



A simple calculus for proteins and cells

Cosimo Laneve, Fabien Tarissan

► To cite this version:

Cosimo Laneve, Fabien Tarissan. A simple calculus for proteins and cells. Theoretical Computer Science, 2008, Membrane Computing and Biologically Inspired Process Calculi, 404 (1-2), pp.127-141. 10.1016/j.tcs.2008.04.011 . hal-01217855

HAL Id: hal-01217855

<https://hal.science/hal-01217855>

Submitted on 20 Oct 2015

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

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

A simple calculus for proteins and cells

Cosimo Laneve, Fabien Tarissan,

Dipartimento di Scienze dell'Informazione, Università di Bologna

In memory of Nadia Busi

Abstract

The use of process calculi to represent biological systems has led to the design of different formalisms such as brane calculi and κ -calculus. Both have proved to be useful to model different types of biological systems. As an attempt to unify the formalisms, we introduce the **bio κ -calculus**, a simple calculus for describing proteins and cells, in which bonds are represented by means of shared names and interactions are modelled at the domain level. In **bio κ -calculus**, protein-protein interactions have to be at most binary and cell interactions have to fit with sort constraints.

In this contribution we define the semantics of **bio κ -calculus**, analyse its properties, discuss the expressivity of the calculus by modelling two significant examples – a signalling pathway and a virus infection –, and study an implementation in Milner's π -calculus.

1 Introduction

One problem when dealing with molecular biology is to extract a functional meaning out of the mass of current knowledge. This problem has pleaded for the development of specific tools that describe biology in a faithful way. Among these tools, process calculi have been proved powerful enough to formalise some of crucial activities of biological systems, to render in a natural way the massive parallelism and concurrency of interactions, and to analyse the overall behaviour.

There are two features of process calculi that are particularly relevant to biological systems. The first one is about the syntax. In process calculi, the syntax of the terms determines their capacity to interact with the environment. This is in strong analogy with molecular biology, where the structure of an object affects its activities. For instance, the three-dimensional structure of a protein, as well as the fact that it is part of an amino-acid sequence, may enable

some reactions and disable other ones. The second feature is about the semantics. Process calculi have models that have been thoroughly studied and that support automatic verification tools. These models may turn out essential in order to formalise the behaviour of biological systems (with qualitative measures, such as rate of reactions) and the reachability of problematic molecular configurations.

Two different process calculi families have brought some interest and various results. A first family is based on Milner’s π -calculus [1], following principles proposed in [2]. According to these principles, bonds between proteins are represented in π -calculus by shared channel names and protein interactions are modelled by process communications. The κ -calculus [3] belongs to this family and it has been showed convenient for representing mechanisms such as signalling pathways or regulatory networks. The second family of calculi relies on Mobile Ambient [4], following ideas proposed by Paun [5]. These calculi use actions and co-actions capabilities located on the surface of cell membranes for representing molecular transports and virus infections. The *brane calculi* [6] and *bio-ambients* [7] belong to this family. As these two families turn out to commit to different paradigms, it is compelling to develop a unique formalism able to handle the two types of systems. The challenge is to enucleate few basic realistic primitives that are expressive enough to cover mechanisms of κ -calculus and brane calculi. Such a formalism is the aim of this contribution.

The calculus presented in this contribution – the **bio** κ -calculus – retains denotations for proteins, cells, and solutions. Protein reactions are *complexations* and *decomplexations* of *two* proteins. These reactions follow the same pattern of those of κ -calculus. Actually, they are similar, but even simpler, to those of $\mathfrak{m}\kappa$ -calculus [3], an intermediate formalism introduced to ease the translation of κ -calculus in π -calculus. Cells are compartments consisting of a *membrane* and a *cytoplasm* – a solution contained into the membrane. Cellular reactions permit to fuse two cells. Fusions open the cell to other cells. In particular, the cytoplasms, whose interactions with the external solution are mediated by the membrane, may directly interact with another solution – the cytoplasm of the fusing cell – after the fusion. Complexations and decomplexations of a reactant that belongs to a cell membrane may change the capability of the membrane, thus preparing the cell to future fusions, endo- or exocytosis. The formal rendering of these phenomena is done in **bio** κ -calculus by admitting protein reactions with side-effects on membrane types. Other relevant cellular reactions are considered in **bio** κ -calculus, such as *translocations*, which transport proteic material inside the cells, and *phagocytosis*, which allows a cell, such as a virus, to enter other cells. Phagocytosis is particularly difficult to model because the entering cell is enclosed into a membrane that is pull out the host cell membrane. This means that the membrane of the host cell must have enough material for the new one. This mechanism is modelled by admitting matches of patterns of proteins.

The formal models of **bio** κ -calculus are labelled transition systems where labels carry information about the reactants and the biological rule that is being used. It is folklore in process calculi that such systems are too descriptive (too *intentional objects*) because they separate solutions that no experiment may distinguish. A number of equivalences over transition systems have been proposed in the literature. In this paper we consider *weak bisimulation* [8] that equates two systems if they simulate each other. We demonstrate that weak bisimulation is a congruence in **bio** κ -calculus: two weakly bisimilar biological systems behave in the same way when put in *every* solution.

The **bio** κ -calculus is simple. It has atomic elements – the proteins – and two syntactic constructors – grouping of solutions and cells. Reactions are interactions between two proteins or between two membranes. Notwithstanding this simplicity, the calculus is expressive enough. It is possible to describe significant biological examples: the RTK-MAPK pathway and a virus infection. This simplicity is also crucial for providing a faithful implementation of core **bio** κ -calculus into π -calculus. The compilation proposed in this paper encodes complexations and decomplexations into exactly one π -calculus interaction; conversely, every interaction corresponds exactly to one reaction.

The contribution is structured as follows. In the next section we define a core version of **bio** κ -calculus where reactions are restricted to be between two proteins. We study the semantics of core **bio** κ -calculus and analyse a signal transduction mechanism. In Section 3 we define an extensional semantics for **bio** κ -calculus and study its properties. In Section 4 we extend the basic interaction mechanisms with fusions that make structural modifications in the hierarchical organisation of a system. In Section 5 we discuss other extensions of the calculus for modelling translocations and phagocytosis. In Section 6 we study an implementation of core **bio** κ -calculus in π -calculus. Section 7 draws few concluding remarks and hints at possible future works.

2 The core of the **bio** κ -calculus

In this section we present a core version of the **bio** κ -calculus where interactions never change the hierarchical organisation of biological solutions. We define its syntax, its operational semantics, and the weak bisimulation. We also analyse its expressiveness by encoding the RTK-MAPK pathway.

2.1 Syntax

Three disjoint countable sets of names will be used in **bio** κ -calculus. A set of *protein names* P , ranged over by A, B, C, \dots to denote the *type* of a protein; a set of *edge names* E , ranged over by x, y, z, \dots to describe the bonds between proteins; a set of *membrane names* M , ranged over by M, N, \dots , to denote the type of a membrane. Protein names are sorted according to the number of *sites* they possess. Let $\mathfrak{s}(\cdot)$ be the function yielding the sites of proteins. The sites of a protein A are indicated by the natural numbers in the set $\{1, \dots, \mathfrak{s}(A)\}$.

Sites may be in three possible states: *bound* to another site, sharing then an element in E with this site, *visible*, i.e. not connected to other sites, or *hidden*, i.e. not connected to other sites but not useful for interactions. Let s be a site, we write s^x if it is connected through the name x , \bar{s} if it is hidden, and simply s if it is visible. Formally, the state of sites are defined by maps, called *interfaces*. Interfaces, ranged over by σ, σ', \dots , are partial functions from naturals to $E \cup \{v, h\}$ (we are assuming that $v, h \notin E$). For instance, $[1 \mapsto x; 2 \mapsto v; 3 \mapsto h]$ is an interface. In order to simplify the reading, we write this map as $1^x + 2 + \bar{3}$. In the following, when we write $\sigma + \sigma'$ we assume that the domains of σ and σ' are disjoint. We always assume that interfaces are injective on E . The sites of a protein A are completely defined by interfaces that are *total* on $\{1, \dots, \mathfrak{s}(A)\}$.

The following figure illustrates the all the possible state of protein domain and the notation for the complexes:

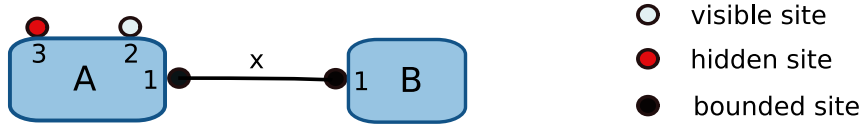


Fig. 1. Graphical notation for the complex $A(1^x + 2 + \bar{3}), B(1^x)$

Biological reactions modify interfaces, thus altering proteins' capabilities. In addition to the creation or suppression of a bond between proteins, reactions may also modify sites that are either hidden or visible to visible or hidden, respectively. In order to model these changes, let a *v-h-map*, ranged over ϕ, ψ, \dots , be every partial interface onto $\{v, h\}$. Let the *swap* of a *v-h-map* ϕ , written $\bar{\phi}$, be the *v-h-map*:

$$\bar{\phi}(i) = \begin{cases} h & \text{if } \phi(i) = v \\ v & \text{if } \phi(i) = h \end{cases}$$

Definition 1 *The syntax of the core-bio* κ -calculus defines (biological) solu-

tions S by the grammar:

$$S, T ::= \begin{array}{cccc} \mathbf{0} & | & A(\sigma) & | & M(S)[T] & | & S, T \\ \text{(empty)} & & \text{(protein)} & & \text{(cell)} & & \text{(group)} \end{array}$$

Solutions can be either empty, or a *protein* $A(\sigma)$ indicating a protein name and its interface, or a *cell* $M(S)[T]$, that is a solution T , called *cytoplasm*¹, surrounded by another solution S , called *membrane*, or a *group* of solutions S , T . The following figure illustrates the syntactical notation for a cell:

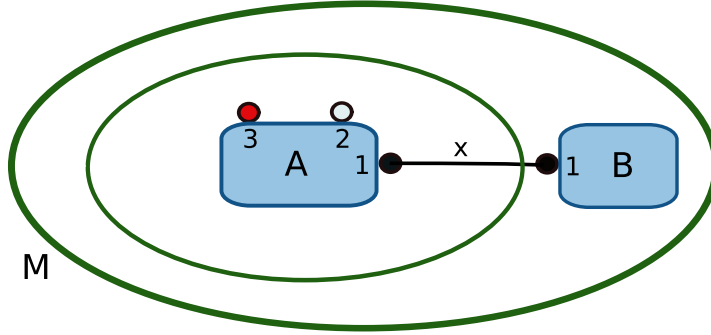


Fig. 2. Graphical notation for the cell $M(B(1^x))[A(1^x + 2 + \bar{3})]$

Three auxiliary functions will be applied to solutions and interfaces:

- $\text{en}(\cdot)$ returns the set of *edge names* occurring in the argument;
- $\text{de}(\cdot)$ returns the set of *dangling edge names* of the argument, namely those names that occur exactly once;
- $\text{be}(\cdot)$ returns the set of *bound edge names* of the argument, namely those names that occur at least twice.

Clearly $\text{de}(S) = \text{en}(S) \setminus \text{be}(S)$ and similarly for σ . For instance, in $S = M(C(1^y + 2))[A(1^x + 2 + 3), B(1 + 2^x)]$ the set $\text{en}(S)$ is $\{y, x\}$ and the set $\text{de}(S)$ is $\{y\}$. We abbreviate the group $A_1(\sigma_1), \dots, A_n(\sigma_n)$ with $\prod_{i \in 1..n} A_i(\sigma_i)$.

In the whole paper, we identify solutions that are equal up to a renaming of edge names that are not dangling (called *alpha-conversion*) and we assume that all solutions meet the following *well-formedness conditions*:

- (*edge-condition*) in every solution, edge names occur at most twice;
- (*membrane-condition*) every membrane is a group of proteins, that is cells do not occur in membranes;

¹ With the term “cytoplasm” we refer to every material surrounded by a structure forbidding any interaction with the external environment. Thus, a nucleus of a cell is a cytoplasm, as well as a DNA strand surrounded by a capsid.

- (*cytoplasm-condition*) the dangling edges of nuclei of cells are connected to the corresponding membrane, that is, for every $M(S)[T]$, $\text{de}(T) \subseteq \text{de}(S)$.

For example $M(C(1^x + 2))[A(1^x + 2 + 3), B(1 + 2^x)]$ does not meet the edge condition because the edge x has three ends (it is a multi-edge). The solution $M(B(1 + 2), C(1 + 2))[A(1 + 2^x + 3)]$ does not meet the cytoplasm condition because the cytoplasm $A(1 + 2^x + 3)$ has a dangling edge that is not connected to the membrane. In the following, solutions that are membranes will be addressed by M, N, \dots .

2.2 Semantics

Biological reactions that we consider in this section are of two types: *complexations*, which create edges between possibly disconnected proteins, and *decomplexations*, which remove edges. For instance a complexation reaction is (we are assuming $\mathfrak{s}(A) = 4$, $\mathfrak{s}(B) = 3$)

$$A(1^x + 2 + \bar{3} + 4), B(1 + \bar{2} + \bar{3}) \longrightarrow A(1^x + 2^y + 3 + 4), B(1^y + 2 + \bar{3})$$

that creates an edge y connecting the site 2 of A and the site 1 of B . These two sites, in order that this reaction be executed, must be visible. This means that the application of a complexation must check whether the sites being connected are visible or not. For example the above reaction cannot be applied to the group $A(1^x + \bar{2} + \bar{3} + 4), B(1 + \bar{2} + \bar{3})$ because the site 2 of A is hidden. Reactions in *bio κ* -calculus may also change the state of sites that are visible or hidden in the reactants, switching them into hidden and visible, respectively. In the example above, this happens to the sites 3 of A and 2 of B . This partial interface is then used as a test to check whether the interactions is possible and at the same time, it is used to describe the effect of the interaction on the internal sites of the proteins. These serve however two different purposes. Then it seems natural to add the possibility to describe a part of the interface used as a test but keeping unchanged by the interaction. This is the case of the site 4 of A if we suppose that the interaction is impossible as soon as this site is hidden. A concise way for defining the above reaction is the schema

$$r : ((A, 2, \bar{3}, 4), (B, 1, \bar{2}, \emptyset))$$

that makes explicit the reactant proteins – the first items of the triples –, the corresponding sites to be complexated – the second items – and the part of the interface which is tested, separating the sub-part which is permuted from the one which keeps unchanged – third and fourth item. For example, the rule r also applies to $A(1 + 2 + \bar{3} + 4), B(1 + \bar{2} + 3)$ or $A(\bar{1} + 2 + \bar{3} + 4), B(1 + \bar{2} + 3^x)$ yielding solutions $A(1 + 2^y + 3 + 4), B(1^y + 2 + 3)$ and $A(\bar{1} + 2^y + 3 + 4), B(1^y +$

$2 + 3^x$), respectively. In general, the shape of a reaction schema is

$$\mathbf{r} : ((A, i, \psi_1, \psi_2), (B, j, \phi_1, \phi_2))$$

that is a reaction name \mathbf{r} and two tuple containing a protein name, a site and two v - h -maps. A generic *application* of the schema \mathbf{r} may be written as

$$\begin{array}{c} A(i + \psi_1 + \psi_2 + \psi') , B(j + \phi_1 + \phi_2 + \phi') \\ \downarrow \\ A(i^x + \overline{\psi_1} + \psi_2 + \psi') , B(j^x + \overline{\phi_1} + \phi_2 + \phi') \end{array}$$

where x is a fresh edge name and $i + \psi_1 + \psi_2 + \psi$ and $j + \phi_1 + \phi_2 + \phi$ are total on $[1 .. \mathfrak{s}(A)]$ and $[1 .. \mathfrak{s}(B)]$, respectively. It is worth to observe that the interfaces σ and σ' are not mentioned in the schema. This means that the interaction occurs *what ever* the states of the sites of these interfaces are.

Decomplexations are complexations in the other way round. For instance a decomplexation reaction is

$$A(1^x + 2^y + 3 + 4) , B(1^y + 2 + \overline{3}) \longrightarrow A(1^x + 2 + \overline{3} + 4) , B(1 + \overline{2} + \overline{3})$$

that removes the edge y . The schema describing decomplexations is similar to that of complexations:

$$\mathbf{r}' : ((A, i, \psi_1, \psi_2), (B, j, \phi_1, \phi_2))$$

The application of the decomplexation rule is different from complexation: in this case the two reactants must be connected by an edge between the site i of A and the site j of B . So, for example, a generic application of \mathbf{r}' is

$$\begin{array}{c} A(i^x + \psi_1 + \psi_2 + \psi') , B(j^x + \phi_1 + \phi_2 + \phi') \\ \downarrow \\ A(i + \overline{\psi_1} + \psi_2 + \psi') , B(j + \overline{\phi_1} + \phi_2 + \phi') \end{array}$$

In order to separate complexations from decomplexations we consider two functions, \mathcal{C} for *complexations* and \mathcal{D} for *decomplexations*, from rule names to tuples $((A, i, \phi_1, \phi_2), (B, j, \psi_1, \psi_2))$. These functions \mathcal{C} and \mathcal{D} are assumed with disjoint domains, therefore a rule name uniquely defines whether it is a complexation or a decomplexation. Let \mathcal{R} range over \mathcal{C} and \mathcal{D} and let also write $(A, i, \phi_1, \phi_2) \in_{\text{left}} \mathcal{R}(\mathbf{r})$ if $\mathcal{R}(\mathbf{r}) = ((A, i, \phi_1, \phi_2), (B, j, \psi_1, \psi_2))$ and $(A, i, \phi_1, \phi_2) \in_{\text{right}} \mathcal{R}(\mathbf{r})$ if $\mathcal{R}(\mathbf{r}) = ((B, j, \psi_1, \psi_2), (A, i, \phi_1, \phi_2))$. Finally, let also $(A, i, \phi_1, \phi_2) \in \mathcal{R}(\mathbf{r})$ if either $(A, i, \phi_1, \phi_2) \in_{\text{left}} \mathcal{R}(\mathbf{r})$ or $(A, i, \phi_1, \phi_2) \in_{\text{right}} \mathcal{R}(\mathbf{r})$.

The operational semantics of core **bio** κ -calculus we are going to define use labels. A label is either a triple (A, x, \mathbf{r}) , also written $A_{\mathbf{r}}^x$, where A is a protein name, x an edge name and \mathbf{r} a rule name, or τ . We use μ to range over labels. Let $\text{diff}(S, S')$ be the set $(\text{en}(S') \setminus \text{en}(S)) \cup (\text{en}(S) \setminus \text{en}(S'))$.

Definition 2 *The transition relation $S \xrightarrow{\mu} T$ is the least relation satisfying the following rules:*

$$\begin{array}{c}
\text{(COM)} \\
\frac{(A, i, \phi_1, \phi_2) \in \mathcal{C}(\mathbf{r}) \quad \mathbf{x} \notin \text{en}(\sigma)}{A(i + \phi_1 + \phi_2 + \sigma) \xrightarrow{A_{\mathbf{r}}^x} A(i^x + \overline{\phi_1} + \phi_2 + \sigma)} \\
\text{(DEC)} \\
\frac{(A, i, \phi_1, \phi_2) \in \mathcal{D}(\mathbf{r})}{A(i^x + \phi_1 + \phi_2 + \sigma) \xrightarrow{A_{\mathbf{r}}^x} A(i + \overline{\phi_1} + \phi_2 + \sigma)} \\
\text{(SOL-L)} \qquad \text{(SOL-R)} \\
\frac{S \xrightarrow{\mu} S' \quad \text{diff}(S, S') \cap \text{en}(T) = \emptyset}{S, T \xrightarrow{\mu} S', T} \qquad \frac{T \xrightarrow{\mu} T' \quad \text{diff}(T, T') \cap \text{en}(S) = \emptyset}{S, T \xrightarrow{\mu} S, T'} \\
\text{(MEM)} \qquad \text{(CYTO)} \\
\frac{M \xrightarrow{\mu} M' \quad \text{diff}(M, M') \cap \text{en}(S) = \emptyset}{M(\lfloor M \rfloor)[S] \xrightarrow{\mu} M(\lfloor M' \rfloor)[S]} \qquad \frac{S \xrightarrow{\tau} S' \quad \text{diff}(S, S') \cap \text{en}(M) = \emptyset}{M(\lfloor M \rfloor)[S] \xrightarrow{\tau} M(\lfloor M \rfloor)[S']} \\
\text{(REACT)} \qquad \text{(MS-REACT)} \\
\frac{S \xrightarrow{A_{\mathbf{r}}^x} S' \quad T \xrightarrow{B_{\mathbf{r}}^x} T' \quad A \neq B}{S, T \xrightarrow{\tau} S', T'} \qquad \frac{M \xrightarrow{A_{\mathbf{r}}^x} M' \quad S \xrightarrow{B_{\mathbf{r}}^x} S' \quad A \neq B}{M(\lfloor M \rfloor)[S] \xrightarrow{\tau} M(\lfloor M' \rfloor)[S']}
\end{array}$$

Rules (COM) and (DEC) respectively define complexations and decomplexations capabilities of proteins by lifting these information to labels of transitions and, at the same time, updating the proteins. Rules (SOL-L), (SOL-R) and (MEM) lift transitions to groups and membranes; it is crucial that edge names created or deleted do not occur elsewhere. Rule (CYTO) lift *internal* transitions of the cytoplasm to the whole cell; as before, edge names created or deleted must not occur elsewhere. We observe that (CYTO) bans complexation or decomplexation between cytoplasm and the solution external to the cells. Rule (REACT) and (MS-REACT) define reactions, both complexations and decomplexations, in groups and cells. In particular (REACT) also accounts for reactions between membranes of different cells. It is worth to notice that the constraint $A \neq B$ bans reactions between same proteins (*cf. self complexation* in [3]). This is for simplicity sake: in case reactants are proteins with a same name we need to carry more information on the labels to separate them.

It is worth to notice that the semantics of **bio** κ -calculus adhere to a style that is different than the one used for defining the transitions of **m** κ -calculus in [3]. The **m** κ -calculus has been equipped with an *unlabelled* reduction relation \rightarrow . This relation corresponds to the foregoing rules (COM), (DEC), (SOL-G), (SOL-D), and (REACT). We have preferred a labelled transition relation for simplicity sake: the number of rules for defining a relation \rightarrow equivalent to $\xrightarrow{\mu}$ should have been larger than that of Definition 2 because of the presence of membranes.

The transition relation preserve well-formedness of solutions.

Proposition 1 *If S is well-formed and $S \xrightarrow{\mu} T$ then T is well-formed as well.*

It is worth to observe that membrane names do not play any role at this stage. They will be relevant in the complete system with complex membrane reactions presented in Section 4.

2.3 Modelisation of the signal transduction

The so-called RTK-MAPK pathway are intensely used and studied in many approaches modelling and simulating biological systems [2,3]. We therefore model the first steps of such a pathway in **bio** κ -calculus, thus providing a touchstone for our calculus.

The *receptors tyrosine kinase* (RTKs) are receptors located at the surface of some cell membranes. After binding with some ligand (insulin, EGF, VEGF, etc ...), they are activated and lead to a specific cellular response (growth of the cell, proliferation, etc ...). There exists actually many different types of RTKs which are gathered into different families. They don't provide the same reaction and the cellular response depends on the type of RTK which is activated. After the binding with a ligand, the RTKs may *self-phosphorylate* (adding a phosphate group to a binding domain). This modification is detected by specific proteins which will bind to the RTK and recruit other proteins which will be activated or inhibited and will lead eventually to the cellular response. Another possibility for the RTKs to mediate their reaction is to phosphorylate other proteins (besides self-phosphorylation). These activated proteins lead to a cascade of signalisations too. This cascade results mainly in the activation or the inhibition of transcription factors which adjusts gene expressions. The generic schema of an RTK-MAPK cascade is then the following: the RTKs, after binding with a ligand (the signal), induce some changes into the cell (ex. gene expression) in order to mediate a biological effect (ex. entering the cell cycle).

The example we are developing here is that of the response induced by the

signal carried by an epidermic growth factor protein (EGF). The dimeric form (1) of EGF can bind to its associated receptor EGFR (2), a transmembrane protein with an extracellular ligand-binding domain located on the plasmic membrane of some cells. This binding activates EGFR by phosphorylating an internal domain of the protein (3 and 4). This activation leads to multiple interactions with cytoplasmic complexes of proteins by successive binding-phosphorylations, starting with the adapter protein SHC (5). The cascade of interactions ends with the activation of the extracellular signal-regulated kinase ERK. Once phosphorylated, this protein can then be transported into the nucleus by means of a translocation mechanism and will thereafter modify the gene expression, stimulating the cell to enter mitosis. This causes the cell to divide and proliferate.

After the biological description EGF, RTK, and SHC have respective arities 3, 4, and 2. We give here the formal rendering of the five first steps described above:

$$\begin{aligned}
r_1 : & \quad ((\text{EGF}, 1, \overline{2}, \emptyset), (\text{EGF}, 1, \overline{2}, \emptyset)) \in \mathcal{C} \\
r_2 : & \quad ((\text{EGF}, 2, \emptyset, \emptyset), (\text{EGFR}, 1, \overline{4}, \emptyset)) \in \mathcal{C} \\
r_3 : & \quad ((\text{EGFR}, 2, \overline{3} + 4, \emptyset), (\text{EGFR}, 2, \overline{3} + 4, \emptyset)) \in \mathcal{C} \\
r_4 : & \quad ((\text{EGFR}, 2, \emptyset, \emptyset), (\text{EGFR}, 2, \emptyset, \emptyset)) \in \mathcal{D} \\
r_5 : & \quad ((\text{EGFR}, 3, \emptyset, \emptyset), (\text{SHC}, 1, \overline{2}, \emptyset)) \in \mathcal{C}
\end{aligned}$$

A possible run of those rules is presented in Figure 3, showing that our language is expressive enough to define the inherent causality involved in the transduction in a precise yet natural way.

A couple of problems of our notation deserve to be discussed though. Consider the initial solution:

$$\begin{aligned}
& \text{EGF}(1 + \overline{2}) , \text{EGF}(1 + \overline{2}) , \text{EGF}(1 + \overline{2}) , \text{EGF}(1 + \overline{2}) , \\
& \mathbf{M} \langle \text{EGFR}(1 + 2^x + \overline{3} + \overline{4}) , \text{EGFR}(1 + 2^x + \overline{3} + \overline{4}) , \mathbf{M} \rangle [\text{SHC}(1 + \overline{2}) , \mathbf{S}]
\end{aligned}$$

After two applications of rule r_1 , we obtain the solution

$$\begin{aligned}
& \text{EGF}(1^x + 2) , \text{EGF}(1^x + 2) , \text{EGF}(1^y + 2) , \text{EGF}(1^y + 2) , \\
& \mathbf{M} \langle \text{EGFR}(1 + 2^x + \overline{3} + \overline{4}) , \text{EGFR}(1 + 2^x + \overline{3} + \overline{4}) , \mathbf{M} \rangle [\text{SHC}(1 + \overline{2}) , \mathbf{S}]
\end{aligned}$$

that reduces, after two application of rule r_2 to the wrong solution

$$\begin{aligned}
& \text{EGF}(1^x + 2) , \text{EGF}(1^x + 2^u) , \text{EGF}(1^y + 2^v) , \text{EGF}(1^y + 2) , \\
& \mathbf{M} \langle \text{EGFR}(1^u + 2^x + \overline{3} + 4) , \text{EGFR}(1^v + 2^x + \overline{3} + 4) , \mathbf{M} \rangle [\text{SHC}(1 + \overline{2}) , \mathbf{S}]
\end{aligned}$$

$$\begin{aligned}
& \text{EGF}(1 + \overline{2}) , \text{EGF}(1 + \overline{2}) , \\
& \quad \mathsf{M}(\text{EGFR}(1 + 2 + \overline{3} + \overline{4}) , \text{EGFR}(1 + 2 + \overline{3} + \overline{4}) , \mathsf{M})[\text{SHC}(1 + \overline{2}) , \mathsf{S}] \\
& \xrightarrow{\tau} \text{EGF}(1^z + 2) , \text{EGF}(1^z + 2) , \\
& \quad \mathsf{M}(\text{EGFR}(1 + 2 + \overline{3} + \overline{4}) , \text{EGFR}(1 + 2 + \overline{3} + \overline{4}) , \mathsf{M})[\text{SHC}(1 + \overline{2}) , \mathsf{S}] \quad (\mathbf{r}_1) \\
& \xrightarrow{\tau} \text{EGF}(1^z + 2^y) , \text{EGF}(1^z + 2) , \\
& \quad \mathsf{M}(\text{EGFR}(1^y + 2 + \overline{3} + \overline{4}) , \text{EGFR}(1 + 2 + \overline{3} + \overline{4}) , \mathsf{M})[\text{SHC}(1 + \overline{2}) , \mathsf{S}] \quad (\mathbf{r}_2) \\
& \xrightarrow{\tau} \text{EGF}(1^z + 2^y) , \text{EGF}(1^z + 2^u) , \\
& \quad \mathsf{M}(\text{EGFR}(1^y + 2 + \overline{3} + \overline{4}) , \text{EGFR}(1^u + 2 + \overline{3} + \overline{4}) , \mathsf{M})[\text{SHC}(1 + \overline{2}) , \mathsf{S}] \quad (\mathbf{r}_2) \\
& \xrightarrow{\tau} \text{EGF}(1^z + 2^y) , \text{EGF}(1^z + 2^u) , \\
& \quad \mathsf{M}(\text{EGFR}(1^y + 2^x + 3 + \overline{4}) , \text{EGFR}(1^u + 2^x + 3 + \overline{4}) , \mathsf{M})[\text{SHC}(1 + \overline{2}) , \mathsf{S}] \quad (\mathbf{r}_3) \\
& \xrightarrow{\tau} \text{EGF}(1^z + 2^y) , \text{EGF}(1^z + 2^u) , \\
& \quad \mathsf{M}(\text{EGFR}(1^y + 2 + 3 + \overline{4}) , \text{EGFR}(1^u + 2 + 3 + \overline{4}) , \mathsf{M})[\text{SHC}(1 + \overline{2}) , \mathsf{S}] \quad (\mathbf{r}_4) \\
& \xrightarrow{\tau} \text{EGF}(1^z + 2^y) , \text{EGF}(1^z + 2^u) , \\
& \quad \mathsf{M}(\text{EGFR}(1^y + 2 + 3^t + \overline{4}) , \text{EGFR}(1^u + 2 + 3 + \overline{4}) , \mathsf{M})[\text{SHC}(1^t + 2) , \mathsf{S}] \quad (\mathbf{r}_5)
\end{aligned}$$

Fig. 3. The first steps of the RTK-MAPK cascade.

where two different dimeric forms of EGF connect to a same pair of EGFR receptors. Our notation is too simple to rule out such configurations. In \mathbf{mk} -calculus, this expressiveness issue is solved by the use of reaction IDs and pattern matching over them. Actually, this issue is related to a more general question named self-assembly problem and worth to be studied independently [9].

The second problem is manifested at the end of the RTK-MAPK pathway. The pathway causes a phosphorylation of a particular protein (ERK) that enters in the nucleus, which is represented as cell, as well. At this stage we have no mechanism that make entities enter in the cell. Such mechanisms will be discussed in detail in Section 5.

3 Extensional semantics: weak bisimulation

The transition relation of Definition 2 associates solutions to graphs where nodes are solutions and μ -labelled edges model the transitions $\mathsf{S} \xrightarrow{\mu} \mathsf{T}$. The induced equivalence on $\mathbf{bio}\kappa$ -calculus is graph isomorphism: two terms are equivalent provided their associated graphs are isomorphic. Graph isomorphism is however too strong as a biological semantics because it dis-

tinguishes solutions that differ for τ -transitions. For instance, let $\mathcal{C}(\mathbf{r}) = ((A, 1, \emptyset), (B, 1, \emptyset)) = \mathcal{D}(\mathbf{r}')$, that is \mathbf{r} and \mathbf{r}' are *reversible* reactions. Then the solutions $A(1^y + \sigma), B(1^y + \sigma')$ and $A(1 + \sigma), B(1 + \sigma')$ have underlying graphs that are not isomorphic. The failure of the isomorphism is due to τ -transitions: every solution may be reduced to the other by means of a τ -transition.

The following equivalence corrects the above criticism. It adapts *weak bisimulation* [1] to our setting. We write $S \xRightarrow{\tau} S'$ if $S \xrightarrow{\tau}^* S'$ and $S \xRightarrow{\mu} S'$, with $\mu \neq \tau$, if $S \xrightarrow{\tau}^* \xrightarrow{\mu} \xrightarrow{\tau}^* S'$.

Definition 3 A (weak) bisimulation is a symmetric binary relation \mathfrak{R} between solutions such that $S \mathfrak{R} T$ implies:

- (1) if $S \xrightarrow{\tau} S'$ then $T \xRightarrow{\tau} T'$ and $S' \mathfrak{R} T'$;
- (2) if $S \xrightarrow{A_r^x} S'$ then $T \xRightarrow{A_r^x} T'$ and $S' \mathfrak{R} T'$.

S is bisimilar to T , written $S \approx T$, if $S \mathfrak{R} T$ for some bisimulation \mathfrak{R} .

With this notion of equivalence, it is possible to show the following properties for the core-bio κ -calculus:

Proposition 2 (1) “ \approx ” is an abelian monoidal operator with identity $\mathbf{0}$. Namely $S, T \approx T, S$ and $(S, T), R \approx S, (T, R)$ and $S, \mathbf{0} \approx S$.
(2) \approx is preserved by injective renamings that are identities on dangling edge names. Namely, let ι be an injective renaming on $\text{en}(S)$ such that ι is the identity on $\text{de}(S)$, then $S \approx \iota(S)$.
(3) \approx is preserved by reversible rules. Namely, let $\mathcal{C}(\mathbf{r}) = ((A, i, \psi), (B, j, \phi))$ and $\mathcal{D}(\mathbf{r}') = ((A, i, \bar{\psi}), (B, j, \bar{\phi}))$ then $A(i + \psi + \sigma), B(j + \phi + \sigma') \approx A(i^x + \bar{\psi} + \sigma), B(j^x + \bar{\phi} + \sigma')$.

Another relevant property of \approx is that every two bisimilar systems behave in the same way when plugged in a same context.

Theorem 1 \approx is a congruence.

Proof: We must prove that, if $S \approx T$ then

- (1) $S, R \approx T, R$,
- (2) $M(S)[R] \approx M(T)[R]$
- (3) $M(M)[S] \approx M(M)[T]$

when these solutions are well-formed. We demonstrate (1), the other cases are similar. Let \mathfrak{R} be the bisimulation containing the pair (S, T) and let $(S', R) \mathfrak{R}' (T', R)$ if

(1) $S' \mathfrak{R} T'$,

(2) and S', R and T', R are well-formed.

To demonstrate that \mathfrak{R}' is a bisimulation one has to prove that if $S', R \xrightarrow{\mu} U$ then there exists U' such that $T', R \xrightarrow{\mu} U'$ and $U \mathfrak{R}' U'$. There are two cases:

$\mu = A_r^x$. Then either $U = S', R'$ with $R \xrightarrow{A_r^x} R'$ or $U = S'', R$ with $S' \xrightarrow{A_r^x} S''$. The first case is immediate. In the second case, $T' \xrightarrow{A_r^x} T''$ and $S' \mathfrak{R} T''$. By definition: $(S'', R) \mathfrak{R}' (T'', R)$.

$\mu = \tau$. The interesting case is when S' and R interact (the other cases are dealt as above) by the rule (REACT). Suppose $U = S'', R'$ with $S' \xrightarrow{A_r^x} S''$ and $R \xrightarrow{B_r^x} R'$. Since $S' \mathfrak{R} T'$, there exists T'' such that $T' \xrightarrow{A_r^x} T''$ and $S' \mathfrak{R} T''$. Then, by (REACT), $T', R \xrightarrow{\tau} T'', R'$ and $(S'', R') \mathfrak{R}' (T'', R')$ by definition of \mathfrak{R}' .

The substitution property of Proposition 2 owned by \approx might be too strong in biology, thus making this equivalence not realistic. In this context, one usually wants to prove that two parts behave in the same way when plugged *under a certain number of contexts*, rather than every possible one. Therefore, a semantics owning a parametric congruence property might fit better with biology. However what these parameters are and what are the properties owned by a “good” extensional semantics for biology remains unclear to us and is left as an open issue.

We finish this section by a series of remarks concerning \approx . It is worth noticing for example that $S \approx T$ does not imply $\text{de}(S) = \text{de}(T)$. For two reasons: First of all, by bisimulation, $S \xrightarrow{A_r^x} S'$ may be matched by $T \xrightarrow{A_r^y} T'$ with $x \neq y$. Second, taking empty biological relations – i.e. $\mathcal{C} = \emptyset$ and $\mathcal{D} = \emptyset$ –, then $A(1^x) \approx \mathbf{0}$ but their dangling edges are different. Nevertheless, one might establish a relation between subsets of dangling edges of two bisimilar solutions. Let $\text{oe}(S)$, called the *observable edges*, be the set

$$\text{oe} = \{x \mid S \xrightarrow{A_r^x} S' \text{ and } r \text{ in domain of } \mathcal{D}\}$$

It is easy to prove that if $S \approx T$ then there is an injective renaming ι such that $\text{oe}(S) = \iota(\text{oe}(T))$.

Finally, as an illustration of Definition 3, one may note that a cell whose membrane cannot interact with elements outside is equivalent to the empty solution. Let M be a group of proteins. M is said to be *inert* if $M \not\xrightarrow{A_r^x}$ for every A_r^x . It is possible to verify that, if M is inert, then $M(\mathbf{M})[S] \approx \mathbf{0}$.

4 Cell interactions

The calculus of Section 2 is not very different from $\mathfrak{m}\kappa$ -calculus. Cells, in particular, do not play any relevant role since their structure is preserved by the transition. In this section we explore an extension with brane primitives, thus being able to model merging and splitting of cells such as the following *endosomes fusion*:

$$\text{ESM}(\langle \mathbf{M} \rangle [\mathbf{S}]) , \text{ESM}(\langle \mathbf{N} \rangle [\mathbf{T}]) \xrightarrow{\tau} \text{ESM}(\langle \mathbf{M} , \mathbf{N} \rangle [\mathbf{S} , \mathbf{T}])$$

The extension of core $\text{bio}\kappa$ -calculus we are going to design retains higher-order mechanisms, following *higher-order π -calculus*, a similar extension already studied for π -calculus [10].

4.1 Mreagents

We begin by augmenting the syntax with membrane reagents:

Definition 4 *The syntax of the $\text{bio}\kappa$ -calculus defines the solutions S by the grammar:*

$$\begin{array}{ccccccc} \mathbf{M}, \mathbf{S}, \mathbf{T} ::= & \mathbf{0} & | & A(\sigma) & | & \mathbf{M}(\langle \mathbf{M} \rangle [\mathbf{S}]) & | & \mathbf{S} , \mathbf{T} & | & \langle \mathbf{M} ; \mathbf{S} \rangle \cdot \mathbf{T} \\ & (empty) & & (protein) & & (cell) & & (group) & & (mreagent) \end{array}$$

An mreagent $\langle \mathbf{M} ; \mathbf{S} \rangle \cdot \mathbf{T}$ is an intermediate (unstable) solution that is used for manifesting the capability of a membrane \mathbf{M} , isolating a solution \mathbf{S} from the environment \mathbf{T} , to fuse with another membrane. This remark leads to consider the elements \mathbf{M} and \mathbf{S} as the components of a cell $\mathbf{M}(\langle \mathbf{M} \rangle [\mathbf{S}])$ and, as such, the following well-formedness constraints are put on their structure:

- \mathbf{M} is a multiset of proteins;
- $\text{de}(\mathbf{S}) \subseteq \text{de}(\mathbf{M})$, namely the dandling edges of \mathbf{S} are connected to proteins located in the membrane \mathbf{M} ;
- mreagents do not contain other mreagents.

Two operations involve membranes: *fusions*, which put together in a unique cell two cytoplasms separated by two membranes, and *activations* which change the type of a membrane as a consequence of a complexation or a decomplexation. Fusions are defined by a function \mathcal{F} from rule names to triples $((\mathbf{M}, \mathbf{M}'), \mathbf{N})$. We write $(\mathbf{M} \otimes \mathbf{M}', \mathbf{N}) = \mathcal{F}(\mathbf{r})$ if either $\mathcal{F}(\mathbf{r}) = ((\mathbf{M}, \mathbf{M}'), \mathbf{N})$ or $\mathcal{F}(\mathbf{r}) = ((\mathbf{M}', \mathbf{M}), \mathbf{N})$. We also write $\mathbf{M} \in \mathcal{F}(\mathbf{r})$ if $(\mathbf{M} \otimes \mathbf{M}', \mathbf{N}) = \mathcal{F}(\mathbf{r})$, for some \mathbf{M}' and \mathbf{N} . We assume that the domains of \mathcal{F} , \mathcal{C} and \mathcal{D} are disjoint. Activations

are defined by a function \mathcal{A} that takes pairs (A_r, M) and returns membrane names. The intuition is that a reaction involving a transmembrane protein may activate the membrane it belongs to, thus preparing the cell to a possible fusion.

With abuse of notation, we use μ to also range over labels of the form M_r .

Definition 5 *The transition relation $\xrightarrow{\mu}$ is the least one that includes the rules in Definition 2 where (MEM) and (MS-REACT) have also the premise “ (A_r, M) not in the domain of \mathcal{A} ” and the following reductions:*

$$\begin{array}{c}
\text{(OPEN)} \\
\frac{M \in \mathcal{F}(r)}{M \langle M \rangle [S] \xrightarrow{M_r} \langle M ; S \rangle \cdot 0}
\end{array}
\qquad
\begin{array}{c}
\text{(GRASP)} \\
\frac{S \xrightarrow{\mu} \langle M ; S'' \rangle \cdot S'}{S, T \xrightarrow{\mu} \langle M ; S'' \rangle \cdot (S', T)}
\end{array}$$

$$\begin{array}{c}
\text{(FUSE-H)} \\
\frac{S \xrightarrow{M_r} \langle M ; S'' \rangle \cdot S' \quad T \xrightarrow{N_r} \langle N ; T'' \rangle \cdot T' \quad \mathcal{F}(r) = (M \otimes N, M')}{S, T \xrightarrow{\tau} S', T', M' \langle M, N \rangle [S'', T'']}
\end{array}
\qquad
\begin{array}{c}
\text{(FUSE-V)} \\
\frac{S \xrightarrow{N_r} \langle N ; T \rangle \cdot S' \quad \mathcal{F}(r) = (M \otimes N, M')}{M \langle M \rangle [S] \xrightarrow{\tau} M' \langle M, N \rangle [S'], T}
\end{array}$$

$$\begin{array}{c}
\text{(MEM-A)} \\
\frac{M \xrightarrow{A_r^x} M' \quad \mathcal{A}(A_r, M) = N \quad \text{diff}(M, M') \cap \text{en}(S) = \emptyset}{M \langle M \rangle [S] \xrightarrow{A_r^x} N \langle M' \rangle [S]}
\end{array}
\qquad
\begin{array}{c}
\text{(MS-AREACT)} \\
\frac{M \xrightarrow{A_r^x} M' \quad S \xrightarrow{B_r^x} S' \quad A \neq B \quad \mathcal{A}(A_r, M) = N}{M \langle M \rangle [S] \xrightarrow{\tau} N \langle M' \rangle [S']}
\end{array}$$

Rule (OPEN) prepares a cell to be fused with an enclosing cell or with a peer cell; the precondition guarantees that the cell may participate to a fusion. Rule (GRASP) lifts the fusion capability to groups by freezing them in a mreagent. Rules (FUSE-H) and (FUSE-V) define fusions between peer and nested cells, respectively. The new cell is created with the membrane name returned by the function \mathcal{F} . Rule (MEM-A) is a refinement of (MEM). It models possible side-effects on the membrane name due to interactions between membrane proteins and proteins outside the cell. Such interactions may activate membranes by changing their fusion capability, which is encoded in our formalism by membrane names. In a similar way, rule (MS-AREACT) refines (MS-REACT).

4.2 Viral infection

A virus is an intracellular parasite that uses the infected cell replication machinery in order to duplicate its own genetic material. Usually, a virus consists of a genetic material (DNA or RNA), a proteic structure called *capsid* providing protection for the genetic material – we use the term of nucleocapsid to denote both the capsid and the genetic material –, and a possible envelope (usually extracted from a former infected cell and used later on to infect other cells).

Below we encode in **bio κ** -calculus an influenza-like virus relying on similar descriptions in [6,11]. We focus on the infection part, as we cannot express any creation of new material. This part consists of the following steps:

- (1) a protein-protein interaction between a virus membrane protein – the hemagglutinin HA – and a receptor – a glycoprotein GLY – on the cell's membrane; this activates GLY and prepares the cell to the phagocytosis;
- (2) the phagocytosis of the virus occurs thus creating a new vesicle VES engulfing the virus;
- (3) a fusion occurs between the new vesicle and an endosomal membrane (EDSM) in the cytoplasm; this makes the virus enter the endosome;
- (4) a further fusion occurs between the endosome and the virus that is now part of the cytoplasm; this leads to an exocytosis that eventually releases the virus nucleocapsid into the cytoplasm.

We consider the following rules:

$$\mathbf{r}_3 : ((\text{EDSM}, \text{VES}), \text{EDSM}) \in \mathcal{F}$$

$$\mathbf{r}_4 : ((\text{EDSM}, \text{VS}), \text{EDSM}) \in \mathcal{F}$$

The initial solution is **Virus** , **Cell** where the components are as follows:

$$\text{Virus} := \text{vs}(\llbracket \text{HA}(1) \rrbracket) [\text{Nucaps}]$$

$$\text{Cell} := \text{cLL}(\llbracket \text{GLY}(1) \rrbracket, \text{M}_c) [\text{Endosome} , \text{Cytosol}]$$

$$\text{Endosome} := \text{EDSM}(\llbracket \text{M}_e \rrbracket) [\text{E}_s]$$

We describe the last part of the infection pathway, assuming that the virus has already been engulfed in the host cell. We skip the first steps because we cannot express the phagocytosis of the virus. In Section 5 we will analyse the missing part. Therefore, let

$$\text{cLL}(\llbracket \text{M}_c \rrbracket) [\text{Endosome} , \text{VES}(\llbracket \text{GLY}(1) \rrbracket) [\text{Virus}] , \text{Cytosol}]$$

be the initial solution. A possible run is:

$$\begin{aligned} & \text{CLL}(\text{M}_c) \text{EDSM}(\text{M}_e) [\text{E}_s] , \text{VES}(\text{GLY}(1)) [\text{Virus}] , \text{Cytosol} \\ & \xrightarrow{\tau} \text{CLL}(\text{M}_c) \text{EDSM}(\text{M}_e, \text{GLY}(1)) [\text{Virus}, \text{E}_s] , \text{Cytosol} \end{aligned} \quad (\text{r}_3)$$

$$\begin{aligned} & \equiv \text{CLL}(\text{M}_c) \text{EDSM}(\text{M}_e, \text{GLY}(1)) [\text{VS}(\text{HA}(1)) [\text{Nucaps}], \text{E}_s] , \text{Cytosol} \\ & \xrightarrow{\tau} \text{CLL}(\text{M}_c) \text{EDSM}(\text{M}_e, \text{GLY}(1), \text{HA}(1)) [\text{E}_s] , \text{Nucaps}, \text{Cytosol} \end{aligned} \quad (\text{r}_4)$$

4.3 An extensional semantics for cells

The extensional semantics of Definition 3 must be refined in order to account with new transitions and mreagents. This refinement should for instance equate solutions such as $A(1^x), M(B(1^x)) [S]$ and $A(1), N(B(1)) [S]$ when $\mathcal{C}(\mathbf{r}) = ((A, 1, \emptyset), (B, 1, \emptyset))$, $\mathcal{D}(\mathbf{r}') = ((A, 1, \emptyset), (B, 1, \emptyset))$ and $\mathcal{A}(B_{\mathbf{r}}, N) = M'$ and $\mathcal{A}(B_{\mathbf{r}'}, M') = N$. This case is very similar to that of reversible reactions discussed in the previous section.

The notations $S \xRightarrow{\tau} T$ and $S \xRightarrow{\mu} T$, $\mu \neq \tau$, are defined in the same way as in Definition 2.

Definition 6 A context bisimulation is a symmetric binary relation \mathfrak{R} between solutions such that it is a bisimulation and if $S \mathfrak{R} T$ and if $S \xrightarrow{M_{\mathbf{r}}} \{M; S''\} \cdot S'$ then $T \xrightarrow{M_{\mathbf{r}}} \{M'; T''\} \cdot T'$ and, for every N, R , and N such that $\mathcal{F}(\mathbf{r}) = (M \otimes N, N')$, we have both:

- $(S'', N' (M, N) [S']) \mathfrak{R} (T'', N' (M', N) [T''])$
- $(S', N' (M, N) [S'', R]) \mathfrak{R} (T', N' (M', N) [T'', R])$.

S is context bisimilar to T , written $S \approx_c T$, if $S \mathfrak{R} T$ for some context bisimulation \mathfrak{R} .

Context bisimulation retains the same substitutivity property of \approx stated in Proposition 2. Besides, we have the following theorem:

Theorem 1 \approx_c is a congruence.

Proof: The proof is similar to that of Theorem 1. We only discuss the case of contexts $[\cdot], R$. Let \mathfrak{R} be a context bisimulation such that $S \mathfrak{R} T$ and let $(S', R) \mathfrak{R}' (T', R)$ if

- (1) $S' \mathfrak{R} T'$,
- (2) and S', R and T', R are well-formed.

To demonstrate that \mathfrak{R}' is a bisimulation one has to prove that if $S', R \xrightarrow{\mu} U$ then $T', R \xrightarrow{\mu} U'$ and $U \mathfrak{R}' U'$. Among the possible cases we discuss the case when $S', R \xrightarrow{\tau} U$ is due to a rule (FUSE) between an mreagent of S' and another of R .

Let $S' \xrightarrow{M_r} \langle M; S''' \rangle \cdot S'', R \xrightarrow{N_r} \langle N; R'' \rangle \cdot R'$, and $\mathcal{F}(r) = (M \otimes N, M')$. Then $U = S'', M'(\langle M, N \rangle[S''', R'']), R'$. Since $S' \mathfrak{R} T'$ then $T' \xrightarrow{M_r} \langle M'; T''' \rangle \cdot T''$ and $U' = T'', M'(\langle M', N \rangle[T''', R'']), R'$. By definition of \mathfrak{R} :

$$(S'', M'(\langle M, N \rangle[S''', R''])) \mathfrak{R} (T'', M'(\langle M', N \rangle[T''', R'']))$$

Therefore $U \mathfrak{R}' U'$ follows because of the \mathfrak{R}' -closure with respect to contexts $[\cdot], R'$.

Context bisimilarity retains a universal quantification that is hard to check in practice. One might wonder whether it is possible to simplify the definition. For example, instead of quantifying on cells, one may simply require the bisimilarity of components of mreagents, namely $M \approx_c M', S'' \approx_c T''$, and $S' \approx_c T'$. It is easy to demonstrate that the induced equivalence, which we note \approx_c^+ , is a congruence and $\approx_c^+ \subseteq \approx_c$. At the time we write this note it is not clear to us whether this containment is strict or not. This issue actually requires further investigations.

5 Translocation and phagocytosis

The $\text{bio}\kappa$ -calculus presented in Sections 2 and 4 has a limited expressive power: mechanisms such as translocation, where a single protein may enter a cell, or phagocytosis, where a cell may enter another cell cannot be described. For this reason we overlooked the first steps of the virus infection in example of Section 4.2. The integration of translocation and phagocytosis in $\text{bio}\kappa$ -calculus is not simple and admits several design choices. We discuss few possible formalisations below.

5.1 Translocation

Translocation is a mechanism enabling the transport of proteins through a membrane. This mechanism is very specific and controlled by particular membrane proteins that are different for each type of membrane.

Translocations usually do not transport the full proteins in one step because they are too big for traversing the membrane. This problem is solved in two ways. One way is that the mRNA containing the genetic code of the protein interacts with a ribosome – big proteic complex in charge of translating the

code. This interaction translates part of the code in the cell and, *at the same time*, creates the encoded protein. Alternatively, ad-hoc proteins, called the *chaperons*, are used to unfold the entering protein during the process (this is actually what happen at the end of the RTK-MAPK cascade described in Section 2).

We abstract from this low level mechanisms and assume that proteins may safely traverse the membrane. A first approximative definition of translocation might only check that the entering protein be disconnected and retains a suitable interface:

$$A(\phi + \psi), M(\llbracket M \rrbracket)[S] \xrightarrow{\tau} M(\llbracket M \rrbracket)[A(\overline{\phi} + \psi), S']$$

According to our notation, both ϕ and ψ are v - h -maps, therefore $\text{en}(\phi + \psi) = \emptyset$. This description is not satisfactory for it makes the membrane play a passive role while this process is on the contrary highly specific. As such, it is impossible to represent the effects provided by the chaperon proteins.

A better way is to model translocation as a decomplexation rule between an external protein already connected to the membrane (and nowhere else) and a membrane protein. To avoid the confusion of two different phenomena – simple decomplexations and decomplexations with translocations – we use a further function \mathcal{T} from rule names to tuples $((A, i, \phi, \phi'), (B, j, \psi, \psi'), M, N)$. As usual we assume that there is no clash between rule names in the domain of \mathcal{T} and the other functions that have been used in the paper. The two rules controlling translocations are:

$$\begin{array}{c} \text{(TRS-P)} \\ \hline \mathcal{T}(\mathbf{r}) = ((A, i, \phi, \phi'), (B, j, \psi, \psi'), M, N) \\ \hline A(i^x + \phi + \phi') \xrightarrow{A_r^x} \mathbf{0} \\ \\ \text{(TRS-M)} \\ \hline M \xrightarrow{B_r^x} M' \\ \mathcal{T}(\mathbf{r}) = ((A, i, \phi, \phi'), (B, j, \psi, \psi'), M, N) \\ \hline M(\llbracket M \rrbracket)[S] \xrightarrow{B_r^x} N(\llbracket M' \rrbracket)[A(i + \overline{\phi} + \phi'), S'] \end{array}$$

Fig. 4. Reduction rules for translocation

In rule (TRS-P) the interface of A has exactly one site bound because the interfaces ψ and ψ' are v - h -maps according our notation. The interface ψ is being turned into $\overline{\psi}$ during the translocation, ψ' is unchanged. The protein A disappears during (TRS-P). Dually, the protein A appears during (TRS-M). The translocation will be the effect of the synchronisation (REACT).

5.2 Phagocytosis.

As stated in the introduction, endocytosis is a mechanism allowing a cell to engulf new material. There exists of course many different types of endocytosis. We focus in particular on the *active endocytosis* which is a phenomenon arising as a response to an activation or, more generally, to an external signal. In order to distinguish it from the general endocytosis and since we intend to use this phenomenon to engulf complex structures (such as a virus), we use the term of *phagocytosis*, after the name of the associated brane primitive. The phagocytosis of $M(M)[S]$ – usually a small cell – by $N(N)[T]$ – usually a big cell – transports the former into the cytoplasm of the latter *by surrounding $M(M)[S]$ with a new membrane that is part of N* . This surrounding mechanism is the problematic one because it amounts to split the host cell membrane in some “not local” way. For example the transition

$$M(M)[S] , N(N)[T] \xrightarrow{\tau} N(N)[N'(\emptyset)[M(M)[S]] , T$$

is not very appropriate because the new membrane N' is empty. As we don't have any mechanism for feeding the membrane yet, the solution ban complexations of the new membrane with proteins.

Actually, phagocytosis should be possible provided the membrane of the host cell had enough material for a new membrane. We therefore model phagocytosis as a decomplexation of two proteins in the membranes of the reactant cells with the side effect of splitting the host cell according to some predefined pattern. As for translocations, we use a new functions from rule names to tuples $((A, i, \phi, \phi'), (B, j, \psi, \psi'), M, N, N', N'', N')$, where $\text{de}(N') = \emptyset$. The meaning of this tuple is the following: (A, i, ϕ, ϕ') and (B, j, ψ, ψ') are the two proteins that decomplexate and are located in two membranes M and N , respectively. The name N' is the one given at the new membrane surrounding the phagocytosed cell, N' is the membrane material of the new cell.

In the following rules, transition labels are extended with $M_{\mathbf{r}}^x$ and still ranged over by μ . Let \uplus denote disjoint union of sets. The rules defining phagocytosis are the one of Figure 5.

Rule (OPEN-P) defines the transition of the phagocytosed cell. The label $A_{\mathbf{r}}^x$ of the membrane transition becomes $M_{\mathbf{r}}^x$ in the cellular transition. This exposes the phagocytosis to the label. Similarly for the rule (OPEN-C). In (OPEN-C) the material needed for the new membrane surrounding the phagocytosed cell is removed from the host cell membrane. This material is restored in the rule (PHAGO).

It is not clear to us how close the above rules are to the biological phagocytosis. It is worth to observe that (OPEN-C) is computationally expensive, at least if

$$\begin{array}{c}
\text{(OPEN-P)} \\
\frac{M \xrightarrow{A_x^x} M' \quad \mathcal{P}(\mathbf{r}) = ((A, \mathbf{i}, \phi, \phi'), (B, \mathbf{j}, \psi, \psi'), M, N, N', N'', N')}{M(\mathbf{M})[S] \xrightarrow{M_x^x} \mathbf{M}' ; S \cdot \mathbf{0}} \\
\\
\text{(OPEN-C)} \\
\frac{M \xrightarrow{B_x^x} \prod_{k \in I \uplus J} B_k(\sigma_k) \quad \mathcal{P}(\mathbf{r}) = ((A, \mathbf{i}, \phi, \phi'), (B, \mathbf{j}, \psi, \psi'), M, N, N', N'', \prod_{j \in J} B_j(\sigma_j))}{N(\mathbf{M})[S] \xrightarrow{N_x^x} \mathbf{N} ; \prod_{i \in I} B_i(\sigma_i) ; S \cdot \mathbf{0}} \\
\\
\text{(PHAGO)} \\
\frac{S \xrightarrow{M_x^x} \mathbf{M} ; S'' \cdot S' \quad T \xrightarrow{N_x^x} \mathbf{N} ; T'' \cdot T' \quad \mathcal{P}(\mathbf{r}) = ((A, \mathbf{i}, \phi, \phi'), (B, \mathbf{j}, \psi, \psi'), M, N, N', N'', N')}{S, T \xrightarrow{\tau} S', T', N'(\mathbf{N})[N''(\mathbf{N}')][M(\mathbf{M})[S'']] , T''}
\end{array}$$

Fig. 5. Reduction rules for phagocytosis

compared to the other operations described in the paper. According to (OPEN-C), extracting a pattern of proteins out of a membrane amounts to a long sequel of checks that lock the membrane, thus inhibiting other interactions. It is an open question whether it is possible or not to design simpler and more basic mechanisms for phagocytosis.

Then we could finish the modelling of the example of Section 4.2. We present Figure 6 the complete set of rules used to make evolve the solution **Virus**, **Cell** to the state in which the nucleocapsid is released into the cytoplasm of the infected cell.

$$\begin{array}{ll}
\mathbf{r}_1 : & ((\text{HA}, 1, \emptyset), (\text{GLY}, 1, \emptyset)) \in \mathcal{C} \\
\mathbf{r}'_1 : & (\text{GLY}_{\mathbf{r}_1}, \text{CLL}) \mapsto \text{ACLL} \in \mathcal{A} \\
\mathbf{r}_2 : & ((\text{HA}, 1, \emptyset), (\text{GLY}, 1, \emptyset)), \text{VS}, \text{ACLL}, \text{CLL}, \text{VES}, (\text{GLY}(1)) \in \mathcal{P} \\
\mathbf{r}_3 : & ((\text{EDSM}, \text{VES}), \text{EDSM}) \in \mathcal{F} \\
\mathbf{r}_4 : & ((\text{EDSM}, \text{VS}), \text{EDSM}) \in \mathcal{F}
\end{array}$$

Fig. 6. Set of interactions rules modelling the virus infection

The formal rendering of the example of Section 4.2 can now be stated from the beginning, namely before that the phagocytosis occurs. This is what is presented in Figure 7.

$$\begin{aligned}
& \text{VS}(\text{HA}(1)) \mid \text{CLL}(\text{GLY}(1), \text{M}_c) \mid [\text{Endosome}, \text{Cytosol}] \\
& \xrightarrow{\tau} \text{VS}(\text{HA}(1^x)) \mid \text{Nucaps} \mid \text{ACLL}(\text{GLY}(1^x), \text{M}_c) \mid [\text{Endosome}, \text{Cytosol}] \quad (\mathbf{r}_1 - \mathbf{r}'_1) \\
& \xrightarrow{\tau} \text{CLL}(\text{M}_c) \mid \text{VES}(\text{GLY}(1)) \mid \text{VS}(\text{HA}(1)) \mid \text{Nucaps} \mid [\text{Endosome}, \text{Cytosol}] \quad (\mathbf{r}_2) \\
& \xrightarrow{\tau} \text{CLL}(\text{M}_c) \mid \text{EDSM}(\text{M}_e, \text{GLY}(1)) \mid \text{Virus}, \text{E}_s \mid [\text{Cytosol}] \quad (\mathbf{r}_3) \\
& \equiv \text{CLL}(\text{M}_c) \mid \text{EDSM}(\text{M}_e, \text{GLY}(1)) \mid \text{VS}(\text{HA}(1)) \mid \text{Nucaps} \mid \text{E}_s \mid [\text{Cytosol}] \\
& \xrightarrow{\tau} \text{CLL}(\text{M}_c) \mid \text{EDSM}(\text{M}_e, \text{GLY}(1), \text{HA}(1)) \mid \text{E}_s \mid \text{Nucaps} \mid [\text{Cytosol}] \quad (\mathbf{r}_4)
\end{aligned}$$

Fig. 7. Run of the first steps of the virus infection

6 Implementing the core bio κ -calculus

In this section we discuss an implementation of core bio κ -calculus in π -calculus [12]. The encoding below is simpler than the one in [3], because rules in core bio κ -calculus have always two interacting proteins. In particular, the implementation below do not change the granularity of reactions: every biological reaction is translated into *single* π -calculus transitions.

We begin with a brief introduction to π -calculus, and then detail the compilation of bio κ -calculus.

6.1 The π -calculus

The π -calculus uses a countable set of *names* \mathcal{N} , ranged over by x, y, z, \dots , and *agent names*, ranged over by $\mathbf{A}, \mathbf{B}, \dots$. Tuples of names are noted \tilde{x} . Processes P are defined by the grammar:

$$\begin{aligned}
P := & \mathbf{0} \mid \bar{x}\tilde{z}.P \mid x(\tilde{z}).P \mid (x)P \mid P + P \mid P \mid P \\
& \mid [x = y]P \mid \mathbf{A}(\tilde{x})
\end{aligned}$$

where process names are defined by a set of equations $\mathbf{A}(\tilde{z}) := P$.

A process can be the inert process $\mathbf{0}$, an output $\bar{x}\tilde{z}.P$ that produces a message $\bar{x}\tilde{z}$ and behaves like P ; an input $x(\tilde{z}).P$ that consumes a message $\bar{x}\tilde{u}$ and behaves like $P\{\tilde{u}/\tilde{z}\}$; a restriction $(x)P$ that behaves as P except that messages on x are prohibited; a choice $P + Q$ that behaves either as P or as Q ; a parallel composition $P \mid Q$ of two processes; a match $[x = y]P$ that behaves as P provided x and y are equal; an agent invocation $\mathbf{A}(\tilde{x})$. The input $x(\tilde{z}).P$, restriction $(x)P$, and agent definition $\mathbf{A}(\tilde{z}) := P$ are binders of names \tilde{z} , x , and \tilde{z} , respectively. The scope of these binders are the processes P . We use the

$$\begin{array}{c}
\text{(INP)} \\
x(\tilde{z}).P \xrightarrow{x(\tilde{z})} P \\
\\
\text{(NEW)} \\
\frac{P \xrightarrow{\mu} Q \quad x \notin \text{fn}(\mu)}{(x)P \xrightarrow{\mu} (x)Q} \\
\\
\text{(SUM)} \\
\frac{P \xrightarrow{\mu} P'}{P + Q \xrightarrow{\mu} P'} \\
\\
\text{(MATCH)} \\
\frac{P \xrightarrow{\mu} Q}{[u = u]P \xrightarrow{\mu} Q} \\
\\
\text{(COM)} \\
\frac{P \xrightarrow{(\tilde{y})\tilde{x}\tilde{z}} P' \quad Q \xrightarrow{x(\tilde{w})} Q' \quad \tilde{y} \cap \text{fn}(Q) = \emptyset}{P|Q \xrightarrow{\tau} (\tilde{y})(P'|Q'\{\tilde{z}/\tilde{w}\})} \\
\\
\text{(OUT)} \\
\tilde{x}\tilde{z}.P \xrightarrow{\tilde{x}\tilde{z}} P \\
\\
\text{(OPEN)} \\
\frac{P \xrightarrow{(\tilde{y})\tilde{x}\tilde{z}} Q \quad z' \neq x \quad z' \in \tilde{z} \setminus \tilde{y}}{(z')P \xrightarrow{(z'\tilde{y})\tilde{x}\tilde{z}} Q} \\
\\
\text{(PAR)} \\
\frac{P \xrightarrow{\mu} P' \quad \text{bn}(\mu) \cap \text{fn}(Q) = \emptyset}{P|Q \xrightarrow{\mu} P'|Q} \\
\\
\text{(APP)} \\
\frac{A(\tilde{x}) := P \quad P\{\tilde{z}/\tilde{x}\} \xrightarrow{\mu} Q}{A(\tilde{z}) \xrightarrow{\mu} Q}
\end{array}$$

Fig. 8. Operational semantics of the π -calculus.

standard notion of α -equivalence, *free* and *bound names* of processes, noted $\text{fn}(P)$ and $\text{bn}(P)$, respectively. In the following we let $\sum_{i \in I} P_i$ be the choice between the processes P_i . We also let $\prod_{i \in I} [x_i = y_i]$ be the sequence of matches $[x_i = y_i]$.

Figure 8 collects the semantics of π -calculus, except for the symmetric forms of rules (SUM) and (PAR) that are omitted. The semantics is described as a transition system on syntactic processes with transitions labelled by *actions*. These actions, written μ , are of three types: *internal* actions τ , *inputs* $x(\tilde{u})$, and *outputs* $(\tilde{y})\tilde{x}\tilde{u}$. Outputs $(\tilde{y})\tilde{x}\tilde{u}$ are *bounded* when \tilde{y} is not empty. Bounded outputs, $(\tilde{y})\tilde{x}\tilde{u}$, generated by the transitions above all satisfy: 1) $\tilde{y} \subseteq \tilde{u}$ and 2) $x \notin \tilde{y}$.

One defines:

$$\begin{array}{ll}
\text{fn}(\tau) = \emptyset & \text{bn}(\tau) = \emptyset \\
\text{fn}(x(\tilde{u})) = \{x\} & \text{bn}(x(\tilde{u})) = \tilde{u} \\
\text{fn}((\tilde{y})\tilde{x}\tilde{u}) = \{x\} \cup (\tilde{u} \setminus \tilde{y}) & \text{bn}((\tilde{y})\tilde{x}\tilde{u}) = \tilde{y}
\end{array}$$

Rules (PAR), (COM), (NEW) and (OPEN) have side-conditions controlling bounded output and involving $\text{fn}(\mu)$ and $\text{bn}(\mu)$. Specifically, these conditions ensure that 1) the exported bounded names do not capture any variables when they

finally appear in the right hand side of the conclusion of rule (COM), 2) it is not possible to input or output on a restricted name.

6.2 The compilation

We distinguish three names, v , h , and b , in the set \mathcal{N} . These names are considered constants. Our compilation is “protein-centric”, that is to say proteins are translated as processes whose behaviour is obtained from all the interactions they participate to. In particular, the behaviour of a protein A is specified by a (unique) process definition

$$A(\tilde{u}; \tilde{x}; \tilde{r}) := P .$$

where the arguments of A have been partitioned using semicolons because they play different roles. The tuples \tilde{u} and \tilde{x} have length $\mathfrak{s}(A)$. Names \tilde{u} are always instantiated by v , h , or b , according to the corresponding site is visible, hidden, or bound, respectively. The tuple \tilde{x} is instantiated by names representing bindings between proteins. These names are meaningful provided the corresponding name in the tuple \tilde{u} is b . When names are meaningless, they are set to h . The tuple \tilde{r} has length twice the number of complexation rules where A is a reactant (\tilde{r} may be seen as a sequence of pairs). Let \mathbf{r} be a complexation with a reactant A and let r_i, r_e be the pair in \tilde{r} corresponding to such a rule. The two names r_i and r_e are used for interacting with proteins outside the cell and with proteins in the cytoplasm by means of \mathbf{r} . This distinction is significant for proteins in membranes; the two names are identical for the other proteins. In the following we assume that complexation rules are enumerated. Therefore, there is a first complexation rule, a second one, etc.

The $A(\tilde{u}; \tilde{x}; \tilde{r}) := P$ requires the following notational conventions:

$\text{prj}_i(\tilde{x})$ is the i -th element of \tilde{x} .

$\sharp_c(\mathbf{r})$ gives the ordinal of the complexation rule \mathbf{r} . It is undefined if \mathbf{r} is a decomplexation.

$\sharp_{d,A}(\mathbf{r})$ gives the bound site of the protein A in the decomplexation rule \mathbf{r} . It is undefined if \mathbf{r} is a complexation or A is not a reactant of \mathbf{r} .

$\text{TEST}(\tilde{u}, \phi)$ gives the sequence of matches $\prod_{j \in \text{dom}(\phi)} [\text{prj}_j(\tilde{u}) = \phi(j)]$.

$\text{SET}(\tilde{u}, i, \phi)$, where $i \notin \text{dom}(\phi)$, is defined elementwise as follows:

$$\text{prj}_j(\text{SET}(\tilde{u}, i, \phi)) := \begin{cases} \text{prj}_j(\tilde{u}) & \text{if } j \notin \{i\} \cup \text{dom}(\phi) \\ b & \text{if } i = j \text{ and } \text{prj}_i(\tilde{u}) = v \\ v & \text{if } i = j \text{ and } \text{prj}_i(\tilde{u}) = b \\ h & \text{if } \phi(j) = v \\ v & \text{if } \phi(j) = h \end{cases}$$

$\text{SET}_x(\tilde{u}, i)$ is defined elementwise as follows:

$$\text{prj}_j(\text{SET}_x(\tilde{u}, i)) := \begin{cases} x & \text{if } i = j \\ \text{prj}_j(\tilde{u}) & \text{otherwise} \end{cases}$$

Every preliminary notion is set for the definition of $\mathbf{A}(\tilde{u}; \tilde{x}; \tilde{r})$.

$$\begin{aligned} \mathbf{A}(\tilde{u}; \tilde{x}; \tilde{r}) := & \sum_{\mathbf{A}(i+\phi+\phi') \in \text{left } \mathcal{C}(\mathbf{r})} \text{TEST}(\tilde{u}, i + \phi + \phi') (z) \Big(\\ & \overline{\text{prj}_{2\#_c(\mathbf{r})}(\tilde{r})} z. \mathbf{A}(\text{SET}(\tilde{u}, i, \phi); \text{SET}_z(\tilde{x}, i); \tilde{r}) \\ & + \overline{\text{prj}_{2\#_c(\mathbf{r})+1}(\tilde{r})} z. \mathbf{A}(\text{SET}(\tilde{u}, i, \phi); \text{SET}_z(\tilde{x}, i); \tilde{r}) \Big) \\ & + \sum_{\mathbf{A}(i+\phi+\phi') \in \text{right } \mathcal{C}(\mathbf{r})} \text{TEST}(\tilde{u}, i + \phi + \phi') \Big(\\ & \text{prj}_{2\#_c(\mathbf{r})}(\tilde{r}) (z). \mathbf{A}(\text{SET}(\tilde{u}, i, \phi); \text{SET}_z(\tilde{x}, i); \tilde{r}) \\ & + \text{prj}_{2\#_c(\mathbf{r})+1}(\tilde{r}) (z). \mathbf{A}(\text{SET}(\tilde{u}, i, \phi); \text{SET}_z(\tilde{x}, i); \tilde{r}) \Big) \\ & + \sum_{\mathbf{A}(i+\phi+\phi') \in \text{left } \mathcal{D}(\mathbf{r})} [\text{prj}_i(\tilde{u}) = b] \text{TEST}(\tilde{u}, \phi + \phi') \\ & \overline{\text{prj}_{\#_{d,\mathbf{A}}(\mathbf{r})}(\tilde{x})} . \mathbf{A}(\text{SET}(\tilde{u}, i, \phi); \text{SET}_h(\tilde{x}, i); \tilde{r}) \\ & + \sum_{\mathbf{A}(i+\phi+\phi') \in \text{right } \mathcal{D}(\mathbf{r})} [\text{prj}_i(\tilde{u}) = b] \text{TEST}(\tilde{u}, \phi + \phi') \\ & \text{prj}_{\#_{d,\mathbf{A}}(\mathbf{r})}(\tilde{x}) (). \mathbf{A}(\text{SET}(\tilde{u}, i, \phi); \text{SET}_h(\tilde{x}, i); \tilde{r}) \end{aligned}$$

Every branch of the choices in $\mathbf{A}(\tilde{u}; \tilde{x}; \tilde{r})$ verifies whether the interface fits with a left-hand side of some rule having \mathbf{A} as a reactant or not. The test is implemented by the operation $\text{TEST}(\tilde{u}, \phi)$, following our encoding of sites in the arguments \tilde{u} of \mathbf{A} . We notice that, the branches of $\mathbf{A}(\tilde{u}; \tilde{x}; \tilde{r})$ corresponding to decomplexation rules have an additional test $[\text{prj}_i(\tilde{u}) = b]$ verifying that the i -th site is actually bound.

Once a branch has been chosen, a biological reaction is compiled into a single π -calculus interaction, where the left reactant plays the role of sender and the

right reactant plays the role of receiver. The sender of a complexation has to create a new name representing the biological edge created by the rule. This name is communicated during the interaction to the other reactant and will be used to decomplexate the two proteins. The new state of the protein A is obtained by updating the sites as prescribed by the rule – see the definition of SET.

Few other notational conventions are useful for the definition of the compilation. Let $\llbracket \sigma \rrbracket_A$ be function yielding a tuple of length $2 \times \mathfrak{s}(A)$ defined elementwise as follows:

$$\text{prj}_i(\llbracket \sigma \rrbracket_A) = \begin{cases} v & \text{if } \sigma(i) = v \\ h & \text{if } \sigma(i) = h \\ b & \text{otherwise} \end{cases}$$

$$\text{prj}_{2i}(\llbracket \sigma \rrbracket_A) = \begin{cases} x & \text{if } \sigma(i) = x \text{ and } x \notin \{v, h\} \\ h & \text{otherwise} \end{cases}$$

For example $\llbracket 1^x + 2 + \overline{3} \rrbracket_A = (b, v, h; x, h, h)$. We notice that $\llbracket \sigma \rrbracket_A$ is undefined if σ is not total on $1..\mathfrak{s}(A)$. Let also $(x_1, y_1) \cdots (x_n, y_n) \mid_A$ be the subsequence containing (x_i, y_i) provided A is a reactant of the i -th complexation rule. Finally, let n be the number of complexation rules of the biological system. The compilation of a solution S is given by $\llbracket S \rrbracket_{(r_1, r_1), \dots, (r_n, r_n)}$, where r_1, \dots, r_n are pairwise different names:

$$\begin{aligned} \llbracket 0 \rrbracket_\ell &:= \mathbf{0} \\ \llbracket A(\sigma) \rrbracket_\ell &:= \mathbf{A}(\llbracket \sigma \rrbracket_A; \ell \mid_A) \\ \llbracket S, T \rrbracket_\ell &:= \llbracket S \rrbracket_\ell \mid \llbracket T \rrbracket_\ell \\ \llbracket \mathbf{M}(\mathbf{M})[S] \rrbracket_{(r_1, r'_1) \cdots (r_n, r'_n)} &:= (r''_1) \cdots (r''_n) \left(\llbracket \mathbf{M} \rrbracket_{(r'_1, r''_1) \cdots (r'_n, r''_n)} \mid \llbracket S \rrbracket_{(r''_1, r'_1) \cdots (r''_n, r'_n)} \right) \end{aligned}$$

The compilation schemas are simple except that of $\mathbf{M}(\mathbf{M})[S]$, which is discussed below. As regards cells, one has to ban reactions between proteins in the cytoplasm and those that are external to the cell. This may be enforced in π -calculus by picking a fresh set of names for complexation rules. For this reason the compilation creates a fresh set of complexation names r''_1, \dots, r''_n and allows them to occur in proteins of S and M. On the contrary, proteins in the membrane M may interact both with the external environment and with the cytoplasm. The compilation supports this feature by encoding proteins in M with the set of complexation names of the external environment and the fresh set of complexation names of the cytoplasm. It is worth to notice that

decomplexations are not an issue because their correctness follows from the well-formedness of the solution.

The following proposition states the correctness of the compilation. The proof is a standard induction on the structure of S .

Proposition 2 *Let $\mathcal{C} = \{\mathbf{r}_1, \dots, \mathbf{r}_n\}$, and let r_1, \dots, r_n and r'_1, \dots, r'_n be two tuples of pairwise different names. Let $\ell = (r_1, r'_1) \dots (r_n, r'_n)$. The following reductions are in core $\mathbf{bio}\kappa$ -calculus.*

- (1) $S \xrightarrow{A_{\mathbf{r}_k}^x} T$, where A is the left reactant (resp. right reactant), if and only if either $\llbracket S \rrbracket_\ell \xrightarrow{(z)\overline{r_k}z} \llbracket T \rrbracket_\ell$ or $\llbracket S \rrbracket_\ell \xrightarrow{(z)r'_kz} \llbracket T \rrbracket_\ell$ (resp. either $\llbracket S \rrbracket_\ell \xrightarrow{r_k(z)} \llbracket T \rrbracket_\ell$ or $\llbracket S \rrbracket_\ell \xrightarrow{r'_k(z)} \llbracket T \rrbracket_\ell$);
- (2) $S \xrightarrow{A_{\mathbf{r}}^x} T$, where \mathbf{r} is a decomplexation rule and A is the left reactant (resp. the right reactant), if and only if $\llbracket S \rrbracket_\ell \xrightarrow{\bar{z}} \llbracket T \rrbracket_\ell$ (resp. $\llbracket S \rrbracket_\ell \xrightarrow{z()} \llbracket T \rrbracket_\ell$);
- (3) $S \xrightarrow{\tau} T$ if and only if $\llbracket S \rrbracket_\ell \xrightarrow{\tau} \llbracket T \rrbracket_\ell$.

We conclude by commenting on the extension of the encoding to the full $\mathbf{bio}\kappa$ -calculus. The problematic rules of $\mathbf{bio}\kappa$ -calculus are fusions because they allow interactions between proteins that were banned. For instance, the fusion between two cells in a same solution allows interactions between the proteins in the two cytoplasms. So, a local interaction has global effects in a given environment. This situation is usually hard to encode into π -calculus. A way out is to use formalisms permitting simple modellings of such global effects, such as the *fusion calculus* [13]. In this calculus, a communication equates names, thus causing communications on names that were not possible before. In the foregoing compilation, a fusion should equate tuples of names r_1, \dots, r_n , thus allowing complexations between proteins using such names. We are currently investigating the soundness of this idea.

7 Conclusions

We have presented a unique framework for modelling proteins and cells interactions – the $\mathbf{bio}\kappa$ -calculus. Protein interactions in $\mathbf{bio}\kappa$ -calculus are of two types: complexations and decomplexations; cell interactions in $\mathbf{bio}\kappa$ -calculus describe fusions. All interactions are “local” in the sense that they always involve two proteins. Fusions have been modelled by using an higher order semantics in the style of [10]. We have studied the operational semantics of $\mathbf{bio}\kappa$ -calculus and an extensional semantics of its – the bisimulation. The expressiveness has been analysed by modelling two significant systems and

comparing them with similar ones that have been proposed in other algebraic approaches.

Some extensions of **bio** κ -calculus rules may be done without difficulties. In this paper we have discussed rules modelling translocations and phagocytosis, even if the latter ones are not very satisfactory. Other biological reactions have not yet been considered and are left to future work, such as those in [6] or in [11].

Extensional semantics of **bio** κ -calculus are an intriguing issue we plan to investigate in the future. In particular we are interested in mathematical tools and techniques that help in assessing properties of biological solutions. Such tools might include stochastic measures in the style of [14] and might be extensively used to predict outputs of experiments in vitro.

Acknowledgements. We thank Gianluigi Zavattaro for the discussions and the suggestions on the encodings of the **bio** κ -calculus into π -calculus.

References

- [1] R. Milner, Communicating and mobile systems: the π -calculus, Cambridge University Press, Cambridge, 1999.
- [2] A. Regev, W. Silverman, E. Shapiro, Representation and simulation of biochemical processes using the π -calculus process algebra, in: R. B. Altman, A. K. Dunker, L. Hunter, T. E. Klein (Eds.), Pacific Symposium on Biocomputing, Vol. 6, World Scientific Press, Singapore, 2001, pp. 459–470.
URL <http://www.smi.stanford.edu/projects/helix/psb01/regev.pdf>
- [3] V. Danos, C. Laneve, Formal molecular biology, Theoretical Computer Science 325 (1) (2004) 69–110.
- [4] L. Cardelli, A. D. Gordon, Mobile ambients., Theoretical Computer Science 240 (1) (2000) 177–213.
- [5] G. Paun, Membrane computing. an introduction., Springer-Verlag, Berlin, 2002.
- [6] L. Cardelli, Brane calculi., in: CMSB, 2004, pp. 257–278.
- [7] A. Regev, E. M. Panina, W. Silverman, L. Cardelli, E. Shapiro, Bioambients: an abstraction for biological compartments, Theoretical Computer Science 325 (1) (2004) 141–167.
- [8] R. Milner, Communication and Concurrency, International Series on Computer Science, Prentice Hall, 1989.

- [9] V. Danos, F. Tarissan, Self-assembling graphs, in: J. Mira, J. R. Álvarez (Eds.), *Mechanisms, Symbols, and Models Underlying Cognition*, Vol. 3561 of *Lecture Notes in Computer Science*, Springer, 2005, pp. 498–507.
- [10] D. Sangiorgi, From π -calculus to Higher-Order π -calculus — and back, in: M.-C. Gaudel, J.-P. Jouannaud (Eds.), *Proc. TAPSOFT'93*, Vol. 668 of *Lecture Notes in Computer Science*, 1993, pp. 151–166.
- [11] V. Danos, S. Pradalier, Projective brane calculus., in: *CMSB*, 2004, pp. 134–148.
- [12] R. Milner, J. Parrow, D. Walker, A calculus of mobile processes, *Journal of Information and Computation* 100 (1992) 1–77.
- [13] J. Parrow, B. Victor, The fusion calculus: Expressiveness and symmetry in mobile processes, in: *Proceedings of LICS '98*, IEEE, Computer Society Press, 1998, pp. 176–185.
- [14] A. Credi, M. Garavelli, C. Laneve, S. Pradalier, S. Silvi, G. Zavattaro, Modelization and simulation of nano devices in nanok calculus, in: *CMSB*, Vol. 4695 of *Lecture Notes in Computer Science*, Springer, 2007, pp. 168–183.