

# From Biochemistry to Stochastic Processes

Cosimo Laneve<sup>1</sup>

*Dipartimento di Scienze dell'Informazione, Università di Bologna*

Sylvain Pradaliér<sup>2</sup>

*Ecole Polytechnique, Paris*

Gianluigi Zavattaro<sup>3</sup>

*Dipartimento di Scienze dell'Informazione, Università di Bologna*

---

## Abstract

The **nanok** calculus is a formalism that models biochemical systems by defining its set of reactions. We study the implementation of **nanok** into the Stochastic Pi Machine where biochemical systems are defined by regarding molecules as processes, and deriving the overall behaviour by means of communication rules. Our implementation complies with the *stochastic behaviors* of systems, thus allowing one to use **nanok** as an intelligible front-end for a process-oriented simulator. This study also permits to reuse, in **nanok**, the theories and tools already developed for process calculi.

**Keywords:** Stochastic process calculi, systems biology, encoding from reactive to process-oriented formalisms.

---

## 1 Introduction

Several stochastic formalisms emerged in the last few years as models for the representation of biological systems (see e.g. [5,12,2,7,9,8,4] just to mention a few). These formalisms usually follow either a *reactive-oriented* (as [5,2,12,9,4] in the list above) or a *process-oriented* approach (as [7,14,8]). According to the former approach – inspired by traditional chemical kinetics – a system is specified as a set of reactions; according to the latter – inspired by process calculi – a system is specified by defining each molecule as a process, and deriving the overall behaviour by means of communication rules.

---

<sup>1</sup> Email:[laneve@cs.unibo.it](mailto:laneve@cs.unibo.it)

<sup>2</sup> Email:[sylvain.pradaliér@lix.polytechnique.fr](mailto:sylvain.pradaliér@lix.polytechnique.fr)

<sup>3</sup> Email:[zavattar@cs.unibo.it](mailto:zavattar@cs.unibo.it)

Process-oriented descriptions depart from ordinary biochemical models because they define the sequences of actions once and for all and use syntaxes usually devoted to computer science. Moreover the modelling of a molecule is a term of size proportional to the number of interactions addressing the molecule. As a consequence, such descriptions are less intelligible to biochemists than reaction-oriented approaches, whose syntax is closer to biochemistry and whose complexity is spread over the reactions. On the other hand, process-oriented calculi retain several simulators and tools, which make them attractive for experiments *in silico* (see for instance [16,7,14]).

In this paper we bridge the gap between the two approaches by implementing the **nanok** calculus [10], a reactive-oriented formalism, into the *Stochastic Pi Machine* [7], SPIM calculus in the following, a simulator for the stochastic  $\pi$ -calculus [19,7]. In **nanok** calculus a molecule is a term  $A[s^1 + r^2](1 + 2^x)$  with fields  $s$  and  $r$  and sites 1 and 2. The fields define the internal state of the molecule – they model its shape or its hydrogen groups or phosphate groups: in the above case, the fields  $s$  and  $r$  are set to 1 and 2, respectively; the sites are the binding capabilities of the molecule: in the above case,  $A$  is bound to another molecule with site 2, the bond is called  $x$ , and is unbound on site 1. Note that only sites can be bound and that only fields can store a value. The **nanok** calculus retains a graphical representation – the above molecule is rendered in Figure 1(a). The dynamics of a **nanok** calculus system is defined by reactions that describe how two reactants may evolve. For example, the reaction

$$\rho_1 \quad A[s^1](1), B[t^0](1) \xrightarrow{\lambda} A[s^0](1^x), B[t^1](1^x)$$

illustrated in Figure 1(b), specifies that every molecule  $A$  with an internal state  $s$

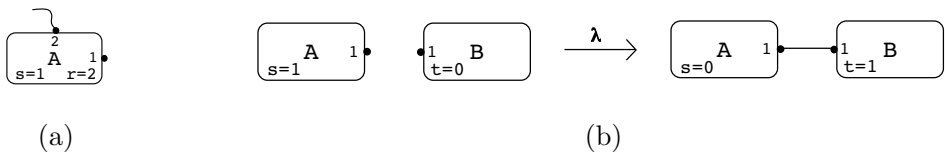


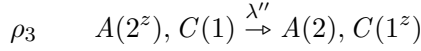
Fig. 1. Molecules and reactions in **nanok** calculus

equal 1 and a free site 1 may react with every  $B$  with internal state  $t$  equal 0 and free site 1. The result is a complex where  $A$  and  $B$  are connected by a bond, called  $x$ , and the two internal states have swapped values. This change in the internal state might represent for instance an exchange of electronical charge. The label  $\lambda$  of the reaction represents its rate. It is worth to notice that this reaction applies to the molecule  $A$  in Figure 1(a), as well as to every other  $A$  with a different value of  $r$  and/or with an unbound site 2.

In **nanok** molecules may react by means of three types of reactions – creations, destructions, and exchanges – and retain a stochastic semantics. The above reaction  $\rho_1$  is an example of creations. The reaction

$$\rho_2 \quad A[s^1](1^x), B[t^0](1^x) \xrightarrow{\lambda'} A[s^0](1), B[t^1](1)$$

defines a destruction that is opposite to the above creation. The rate  $\lambda'$  may be different from  $\lambda$ , thus yielding *equilibria* in accordance with them. The reaction



is an exchange rule defining a bond flipping from the site 2 of  $A$  to 1 of  $C$ . Bond flipping is of a peculiar interest to model nano-machines where links between components are dynamically rearranged [10]. Despite the simplicity of these reactions, their implementation in SPIM calculus is complex if the stochastic semantics must be preserved. Let us discuss the problems through a number of examples.

An implementation of **nanok** into SPIM should project the behaviour of each molecule out of the set of reactions and collect them into a process definition. For example, the SPIM process  $\hat{A}$  of the molecule  $A$  in the above examples is

$$\begin{aligned} \hat{A}(s, t, a_1, a_2) = & \text{behaviour-of } A \text{ in } \rho_1 \\ & + \text{behaviour-of } A \text{ in } \rho_2 \\ & + \text{behaviour-of } A \text{ in } \rho_3 \end{aligned}$$

That is, a molecule is implemented by a parametric process definition, where the parameters  $s, t, a_1, a_2$  define the values of fields and sites. The “behaviour-of  $A$  in  $\rho_1$ ” might be defined as

$$[a_1 = \varepsilon, s = 1] \overline{\rho_1}(x). \hat{A}(0, t, x, a_2)$$

where

- $[a_1 = \varepsilon, s = 1]$  means that such a behaviour may be triggered provided the site  $a_1$  is unbound (has value  $\varepsilon$ ) and the field  $s$  is equal 1;
- in this case the channel  $\rho_1$  is used to output a fresh name (modelling the bond). We expect that  $\hat{B}$  will perform a corresponding input when the field  $r$  is 1 and the site 1 is unbound. We also expect that the rate of the channel  $\rho_1$  has been declared to be  $\lambda$ ;
- then  $\hat{A}$  will continue as the process  $\hat{A}(0, t, x, a_2)$ .

However the analogy *reaction-names as channels* cannot be pushed forward to destruction. In facts, if the “behaviour-of  $A$  in  $\rho_2$ ” was modelled as

$$[a_1 = \neg\varepsilon, s = 0] \overline{\rho_2}(). \hat{A}(1, t, \varepsilon, a_2)$$

then  $\hat{A}$  might interact with the wrong  $\hat{B}$  in the SPIM implementation ( $a_1 = \neg\varepsilon$  means that  $a_1$  is bound). This difficulty may be circumvented by using the names encoding bonds – that are shared exactly by the two connected molecules – in order to send the disconnection signal. So the “behaviour-of  $A$  in  $\rho_2$ ” becomes

$$[a_1 = \neg\varepsilon, s = 0] \overline{a_1}(). \hat{A}(1, t, \varepsilon, a_2)$$

and we are assuming that  $\hat{B}$  will input on  $a_1$  and that the rate of this name is  $\lambda'$ .

Unfortunately, this solution is also defective. Consider the rule

$$\rho_4 \quad A(1^x + 2^z), B(1^x + 2) \xrightarrow{\lambda'''} A(1^x + 2), B(1^x + 2^z)$$

with  $\lambda''' \neq \lambda'$ . In this case, a bond may flip from  $A$  to  $B$  (with an exchange rule) *provided they are connected* through the site 1. A simple solution is to implement this rule with exactly one SPIM interaction by using the channel  $x$  for exchanging the bond  $z$ . But the mismatch between the rates  $\lambda'$  and  $\lambda'''$  makes this solution (stochastically) unfeasible. However, this last difficulty may be overcome by introducing, at creation time, a tuple of channels, one for every use. So in case of a destruction, a channel of the tuple is used (and the whole tuple is destroyed), in case of exchange, another one is used, which retains a different rate. Said otherwise, a bond in **nanok** calculus corresponds to a tuple of channels in SPIM with possibly different rates.

In facts, the above one is the implementation we describe in this contribution. In an accurate solution, a bond in **nanok** should be represented in SPIM by a tuple whose length is the number of reactions that address that bond directly or indirectly through sequences of exchanges. Actually, for simplicity, our solution over-approximates the “precise” solution, by representing bonds with tuples whose length is the size of the set of reactions – the *gangs*, in our terminology.

The encoding of **nanok** into SPIM defined in this paper, let it be  $\llbracket \cdot \rrbracket$ , is such that  $S \xrightarrow{\lambda}_{\text{nanok}} T$  if and only if  $\llbracket S \rrbracket \xrightarrow{\lambda}_{\text{spim}} \llbracket T \rrbracket$  (the arrow is subscribed to ease the reading). It follows that  $S$  and  $\llbracket S \rrbracket$  are *strongly stochastic bisimilar* [3]. This is different from usual implementations that almost never preserve the granularity of transitions. Of course, this strong relationship may be weakened by defining suitable protocols, in the style of [13]. However such a solution might hardly preserve the stochastic semantics. In facts, the stochastic semantics defines an exponential law controlling the waiting time before a transition can be fired (the *sojourn time*). In turn, the match of one reaction with a sequence of transitions amounts to matching an exponential distribution with a sum of exponential distributions, which is not possible.

## Related works.

In [20], it has been shown that systems of molecular interactions with explicit bonds might be represented and simulated using the stochastic  $\pi$ -calculus. Our encoding corroborates this result since the SPIM calculus is a subset of the stochastic  $\pi$ -calculus. We remark that the example provided in [20] and, we believe, the descriptions done in this approach, can easily be rewritten in SPIM calculus and even in a sub-calculus of it, since our encodings doesn't use its full power.

In [6], Cardelli has encoded chemical systems into process algebra and back preserving both the stochastic and the ODE semantics. Our encoding extends these encodings because the CGF process algebra used in [6] is a subset of the SPIM calculus and because the **nanok** calculus extends the language of chemical reactions of [6] with explicit bonds between molecules and with internal states. However,

our results are weaker than those in [6], since we only assert the correctness of the encoding with respect to the stochastic semantics.

Another stochastic process calculus that has been used also for the modeling of biochemical systems is PEPA [14]. For instance, in [4], PEPA has been exploited to examine the influence of the RAF Kinase Inhibitor Protein (RKIP) on the Extracellular signal Regulated Kinase (ERK) signalling pathway. Also in that paper, as in the present one, a reagent-centric view and a pathway-centric (process-centric in our terminology) view are studied. Our analysis of the two approaches is different for two main reasons. First of all, in the PEPA-based approach one process is used to represent the concentration of one species while we follow the Cardelli's approach considering one process for each molecule. In fact, we have found this approach appropriate for a compositional model of discrete state systems (in which we count the number of molecules instead of considering their concentrations). The second difference is that in [4] only finitely many different species are considered, thus the translation from the reagent-centric to the pathway-centric views can be obtained using an intermediate matrix representation that quantifies the impact of each reaction on each reagent in a manner analogous to the stoichiometry matrix of the chemical reactions. We cannot exploit this approach as we do not impose any bound to the number of different complexes that can be produced in a **nanok** calculus system. More recently, the PEPA approach has been also extended with a reaction-centric model called Bio-PEPA [9]. Also in this case, one process is used to represent the concentration of one species.

Encodings from the full  $\kappa$  calculus to **nanok** calculus, or to  $\pi$ -calculus are presented in [13] and [11]. Yet, they only preserve non-stochastic semantics. In facts, encodings preserving the stochastic semantics do not exist, due to the negative results of [17].

## Structure of the article.

The rest of the article is organized as follows. First we recall the syntax of **nanok** and SPIM, and present their basic stochastic semantics. In part 3, we first introduce the gangs and then complete the definition of the encoding. Its correctness is asserted by theorem 3.2 and 3.4. In part 4 we present the collective stochastic semantics and state the correctness of the encoding with respect to the collective stochastic semantics in theorem 4.3.

## 2 The stochastic calculi

We shortly present the two stochastic calculi we analyze in this paper: a subcalculus of **nanok** calculus, where reactants share at most one bond, and a subcalculus of SPIM. Examples and additional details can be found in [10] and [7].

### 2.1 The **nanok** calculus

Terms, called *solutions*, are sequences of *molecules*. Each molecule belongs to a *species* and retain an *internal state*, which is determined by a tuple of *fields*, and

an *interface*, which is a tuple of *sites* that may be bound to other sites. Formally, a molecule will be written  $A[u](\sigma)$ , where

$A$  is the species. The molecules of a species retain the same set of fields and the same set of sites that are finitely many; fields and sites will be addressed by numbers  $0, 1, 2, \dots$ ;

$u$  – called the *evaluation* – is a total map from fields of  $A$  to *finite sets* (the internal states of molecules are always finitely many);

$\sigma$  – called the *interface* – is a total map from sites of  $A$  to either *bonds*, which are names of a *totally ordered countable set* ranged over by  $x, y, z, \dots$ , or  $\varepsilon$ , a special value indicating that the site is not bound.

For example,  $A[1 \mapsto 0; 2 \mapsto 1](1 \mapsto \varepsilon; 2 \mapsto x; 3 \mapsto \varepsilon)$  is a molecule with two fields 1 and 2 and three sites 1, 2, and 3. The fields 1 and 2 have values 0 and 1, respectively; the site 2 is the only one that is bound and the bond is  $x$ . In order to ease the reading, we write this molecule as  $A[1^0 + 2^1](1 + 2^x + 3)$  (the value  $\varepsilon$  is always omitted). Let  $\emptyset$  be the empty map. We write  $A(\sigma)$  instead of  $A[\emptyset](\sigma)$ ,  $A[u]$  instead of  $A[u](\emptyset)$ , and simply  $A$  instead of  $A[\emptyset](\emptyset)$ . We denote by  $\text{ran}(\sigma)$  the range of an interface  $\sigma$  deprived of  $\varepsilon$  and by  $\text{bonds}(S)$  the set of the bonds appearing in the solution  $S$ .

**Definition 2.1** A *solution* is a term defined by the grammar

$$S ::= A[u](\sigma) \mid S, S$$

The operator “,” is assumed to be associative, so  $(S, T), R$  is equal to  $S, (T, R)$  (therefore parentheses are always omitted).

Bonds always occur at most twice in solutions. A solution is *proper* if every bond therein occurs exactly twice.

The **nanok** calculus semantics is defined by means of reaction rules. A few preliminary definitions are in order:

- we write  $\sigma \leq \sigma'$  if  $\text{dom}(\sigma) = \text{dom}(\sigma')$  and, for every  $i$ , if  $\sigma(i) \neq \varepsilon$  then  $\sigma(i) = \sigma'(i)$  (the two interfaces may differ on sites mapped to the empty value  $\varepsilon$  by  $\sigma$ :  $\sigma'$  may map such sites to bonds);
- a *pre-solution* is a sequence of terms  $A[u](\sigma)$  where  $u$  and  $\sigma$  are partial functions (with an abuse of notation, we denote partial and total functions in the same way);
- a pre-solution is *proper* when every bond therein occurs exactly twice.

In the following when we write  $u + u'$  and  $\sigma + \sigma'$  we assume that  $\text{dom}(u) \cap \text{dom}(u') = \emptyset$  and  $\text{dom}(\sigma) \cap \text{dom}(\sigma') = \emptyset$ .

**Definition 2.2** Reactions of **nanok** calculus are either *creations*, *destructions*, or *exchanges* that are labelled by *rates*, which are positive real numbers or  $\infty$ . Cre-

ations have format

$$A[u](\sigma), B[v](\rho) \xrightarrow{\lambda} A[u'](\sigma'), B[v'](\rho'), C_1[w_1](\eta_1), \dots, C_n[w_n](\eta_n)$$

where  $\sigma \leq \sigma'$ ,  $\rho \leq \rho'$ ,  $\text{dom}(u) = \text{dom}(u')$ ,  $\text{dom}(v) = \text{dom}(v')$ , and  $w_i$  and  $\eta_i$  are total. Destructions have formats

$$\begin{aligned} A[u](\sigma), B[v](\rho) &\xrightarrow{\lambda} A[u'](\sigma'), B[v'](\rho') \\ A[u](\sigma), B[v](\rho) &\xrightarrow{\lambda} A[u'](\sigma') \end{aligned}$$

where  $\sigma \geq \sigma'$ ,  $\text{dom}(u) = \text{dom}(u')$ , and, in the first case,  $\rho \geq \rho'$ ,  $\text{dom}(v) = \text{dom}(v')$  and, in the second case,  $\rho$  has to be total. Exchanges have one of the formats:

$$\begin{aligned} A[u](\sigma), B[v](\rho) &\xrightarrow{\lambda} A[u'](\sigma), B[v'](\rho) \\ A[u](a^x + \sigma), B[v](b + \rho) &\xrightarrow{\lambda} A[u'](a + \sigma), B[v'](b^x + \rho) \end{aligned}$$

where the pre-solutions  $A[u](\sigma), B[v](\rho)$  and  $A[u](a + \sigma), B[v](b + \rho)$  are proper and  $\text{dom}(u) = \text{dom}(u')$  and  $\text{dom}(v) = \text{dom}(v')$ .

In the rest of the paper we assume that reactants share at most one bond, i.e.  $\text{ran}(\sigma) \cap \text{ran}(\rho)$  is either an empty set or a singleton.

Creations produce new bonds between two unbound sites and/or synthesize new molecules. Destructions behave in the other way around. Exchanges either leave the interfaces unchanged or move one bond from a reactant to the other (bond-flipping exchange).

It is worthwhile to remark that reactions do not address every field and site of the reactants (evaluations and interfaces are partial). The intended meaning is that two molecules reacts if they are *instances* of the left-hand side of a reaction. We will formalize this notion later on in the section.

## 2.2 The SPIM calculus

The SPIM calculus uses two sets of identifiers: *names*, which is totally ordered and ranged over by  $x, y, u, \dots$ , *agents*, ranged over by  $A, B, \dots$ . Names have a rate that is a positive real number or  $\infty$ . This rate may be explicitly declared in the process or globally defined (for free names). The following syntactic categories are used in SPIM calculus:

$$\begin{array}{lll} M ::= & [u = v] \mid M M & \text{matches} \\ \alpha ::= & x(\tilde{u}) \mid \bar{x}\tilde{u} \mid \bar{x}(\tilde{u} : \tilde{\lambda}) & \text{actions} \\ P ::= & \mathbf{0} \mid A(\tilde{u})|P & \text{terms} \end{array}$$

Matches are sequences of equalities between values. Actions are either *input*  $x(\tilde{u})$  on  $x$  of a tuple  $\tilde{u}$ , or *output*  $\bar{x}\tilde{u}$  on  $x$  of a tuple  $\tilde{u}$ , or *bound output*  $\bar{x}(\tilde{u} : \tilde{\lambda})$  on

$x$  of a tuple  $\tilde{u}$  with rates  $\tilde{\lambda}$ . Terms can be the inert  $\mathbf{0}$  or a parallel composition of agent invocations. The parallel operator  $|$  is assumed to be associative. Agent declarations have the form:

$$A(\tilde{x}) \triangleq \sum_{i \in I} M_i \alpha_i . P_i$$

### Notation.

Whenever a match has the form  $[u = u]$ , or a sum has only one branch we omit to write them explicitly. For instance  $A(\tilde{x}) \triangleq \sum_{i \in \{1\}} [\tilde{x}_i = \tilde{x}_i] \alpha . P$  is written  $A(\tilde{x}) \triangleq \alpha . P$ .

A process is a term  $(\tilde{x} : \tilde{\lambda}) P$  – the set of agent definitions is kept implicit – where  $\tilde{\lambda}$  are rates. The term  $\tilde{x} : \tilde{\lambda}$  has to be considered a set with the constraint that every two different elements have different names. Processes are ranged over by  $P, Q, \dots$ .

Scope restrictions bind names, that is in  $(x : \lambda) P$  the  $x$  free in  $P$  is bound by  $x : \lambda$ . Likewise, input  $x(\tilde{u}).P$  and bound output  $\bar{x}(\tilde{u} : \tilde{\lambda}).P$  bind  $\tilde{u}$  with scope  $P$ . The agent definition  $A(\tilde{u}) \triangleq \sum_{i \in I} M_i \alpha_i . P_i$  binds  $\tilde{u}$  with scope to the right hand side of the definition. Names that are not bound are called *free* and we write  $\mathbf{fn}(T)$  for the set of such names in  $T$ .

We assume that all terms meet the following well formed properties:

- in  $(\tilde{x} : \tilde{\lambda})P$ ,  $\tilde{x} \subseteq \mathbf{fn}(P)$  (there is no garbage);
- bound names in agent definitions never clash with free names (this allows us to avoid alpha-conversions).

The reductions of SPIM calculus are communications on a channel. Since they are fixed, they will not play any relevant role in the following Definition 2.4. (They will be embodied in the (*init*) item of the definition.) Therefore we omit the formal definition here.

### 2.3 Basic transition relations

Reactions only define the (biochemical) changes of the reactants. These descriptions are used to infer transitions of solutions consisting of several possible reactants. Such transition relations are given in two steps: a first one, called *basic transition relation*, that records the position of the reactants in the whole solution; a second one, called *collective transition relation*, that computes the rate of a transition by summing the rate of the basic transitions that produce the same solution (regardless the position of the molecules/agents). Below we define the basic transition relation for **nanok** and SPIM calculi. (The two definitions are very close, this is why they have been collected in this subsection.) Our result of correctness of the encoding of **nanok** in SPIM regards the basic transition relation. It follows that this correctness also holds for the collective semantics, since it is derived in the same way from the basic transition relation (see theorem 4.3 in Section 4).

The definition of the basic transition relation of the **nanok** calculus requires few



notations. Let  $\mu$  range over  $\rho_L$  and  $\rho_R$  and let  $\overline{\rho}_L = \rho_R$  and  $\overline{\rho}_R = \rho_L$  (notice that  $\overline{\overline{\mu}} = \mu$ ). The **nanok** reactions may be addressed by:

$$A[u](\sigma), B[v](\rho) \xrightarrow{\lambda} A[u'](\sigma'), S$$

where  $S$  may also be empty (denoted by  $\_$ ). The special term  $\_$  is considered a unit for the “,” operator (the solutions  $\_, S, S, \_$  and  $S$  are equal). With an abuse of notation we lift a renaming  $\iota$  to a solution by applying it pointwise.

**Definition 2.3** The basic transition relation of **nanok**, written either  $\xrightarrow{\rho, \iota}_{l, l'}$  or  $\xrightarrow{\mu, \iota}_l$ , is the least relation that satisfies the following rules:

- (init) let  $\rho = A[u](\sigma), B[v](\phi) \xrightarrow{\lambda} A[u'](\sigma'), S$ . Then both  $A[u + w](\sigma \circ \iota + \nu) \xrightarrow{\rho_L, \iota}_1 A[u' + w](\sigma' \circ \iota + \nu)$  and  $B[v + w](\phi \circ \iota + \nu) \xrightarrow{\rho_R, \iota}_1 T$ , where  $T$  is either  $B[v' + w](\phi' \circ \iota + \nu), \iota(S)$  or  $\iota(S)$ , according to the shape of the right hand side, and  $\iota$  is an injective renaming with  $\text{ran}(\iota) \cap \text{ran}(\nu) = \emptyset$ ;
- (lifts) if  $S \xrightarrow{\mu, \iota}_l S'$  and  $(\text{bonds}(S') \setminus \text{bonds}(S)) \cap \text{bonds}(T) = \emptyset$ , then both  $S, T \xrightarrow{\mu, \iota}_l S', T$  and  $T, S \xrightarrow{\mu, \iota}_{l' + l} T, S'$ , where  $T$  has  $l'$  molecules;
- (communications) if  $S \xrightarrow{\mu, \iota}_l S'$  and  $T \xrightarrow{\overline{\mu}, \iota}_{l'} T'$  and  $\iota$  is an order-preserving injection that map bonds into the least ones not used in  $S, T$  then  $S, T \xrightarrow{\rho}_{l, l'' + l'} S', T'$ , where  $\rho$  is the rule of  $\mu$  and  $S$  has  $l''$  molecules.

The indexes of the basic transition relation identify the position of the reactants since solutions are sequences of molecules. In the case (*init*), the position is always 1 because the solution consists of one molecule. In the case of (*lifts*), the index is increased by the number of the molecules on the left, if any. The last case models a reaction: the solution is split into two parts  $S$  and  $T$  containing the reactants at positions  $l$  and  $l'$ , respectively. In the composite solution  $S, T$ , the reactants are at position  $l$  and  $l'' + l'$ , where  $l''$  is the number of molecules of  $S$ . For example let  $kM$  be  $\underbrace{M, \dots, M}_{k \text{ times}}$  and let  $\rho : H(1), H(1) \xrightarrow{\lambda} H(1^u), H(1^u)$  be the hydrogen gas reaction.

Then the following three transitions are possible

$$3H(1) \xrightarrow{\rho}_{1,2} 2H(1^x), H(1)$$

$$3H(1) \xrightarrow{\rho}_{1,3} H(1^x), H(1), H(1^x)$$

$$3H(1) \xrightarrow{\rho}_{2,3} H(1), 2H(1^x)$$

The basic transition relation is labelled by finite injective renamings. To clarify this point, consider the creation  $\varrho = Na(1^x + 2), Na(1^x + 2) \xrightarrow{10} Na(1^x + 2^y), Na(1^x + 2^y)$  (a bond is created between two sodium molecules provided they are already bound). Then take the solution  $Na[ion^0](1^z + 2), Na[ion^0](1^v + 2), Na[ion^1](1^z + 2)$ ,

$Na[ion^0](1^v + 2)$ . We derive the expected transition

$$Na[ion^0](1^z + 2), Na[ion^0](1^v + 2), Na[ion^1](1^z + 2), Na[ion^0](1^v + 2) \\ \xrightarrow{\varrho}_{1,3} Na[ion^0](1^z + 2^w), Na[ion^0](1^v + 2), Na[ion^1](1^z + 2^w), Na[ion^0](1^v + 2)$$

following a structured operational semantics approach [18]. Namely, we focus on the single reactants and lift the transitions to “,”-contexts. This is correct inasmuch as one records the instantiation of bonds in the left-hand sides of reactions with the actual names of the molecules: the two reactants must instantiate bonds in the same way. This is the reason why the first two molecules of the above solution cannot react with  $\varrho$ . More precisely,  $Na[ion^0](1^z + 2) \xrightarrow{\varrho_{L,\iota}}_1 Na[ion^0](1^z + 2^w)$ , where  $\iota = [x \mapsto z, y \mapsto w]$ , and  $Na[ion^0](1^v + 2) \not\xrightarrow{\varrho_{R,\iota}}_1$ .

Our final remarks regard the rule (*communications*). There are possibly infinitely many transitions  $S \xrightarrow{\rho,\iota}_l T$  because there are infinitely many renamings  $\iota$  which satisfy the conditions of the (*init*) rule. However this nondeterminism is removed when the reaction occurs because the created bonds have to be the least names not occurring in  $S$ , and because the renaming has to be order-preserving. Said otherwise, the relation  $\xrightarrow{\mu}_{l,l'}$ , which models the evolution of a solution, is finitely branching, while the auxiliary relation  $\xrightarrow{\mu}_l$  is not finitely branching. It is also worth to notice that there is no rule lifting a transition  $\xrightarrow{\mu}_{l,l'}$  to a context “,”: we use the associativity of , to partition a solution  $S$  into  $S', S''$  such that the reactants are in  $S'$  and  $S''$ .

The basic transition relation of the SPIM calculus requires few definitions:

- $M$  is true if  $M$  is a sequence of  $[x = x]$ ;
- $length(A_1(\tilde{u}_1) \mid \cdots \mid A_n(\tilde{u}_n))$  returns  $n$ ;
- $\tilde{x} : \tilde{\lambda} + \tilde{y} : \tilde{\lambda}'$  is the sequence  $z_1 : \lambda_1, \dots, z_n : \lambda_n$  where  $z_1, \dots, z_n$  are pairwise different names,  $\{z_1, \dots, z_n\} = \tilde{x} \cup \tilde{y}$ , and  $z_i : \lambda_i$  if either  $z_i : \lambda_i \in \tilde{y} : \tilde{\lambda}'$  or  $z_i \notin \tilde{y}$  and  $z_i : \lambda_i \in \tilde{x} : \tilde{\lambda}$ .
- $[(\tilde{x} : \tilde{\lambda})P]_{gc} = (\tilde{z} : \tilde{\lambda}')P$  such that  $y : \lambda''$  is in  $\tilde{z} : \tilde{\lambda}'$  if  $y \in \mathbf{fn}(P)$  and  $y : \lambda''$  is in  $\tilde{x} : \tilde{\lambda}$ .
- with an abuse of notation we lift a renaming  $\iota$  to a tuple of names or to a process by applying it pointwise.

Let  $\mathbf{bn}(\alpha)$  be  $\tilde{u}$  if  $\mu$  is either  $x(\tilde{u})$  or  $\bar{x}(\tilde{u} : \tilde{\lambda})$ ; it is  $\emptyset$  if  $\alpha = \bar{x}\tilde{u}$ .

**Definition 2.4** The basic transition relation of the SPIM calculus, written either  $\xrightarrow{\alpha}_{l,i,l',j}$  or  $\xrightarrow{\tau_\lambda}_{l,i,l',j}$  or  $\xrightarrow{\alpha}_{l,i}$ , is the least one satisfying the following rules:

- (init) let  $A(\tilde{u}) = \sum_{i \in I} M_i \alpha_i . P_i$  and let  $M_j \{\tilde{v}/\tilde{u}\}$  be true. If  $\alpha_j \{\tilde{v}/\tilde{u}\} = \bar{x}\tilde{w}$  then  $A(\tilde{v}) \xrightarrow{\bar{x}\tilde{w}}_{1,j} P_j$ ; if  $\alpha_j \{\tilde{v}/\tilde{u}\} = \bar{x}(\tilde{w} : \tilde{\lambda})$  then  $A(\tilde{v}) \xrightarrow{\bar{x}(\iota(\tilde{w}) : \tilde{\lambda})}_{1,j} \iota(P_j)$ ; if  $\alpha_j \{\tilde{v}/\tilde{u}\} = x(\tilde{w})$  then  $A(\tilde{v}) \xrightarrow{x(\iota(\tilde{w}))}_{1,j} \iota(P_j)$ , where  $\iota$  is an injective order-preserving renaming;
- (lifts) if  $P \xrightarrow{\alpha}_{l,i} P'$  and  $\mathbf{bn}(\alpha) \cap \mathbf{fn}(Q) = \emptyset$  and  $l' = length(Q)$ , then both

$P \mid Q \xrightarrow{l.i} P' \mid Q$  and  $Q \mid P \xrightarrow{l'.i} Q \mid P'$ ;

- (communications) let  $l'' = \text{length}(P)$ ,  $\lambda$  be the rate of  $x$ , and  $Q \xrightarrow{x(\tilde{u})}_{l'.i'} Q'$ . If  $P \xrightarrow{\tilde{x}\tilde{v}}_{l.i} P'$  then  $(\tilde{z} : \tilde{\lambda}')(P \mid Q) \xrightarrow{\tau_\lambda}_{l.i,l'+l'',i'} [(\tilde{z} : \tilde{\lambda}')(P \mid Q\{\tilde{v}/\tilde{u}\})]_{GC}$ ; if  $P \xrightarrow{\tilde{x}(\tilde{v}:\tilde{\lambda}'')}_{l.i} P'$  then  $(\tilde{z} : \tilde{\lambda}')(P \mid Q) \xrightarrow{\tau_\lambda}_{l.i,l'+l'',i'} [(\tilde{z} : \tilde{\lambda}' + \tilde{v} : \tilde{\lambda}'')(P' \mid Q'\{\tilde{v}/\tilde{u}\})]_{GC}$  where  $\tilde{v}$  are the least names not occurring in  $P \mid Q$ . Symmetrically when  $P$  performs an input and  $Q$  performs an output.

As for **nanok**, there is always at most one  $(\tilde{z} : \tilde{\lambda}')P'$  such that  $(\tilde{x} : \tilde{\lambda})P \xrightarrow{\tau_{\lambda''}}_{l,i,l',i'} (\tilde{z} : \tilde{\lambda}')P'$  because alpha-conversion is never considered in the basic transition relation, because created names are the least possible ones and because the renamings are order-preserving.

### 3 Encoding the nanok calculus into the SPIM calculus

The definition of the encoding of **nanok** calculus into SPIM calculus is presented in two steps. The first one defines an internal translation of **nanok** calculus that expands every bond into tuples of bonds. The bonds in the tuple are an over-approximation of the reactions that use the bond. We call these tuples of newly generated names *gangs*. The second step defines a translation from **nanok** (with gangs) to the SPIM calculus.

An example illustrating the encoding is postponed to the appendix.

#### 3.1 Gangs: a dedicated name for every reaction

In the following we use tuples that will be ordered as follows:  $(x_1, \dots, x_m) \leq (y_1, \dots, y_m)$  if and only if, for every  $i$ ,  $x_i \leq y_i$ . Let  $\varepsilon^m$  be a tuple of  $m$  elements  $\varepsilon$ .

**Definition 3.1** Let  $\mathcal{R} = \{\rho_i : L_i \xrightarrow{\lambda_i} R_i \mid i \in 1..n\}$  be a set of **nanok** reaction rules and let  $j$  be a bijective function that maps  $\varepsilon$  to  $\varepsilon^n$  and bonds to  $n$ -tuples of bonds such that if  $x \leq y$  then  $j(x) \leq j(y)$  (such a  $j$  exists because the set of names is countable).

The solution  $\llbracket S \rrbracket_j$  is  $S$  where every  $z$  being either a bond or  $\varepsilon$  is replaced by  $j(z)$ . The set of reactions  $\llbracket \mathcal{R} \rrbracket_j$  is  $\{\rho_i : \llbracket L_i \rrbracket_j \xrightarrow{\lambda_i} \llbracket R_i \rrbracket_j \mid i \in 1..n\}$ .

Namely  $\llbracket \mathcal{R} \rrbracket_j$  and  $\llbracket S \rrbracket_j$  are such that

- interfaces map sites to tuples of bonds of length  $n$  – a *gang*;
- two distinct tuples do not contain the same name;
- tuples preserve the order of bonds in  $\mathcal{R}$  and  $S$ .

We let  $\llbracket \iota \rrbracket_j = j \circ \iota \circ j^{-1}$  and  $\llbracket \mu \rrbracket_j$  be either  $\rho_L, \llbracket \iota \rrbracket_j$  or  $\rho_R, \llbracket \iota \rrbracket_j$ , depending on whether  $\mu$  is  $\rho_L, \iota$  or  $\rho_R, \iota$ . The correctness of the encoding of Definition 3.1 is stated in the following theorem.

**Theorem 3.2** (i) if  $S \xrightarrow{\mu}_t T$  then  $\llbracket S \rrbracket_j \xrightarrow{\llbracket \mu \rrbracket_j}_t \llbracket T \rrbracket_j$  (similarly for  $S \xrightarrow{\rho}_{l,l'} T$ );

- (ii) if  $\llbracket S \rrbracket_j \xrightarrow{\mu}_l T$  then there exists  $S'$  and  $\mu'$  such that:  $\llbracket S' \rrbracket_j = T$ ,  $S \xrightarrow{\mu'}_l S'$ , and  $\mu = \llbracket \mu' \rrbracket_j$  (similarly for  $\llbracket S \rrbracket_j \xrightarrow{\rho}_{l,\nu'} T$ ).

**Proof.** We first consider the transitions  $\xrightarrow{\mu}_l$  and then we lift the results to  $\xrightarrow{\rho}_{l,\nu'}$ .

We consider  $S \xrightarrow{\mu}_l T$ . We prove that the first item holds by induction on  $\text{length}(S)$ . If  $l = 1$ , the transition has been obtained by the *(init)* rule of the **nanok** calculus, and we remark that  $\llbracket S \rrbracket$  and  $\llbracket l \rrbracket$  fit the requirements of the *(init)* rule of the SPIM calculus thus  $\llbracket S \rrbracket \xrightarrow{\llbracket \mu \rrbracket}_l \llbracket T \rrbracket$ . If  $l > 1$ , the transition has been obtained by the *(lift)* rule of the **nanok** calculus and  $S = S_1, S_2$  with either  $S_1 \xrightarrow{\mu}_l T_1$  (and  $T = T_1, S_2$ ) or  $S_2 \xrightarrow{\mu}_{l-\text{length}(S_1)} T_2$  (and  $T = S_1, T_2$ ). By induction hypothesis we have that either  $\llbracket S_1 \rrbracket \xrightarrow{\llbracket \mu \rrbracket}_l \llbracket T_1 \rrbracket$  or  $\llbracket S_2 \rrbracket \xrightarrow{\llbracket \mu \rrbracket}_{l-\text{length}(S_1)} \llbracket T_2 \rrbracket$ . By application of the *(lift)* rule of the SPIM calculus we have that  $\llbracket S \rrbracket \xrightarrow{\llbracket \mu \rrbracket}_l \llbracket T \rrbracket$ .

The second item for the case  $\llbracket S \rrbracket \xrightarrow{\mu}_l T$  is proved similarly.

We now consider the first item for the case  $S \xrightarrow{\rho}_{l,\nu'} T$ . This transition is obtained by application of the *(communications)* rule, thus  $S = S_1, S_2$  where  $S_1$  and  $S_2$  performs two complementary transitions of the form  $\xrightarrow{\mu}_l$ . As we have already proved that the theorem holds for such transitions, we have that also  $\llbracket S_1 \rrbracket$  and  $\llbracket S_2 \rrbracket$  can perform complementary transitions. By application of the *(communications)* rule it is easy to see that also  $\llbracket S \rrbracket \xrightarrow{\llbracket \rho \rrbracket}_{l,\nu'} \llbracket T \rrbracket$ .

The second item for the case  $\llbracket S \rrbracket \xrightarrow{\rho}_{l,\nu'} T$  is proved similarly.  $\square$

### 3.2 From gangs to the SPIM calculus: agents as molecules

The second step of our translation encodes the **nanok** calculus with gangs of bonds into processes of SPIM. As discussed in the Introduction, we encode a species  $A$  by a parametric agent definition  $\widehat{A}(\tilde{x}) = P$ , whose parameters  $\tilde{x}$  represent the possible values of fields and sites of the molecules of that species. The body  $P$  is a choice with a branch for every reaction involving the species  $A$ . A molecule  $A[u](\sigma)$  is an invocation  $\widehat{A}(\{u, \sigma\})$ .

We begin by defining  $\{u, \sigma\}$ . Let  $\varepsilon$  and  $\neg\varepsilon$  be two distinguished channels. Then  $\{u, \sigma\}$  is equal to  $\{u\}_0, \{\sigma\}_1, \{\sigma\}_2$ , where

- $\{u\}_0$  yields the tuple of the values of the fields in  $u$ ;
- $\{\sigma\}_1$  yields the concatenation of the gangs in the range of  $\sigma$ ;
- $\{\sigma\}_2$  yields a tuple of length equal to the size of  $\text{dom}(\sigma)$ , whose  $i$ -th element is  $\varepsilon$  if all the element of the tuple  $\sigma(i)$  are  $\varepsilon$  and  $\neg\varepsilon$  if not.

Then we continue with a sequence of definitions. We assume given a set of  $n$  reactions  $\mathcal{R}$ .

- $[x_1, \dots, x_m]_u$  is the sequence of matches  $[x_i = u(i)]_{i \in \text{dom}(u)}$  ( $u$  is a partial map);
- $[x_1, \dots, x_m]_\sigma$  is the sequence of matches  $(M_i)_{i \in \text{dom}(\sigma)}$  where  $M_i = [x_i = \varepsilon]$  if all the element of the tuple  $\sigma(i)$  are  $\varepsilon$ , and  $M_i = [x_i = \neg\varepsilon]$  if not;

- $set(\tilde{x}, u)$  is the tuple where the  $i$ -th element is  $u(i)$  whenever  $i \in \text{dom}(u)$ , it is the  $i$ -th element of  $\tilde{x}$ , otherwise;
- $set_1(\tilde{x}, \sigma)$  is the tuple where the element  $n * (i - 1) + j$  is the  $j$ -th element of the tuple  $\sigma(i)$ , when  $i \in \text{dom}(\sigma)$ , and  $x_{n*(i-1)+j}$  otherwise;
- $set_2(\tilde{x}, \sigma)$  is the tuple where the element  $n * (i - 1) + j$  is  $\varepsilon$  if  $\sigma(i) = \varepsilon^n$ , it is  $\neg\varepsilon$  if  $i \in \text{dom}(\sigma)$  and  $\sigma(i) \neq \varepsilon^n$ , and it is  $x_{n*(i-1)+j}$  otherwise;
- $proj(\tilde{x}, a)$  is the tuple  $(x_{n*(a-1)+i})_{i \leq n}$  and  $proj(\tilde{x}, a, i)$  is  $x_{n*(a-1)+i}$ ;
- if  $A[u](\sigma), B[v](\phi) \xrightarrow{\lambda} A[u'](\sigma'), S \in \mathcal{R}$  then both  $A[u](\sigma) \xrightarrow{\lambda} A[u'](\sigma') \in_L \mathcal{R}$  and  $B[v](\phi) \xrightarrow{\lambda} S \in_R \mathcal{R}$ ;
- If  $\rho$  is a creation,  $\text{CR}(\rho, \mathcal{R})$  is a sequence  $(x_1 : \lambda_1, \dots, x_m : \lambda_m)$  where every subsequence  $(x_{i \times n} : \lambda_{i \times n}, x_{i \times n + 1} : \lambda_{i \times n + 1}, \dots, x_{i \times n + n - 1} : \lambda_{i \times n + n - 1})$  corresponds to the  $i$ -th bond created by  $\rho$  and  $\lambda_{i \times n}, \dots, \lambda_{i \times n + n - 1}$  are the rates of the reactions in  $\mathcal{R}$ .

Every preliminary notation is in place for the definition of the encoding from **nanok** with gangs to SPIM.

**Definition 3.3** Let  $\mathcal{R}$  be a set of  $n$  reactions in **nanok**. The SPIM agent corresponding to the species  $A$  is:

$$\begin{aligned} \hat{A}(\tilde{x}, \tilde{y}, \tilde{z}) = & \sum_{\rho: A[u](\sigma) \xrightarrow{\lambda} A[u'](\sigma') \in_L \mathcal{R}} [\tilde{x}]_u [\tilde{z}]_\sigma \alpha_{\rho, L} \cdot P_{\rho, L} \\ & + \sum_{\rho: A[u](\sigma) \xrightarrow{\lambda} S \in_R \mathcal{R}} [\tilde{x}]_u [\tilde{z}]_\sigma \alpha_{\rho, R} \cdot P_{\rho, R} \end{aligned}$$

where the length of  $\tilde{x}$  is the number of fields of  $A$ , and the lengths of  $\tilde{y}$  and  $\tilde{z}$  are the number of sites of  $A$  times  $n$ . In addition:

- if  $\rho$  is a creation with an empty set of bonds in the left-hand side then  $\alpha_{\rho, L} = \bar{\rho}(u : \lambda)$  and  $\alpha_{\rho, R} = \rho(\tilde{u})$  and  $(u : \lambda) = \text{CR}(\rho, \mathcal{R})$ ;
- if  $\rho$  is a creation with a bond  $x$  in the left-hand side then  $\alpha_{\rho, L} = \overline{proj(\tilde{y}, a, i)}(u : \lambda)$  and  $\alpha_{\rho, R} = proj(\tilde{y}, a, i)(\tilde{u})$ , where  $a$  is the site of  $A$  bound by  $x$ ,  $i$  is the index of  $\rho$  in  $\mathcal{R}$  and  $(u : \lambda) = \text{CR}(\rho, \mathcal{R})$ ;
- if  $\rho$  is a destruction with a bond  $x$  in the left-hand side then  $\alpha_{\rho, L} = \overline{proj(\tilde{y}, a, i)}( )$  and  $\alpha_{\rho, R} = proj(\tilde{y}, a, i)( )$ , where  $a$  is the site of  $A$  bound by  $x$  and  $i$  is the index of  $\rho$  in  $\mathcal{R}$ ;
- if  $\rho$  is an exchange with an empty set of bonds in the left-hand side or with a bond occurring once and in  $A$  then  $\alpha_{\rho, L} = \bar{\rho}\tilde{u}$  and  $\alpha_{\rho, R} = \rho(\tilde{u})$ , where  $\tilde{u}$  is either empty, if there is no bond in the left-hand side, or  $proj(\tilde{y}, A, a)$  if the site with the bond is  $a$ ;
- if  $\rho$  is an exchange with a bond  $x$  shared by the reactants then one defines  $\alpha_{\rho, L} = proj(\tilde{y}, A, a, i)(\tilde{u})$  and  $\alpha_{\rho, R} = proj(\tilde{y}, A, a, i)(\tilde{u})$ , where  $a$  is the site of  $A$  bound by  $x$ ,  $i$  is the index of  $\rho$  in  $\mathcal{R}$  and  $\tilde{u}$  is either empty, if there is no bond in the left-hand side apart  $x$ , or  $proj(\tilde{y}, A, a')$  if  $A$  has a further bond on the site  $a'$ .

As regards continuations,  $P_{\rho,L} = \widehat{A}(\text{set}(\tilde{x}, u'), \text{set}_1(\tilde{y}, \sigma'), \text{set}_2(\tilde{z}, \sigma'))$  and  $P_{\rho,R}$  is either  $\mathbf{0}$ , if  $S = -$ , or  $\widehat{A}(\text{set}(\tilde{x}, u'), \text{set}_1(\tilde{y}, \sigma'), \text{set}_2(\tilde{z}, \sigma')), \widehat{C}_1(\{\{u_1\}\}_1, \{\{\phi_1\}\}_2, \{\{\phi_1\}\}_3), \dots, \widehat{C}_n(\{\{u_h\}\}_1, \{\{\phi_h\}\}_2, \{\{\phi_h\}\}_3)$  if  $S = A[u'](\sigma'), C_I[v_1](\phi_1), \dots, C_n[v_n](\phi_n)$ .

The encoding of a **nano** $\kappa$  calculus solution with gangs is:

$$\{[A_I[u_1](\sigma_1), \dots, A_m[u_m](\sigma_m)]\} \triangleq (\delta_S)(\widehat{A}_1\{\{u_1, \sigma_1\}\}, \dots, \widehat{A}_m\{\{u_m, \sigma_m\}\})$$

where  $\delta_S$  is the minimal set that contains

- $(\rho : \lambda)$ , if  $\rho$  has no bond between reactants and has rate  $\lambda$ ,
- $(x_{n \times (i-1)} : \lambda_1, \dots, x_{n \times (i-1) + n-1} : \lambda_n)$ , if there is an agent invocation  $\widehat{A}\{\{u_1, \sigma_1\}\}$  and  $\{\{\sigma_1\}\}_2 = (\dots, x_{n \times (i-1)}, \dots, x_{n \times (i-1) + n-1}, \dots)$ , with  $x_i \neq \varepsilon$  and  $\lambda_1, \dots, \lambda_n$  being the rates of the reactions in  $\mathcal{R}$ .

To illustrate our encoding we chose a toy-modelling of the transcription of a gene. There are four species:

- *Gn* models a gene. It has one field *tr* and two sites *pr* and *rnap*; *tr* is 1 when the gene is being transcribed by the RNA polymerase and 0 if not; *pr* and *rnap* are used to link to *Pr* and *RNAp*, respectively;
- *Pr* models the various promoter-sequences of the gene. It has one field *act* and two sites *rnap* and *gn*. The activation of the promoters by the transcription-factors is represented by switching *act* from 0 to 1. The sites *rnap* and *gn* are used to link *RNAp* and *Gn*, respectively;
- *RNAp* models the RNA polymerase. It has one field *act* that is set to 1 when the molecule is activated by the complexations with the promoters and the transcription factors, it is set to 0 otherwise. It has a site *link* that may be bound either to *Pr* or to *Gn*, according to the stage of the transcription;
- The species *mRNA* models the RNA messenger corresponding to the gene. It has neither fields nor sites.

There are three reactions. The creation  $\rho_1$  models the binding of the RNA polymerase to the promoters (between sites *rnap* of *Pr* and *link* of *RNAp*) and their activation (update of the fields *act* of *Pr* and *RNAp*). The exchange  $\rho_2$  models the movement of the RNA-polymerase to the gene and the beginning of the transcription itself (update of the field *tr* of *Gn* to 1). The destruction  $\rho_3$  models the termination of the transcription.

$$\begin{aligned}
\rho_1 : & Pr[act^0](rnap) , RNAP[act^0](link) \xrightarrow{\lambda_1} Pr[act^1](rnap^x) , RNAP[act^1](link^x) \\
\rho_2 : & Pr[act^1](rnap^y + gn^x) , Gn[tr^0](pr^x + rnap) \\
& \xrightarrow{\lambda_2} Pr[act^0](rnap + gn^x) , Gn[tr^1](pr^x + rnap^y) \\
\rho_3 : & RNAP[act^1](link^x) , Gn[tr^1](rnap^x) \\
& \xrightarrow{\lambda_3} RNAP[act^0](link) , Gn[tr^0](rnap) , mRNA
\end{aligned}$$

The encoding of this **nanok** systems yields the following four recursive definitions in SPIM. We notice that, in the encoding of  $Gn$ , the parameters  $p\rho_1, p\rho_2, p\rho_3, ?p$  correspond to the gang of the site  $pr$  (the three former are yielded by  $\{[.]_1$  and the latter by  $\{[.]_2$ ), similarly the parameters  $r\rho_1, r\rho_2, r\rho_3, ?r$  correspond to the gang of the site  $rnap$ .

$$\begin{aligned}
& \widehat{Gn}(tr, p\rho_1, p\rho_2, p\rho_3, ?p, r\rho_1, r\rho_2, r\rho_3, ?r) \triangleq \\
& [tr = 0, ?p = \neg\epsilon, ?r = \epsilon] p\rho_2(r_1, r_2, r_3) . \widehat{Gn}(1, p\rho_1, p\rho_2, p\rho_3, ?p, r_1, r_2, r_3, \neg\epsilon) \\
& + [tr = 1, ?r = \epsilon] r\rho_3() . \widehat{Gn}(1, p\rho_1, p\rho_2, p\rho_3, ?p, \epsilon, \epsilon, \epsilon, \epsilon)
\end{aligned}$$

In the encoding of  $Pr$ , the parameters  $r\rho_1, r\rho_2, r\rho_3, ?r$  correspond to the gang of the site  $rnap$  and the parameters  $g\rho_1, g\rho_2, g\rho_3, ?g$  correspond to the gang of the site  $gn$ .

$$\begin{aligned}
& \widehat{Pr}(act, r\rho_1, r\rho_2, r\rho_3, ?r, g\rho_1, g\rho_2, g\rho_3, ?g) \triangleq \\
& [act = 0, ?r = \epsilon] \overline{\rho_1}(r_1 : \lambda_1, r_2 : \lambda_1, r_3 : \lambda_3) . \widehat{Pr}(1, r_1, r_2, r_3, \neg\epsilon, g\rho_1, g\rho_2, g\rho_3, ?g) \\
& + [act = 1, ?r = \neg\epsilon, ?g = \neg\epsilon] \overline{g\rho_2}(r\rho_1, r\rho_2, r\rho_3) . \widehat{Pr}(0, \epsilon, \epsilon, \epsilon, \epsilon, g\rho_1, g\rho_2, g\rho_3, ?g)
\end{aligned}$$

In the encoding of the species  $RNAP$ , the parameters  $l\rho_1, l\rho_2, l\rho_3, ?l$  correspond to the gang of the site  $link$ .

$$\begin{aligned}
& \widehat{RNAP}(act, l\rho_1, l\rho_2, l\rho_3, ?l) \triangleq \\
& [act = 0, ?l = \epsilon] \rho_1(r_1, r_2, r_3) . \widehat{RNAP}(1, r_1, r_2, r_3, \neg\epsilon) \\
& + [act = 1, ?l = \neg\epsilon] \overline{l\rho_3} . (\widehat{RNAP}(0, \epsilon, \epsilon, \epsilon, \epsilon) \mid \widehat{mRNA})
\end{aligned}$$

Since the molecule  $mRNA$  is not involved in any reaction, the corresponding process is:

$$\widehat{mRNA}() \triangleq 0$$

Finally, the encoding of the solution  $Pr[act^0](rnap+gn^x), Gn[tr^0](prev^x+rnap)$ ,

$RN\mathit{Ap}[act](link)$  is the term

$$\begin{aligned}
 &(\rho_1 : \lambda_1, pr_1 : \lambda_1, pr_2 : \lambda_2, pr_3 : \lambda_3)(\widehat{Pr}(0, \epsilon, \epsilon, \epsilon, \epsilon, pr_1, pr_2, pr_3, \neg\epsilon) \\
 &\quad | \widehat{Gn}(0, pr_1, pr_2, pr_3, \neg\epsilon, \epsilon, \epsilon, \epsilon, \epsilon) \\
 &\quad | \widehat{RN\mathit{Ap}}(act, \epsilon, \epsilon, \epsilon, \neg\epsilon) \\
 &\quad )
 \end{aligned}$$

Remarkably, our encoding retains some interesting compositional properties. First of all, consider two solutions  $S$  and  $S'$  and their encodings  $(\delta_S)P$  and  $(\delta_{S'})P'$ , respectively. The encoding of the parallel composition  $S|S'$  is  $(\delta_S \cup \delta_{S'})P|P'$  where  $(\delta_S \cup \delta_{S'})$  is the union of the name declarations  $\delta_S$  and  $\delta_{S'}$  assuming that local names are disjoint (this assumption can be guaranteed exploiting  $\alpha$ -conversion). Furthermore, if we add a new reaction rule to a system that we have already encoded, there are very few changes to be made. First, we need to expand the length of the gangs by one. After, we have to add one line to the definitions corresponding to the two species occurring in the left hand side of the new reaction, in order to describe this additional behaviour for the molecules belonging to these two species.

The next theorem states the correctness of  $\{\llbracket \cdot \rrbracket\}$ . If  $A$  is the  $l$ -th molecule in  $S$  and if  $\rho$  corresponds to the  $i$ -th branch of the choice in  $\widehat{A}$ , we let  $\{\llbracket l \rrbracket\}_\rho$  be the pair  $(l, i)$ . We also let  $\{\llbracket \rho_L, i \rrbracket\}$  and  $\{\llbracket \rho_R, i \rrbracket\}$  to be respectively  $\alpha_{\rho, L}$  and  $\alpha_{\rho, R}$  as defined in Definition 3.3.

**Theorem 3.4** (i) If  $S \xrightarrow{\mu}_l T$  then  $\{\llbracket S \rrbracket\} \xrightarrow{\{\llbracket \mu \rrbracket\}_{\{\llbracket l \rrbracket\}_\rho}} \{\llbracket T \rrbracket\}$  (similarly for  $S \xrightarrow{\rho}_{l, i'} T$ );  
(ii) if  $\{\llbracket S \rrbracket\} \xrightarrow{\mu}_{l, i} T$  (resp.  $\{\llbracket S \rrbracket\} \xrightarrow{\rho}_{l, i, i', j} T$ ) then there exist  $S'$  and  $\mu'$  such that:  
 $\{\llbracket S' \rrbracket\} = T$ ,  $S \xrightarrow{\mu'}_{l'} S'$ , and  $\mu = \{\llbracket \mu' \rrbracket\}$  (similarly for  $\{\llbracket S \rrbracket\} \xrightarrow{\rho}_{l, i, i', i'} T$ ).

**Proof.** The theorem is proven similarly to theorem 3.2 (this because we have called the three items of the definition 2.4 of the basic transition relation of the SPIM calculus with the same names (*init*), (*lifts*), and (*communications*) of the corresponding definition 2.3 for the **nanok** calculus).  $\square$

## 4 The stochastic collective semantics

The basic transition relation we considered keeps track of all the possible transitions that the molecules in a solution can perform. However, some of these transitions are somehow “equivalent” because, for instance, they have the same source and the targets are indistinguishable. This is the case when the solution contains several copies of a molecule and the reaction is an homeodimerization, i.e. two identical molecules get bound.

The following collective semantics merges “equivalent” transitions into one transition with an associated rate obtained as the sum of the rates of the merged transitions. It uses the structural equivalence to formalize the indistinguishability of solutions.



**Definition 4.1** The structural equivalence of the **nanok** calculus, noted  $\equiv$ , is the least equivalence satisfying the following rules (solutions are already identified by associativity of “,”):

- (i)  $S, T \equiv T, S$ ;
- (ii)  $S \equiv T$  if there exists an injective renaming  $\iota$  on bonds such that  $S = \iota(T)$ .

The structural equivalence of the SPIM calculus, that, with an abuse of notation, we also note  $\equiv$ , is the least equivalence satisfying the following rules:

- $P|Q \equiv Q|P$
- if  $P$  is  $\alpha$ -equivalent to  $Q$  then  $P \equiv Q$

In order to give a unique definition of the collective semantics, we introduce a few notations. The letters  $F, G$  are used to range over solutions or processes; we assume that transitions of the basic transition systems have shape  $\xrightarrow{\rho}_{\partial}$ , where  $\partial$  is a pair (we are considering evolutions of closed systems). Let also

- $next(F) = \{((\rho, \partial), G) \mid F \xrightarrow{\rho}_{\partial} G\}$ ;
- $next_{\infty}(F) = \{((\rho, \partial), G) \mid F \xrightarrow{\rho}_{\partial} G \text{ and } rate(\rho) = \infty\}$
- $F$  has finite rates if and only if  $next_{\infty}(F) = \emptyset$
- given a set  $\mathcal{F}$  of pairs  $(X, T)$ , where  $T$  is a term, let  $[\mathcal{F}]_G$ , where  $G$  is a fixed term, be  $\{(X, G') \mid (X, G') \in \mathcal{F} \text{ and } G' \equiv G\}$ ;
- $can(\mathcal{F})$  is defined over sets of pairs  $(X, G)$  (the second element is a term, the first one is left unspecified), such that the terms occurring as second element of the pairs are all structurally equivalent. It returns a term  $G'$  such that there is  $X$  with  $(X, G') \in \mathcal{F}$ .

**Definition 4.2** [Stochastic collective transition relation] The stochastic transition relation  $\mapsto$  induced by a basic transition relation  $\xrightarrow{\rho}_{\partial}$  ( $\partial$  is a pair of indexes) and structural equivalence  $\equiv$  on a language is the least relation satisfying the following rules:

- if  $F \xrightarrow{\rho}_{\partial} G$  and  $rate(\rho) = \infty$  then  $F \mapsto^{\infty} can([next_{\infty}(F)]_G)$ ;
- if  $F \xrightarrow{\rho}_{\partial} G$  and  $F$  has finite rates then  $F \mapsto^{\lambda} can([next(F)]_G)$ , where

$$\lambda = \sum_{((\rho, \partial), G') \in [next(F)]_G} rate(\rho)$$

The correctness result of the collective transition relation is stated below.

**Theorem 4.3**  $S \mapsto^{\lambda} T$  in **nanok** if and only if there exists  $P$  such that  $\{\llbracket S \rrbracket_j\} \mapsto^{\lambda} P$  and  $P \equiv \{\llbracket T \rrbracket_j\}$  in SPIM.

Our correctness notion corresponds to the subcase of the strong stochastic bisimulation [3] where the bisimulation relation is a bijection.

## 5 Future works

Our current interests are mainly about simulators and analysis tools for SPIM calculus. In fact, this contribution allows us to simulate **nanok** systems. However, the same encoding also makes it possible to model-check **nanok** formalizations in the PRISM platform [1], since it supports verifications of probabilistic and stochastic extensions of  $\pi$ -calculus [15]. More precisely, it should be possible to wire our encoding from the **nanok** calculus to SPIM – a subset of stochastic  $\pi$ -calculus – with the implementation in [15]. This requires that reactions do not create molecules, otherwise the state space would be infinite, and PRISM cannot handle such systems.

There are two questions to bother with. Firstly, our encoding uses polyadic communications, which is still not considered in [15]. However this should be one of the next extensions of this work. The second issue is more problematic. A relevant constraint for the efficiency of the encoding in [15] is the absence of name creations within agent definition. This is not the case for our encoding, because agents may perform bounded outputs. Yet, in **nanok** subsystems where the creation of new molecules is finite, the number of names used at every stage of the computation is finite. So, a clever algorithm might compute this number statically (an over-approximation is  $k \times h$ , where  $k$  is the maximal number of molecules and  $h$  is the maximal length of the arguments of an agent) and use a garbage-collection mechanism to recycle names. This should allow the static allocation of variables in the PRISM language to handle all the private names.

## References

- [1] L.de Alfaro, M.Z.Kwiatkowska, G.Norman, D.Parker, and R.Segala. Symbolic model checking of probabilistic processes using mtbdds and the kronecker representation. In *TACAS*, pages 395–410, 2000.
- [2] R. Barbuti, A. Maggiolo-Schettini, P. Milazzo, and A. Troina. A calculus of looping sequences for modelling microbiological systems. *Fundamenta Informaticae*, 72(1-3):21–35, 2006.
- [3] M. Bernardo. A survey of markovian behavioral equivalences. In *Proc. of International School on Formal Methods for the Design of Computer, Communication, and Software Systems 2007*, volume 4486 of *LNCS*, pages 180–219, 2007.
- [4] M. Calder, S. Gilmore, and J. Hillston. Modelling the influence of rkip on the erk signalling pathway using the stochastic process algebra pepa. In *Transactions on Computational Systems Biology VII*, volume 4230 of *Lecture Notes in Computer Science*, pages 1–23, 2006.
- [5] L. Calzone, F. Fages, and S. Soliman. Biocham: an environment for modeling biological systems and formalizing experimental knowledge. *Bioinformatics*, 22(14):1805–1807, 2006.
- [6] L. Cardelli. On process rate semantics. *Theoretical Computer Science*, 391(3):190–215, 2008.
- [7] L. Cardelli and A. Phillips. A corret abstract machine for the stochastic pi-calculus. In *Proc. of Workshop on Concurrent Models in Molecular Biology*, 2004.
- [8] L. Cardelli and G. Zavattaro. On the computational power of biochemistry. In *Proc. of Algebraic Biology 2008*, volume to appear of *LNCS*, 2008.
- [9] F. Ciocchetta and J. Hillston. Bio-pepa: a framework for modelling and analysis of biological systems. *Theoretical Computer Science*, to appear.
- [10] A. Credi, M. Garavelli, C. Laneve, S. Pradaliere, S. Silvi, and G. Zavattaro. Modelization and simulation of nano devices in the nano-k calculus. In *Proc. of Computational Methods in Systems Biology 2007*, volume 4695 of *LNCS*, pages 168–183, 2007.

- [11] Pierre-Louis Curien, Vincent Danos, Jean Krivine, and Min Zhang. Computational self-assembly. *Theoretical Computer Science*, 404(1-2):61–75, 2008.
- [12] V. Danos, J. Feret, W. Fontana, and J. Krivine. Scalable simulation of cellular signaling networks. In *Proc. of Asian Symposium on Programming Languages and Systems 2007*, volume 4807 of *LNCS*, pages 139–157, 2007.
- [13] V. Danos and C. Laneve. Formal molecular biology. *Theoretical Computer Science*, 325(1):69–110, 2004.
- [14] Stephen Gilmore and Jane Hillston. The pepa workbench: a tool to support a process algebra-based approach to performance modelling. In *Proceedings of the 7th international conference on Computer performance evaluation : modelling techniques and tools*, pages 353–368, Secaucus, NJ, USA, 1994. Springer-Verlag New York, Inc.
- [15] D.Parker G.Norman, C.Palamidessi and P.Wu. Model-checking probabilistic and stochastic extensions of the pi-calculus. In *IEEE Transactions on Software engineering*, 2007.
- [16] J.Heath, M.Z.Kwiatkowska, G.Norman, D.Parker, and O.Tymchyshyn. Probabilistic model checking of complex biological pathways. In *CMSB*, pages 32–47, 2006.
- [17] C. Laneve and A.Vitale. Expressivness in the  $\kappa$ -family. In *Proc. of MFPS, ENTCS*, 2008.
- [18] G. D. Plotkin. A Structural Approach to Operational Semantics. Technical Report DAIMI FN-19, University of Aarhus, 1981.
- [19] Corrado Priami. Stochastic pi-calculus. *Computer Journal*, 38(7):578–589, 1995.
- [20] Corrado Priami, Aviv Regev, Ehud Shapiro, and William Silverman. Application of a stochastic name-passing calculus to representation and simulation of molecular processes. *Information Processing Letters*, 80:25–31, 2001.