



**HAL**  
open science

## Metric Reasoning About $\lambda$ -Terms: The Affine Case

Raphaëlle Crubillé, Ugo Dal Lago

► **To cite this version:**

Raphaëlle Crubillé, Ugo Dal Lago. Metric Reasoning About  $\lambda$ -Terms: The Affine Case. LICS 2015, Jul 2015, Kyoto, Japan. 10.1109/LICS.2015.64 . hal-01231814

**HAL Id: hal-01231814**

**<https://inria.hal.science/hal-01231814>**

Submitted on 20 Nov 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.

# Metric Reasoning About $\lambda$ -Terms: The Affine Case

Raphaëlle Crubillé

ENS Lyon

Email: raphaelle.crubille@ens-lyon.fr

Ugo Dal Lago

Università di Bologna & INRIA

Email: ugo.dallago@unibo.it

**Abstract**—Terms of Church’s  $\lambda$ -calculus can be considered *equivalent* along many different definitions, but context equivalence is certainly the most direct and universally accepted one. If the underlying calculus becomes probabilistic, however, equivalence is too discriminating: terms which have totally unrelated behaviours are treated the same as terms which behave very similarly. We study the problem of evaluating the *distance* between affine  $\lambda$ -terms. A natural generalisation of context equivalence, is shown to be characterised by a notion of *trace* distance, and to be bounded from above by a coinductively defined distance based on the Kantorovich metric on distributions. A different, again fully-abstract, tuple-based notion of trace distance is shown to be able to handle nontrivial examples.

## I. INTRODUCTION

Probabilistic models are formidable tools when abstracting the behaviour of complicated, intractable systems by simpler ones, at the price of introducing uncertainty. But there is more: randomness can be seen as *a way to compute*; in modern cryptography, as an example, having access to a source of uniform randomness is essential to achieve security in an asymmetric setting [14]. Other domains where probabilistic models play a key role include machine learning [24], robotics [27], and linguistics [21].

Probabilistic models of computation have been studied not only directly, but also through concrete or abstract programming languages, which most often are extensions of their deterministic siblings. Among the many ways probabilistic choice can be captured in programming, the simplest one consists in endowing the language of programs with an operator modelling the flipping of a fair coin. This renders program evaluation a probabilistic process, and under mild assumptions the language becomes universal for probabilistic computation. Particularly fruitful in this sense has been the line of work on the functional paradigm, both at a theoretical [17], [26], [23] and at a more practical level [15].

In presence of higher-order functions, program equivalence can be captured by so-called *context equivalence*: two programs  $M$  and  $N$  are considered equivalent if they behave the same no matter how the environment interacts with them: for every context  $\mathcal{C}$ , it holds that  $\text{Obs}(\mathcal{C}[M]) = \text{Obs}(\mathcal{C}[N])$ . However, this definition has the drawback of being based on an universal quantification over *all* contexts: showing that two programs are equivalent requires considering their interaction with every possible context. The problem of giving handier characterisations of context equivalence can be approached in

many different ways. As an example, coinductive methodologies for program equivalence have been studied thoroughly in deterministic [1], [25] and non-deterministic [19] computation, with new and exciting results appearing recently also for probabilistic languages: applicative bisimilarity, a coinductively defined notion of equivalence for functional programs, has been shown to be sound, and sometime even fully abstract, for probabilistic  $\lambda$ -calculi [6], [4].

In a probabilistic setting, however, equivalences are too discriminating if defined as above. Indeed, two programs are equivalent if their probabilistic behaviour is *exactly* the same (in every context). The actual value of probabilities in a probabilistic model often comes from statistical measurements, and should be considered more as an approximation to the actual probability law. Consequently, we would like to compare programs by appropriately reflecting small variations in them. Another scenario in which a richer, more informative way of comparing programs is needed is cryptography, where a central notion of equivalence, called *computational indistinguishability* [13] is indeed based on statistical distance rather than equality: the adversary *can* win the game, but with a *small* probability. Summing up, equivalences should be refined into metrics, and this is the path we will follow in this paper.

In probabilistic  $\lambda$ -calculi, the notion of observation  $\text{Obs}(\cdot)$  is quantitative: it is either the *probability* of convergence to a certain observable base value (e.g. the empty string), or the probability of convergence *tout court*. One can then easily define a notion of *context distance* as the maximal distance contexts can achieve when separating two terms:

$$\delta^{\text{ctx}}(M, N) = \sup_{\mathcal{C}} |\text{Obs}(\mathcal{C}[M]) - \text{Obs}(\mathcal{C}[N])|.$$

This looks very close to computational indistinguishability, except for the absence of a security parameter: a scheme is secure if the advantage of any adversary in a given game (e.g., consisting in distinguishing between the case where the scheme is used, and the case where it is replaced by a truly random process) is “small” (e.g., negligible). Again, however, we find ourselves in front of a definition which risks to be useless in proofs, given that all contexts must be taken into account. But how difficult evaluating the distance between concrete higher-order terms really is? Are there ways to alleviate the burden of dealing with all contexts, like for equivalences? These are the questions we address in this

paper, and which have to the authors’ knowledge not been investigated before.

As we will discuss in Section II below, finding handier characterisations of the context distance poses challenges which are simply different (and often harder) than the ones encountered in context equivalence. In particular, the context distance tends to *trivialise* and, perhaps worse, naively applying the natural generalisation of techniques known for equivalence is bound to lead to unsound methodologies. Indeed, one immediately realises that the number of times contexts access their argument is a crucial parameter, which must necessarily be dealt with. This is the reason why we work with an affine  $\lambda$ -calculus in this paper: this is a necessary first step, but also points to the right way to tame the general, non-linear case, as we hinge in Section VI-C.

An extended version of this paper with more details is available [5]. The authors are partially supported by the ANR project 12IS02001 PACE.

### Contributions

We introduce in this paper *three* distinct notions of distance for terms in an untyped, probabilistic, and affine  $\lambda$ -calculus. The first one is a notion of *trace distance*, in which terms are faced with *linear* tests, i.e. sequences of arguments. The distance between two terms is then defined as the greatest separation any linear test achieves. The first results of this paper are the non-expansiveness of the trace distance, which implies (given that any linear test can easily be implemented by an affine context) that the trace and context distances coincide. This is the topic of Section IV below.

Section V, instead, focuses on another notion of distance, which is coinductively defined following the well-known Kantorovich metric [18] for distributions of states in any labelled Markov chain (LMC in the following), and that we dub the *bisimulation distance*. This second notion of a distance is not only at least as discriminating as the trace distance, which is well expected, but non-expansive itself. This is proved by a variation on the Howe’s method [16], a well-known technique for proving that bisimilarity is a congruence in an higher-order setting, and which has never been used for metrics before. On the other hand, the bisimulation distance *does not* coincide with the context distance, a fact that we do not only prove by giving a counterexample, but that we justify by relying on a test-based characterisation of the bisimulation distance known from the literature.

For the sake of simplicity, the trace and bisimulation distances are analysed on a purely applicative  $\lambda$ -calculus, keeping in mind that pairs could be very easily handled, and can even be encoded in the applicative fragment, as discussed in Section IV-C. The presence of pairs, however, allows us to form very interesting examples of distance problems, one of which will drive us throughout the paper but unfortunately turns out hard to handle neither by the trace distance nor by the bisimulation distance. This is the starting point for the third notion of distance introduced in this paper, which is the subject of Section VI, and which we call the *tuple distance*.

Our third notion of distance can be proved to coincide with the trace distance, and thus with the context distance. But this is not the end of the story: in the tuple distance, not a single but *many* terms are compared, and this makes the distance between concrete terms much easier to evaluate: interaction is somehow internalised. In particular, our running example can be handled quite easily. The way the tuple distance is defined makes it adaptable to non-affine calculi, a topic which is outside the scope of this paper, but which we briefly discuss in Section VI-C.

### Related Work

This is definitely not the first work on metrics for probabilistic systems: notions of coinductively defined metrics for LMCs, as an example, have been extensively studied (e.g. [10], [9], [28]). There have been, to our knowledge, not so many investigations on the meaning of metrics for concrete programming languages [12], and almost nothing on metric for higher-order languages.

If the key property notions of *equivalences* are required to satisfy consists in being *congruences*, the corresponding property for metrics has traditionally been taken as *non-expansiveness*. Indeed, many results from the literature (e.g. [10], [22]) have precisely the form of non-expansiveness results for metrics defined in various forms. The underlying language, however, invariably take the form of a process algebra without any higher-order feature. The work by Gebler, Tini, and co-authors shows that one could go beyond non-expansiveness and towards uniform continuity [12] but, again, higher-order functions remain out of scope.

Notions of *equivalence* for various forms of probabilistic  $\lambda$ -calculi have also been extensively studied, starting from the pioneering work by Plotkin and Jones [17], down to recent results on probabilistic applicative bisimulation [6], [4], logical relations [3], and probabilistic coherent spaces [7], [11]. None of the works above, however, go beyond equivalences and deals with notions of distance between terms.

## II. THE ANATOMY OF A DISTANCE

In this section, we describe the difficulties one encounters when trying to characterise the context distance with either bisimulation or trace metrics.

Suppose we have two terms  $M$  and  $N$  of boolean type written in a probabilistic  $\lambda$ -calculus. As such,  $M$  and  $N$  do not evaluate to a value in the domain of booleans but to a *distribution* over the same domain.  $M$  evaluates to the distribution assigning **true** probability 1, while  $N$  evaluates to the uniform distribution over booleans, (i.e. the distribution which attributes probability  $\frac{1}{2}$  to both **true** and **false**). Figure 1 depicts the relevant fragment of a LMC, whose induced notion of probabilistic bisimilarity has been proved to be sound for context equivalence [4].  $M$  and  $N$  are not bisimilar. Indeed, **true** and **false** are trivially not bisimilar, while  $M$  and  $N$  go to equivalent states with different probabilities. The two terms are non-equivalent also contextually. But what *should be* the distance between  $M$  and  $N$ ?

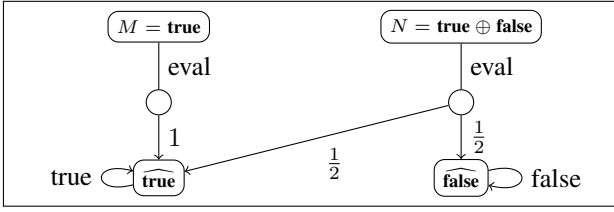


Fig. 1.  $M$  and  $N$  as States of a LMC.

For the moment, let us forget about the context distance, and concentrate on the notions of distance for LMCs we mentioned in Section I. In all cases we are aware of, we obtain that  $M$  and  $N$  are at distance  $\frac{1}{2}$ . As an example, if we consider a trace metric, we have to compare the success probability of linear tests, starting from  $M$  and  $N$ . More precisely, the tests of interest with respect to these two terms are:

$$t := \text{eval}; \quad s := \text{eval} \cdot \text{true}; \quad r := \text{eval} \cdot \text{false}.$$

Since neither  $M$  nor  $N$  has a non-zero divergence probability, they both pass the test  $t$  with probability 1. The success probability of the test  $s$  corresponds to the probability of evaluating to **true**: it is 1 for  $M$  and  $\frac{1}{2}$  for  $N$ . Similarly, the success probability of  $r$  corresponds to the probability to obtain **false** after evaluation: it is 0 for  $M$  and  $\frac{1}{2}$  for  $N$ . So we can see that the maximal separation linear tests can obtain is  $\frac{1}{2}$ . The situation is quite similar for bisimulation metrics [10], which attribute distance  $\frac{1}{2}$  to  $M$  and  $N$ .

It is easy, however, to find a family of contexts  $\{C_n\}_{n \in \mathbb{N}}$  such that  $C_n[M]$  evaluates to **true** with probability 1, and  $C_n[N]$  evaluates to **false** with probability  $1 - \frac{1}{2^n}$ : define  $C_n$  as a context that copies its argument  $n$  times, returning **false** if *at least one* of the  $n$  copies evaluates to **false**, and otherwise returns **true**. As a consequence, the context distance between  $M$  and  $N$  is 1. In fact, this reasoning can be extended to any pair of programs which are not equivalent but whose probability of convergence is 1: out of a context which separates them of  $\varepsilon > 0$ , with  $\varepsilon$  very small, we can construct a context that separates them of 1 performing some statistical reasoning. The situation is more complicated if we take the probability of convergence as an observable: we cannot *always* construct contexts that discriminate terms based on their probability of convergence, although something can be done if the terms' probabilities of convergence are different but close to 1. The context metric, in other words, risks to be not continuous and close to trivial if contexts are too powerful. What the example above shows, however, is something even worse: if contexts are allowed to copy their arguments, then any metric defined upon the usual presentation of probabilistic  $\lambda$ -calculus as an LMC (a fragment of which is depicted in Figure 1) is bound to be *unsound* w.r.t. the context metric.

Whether bisimulation metrics are sound, how close they are to the context distance, and whether they are useful in relieving the burden of evaluating it, are however open and interesting questions even in absence of the copying capability, i.e., when the underlying language is *affine*. This is the main reason why

we focus in this work on such a  $\lambda$ -calculus, whose expressive power is limited (although definitely non-trivial [20]) but which is anyway higher-order. We discover this way an elegant and deep theory in which trace and bisimulation metrics are indeed sound. At the end of this paper, some hints will be given about how the case of the untyped  $\lambda$ -calculus can be handled, a problem which we leave for future work.

Evaluating the context distance between affine terms is already an interesting and nontrivial problem. Consider, as an example, a sequence of terms  $\{M_n\}_{n \in \mathbb{N}}$  defined inductively as follows (where  $\Omega$  stands for a term with zero probability of converging):

$$M_0 = \langle \lambda x. \Omega, \lambda x. \Omega \rangle; \quad M_{n+1} = \langle \lambda x. M_n, \lambda x. \Omega \rangle.$$

$M_0$  is the pair whose components are both equal to  $\lambda x. \Omega$ , and  $M_{n+1}$  is defined as a pair whose first component is the function which returns  $M_n$  whatever its argument is, and the second component is again  $\lambda x. \Omega$ . We are now going to define another sequence of terms  $\{N_n\}_{n \in \mathbb{N}}$ , which can be seen as a noisy variation on  $\{M_n\}_{n \in \mathbb{N}}$ . More precisely,  $N_0$  is the same as  $M_0$ , and for each  $n \in \mathbb{N}$ ,  $N_{n+1}$  is constructed similarly to  $M_{n+1}$ , but adding some negligible noise in *both* components:

$$N_0 = \langle \lambda x. \Omega, \lambda x. \Omega \rangle; \\ N_{n+1} = \langle \lambda x. N_n \oplus \frac{1}{2^{n+1}} \Omega, \lambda x. \Omega \oplus \frac{1}{2^{n+1}} I \rangle.$$

( $I$  stands for the identity:  $\lambda x. x$ , while the term  $L \oplus^p K$  has the same behaviour as  $L$  with probability  $(1 - p)$ , and the same behaviour as  $K$  with probability  $p$ .) We would like to study how the distance between  $M_n$  and  $N_n$  evolves when  $n$  tends to infinity: do the little differences we apply at each step  $n$  accumulate, and how can we express this accumulation quantitatively?

Intuitively, it is easy for the environment to separate  $M_n$  and  $N_n$  of  $\frac{1}{2^n}$ : it is enough to consider a context  $\mathcal{C}$  which simply takes the second component of the pair, passes any argument to it, and evaluates it: the convergence probability of  $\mathcal{C}[M_n]$  is 0, while the convergence probability of  $\mathcal{C}[N_n]$  is  $\frac{1}{2^n}$ . But the environment can also decide to take the *first* component of the pair, in order to use the fact that  $M_{n-1}$  and  $N_{n-1}$  can be distinguished: more precisely, let us suppose that we have a context  $\mathcal{C}$  which separates  $M_{n-1}$  and  $N_{n-1}$ . Then we can construct a context  $\mathcal{D}$  which takes the first element of the pair, passes any argument to it, tries to evaluate it, and if it succeeds, gives the result as an argument to  $\mathcal{C}$ . We would like to express the supremum of the separation that such a context can obtain as a function of the distance between  $M_{n-1}$  and  $N_{n-1}$ . Unfortunately, this is not so simple: if  $\mathcal{C}$  is such that the convergence probability of  $\mathcal{C}[M_{n-1}]$  is  $\varepsilon$  and the convergence probability of  $\mathcal{C}[N_{n-1}]$  is  $\iota$ , we can see that the convergence probability of  $\mathcal{D}[M_n]$  is  $\varepsilon$ , whereas the convergence probability of  $\mathcal{D}[N_n]$  is  $(\iota \cdot (1 - \frac{1}{2^n}))$ . But it is not possible to express  $|\varepsilon - \iota \cdot (1 - \frac{1}{2^n})|$  as a function of  $|\varepsilon - \iota|$  and  $n$ : intuitively, the separation that the context  $\mathcal{D}$  can achieve depends not only on the separation that the context  $\mathcal{C}$  can achieve, but also on *how*  $\mathcal{C}$  achieves it. And moreover, the environment may of course

decide to use the two components of the pair, and to make them interact in an arbitrary way. Summing up, although the mechanism of construction of these terms seems to be locally easy to measure, it is complicated to have any idea about how the distance between them evolves when  $n$  tends to infinity.

### III. PRELIMINARIES

In this section, an affine and probabilistic  $\lambda$ -calculus, which is the object of study of this paper, will be introduced formally, together with a notion of context distance for it.

#### A. An Affine, Untyped, Probabilistic $\lambda$ -Calculus

We endow the  $\lambda$ -calculus with a probabilistic choice operator  $\oplus$ , which corresponds to the possibility for the program to choose one between two arguments, each with the same probability. *Terms* are expressions generated by the following grammar:

$$M ::= x \mid \lambda x.M \mid MM \mid M \oplus M \mid \Omega,$$

where  $\Omega$  models divergence<sup>1</sup>, and  $x$  ranges over a countable set  $\mathbf{X}$  of variables.

The class of affine terms, which model functions using their arguments at most once, can be isolated by way of a formal system, whose judgements are in the form  $\Gamma \vdash M$  (where  $\Gamma$  is any finite set of variables) and whose rules are the following (where  $\Gamma, \Delta$  stands for the union of two disjoint contexts):

$$\frac{}{\Gamma, x \vdash x} \quad \frac{\Gamma, x \vdash M}{\Gamma \vdash \lambda x.M} \quad \frac{\Gamma \vdash M \quad \Delta \vdash N}{\Gamma, \Delta \vdash MN} \\ \frac{\Gamma \vdash M \quad \Gamma \vdash N}{\Gamma \vdash M \oplus N} \quad \frac{}{\Gamma \vdash \Omega}$$

A *program* is a term such that  $\emptyset \vdash M$ , and  $\mathbf{P}$  is the set of all such terms. We will call them *closed terms*. We say that a program is a *value* if it is of the form  $\lambda x.M$ , and  $\mathbf{V}$  is the set of all such programs. The semantics of the just defined calculus is expressed as a binary relation  $\Downarrow$  between programs and *value subdistributions* (or simply *value distributions*), i.e. functions from values to real numbers whose sum is *smaller* or equal to 1. The relation  $\Downarrow$  is inductively defined by the following rules:

$$\frac{}{\Omega \Downarrow \emptyset} \quad \frac{V \in \mathbf{V}}{V \Downarrow \{V^1\}} \quad \frac{M \Downarrow \mathcal{D} \quad N \Downarrow \mathcal{E}}{M \oplus N \Downarrow \frac{1}{2}\mathcal{D} + \frac{1}{2}\mathcal{E}} \\ \frac{M \Downarrow \mathcal{D} \quad N \Downarrow \mathcal{E}}{\{L\{V/x\} \Downarrow \mathcal{F}_{L,V}\}_{\lambda x.L \in \mathcal{S}(\mathcal{D}), V \in \mathcal{S}(\mathcal{E})}} \\ MN \Downarrow \sum \mathcal{D}(\lambda x.L) \cdot \mathcal{E}(V) \cdot \mathcal{F}_{L,V}$$

where  $\mathcal{S}(\mathcal{D})$  stands for the support of the distribution  $\mathcal{D}$ . The divergent program  $\Omega$ , as expected, evaluates to the *empty* value distribution  $\emptyset$  which assigns 0 to any value. The expression  $\{V^1\}$  stands for the Dirac's value distribution on  $V$ ; more generally the expression  $\{V_1^{p_1}, \dots, V_n^{p_n}\}$  indicates the value distribution assigning probability  $p_i$  to each  $V_i$  (and 0 to any other value).

<sup>1</sup>Since we only consider affine terms, we cannot encode divergence by the usual constructions of  $\lambda$ -calculus.

For every program  $M$ , there exists precisely *one* value distribution  $\mathcal{D}$  such that  $M \Downarrow \mathcal{D}$ , that we note  $\llbracket M \rrbracket$ . This holds only because we restrict ourselves to affine terms. Moreover,  $\llbracket M \rrbracket$  is always a finite distribution. The rule for application expresses the fact that the semantics is call-by-value: the argument is evaluated before being passed to the function. There is no special reason why we adopt call-by-value here, and all we are going to say also holds for (weak) call-by-name evaluation.

The way  $\Downarrow$  is defined means that it is a *big-step* notion of semantics. In some circumstances, we need to have a more local view of how programs behave. We can define *small-step semantics*, again as a relation  $\Rightarrow$  between programs and value distributions (itself defined on top of a relation  $\rightarrow$  between programs and *program* distributions capturing a single evaluation step). The definition poses no significant problems [5]. Big-step and small-step semantics are equivalent: for every program  $M$ , there exists a unique distribution  $\mathcal{D}$  such that  $M \Rightarrow \mathcal{D}$ , and moreover  $\mathcal{D} = \llbracket M \rrbracket$ .

#### B. Context Distance

We now want to define a notion of observation for programs which somehow measures the convergence probability of a program. We will do that following the previous literature on this subject. For any distribution  $\mathcal{D}$  over a set  $A$ , its sum  $\sum_{a \in A} \mathcal{D}(a)$  is indicated as  $\sum \mathcal{D}$  and is said to be the *weight* of  $\mathcal{D}$ . The *convergence probability* of a term  $M$ , that we note  $\mathcal{P}^{\text{cv}}(M)$ , is simply  $\sum \llbracket M \rrbracket$ , i.e., the weight of its semantics. For instance, the convergence probability of  $\Omega$  is zero.

The environment, as usual, is modelled by the notion of a *context*, which is nothing more than a term with at most one occurrence of the *hole*  $[\cdot]$ . Contexts are generated by the following grammar:

$$\mathcal{C} ::= [\cdot] \mid x \mid \lambda x.C \mid \mathcal{C}M \mid M\mathcal{C} \mid \mathcal{C} \oplus \mathcal{C} \mid \Omega.$$

Affine contexts can be identified by a formal system akin to the one for terms, that we elide due to lack of space (see [5]). We note as  $\mathcal{C}[M]$  the program obtained by replacing  $[\cdot]$  by the closed term  $M$  in  $\mathcal{C}$ . The interaction of a program  $M$  with a context  $\mathcal{C}$  is the execution of the program  $\mathcal{C}[M]$ .

We now consider three different ways of comparing programs, based on their behaviour when interacting with the environment: a preorder  $\leq^{\text{ctx}}$ , an equivalence relation  $\equiv^{\text{ctx}}$ , and a map  $\delta^{\text{ctx}}$ :

*Definition 1 (Context Equivalence, Context Distance):* Let  $M$  and  $N$  be two programs. Then we write that  $M \leq^{\text{ctx}} N$  if and only if for every context  $\mathcal{C}$ , it holds that  $\mathcal{P}^{\text{cv}}(\mathcal{C}[M]) \leq \mathcal{P}^{\text{cv}}(\mathcal{C}[N])$ . If  $M \leq^{\text{ctx}} N$  and  $N \leq^{\text{ctx}} M$ , then we say that the two terms are *context equivalent*, and we write  $M \equiv^{\text{ctx}} N$ . With the same hypotheses, we say the *context distance* between  $M$  and  $N$  is the real number  $\delta^{\text{ctx}}(M, N)$  defined as  $\sup_{\mathcal{C}} |\mathcal{P}^{\text{cv}}(\mathcal{C}[M]) - \mathcal{P}^{\text{cv}}(\mathcal{C}[N])|$ .

Please observe that, following [8], we only compare programs and not arbitrary terms. This is anyway harmless in an affine setting.

*Example 1:* Let  $I$  be the identity  $\lambda x.x$ .  $I$  and  $\Omega$  are as far as two programs can be:  $\delta^{\text{ctx}}(I, \Omega) = 1$ . To prove that, finding a context which always converges for one of the terms, and always diverges for the other one, suffices. We can take  $\mathcal{C} = [\cdot]$ , and we have that  $\mathcal{P}^{\text{cv}}(\mathcal{C}[I]) = 1$  and  $\mathcal{P}^{\text{cv}}(\mathcal{C}[\Omega]) = 0$ . Of course,  $I$  and  $\Omega$  are not context equivalent. Throwing in probabilistic choice can complicate matters a bit. Consider the two terms  $I \oplus \Omega$  and  $I$ . One can easily prove that  $\delta^{\text{ctx}}(I \oplus \Omega, I) \geq \frac{1}{2}$ ; just consider  $\mathcal{C} = [\cdot]$ . However, showing that the above inequality is in fact an equality, requires showing that there cannot exist any context that separates more, which is possible, but definitely harder. This will be shown in the next section, using a trace-based characterisation of context distance.

### C. On Pseudometrics

Which properties does the context distance satisfy, and which structure it then gives to the set of programs? This section answers these questions, and prepares the ground for the sequel by fixing some terminology.

*Definition 2 (Pseudometrics):* Let  $S$  be a set. A *premetric* on  $S$  is any function  $\mu : S \rightarrow S$  such that  $0 \leq \mu(s, t) \leq 1$  and  $\mu(s, s) = 0$ . A *pseudometric* on  $S$  is any premetric such that for every  $s, t, u \in S$ , it holds that  $\mu(s, t) = \mu(t, s)$  and  $\mu(s, t) \leq \mu(s, u) + \mu(u, t)$ . The set of all pseudometrics on  $S$  is indicated with  $\Delta(S)$ .

Please observe that pseudometrics are not metrics in the usual sense, since  $\mu(s, t) = 0$  does not necessarily imply that  $s = t$ . If we have a pseudometric  $\mu$ , we can construct an equivalence relation by considering the *kernel* of  $\mu$ , that is the set of those pairs  $(s, t)$  such that  $\mu(s, t) = 0$ . It is easy to prove that the context distance is indeed a pseudometric, and that its kernel is context equivalence. We would now want to define a preorder  $\leq^{\text{metr}}$  on pseudometrics in such a way that if  $\mu \leq^{\text{metr}} \rho$ , then the kernel of  $\mu$  is included in the kernel of  $\rho$ . The natural choice, then, is to take the following definition, which is the reverse of the pointwise order on  $[0, 1]$ :

*Definition 3 (Pseudometric Ordering):* Let  $S$  be any set, and let  $\mu$  and  $\rho$  be two metrics in  $\Delta(S)$ . Then we stipulate that  $\mu \leq^{\text{metr}} \rho$  if and only if, for every  $s, t \in S$  we have that  $\rho(s, t) \leq \mu(s, t)$ .

As is well-known, the relation  $\leq^{\text{metr}}$  gives  $\Delta(S)$  the structure of a complete lattice.

But when, precisely, can a pseudometric on programs be considered a sound notion of distance? First of all, we would like it to put two programs at least as far as the difference between their convergence probabilities, since this is precisely our notion of observation:

*Definition 4 (Adequacy):* Let  $\mu$  be a pseudometric on the set of programs. Then  $\mu$  is an *adequate pseudometric* if for any programs  $M$  and  $N$ , we have that  $|\sum \llbracket M \rrbracket - \sum \llbracket N \rrbracket| \leq \mu(M, N)$ .

Secondly, we are interested in how programs behave when interacting with the environment. Especially, if we have two terms  $M$  and  $N$  at a given distance  $\varepsilon$ , and we put them in an environment  $\mathcal{C}$ , we would like a pseudometric  $\mu$  to give us

some information about the distance between  $\mathcal{C}[M]$  and  $\mathcal{C}[N]$ . This is the idea behind the following, standard, definition:

*Definition 5 (Non-Expansiveness):* Let  $\mu$  be a pseudometric on programs. We say  $\mu$  is *non-expansive* if for every pair of programs  $M$  and  $N$  and for every context  $\mathcal{C}$ , we have that  $\mu(\mathcal{C}[M], \mathcal{C}[N]) \leq \mu(M, N)$ .

Non-expansiveness is the natural generalisation of the notion of a *congruence*. By construction,  $\delta^{\text{ctx}}$  is a non-expansive pseudometric. We can also adapt the notion of soundness to pseudometric;  $\mu$  is said to be a *sound* pseudometric on programs if  $\mu \leq^{\text{metr}} \delta^{\text{ctx}}$ . Clearly, any adequate and non-expansive pseudometric is sound. In the rest of this paper, we will only deal with pseudometrics, but for the sake of simplicity we will refer to them simply as metrics.

## IV. THE TRACE DISTANCE

The first notion of metric we study is based on traces, i.e., linear tests. This is handier than the context distance, since contexts are replaced by objects with a simpler structure.

### A. Definition

A *trace*  $s$  is a sequence in the form  $@V_1 \cdots @V_n$ , where  $V_1, \dots, V_n$  are values, and we note  $\mathcal{T}r$  the set of traces. We define the probability that a program accepts a trace inductively on the length of the trace, as follows:

$$\begin{aligned} Pr(\lambda x.M, \epsilon) &= 1; \\ Pr(\lambda x.M, @V \cdot s) &= Pr(M\{V/x\}, s); \\ Pr(M, s) &= \sum_V \llbracket M \rrbracket(V) \cdot Pr(V, s) \quad \text{if } M \notin \mathbf{V}. \end{aligned}$$

Please observe that the probability that a term  $M$  accepts a trace  $s = @V_1 \cdots @V_n$  is the probability of convergence of  $MV_1 \cdots V_n$ . We are now going to define a metric, based on the probability that programs accept arbitrary traces:

*Definition 6:* Let  $M, N$  be two programs. Then we define the *trace distance* between them as  $\delta^{\text{tr}}(M, N) = \sup_s |Pr(M, s) - Pr(N, s)|$ . One can then define *trace equivalence* and the *trace preorder*, in the expected way [5].

Please observe that  $\delta^{\text{tr}}$  is a pseudometric on programs in the sense of Definition 2, and that it is an adequate one. The kernel of  $\delta^{\text{tr}}$  is nothing more than trace equivalence.

*Example 2:* The trace distance  $\delta^{\text{tr}}(I \oplus \Omega, I)$  between  $I \oplus \Omega$  and  $I$  is  $\frac{1}{2}$ . Showing that it is greater than  $\frac{1}{2}$  is easy: it is sufficient to consider the empty trace. The other inequality, requires evaluating, for any trace  $s$ , the probability of accepting it. This is however much easier than dealing with all contexts, because we can now control the structure of the overall program we obtain: for any trace  $s = @V_1 \cdots @V_n$ , we can see that:  $Pr(I \oplus \Omega, s) = \frac{1}{2} \cdot \sum \llbracket V_1 \cdots V_n \rrbracket$ , and  $Pr(I, s) = \sum \llbracket V_1 \cdots V_n \rrbracket$ . The difference (in absolute value) between  $Pr(I \oplus \Omega, s)$  and  $Pr(I, s)$ , then, cannot be greater than  $\frac{1}{2}$ .

The trace distance and the context distance indeed *coincide*, as well as the trace and context preorder, and the trace and

context equivalence. In the rest of this section, we will give the details of the proof for the pseudometric case, but the proof is similar for  $\equiv^{\text{ctx}}$  and  $\leq^{\text{ctx}}$ . It is easy to realise that the context distance is a lower bound on the trace distance, since any trace  $@V_1 \cdots @V_n$  can be seen as the context  $[\cdot]V_1 \cdots V_n$ . We thus focus on non-expansiveness.

### B. Non-Expansiveness

Are there contexts that can separate strictly more than traces? In order to show that it is not the case, it is enough to show that  $\delta^{\text{tr}}$  is non-expansive:

*Theorem 1:* Let  $M$  and  $N$  be two programs, and let  $\mathcal{C}$  be a context. Then  $\delta^{\text{tr}}(\mathcal{C}[M], \mathcal{C}[N]) \leq \delta^{\text{tr}}(M, N)$ .

Since  $\delta^{\text{tr}}$  is adequate, we can conclude that trace metric and context metric actually coincide:

*Theorem 2:*  $\delta^{\text{ctx}} = \delta^{\text{tr}}$ .

The rest of this section is devoted to an outline of the proof of Theorem 1 (see [5] for more details). The proof we give here is roughly inspired by the proof of congruence of trace equivalence for a non-deterministic  $\lambda$ -calculus [8]. The overall structure of the proof is the following: we first express the capacity of a program to do a trace by means of a labelled transition system (LTS in the following) whose states are distributions over programs. Then we consider *another* LTS, where the states are distributions over pairs of contexts and programs, that intuitively models the execution of  $\mathcal{C}[M]$ , but keeps the evolution of  $\mathcal{C}$  and  $M$  apart. The first LTS, called  $\mathcal{L}^{\text{tr}}$ , has distributions over programs as states, and traces as actions. We indicate with  $\dot{\rightarrow}$  the transition relation associated to  $\mathcal{L}^{\text{tr}}$ . This LTS is itself defined on top of a labelled transition relation, noted  $\overset{a}{\rightarrow}$ , where  $a \in \{\tau\} \cup \{@V \mid V \in \mathbf{V}\}$ . Intuitively, a  $\tau$ -step corresponds to an internal computation step for any term in the support of the distribution, while a  $@V$ -step corresponds to an interaction with the environment, which provides  $V$  as an argument. The relation  $\dot{\rightarrow}$  is defined as the accumulation of several steps of  $\overset{a}{\rightarrow}$ . Both the relations  $\dot{\rightarrow}$  and  $\overset{a}{\rightarrow}$  are given in Figure 2. We write  $\mathcal{D} \dot{+} \mathcal{E}$  for  $\mathcal{D} + \mathcal{E}$  when we want to insist on  $\mathcal{D}$  and  $\mathcal{E}$  to have disjoint supports.

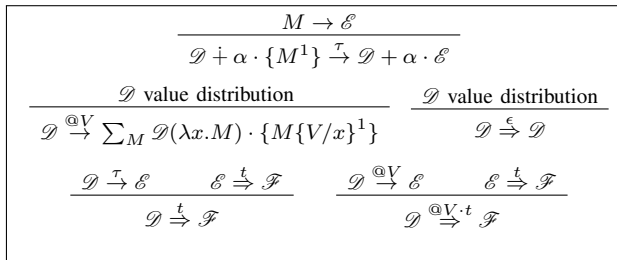


Fig. 2. Small-step Trace Relations on Program Distributions.

Please observe that these relations are *not* probabilistic. The relation  $\xrightarrow{\tau}$  is non-deterministic, since at each step we can decide which term of the distribution we want to reduce. However,  $\xrightarrow{\tau}$  is strongly normalising and confluent. Moreover, we can show that  $\dot{\rightarrow}$  is in fact deterministic: for every program

distribution  $\mathcal{D}$ , and every trace  $s$ , there exists a unique  $\mathcal{E}$  such that  $\mathcal{D} \xrightarrow{s} \mathcal{E}$ . The interest of the relation  $\dot{\rightarrow}$  is that it gives an alternate formulation for the probability of success in doing a trace:  $Pr(M, s) = p$  if and only if there is a distribution  $\mathcal{D}$  such that  $p = \sum \mathcal{D}$  and  $\{M^1\} \xrightarrow{s} \mathcal{D}$ . We can now use the relation  $\dot{\rightarrow}$  to give an equivalent formulation of Theorem 1: if  $M$  and  $N$  are such that  $\delta^{\text{tr}}(M, N) \leq \epsilon$ , then for every trace  $s$  and context  $\mathcal{C}$ , if  $\{\mathcal{C}[M]^1\} \xrightarrow{s} \mathcal{D}$  and  $\{\mathcal{C}[N]^1\} \xrightarrow{s} \mathcal{E}$ , then it holds that  $|\sum \mathcal{D} - \sum \mathcal{E}| \leq \epsilon$ . This statement, however, cannot be proved directly, yet, because the way  $\mathcal{C}$  and the argument terms interact is lost. It is then time to introduce our second LTS, called  $\mathcal{L}_{\mathbf{C} \times \mathbf{P}}^{\text{tr}}$ , which will allow us to relate  $\{\mathcal{C}[M]^1\} \dot{\rightarrow} \cdot$  to the behaviour of  $M$ : we want to talk about the evolution of a system consisting of the program  $M$  and the environment  $\mathcal{C}$ , while keeping the system and the environment as separate as possible.  $\mathbf{C} \times \mathbf{P}$  is the set of pairs of the form  $(\mathcal{C}, M)$ , where  $\mathcal{C}$  is a context and  $M$  is a program. We say that  $(\mathcal{C}, M) \in \mathbf{C} \times \mathbf{P}$  is a value if  $\mathcal{C}[M] \in \mathbf{V}$ . The states of  $\mathcal{L}_{\mathbf{C} \times \mathbf{P}}^{\text{tr}}$  are distributions over  $\mathbf{C} \times \mathbf{P}$ , while its labels are traces. We will note as  $\dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}}$  the transition relation of  $\mathcal{L}_{\mathbf{C} \times \mathbf{P}}^{\text{tr}}$ . Similarly to what we have done for  $\mathcal{L}^{\text{tr}}$ , we first define an auxiliary relation  $\overset{a}{\rightarrow}_{\mathbf{C} \times \mathbf{P}}$ . The rules defining  $\dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}}$  and  $\overset{a}{\rightarrow}_{\mathbf{C} \times \mathbf{P}}$  are in Figure 3. The relation  $\dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}}$  is strongly normalising. As a consequence,  $\dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}}$  is deterministic.

Intuitively, if  $s$  is a trace, the transition relation  $\dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}}$  corresponds to the transition relation  $\xrightarrow{s}$  starting from  $\{\mathcal{C}[M]^1\}$ . More precisely, if  $\{\mathcal{C}[M]^1\} \xrightarrow{s} \mathcal{D}$  and  $\{(\mathcal{C}, M)^1\} \dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}} \mathcal{E}$ , then  $\sum \mathcal{D} = \sum \mathcal{E}$ . This allows us to give yet another equivalent formulation of Theorem 1: if  $\delta^{\text{tr}}(M, N) \leq \epsilon$ , then if  $\{(\mathcal{C}, M)^1\} \dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}} \mathcal{D}$  and  $\{(\mathcal{C}, N)^1\} \dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}} \mathcal{E}$ , it holds that  $|\sum \mathcal{D} - \sum \mathcal{E}| \leq \epsilon$ . We are in fact going to show a stronger result, which uses the notion of  $\epsilon$ -related distributions:

*Definition 7:* Let  $\epsilon \in [0, 1]$ . We say that two distributions  $\mathcal{D}$  and  $\mathcal{E}$  over  $\mathbf{C} \times \mathbf{P}$  are  $\epsilon$ -related, and we note  $\mathcal{D} \equiv_{\epsilon}^{\text{par}} \mathcal{E}$  if there exist  $n \in \mathbb{N}$ ,  $\mathcal{C}_1, \dots, \mathcal{C}_n$  distinct contexts,  $p_1, \dots, p_n$  positive real numbers with  $\sum_i p_i \leq 1$ , and  $\mathcal{F}_1, \dots, \mathcal{F}_n$ , and  $\mathcal{G}_1, \dots, \mathcal{G}_n$  distributions over  $\mathbf{P}$ , such that:

- $\mathcal{D} = \sum_{1 \leq i \leq n} p_i \cdot (\mathcal{C}_i, \mathcal{F}_i)$ ;
- $\mathcal{E} = \sum_{1 \leq i \leq n} p_i \cdot (\mathcal{C}_i, \mathcal{G}_i)$ ;
- $\forall i, \delta^{\text{tr}}(\mathcal{F}_i, \mathcal{G}_i) \leq \epsilon$ .

Please observe that, if  $\delta^{\text{tr}}(M, N) \leq \epsilon$ , then for every context  $\mathcal{C}$ , the distributions  $\{(\mathcal{C}, M)^1\}$  and  $\{(\mathcal{C}, N)^1\}$  are  $\epsilon$ -related. In fact, the notion of  $\epsilon$ -relatedness is a way to capture pairs of distributions over  $\mathbf{C} \times \mathbf{P}$  representing *the same* environment, in which we put programs which are close for the trace pseudometric. The following can be seen as a stability result: if we start from  $\epsilon$ -related distributions, and we do a trace  $s$ , we end up in two distributions which are still  $\epsilon$ -related.

*Lemma 1:* Let  $\mathcal{D}, \mathcal{E}$  be distributions over  $\mathbf{C} \times \mathbf{P}$ , and  $\epsilon \in [0, 1]$  such that  $\mathcal{D} \equiv_{\epsilon}^{\text{par}} \mathcal{E}$ . Let  $s$  be a trace. Let  $\mathcal{F}$  and  $\mathcal{G}$  be such that:  $\mathcal{D} \dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}} \mathcal{F}$ , and  $\mathcal{E} \dot{\rightarrow}_{\mathbf{C} \times \mathbf{P}} \mathcal{G}$ . Then  $\mathcal{F} \equiv_{\epsilon}^{\text{par}} \mathcal{G}$ .

We can now see Theorem 1 as a direct consequence of Lemma 1. Indeed, let  $M$  and  $N$  be two programs at distance at most  $\epsilon$  for the trace metric, and let  $\mathcal{D}$  and  $\mathcal{E}$  be such that

$\frac{\{(C, M)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{E}}{\mathcal{D} \dot{+} p \cdot \{(CN, M)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{D} + p \cdot \sum_{\mathcal{D}, L} \mathcal{E}(\mathcal{D}, L) \cdot \{(\mathcal{D}N, L)^1\}}$	$\frac{C[M] \in \mathbf{V} \quad N \rightarrow \mathcal{E}}{\mathcal{D} \dot{+} p \cdot \{(CN, M)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{D} + p \cdot \sum_L \mathcal{E}(L) \cdot \{(CL, M)^1\}}$
$\frac{N \rightarrow \mathcal{E}}{\mathcal{D} \dot{+} p \cdot \{(NC, M)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{D} + p \cdot \sum_L \mathcal{E}(L) \cdot \{(LC, M)^1\}}$	$\frac{\{(C, M)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{E} \quad V \in \mathbf{V}}{\mathcal{D} \dot{+} p \cdot \{(VC, M)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{D} + p \cdot \sum_{\mathcal{D}, L} \mathcal{E}(\mathcal{D}, L) \cdot \{(V\mathcal{D}, L)^1\}}$
$\frac{M \rightarrow \mathcal{E}}{\mathcal{D} \dot{+} p \cdot \{([\cdot], M)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{D} + p \cdot \sum_L \mathcal{E}(L) \cdot \{([\cdot], L)^1\}}$	$\frac{}{\mathcal{D} \dot{+} p \cdot \{([\cdot]V, \lambda x.N)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{D} + p \cdot \{([\cdot], N\{V/x\})^1\}}$
$\frac{C[M] \in \mathbf{V}}{\mathcal{D} \dot{+} p \cdot \{((\lambda x.N)C, M)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{D} + p \cdot \{(N\{C/x\}, M)^1\}}$	$\frac{}{\mathcal{D} \dot{+} p \cdot \{(\Omega, M)^1\} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{D}}$
$p \cdot \{(\lambda x.C, M)^1\} \xrightarrow{\textcircled{V}}_{\mathbf{C} \times \mathbf{P}} p \cdot \{(C\{V/x\}, M)^1\}$	$p \cdot \{([\cdot], \lambda x.M)^1\} \xrightarrow{\textcircled{V}}_{\mathbf{C} \times \mathbf{P}} p \cdot \{([\cdot], M\{V/x\})^1\}$
$\frac{(\mathcal{D}_i \xrightarrow{\textcircled{V}}_{\mathbf{C} \times \mathbf{P}} \mathcal{E}_i)_{1 \leq i \leq n}}{\sum_{1 \leq i \leq n} \mathcal{D}_i \xrightarrow{\textcircled{V}}_{\mathbf{C} \times \mathbf{P}} \sum_{1 \leq i \leq n} \mathcal{E}_i}$	
$\frac{\mathcal{D} \text{ value distribution}}{\mathcal{D} \xrightarrow{\textcircled{V}}_{\mathbf{C} \times \mathbf{P}} \mathcal{D}}$	$\frac{\mathcal{D} \xrightarrow{\tau}_{\mathbf{C} \times \mathbf{P}} \mathcal{E} \quad \mathcal{E} \xrightarrow{t}_{\mathbf{C} \times \mathbf{P}} \mathcal{F}}{\mathcal{D} \xrightarrow{t}_{\mathbf{C} \times \mathbf{P}} \mathcal{F}}$
$\frac{\mathcal{D} \xrightarrow{\textcircled{V}}_{\mathbf{C} \times \mathbf{P}} \mathcal{E} \quad \mathcal{E} \xrightarrow{t}_{\mathbf{C} \times \mathbf{P}} \mathcal{F}}{\mathcal{D} \xrightarrow{\textcircled{V}, t}_{\mathbf{C} \times \mathbf{P}} \mathcal{F}}$	

Fig. 3. Small-Step Trace Relations on Distributions over  $\mathbf{C} \times \mathbf{P}$ .

$\{(C, M)^1\} \xrightarrow{\textcircled{V}}_{\mathbf{C} \times \mathbf{P}} \mathcal{D}$ , and  $\{(C, N)^1\} \xrightarrow{\textcircled{V}}_{\mathbf{C} \times \mathbf{P}} \mathcal{E}$ . Then, as we have already observed,  $\{(C, M)^1\}$  and  $\{(C, N)^1\}$  are  $\varepsilon$ -related. By Lemma 1, we can deduce that  $\mathcal{D}$  and  $\mathcal{E}$  are  $\varepsilon$ -related. And it is easy to see that this implies  $|\sum \mathcal{D} - \sum \mathcal{E}| \leq \varepsilon$ .

### C. Adding Pairs to the Calculus

The trace distance and the results we have just presented about it can be extended to an affine  $\lambda$ -calculus *with pairs*, namely a calculus whose language of terms also includes the following two constructs:

$$M ::= \langle M, N \rangle \mid \text{let } \langle x, y \rangle = M \text{ in } N.$$

We assume that terms are typed in any linear type system guaranteeing the absence of deadlocks (e.g., simple or recursive types), and we generalize the operational semantics in a natural way [5]. We would now like to extend the definition of a trace accordingly: which action should we perform on a term in the form  $\langle M, N \rangle$ ? The naïve solution would be to add projections to the trace language:  $s ::= \pi_1 \cdot s \mid \pi_2 \cdot s$ , with trace interpretation extended in the expected way. However, this way the trace distance would not be sound for the context distance, anymore. Indeed, let us consider the following example:

*Example 3:* We are going to compare the following terms:

$$M := \langle \lambda z.(I \oplus \Omega), \lambda z.(I \oplus \Omega) \rangle; \quad N := \langle \lambda z.I, \lambda z.I \rangle.$$

These two terms are at context distance at least  $\frac{3}{4}$ , since we can consider the context  $\mathcal{C} := \text{let } \langle x, y \rangle = [\cdot] \text{ in } (xI)(yI)$ , and we can see that  $\sum \llbracket \mathcal{C}[M] \rrbracket = \frac{1}{4}$ , while  $\sum \llbracket \mathcal{C}[N] \rrbracket = 1$ . But we cannot find any trace that separates them more than  $\frac{1}{2}$ . The interesting case is when  $s = \pi_i \cdot t$ . But then:

$$\begin{aligned} |Pr(M, s) - Pr(N, s)| &= |Pr(\lambda z.(\Omega \oplus I), t) - Pr(\lambda z.I, t)| \\ &\leq \delta^{\text{tr}}(\lambda z.(\Omega \oplus I), \lambda z.I). \end{aligned}$$

And it is easy to see that in the calculus with pairs we still have  $\delta^{\text{tr}}(\lambda z.(\Omega \oplus I), \lambda z.I) = \frac{1}{2}$ .

The reason why we cannot recover the context distance by way of projections is that the `let` construct above allows us to access *both* components of a pair, and the distances each of them induce can *add up*. A way out consists in extending the trace language to pairs really following linearity, and considering a new action in the form  $\otimes L$  with the following trace interpretation:

$$Pr(\langle M, N \rangle, \otimes L \cdot t) = \sum_{V, W} \llbracket M \rrbracket(V) \cdot \llbracket N \rrbracket(W) \cdot Pr(L\{V, W/x, y\}, t).$$

Remarkably, this is consistent with the well-known embedding of pairs into the applicative fragment (see, e.g., [2], and of course [5]).

This way of handling pairs allows the trace distance and the context distance to coincide, again. However, the trace distance loses its grip with respect to the context distance. Consider, for instance, the terms  $M$  and  $N$  from Example 3. Showing an upper bound on the distance between  $M$  and  $N$  is the same thing as showing an upper bound on  $\delta^{\text{tr}}(L\{\lambda z.(\Omega \oplus I), \lambda z.(\Omega \oplus I)/x, y\}, L\{\lambda z.I, \lambda z.I/x, y\})$  for all terms  $L$  such that  $x, y \vdash L$ , which is in fact not far away from what we should show if we were considering the context distance directly.

## V. THE BISIMULATION DISTANCE

As we realised in the last section, the trace metric can be a way to alleviate the burden of evaluating the context distance between terms but, in particular in presence of pairs, its usefulness is limited. In this section, we will look at another way to define the distance between programs which is genuinely coinductive, being based on the Kantorovich metric for distributions.

### A. Definition

A labelled Markov chain (LMC) is a triple  $\mathcal{M} = (S, \mathcal{L}, \mathcal{P})$ , where  $S$  is a countable set of states,  $\mathcal{L}$  is a



countable set of labels, and  $\mathcal{P}$  is a transition probability matrix, that is a function:  $\mathcal{P} : S \times \mathcal{L} \rightarrow \text{Distr}(S)$ . Moreover, if the image of  $\mathcal{P}$  only consists of distributions with *finite* support, we call  $\mathcal{M}$  an *image-finite* LMC. We are now going to define, similarly to [10] (but in absence of non-determinism), the metric analogue to bisimulation. The idea is to define a metric on the set  $S$  of states of the LMC as the greatest fixed point of some monotone operator on metrics. Please recall that  $(\Delta(S), \leq^{\text{metr}})$  is a complete lattice, and so any monotone operator has indeed a greatest fixed point.

*Lifting Metrics to Distributions:* We are going to define a way to turn any premetric over a set  $S$  into a metric over finite distribution over  $S$ .

*Definition 8:* Let  $\mu$  be a premetric on a set  $S$ . We define the *lifting of  $\mu$*  as the metric on the set of finite distributions over  $S$  defined as follows: for every  $\mathcal{D}, \mathcal{E}$  finite distributions over  $S$ ,  $\mu(\mathcal{D}, \mathcal{E})$  is the optimum solution to the following linear program:

$$\begin{aligned} \min \quad & \sum_{(s,t) \in \mathcal{S}(\mathcal{D}) \times \mathcal{S}(\mathcal{E})} h_{s,t} \cdot \mu(s,t) + \sum_s w_s + \sum_t z_t \\ \text{subject to} \quad & \sum_{s \in \mathcal{S}(\mathcal{D})} h_{s,t} + z_t = \mathcal{E}(t); \\ & \sum_{t \in \mathcal{S}(\mathcal{E})} h_{s,t} + w_s = \mathcal{D}(s); \\ & \forall (s,t) \in \mathcal{S}(\mathcal{D}) \times \mathcal{S}(\mathcal{E}), h_{s,t}, z_t, w_s \geq 0. \end{aligned}$$

Please observe that this linear program has an optimal solution. We can make use of the notion of duality from linear programming, and obtain an alternative characterisation of lifting:

*Theorem 3:* Let  $\mu$  be a premetric on  $S$  and Let  $\mathcal{D}, \mathcal{E}$  be finite distributions over  $S$ . Then:

$$\begin{aligned} \mu(\mathcal{D}, \mathcal{E}) &= \max \sum_s a_s \cdot \mathcal{D}(s) + b_s \cdot \mathcal{E}(s) \\ \text{subject to} \quad & \forall s \in S, a_s \leq 1; \\ & \forall s \in S, b_s \leq 1; \\ & \forall s, t \in S, a_s + b_t \leq \mu(s,t). \end{aligned}$$

The interest of the lifting construction comes from the fact that the lifting of a metric  $\mu$  behaves coherently with the original metric  $\mu$ . In particular, if we know the lifting of  $\mu$ , we are able to recover  $\mu$  by considering Dirac distributions:

*Lemma 2:* Let  $\mu$  be a premetric on  $S$ , and  $s, t \in S$ . Then  $\mu(\{s^1\}, \{t^1\}) = \mu(s, t)$ .

Moreover, the metric's structure is preserved when we consider the lifting: indeed, if a premetric on states is symmetric and verifies the triangular inequality, the same holds for its lifting.

*Metric as Fixpoints:* In a non-probabilistic setting, a relation  $R$  is a bisimulation if every pair of states  $s, t$  such that  $s R t$  can do the same actions and end up into states which are themselves in the relation. In order to obtain a quantitative counterpart of the scheme above, we define an operator  $F$  on the set of metrics over the states of a LMC: intuitively, given a metric  $\mu$ , we define a new metric  $F(\mu)$  which corresponds to the distance obtained by first doing a step of the transition

relation, and then applying the lifting of  $\mu$  to the resulting distributions.

*Definition 9:* Let  $\mathcal{M} = (S, \mathcal{L}, \mathcal{P})$  be an image-finite LMC. We define an operator  $F$  on  $\Delta(S)$  as

$$F(\mu)(s, t) = \sup\{\mu(\mathcal{P}(s, a), \mathcal{P}(t, a)) \mid a \in \mathcal{L}\}.$$

Please observe that  $F$  is a monotone operator on the complete lattice of metrics. We define the *bisimulation metric*  $\delta_{\mathcal{M}}^b$  on  $\mathcal{M}$  as the greatest fixpoint of  $F$ .

*Bisimulation Metric and the Affine  $\lambda$ -Calculus:* We are now going to consider a specific LMC  $\mathcal{M}^\Lambda$ , which captures the interactive behaviour of our calculus.

*Definition 10:* We define the LMC  $\mathcal{M}^\Lambda = (S^\Lambda, \mathcal{L}^\Lambda, \mathcal{P}^\Lambda)$  where:

- The set of states  $S^\Lambda$  is defined as follows:  $\mathbf{P} \uplus \mathbf{V}$ . A value  $V$  in the second component of  $S^\Lambda$  is distinguished from one in the first by using the notation  $\hat{V}$ .
- The set of labels  $\mathcal{L}^\Lambda$  is taken to be

$$\mathcal{L}^\Lambda = \{\text{@}V \mid V \in \mathbf{V}\} \cup \{\text{eval}\}.$$

- The transition probability matrix  $\mathcal{P}^\Lambda$  is such that: for every  $M \in \mathbf{P}$ , and any value  $V \in S(\llbracket M \rrbracket)$ , it holds that  $\mathcal{P}^\Lambda(M, \text{eval})(\hat{V}) = \llbracket M \rrbracket(V)$ , and that for every term  $M$  such that  $\lambda x.M \in \mathbf{P}$ , and  $V \in \mathbf{V}$ , it holds that  $\mathcal{P}^\Lambda(\widehat{\lambda x.M}, \text{@}V)(M\{V/x\}) = 1$ .

The results we have proved previously in this section apply to  $\mathcal{M}^\Lambda$ . In particular, one can define the bisimulation metric on  $\mathcal{M}^\Lambda$ . The *bisimulation distance* on programs, which we indicate as  $\delta^b$ , is defined to be the restriction of  $\delta_{\mathcal{M}^\Lambda}^b$  to programs.

We can see easily that  $\delta^b$  is an adequate metric. But there is more, since the bisimulation metric is well-known to be a lower bound on the trace distance: the bisimulation distance is a *sound* metric. In the next section, we anyway show non-expansiveness for it, which is stronger.

## B. Non-Expansiveness

Proving the non-expansiveness of  $\delta^b$  cannot be done directly, by a plain induction on contexts. Our strategy towards the result is the Howe's technique [16], a way of proving congruence of coinductively-defined equivalences which has been widely used for deterministic and non-deterministic languages, and that we here adapt to metrics.

The idea, then, is to start from  $\delta^b$ , construct another metric  $\delta^{\text{b}^H}$  on top of  $\delta^b$  (which turns out to be non-expansive by construction), and show that  $\delta^{\text{b}^H} = \delta^b$ . We first need to transform our metric  $\delta^b$  on programs into a metric on (potentially open) terms. Any metric  $\mu$  on programs can be extended into a metric on open terms, which by abuse of notation we continue to call  $\mu$  and which is defined as follows

$$\begin{aligned} \mu(M, N) &= \sup_{V_1, \dots, V_n \in \mathbf{V}} \mu(M\{V_1, \dots, V_n/x_1, \dots, x_n\}, \\ &\quad N\{V_1, \dots, V_n/x_1, \dots, x_n\}), \end{aligned}$$

where  $x_1, \dots, x_n$  are the variables occurring free in either  $M$  or  $N$ . The Howe's lifting of a metric  $\mu$  is defined based on a formal system for judgements in the form  $\Gamma \vdash \mu^H(M, N) \leq \varepsilon$ , where  $\varepsilon$  is a real number between 0 and 1. Its rules are given in Figure 4. Please observe that, potentially, there are several different  $\varepsilon$  such that  $\Gamma \vdash \mu^H(M, N) \leq \varepsilon$ . We are finally in a

$$\begin{array}{c}
\frac{\mu(x, M) \leq \varepsilon \quad x, \Gamma \vdash M}{x, \Gamma \vdash \mu^H(x, M) \leq \varepsilon} \\
\\
\frac{\Gamma \vdash \mu^H(M, K) \leq \varepsilon \quad \mu(K \oplus T, L) \leq \gamma}{\Gamma \vdash \mu^H(N, T) \leq \iota} \quad \frac{\mu(K \oplus T, L) \leq \gamma}{\Gamma \vdash L} \\
\\
\frac{\Gamma \vdash \mu^H(M \oplus N, L) \leq \frac{\varepsilon + \iota}{2} + \gamma}{\Gamma \vdash \mu^H(M, K) \leq \varepsilon \quad \mu(\lambda x. K, L) \leq \iota} \\
\\
\frac{x, \Gamma \vdash \mu^H(M, K) \leq \varepsilon \quad \mu(\lambda x. K, L) \leq \iota}{\Gamma \vdash \mu^H(\lambda x. M, L) \leq \varepsilon + \iota} \\
\\
\frac{\Gamma \vdash \mu^H(M, K) \leq \varepsilon \quad \mu(KT, L) \leq \gamma}{\Delta \vdash \mu^H(N, T) \leq \iota} \quad \frac{\mu(KT, L) \leq \gamma}{\Gamma, \Delta \vdash L} \\
\\
\frac{\Gamma \vdash \mu^H(M, K) \leq \varepsilon \quad \mu(KT, L) \leq \gamma}{\Gamma, \Delta \vdash \mu^H(MN, L) \leq \varepsilon + \iota + \gamma}
\end{array}$$

Fig. 4. Howe's Rules.

position to define the Howe's lifting of  $\mu$ :

*Definition 11:* Let  $\mu$  be a metric on terms. We define a premetric  $\mu^H$  on terms by:

$$\mu^H(M, N) = \inf \left( \{ \varepsilon \mid \exists \Gamma, \Gamma \vdash \mu^H(M, N) \leq \varepsilon \} \cup \{1\} \right).$$

We can see that  $\delta^{bH}$  is a premetric on open terms. Please observe that it is not necessarily a metric, since its construction entails neither symmetry nor the triangular inequality. The interest of this construction is that the metric  $\delta^{bH}$  is (more or less by construction) non-expansive:

*Lemma 3 (Non-expansiveness of  $\delta^{bH}$ ):* For every context  $\mathcal{C}$  and for every terms  $M, N$  it holds that  $\delta^{bH}(\mathcal{C}[M], \mathcal{C}[N]) \leq \delta^{bH}(M, N)$ .

The goal now is to show that  $\delta^{bH} \leq^{\text{metr}} \delta^b$ . Since  $\delta^b$  is the greatest fixed point of  $F$  for our LMC  $\mathcal{M}^\Lambda$ , we are going to show that  $\delta^{bH}$  can be extended into a metric on the states of  $\mathcal{M}^\Lambda$ , obtaining a fixed point for the operator  $F$ . First we extend  $\delta^{bH}$  to a premetric on  $S^\Lambda$ :

*Definition 12:* We define the extension of  $\delta^{bH}$  to  $S^\Lambda$  (that we note still  $\delta^{bH}$  by abuse of notation), by:

$$\begin{aligned}
\delta^{bH}(M, N) &:= \delta^{bH}(M, N); \\
\delta^{bH}(\widehat{V}, \widehat{W}) &:= \delta^{bH}(V, W); \\
\delta^{bH}(M, \widehat{W}) &:= 1.
\end{aligned}$$

Since  $\delta^{bH}$  isn't guaranteed to be a metric, we are forced to further refine it, by adding rules corresponding to symmetry and to the triangular inequality: we define  $\delta_{\Delta}^{bH}$  over  $S^\Lambda$  by the rules of Figure 5.

We can see easily that  $\delta^b \leq^{\text{metr}} \delta^{bH} \leq^{\text{metr}} \delta_{\Delta}^{bH}$  with respect to the preorder on terms. We want to show that  $\delta^{bH} = \delta^b$ . In

$$\begin{array}{c}
\frac{\delta^{bH}(s, t) \leq \varepsilon}{\vdash \delta_{\Delta}^{bH}(s, t) \leq \varepsilon} \\
\\
\frac{\vdash \delta_{\Delta}^{bH}(s, t) \leq \varepsilon \quad \vdash \delta_{\Delta}^{bH}(t, s) \leq \iota}{\vdash \delta_{\Delta}^{bH}(s, t) \leq \min(\varepsilon, \iota)} \\
\\
\frac{\vdash \delta_{\Delta}^{bH}(s, t) \leq \varepsilon \quad \vdash \delta_{\Delta}^{bH}(t, u) \leq \iota}{\vdash \delta_{\Delta}^{bH}(s, u) \leq \varepsilon + \iota}
\end{array}$$

Fig. 5. Symmetric and Transitive Closure of  $\delta^{bH}$ .

order to have that, we will show that  $\delta_{\Delta}^{bH} \leq^{\text{metr}} \delta^b$ . That is a direct consequence of the following theorem:

*Theorem 4:*  $\delta_{\Delta}^{bH}$  is a pre-fixpoint of  $F$ .

**Proof.** We need to show that  $\delta_{\Delta}^{bH} \leq^{\text{metr}} F(\delta_{\Delta}^{bH})$ . Please remember that the preorder on metrics corresponds to the reverse of the point-wise preorder for states. So if we read this inequality on metrics as an inequality on the states of  $\mathcal{M}^\Lambda$ , we see that it is equivalent to: for every  $s, t \in S^\Lambda$ ,  $F(\delta_{\Delta}^{bH})(s, t) \leq \delta_{\Delta}^{bH}(s, t)$ . If we unfold the definition of the operator  $F$  on metrics, we can see that it means that for every  $a \in \mathcal{L}^\Lambda$ ,  $\delta_{\Delta}^{bH}(\mathcal{P}^\Lambda(s, a), \mathcal{P}^\Lambda(t, a)) \leq \delta_{\Delta}^{bH}(s, t)$ . Please remember that there are two kinds of actions in our LMC: the action *eval* of evaluating a program to obtain a value distribution, and the action  $@V$ , which corresponds to passing the value  $V$  to a distinguished value. If we consider separately each of these actions, we see that the result we want to have is equivalent to:

- Let  $M, N$  be closed terms. Then  $\delta_{\Delta}^{bH}(\widehat{[M]}, \widehat{[N]}) \leq \delta_{\Delta}^{bH}(M, N)$ ;
- Let  $M, N$  be such that  $x \vdash M$  and  $x \vdash N$ , and let  $V$  be a value. Then it holds that:  $\delta_{\Delta}^{bH}(M\{V/x\}, N\{V/x\}) \leq \delta_{\Delta}^{bH}(\lambda x. M, \lambda x. N)$ .

The first point is the so-called Key Lemma, while the second one is a form of substitutivity. Both are non-trivial to prove (see [5] for more details). Remarkably, the Key Lemma can be proved exploiting the duality between the two presentations of Kantorovich lifting (see Theorem 3).  $\square$

Since  $\delta^{bH}$  is non-expansive by construction, we now have the result we were aiming for:

*Theorem 5:*  $\delta^b$  is non-expansive.

### C. On Full-Abstraction and Pairs

The bisimulation distance is a sound approximation of the context distance. But how about full-abstraction? Is there any hope to prove that the two coincide? The answer is negative: there are terms whose distance is *strictly* higher in the bisimulation metric than in the context (or trace) metric.

*Example 4:* Consider the following terms:  $M$  corresponds to the program that takes an argument, and then returns  $I$  with probability  $\frac{1}{2}$ , and diverges with probability  $\frac{1}{2}$ .  $N$  corresponds to the program which chooses first between the function which

return  $I$  whenever it is called, and the function which diverges whenever called. Formally:

$$M := \lambda x.(I \oplus \Omega); \quad N := (\lambda x.I) \oplus (\lambda x.\Omega).$$

These two terms are at distance 0 for the context distance: since the calculus is linear, the step where the choice is done is irrelevant. However,  $\delta^b(M, N) = \frac{1}{2}$ : the proof, whose details can be found in [5] use the characterisation of bisimulation distance by testing from [9], in which not only linear tests, but also more complicated tests (like *threshold tests*) are available.

But how about pairs? Indeed, for the sake of simplicity, we have presented the metatheory of the bisimulation metric for a purely applicative  $\lambda$ -calculus. Following the lines of our discussion in Section IV-C, however, the LMC  $\mathcal{M}^\Lambda$  can be extended into one handling pairs in a relatively simple way. The difficulties we encountered when trying to evaluate the (trace, or context) distance between pairs of terms unfortunately remain: it is not clear whether coinduction could provide any additional *advantage* over contextual distance.

## VI. THE TUPLE DISTANCE

The two metrics we have just defined have been shown to be non-expansive, even if the calculus is extended with pairs. In that case, however, they do not represent so much of an improvement with respect to the context distance. Please recall *where* the problem comes from: we would like to define actions starting from  $\langle M, N \rangle$ , and respecting the affine paradigm. We have seen that taking *projections* as actions lead to an unsound metric, and we have circumvented the problem by considering an action  $\otimes L$ , following [8]. Intuitively the action  $\otimes L$  corresponds to replacing the free variables of  $L$  (which are supposed to be included in  $\{x, y\}$ ) by the components of the pair: if for instance  $V$  and  $W$  are values, we have that  $\langle V, W \rangle \xrightarrow{\otimes L} \{L\{V, W/x, y\}^1\}$ . But what can any environment  $L$  do if we give it  $V$  and  $W$  as two values to interact with? Let us suppose that both  $V$  and  $W$  are functions, and remember that we are in an affine setting. The environment can (probabilistically) pass some arguments to  $V$ , and independently some other arguments to  $W$ , and then possibly pass to one of the two programs an argument that contains the other one. The idea behind the construction we present in this section, then, is to *keep* the information about the two components of the pairs in the states until they really interact with each other.

Our idea can be made concrete by introducing another LMC, whose states are not closed terms anymore, but *tuples* in the form  $[V_1, \dots, V_n]$ , where  $V_1, \dots, V_n$  are values. The possible actions the environment can perform on a tuple  $[V_1, \dots, V_n]$  correspond to the choice of an index  $i \in \{1, \dots, n\}$  and of an action to apply to the value  $V_i$ . If  $V_i$  is a pair, the only possible action is to split it into two components. We call this action *unfold* <sup>$i$</sup> . If  $V_i$  is a function, the environment can pass it an argument, which can possibly be constructed using other  $V_j$ 's. More precisely, the argument is built by way of an open term  $\mathcal{C}$ , and a typing context  $\Gamma$ ,

such that  $\Gamma \vdash \mathcal{C}$ , and  $\Gamma$  is a subset of  $\{x_j \mid j \neq i\}$ : the free variables of  $\Gamma$  represent the places where other values  $V_j$ , with  $j \neq i$ , are used. Moreover, we ask that for any values  $W_1, \dots, W_n$ , the term obtained in substituting  $x_j$  by  $W_j$  is a value: it means that  $\mathcal{C}$  is one of the  $x_j$ , or of the form  $\lambda y.\mathcal{D}$ . We call a pair  $(\Gamma, \mathcal{C})$  which verifies these conditions a  $(n, i)$ -*open value*. Formally, the LMC  $\mathcal{M}_{\text{mul}}^\Lambda = (S_{\text{mul}}^\Lambda, \mathcal{A}_{\text{mul}}^\Lambda, \mathcal{P}_{\text{mul}}^\Lambda)$  is defined in Figure 6.

### A. The Metric

We are going to define a metric on closed terms which corresponds to linear tests in  $\mathcal{M}_{\text{mul}}^\Lambda$ . First, we define *tuple traces* simply as words over  $\mathcal{A}_{\text{mul}}^\Lambda$ . The probability to succeed in doing a trace  $s$  starting from a tuple  $K \in S_{\text{mul}}^\Lambda$  can be naturally defined, and paves the way to defining a metric on tuples of values:

$$\begin{aligned} Pr^{\text{mul}}(K, \varepsilon) &= 1; \\ Pr^{\text{mul}}(K, a \cdot s) &= \sum_H \mathcal{P}(K, a)(H) \cdot Pr^{\text{mul}}(H, s); \\ \delta^{\text{mul}}(K, H) &= \sup_s |Pr^{\text{mul}}(K, s) - Pr^{\text{mul}}(H, s)|. \end{aligned}$$

What we need, however, is a metric on *programs*. Please remember that states of the LMC  $\mathcal{M}_{\text{mul}}^\Lambda$  are tuples of *values*. Any program  $M$ , however, can be viewed as the distribution of values obtained by evaluating it, i.e. its semantics  $\llbracket M \rrbracket$ :

$$\begin{aligned} \delta^{\text{mul}}(M, N) &= \sup_s \left| \sum [M](V) \cdot Pr^{\text{mul}}([V], s) \right. \\ &\quad \left. - \sum [N](W) \cdot Pr^{\text{mul}}([W], s) \right|. \end{aligned}$$

The just introduced metric should at least be put in relation to the context metric for it to be useful. We know from Section IV that the context metric coincides with the trace metric. The following theorem relates the trace metric  $\delta^{\text{tr}}$  and the metric  $\delta^{\text{mul}}$ :

*Theorem 6:* Let  $I$  be any finite set of variables, and  $\{V_x\}_{x \in I}$  and  $\{W_x\}_{x \in I}$  any two collections of values. For any open term  $\mathcal{C}$  such that  $I \vdash \mathcal{C}$ , it holds that:

$$\begin{aligned} \delta^{\text{tr}}(\mathcal{C}\{V_x/x\}_{(x \in I)}, \mathcal{C}\{W_x/x\}_{(x \in I)}) \\ \leq \delta^{\text{mul}}([V_x]_{(x \in I)}, [W_x]_{(x \in I)}). \end{aligned}$$

Theorem 6 can be read as a non-expansiveness result: if we have a system  $\mathcal{E}$ , playing the role of the environment, and which is prepared to interact with  $n$  components, and moreover we have two tuples  $K$  and  $H$  of length  $n$ , then the tuple distance between  $K$  and  $H$  gives us an upper bound on the trace distance between the system composed of  $\mathcal{E}$  interacting with  $K$ , and the system composed of  $\mathcal{E}$  interacting with  $H$ .

We can now see that  $\delta^{\text{mul}}$  coincides with the context metric: one inequality comes from Theorem 6, the other comes from the fact that any trace  $s$  over  $\mathcal{A}_{\text{mul}}^\Lambda$  and designed to start from a single value, can be simulated by a context.

*Theorem 7:* On programs,  $\delta^{\text{mul}} = \delta^{\text{ctx}}$ .

$$\begin{aligned}
S_{\text{mul}}^\Lambda &= \{[V_1, \dots, V_n] \mid V_1, \dots, V_n \text{ closed values}\}; \\
\mathcal{A}_{\text{mul}}^\Lambda &= \{\text{unfold}^i \mid i \in \mathbb{N}\} \cup \left\{ @(\Gamma, \mathcal{C})^i \mid i \in \mathbb{N}, (\Gamma, \mathcal{C}) \text{ a } (n, i)\text{-open-value} \right\}; \\
\mathcal{P}_{\text{mul}}^\Lambda([s_1, \dots, \langle N, L \rangle, \dots, s_n], \text{unfold}^i)([s_1, \dots, s_{i-1}, V, W, s_{i+1}, \dots, s_n]) &= \llbracket N \rrbracket(V) \cdot \llbracket L \rrbracket(W); \\
\mathcal{P}_{\text{mul}}^\Lambda([s_1, \dots, \lambda y. N, \dots, s_n], @(\Gamma, \mathcal{C})^i)([s_{h_1}, \dots, W, \dots, s_{h_m}]) &= \llbracket N \{ \mathcal{C} \{ s_{j_i} / x_{j_i} \}_{(1 \leq i \leq k)} / y \} \rrbracket(W); \\
&\text{with } \{1, \dots, n\} = \{i\} \cup \{j_1, \dots, j_k\} \cup \{h_1, \dots, h_m\}, n = 1 + k + m, \text{ and } \Gamma = x_{j_1}, \dots, x_{j_k}.
\end{aligned}$$

Fig. 6. The Tuple LMC.

## B. Examples

The tuple distance, that we have just proved to be fully-abstract, can be seen as yet another presentation of the context distance. But there is much more: it allows to evaluate the distance between concrete programs, even when the latter contain pairs, in a relatively easy way. In this section, we will give two examples.

1) *A Simple Example:* Consider the terms  $M$  and  $N$  defined in Example 3. We can prove that  $\delta^{\text{mul}}(M, N) = \frac{3}{4}$ . The proof of the fact that  $\delta^{\text{mul}}(M, N) \leq \frac{3}{4}$  can be found in [5]. Here we are only going to show that  $\delta^{\text{mul}}(M, N) \geq \frac{3}{4}$ . In order to show that, we are going to present a particular trace  $s$  such that  $|Pr^{\text{mul}}([M], s) - Pr^{\text{mul}}([N], s)| = \frac{3}{4}$ . More precisely, we take  $s = \text{unfold}^1 \cdot @(\emptyset, I)^1 \cdot @(\emptyset, I)^2$ : it corresponds to first separating the two components of the pair, and then passing  $I$  as an argument to the first and to the second component. The relevant fragment of  $\mathcal{M}_{\text{mul}}^\Lambda$  can be found in Figure 7. In particular, we can see that  $Pr([M], s) = 1$ , and  $Pr([N], s) = \frac{1}{4}$ .

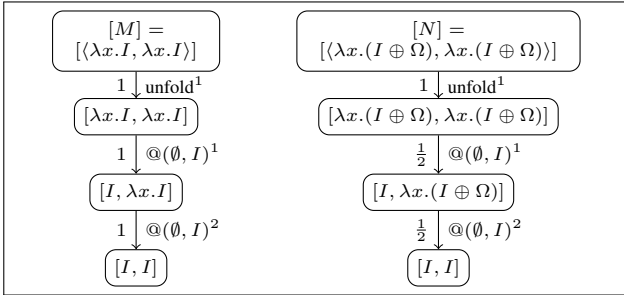


Fig. 7. A Fragment of the Tuple LMC.

2) *A More Complicated Example:* Please recall the example we presented in Section II. We note  $\{u_n\}_{n \in \mathbb{N}}$  the sequence defined as:  $u_n = \prod_{1 \leq i \leq n} (1 - \frac{1}{2^i})$ .

*Theorem 8:* For every  $n \in \mathbb{N}$ ,  $\delta^{\text{mul}}(M_n, N_n) = 1 - u_n$ .

**Proof.** We first show that  $\delta^{\text{mul}}(M_n, N_n) \geq 1 - u_n$ . As in the previous example, we do that by finding, for each  $n \in \mathbb{N}$ , a trace  $s_n$  such that  $|Pr([M_n], s_n) - Pr([N_n], s_n)| = 1 - u_n$ . We define the sequence  $(s_n)_{n \in \mathbb{N}}$  inductively as follows:

$$s_0 = \epsilon; \quad s_{n+1} = \text{unfold}^1 \cdot @(\emptyset, I)^1 \cdot s_n.$$

$s_0$  is the trace which always succeeds, whatever the starting state is.  $s_{n+1}$  corresponds to separating the two components of the pair which is in first position in the tuple, then passing the identity as an argument to the first component of this pair, and then executing  $s_n$ . For this sequence of traces, the recursive equations of Figure 8 are verified (the proof can be found in [5]). We can see by solving these equations that for every

$$\begin{aligned}
Pr([M_0], s_0) &= 1; & Pr([N_0], s_0) &= 1; \\
Pr([M_{n+1}], s_{n+1}) &= Pr([M_n], s_n); \\
Pr([N_{n+1}], s_{n+1}) &= (1 - \frac{1}{2^{n+1}}) \cdot Pr([N_n], s_n).
\end{aligned}$$

Fig. 8. Recursive Equations Verified by  $s_n$ .

$n \in \mathbb{N}$ ,  $Pr(M_n, s_n) = 1$  and  $Pr(N_n, s_n) = u_n$ . As a direct consequence, we obtain the result. We want now to show that  $\delta^{\text{mul}}(M_n, N_n) \leq 1 - u_n$ . To do that, we need to establish that there doesn't exist a trace  $t$  such that  $|Pr([M_n], t) - Pr([N_n], t)| > 1 - u_n$ . We're in fact going to show something stronger: for every  $n \in \mathbb{N}$ , we're going to define a set  $A_n$  of pairs of tuple, which contains the pair  $([M_n], [N_n])$ , and such that for every  $(K, H) \in A_n$ , for every trace  $t$ ,  $|Pr(K, t) - Pr(H, t)| \leq 1 - u_n$ . Intuitively, the idea behind the sequence  $\{A_n\}_{n \in \mathbb{N}}$  is the following: if we start from  $[M_n]$ , follow a trace of even length, and end up in a tuple  $K$  with a non-zero probability, and when we follow *the same trace* starting from  $[N_n]$  end up in the tuple  $H$ , then the pair of tuples  $(K, H)$  is in one of the  $A_j$ , with  $j$  smaller than  $n$ .

*Definition 13:* Let be  $n \in \mathbb{N}$ . Let  $A_n$  be the set of  $(K, H)$  such that: there exist  $m \in \mathbb{N}$ , and  $k_i \geq n + 1$  (for  $1 \leq i \leq m$ ), where:

$$\begin{aligned}
K &= [M_n, [\lambda x. \Omega]^m]; \\
H &= [N_n, [\lambda x. \Omega \oplus \frac{1}{2^{k_i}} I]_{1 \leq i \leq m}].
\end{aligned}$$

We want now to give an upper bound to the separation between  $K$  and  $H$  any trace can induce, if  $(K, H) \in A_n$ .

*Lemma 4:* For every  $n \in \mathbb{N}$ , for every  $(K, H) \in A_n$ , we can partition the set of traces as:

$$\mathcal{T}r = \{s \mid Pr(K, s) = 0 \text{ and } Pr(H, s) \leq \frac{1}{2}\} \\ \cup \{s \mid Pr(K, s) = 1 \text{ and } Pr(H, s) \geq u_n\}.$$

The proof of Lemma 4 can be found in [5]. The result we want to show is a direct consequence of Lemma 4: we can see easily that for any trace  $s$ , if  $(K, H) \in A_n$ , the separation  $s$  can induce is smaller or equal to  $1 - u_n$ .  $\square$

### C. On Tuples and Copying

The tuple distance naturally suggests a way to handle  $\lambda$ -calculi in which copying is indeed allowed. Although the details are clearly outside the scope of this paper, we anyway want to give some hints about why this is the case.

What makes the trace and behavioural distances unsound in presence of copying is their inability to capture an environment which can access the program at hand *more than once*. In our view, however, the problem does not come from the way those distances are defined in the abstract, but rather in the way *the underlying LMC* reflects the operational semantics of the calculus at hand. In a sense, it is in the responsibility of the LMC to guarantee that the environment can access terms multiple times. The LMC  $\mathcal{M}^\wedge$  we introduced in this paper (which is close to the ones from the literature [4], [6], [8]), as an example, is not adequate.

Suppose, however, to extend  $\mathcal{M}_{mul}^\wedge$  to an LMC for a  $\lambda$ -calculus in the style of Wadler’s linear  $\lambda$ -calculus [29]: there, the grammar of terms includes a construct  $!M$  whose purpose is marking those subterms which can indeed be duplicated. The actions the environment can perform on a term in the form  $!M$  simply reflects the above: the environment can create a *new copy* of  $!M$ , *but also keeps* the possibility to access  $!M$  in the future. One immediately realises that tuples are indeed the right way to model the access to both  $!M$  and  $M$ .

## VII. CONCLUSIONS

We have initiated the study of metrics in higher-order languages, starting with the relatively easy case of affine  $\lambda$ -terms, where copying capabilities are simply not available. We showed that three different notions of distance are sound (and sometime fully-abstract) for the context distance, the natural generalisation of Morris’ observational equivalence. One of them, the tuple distance, reflects the inherently monoidal structure of the underlying calculus, this way allowing to solve some nontrivial distance problems.

We are actively working on extending the results described here to the non-affine case, which for various reasons turns out to be more difficult, as discussed in Section II. We are in particular quite optimistic about the possibility of generalising the tuple distance to a metric reflecting copying. The real challenge, however, consists in handling the case in which copying is indeed available, but the number of copies of a given term the environment can have access to is somehow bounded, maybe polynomially in the value of a security

parameter. That would indeed be a way to get closer to computational indistinguishability, a central notion in modern cryptography.

## REFERENCES

- [1] S. Abramsky. The Lazy  $\lambda$ -Calculus. In D. Turner, editor, *Research Topics in Functional Programming*, pages 65–117. Addison Wesley, 1990.
- [2] A. Asperti and L. Roversi. Intuitionistic light affine logic. *ACM Trans. Comput. Log.*, 3(1):137–175, 2002.
- [3] A. Bizjak and L. Birkedal. Step-indexed logical relations for probability. In *FoSSaCS*, pages 279–294, 2015.
- [4] R. Crubillé and U. Dal Lago. On probabilistic applicative bisimulation and call-by-value  $\lambda$ -calculi. In *ESOP*, pages 209–228, 2014.
- [5] R. Crubillé and U. Dal Lago. Metric reasoning about  $\lambda$ -terms: the affine case (long version). Available at <http://arxiv.org/abs/1505.03638>, 2015.
- [6] U. Dal Lago, D. Sangiorgi, and M. Alberti. On coinductive equivalences for higher-order probabilistic functional programs. In *POPL*, pages 297–308, 2014.
- [7] V. Danos and T. Ehrhard. Probabilistic coherence spaces as a model of higher-order probabilistic computation. *Inf. Comput.*, 209(6):966–991, 2011.
- [8] Y. Deng and Y. Zhang. Program equivalence in linear contexts. To appear in *Theoretical Computer Science*, 2014.
- [9] J. Desharnais, V. Gupta, R. Jagadeesan, and P. Panangaden. Metrics for labeled markov systems. In *CONCUR*, 1999.
- [10] J. Desharnais, R. Jagadeesan, V. Gupta, and P. Panangaden. The metric analogue of weak bisimulation for probabilistic processes. In *LICS*, pages 413–422, 2002.
- [11] T. Ehrhard, C. Tasson, and M. Pagani. Probabilistic coherence spaces are fully abstract for probabilistic PCF. In *POPL*, pages 309–320, 2014.
- [12] D. Gebler and S. Tini. Fixed-point characterization of compositionality properties of probabilistic processes combinators. In *EXPRESS-SOS*, pages 63–78, 2014.
- [13] O. Goldreich. *Modern Cryptography, Probabilistic Proofs and Pseudorandomness*, volume 17 of *Algorithms and Combinatorics*. Springer, 1998.
- [14] S. Goldwasser and S. Micali. Probabilistic encryption. *J. Comput. Syst. Sci.*, 28(2):270–299, 1984.
- [15] N. D. Goodman, V. K. Mansinghka, D. M. Roy, K. Bonawitz, and J. B. Tenenbaum. Church: a language for generative models. In *UAI 2008*, pages 220–229, 2008.
- [16] D. J. Howe. Proving congruence of bisimulation in functional programming languages. *Inf. Comput.*, 124(2):103–112, 1996.
- [17] C. Jones and G. D. Plotkin. A probabilistic powerdomain of evaluations. In *LICS*, pages 186–195, 1989.
- [18] L. V. Kantorovich. On the transfer of masses. In *Dokl. Akad. Nauk. SSSR*, volume 37, pages 227–229, 1942.
- [19] S. B. Lassen. Relational reasoning about contexts. In *Higher Order Operational Techniques in Semantics, Publications of the Newton Institute*, pages 91–135. Cambridge University Press, 1998.
- [20] H. G. Mairson. Linear lambda calculus and ptime-completeness. *J. Funct. Program.*, 14(6):623–633, 2004.
- [21] C. D. Manning and H. Schütze. *Foundations of statistical natural language processing*, volume 999. MIT Press, 1999.
- [22] J. C. Mitchell, A. Ramanathan, A. Scedrov, and V. Teague. A probabilistic polynomial-time process calculus for the analysis of cryptographic protocols. *Theor. Comput. Sci.*, 353(1-3):118–164, 2006.
- [23] S. Park, F. Pfenning, and S. Thrun. A probabilistic language based on sampling functions. *ACM Trans. Program. Lang. Syst.*, 31(1), 2008.
- [24] J. Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, 1988.
- [25] A. M. Pitts. Operationally-based theories of program equivalence. In *Semantics and Logics of Computation*, pages 241–298. Cambridge University Press, 1997.
- [26] N. Ramsey and A. Pfeffer. Stochastic lambda calculus and monads of probability distributions. In *POPL*, pages 154–165, 2002.
- [27] S. Thrun. Robotic mapping: A survey. *Exploring artificial intelligence in the new millennium*, pages 1–35, 2002.
- [28] F. van Breugel and J. Worrell. A behavioural pseudometric for probabilistic transition systems. *Theor. Comput. Sci.*, 331(1):115–142, 2005.
- [29] P. Wadler. A syntax for linear logic. In *MFPS*, pages 513–529, 1993.