

# On approximation metrics for linear temporal model-checking of stochastic systems<sup>\*</sup>

Ilya Tkachev  
Delft University of Technology  
i.tkachev@tudelft.nl

Alessandro Abate  
University of Oxford  
aabate@cs.ox.ac.uk

## ABSTRACT

This paper proposes criteria for metrics between stochastic systems with a focus on the task of linear temporal model-checking. It explicitly puts forward two metrics which partially satisfy those criteria, and discusses their connection with other metrics studied in the literature. In particular, the notion of coupling between stochastic processes is shown to be crucial: omitting the explicit choice of coupling may lead to conservative results. The theoretical claims in the paper are supported by numerical examples.

## 1. INTRODUCTION

Stochastic systems have found broad applications in diverse areas where uncertainty can be quantified (cf. references in [31]). One particularly interesting class of problems concerns linear temporal (LT) model-checking of stochastic systems [8], which seeks to find the expectation of a path-dependent reward (or cost) functional, e.g. the probability that a realization of a system satisfies a given specification. Among the specifications of interest are reachability, safety, reach-avoid, and richer properties over a trajectory. If the system allows for a control input, one may further be interested in optimizing such an expectation or probability over all the admissible control policies [32]. Clearly, the simpler the system the easier the resulting model-checking procedure. In particular, if the state space of the process is finite, then model-checking can be performed algorithmically, by means of a specialized software [18, 21]. It is thus of interest to develop metrics<sup>1</sup> between stochastic systems

<sup>\*</sup>A. Abate is also affiliated with the Delft University of Technology. This work is supported by the European Commission STREP project MoVeS 257005, by the European Commission Marie Curie grant MANTRAS 249295, by the European Commission IAPP project AMBI 324432, by the European Commission NoE Hycon2 257462, and by the NWO VENI grant 016.103.020.

<sup>1</sup> In this paper we mostly deal with pseudo-metrics, which are allowed to be equal to zero for two different elements that are equivalent, in a certain sense. Since the difference between a metric and a pseudo-metric is not of importance to us, we refer to both objects simply as metrics.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

HSCC'14, April 14–17, 2014, Berlin, Germany.

Copyright 2014 ACM 978-1-4503-1567-8/13/04 ...\$15.00.

that allow quantifying an error over LT properties when substituting a complex, concrete system with a simpler, abstract one.

A recent survey on stochastic metrics can be found in [3]. Here we briefly recapitulate two main approaches to defining and using such metrics. The first one has been taken in the Computer Science community, whereas the second has been proposed in the Systems & Control area.

The work in [11] has been the first to emphasize the necessity for metrics between stochastic systems, justifying their relevance by the conservative nature and lack of robustness of the notion of exact probabilistic bisimulation relation [22]. According to the notion in [22], an infinite system is bisimilar to a finite one only if the former has a certain “lumpable” structure relating the two models, and any small perturbation of the models can violate such a structure, thus resulting in systems that originally were precisely related and that become now completely unrelated. In contrast, the metric introduced in [11] both admits the bisimulation as its zero level set, and varies continuously with respect to perturbations of the system dynamics. A few other metrics have been developed along the same lines (see e.g. [35] and references therein), however in all these cases the applicability to LT model-checking has not been discussed.

Similar ideas have been later applied to non-stochastic systems [16]. The approach there is different from the one in [11, 35], as it starts from trace-like metrics over the global dynamics (paths), and then relates the former metrics on path spaces to a bisimulation-like metric over the local dynamics (transitions), the latter being more conservative but easier to compute. The developed metrics were designed to quantify errors over LT properties. The extension of this technique to stochastic systems has been proposed in [19], however the obtained results appear to be rather conservative, possibly due to the fact that the work with stochastic systems has been done at the level of random elements (rather than distributions), and as their coupling procedure is not addressed explicitly (cf. Sections 2.1 and 3.2).

In summary, there are two main approaches to metrics between stochastic systems that are currently available. The first, taken in the Computer Science community [35, 11], does not provide an explicit way of using metrics for LT model-checking. The second, originated in the area of Systems & Control [19], although applicable to the aforementioned task, leads to conservative results. These reasons suggest to look into alternative approaches to define metrics between stochastic systems. We propose the following criteria for such metrics:

- i. a metric between two systems quantifies the difference between LT specifications computed over such systems;
- ii. given any two systems the metric between them is com-

putable, either analytically or at least by means of Monte Carlo simulations;

- iii. there is a class of “nice” infinite systems that allow for finite abstractions that are arbitrarily close in a given metric.

With focus on the criteria above, in this work we provide two new metrics. Unlike [19], we do not assume any form of stability over the systems under study, hence most of the results in this work hold true over a finite time horizon  $n$ , the dependence on which we then explicitly mention in the notation for the metrics. The first metric  $d_{TV^n}$  is based on the total variation distance between measures, whereas the second metric  $d_{W^n}$  is inspired by the Wasserstein distance.<sup>2</sup> Similar to [16, 19], we start with trace-like metrics to assure that the first criterion is satisfied. We also provide Monte Carlo methods to quantify both metrics in order to meet the requirements of the second criterion. Finally, we provide precise (rather than simulation-based) bounds on  $d_{TV^n}$  in terms of a bisimulation-like analogue, here called  $d_B$ : this allows us to characterize a class of systems that satisfy the last criterion with respect to  $d_{TV^n}$ , which in turn aligns with results on formula-free abstractions proposed in [31]. Unfortunately, such analytic results are much harder to obtain for the metric  $d_{W^n}$ .

The rest of the work is structured as follows. The majority of our ideas and methods rely upon the notion of random elements, their distributions, and their couplings. We briefly go over this theory in Section 2, which also introduces labelled discrete-time Markov processes (ldt-MP), the class of systems we focus on in this paper. The metrics  $d_{TV^n}$  and  $d_{W^n}$ , and relations between them, are discussed in Section 3. Theoretical and computational examples are further provided in Section 4, and the paper is concluded in Section 5. The notation can be found in Section 7.

## 2. PRELIMINARIES

### 2.1 Coupling

The notion of coupling [23] is crucial when considering an interplay between two stochastic processes, particularly in case one wants to quantify the difference between them. To introduce this notion, we need to elaborate on the definition of stochastic processes and on their representations. Recall that a stochastic process is a special case of a random element – this concept is often used to model probabilistic phenomena.

**DEFINITION 1.** *A random element on a measurable space  $(E, \mathcal{E})$  is a tuple  $(\Omega, \mathcal{F}, P, f)$ , where  $(\Omega, \mathcal{F}, P)$  is a probability space and  $f : \Omega \rightarrow E$  is a measurable map. We say that  $(\Omega, \mathcal{F}, P)$  is the sample space of the random element, and that  $(E, \mathcal{E})$  is its range. The distribution of the random element  $(\Omega, \mathcal{F}, P, f)$  is the probability measure  $f_*P \in \mathcal{P}(E, \mathcal{E})$  on its range.*

We will shortly discuss that distributions can be induced by different random elements, thus one can think of the latter as particular *representations* of distributions. More precisely, we say that a random element  $(\Omega, \mathcal{F}, P, f)$  is the representation of a probability measure  $Q$  whenever it holds that  $f_*P = Q$ . Note that, for any distribution  $Q \in \mathcal{P}(E, \mathcal{E})$ , there always exists at least one representation over its range, given by  $(E, \mathcal{E}, Q, \text{id}_E)$

<sup>2</sup> The Wasserstein metric is also known as Kantorovich or Hutchinson metric [35],[14, footnote in Section 2.2].

since  $(\text{id}_E)_*Q = Q$ . This is also known as the *canonical representation* of the distribution  $Q$ <sup>3</sup>. Let us provide some examples to further clarify the concept of representation. The first example shows that for any distribution a representation is never unique.

**EXAMPLE 1.** *Consider a distribution  $Q$  and let  $(\Omega, \mathcal{F}, P, f)$  be its arbitrary representation. For any probability space  $(\Omega', \mathcal{F}', P')$  it holds that  $(\Omega \times \Omega', \mathcal{F} \otimes \mathcal{F}', P \otimes P', f \circ \pi)$  is another representation of  $Q$ , where  $\pi : \Omega \times \Omega' \rightarrow \Omega$  denotes the projection map.*

In the previous example it is crucial that the sample space is modified. The next example emphasizes this fact by showing that if the sample space is fixed, then there may exist distributions that admit a unique representation over this sample space.

**EXAMPLE 2.** *Let  $E := \{-1, 1\}$ ,  $\mathcal{E} := 2^E$ , and let  $Q(\{1\}) = \frac{1}{3}$ . Note that there exist four distinct maps  $f : E \rightarrow E$ , that is  $\text{id}_E, -\text{id}_E$ , and the constant maps  $-1$  and  $1$ . Clearly,  $Q = (\text{id}_E)_*Q$  but  $Q \neq f_*Q$  if  $f \neq \text{id}_E$ , so that  $Q$  can be represented in the unique (canonical) way if the sample space is  $(E, \mathcal{E})$ . On the other hand, a symmetrical distribution  $Q'$  given by  $Q'(\{1\}) = \frac{1}{2}$  can be represented in two ways since  $Q' = (-\text{id}_E)_*Q'$ . More generally, if a distribution admits distinct representations over its range, they can be considered as symmetries of this distribution.*

According to [24, Section 2.1], a stochastic process on a measurable space  $(X, \mathcal{X})$  is a parameterized collection of  $X$ -valued random elements  $(x_t)_{t \in T}$ , all defined over the same probability space  $(\Omega, \mathcal{F}, P)$ . Clearly, one can equivalently consider a stochastic process as a single random element  $(\Omega, \mathcal{F}, P, f)$  where the map  $f : \Omega \rightarrow X^T$  is uniquely determined by  $f \circ \pi_t = x_t$  for any  $t$ ,  $\pi_t : X^T \rightarrow X$  being obvious projection maps. Recall that here  $X^T$  is the product measurable space consisting of all maps from  $T$  to  $X$ . As a consequence, all the results on random elements above now apply to stochastic processes. In particular, it is important for us that any stochastic process can be understood also in a “weak” form – as a distribution over  $X^T$ , rather than only in a perhaps more intuitive “strong” form – as a random element inducing such a distribution. This approach is important in the context of this work, since only distributions matter for LT model checking, rather than particular representations of a stochastic process. A representation, especially if it is not canonical, is sometimes provided to show a constructive definition of a distribution, as in the following example.

**EXAMPLE 3.** *A one-dimensional Itô diffusion [24, Chapter 7]*

$$dx_t = a(x_t)dt + b(x_t)dB_t, \quad x_0 = x \in \mathbb{R}, \quad (1)$$

*is a stochastic process with trajectories in  $E = \mathcal{C}([0, \infty))$ . Since the diffusion (1) is Markovian, it can be considered from two perspectives: either as a solution of the stochastic differential equation (SDE) (1), or as a Markov process on  $X = \mathbb{R}$ .*

*In order to treat a diffusion as a solution of an SDE, we consider a probability space  $(\Omega, \mathcal{B}(\Omega), P)$ , where  $\Omega = \mathcal{C}_0([0, \infty))$  is a set of all  $\omega \in E$  satisfying  $\omega(0) = 0$  and where  $P$  is the Wiener measure (the distribution of the Brownian motion  $(B_t)_{t \in [0, \infty)}$ ). In this setting a diffusion is constructed as a random element  $(\Omega, \mathcal{F}, P, f)$ , where  $f$  is the strong solution of an SDE [24, Section 5.3] with the initial condition  $x$ . Let us denote the induced distribution by  $Q := f_*P$ .*

<sup>3</sup> The existence of a canonical representation in particular means that in order to represent any given distribution its range suffices and there is no necessity in coming up with a “bigger” sample space  $(\Omega, \mathcal{F})$ .

Alternatively, one can treat the diffusion in (1) as a Markov process, and construct a distribution  $Q$  directly on  $E$  using the transition function associated with the diffusion [13, Section 4.1]. The stochastic process for the diffusion (in a strong sense) would be a canonical representation of  $Q$ . In this Markovian setting there is no need to define an auxiliary probability space  $(\Omega, \mathcal{F}, P)$  for the Brownian motion, so the current Markovian approach is more direct. On the other hand, such a construction is perhaps less intuitive as it does not emphasize the connection between the diffusion and the Brownian motion it is driven by.

Note that both methods lead to the same stochastic process in a weak sense (at the level of distributions), but to different stochastic processes in a strong sense (at the level of random elements). At the same time, if one needs to compute the probability that a trajectory of (1) reaches a target set, it does not matter which representation of a diffusion is used, since such a probability is uniquely determined by the distribution of the diffusion.

The motivation behind approximate abstractions of stochastic systems is to use the information obtained over a simpler abstract system to deduce properties of a more complicated concrete system. Assume that the range  $(E, \mathcal{E})$  is endowed with a metric  $d_E$ . Given a concrete random element  $(\Omega, \mathcal{F}, P, f)$  and its approximation  $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{P}, \tilde{f})$  one wants to obtain an inequality

$$\mathbb{P}(d_E(f, \tilde{f}) > \delta) \leq \varepsilon \quad (2)$$

to be able to interpret over  $f$  the results obtained for  $\tilde{f}$ . As an example [6], if we know the probability  $p = \tilde{P}(\tilde{f} \in A)$  for some set  $A \in \tilde{\mathcal{F}}$ , we can use (2) to provide bounds on  $P(f \in A_\delta)$ , where  $A_\delta = \{x \in E : d_E(x, A) \leq \delta\}$  is the  $\delta$ -inflation of the set  $A$ :

$$\begin{aligned} P(f \in A_\delta) &= \mathbb{P}(f \in A_\delta) \geq \mathbb{P}(\{\tilde{f} \in A\} \cap \{d_E(f, \tilde{f}) \leq \delta\}) \\ &\geq \mathbb{P}(\tilde{f} \in A) - \mathbb{P}(d_E(f, \tilde{f}) > \delta) \geq p - \varepsilon. \end{aligned}$$

Whenever  $f$  and  $\tilde{f}$  are stochastic processes and  $A$  is a reachability specification, the result allows using the probabilistic reachability analysis over  $\tilde{f}$  to study that over  $f$ . Recall that the abstraction is a random element on its own, and that the analysis over the abstraction can be carried out regardless of its relation to the concrete system: such a relation only matters when extrapolating results of this analysis from the abstraction back to the concrete system. For this purpose one has to define a common sample space for both systems in order to specify the probability measure in (2). The procedure of building a common sample space for *a priori* unrelated distributions or random elements is called coupling.

**DEFINITION 2.** A coupling of  $Q, \tilde{Q} \in \mathcal{P}(E, \mathcal{E})$  is a probability measure  $\mathbb{Q} \in \mathcal{P}(E^2, \mathcal{E}^2)$  satisfying the following equalities:

$$\pi_* \mathbb{Q} = Q, \quad \tilde{\pi}_* \mathbb{Q} = \tilde{Q}, \quad (3)$$

where  $\pi(x, \tilde{x}) = x$  and  $\tilde{\pi}(x, \tilde{x}) = \tilde{x}$  for all  $(x, \tilde{x}) \in E^2$  are obvious projection maps. A coupling of two random elements  $(\Omega, \mathcal{F}, P, f)$  and  $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{P}, \tilde{f})$  is any random element  $(\Xi, \mathcal{G}, \mathbb{P}, (g, \tilde{g}))$  such that  $(g, \tilde{g})_* \mathbb{P}$  is a coupling of  $f_* P$  and  $\tilde{f}_* \tilde{P}$ , and  $(\Xi, \mathcal{G})$  is an arbitrary measurable space.

Whenever (3) holds true, we say that  $Q$  and  $\tilde{Q}$  are marginals of  $\mathbb{Q}$ . It shall be clear that given any two marginal distributions, there always exists at least one coupling of them, called the *independence coupling*, which is given by  $\mathbb{Q} := Q \otimes \tilde{Q}$ . The most important point about the coupling is that it is only unique when

one of the marginals  $Q$  or  $\tilde{Q}$  is the Dirac measure, or equivalently when one of the random elements is deterministic. In particular, since the concrete system and the abstraction are almost never coupled *a priori* (neither at the level of random elements, nor at the level of distributions), one can e.g. optimize over all admissible couplings to choose the best for inequality (2). This idea constitutes to the core of our method. Before going into the details of it, let us provide examples of couplings.

**EXAMPLE 4.** Consider the diffusion  $x$  as per (1), and let

$$d\tilde{x}_t = \tilde{a}(\tilde{x}_t)dt + \tilde{b}(\tilde{x}_t)d\tilde{B}_t, \quad \tilde{x}_0 = \tilde{x} \in \mathbb{R}, \quad (4)$$

be another diffusion. Let us represent both of them as solutions of SDEs, that is as random elements on the sample spaces  $(\Omega, \mathcal{F}, P)$  and  $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{P})$ , where the former is defined in Example 3, and the latter is defined for (4) analogously. We provide three versions of couplings obtained via the Brownian motions driving these diffusions:

- $\Xi = \Omega$ ,  $\mathcal{G} = \mathcal{F}$ ,  $\mathbb{P} = P$  and  $(g, \tilde{g}) = (f, \tilde{f})$ . This coupling corresponds to the case  $B_t = \tilde{B}_t$  when diffusions are “sharing” the same noise [19].
- $\Xi = \Omega \times \tilde{\Omega}$ ,  $\mathcal{G} = \mathcal{F} \otimes \tilde{\mathcal{F}}$ ,  $\mathbb{P} = P \otimes \tilde{P}$  and  $(g, \tilde{g}) = (f \circ \pi, \tilde{f} \circ \tilde{\pi})$  where the projection maps  $\pi : \Xi \rightarrow \Omega$  and  $\tilde{\pi} : \Xi \rightarrow \tilde{\Omega}$  are as in Definition 2. This clearly corresponds to the case of the independence coupling, that is  $B_t \perp \tilde{B}_t$  [2].
- $\Xi = \Omega$ ,  $\mathcal{G} = \mathcal{F}$ ,  $\mathbb{P} = P$  and  $(g, \tilde{g}) = (f, \tilde{f} \circ n)$  where  $n(\omega) = -\omega$  for any  $\omega \in \Omega$ . In this case  $B_t = -\tilde{B}_t$ : noises driving diffusions are “reflected”. Such a construction is possible thanks to the fact that  $n$  is a symmetry of  $P$  (i.e.  $n_* P = P$ ), or in other words  $-\tilde{B}_t$  is a Brownian motion whenever  $\tilde{B}_t$  is.

Finally, let us mention that as much as any distribution of a single random element admits a representation over its range, the distribution of any coupling of two random elements admits a representation over the product of the ranges: this follows directly from Definition 2. In particular, if  $(\Omega, \mathcal{F}, P, f)$  and  $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{P}, \tilde{f})$  are any random elements with a range  $(E, \mathcal{E})$  endowed with a metric  $d$ , and  $(\Xi, \mathcal{G}, \mathbb{P}, (g, \tilde{g}))$  is one of their couplings with a distribution  $\mathbb{Q} = (g, \tilde{g})_* \mathbb{P}$ , then the probability in (2) can be expressed as

$$\mathbb{P}(d(g, \tilde{g}) > \delta) = \mathbb{Q}(E^2 \setminus \Delta_E^\delta),$$

where  $\Delta_E^\delta := \{(x, \tilde{x}) : d(x, \tilde{x}) \leq \delta\}$  is a diagonal  $\delta$ -strip in  $E^2$ . The benefit of dealing with stochastic systems (and their couplings) at the level of distributions lies the easy introduction of the distance between the concrete and the abstract systems. To support this point, let us mention that in probability theory distances are classically introduced between distributions, rather than between random elements [15, 25].

## 2.2 Behaviors of ldt-MP

Although most of the concepts and methods we have introduced apply to arbitrary stochastic processes, a focus on discrete-time Markov processes (dt-MP) allows us to provide a more detailed analysis. A dt-MP is a pair  $(X, P)$ , where  $X$  is a Borel space referred to as the state space, and  $P : X \rightarrow \mathcal{P}(X)$  is a stochastic kernel on  $X$ . It follows from [26, Theorem 2.8] that given any initial state  $x \in X$  there exists a unique probability measure  $P_x$  on the state-path space  $X^\omega$  satisfying

$$P_x \left( \prod_{k=0}^{\infty} dx_k \right) = 1_{\{x\}}(x) \cdot \prod_{k=0}^{\infty} P(x_k, dx_k). \quad (5)$$

We say that  $(X, P)$  is finite whenever  $X$  is a finite set.

REMARK 1. *Similar to Example 3, for a given dt-MP and a fixed initial state, (5) provides a weak stochastic process. The corresponding canonical strong stochastic process can be introduced by defining the coordinate maps  $\mathbf{x}_k : X^\omega \rightarrow X$  as obvious projections. As an alternative, it follows from [20, Proposition 7.6] that any dt-MP can be expressed as a stochastic difference equation  $\mathbf{x}_{k+1} = F(\mathbf{x}_k, w_k)$ , where  $w_k$  is a sequence of iid uniform random variables on  $[0, 1]$ . With focus on Section 2.1, this means that any dt-MP with a fixed initial state  $x \in X$  admits the following non-canonical representation:  $([0, 1]^\omega, \mathcal{B}([0, 1]^\omega), \lambda^\omega, f)$ , where  $\lambda$  is the Lebesgue measure on  $[0, 1]$  and the map  $f$  is obtained by iterating  $F$  starting from  $x$ .*

For a non-stochastic system on a state space  $X$ , its internal behavior is any element of  $X^\omega$  where each transition  $\mathbf{x}_k \rightarrow \mathbf{x}_{k+1}$  is allowed in the system [29, Section 1.2]. For stochastic systems such a definition is not suitable since formally all transitions are allowed, so that any single element of  $X^\omega$  is a possible internal behavior of a dt-MP, albeit possibly of zero probability. Due to this reason, it is more appropriate to talk about *behavioral distributions*, that is  $P_x$  shall be understood as a distribution of internal behaviors that indicates which behaviors are more likely appear as trajectories of a dt-MP.

Often one is not interested in each single behavior, but rather in collections thereof that can be characterised by means of observations (or labels). A labeled dt-MP (ldt-MP for short) is a tuple  $(X, P, Y, L)$  where  $(X, P)$  is a dt-MP,  $Y$  is a Borel space and  $L : X \rightarrow Y$  is a measurable map [31]. In an ldt-MP, to any internal behavior  $(x_0, x_1, \dots)$  there corresponds an *external behavior*, or a *trace*, given by an output of the trace map:

$$L_\omega(x_0, x_1, \dots) := (L(x_0), L(x_1), \dots).$$

For non-stochastic system a collection of its all admissible external behaviors is also called a *language* of a system [17]. Similar to the case of internal behaviors, in our setting it is more appropriate to talk about distributions of external behaviors as all of them are allowed in the ldt-MP. We refer to such distributions as *trace distributions* [27]: let us emphasize again that trace distributions are analogues of languages for non-stochastic systems. Since the map  $L_\omega : X^\omega \rightarrow Y^\omega$  is measurable [31, Theorem 1], for any ldt-MP its trace distribution can be expressed as  $(L_\omega)_* P_x$ . We can formulate the LT model-checking problem for ldt-MP as follows:

PROBLEM 1. *Given an ldt-MP  $(X, P, Y, L)$  and the observation-path dependent functional  $f \in \mathcal{B}(Y^\omega)$ , compute the expectation  $((L_\omega)_* P_x) f$  for any initial state  $x \in X$ .*

An important case of a functional  $f$  in Problem 1 is the indicator function  $f = 1_A$  of some event of interest  $A \in \mathcal{B}(Y^\omega)$ : in such a case the expectation to be computed turns out to be the probability  $((L_\omega)_* P_x)(A) = P_x(L_\omega(\mathbf{x}) \in A)$ . A common example of an observation space is given by finite sets  $Y$ , also called *alphabets*. Over alphabets the event  $A$  can be for instance an  $\omega$ -regular language<sup>4</sup> expressed as an automaton, or an LTL formula – for a detailed exposition see e.g. [8, Chapters 4, 5].

Dealing with stochastic systems at the level of observations allows one to compare systems possibly endowed with different state spaces. This feature is extremely important since a

<sup>4</sup> Indeed,  $\omega$ -regular languages over a finite  $Y$  are always elements of  $\mathcal{B}(Y^\omega)$ , for the proof see e.g. [36, Proposition 2.3].

complex ldt-MP can be approximated by a simpler one over a smaller state space – for example, by a finite ldt-MP: for the latter LT model-checking allows for analytical solutions, and numerical procedures can be computationally efficient. In order to be able to quantitatively argue about the original ldt-MP using results obtained over its abstract approximation, it is useful to endow the trace space with some metrics [16]. The choice of the latter depends on how one wants to interpret over the original system the results obtained over the abstraction. Along the lines of the discussion in Section 2.1, for ldt-MP we define these metrics to measure the difference (or similarity) between trace distributions, regardless of the way the latter are represented. The trace equivalence for any two ldt-MP can be defined by requiring them to have the same trace distributions, however there may be several choices of metrics whose zero level sets coincide with such trace equivalence. The next section proposes two such metrics for the trace distributions over the ldt-MP with the same observation spaces based on the *total variation* distance TV and on the *Wasserstein* distance W.

## 3. METRICS FOR LDT-MP

### 3.1 Total variation distance

Perhaps the most direct way to define a distance between two probability measures that fits the purposes of Problem 1 is to maximize over all functionals (or events) the difference between the corresponding expectations. Interestingly, such seemingly naïve approach yields a useful metric called the total variation distance (see Section 7.2).

Let us consider two ldt-MP  $\mathcal{D} = (X, P, Y, L)$  and  $\tilde{\mathcal{D}} = (\tilde{X}, \tilde{P}, Y, \tilde{L})$  over the same observation space  $Y$ . Given initial states  $x \in X$  and  $\tilde{x} \in \tilde{X}$ , we denote the trace distributions of  $\mathcal{D}$  and  $\tilde{\mathcal{D}}$  by  $Q_x$  and  $\tilde{Q}_{\tilde{x}}$ , respectively. Suppose for example that  $\tilde{\mathcal{D}}$  is finite,  $Y$  is a finite alphabet and  $\Phi$  is some LTL formula over  $Y$ . If we know the value of  $\text{TV}(Q_x, \tilde{Q}_{\tilde{x}})$ , we can compute  $\tilde{Q}_{\tilde{x}}(\Phi)$  and use it to estimate  $Q_x(\Phi)$ , since by definition of TV (12)

$$|Q_x(\Phi) - \tilde{Q}_{\tilde{x}}(\Phi)| \leq \frac{1}{2} \text{TV}(Q_x, \tilde{Q}_{\tilde{x}}).$$

From the definition of TV (for details see Section 7.2), it follows that similar bounds in terms of  $\text{TV}(Q_x, \tilde{Q}_{\tilde{x}})$  can be also obtained on the difference of expectations for more general cost functionals, rather than just indicator functions of LTL formulae. As a result, the TV satisfies the first criterion from Section 1 and can thus represent a good candidate for a metric. However, notice that  $\text{TV}(Q_x, \tilde{Q}_{\tilde{x}})$  quantifies the distance between probabilities on the infinite time horizon: this may be a too conservative requirement, as mentioned in [33, Section 3.1] and as the following example shows.

EXAMPLE 5. *Let us consider two simple ldt-MP with only two states:  $\tilde{X} = X = \{0, 1\}$ ,  $Y = X$ ,  $L = \tilde{L} = \text{id}_X$ , and suppose that transition matrices  $P \neq \tilde{P}$  have only positive entries. It follows that these ldt-MP are ergodic; we denote their unique invariant distributions by  $\mu$  and  $\tilde{\mu}$ , respectively. For any  $h : X^2 \rightarrow \mathbb{R}$  define*

$$A_h := \left\{ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} h(\mathbf{x}_k, \mathbf{x}_{k+1}) = \int_{X^2} h(x, y) P(x, dy) \mu(dy) \right\}.$$

*It holds that  $P_x(A_h) = 1$  for any initial state  $x$ . A similar fact can be obtained for  $\tilde{P}_x$ . Clearly, we can always find an  $h$  such that*

$$\int_{X^2} h(x, y) P(x, dy) \mu(dy) \neq \int_{X^2} h(x, y) \tilde{P}(x, dy) \tilde{\mu}(dy).$$

Since  $Q_x = P_x$  and  $\tilde{Q}_x = \tilde{P}_x$  we get  $\text{TV}(Q_x, \tilde{Q}_x) = 2$  by putting  $A_h$  in (12), no matter how small the difference between  $P$  and  $\tilde{P}$  is.

Due to the reasons discussed above, we focus our attention to finite-horizon behaviors of ldt-MP, characterized by restrictions of the trace distribution  $Q_x$  and  $\tilde{Q}_x$  to the set  $Y^{n+1}$ : we denote them by  $Q_x^n$  and  $\tilde{Q}_x^n$ , respectively. Since given an initial state over the concrete ldt-MP one has the freedom of choosing that over the abstraction, let us define the TV-like distance between ldt-MPs with the same observation spaces by

$$d_{\text{TV}^n}(\mathcal{D}, \tilde{\mathcal{D}}) := \sup_{x \in X} \inf_{\tilde{x} \in \tilde{X}} \text{TV}(Q_x^n, \tilde{Q}_{\tilde{x}}^n). \quad (6)$$

**DEFINITION 3.** We say that two ldt-MP  $\mathcal{D}$  and  $\tilde{\mathcal{D}}$  are  $\varepsilon$ -trace equivalent in the  $d_{\text{TV}^n}$  metrics if  $d_{\text{TV}^n}(\mathcal{D}, \tilde{\mathcal{D}}) \leq \varepsilon$ .

Although the approximate trace equivalence of ldt-MP in  $d_{\text{TV}^n}$  for  $n < \infty$  is not in general useful in infinite-horizon LT model-checking, it is sufficient to argue about finite-horizon properties, for example those expressed as BLTL formulae [31, Section 2.4]. Until the end of Section 3.1 we focus exclusively on the case when  $Y$  is a finite alphabet. A finite alphabet on a finite time horizon contains only a finite number of elements, which justifies the use of the following formula for the total variation [9]:

$$\text{TV}(Q_x^n, \tilde{Q}_{\tilde{x}}^n) = \sum_{y \in Y^{n+1}} |Q_x(\{y\}) - \tilde{Q}_{\tilde{x}}(\{y\})|. \quad (7)$$

Though the expression above is still challenging to compute precisely even over finite-state ldt-MPs, it can be computed by means of Monte Carlo simulations, which shows that  $d_{\text{TV}^n}$  satisfies the second criterion mentioned in Section 1.

Let us fix the time horizon  $n$ , and let us sample independently  $N$  copies of the observations of  $\mathcal{D}$  and  $\tilde{\mathcal{D}}$  over the given time horizon, which we further denote by  $(\mathbf{y}^i)_{i=1}^N$  and  $(\tilde{\mathbf{y}}^i)_{i=1}^N$  respectively.<sup>5</sup> The index  $i$  refers to different runs: each of them can be obtained e.g. by sampling the state-path of  $\mathcal{D}$ , say  $\mathbf{x}^i$ , on the time horizon  $n$ , and then by mapping  $\mathbf{x}^i$  into  $\mathbf{y}^i$  by means of map  $L$ . Define

$$Q_x^{n,N}(\cdot) := \frac{1}{N} \sum_{i=1}^N 1_{\{\mathbf{y}^i \in \cdot\}}, \quad \tilde{Q}_{\tilde{x}}^{n,N}(\cdot) := \frac{1}{N} \sum_{i=1}^N 1_{\{\tilde{\mathbf{y}}^i \in \cdot\}} \quad (8)$$

to be empirical distributions. It is fairly simple to compute the approximate distance  $\hat{\nu}_N := \text{TV}(Q_x^{n,N}, \tilde{Q}_{\tilde{x}}^{n,N})$  using (7), which leads to assess how good such an approximation of the original distance  $\nu := \text{TV}(Q_x^n, \tilde{Q}_{\tilde{x}}^n)$  is. Let us denote by  $\Pi$  the joint probability distribution of the two iid sequences  $(\mathbf{y}^i)_{i=1}^N$  and  $(\tilde{\mathbf{y}}^i)_{i=1}^N$ .

**THEOREM 1.** For any  $\delta > 0$  it holds that

$$\Pi(|\nu - \hat{\nu}_N| \leq 2\delta) \geq \left(1 - (N+1)^{|Y|^{n+1}} e^{-N\delta^2}\right)^2, \quad (9)$$

provided that  $(N+1)^{|Y|^{n+1}} e^{-N\delta^2} \leq 1$ .

**PROOF.** The idea of the proof is to start with  $\text{TV}(Q_x^n, Q_x^{n,N})$ :

$$\Pi\left(\text{TV}(Q_x^n, Q_x^{n,N}) \geq \delta\right) = \sum_{\nu \in L_N^\delta} \Pi\left(Q_x^{n,N} = \nu\right) \leq \sum_{\nu \in L_N^\delta} e^{-N \cdot \text{KL}(\nu, Q_x^n)},$$

<sup>5</sup> Note that for any  $i \in [1; N]$ , we have that  $\mathbf{y}^i$  is a vector of elements of  $Y$  of length  $n+1$ : that is a vector of observations of  $\mathcal{D}$  over the time horizon  $n$ .

where by KL we denote the Kullback-Leibler divergence between probability measures [15],  $L_N^\delta$  is the set of all empirical measures  $\nu$  satisfying  $\text{TV}(Q_x^n, \nu) \geq \delta$ , and the last inequality immediately follows from [9, Lemma 2.1.9]. Since for all  $\nu \in L_N^\delta$  it holds that  $\text{KL}(\nu, Q_x^n) \geq (\text{TV}(\nu, Q_x^n))^2$ , and [9, Lemma 2.1.2(a)] implies that  $|L_N^\delta| \leq (N+1)^{|Y|^{n+1}}$ , we obtain:

$$\Pi\left(\text{TV}(Q_x^n, Q_x^{n,N}) \geq \delta\right) \leq (N+1)^{|Y|^{n+1}} e^{-N\delta^2}. \quad (10)$$

A similar estimate obviously can be derived for  $\text{TV}(\tilde{Q}_{\tilde{x}}^n, \tilde{Q}_{\tilde{x}}^{n,N})$ . By triangular inequality we further get

$$|\nu - \hat{\nu}_N| \leq \text{TV}(Q_x^n, Q_x^{n,N}) + \text{TV}(\tilde{Q}_{\tilde{x}}^n, \tilde{Q}_{\tilde{x}}^{n,N}),$$

which together with (10) yields (9) as desired.  $\square$

Clearly, (9) implies that the estimator  $\hat{\nu}_N$  converges in  $\Pi$  to the true distance  $\nu$ , and further provides an explicit bound on the confidence level. Note also that for fixed sizes of the alphabet and of the time horizon, the bound in (9) depends quadratically on the precision  $\delta$ , and logarithmically on the confidence level. Unfortunately, such a dependence is only asymptotical, and for  $N$  not sufficiently large the polynomial term  $(N+1)^{|Y|^{n+1}}$  dominates the exponential term  $e^{-N\delta^2}$ . Moreover, such polynomial term grows very fast with respect to the size of the alphabet, and even faster with respect to the time horizon. For example, even if  $|Y| = 2$ ,  $n = 9$  and  $\delta = 0.1$ , one needs approximately an order of  $N = 1.5 \times 10^6$  samples to obtain good confidence levels, whereas for  $n = 10$  around  $N = 3.5 \times 10^6$  samples are required. In addition, the precision of the estimator  $\hat{\nu}_N$  can be only guaranteed with some confidence as its computation relies on randomized methods. Finally, Theorem 1 only gives a way to estimate  $\text{TV}(Q_x^n, \tilde{Q}_{\tilde{x}}^n)$  for given initial states, whereas (6) requires solving an optimization problem. All of this motivates looking into alternative methods for computing the  $d_{\text{TV}^n}$  metric.

We start with the case when both ldt-MPs  $\mathcal{D}$  and  $\tilde{\mathcal{D}}$  have the same state space, say  $X$ ; let  $P$  and  $\tilde{P}$  be the corresponding stochastic kernels. Recall that each of them acts as a linear operator on the Banach space  $\mathfrak{b}\mathcal{B}(X)$ , e.g.  $Pf(x) = \int_X f(y)P(x, dy)$ . Since  $\mathfrak{b}\mathcal{B}(X)$  is endowed with a sup-norm, one can define an induced norm on operators, as follows:

$$\|P - \tilde{P}\| := \sup_{f \in \mathfrak{b}\mathcal{B}(X), \|f\| \leq 1} \|(P - \tilde{P})f\|.$$

We obtain that  $\|P - \tilde{P}\| = \sup_{x \in X} \text{TV}(P(x, \cdot), \tilde{P}(x, \cdot))$  and as a result  $\|P - \tilde{P}\| \geq \text{TV}^1(Q_x, \tilde{Q}_{\tilde{x}})$ . We define a new metric as

$$d_B(\mathcal{D}, \tilde{\mathcal{D}}) := \|P - \tilde{P}\|.$$

The next theorem shows how the latter metric can be used in order to derive upper bounds on  $d_{\text{TV}^n}$ .

**THEOREM 2.** For any  $n \in \mathbb{N}_0$  it holds that

$$d_{\text{TV}^n}(\mathcal{D}, \tilde{\mathcal{D}}) \leq 2 - 2 \left(1 - \frac{1}{2} d_B(\mathcal{D}, \tilde{\mathcal{D}})\right)^n. \quad (11)$$

**PROOF.** From [7, Theorem 2] it follows that

$$\text{TV}(P_x^n, \tilde{P}_x^n) \leq 2 - 2 \left(1 - \frac{1}{2} d_B(\mathcal{D}, \tilde{\mathcal{D}})\right)^n,$$

for any initial state  $x \in X$ . Since  $Q_x^n$  and  $\tilde{Q}_x^n$  are image measures generated by  $P_x^n$  and  $\tilde{P}_x^n$  respectively, it holds that  $\text{TV}(Q_x^n, \tilde{Q}_x^n) \leq \text{TV}(P_x^n, \tilde{P}_x^n)$ . Finally, it clearly holds that

$$d_{\text{TV}^n}(\mathcal{D}, \tilde{\mathcal{D}}) \leq \sup_{x \in X} \text{TV}(Q_x^n, \tilde{Q}_x^n),$$

which further leads to (11).  $\square$

Before we elaborate on Theorem 2, let us first show an example of how it can be applied to construct finite abstractions with any given precision: that would assure that  $d_{\text{TVP}^n}$  satisfies the third criterion in Section 1. Suppose that we are given an ldt-MP  $\mathcal{D} = (X, Y, P, L)$ , where  $Y$  is a finite alphabet and  $P$  is an integral kernel, that is  $P(x, dy) = p(x, y)\mu(dy)$ . Here  $p$  is a jointly measurable function and  $\mu$  is a  $\sigma$ -additive measure on  $X$ . We say that  $p$  is a density of  $P$  with respect to the measure  $\mu$ : for example if  $X = \mathbb{R}^m$  and  $\mu$  is the Lebesgue measure,  $p$  is a common density function. Assume further that the state space  $X$  is endowed with some metric  $d_X$  compatible with its topology, and is bounded with respect to this metric. Let  $(X_i)_{i=1}^m$  be a finite partition of  $X$  such that  $L$  is constant when restricted to any  $X_i$ , and denote by  $\delta_i$  the diameter of  $X_i$  in the metric  $d_X$ . Let  $x_i \in X_i$  be arbitrary points, and define the finite ldt-MP  $\tilde{\mathcal{D}} = (\tilde{X}, \tilde{P}, Y, \tilde{L})$  as follows:  $\tilde{X} = [1; m]$ ,  $\tilde{P}(i, \{j\}) := P(x_i, X_j)$  and  $\tilde{L}(i) = L(x_i)$ . Notice that we cannot compare  $\mathcal{D}$  and  $\tilde{\mathcal{D}}$  directly using the metric  $d_B$ : the latter is only defined for ldt-MP that share the same state space. Due to this reason, below we construct an auxiliary ldt-MP  $\hat{\mathcal{D}}$  which has an infinite state space  $X$  so that  $d_B(\mathcal{D}, \hat{\mathcal{D}})$  is well-defined, but whose structure makes it trace-equivalent to the finite ldt-MP  $\tilde{\mathcal{D}}$ . We require the following Lipschitz-like condition.

**ASSUMPTION 1.** *There exist measurable non-negative functions  $\kappa_i : X \rightarrow [0, \infty)$ , such that  $K_i := \int_X \kappa_i(y)\mu(dy) < \infty$  for all  $i \in [1; m]$ , and such that*

$$|p(x', y) - p(x'', y)| \leq \kappa_i(y)d_X(x', x'') \quad \forall x', x'' \in X_i, \quad \forall y \in X.$$

**THEOREM 3.** *Under Assumption 1 for any  $n \in \mathbb{N}_0$  it holds that*

$$d_{\text{TVP}^n}(\mathcal{D}, \tilde{\mathcal{D}}) \leq 2 - 2 \left(1 - \max_{i \in [1; m]} K_i \delta_i\right)^n.$$

**PROOF.** We construct an auxiliary ldt-MP  $\hat{\mathcal{D}} = (X, \hat{P}, Y, L)$  with a lumpable structure to assure that  $d_{\text{TVP}^n}(\hat{\mathcal{D}}, \tilde{\mathcal{D}}) = 0$ , and to compare  $d_{\text{TVP}^n}(\mathcal{D}, \hat{\mathcal{D}})$  using Theorem 2. For this purpose we let  $\hat{P}(x, dy) = \hat{p}(x, y)\mu(dy)$  to be an integral kernel where the density  $\hat{p}$  is defined as  $\hat{p}(x, y) = p(x_i, y)$  for all points  $x \in X_i$  and  $y \in X$ . For any  $f \in \text{b}\mathcal{B}(X)$ ,  $\|f\| \leq 1$  and  $x \in X_i$  it holds that

$$\begin{aligned} |(P - \hat{P})f(x)| &= \left| \int_X f(y)(p(x, y) - p(x_i, y))\mu(dy) \right| \\ &\leq \int_X |f(y)|\delta_i \kappa_i(y)\mu(dy) \leq K_i \delta_i, \end{aligned}$$

and hence  $\|P - \hat{P}\| \leq \max_{i \in [1; m]} K_i \delta_i$ , so that

$$d_{\text{TVP}^n}(\mathcal{D}, \hat{\mathcal{D}}) \leq 2 - 2 \left(1 - \max_{i \in [1; m]} K_i \delta_i\right)^n$$

by Theorem 2. Let  $\iota : X \rightarrow \tilde{X}$  be the map defined by  $\iota(x) = i$  iff  $x \in X_i$ . For any initial state  $x \in X$  it holds that  $\mathbb{Q}_x^n = \tilde{\mathbb{Q}}_{\iota(x)}^n$ : indeed, these measures coincide on each element of the finite set  $Y^{n+1}$  by construction. As a result,  $d_{\text{TVP}^n}(\mathcal{D}, \tilde{\mathcal{D}}) = 0$ , which completes the proof of the theorem.  $\square$

This theorem shows that in case an original ldt-MP satisfies Assumption 1, for any  $\varepsilon > 0$  it is possible to construct a finite ldt-MP that is  $\varepsilon$ -trace equivalent to the original one in the  $d_{\text{TVP}^n}$  metric. As a result, the metric  $d_{\text{TVP}^n}$  satisfies all the criteria in Section 1. Note also that to estimate the metric  $d_{\text{TVP}^n}$ , which quantifies the difference between path measures, namely the *global* (in

time) dynamics, we have used the metric  $d_B$ , which compares the stochastic kernels, hence measuring the difference between transitions of ldt-MP – their *local* dynamics. The difference in local dynamics is often easier to assess, which in turn provides a method to argue about similarity of the global dynamics – see also the discussion in [16, Section II.B]. To formally speak about the similarity of the local dynamics for the processes with different state spaces, we would need to introduce notions of approximate bisimulations. For example, in such a case we could compare  $\mathcal{D}$  and  $\tilde{\mathcal{D}}$  in the proof of Theorem 3 directly, without the need to introduce the auxiliary ldt-MP  $\hat{\mathcal{D}}$ . However the discussion on approximate bisimulation goes beyond the scope of this contribution.

Theorems 2 and 3 further improve the formula-free abstraction procedure for ldt-MP introduced in [31]: the bounds on the path measures there increase linearly in the time horizon, whereas (11) provides bounds which never exceed 2 – the maximal meaningful value of the total variation distance. Moreover, such bounds are tight, that is in some special cases the equality in (11) holds.

## 3.2 Wasserstein distance

Above we have shown the usefulness of the metric  $d_{\text{TVP}^n}$  for LT model-checking. Although this metric is based on the total variation distance, which has an important characterization through coupling (13), we have not used extensively the notion of coupling in the discussion on  $d_{\text{TVP}^n}$  (though it is important in the proof [7, Theorem 2] used in Theorem 2).

The notion of coupling appears to be much more important for another metric that we consider next. Suppose that we are given two ldt-MP  $\mathcal{D} = (X, Y, P, L)$  and  $\tilde{\mathcal{D}} = (\tilde{X}, \tilde{P}, Y, \tilde{L})$  expressed as

$$\begin{cases} \mathbf{x}_{k+1} = F(\mathbf{x}_k, w_k), \\ \mathbf{y}_k = L(\mathbf{x}_k) \end{cases} \quad \begin{cases} \tilde{\mathbf{x}}_{k+1} = \tilde{F}(\tilde{\mathbf{x}}_k, \tilde{w}_k), \\ \tilde{\mathbf{y}}_k = \tilde{L}(\tilde{\mathbf{x}}_k) \end{cases}$$

where each of the  $(w_k)_{k \in \mathbb{N}_0}, (\tilde{w}_k)_{k \in \mathbb{N}_0}$  is a sequence of iid random variables. Suppose further that we are interested in approximating  $\mathcal{D}$  with  $\tilde{\mathcal{D}}$ . If the observation space  $Y$  is endowed with some metric  $d_Y$ , we can endow  $Y^{n+1}$  with the product metric

$$d_{Y^{n+1}}((y_0, \dots, y_n), (\tilde{y}_0, \dots, \tilde{y}_n)) := \max_{k \in [0; n]} d_Y(y_k, \tilde{y}_k).$$

If we are able to assure that  $d_{Y^{n+1}}(\mathbf{y}, \tilde{\mathbf{y}}) \leq \delta$  with a high probability, the fact that  $\tilde{\mathbf{y}}$  satisfies some property  $\Phi \in \text{b}\mathcal{B}(Y^{n+1})$  would imply that  $\mathbf{y}$  satisfies the modified (inflated) property

$$\Phi_\delta = \{y \in Y^{n+1} : d_{Y^{n+1}}(y, \Phi) \leq \delta\}.$$

See [19, Theorem 7] for the case when  $\Phi$  is a reachability specification. However, in order to talk about the probability of the value of  $d_{Y^{n+1}}(\mathbf{y}, \tilde{\mathbf{y}})$  being less than  $\delta$ , we need to consider a coupling between  $\mathbf{y}$  and  $\tilde{\mathbf{y}}$ . The work in [19] has considered the case of linear jump-diffusions, and implicitly suggested to use a coupling that matches the noises, that is  $\tilde{w}_k = w_k$ . In the follow-up paper [2] the coupling was considered to be the independent one (cf. Example 4). However, one has a complete freedom in choosing the coupling, so we can define a metric as follows:

$$d_{W^n}(\mathcal{D}, \tilde{\mathcal{D}}) := \sup_{x \in X} \inf_{\tilde{x} \in \tilde{X}} W(\mathbb{Q}_x^n, \tilde{\mathbb{Q}}_{\tilde{x}}^n),$$

thanks to the coupling characterization of the Wasserstein distance in (14). In case we know  $d_{W^n}(\mathcal{D}, \tilde{\mathcal{D}})$ , for any initial state  $x \in X$  we can find a corresponding  $\tilde{x} \in \tilde{X}$  such that

$$\mathbb{Q}(d_{Y^{n+1}}(\mathbf{y}, \tilde{\mathbf{y}}) \geq \delta) \leq \frac{d_{W^n}(\mathcal{D}, \tilde{\mathcal{D}})}{\delta},$$

thanks to Markov’s inequality, where  $\mathbb{Q} \in \Gamma(\mathbb{Q}_x^n, \tilde{\mathbb{Q}}_{\tilde{x}}^n)$  is any among the optimal couplings in (14). As a result, the metric  $d_{W^n}$  as well satisfies the first criterion in Section 1. Randomized methods to compute this metric can be found in [14, Sections 2.2, 2.3] and [28], thus assuring that  $d_{W^n}$  satisfies the second criterion we have raised. These methods are based on the representation of  $W$  as an integral metric (see Section 7.2), for example the one in [28] suggests solving a linear programming problem to evaluate  $W(\mathbb{Q}_x^{n,N}, \tilde{\mathbb{Q}}_{\tilde{x}}^{n,N})$  over empirical distributions (8).

Both  $TV$  and  $W$  are integral metrics, and for bounded metric spaces  $TV$  is always bigger than  $W$ , multiplied by a constant. Hence, one may expect that  $d_{W^n}$  provides less conservative estimates than  $d_{TV^n}$ . Unfortunately, although one can try defining an analogue of the transition-based metric  $d_B$  for  $d_{W^n}$ , obtaining estimates similar to Theorem 2 is not an easy task. Thus we are not able to claim that  $d_{W^n}$  satisfies the third criterion.

Let us also mention that the Wasserstein metric has also been applied for purposes of abstractions of controlled ldt-MP in [14] with focus on discounted additive cost functionals. Furthermore, [30] has employed the Wasserstein metric to compare outputs of stochastic processes. However, the results in both papers are not suitable for Problem 1.

## 4. COMPARISON WITH THE LITERATURE

### 4.1 Total Variation distance

Recent literature on approximate probabilistic bisimulations has led to several metrics for discrete-time stochastic systems that are related to the metrics  $d_B$  and  $d_{TV^n}$  considered here. Before we elaborate on their similarity, let us first recall a few of these metrics.

One of the first metrics for stochastic systems is introduced in [11] and [12] for finite and infinite systems, respectively. As we discussed in Section 1, the notion of exact probabilistic bisimulation, introduced for finite systems in [22] and extended to infinite systems in [10], appears to be too restrictive and to lack robustness: this motivates looking for more flexible relations between systems. The notion of exact probabilistic bisimulation is characterized via a binary logic  $\mathcal{L}$ , so that two states are bisimilar if and only if they satisfy the same formulae in  $\mathcal{L}$ . The extension of  $\mathcal{L}$  to some *real-valued* logic  $\mathcal{L}_r$  leads to the definition of a distance between two states  $x$  and  $\tilde{x}$  as the maximal difference between all formulae of  $\mathcal{L}_r$  computed over such states:  $d_{\mathcal{L}_r}(x, \tilde{x}) := \sup_{f \in \mathcal{L}_r} |f(x) - f(\tilde{x})|$ . The fact that  $\mathcal{L}_r$  is an extension of  $\mathcal{L}$  requires that  $d_{\mathcal{L}_r}$  is equal to zero over bisimilar states. The value of the metric between two stochastic systems is defined (over the disjoint union of the two systems) as the  $d_{\mathcal{L}_r}$  distance between their initial conditions.

A different approach is taken in [35], where stochastic systems are represented as co-algebras, namely pairs  $(X, P)$ , where  $X$  is a state space and  $P : X \rightarrow \mathcal{P}(X)$  (note the similarity with dt-MP models considered in this work). The state spaces are assumed to be endowed with a metric, and the functor  $\mathcal{P}$  is assumed to push the original metric of  $X$  into a Kantorovich metric over  $\mathcal{P}(X)$ .<sup>6</sup> It is further shown that a final co-algebra exists for such a functor, and since the final co-algebra is often thought of as a collection of behaviors, [35] introduces a metric between two states  $x, \tilde{x} \in X$  as  $d_{co}(x, \tilde{x}) = d_{fin}(!x, !\tilde{x})$ , where  $!$  is the unique map from  $X$  to the final co-algebra.

<sup>6</sup> Note however, that despite a strong connection between Kantorovich and Wasserstein metrics (see e.g. [14, Section 2.4.6]), this approach differs from the one we apply here in Section 3.2.

Although  $d_{\mathcal{L}_r}$  and  $d_{co}$  are introduced in completely different ways, it is shown in [35] that under certain conditions  $d_{co} = d_{\mathcal{L}_r}$ . This could be considered as evidence that any of these two (equal) metrics is the natural metric for the intended applications. However, it is of interest whether either of the two approaches is *the* natural one. Indeed, the definition of  $d_{\mathcal{L}_r}$  depends on the choice of the real-valued logic  $\mathcal{L}_r$ , which in [12] is chosen to be just one of many possible extensions of the binary logic  $\mathcal{L}$ . Similarly, given a metric on the Borel space  $X$ , there are many ways to introduce a related metric on  $\mathcal{P}(X)$ : besides the Kantorovich metric proposed in [35] one could employ a Prokhorov metric [15] which would lead to a different value for  $d_{co}$ . As convergence in Prokhorov metric is equivalent to weak convergence of measures, and hence to the topology of  $\mathcal{P}(X)$ , arguably the latter may be a more natural choice for a metric on  $\mathcal{P}(X)$  than the Kantorovich one. In conclusion, we contend that it is in general not possible to assert how good or natural the choice of a particular metric is, based on the way this metric is defined: the usefulness of a metric is rather determined by the applications it is meant to be used for. From this more practical perspective, it is unfortunately not clear which applications the models studied in [12, 35] are suitable for, besides their use in testing [34]: in fact, their semantics has not been defined explicitly and constructively. Besides, up to the best of the authors knowledge, in the case of testing there exist no precise bounds on the difference between testing probabilities over two system given their  $d_{co}$  or  $d_{\mathcal{L}_r}$  distances.

From the perspective of the preceding discussion, we aim next to elaborate on the similarity between  $d_{\mathcal{L}_r}$  and  $d_{TV^n}$ . Let us in particular show how  $d_{\mathcal{L}_r}$  can be applied on BLTL model-checking over ldt-MP (which are the models we have defined  $d_{TV^n}$  over). Since  $d_{\mathcal{L}_r}$  has been introduced over Labelled Markov Processes (LMP) [10], we need to provide a way of transform an ldt-MP into a corresponding LMP<sup>7</sup>. Consider a ldt-MP  $\mathcal{D} = (X, Y, P, L)$  with a finite observation space  $Y$ ; we define a corresponding LMP  $(X, Y, \{\tau_y\}_{y \in Y})$ , endowed with a state space  $X$ , a labels set  $Y$ , and the transition kernel

$$\tau_y(x, dx') := \mathbf{1}_{L^{-1}(y)}(x) \cdot P(x, dx'),$$

that is, we enable the label  $y$  in the LMP exactly in those states  $x$  that are labeled with  $y$  in  $\mathcal{D}$ . Suppose we would like to compute over  $\mathcal{D}$  a probability of some basic BLTL formula, that is a finite word  $y_0 \dots y_n \in Y^{n+1}$ : i.e. we are to find  $\mathbb{Q}_x^n(y_0 \dots y_n)$ , where  $\mathbb{Q}_x^n$  is the trace distribution  $\mathcal{D}$  induces on  $Y^{n+1}$ . As LMP do not have explicitly-defined trace semantics, we can interpret  $\mathbb{Q}_x^n(y_0 \dots y_n)$  as a value of the  $\mathcal{L}_r$  formula  $\langle y_0 \rangle \dots \langle y_n \rangle 1(x)$ .<sup>8</sup> The latter interpretation tells us that the LMP can be applied to compute probabilities of basic BLTL formulae over a corresponding ldt-MP, and hence it can be applied to any BLTL formula [31, Section 2.4]. Note however, that probabilities of more general BLTL formulae over ldt-MP are not necessarily elements of  $\mathcal{L}_r$  for the corresponding LMP: the disjunction and conjunction in BLTL are not related to max and min in  $\mathcal{L}_r$  at all. As a result, although any BLTL formula over an ldt-MP can be computed via functions in  $\mathcal{L}_r$  over the corresponding LMP, such computations may require summation of  $\mathcal{L}_r$  functions. Hence, if any function in  $\mathcal{L}_r$  is known only with some error, this operation will require accumulating errors – this issue was already discussed in [31, Section 3] concerning the possible extensions of safety-focused

<sup>7</sup> Up to our knowledge, this paper is the first to explicitly elaborate on such a transformation.

<sup>8</sup> We assume here a version of  $\mathcal{L}_r$  without discounting.

results in [4] to the whole BLTL. Thus, although [12, Proposition 7.5] concerning  $d_{\mathcal{L}}$  provides an analysis similar to Theorem 2, the bounds from the former result when interpreted over ldt-MP only apply to basic BLTL formulae, whereas the latter result provides bounds for any BLTL formula.

Let us also mention that Theorems 2 and 3 concerning the  $d_{TV^n}$  metric improve formula-free guarantees for ldt-MP introduced in [31] by providing tighter bounds on the difference between trace distributions: indeed, [31, Lemma 1, Theorem 3] derives similar bounds as Theorem 2, but provide more conservative results. An alternative approach [5] suggests expressing an LT property as an automaton, and solving a safety problem over the product system. Unfortunately, the overall error associated to the abstraction needed for the solution of the safety problem depends on the size of the automaton, which is particularly crucial in case of BLTL formulae for which automata can be large [31, Section 3.4]. The formula-free guarantees do not depend on a particular BLTL formula (as the name suggests), and provide arguably less conservative bounds than safety-based approximate model-checking of complex LT properties (cf. [31, Section 5]).

## 4.2 Wasserstein distance

As we have mentioned in Section 3.2, the definition of the  $d_{W^n}$  metric is inspired by the work [19], which introduces an analogous metric restricted to a fixed coupling structure. We contend that this feature substantially increases the conservatism of the results, since in general one has the freedom to optimize over the coupling. The notions in [19] hinge upon similar ones developed over non-stochastic systems [16], where the coupling is necessarily unique so that its choice does not play a role. The choice of coupling for stochastic processes is especially important when it is not meaningful to claim that both processes are driven by the same noise. For instance, if it is clear how to compare noises that drive two diffusions (cf. Example 4), it is much less clear how to couple a diffusion with a finite-space continuous-time Markov Chain serving as its approximation [2]: in the latter case the only natural coupling seems to be the independent one, and it is rather unlikely that this coupling is the optimal one.

The work in [19] suggests how to compute a metric over a fixed coupling under rather strong assumptions on the stability of the models (over a fixed coupling). In particular, these techniques require Lyapunov-like functions to synthesize a metric, and provide practically relevant sufficient conditions for these functions only in the linear case. In contrast, [6] introduces randomized methods to compute fixed coupling metrics without any stability assumptions – as in our case the results are valid over a finite time horizon. It is of interest to compare those results with ours, to see whether the choice of the coupling leads to a reduced conservatism, as expected.

Let us consider an example drawn from [1] and also studied in [32]. We consider a regional power network consisting of two local subnetworks: each of them has its own energy storage capacity. There are two sources of the energy: a coal plant shared by both networks, and two separate wind farms producing renewable energy. [32] studied the approximate abstraction of this model, with the objective of coal plant energy production and distribution optimizing some desired property over the whole network. On the other hand, here we are interested in studying the effect of the correlation over the energy produced by the two wind farms – it practically makes sense to consider this correlation when the two local subnetworks are close geo-

graphically. We assume that the coal plant energy is evenly distributed between both networks, which leads to the following (autonomous) dt-MP model of the network:

$$\begin{cases} \mathbf{x}_{k+1}^1 = \beta \left( \mathbf{x}_k^1 + \frac{1}{2} p_k + r_k^1 - d_k^1 \right) \vee 0 \wedge M, \\ \mathbf{x}_{k+1}^2 = \beta \left( \mathbf{x}_k^2 + \frac{1}{2} p_k + r_k^2 - d_k^2 \right) \vee 0 \wedge M, \end{cases}$$

where  $\mathbf{x}^i$  is the amount of the stored energy in the  $i$ -th subnetwork,  $\beta = 0.8$  is the loss rate of the stored energy,  $p_k$  is the energy produced by the coal plant,  $r_k^i$  is a renewable energy produced by the  $i$ -th wind farm,  $d_k^i$  is the energy demand in the  $i$ -th network, and  $M = 30$  is the max storage capacity. We assume that  $p_k \sim U([3, 5])$  and  $d_k^i \sim U([1, 2])$  are distributed uniformly. The distribution of the renewable energy instead consists of two components:  $r_k^i = (1 - \rho)e_k^i + \rho r_k$ , where  $e_k^i, r_k \in \mathcal{E}(1)$  are distributed exponentially and  $\rho \in [0, 1]$  measures the correlation level of the distributions  $r_k^1$  and  $r_k^2$ . The initial conditions  $\mathbf{x}_0^i \sim U([5, 10])$  are distributed uniformly.

We study two copies of the model distinguished by the choice of parameter  $\rho = 0.2$  and  $\rho = 0$  respectively. Recall that we only know the distribution of the stochastic process associated to each of the two models, but not their joint distribution: this allows to choose the coupling. We then consider three possible choices: coupling via the same noise, independence coupling, and reflected coupling. In the former case we use the same samples of  $\tilde{d}_k^i = d_k^i$ ,  $\tilde{p}_k = p_k$  and  $\tilde{\mathbf{x}}_0^i = \mathbf{x}_0^i$  for both systems, whereas in the latter case we use the symmetry of the uniform distribution to define  $\tilde{d}_k^i, \tilde{p}_k$  and  $\tilde{\mathbf{x}}_0^i$  as  $d_k^i, p_k$  and  $\mathbf{x}_0^i$  reflect with respect to their mean values. Finally, in case of the independence coupling all samples are assumed to be independent, that is  $\tilde{d}_k^i \perp d_k^i$ ,  $\tilde{p}_k \perp p_k$  and  $\tilde{\mathbf{x}}_0^i \perp \mathbf{x}_0^i$ . We assume that the observation space corresponds to the state space endowed with the Euclidean metric, thus the observation map is the identity function.

The results are presented in Figure 1. The computations are run over a time horizon  $n = 100$  for  $N \in [2; 1000]$  samples. From the figure it can be seen that over the number of samples considered, for each estimator the convergence does hold. One can also observe that although the coupling via the same noise performs better than the independence coupling, or than that via reflected noise, their associated distance is greater than the distance  $d_{W^{100}}$  obtained by optimizing over all the possible couplings. This optimization is performed using the method in [28], by means of solving a linear programming problem: the latter clearly represents the computational bottleneck of the method, as it requires  $2N$  variables and  $N(N - 1)$  constraints. The complexity then depends on the number of samples, and although there is no explicit connection between the complexity of the linear programming problem and the dimension of original systems (or the time horizon), the latter factor may affect the number of samples needed for convergence to the true solution. Unfortunately, up to our knowledge there are no explicit results on the convergence for the W metric, similar to Theorem 1 for the TV metric. Methods mentioned in [14, Sections 2.2, 2.3] can be explored as possible computationally efficient alternatives. Let us finally mention that we do not provide here a comparison over the  $d_{TV^n}$  metric as it would not enlighten the importance of coupling as significantly as the  $d_{W^n}$  metric does.

## 5. CONCLUSIONS

This paper has discussed the formulation of two new metrics for stochastic transition systems, with focus on their applications for linear-time model-checking. Such metrics are shown to be useful according to three criteria defined in the article, and are

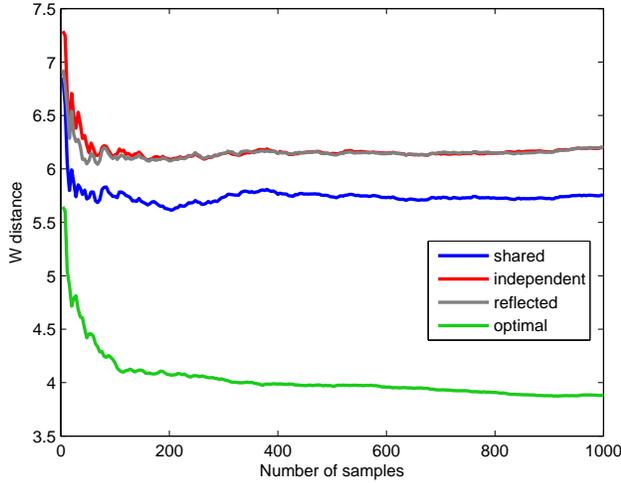


Figure 1: Evaluation of expected distances between  $x$  and  $\bar{x}$  depending on the number of samples, over different couplings: shared (blue line in the middle), independent and reflected noises (resp. red and gray lines on the top), optimal coupling (green line at the bottom).

claimed to perform better compared to other metrics proposed in the literature. Table 1 further shows which of those criteria the metrics introduced in this paper satisfy, and according to which result. A promising direction of future research concerns the connections of the metrics proposed in this work with notions of approximate probabilistic simulation and bisimulation.

Criterion	$TV^n$	$W^n$
[i.]	YES. By definition.	YES. By definition.
[ii.]	YES. By Theorem 1.	YES. See e.g. [28].
[iii.]	YES. By Theorems 2, 3.	???

Table 1: Satisfaction of criteria by metrics  $TV^n$  and  $W^n$ .

## 6. REFERENCES

- [1] MoVeS website. <http://www.movesproject.eu>.
- [2] A. Abate. A contractivity approach for probabilistic bisimulations of diffusion processes. In *Proceedings of the 48th IEEE Conference on Decision and Control and 28th Chinese Control Conference*, pages 2230–2235, 2009.
- [3] A. Abate. Approximation metrics based on probabilistic bisimulations for general state-space markov processes: a survey. *Electronic Notes in Theoretical Computer Sciences*, 2013. In Print.
- [4] A. Abate, J.-P. Katoen, J. Lygeros, and M. Prandini. Approximate model checking of stochastic hybrid systems. *European Journal of Control*, 16:624–641, 2010.
- [5] A. Abate, J.-P. Katoen, and A. Mereacre. Quantitative automata model checking of autonomous stochastic hybrid systems. In *Proceedings of the 14th international conference on Hybrid Systems: Computation and Control, HSCC '11*, pages 83–92, New York, NY, USA, 2011. ACM.
- [6] A. Abate and M. Prandini. Approximate abstractions of stochastic systems: A randomized method. In *Proceedings of the 50th IEEE Conference on Decision and Control and European Control Conference*, pages 4861–4866, 2011.
- [7] A. Abate, F. Redig, and I. Tkachev. On the effect of perturbation of conditional probabilities in total variation. *Statistics & Probability Letters*, 2014. In Press.
- [8] C. Baier and J.-P. Katoen. *Principles of model checking*. The MIT Press, 2008.
- [9] A. Dembo and O. Zeitouni. *Large deviations techniques and applications*. Springer-Verlag, Berlin, 2010.
- [10] J. Desharnais, A. Edalat, and P. Panangaden. Bisimulation for Labelled Markov Processes. *Information and Computation*, 179(2):163 – 193, 2002.
- [11] J. Desharnais, V. Gupta, R. Jagadeesan, and P. Panangaden. Metrics for Labeled Markov Systems. In *Proceedings of the 10th International Conference on Concurrency Theory, CONCUR '99*, pages 258–273, London, UK, 1999. Springer-Verlag.
- [12] J. Desharnais, V. Gupta, R. Jagadeesan, and P. Panangaden. Metrics for labelled Markov processes. *Theoretical Computer Science*, 318(3):323–354, 2004.
- [13] S. N. Ethier and T. G. Kurtz. *Markov processes: characterization and convergence*. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. John Wiley & Sons Inc., New York, 1986.
- [14] N. Ferns, P. Panangaden, and D. Precup. Bisimulation metrics for continuous Markov decision processes. *SIAM Journal on Computing*, 40(6):1662–1714, 2011.
- [15] A. L. Gibbs and F. E. Su. On choosing and bounding probability metrics. *International Statistical Review*, 70(3):419–435, 2002.
- [16] A. Girard and G. J. Pappas. Approximation metrics for discrete and continuous systems. *IEEE Transactions on Automatic Control*, 52(5):782 –798, 2007.
- [17] A. Girard and G. J. Pappas. Approximate bisimulation: A bridge between computer science and control theory. *European Journal of Control*, 17(5-6):568–578, 2011.
- [18] A. Hinton, M. Kwiatkowska, G. Norman, and D. Parker. PRISM: A tool for automatic verification of probabilistic systems. In H. Hermanns and J. Palsberg, editors, *Tools and Algorithms for the Construction and Analysis of Systems*, volume 3920 of *Lecture Notes in Computer Science*, pages 441–444. Springer Verlag, 2006.
- [19] A.A. Julius and G.J. Pappas. Approximations of stochastic hybrid systems. *IEEE Transactions on Automatic Control*, 54(6):1193–1203, 2009.
- [20] O. Kallenberg. *Foundations of modern probability*. Probability and its Applications. Springer-Verlag, New York, 1997.
- [21] J.-P. Katoen, M. Khattri, and I. S. Zapreev. A Markov reward model checker. In *Proceedings of the Second International Conference on the Quantitative Evaluation of Systems, QEST '05*, pages 243–244, Washington, DC, USA, 2005. IEEE Computer Society.
- [22] K. G. Larsen and A. Skou. Bisimulation through probabilistic testing. *Information and Computation*, 94(1):1–28, 1991.
- [23] T. Lindvall. *Lectures on the coupling method*. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. John Wiley & Sons Inc., New York, 1992.

- [24] B. Øksendal. *Stochastic differential equations: an introduction with applications*. Universitext. Springer-Verlag, Berlin, sixth edition, 2003.
- [25] K. R. Parthasarathy. *Probability measures on metric spaces*. Probability and Mathematical Statistics, No. 3. Academic Press Inc., New York, 1967.
- [26] D. Revuz. *Markov chains*. North-Holland Publishing, Amsterdam, second edition, 1984.
- [27] R. Segala. A compositional trace-based semantics for probabilistic automata. In I. Lee and S. A. Smolka, editors, *CONCUR '95: Concurrency Theory*, volume 962 of *Lecture Notes in Computer Science*, pages 234–248. Springer Berlin Heidelberg, 1995.
- [28] B. K. Sriperumbudur, K. Fukumizu, A. Gretton, B. Schölkopf, and G. Lanckriet. On the empirical estimation of integral probability metrics. *Electronic Journal of Statistics*, 6:1550–1599, 2012.
- [29] P. Tabuada. *Verification and control of hybrid systems: A symbolic approach*. Springer Verlag, New York, 2009.
- [30] D. Thorsley and E. Klavins. Approximating stochastic biochemical processes with Wasserstein pseudometrics. *Systems Biology, IET*, 4(3):193–211, May 2010.
- [31] I. Tkachev and A. Abate. Formula-free finite abstractions for linear temporal verification of stochastic hybrid systems. In *Proceedings of the 16th international conference on Hybrid Systems: Computation and Control, HSCC '13*, pages 283–292, New York, NY, USA, 2013. ACM.
- [32] I. Tkachev, A. Mereacre, J.-P. Katoen, and A. Abate. Quantitative automata-based controller synthesis for non-autonomous stochastic hybrid systems. In *Proceedings of the 16th international conference on Hybrid Systems: Computation and Control, HSCC '13*, pages 293–302, New York, NY, USA, 2013. ACM.
- [33] Ilya Tkachev and Alessandro Abate. Characterization and computation of infinite-horizon specifications over markov processes. *Theoretical Computer Science*, 515(0):1–18, 2014.
- [34] F. van Breugel, S. Shalit, and J. Worrell. Testing labelled Markov processes. In *Proceedings of the 29th International Colloquium on Automata, Languages and Programming, ICALP '02*, pages 537–548, London, UK, 2002. Springer-Verlag.
- [35] F. van Breugel and J. Worrell. Towards quantitative verification of probabilistic transition systems. In *Proceedings of the 28th International Colloquium on Automata, Languages and Programming, ICALP '01*, pages 421–432, London, UK, 2001. Springer-Verlag.
- [36] M.Y. Vardi. Automatic verification of probabilistic concurrent finite state programs. In *26th Annual Symposium on Foundations of Computer Science*, pages 327–338, 1985.

## 7. APPENDIX

### 7.1 General notation

We use extensively standard notions from measure and probability theory, for precise definitions see e.g. [20]. A measurable space is a pair  $(E, \mathcal{E})$  where  $E$  is an arbitrary set and  $\mathcal{E}$  is a  $\sigma$ -algebra on  $E$ . The space of all probability measures on  $(E, \mathcal{E})$  is denoted by  $\mathcal{P}(E, \mathcal{E})$ ; it is assumed to be endowed with the smallest  $\sigma$ -algebra that makes any evaluation map  $\theta_A : \mathcal{P}(E, \mathcal{E}) \rightarrow \mathbb{R}$ , defined by  $\theta_A(p) = p(A)$  for all  $p \in \mathcal{P}(E, \mathcal{E})$  and

$A \in \mathcal{E}$ , measurable. A stochastic kernel on  $(E, \mathcal{E})$  is a measurable map  $K : E \rightarrow \mathcal{P}(E, \mathcal{E})$ . A probability space is a triple  $(E, \mathcal{E}, p)$ , where  $(E, \mathcal{E})$  is a measurable space and  $p \in \mathcal{P}(E, \mathcal{E})$ . If  $(\Omega, \mathcal{F}, P)$  is a probability space and  $f : (\Omega, \mathcal{F}) \rightarrow (E, \mathcal{E})$  is a measurable map then

$$(f_*P)(A) := P(f^{-1}(A)), \quad \forall A \in \mathcal{E}$$

defines the probability measure in  $\mathcal{P}(E, \mathcal{E})$ . We say that  $f$  induces the measure  $f_*P$ . We further denote by  $\text{id}_\Omega$  the identity map on  $\Omega$ ; clearly, it holds that  $(\text{id}_\Omega)_*P = P$  for any  $P \in \mathcal{P}(\Omega, \mathcal{F})$ . By  $\otimes$  we denote the product of measures or  $\sigma$ -algebras. The set of real numbers is denoted by  $\mathbb{R}$ , the set of natural numbers by  $\mathbb{N}$  and we write  $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$ . We write  $E^\omega$  in place of  $E^{\mathbb{N}_0}$ .

Any topological space  $X$  is assumed to be endowed with its Borel  $\sigma$ -algebra  $\mathcal{B}(X)$  generated by all open subsets of  $X$ . For topological spaces we simply write  $\mathcal{P}(X)$  in place of  $\mathcal{P}(X, \mathcal{B}(X))$ . We refer to elements of  $\mathcal{B}(X)$  as Borel subsets of  $X$ . A (standard) Borel space is a topological space that is homeomorphic to a Borel subset of complete separable metric space. By  $\mathcal{C}(X)$  we denote the linear space of all continuous real-valued functions on  $X$ .

Given a measurable space  $(E, \mathcal{E})$ , we denote by  $\text{b}\mathcal{E}$  the linear space of all bounded functions  $f : E \rightarrow \mathbb{R}$ . It is a Banach space endowed with a sup-norm  $\|f\| := \sup_{x \in E} |f(x)|$ .

### 7.2 Distances between probability measures

For a brief overview on distance between probability measures see [15]; a detailed exposition can be found in [25]. In this paper we focus on two metrics: the total variation distance TV that applies to measures over any measurable spaces, and the Wasserstein distance W that requires the underlying space to be a metric space. Both metrics fall into the class of integral probability metrics: let  $(E, \mathcal{E})$  be an some measurable space and let  $\mathcal{G}$  be any collection real-valued measurable functions on  $E$ . The corresponding metric on  $\mathcal{P}(E, \mathcal{E})$  is defined as

$$I_{\mathcal{G}}(\mu, \tilde{\mu}) := \sup_{g \in \mathcal{G}} \left| \int_E g d(\mu - \tilde{\mu}) \right|.$$

The total variation metric is given by  $\text{TV} = I_{\text{b}\mathcal{E}_1}$  where  $g \in \text{b}\mathcal{E}_1$  iff  $g \in \text{b}\mathcal{E}$  and  $\|g\| \leq 1$ . Alternatively, the metric TV can be introduced via the difference of  $\mu$  and  $\tilde{\mu}$  over sets as follows:

$$\text{TV}(\mu, \tilde{\mu}) = 2 \cdot \sup_{A \in \mathcal{E}} |\mu(A) - \tilde{\mu}(A)|. \quad (12)$$

Unlike the total variation distance that depends exclusively on the measurability structure of the set  $E$ , the Wasserstein distance requires  $E$  to be a metric space. Let  $E$  be endowed with the metric  $d_E$  such that  $\mathcal{E} = \mathcal{B}(E)$ . For any  $g \in \text{b}\mathcal{E}$  let

$$\text{Lip}(g) := \sup_{x \neq x'} \frac{|g(x) - g(x')|}{d_E(x, x')}$$

denote the Lipschitz constant of  $g$ . The Wasserstein distance is defined as  $W = I_{\mathcal{G}_1}$ , where  $g \in \mathcal{G}_1$  iff  $\text{Lip}(g) \leq 1$ . Provided that  $\text{diam}(E) := \sup_{x, x' \in E} d_E(x, x') < \infty$ , the following inequality relates the two metrics:  $\text{TV} \geq \frac{1}{\text{diam}(E)} W$  [15], so that the convergence in TV implies the convergence in W in case the latter is defined using a metric in which  $E$  is bounded.

Another relation between TV and W can be derived via recalling that each of the metric can be characterized by a coupling, optimal in a certain sense. For  $\mu, \tilde{\mu} \in \mathcal{P}(E, \mathcal{E})$  let us denote by

$\Gamma(\mu, \tilde{\mu}) \subset \mathcal{P}(E^2, \mathcal{E}^2)$  the collection of all coupling measures for  $\mu$  and  $\tilde{\mu}$ . Provided that  $E$  is a Borel space and  $\mathcal{E} = \mathcal{B}(E)$ :

$$\text{TV}(\mu, \tilde{\mu}) = 2 \cdot \sup_{M \in \Gamma(\mu, \tilde{\mu})} M(\Delta_E) \quad (13)$$

and if  $(E, d_E)$  is a bounded complete separable metric space:

$$W(\mu, \tilde{\mu}) = \inf_{M \in \Gamma(\mu, \tilde{\mu})} \int_{E \times E} d_E(x, \tilde{x}) M(dx \times d\tilde{x}). \quad (14)$$

If we consider any two random elements  $f$  and  $\tilde{f}$  that represent distributions  $\mu$  and  $\tilde{\mu}$ , then the total variation distance can be obtained via the coupling that maximizes the probability that  $f = \tilde{f}$ , whereas the Wasserstein distance is exactly the minimal expected value of  $d_E(f, \tilde{f})$  over all possible couplings.