(R4)} \left( SR_{E1} \cdot R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot SR_{E1} \cdot R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot EGFR \cdot EGFR \right)^L \rfloor (SHC \mid SHC)$$

$$\xrightarrow{(R5)} \left( SR_{E1} \cdot R_{E2}^1 \cdot PR_{I1} \cdot R_{I2}^2 \cdot SR_{E1} \cdot R_{E2}^1 \cdot PR_{I1} \cdot R_{I2} \cdot EGFR \cdot EGFR \right)^L \rfloor (SHC^2 \mid SHC)$$

# Outline of the talk

# Current and future work

We developed a prototype simulator based on Stochastic CLS to run the *lac* operon example

- currently, we are developing a complete and efficient simulator

In order to model cell divisions and differentiations, tissues, etc...

- we are developing a spatial extension of CLS in which terms are placed and can move in a 2D/3D space

Moreover,

- we are developing a translation of Kohn Molecular Interaction Maps into CLS

As future work:

- we plan to develop a symbolic semantics of CLS, and a symbolic bisimulation relation to allow the development of a verification tool
- we plan to use CLS to study (in collaboration with biologists) retinal cell develpment and differentiation

# References

P. Milazzo. *Qualitative and Quantitative Formal Modeling of Biological Systems*, PhD Thesis, Università di Pisa.

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.

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.

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.

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.