# TRACE REFINEMENT IN
# LABELLED MARKOV DECISION PROCESSES

NATHANAËL FIJALKOW [a,b], STEFAN KIEFER [c], AND MAHSA SHIRMOHAMMADI [c,d]

[a] CNRS, LaBRI, Bordeaux

[b] Alan Turing Institute of data science, London

[c] University of Oxford, UK

[d] CNRS, IRIF, Paris

ABSTRACT. Given two labelled Markov decision processes (MDPs), the trace-refinement problem asks whether for all strategies of the first MDP there exists a strategy of the second MDP such that the induced labelled Markov chains are trace-equivalent. We show that this problem is decidable in polynomial time if the second MDP is a Markov chain. The algorithm is based on new results on a particular notion of bisimulation between distributions over the states. However, we show that the general trace-refinement problem is undecidable, even if the first MDP is a Markov chain. Decidability of those problems was stated as open in 2008. We further study the decidability and complexity of the trace-refinement problem provided that the strategies are restricted to be memoryless.

## 1. INTRODUCTION

We consider labelled Markov chains (MCs) whose transitions are labelled with symbols from an alphabet $\mathsf{L}$. Upon taking a transition, the MC emits the associated label. In this way, an MC defines a *trace-probability* function $Tr : \mathsf{L}^* \to [0,1]$ which assigns to each finite trace $w \in \mathsf{L}^*$ the probability that the MC emits $w$ during its first $|w|$ transitions. Consider the MC depicted in Figure 1 with initial state $p_0$. For example, in state $p_0$, with probability $\frac{1}{4}$, a transition to state $p_c$ is taken and $c$ is emitted. We have, e.g., $Tr(abc) = \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4}$. Two MCs over the same alphabet $\mathsf{L}$ are called *equivalent* if their trace-probability functions are equal.

The study of labelled MCs and their equivalence has a long history, going back to Schützenberger [21] and Paz [18] who studied *weighted* and *probabilistic* automata, respectively. Those models generalize labelled MCs, but the respective equivalence problems are essentially the same. It can be extracted from [21] that equivalence is decidable in polynomial

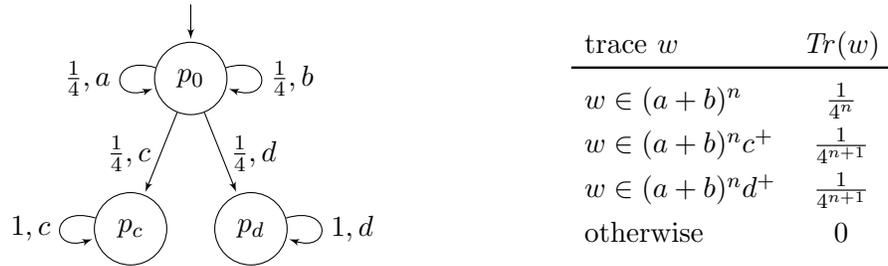| trace $w$ | $Tr(w)$ |
|---|---|
| $w \in (a+b)^n$ | $\frac{1}{4^n}$ |
| $w \in (a+b)^n c^+$ | $\frac{1}{4^{n+1}}$ |
| $w \in (a+b)^n d^+$ | $\frac{1}{4^{n+1}}$ |
| otherwise | $0$ |

Figure 1: An MC with its trace-probability function. This MC, denoted by $\mathcal{C}(\mathcal{A})$, will also be used in Section 3 for the reduction from universality of probabilistic automata to the trace-refinement problem.

time, using a technique based on linear algebra. Variants of this technique were developed in [23, 7]. Tzeng [24] considered the path-equivalence problem for nondeterministic automata which asks, given nondeterministic automata $A$ and $B$, whether each word has the same number of accepting paths in $A$ as in $B$. He gives an NC algorithm[1] for deciding path equivalence which can be straightforwardly adapted to yield an NC algorithm for equivalence of MCs.

More recently, the efficient decidability of the equivalence problem was exploited, both theoretically and practically, for the verification of probabilistic systems, see, e.g., [13, 14, 19, 17, 16]. In those works, equivalence naturally expresses properties such as obliviousness and anonymity, which are difficult to formalize in temporal logic. The *inclusion problem* for two probabilistic automata asks whether for each word the acceptance probability in the first automaton is less than or equal to the acceptance probability in the second automaton. Despite its semblance to the equivalence problem, the inclusion problem is undecidable [5], even for automata of fixed dimension [2]. This is unfortunate, especially because deciding language inclusion is often at the heart of verification algorithms.

We study another "inclusion-like" generalization of the equivalence problem: trace refinement in labelled Markov decision processes (MDPs). MDPs extend MCs by nondeterminism, and labelled MDPs generate outputs (labels); thus labelled MDPs are a generative model with nondeterminism. In each state, a controller chooses, possibly randomly and possibly depending on the history, one out of finitely many *moves*[2]. A move determines a probability distribution over the emitted label and the successor state. In this way, an MDP and a strategy of the controller induce an MC.

The *trace-refinement problem* asks, given two MDPs $\mathcal{D}$ and $\mathcal{E}$, whether for all strategies for $\mathcal{D}$ there is a strategy for $\mathcal{E}$ such that the induced MCs are equivalent. Consider the MDP depicted in Figure 2 where in state $q_1$ there are two available moves; one move generates the label $c$ with probability 1, the other move generates $d$ with probability 1. A strategy of the controller that, in state $q_1$, chooses the last generated label (either $c$ or $d$) with probability 1, induces the same trace-probability function as the MC shown in Figure 1; the MDP thus refines that MC. The described strategy needs one bit of memory to keep track of the last generated label. It was shown in [7] that the strategy for $\mathcal{E}$ may require *infinite memory*,

---

[1]The complexity class NC is the subclass of P containing those problems that can be solved in polylogarithmic parallel time (see, e.g., [10]).

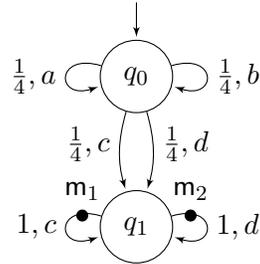[2]As in [7] we speak of moves rather than of actions, to avoid possible confusion with the label alphabet L.

Figure 2: An MDP where the choice of controller is relevant only in $q_1$. Two available moves $\mathsf{m}_1, \mathsf{m}_2$ are shown with small black circles.

even if $\mathcal{D}$ is an MC. The decidability of trace refinement was posed as an open problem, both in the introduction and in the conclusion of [7]. The authors of [7] also ask about the decidability of subcases, where $\mathcal{D}$ or $\mathcal{E}$ are restricted to be MCs. In this paper we answer all those questions. We show that trace refinement is undecidable, even if $\mathcal{D}$ is an MC. In contrast, we show that trace refinement is decidable efficiently (in $\mathsf{NC}$, hence in $\mathsf{P}$), if $\mathcal{E}$ is an MC. Moreover, we prove that the trace-refinement problem becomes decidable if one imposes suitable *restrictions on the strategies* for $\mathcal{D}$ and $\mathcal{E}$, respectively. More specifically, we consider *memoryless* (i.e., no dependence on the history) and *pure memoryless* (i.e., no randomization and no dependence on the history) strategies, establishing various complexity results between $\mathsf{NP}$ and $\mathsf{PSPACE}$.

To obtain the aforementioned $\mathsf{NC}$ result, we demonstrate a link between trace refinement and a particular notion of *bisimulation* between two MDPs that was studied in [11]. This variant of bisimulation is not defined between two states as in the usual notion, but between two *distributions* on states. An exponential-time algorithm that decides (this notion of) bisimulation was provided in [11]. We sharpen this result by exhibiting a $\mathsf{coNP}$ algorithm that decides bisimulation between two MDPs, and an $\mathsf{NC}$ algorithm for the case where one of the MDPs is an MC. For that we refine the arguments devised in [11]. The model considered in [11] is more general than ours in that they also consider continuous state spaces, but more restricted than ours in that the label is determined by the move.

## 2. Preliminaries

A trace over a finite set $\mathsf{L}$ of labels is a finite sequence $w = a_1 \cdots a_n$ of labels where the length of the trace is $|w| = n$. The empty trace $\epsilon$ has length zero. For $n \geq 0$, let $\mathsf{L}^n$ be the set of all traces with length $n$; we denote by $\mathsf{L}^*$ the set of all (finite) traces over $\mathsf{L}$.

For a function $d : S \to [0, 1]$ over a countable set $S$, define the *norm* $\|d\| := \sum_{s \in S} d(s)$. The *support* of $d$ is the set $\mathsf{Supp}(d) = \{s \in S \mid d(s) > 0\}$. The function $d$ is a *probability subdistribution* over $S$ if $\|d\| \leq 1$; it is a *probability distribution* if $\|d\| = 1$. We denote by $\mathsf{subDist}(S)$ (resp. $\mathsf{Dist}(S)$) the set of all probability subdistributions (resp. distributions) over $S$. Given $s \in S$, the *Dirac distribution* on $s$ assigns probability 1 to $s$; we denote it by $d_s$. For a non-empty finite subset $T \subseteq S$, the *uniform distribution* over $T$ assigns probability $\frac{1}{|T|}$ to every element in $T$.

2.1. **Labelled Markov Decision Processes.** In this article a *labelled Markov decision process* (MDP) is a generative probabilistic model with nondeterminism. Formally, an MDP is a quadruple $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$, consisting of a finite set $Q$ of states, an initial distribution $\mu_0 \in \mathsf{Dist}(Q)$, a finite set $\mathsf{L}$ of labels, and a finite probabilistic transition relation $\delta \subseteq Q \times \mathsf{Dist}(\mathsf{L} \times Q)$ where states are in relation with distributions over pairs of labels and successors. We assume that for each state $q \in Q$ there exists some distribution $d \in \mathsf{Dist}(\mathsf{L} \times Q)$ where $\langle q, d \rangle \in \delta$. The set of *moves* in $q$ is $\mathsf{moves}(q) = \{d \in \mathsf{Dist}(\mathsf{L} \times Q) \mid \langle q, d \rangle \in \delta\}$; denote by $\mathsf{moves} = \bigcup_{q \in Q} \mathsf{moves}(q)$ the set of all moves.

For the complexity results, we assume that probabilities of transitions are rational and given as fractions of integers represented in binary.

We describe the behaviour of an MDP as a trace generator running in steps. The MDP starts in the first step in state $q$ with probability $\mu_0(q)$. In each step, if the MDP is in state $q$ the controller chooses $\mathsf{m} \in \mathsf{moves}(q)$; then, with probability $\mathsf{m}(a, q')$, the label $a$ is generated and the next step starts in the successor state $q'$.

Given $q \in Q$, denote by $\mathsf{post}(q)$ the set $\{(a, q') \in \mathsf{Supp}(\mathsf{m}) \mid \mathsf{m} \in \mathsf{moves}(q)\}$. A *path* in $\mathcal{D}$ is a sequence $\rho = q_0 a_1 q_1 \ldots a_n q_n$ such that $(a_{i+1}, q_{i+1}) \in \mathsf{post}(q_i)$ for all $0 \le i < n$. The last state of $\rho$ is $\mathsf{last}(\rho) = q_n$. The trace $\mathsf{trace}(\rho)$ generated by $\rho$ is $a_1 a_2 \cdots a_n$. We let $\mathsf{Paths}(\mathcal{D})$ denote the set of paths in $\mathcal{D}$, $\mathsf{Paths}(w)$ denote $\{\rho \in \mathsf{Paths}(\mathcal{D}) \mid \mathsf{trace}(\rho) = w\}$ the set of paths generating $w$, and $\mathsf{Paths}(w, q)$ denote $\{\rho \in \mathsf{Paths}(\mathcal{D}) \mid \mathsf{trace}(\rho) = w \text{ and } \mathsf{last}(\rho) = q\}$ the set of paths generating $w$ ending in $q$.

**Strategies.** A *strategy* for an MDP $\mathcal{D}$ is a function $\alpha : \mathsf{Paths}(\mathcal{D}) \to \mathsf{Dist}(\mathsf{moves})$ that, given a path $\rho$, returns a probability distribution $\alpha(\rho) \in \mathsf{Dist}(\mathsf{moves}(\mathsf{last}(\rho)))$. Let $q = \mathsf{last}(\rho)$, then $\alpha(\rho)$ generates a label $a$ and selects a successor state $q'$ with probability

$$\sum_{\mathsf{m} \in \mathsf{moves}(q)} \alpha(\rho)(\mathsf{m}) \cdot \mathsf{m}(a, q').$$

Abusing notation slightly we write $\alpha(\rho)(a, q')$ for $\sum_{\mathsf{m} \in \mathsf{moves}} \alpha(\rho)(\mathsf{m}) \cdot \mathsf{m}(a, q')$.

A strategy $\alpha$ is *pure* if for all $\rho \in \mathsf{Paths}(\mathcal{D})$, there exists $\mathsf{m} \in \mathsf{moves}$ such that $\alpha(\rho)(\mathsf{m}) = 1$; we thus view pure strategies as functions $\alpha : \mathsf{Paths}(\mathcal{D}) \to \mathsf{moves}$. A strategy $\alpha$ is *memoryless* if $\alpha(\rho) = \alpha(\rho')$ for all paths $\rho, \rho'$ with $\mathsf{last}(\rho) = \mathsf{last}(\rho')$; we thus view memoryless strategies as functions $\alpha : Q \to \mathsf{Dist}(\mathsf{moves})$. A strategy $\alpha$ is *trace-based* if $\alpha(\rho) = \alpha(\rho')$ for all $\rho, \rho'$ where $\mathsf{trace}(\rho) = \mathsf{trace}(\rho')$ and $\mathsf{last}(\rho) = \mathsf{last}(\rho')$; we view trace-based strategies as functions $\alpha : \mathsf{L}^* \times Q \to \mathsf{Dist}(\mathsf{moves})$. For a traced-based strategy $\alpha$ we write $\alpha(w, q)(a, q')$ for $\sum_{\mathsf{m} \in \mathsf{moves}} \alpha(w, q)(\mathsf{m}) \cdot \mathsf{m}(a, q')$.

**Trace-probability function.** For an MDP $\mathcal{D}$ and a strategy $\alpha$, the probability of a single path is inductively defined by $\mathsf{Pr}_{\mathcal{D}, \alpha}(q) = \mu_0(q)$ and

$$\mathsf{Pr}_{\mathcal{D}, \alpha}(\rho a q') = \mathsf{Pr}_{\mathcal{D}, \alpha}(\rho) \cdot \alpha(\rho)(a, q').$$

This induces $\mathsf{Pr}_{\mathcal{D}, \alpha}(w, q) = \sum_{\rho \in \mathsf{Paths}(w, q)} \mathsf{Pr}_{\mathcal{D}, \alpha}(\rho)$ and $\mathsf{Pr}_{\mathcal{D}, \alpha}(w) = \sum_{\rho \in \mathsf{Paths}(w)} \mathsf{Pr}_{\mathcal{D}, \alpha}(\rho)$.

The *trace-probability* function $Tr_{\mathcal{D}, \alpha} : \mathsf{L}^* \to [0, 1]$ is, given a trace $w$, defined by

$$Tr_{\mathcal{D}, \alpha}(w) = \mathsf{Pr}_{\mathcal{D}, \alpha}(w).$$

We may drop the subscript $\mathcal{D}$ or $\alpha$ from $Tr_{\mathcal{D}, \alpha}$ and from $\mathsf{Pr}_{\mathcal{D}, \alpha}$ if it is understood. We let $subDist_{\mathcal{D}, \alpha}(w) \in \mathsf{subDist}(Q)$ denote the subdistribution after generating a trace $w$, that is

$$subDist_{\mathcal{D}, \alpha}(w)(q) = \mathsf{Pr}_{\mathcal{D}, \alpha}(w, q).$$

We have:
$$Tr_{\mathcal{D},\alpha}(w) = \|subDist_{\mathcal{D},\alpha}(w)\| \tag{2.1}$$

Let $\mathcal{D}$ be an MDP and $\alpha, \beta$ be two strategies; we denote by $Tr_\alpha = Tr_\beta$ when the equality $Tr_\alpha(w) = Tr_\beta(w)$ holds for all traces $w \in \mathsf{L}^*$. A version of the following lemma was proved in [7, Lemma 1]:

**Lemma 2.1.** *Let $\mathcal{D}$ be an MDP and $\alpha$ be a strategy. There exists a trace-based strategy $\beta$ such that $Tr_\alpha = Tr_\beta$.*

*Proof.* Let $\alpha$ be a strategy $\alpha : \mathsf{Paths} \to \mathsf{Dist}(\mathsf{moves})$ of the MDP $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$. We define a trace-based strategy $\beta : \mathsf{L}^* \times Q \to \mathsf{Dist}(\mathsf{moves})$ as follows: given a pair of state $q$ and trace $w$, let
$$\beta(w,q) = \frac{\sum_{\rho \in \mathsf{Paths}(w,q)} \mathsf{Pr}_\alpha(\rho) \cdot \alpha(\rho)}{\mathsf{Pr}_\alpha(w,q)} .$$
We prove by induction that $\mathsf{Pr}_\beta(w,q) = \mathsf{Pr}_\alpha(w,q)$. The induction base, $w = \epsilon$, is simple.

For any $w, q'$ we have:
$$\begin{aligned}
\mathsf{Pr}_\beta(wa, q') &= \sum_{q \in Q} \mathsf{Pr}_\beta(w,q) \cdot \beta(w,q)(a,q') \\
&= \sum_{q \in Q} \mathsf{Pr}_\alpha(w,q) \cdot \frac{\sum_{\rho \in \mathsf{Paths}(w,q)} \mathsf{Pr}_\alpha(\rho) \cdot \alpha(\rho)(a,q')}{\mathsf{Pr}_\alpha(w,q)} \\
&= \sum_{q \in Q} \sum_{\rho \in \mathsf{Paths}(w,q)} \mathsf{Pr}_\alpha(\rho) \cdot \alpha(\rho)(a,q') \\
&= \sum_{\rho \in \mathsf{Paths}(w)} \mathsf{Pr}_\alpha(\rho) \cdot \alpha(\rho)(a,q') \\
&= \mathsf{Pr}_\alpha(wa, q'). \qquad\qquad \square
\end{aligned}$$

A strategy can be implemented by means of memory; to make this explicit, we define a variant of the notion of strategies. We instrument strategies with a countable set $\mathsf{mem}$ of *memory modes*. For an MDP $\mathcal{D}$, a *generalized strategy* with memory $\mathsf{mem}$ is defined by $d_0 \in \mathsf{Dist}(\mathsf{mem})$ and $\alpha : \mathsf{mem} \times \mathsf{Paths}(\mathcal{D}) \to \mathsf{Dist}(\mathsf{moves} \times \mathsf{mem})$. The strategy $\alpha$ returns a probability distribution over the next moves and next memory modes based on the taken path to the current state and the current memory mode. We show in Lemma 2.2 that for each generalized strategy $\alpha$ there is a strategy $\beta$ such that each path in the MDP is equally probable under both strategies $\alpha$ and $\beta$, implying that the trace-probability functions induced by $\alpha$ and $\beta$ are also equal.

To formalize generalized strategies, we extend the definitions: An extended path is a sequence $\pi = q_0 M_0 a_1 q_1 M_1 \ldots a_n q_n M_n$ such that $(a_{i+1}, q_{i+1}) \in \mathsf{post}(q_i)$ for all $0 \le i < n$. The last memory state is $\mathsf{last}_{\mathsf{mem}}(\pi) = M_n$. All notions for paths are naturally transferred to extended paths. For an extended path $\pi = q_0 M_0 a_1 q_1 M_1 \ldots a_n q_n M_n$, let $\rho_\pi$ denote its projection to a path, that is $\rho_\pi = q_0 a_1 q_1 \ldots a_n q_n$. We let $\mathsf{ExtPaths}$ be the set of extended paths, $\mathsf{ExtPaths}(\rho)$ be the set of extended paths projecting to $\rho$, and $\mathsf{ExtPaths}(\rho, M)$ be the set of extended paths $\pi$ projecting to $\rho$ and such that $\mathsf{last}_{\mathsf{mem}}(\pi) = M$.

Given a memory mode $M$ and a path $\rho \in \mathsf{Paths}(\mathcal{D})$, we write $\alpha(M, \rho)(a, q', M')$ for $\sum_{\mathsf{m} \in \mathsf{moves}} \alpha(M, \rho)(\mathsf{m}, M') \cdot \mathsf{m}(a, q')$. The probability of an extended path $\pi$ is defined

inductively by $\mathsf{Pr}_{\mathcal{D},\alpha}(qm) = \mu_0(q) \cdot d_0(m)$ and

$$\mathsf{Pr}_{\mathcal{D},\alpha}(\pi a q' M') = \mathsf{Pr}_{\mathcal{D},\alpha}(\pi) \cdot \alpha(\mathsf{last}_{\mathsf{mem}}(\pi), \rho_\pi)(a, q', M').$$

We then let

$$\mathsf{Pr}_{\mathcal{D},\alpha}(\rho) = \sum_{\pi \in \mathsf{ExtPaths}(\rho)} \mathsf{Pr}_{\mathcal{D},\alpha}(\pi) \quad \text{and} \quad \mathsf{Pr}_{\mathcal{D},\alpha}(\rho, M) = \sum_{\pi \in \mathsf{ExtPaths}(\rho,M)} \mathsf{Pr}_{\mathcal{D},\alpha}(\pi).$$

It easily follows from these definitions that

$$\mathsf{Pr}_{\mathcal{D},\alpha}(\rho a q') = \sum_{M \in \mathsf{mem}} \mathsf{Pr}_{\mathcal{D},\alpha}(\rho, M) \cdot \alpha(M, \rho)(a, q').$$

**Lemma 2.2.** *Let $\mathcal{D}$ be an MDP and $\alpha$ be a generalized strategy. There exists a strategy $\beta$ such that $Tr_\alpha = Tr_\beta$.*

*Proof.* Let $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ be an MDP and $\alpha : \mathsf{mem} \times \mathsf{Paths}(\mathcal{D}) \to \mathsf{Dist}(\mathsf{moves} \times \mathsf{mem})$ be a generalized strategy of $\mathcal{D}$. We drop the subscript $\mathcal{D}$ in the rest of the proof.

Define a strategy $\beta : \mathsf{Paths} \to \mathsf{Dist}(\mathsf{moves})$ from $\alpha$ as follows. Given a path $\rho$, let

$$\beta(\rho) = \frac{\sum_{M \in \mathsf{mem}} \mathsf{Pr}_\alpha(\rho, M) \cdot \alpha(M, \rho)}{\mathsf{Pr}_\alpha(\rho)}.$$

We prove by induction that $\mathsf{Pr}_\beta(\rho) = \mathsf{Pr}_\alpha(\rho)$. The induction base, $\rho = q$ for $q \in Q$, is simple. For all $\rho, a, q'$ we have:

$$\begin{aligned}
\mathsf{Pr}_\beta(\rho a q') &= \mathsf{Pr}_\beta(\rho) \cdot \beta(\rho)(a, q') \\
&= \mathsf{Pr}_\alpha(\rho) \cdot \frac{\sum_{M \in \mathsf{mem}} \mathsf{Pr}_\alpha(\rho, M) \cdot \alpha(M, \rho)(a, q')}{\mathsf{Pr}_\alpha(\rho)} \\
&= \sum_{M \in \mathsf{mem}} \mathsf{Pr}_\alpha(\rho, M) \cdot \alpha(M, \rho)(a, q') \\
&= \mathsf{Pr}_\alpha(\rho a q').
\end{aligned}$$

Since the probability of a trace $w$ is a summation over the probability of all paths emitting $w$, having $\mathsf{Pr}_\alpha(\rho) = \mathsf{Pr}_\beta(\rho)$ for all $\rho$ implies that $Tr_\alpha = Tr_\beta$. $\qquad\square$

**Labelled Markov Chains.** A finite-state labelled Markov chain (MC for short) is an MDP where only a single move is available in each state, and thus controller's choice plays no role. An MC $\mathcal{C} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ is an MDP where $\delta : Q \to \mathsf{Dist}(\mathsf{L} \times Q)$ is a probabilistic transition function. Since MCs are MDPs, we analogously define paths, and the probability of a single path inductively as follows: $\mathsf{Pr}_{\mathcal{C}}(q) = \mu_0(q)$ and $\mathsf{Pr}_{\mathcal{C}}(\rho a q) = \mathsf{Pr}_{\mathcal{C}}(\rho) \cdot \delta(q')(a, q)$ where $q' = \mathsf{last}(\rho)$. The notations $subDist_{\mathcal{C}}(w)$ and $Tr_{\mathcal{C}}$ are defined analogously.

2.2. **Trace Refinement.** Given two MDPs $\mathcal{D}$ and $\mathcal{E}$ with the same set $\mathsf{L}$ of labels, we say that $\mathcal{E}$ *refines* $\mathcal{D}$, denoted by $\mathcal{D} \sqsubseteq \mathcal{E}$, if for all strategies $\alpha$ for $\mathcal{D}$ there exists some strategy $\beta$ for $\mathcal{E}$ such that $Tr_{\mathcal{D}} = Tr_{\mathcal{E}}$. We are interested in the problem $\mathsf{MDP} \sqsubseteq \mathsf{MDP}$, which asks, for two given MDPs $\mathcal{D}$ and $\mathcal{E}$, whether $\mathcal{D} \sqsubseteq \mathcal{E}$. The decidability of this problem was posed as an open question in [7]. We show in Theorem 3.1 that the problem $\mathsf{MDP} \sqsubseteq \mathsf{MDP}$ is undecidable.

We consider various subproblems of $\mathsf{MDP} \sqsubseteq \mathsf{MDP}$, which asks whether $\mathcal{D} \sqsubseteq \mathcal{E}$ holds. Specifically, we speak of the problem

- $\mathsf{MDP} \sqsubseteq \mathsf{MC}$ when $\mathcal{E}$ is restricted to be an MC;

- MC $\sqsubseteq$ MDP when $\mathcal{D}$ is restricted to be an MC;
- MC $\sqsubseteq$ MC when both $\mathcal{D}$ and $\mathcal{E}$ are restricted to be MCs.

We show in Theorem 3.1 that even the problem MC $\sqsubseteq$ MDP is undecidable. Hence we consider further subproblems. Specifically, we denote by MC $\sqsubseteq_m$ MDP the problem where the MDP is restricted to use only memoryless strategies, and by MC $\sqsubseteq_{pm}$ MDP the problem where the MDP is restricted to use only pure memoryless strategies. When both MDPs $\mathcal{D}$ and $\mathcal{E}$ are restricted to use only pure memoryless strategies, the trace-refinement problem is denoted by $\mathsf{MDP_{pm}} \sqsubseteq_{pm} \mathsf{MDP_{pm}}$. The problem MC $\sqsubseteq$ MC equals the *trace-equivalence problem* for MCs: given two MCs $\mathcal{C}_1, \mathcal{C}_2$ we have $\mathcal{C}_1 \sqsubseteq \mathcal{C}_2$ if and only if $Tr_{\mathcal{C}_1} = Tr_{\mathcal{C}_2}$ if and only if $\mathcal{C}_2 \sqsubseteq \mathcal{C}_1$. This problem is known to be in NC [24], hence in P.

## 3. Undecidability Results

In this section we show:

**Theorem 3.1.** *The problem* MC $\sqsubseteq$ MDP *is undecidable. Hence a fortiori,* MDP $\sqsubseteq$ MDP *is undecidable.*

*Proof.* To show that the problem MC $\sqsubseteq$ MDP is undecidable, we establish a reduction from the universality problem for probabilistic automata. A *probabilistic automaton* is a tuple $\mathcal{A} = \langle Q, \mu_0, \mathsf{L}, \delta, F \rangle$ consisting of a finite set $Q$ of states, an initial distribution $\mu_0 \in \mathsf{Dist}(Q)$, a finite set $\mathsf{L}$ of letters, a transition function $\delta : Q \times \mathsf{L} \to \mathsf{Dist}(Q)$ assigning to every state and letter a distribution over states, and a set $F$ of final states. For a word $w \in \mathsf{L}^*$ we write $dis_\mathcal{A}(w) \in \mathsf{Dist}(Q)$ for the distribution such that, for all $q \in Q$, we have that $dis_\mathcal{A}(w)(q)$ is the probability that after inputting $w$ the automaton $\mathcal{A}$ is in state $q$. We write $\mathsf{Pr}_\mathcal{A}(w) = \sum_{q \in F} dis_\mathcal{A}(w)(q)$ to denote the probability that $\mathcal{A}$ accepts $w$. The *universality problem* asks, given a probabilistic automaton $\mathcal{A}$, whether $\mathsf{Pr}_\mathcal{A}(w) \geq \frac{1}{2}$ holds for all words $w$. This problem is known to be undecidable [18, 9].

Let $\mathcal{A} = \langle Q, \mu_0, \mathsf{L}, \delta, F \rangle$ be a probabilistic automaton; without loss of generality we assume that $\mathsf{L} = \{a, b\}$. We construct an MDP $\mathcal{D}$ such that $\mathcal{A}$ is universal if and only if $\mathcal{C} \sqsubseteq \mathcal{D}$ where $\mathcal{C}$ is the MC shown in Figure 1.

The MDP $\mathcal{D}$ is constructed from $\mathcal{A}$ as follows; see Figure 3. Its set of states is $Q \cup \{q_c, q_d\}$, and its initial distribution is $\mu_0$. (Here and in the following we identify subdistributions $\mu \in \mathsf{subDist}(Q)$ and $\mu \in \mathsf{subDist}(Q \cup \{q_c, q_d\})$ if $\mu(q_c) = \mu(q_d) = 0$.) We describe the transitions of $\mathcal{D}$ using the transition function $\delta$ of $\mathcal{A}$. Consider a state $q \in Q$:

- If $q \in F$, there are two available moves $\mathsf{m}_c, \mathsf{m}_d$; both emit $a$ with probability $\frac{1}{4}$ and simulate the probabilistic automaton $\mathcal{A}$ reading the letter $a$, or emit $b$ with probability $\frac{1}{4}$ and simulate the probabilistic automaton $\mathcal{A}$ reading the letter $b$. With the remaining probability of $\frac{1}{2}$, $\mathsf{m}_c$ emits $c$ and leads to $q_c$ and $\mathsf{m}_d$ emits $d$ and leads to $q_d$. Formally, $\mathsf{m}_c(c, q_c) = \frac{1}{2}$, $\mathsf{m}_d(d, q_d) = \frac{1}{2}$ and $\mathsf{m}_c(e, q') = \mathsf{m}_d(e, q') = \frac{1}{4}\delta(q, e)(q')$ where $q' \in Q$ and $e \in \{a, b\}$.
- If $q \notin F$, there is a single available move $\mathsf{m}$ such that $\mathsf{m}(d, q_d) = \frac{1}{2}$ and $\mathsf{m}(e, q') = \frac{1}{4}\delta(q, e)(q')$ where $q' \in Q$ and $e \in \{a, b\}$.
- The only move from $q_c$ is the Dirac distribution on $(c, q_c)$; likewise the only move from $q_d$ is the Dirac distribution on $(d, q_d)$.

This MDP $\mathcal{D}$ "is almost" an MC, in the sense that a strategy $\alpha$ does not influence its behaviour until eventually a transition to $q_c$ or $q_d$ is taken. Indeed, for all $\alpha$ and for all
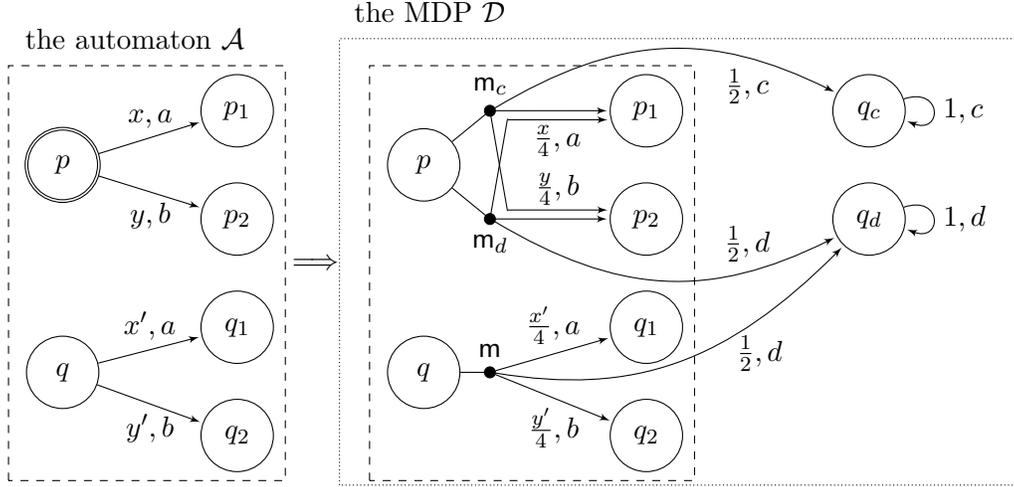
the automaton $\mathcal{A}$

the MDP $\mathcal{D}$



Figure 3: Sketch of the construction of the MDP $\mathcal{D}$ from the probabilistic automaton $\mathcal{A}$, for the undecidability result of $\mathsf{MC} \sqsubseteq \mathsf{MDP}$. Here, $p$ is an accepting state whereas $q$ is not. To read the picture, note that in $p$ there is a transition to the state $p_1$ with probability $x$ and label $a$: $\delta(p, a)(p_1) = x$.

$w \in \{a, b\}^*$ we have $subDist_{\mathcal{D},\alpha}(w) = \frac{1}{4^{|w|}} dis_{\mathcal{A}}(w)$. In particular, it follows $Tr_{\mathcal{D},\alpha}(w) = \|subDist_{\mathcal{D},\alpha}(w)\| = \frac{1}{4^{|w|}} \|dis_{\mathcal{A}}(w)\| = \frac{1}{4^{|w|}}$. Further, if $\alpha$ is trace-based we have:

$$
\begin{aligned}
Tr_{\mathcal{D},\alpha}(wc) &= \|subDist_{\mathcal{D},\alpha}(wc)\| & \text{by (2.1)} \\
&= subDist_{\mathcal{D},\alpha}(wc)(q_c) & \text{structure of } \mathcal{D} \\
&= \sum_{q \in F} subDist_{\mathcal{D},\alpha}(w)(q) \cdot \alpha(w,q)(\mathsf{m}_c) \cdot \frac{1}{2} & \text{structure of } \mathcal{D} \\
&= \frac{1}{4^{|w|}} \sum_{q \in F} dis_{\mathcal{A}}(w)(q) \cdot \alpha(w,q)(\mathsf{m}_c) \cdot \frac{1}{2} & \text{as argued above}
\end{aligned}
\tag{3.1}
$$

We show that $\mathcal{A}$ is universal if and only if $\mathcal{C} \sqsubseteq \mathcal{D}$. Let $\mathcal{A}$ be universal. Define a trace-based strategy $\alpha$ with $\alpha(w,q)(\mathsf{m}_c) = \frac{1}{2\mathsf{Pr}_{\mathcal{A}}(w)}$ for all $w \in \{a, b\}^*$ and $q \in F$. Note that $\alpha(w,q)(\mathsf{m}_c)$ is a probability as $\mathsf{Pr}_{\mathcal{A}}(w) \geq \frac{1}{2}$. Let $w \in \{a, b\}^*$. We have:

$$
\begin{aligned}
Tr_{\mathcal{D},\alpha}(w) &= \frac{1}{4^{|w|}} & \text{as argued above} \\
&= Tr_{\mathcal{C}}(w) & \text{Figure 1}
\end{aligned}
$$

Further we have:

$$Tr_{\mathcal{D},\alpha}(wc) = \frac{1}{4^{|w|}} \sum_{q \in F} dis_{\mathcal{A}}(w)(q) \cdot \alpha(w,q)(\mathsf{m}_c) \cdot \frac{1}{2} \qquad \text{by (3.1)}$$

$$= \frac{1}{4^{|w|}} \sum_{q \in F} dis_{\mathcal{A}}(w)(q) \cdot \frac{1}{\mathsf{Pr}_{\mathcal{A}}(w)} \cdot \frac{1}{4} \qquad \text{definition of } \alpha$$

$$= \frac{1}{4^{|w|+1}} \qquad\qquad \mathsf{Pr}_{\mathcal{A}}(w) = \sum_{q \in F} dis_{\mathcal{A}}(w)(q)$$

$$= Tr_{\mathcal{C}}(wc) \qquad\qquad \text{Figure 1}$$

It follows from the definitions of $\mathcal{D}$ and $\mathcal{C}$ that for all $k \geq 1$, we have $Tr_{\mathcal{D},\alpha}(wc^k) = Tr_{\mathcal{D},\alpha}(wc) = Tr_{\mathcal{C}}(wc) = Tr_{\mathcal{C}}(wc^k)$. We have $\sum_{e \in \{a,b,c,d\}} Tr_{\mathcal{D},\alpha}(we) = Tr_{\mathcal{D},\alpha}(w) = Tr_{\mathcal{C}}(w) = \sum_{e \in \{a,b,c,d\}} Tr_{\mathcal{C}}(we)$. Since for $e \in \{a,b,c\}$ we also proved that $Tr_{\mathcal{D},\alpha}(we) = Tr_{\mathcal{C}}(we)$ it follows that $Tr_{\mathcal{D},\alpha}(wd) = Tr_{\mathcal{C}}(wd)$. Hence, as above, $Tr_{\mathcal{D},\alpha}(wd^k) = Tr_{\mathcal{C}}(wd^k)$ for all $k \geq 1$. Finally, if $w \notin (a+b)^* \cdot (c^* + d^*)$ then $Tr_{\mathcal{D},\alpha}(w) = 0 = Tr_{\mathcal{C}}(w)$.

For the converse, assume that $\mathcal{A}$ is not universal. Then there is $w \in \{a,b\}^*$ with $\mathsf{Pr}_{\mathcal{A}}(w) < \frac{1}{2}$. Let $\alpha$ be a trace-based strategy. Then we have:

$$Tr_{\mathcal{D},\alpha}(wc) = \frac{1}{4^{|w|}} \sum_{q \in F} dis_{\mathcal{A}}(w)(q) \cdot \alpha(w,q)(\mathsf{m}_c) \cdot \frac{1}{2} \qquad \text{by (3.1)}$$

$$\leq \frac{1}{4^{|w|}} \cdot \frac{1}{2} \cdot \sum_{q \in F} dis_{\mathcal{A}}(w)(q) \qquad\qquad \alpha(w,q)(\mathsf{m}_c) \leq 1$$

$$= \frac{1}{4^{|w|}} \cdot \frac{1}{2} \cdot \mathsf{Pr}_{\mathcal{A}}(w) \qquad\qquad \mathsf{Pr}_{\mathcal{A}}(w) = \sum_{q \in F} dis_{\mathcal{A}}(w)(q)$$

$$< \frac{1}{4^{|w|}} \cdot \frac{1}{2} \cdot \frac{1}{2} \qquad\qquad \text{definition of } w$$

$$= Tr_{\mathcal{C}}(wc) \qquad\qquad \text{Figure 1}$$

We conclude that there is no trace-based strategy $\alpha$ with $Tr_{\mathcal{D},\alpha} = Tr_{\mathcal{C}}$. By Lemma 2.1 there is *no* strategy $\alpha$ with $Tr_{\mathcal{D},\alpha} = Tr_{\mathcal{C}}$. Hence $\mathcal{C} \not\sqsubseteq \mathcal{D}$. □

A straightforward reduction from $\mathsf{MDP} \sqsubseteq \mathsf{MDP}$ now establishes:

**Theorem 3.2.** *The problem that, given two MDPs $\mathcal{D}$ and $\mathcal{E}$, asks whether $\mathcal{D} \sqsubseteq \mathcal{E}$ and $\mathcal{E} \sqsubseteq \mathcal{D}$ is undecidable.*

*Proof.* We give a reduction from the problem $\mathsf{MDP} \sqsubseteq \mathsf{MDP}$. Given two MDPs $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ and $\mathcal{E} = \langle Q', \mu_0', \mathsf{L}, \delta' \rangle$, we construct two MDPs called $\mathcal{D} + \mathcal{E}$ and $\mathcal{E}_2$ such that $\mathcal{D} \sqsubseteq \mathcal{E}$ if, and only if, $\mathcal{E}_2 \sqsubseteq \mathcal{D} + \mathcal{E}$ and $\mathcal{D} + \mathcal{E} \sqsubseteq \mathcal{E}_2$. See Figure 4 for an illustration of the construction.

We first construct the MDP $\mathcal{D} + \mathcal{E}$ by simply having a copy of each MDP, adding a new label $\#$ and a new state $p_0$. The initial distribution of $\mathcal{D} + \mathcal{E}$ is the Dirac distribution on $p_0$, where there are two available moves $\mathsf{m}_{\mathcal{D}}$ and $\mathsf{m}_{\mathcal{E}}$. Let $\mathsf{m}_{\mathcal{D}}(\#, q) = \mu_0(q)$ for all $q \in Q$ and $\mathsf{m}_{\mathcal{D}}(\#, q) = 0$ otherwise; and let $\mathsf{m}_{\mathcal{E}}(\#, q) = \mu_0'(q)$ for all $q \in Q'$ and $\mathsf{m}_{\mathcal{D}}(\#, q) = 0$ otherwise. We see that $\mathcal{D} + \mathcal{E}$ always starts by generating label $\#$ with probability 1; next, using memory, $\mathcal{D} + \mathcal{E}$ can commit to either simulating $\mathcal{D}$ or $\mathcal{E}$.
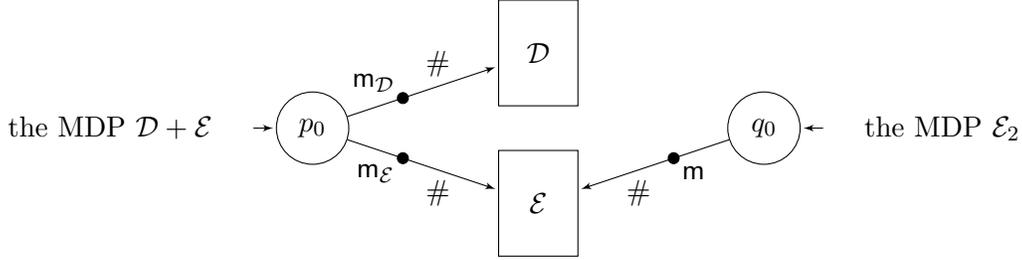
Figure 4: The construction of MDPs $\mathcal{D} + \mathcal{E}$ and $\mathcal{E}_2$ in the proof of Theorem 3.2.

We construct $\mathcal{E}_2$ from $\mathcal{E}$ as follows. We extend the set of labels with $\#$, and the set $Q'$ of states with a new state $q_0$. The initial distribution of $\mathcal{E}_2$ is the Dirac distribution on $q_0$, where there is only one available move $\mathsf{m}$ such that $\mathsf{m}(\#, q) = \mu'_0(q)$. We see that $\mathcal{E}_2$ always starts by generating label $\#$ with probability 1, and then simply behaves as $\mathcal{E}$.

We now argue that $\mathcal{D} \sqsubseteq \mathcal{E}$ if, and only if, $\mathcal{E}_2 \sqsubseteq \mathcal{D} + \mathcal{E}$ and $\mathcal{D} + \mathcal{E} \sqsubseteq \mathcal{E}_2$. This follows from three simple observations:

- The relation $\mathcal{E}_2 \sqsubseteq \mathcal{D} + \mathcal{E}$ always holds. Strategies of $\mathcal{D} + \mathcal{E}$ can choose to simulate $\mathcal{E}_2$ by playing $\mathsf{m}_\mathcal{E}$ with probability 1.
- If $\mathcal{D} \sqsubseteq \mathcal{E}$ then $\mathcal{D} + \mathcal{E} \sqsubseteq \mathcal{E}_2$: a strategy $\gamma$ of $\mathcal{D} + \mathcal{E}$, in the first step, plays $\mathsf{m}_\mathcal{D}$ and $\mathsf{m}_\mathcal{E}$ with probabilities $\gamma(p_0)(\mathsf{m}_\mathcal{D})$ and $\gamma(p_0)(\mathsf{m}_\mathcal{E})$. Next, it follows a strategy $\alpha$ for the copy of $\mathcal{D}$ and a strategy $\beta$ for the copy of $\mathcal{E}$. Since $\mathcal{D} \sqsubseteq \mathcal{E}$, there exists some strategy $\beta'$ of $\mathcal{E}$ such that $Tr_{\mathcal{D},\alpha} = Tr_{\mathcal{E},\beta'}$. Let $\gamma_2$ be the *generalized* (recall the definition preceding Lemma 2.2) strategy for $\mathcal{E}_2$ that first plays $\mathsf{m}$, and then plays $\beta'$ with probability $\gamma(p_0)(\mathsf{m}_\mathcal{D})$ and $\beta$ with probability $\gamma(p_0)(\mathsf{m}_\mathcal{E})$. Then $Tr_{\mathcal{D}+\mathcal{E},\gamma} = Tr_{\mathcal{E}_2,\gamma_2}$. By Lemma 2.2 there is a (non-generalized) strategy, $\gamma'_2$ such that $Tr_{\mathcal{E}_2,\gamma'_2} = Tr_{\mathcal{E}_2,\gamma_2} = Tr_{\mathcal{D}+\mathcal{E},\gamma}$. Thus $\mathcal{D} + \mathcal{E} \sqsubseteq \mathcal{E}_2$.
- If $\mathcal{D} + \mathcal{E} \sqsubseteq \mathcal{E}_2$ then $\mathcal{D} \sqsubseteq \mathcal{E}$: consider a strategy $\alpha$ of $\mathcal{D}$. Construct strategy $\alpha'$ of $\mathcal{D} + \mathcal{E}$ such that $\alpha'(p_0)(\mathsf{m}_\mathcal{D}) = 1$ and $\alpha'(p_0\#\rho) = \alpha(\rho)$ for all paths $\rho$. Since $\mathcal{D} + \mathcal{E} \sqsubseteq \mathcal{E}_2$, there must be some strategy $\beta'$ such that $Tr_{\mathcal{D}+\mathcal{E},\alpha'} = Tr_{\mathcal{E}_2,\beta'}$. For the strategy $\beta$ where $\beta(\rho) = \beta'(q_0\#\rho)$ for all paths $\rho$, we have $Tr_{\mathcal{D},\alpha} = Tr_{\mathcal{E},\beta}$, and thus $\mathcal{D} \sqsubseteq \mathcal{E}$. □

## 4. Decidability for Memoryless Strategies

Given two MCs $\mathcal{C}_1$ and $\mathcal{C}_2$, the (symmetric) trace-equivalence relation $\mathcal{C}_1 \sqsubseteq \mathcal{C}_2$ is polynomial-time decidable [24]. An MDP $\mathcal{D}$ under a memoryless strategy $\alpha$ induces a finite MC $\mathcal{D}(\alpha)$, and thus once a memoryless strategy is fixed for the MDP, its relation to another given MC in the trace-equivalence relation $\sqsubseteq$ can be decided in P. Theorems 4.1 and 4.2 provide tight complexity bounds of the trace-refinement problems for MDPs that are restricted to use only pure memoryless strategies. In Theorems 4.3 and 4.5 we establish bounds on the complexity of the problem when randomization is allowed for memoryless strategies.

4.1. **Pure Memoryless Strategies.** In this subsection, we show that the problems $\mathsf{MC} \sqsubseteq_{\mathsf{pm}} \mathsf{MDP}$ and $\mathsf{MDP} \sqsubseteq_{\mathsf{pm}} \mathsf{MDP}$ are NP-complete and $\Pi^p_2$-complete, respectively. The hardness results are by reductions from the *subset-sum problem* and a variant of the *quantified subset-sum* problem.
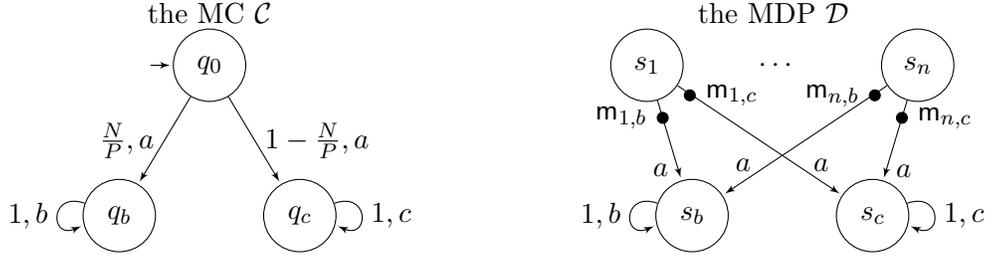
Figure 5: The MC $\mathcal{C}$ and the MDP $\mathcal{D}$ in the reduction for NP-hardness of MC $\sqsubseteq_{\mathsf{m}}$ MDP.

Given a set $\{s_1, s_2, \ldots, s_n\}$ of natural numbers and $N \in \mathbb{N}$, the subset-sum problem asks whether there exists a subset $S \subseteq \{s_1, \ldots, s_n\}$ such that $\sum_{s \in S} s = N$. The subset-sum problem is known to be NP-complete [6]. The quantified version of subset sum is a game between a *universal player* and an *existential player*. Given $k, N \in \mathbb{N}$ and two sets $\{s_1, s_2, \ldots, s_n\}$ and $\{t_1, t_2, \ldots, t_m\}$ of natural numbers, the game is played turn-based for $k$ rounds. In each round $i$ ($1 \leq i \leq k$), the universal player first chooses $S_i \subseteq \{s_1, \ldots, s_n\}$ and then the existential player chooses $T_i \subseteq \{t_1, \ldots, t_m\}$. The existential player wins if and only if

$$\sum_{s \in S_1} s + \sum_{t \in T_1} t + \cdots + \sum_{s \in S_k} s + \sum_{t \in T_k} t = N.$$

The *quantified subset-sum* problem is to check whether the existential player has a winning strategy. The problem is known to be PSPACE-complete [8]. The proof therein implies that the variant of the problem with a fixed number $k$ of rounds is $\Pi^p_{2k}$-complete.

**Theorem 4.1.** *The problem* MC $\sqsubseteq_{\mathsf{pm}}$ MDP *is* NP-*complete.*

*Proof.* Membership of MC $\sqsubseteq_{\mathsf{pm}}$ MDP in NP is obtained as follows. Given an MC $\mathcal{C}$ and an MDP $\mathcal{D}$, the polynomial-time verifiable witness of $\mathcal{C} \sqsubseteq \mathcal{D}$ is a pure memoryless strategy $\alpha$ for $\mathcal{D}$. Once $\alpha$ is fixed, then $\mathcal{C} \sqsubseteq \mathcal{D}(\alpha)$ can be decided in P.

To establish NP-hardness of MC $\sqsubseteq_{\mathsf{pm}}$ MDP, consider an instance of subset sum, i.e., a set $\{s_1, \ldots, s_n\}$ and $N \in \mathbb{N}$. We can assume without loss of generality that $N \leq P$, where $P = s_1 + \cdots + s_n$. We construct an MC $\mathcal{C}$ and an MDP $\mathcal{D}$ such that there exists $S \subseteq \{s_1, \ldots, s_n\}$ with $\sum_{s \in S} s = N$ if and only if $\mathcal{C} \sqsubseteq \mathcal{D}$ when $\mathcal{D}$ uses only pure memoryless strategies.

The MC $\mathcal{C}$ is shown in Figure 5 on the left. The initial distribution is the Dirac distribution on $q_0$; $\mathcal{C}$ generates traces in $ab^+$ with probability $\frac{N}{P}$ and traces in $ac^+$ with probability $1 - \frac{N}{P}$.

The MDP $\mathcal{D}$ is shown in Figure 5 on the right. For all states $s_i$, two moves $\mathsf{m}_{i,b}$ and $\mathsf{m}_{i,c}$ are available, the Dirac distributions on $(a, s_b)$ and $(a, s_c)$. The states $s_b, s_c$ emit only the single labels $b$ and $c$. The initial distribution $\mu_0$ is such that $\mu_0(s_i) = \frac{s_i}{P}$ for all $1 \leq i \leq n$. Intuitively, choosing $b$ in $s_i$ simulates the membership of $s_i$ in $S$ by adding $\frac{s_i}{P}$ to the probability of generating $ab^+$.

For a pure strategy $\alpha$ for $\mathcal{D}$, let $S_\alpha$ be the set of states $s_i$ where $\alpha(s_i) = m_{i,b}$. Then, $Tr_D(ab^+) = \sum_{s \in S_\alpha} \frac{s}{P}$ and $Tr_D(ac^+) = 1 - Tr_D(ab^+)$. Hence $\mathcal{C} \sqsubseteq \mathcal{D}$ holds if and only if there exists a strategy $\alpha$ for $\mathcal{D}$ such that $\sum_{s \in S_\alpha} \frac{s}{P} = \frac{N}{P}$. It implies that the instance of subset problem is positive, meaning that there exists a subset $S \subseteq \{s_1, s_2, \ldots, s_n\}$ such
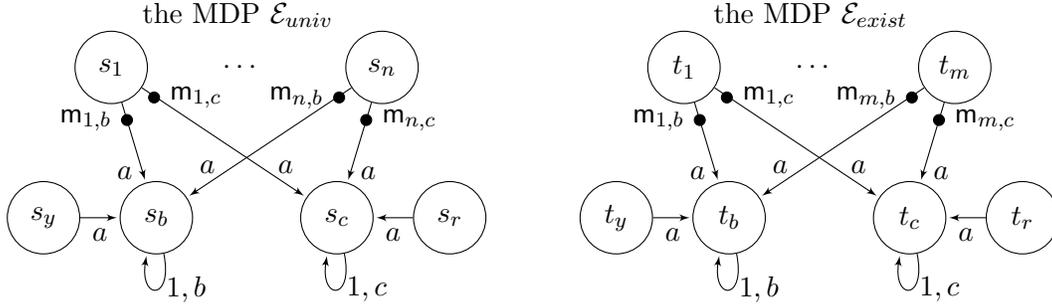
Figure 6: The MDPs $\mathcal{E}_{univ}$ and $\mathcal{E}_{exist}$ in the reduction for $\Pi_2^p$-hardness of MDP $\sqsubseteq_{pm}$ MDP.

that $\sum_{s \in S} s = N$, if and only if $\mathcal{C} \sqsubseteq \mathcal{D}$ when $\mathcal{D}$ uses only pure memoryless strategies. The NP-hardness results follows. □

In the next theorem, we show that MDP $\sqsubseteq_{pm}$ MDP is $\Pi_2^p$-complete. The hardness is by reduction from the quantified subset-sum problem with $k = 1$ (one alternation).

**Theorem 4.2.** *The problem* MDP $\sqsubseteq_{pm}$ MDP *is* $\Pi_2^p$*-complete.*

*Proof.* Membership of MDP $\sqsubseteq_{pm}$ MDP in $\Pi_2^p$ is obtained as follows. Let $\mathcal{D}$ and $\mathcal{E}$ be two MDPs. To check $\mathcal{E} \sqsubseteq \mathcal{D}$, for all pure memoryless strategies $\beta$ of $\mathcal{E}$ one can guess a polynomial-time verifiable witness $\alpha$, a strategy of $\mathcal{D}$. Once $\alpha$ and $\beta$ are fixed in $\mathcal{D}$ and $\mathcal{E}$ respectively, checking $\mathcal{E}(\beta) \sqsubseteq \mathcal{D}(\alpha)$ can be done in P.

To establish the hardness, consider an instance of quantified subset sum, i.e., $N \in \mathbb{N}$ and two sets $\{s_1, \ldots, s_n\}$ and $\{t_1, \ldots, t_m\}$. We construct MDPs $\mathcal{E}_{univ}$ and $\mathcal{E}_{exist}$ such that the existential player wins in one round if and only if $\mathcal{E}_{univ} \sqsubseteq \mathcal{E}_{exist}$ holds, where the MDPs use only pure memoryless strategies.

Let $P = s_1 + \cdots + s_n$ and $R = t_1 + \cdots + t_m$. Pick a small real number $0 < x < 1$ so that $0 < xP, xR, xN < 1$. Pick real numbers $0 \leq y_1, y_2 < 1$ such that $y_1 + xN < 1$ and $y_1 + xN = y_2 + xR$.

The MDPs $\mathcal{E}_{univ}$ and $\mathcal{E}_{exist}$ have symmetric constructions. The MDP $\mathcal{E}_{univ}$ simulates choices of the universal player and is drawn in Figure 6 on the left. For all states $s_i$, two moves $\mathsf{m}_{i,b}$ and $\mathsf{m}_{i,c}$ are available, the Dirac distributions on $(a, s_b)$ and $(a, s_c)$. The initial distribution $\mu_0$ for $\mathcal{E}_{univ}$ is such that $\mu_0(s_y) = \frac{1}{2}y_1$ and $\mu_0(s_r) = 1 - \frac{1}{2}(xP + y_1)$, and $\mu_0(s_i) = \frac{1}{2}xs_i$ for all $1 \leq i \leq n$. The MDP $\mathcal{E}_{exist}$ simulates choices of the existential player and is drawn in Figure 6 on the right. For all states $t_i$, two moves $\mathsf{m}_{i,b}$ and $\mathsf{m}_{i,c}$ are available, the Dirac distributions on $(a, t_b)$ and $(a, t_c)$, similar to $\mathcal{E}_{univ}$. The initial distribution $\mu_0'$ for $\mathcal{E}_{exist}$ is such that $\mu_0'(t_y) = \frac{1}{2}y_2$ and $\mu_0'(t_r) = 1 - \frac{1}{2}(xR + y_2)$, and $\mu_0'(t_j) = \frac{1}{2}xt_j$ for all $1 \leq j \leq m$. Choosing $b$ in a set of states $s_i$ by the universal player is responded by choosing $c$ in a right set of states $t_j$ by the existential player such that the probabilities of emitting $ab^+$ in the MDPs are equal.

For a pure strategy $\alpha$ of $\mathcal{E}_{univ}$, let $S_\alpha$ be the set of states $s_i$ where $\alpha(s_i) = \mathsf{m}_{i,b}$. We therefore have $Tr_{\mathcal{E}_{univ}}(ab^+) = \frac{1}{2}y_1 + \frac{1}{2}\sum_{s \in S_\alpha} xs$. For a pure strategy $\beta$ of $\mathcal{E}_{exist}$, let $T_\beta$ be the set of states $t_j$ where $\beta(t_j) = \mathsf{m}_{i,c}$. Then, $Tr_{\mathcal{E}_{exist}}(ac^+) = 1 - \frac{1}{2}(xR + y_2) + \frac{1}{2}\sum_{t \in T_\beta} xt$. It implies that $Tr_{\mathcal{E}_{exist}}(ab^+) = \frac{1}{2}(xR + y_2) - \frac{1}{2}\sum_{t \in T_\beta} xt$.

Since $y_1 + xN = y_2 + xR$, to achieve $Tr_{\mathcal{E}_{univ}} = Tr_{\mathcal{E}_{exist}}$ the equality $\sum_{s \in S_\alpha} s = N - \sum_{t \in T_\beta} t$ must be guaranteed. It shows that the existential player wins in one round, meaning that for all subsets $S \subseteq \{s_1, s_2, \ldots, s_n\}$ there exists a subset $T \subseteq \{t_1, t_2, \ldots, t_m\}$ such that $\sum_{s \in S} s + \sum_{t \in T} t = N$, if and only if for all pure and memoryless strategies $\alpha$ of $\mathcal{D}$ there exists some pure and memoryless strategy $\beta$ for $\mathcal{E}$ such that $Tr_{\mathcal{E}_{univ}} = Tr_{\mathcal{E}_{exist}}$. The $\Pi_2^p$-hardness result follows.      $\square$

4.2. **Memoryless Strategies.** In this subsection, we provide upper and lower complexity bounds for the problem $\mathsf{MC} \sqsubseteq_\mathsf{m} \mathsf{MDP}$: a reduction to the existential theory of the reals and a reduction from nonnegative matrix factorization.

A formula of the *existential theory of the reals* is of the form $\exists x_1 \ldots \exists x_m \, R(x_1, \ldots, x_n)$, where $R(x_1, \ldots, x_n)$ is a boolean combination of comparisons of the form $p(x_1, \ldots, x_n) \sim 0$, where $p(x_1, \ldots, x_n)$ is a multivariate polynomial and $\sim \in \{<, >, \leq, \geq, =, \neq\}$. The validity of closed formulas (i.e., when $m = n$) is decidable in $\mathsf{PSPACE}$ [3, 20], and is not known to be $\mathsf{PSPACE}$-hard.

**Theorem 4.3.** *The problem* $\mathsf{MC} \sqsubseteq_\mathsf{m} \mathsf{MDP}$ *is polynomial-time reducible to the existential theory of the reals, hence in* $\mathsf{PSPACE}$.

Given an MC $\mathcal{C} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$, to each label $a \in \mathsf{L}$ we associate a *transition matrix* $\Delta(a) \in [0,1]^{Q \times Q}$ with $\Delta(a)[q, q'] = \delta(q)(a, q')$. We view subdistributions $\mu_0$ over states as row vectors $\mu_0 \in [0,1]^Q$. We denote column vectors in boldface; in particular, $\mathbf{1} \in \{1\}^Q$ and $\mathbf{0} \in \{0\}^Q$ are column vectors all whose entries are 1 and 0, respectively. We build on [15, Proposition 10] which reads—translated to our framework—as follows:

**Proposition 4.4.** *Let* $\mathcal{C}_1 = \langle Q_1, \mu_0, \mathsf{L}, \delta \rangle$ *and* $\mathcal{C}_2 = \langle Q_2, \mu_0', \mathsf{L}, \delta' \rangle$ *be MCs with* $Q$ *as the disjoint union of* $Q_1, Q_2$. *Then* $Tr_{\mathcal{C}_1} = Tr_{\mathcal{C}_2}$ *if and only if there exists a matrix* $F \in \mathbb{R}^{Q \times Q}$ *such that*
- *the first row of* $F$ *equals* $(\mu_0, -\mu_0')$,
- $F\mathbf{1} = \mathbf{0}$,

*and, moreover, for all labels* $a \in \mathsf{L}$ *there exist matrices* $M(a) \in \mathbb{R}^{Q \times Q}$ *such that*

$$F \begin{pmatrix} \Delta(a) & 0 \\ 0 & \Delta'(a) \end{pmatrix} = M(a)F$$

*where* $\Delta(a), \Delta(a)'$ *are the transition matrices of* $\mathcal{C}_1$ *and* $\mathcal{C}_2$ *for the label* $a$.

With this at hand we prove Theorem 4.3:

*Proof of Theorem 4.3.* Let $\mathcal{C} = \langle Q_1, \mu_0, \mathsf{L}, \delta \rangle$ be an MC and $\mathcal{D} = \langle Q_2, \mu_0', \mathsf{L}, \delta' \rangle$ be an MDP with $Q$ as the disjoint union of $Q_1, Q_2$. A memoryless strategy $\alpha$ of $\mathcal{D}$ can be characterized by numbers $x_{q,\mathsf{m}} \in [0,1]$ where $q \in Q_2$ and $\mathsf{m} \in \mathsf{moves}(q)$, such that $x_{q,\mathsf{m}} = \alpha(q)(\mathsf{m})$. We have $\sum_{\mathsf{m} \in \mathsf{moves}(q)} x_{q,\mathsf{m}} = 1$ for all states $q$. We write $\overline{x}$ for the collection $(x_{q,\mathsf{m}})_{q \in Q_2, \, \mathsf{m} \in \mathsf{moves}(q)}$, and $\alpha(\overline{x})$ for the memoryless strategy characterized by $\overline{x}$. We have:

$$\mathcal{C} \sqsubseteq_\mathsf{m} \mathcal{D} \iff \exists \text{ memoryless strategy } \alpha : Tr_\mathcal{C} = Tr_{\mathcal{D}, \alpha} \qquad \text{definition}$$
$$\iff Cond \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{Proposition 4.4,}$$

where $Cond$ is the following condition:

There exist

- $x_{q,\mathsf{m}} \in [0,1]$ for all $q \in Q_2$ and all $\mathsf{m} \in \mathsf{moves}(q)$
- matrices $M(a) \in \mathbb{R}^{Q \times Q}$ for all labels $a \in \mathsf{L}$,
- a matrix $F \in \mathbb{R}^{Q \times Q}$

such that

- $\sum_{\mathsf{m} \in \mathsf{moves}(q)} x_{q,\mathsf{m}} = 1$ for all $q \in Q_2$,
- the first row of $F$ equals $(\mu_0, -\mu_0')$,
- $F\mathbf{1} = \mathbf{0}$,
- for all labels $a \in \mathsf{L}$,

$$F \begin{pmatrix} \Delta(a) & 0 \\ 0 & \Delta'(a) \end{pmatrix} = M(a)F$$

where $\Delta(a), \Delta'(a)$ are the transition matrices of $C$ and the finite MC $\mathcal{D}(\alpha(\overline{x}))$ induced by $\mathcal{D}$ under the strategy $\alpha(\overline{x})$.

This condition *Cond* is a closed formula in the existential theory of the reals. $\qquad\square$

Given a nonnegative matrix $M \in \mathbb{R}^{n \times m}$, a *nonnegative factorization* of $M$ with *inner dimension* $r$ is a decomposition of the form $M = A \cdot W$ where $A \in \mathbb{R}^{n \times r}$ and $W \in \mathbb{R}^{r \times m}$ are nonnegative matrices (see [4, 25, 1] for more details). The *NMF problem* asks, given a nonnegative matrix $M \in \mathbb{R}^{n \times m}$ and a number $r \in \mathbb{N}$, whether there exists a factorization $M = A \cdot W$ with nonnegative matrices $A \in \mathbb{R}^{n \times r}$ and $W \in \mathbb{R}^{r \times m}$. The NMF problem is known to be NP-hard [25].

**Theorem 4.5.** *The NMF problem is polynomial-time reducible to* $\mathsf{MC} \sqsubseteq_\mathsf{m} \mathsf{MDP}$*, hence* $\mathsf{MC} \sqsubseteq_\mathsf{m} \mathsf{MDP}$ *is* NP*-hard.*
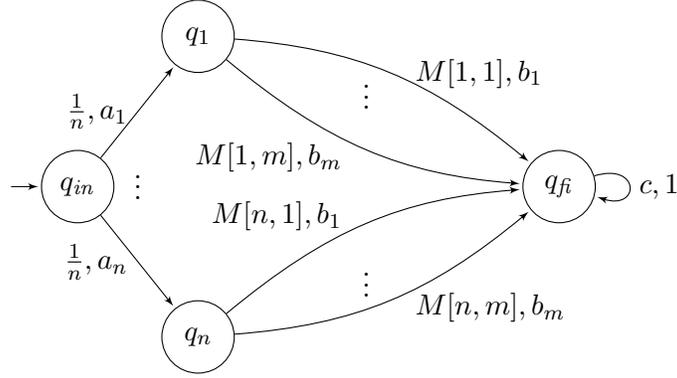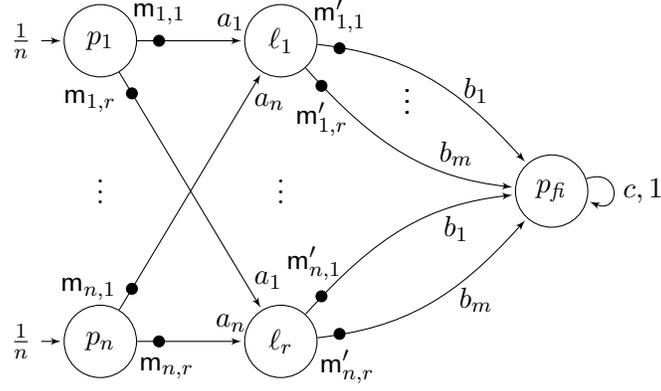
*Proof.* To establish the reduction, consider an instance of the NMF problem, a nonnegative matrix $M \in \mathbb{R}^{n \times m}$ and a number $r \in \mathbb{N}$. We construct an MC $\mathcal{C}$ and an MDP $\mathcal{D}$ such that the NMF instance is a yes-instance if and only if $\mathcal{C} \sqsubseteq \mathcal{D}$ where $\mathcal{D}$ is restricted to use only memoryless strategies.

We can assume, without loss of generality [4, Section 3], that $M$ is a stochastic matrix, that is $\sum_{j=1}^{m} M[i,j] = 1$ for all rows $1 \le i \le n$. Also by [4, Section 3] we know that there exists a nonnegative factorization of $M$ with inner dimension $r$ if and only if there exist two stochastic matrices $A \in \mathbb{R}^{n \times r}$ and $W \in \mathbb{R}^{r \times m}$ such that $M = A \cdot W$.

The transition probabilities in the MC $\mathcal{C}$ encode the entries of matrix $M$. The initial distribution of the MC is the Dirac distribution on $q_{in}$; see Figure 7. There are $n + m + 1$ labels $a_1, \dots, a_n, b_1, \dots, b_m, c$. The transition in $q_{in}$ is the uniform distribution over $\{(a_i, q_i) \mid 1 \le i \le n\}$. In each state $q_i$, each label $b_j$ is emitted with probability $M[i,j]$, and a transition to $q_{fi}$ is taken. In state $q_{fi}$ only $c$ is emitted. Observe that for all $1 \le i \le n$ and $1 \le j \le m$ we have $Tr_{\mathcal{C}}(a_i) = \frac{1}{n}$ and $Tr_{\mathcal{C}}(a_i \cdot b_j \cdot c^*) = \frac{1}{n} M[i,j]$.

The initial distribution of the MDP $\mathcal{D}$ is the uniform distribution over $\{p_1, \dots, p_n\}$; see Figure 8. In each $p_i$ (where $1 \le i \le n$), there are $r$ moves $\mathsf{m}_{i,1}, \mathsf{m}_{i,2}, \dots, \mathsf{m}_{i,r}$ where $\mathsf{m}_{i,k}(a_i, \ell_k) = 1$ and $1 \le k \le r$. In each $\ell_k$, there are $m$ moves $\mathsf{m}'_{k,1}, \mathsf{m}'_{k,2}, \dots, \mathsf{m}'_{k,m}$ where $\mathsf{m}'_{k,j}(b_j, p_{fi}) = 1$ where $1 \le j \le m$. In state $p_{fi}$, only $c$ is emitted. The probabilities of choosing the move $\mathsf{m}_{i,k}$ in $p_i$ and choosing $\mathsf{m}'_{k,j}$ in $\ell_k$ simulate the entries of $A[i,k]$ and $W[k,j]$.

We prove that there is a nonnegative factorization for $M = A \cdot W$ such that $A \in \mathbb{R}^{n \times r}$ and $W \in \mathbb{R}^{r \times m}$ if and only if $\mathcal{C} \sqsubseteq \mathcal{D}$ where $\mathcal{D}$ is restricted to memoryless strategies.

Figure 7: The MC $\mathcal{C}$ of the reduction from NMF to $\mathsf{MC} \sqsubseteq_{\mathsf{m}} \mathsf{MDP}$.



Figure 8: The MDP $\mathcal{D}$ of the reduction from NMF to $\mathsf{MC} \sqsubseteq_{\mathsf{m}} \mathsf{MDP}$.

Suppose $M$ has a nonnegative factorization, i.e., there are stochastic matrices $A \in \mathbb{R}^{n \times r}$ and $W \in \mathbb{R}^{r \times m}$ such that $M = A \cdot W$. To prove that $\mathcal{C} \sqsubseteq \mathcal{D}$, we construct a memoryless strategy $\alpha$ such that $Tr_{\mathcal{C}} = Tr_{\mathcal{D},\alpha}$. For all states $q$ of $\mathcal{D}$, strategy $\alpha$ is defined by

$$\alpha(q) = \begin{cases} d \in \mathsf{Dist}(\mathsf{moves}(p_i)) & \text{if } q = p_i \text{ and } 1 \leq i \leq n, \\ \text{where } d(\mathsf{m}_{i,k}) = A[i,k] \text{ for all } 1 \leq k \leq r \\[2mm] d \in \mathsf{Dist}(\mathsf{moves}(\ell_k)) & \text{if } q = \ell_k \text{ and } 1 \leq k \leq r, \\ \text{where } d(\mathsf{m}'_{k,j}) = W[k,j] \text{ for all } 1 \leq j \leq m \\[2mm] \text{the Dirac distribution on } (c, p_{fi}) & \text{if } q = p_{fi}. \end{cases}$$

The trace-probability function for $\mathcal{D}$ and $\alpha$ is such that for all $1 \leq i \leq n$ and all $1 \leq j \leq m$, we have $Tr_{\mathcal{D},\alpha}(a_i) = \frac{1}{n}$, and

$$Tr_{\mathcal{D},\alpha}(a_i \cdot b_j \cdot c^*) = \frac{1}{n} \sum_{k=1}^{r} \alpha(p_i)(\mathsf{m}_{i,k}) \cdot \alpha(\ell_k)(\mathsf{m}'_{k,j})$$

$$= \frac{1}{n} \sum_{k=1}^{r} A[i,k] \cdot W[k,j] = \frac{1}{n} M[i,j].$$

This gives $Tr_{\mathcal{D},\alpha} = Tr_{\mathcal{C}}$, and thus $\mathcal{C} \sqsubseteq \mathcal{D}$ where $\mathcal{D}$ uses a memoryless strategy.

Conversely, suppose that there exists a memoryless strategy $\beta$ for the MDP $\mathcal{D}$ such that $Tr_{\mathcal{C}} = Tr_{\mathcal{D},\beta}$. We exhibit a factorization $M = A \cdot W$ where $A \in \mathbb{R}^{n \times r}$ and $W \in \mathbb{R}^{r \times m}$. For all $1 \leq i \leq n$, $1 \leq k \leq r$ and $1 \leq j \leq m$, let

$$A[i,k] = \beta(p_i)(\mathsf{m}_{i,k}) \qquad \text{and} \qquad W[k,j] = \beta(\ell_k)(\mathsf{m}'_{k,j}).$$

Since $\mathcal{D}$ under the strategy $\beta$ refines $\mathcal{C}$, then for all $1 \leq i \leq n$ and all $1 \leq j \leq m$

$$Tr_{\mathcal{D},\beta}(a_i \cdot b_j \cdot c^*) = Tr_{\mathcal{C}}(a_i \cdot b_j \cdot c^*) = \frac{1}{n} M[i,j].$$

Since the probability of generating $a_i \cdot b_j \cdot c^*$ is $\frac{1}{n} \sum_{k=1}^{r} \beta(p_i)(\mathsf{m}_{i,k}) \cdot \beta(\ell_k)(\mathsf{m}'_{k,j})$ then we have

$$\sum_{k=1}^{r} A[i,k] \cdot W[k,j] = M[i,j].$$

This completes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

Shitov [22] recently claimed that the NMF problem is complete for the existential theory of the reals. This claim, combined with Theorems 4.3 and 4.5, implies that $\mathsf{MC} \sqsubseteq_{\mathsf{m}} \mathsf{MDP}$ is complete for the existential theory of the reals.

## 5. Bisimulation

In this section we show that the problem $\mathsf{MDP} \sqsubseteq \mathsf{MC}$ is in $\mathsf{NC}$, hence in $\mathsf{P}$.

First, in Subsection 5.1, we establish a link between trace refinement and a notion of bisimulation between distributions that was studied in [11].

Second, in Subsection 5.2 we give a necessary and sufficient condition for the MDPs to be bisimilar. It resembles the properties developed in [11], but we rebuild a detailed proof from scratch, as the authors were unable to verify some of the technical claims made in [11].

As corollaries, we show in Subsection 5.3 that bisimulation between two MDPs can be decided in $\mathsf{coNP}$, improving the exponential-time result from [11], and in Subsection 5.4 that the problem $\mathsf{MDP} \sqsubseteq \mathsf{MC}$ is in $\mathsf{NC}$, hence in $\mathsf{P}$.

5.1. **A Link between Trace Refinement and Bisimulation.** A *local strategy* for an MDP $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ is a function $\alpha : Q \to \mathsf{Dist}(\mathsf{moves})$ that maps each state $q$ to a distribution $\alpha(q) \in \mathsf{Dist}(\mathsf{moves}(q))$ over moves in $q$. We call $\alpha$ *pure* if for all states $q$ there is a move $\mathsf{m}$ such that $\alpha(q)(\mathsf{m}) = 1$. For a subdistribution $\mu \in \mathsf{subDist}(Q)$, a local strategy $\alpha$, and a label $a \in \mathsf{L}$, define the *successor* subdistribution $\mathsf{Succ}(\mu, \alpha, a)$ with

$$\mathsf{Succ}(\mu, \alpha, a)(q') = \sum_{q \in Q} \mu(q) \cdot \sum_{\mathsf{m} \in \mathsf{moves}(q)} \alpha(q)(\mathsf{m}) \cdot \mathsf{m}(a, q')$$

for all $q' \in Q$. We often view a subdistribution $d \in \mathsf{subDist}(Q)$ as a row vector $d \in [0, 1]^Q$. For a local strategy $\alpha$ and a label $a$, define the *transition matrix* $\Delta_\alpha(a) \in [0, 1]^{Q \times Q}$ with $\Delta_\alpha(a)[q, q'] = \sum_{\mathsf{m} \in \mathsf{moves}(q)} \alpha(q)(\mathsf{m}) \cdot \mathsf{m}(a, q')$. Viewing subdistributions $\mu$ as row vectors, we have:

$$\mathsf{Succ}(\mu, \alpha, a) = \mu \cdot \Delta_\alpha(a) \tag{5.1}$$

For a trace-based strategy $\alpha : \mathsf{L}^* \times Q \to \mathsf{Dist}(\mathsf{moves})$ and a trace $w \in \mathsf{L}^*$, define the local strategy $\alpha[w] : Q \to \mathsf{Dist}(\mathsf{moves})$ with $\alpha[w](q) = \alpha(w, q)$ for all $q \in Q$. We have the following lemma.

**Lemma 5.1.** *Let* $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ *be an MDP. Let* $\alpha : \mathsf{L}^* \times Q \to \mathsf{Dist}(\mathsf{moves})$ *be a trace-based strategy. Let* $w \in \mathsf{L}^*$ *and* $a \in \mathsf{L}$*. Then:*

$$subDist_{\mathcal{D},\alpha}(wa) = subDist_{\mathcal{D},\alpha}(w) \cdot \Delta_{\alpha[w]}(a)$$

*Proof.* Let $q' \in Q$. We have:

$$subDist_{\mathcal{D},\alpha}(wa)(q')$$

$$= \sum_{\rho \in \mathsf{Paths}(w)} \mathsf{Pr}_{\mathcal{D},\alpha}(\rho a q') \qquad\qquad \text{definition of } subDist$$

$$= \sum_{q \in Q} \sum_{\rho \in \mathsf{Paths}(w,q)} \mathsf{Pr}_{\mathcal{D},\alpha}(\rho) \cdot \sum_{\mathsf{m} \in \mathsf{moves}(q)} \alpha(\rho)(\mathsf{m}) \cdot \mathsf{m}(a, q') \qquad \text{definition of } \mathsf{Pr}$$

$$= \sum_{q \in Q} \sum_{\rho \in \mathsf{Paths}(w,q)} \mathsf{Pr}_{\mathcal{D},\alpha}(\rho) \cdot \sum_{\mathsf{m} \in \mathsf{moves}(q)} \alpha(w, q)(\mathsf{m}) \cdot \mathsf{m}(a, q') \qquad \alpha \text{ is trace-based}$$

$$= \sum_{q \in Q} subDist_{\mathcal{D},\alpha}(w)(q) \cdot \sum_{\mathsf{m} \in \mathsf{moves}(q)} \alpha(w, q)(\mathsf{m}) \cdot \mathsf{m}(a, q') \qquad \text{definition of } subDist$$

$$= \sum_{q \in Q} subDist_{\mathcal{D},\alpha}(w)(q) \cdot \sum_{\mathsf{m} \in \mathsf{moves}(q)} \alpha[w](q)(\mathsf{m}) \cdot \mathsf{m}(a, q') \qquad \text{definition of } \alpha[w]$$

$$= \sum_{q \in Q} subDist_{\mathcal{D},\alpha}(w)(q) \cdot \Delta_{\alpha[w]}(a)[q, q'] \qquad\qquad \text{definition of } \Delta_{\alpha[w]}(a)$$

$$= \left( subDist_{\mathcal{D},\alpha}(w) \cdot \Delta_{\alpha[w]}(a) \right)(q') \qquad\qquad\qquad \square$$

The following lemma is based on the idea that, using Lemma 5.1, we can "slice" a strategy into local strategies, and conversely we can compose local strategies to a strategy.

**Lemma 5.2.** *Let* $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ *be an MDP. Let* $w = a_1 a_2 \cdots a_n \in \mathsf{L}^*$*. Let* $\mu_1, \mu_2, \ldots, \mu_n$ *be subdistributions over* $Q$*. Then there is a strategy* $\alpha : \mathsf{Paths}(\mathcal{D}) \to \mathsf{Dist}(\mathsf{moves})$ *with*

$$\mu_i = subDist_{\mathcal{D},\alpha}(a_1 a_2 \cdots a_i) \qquad \text{for all } i \in \{0, 1, \ldots, n\}$$

*if and only if there are local strategies* $\alpha_0, \alpha_1, \ldots, \alpha_{n-1} : Q \to \mathsf{Dist}(\text{moves})$ *with*

$$\mu_{i+1} = \mathsf{Succ}(\mu_i, \alpha_i, a_{i+1}) \qquad \text{for all } i \in \{0, 1, \ldots, n-1\}.$$

*Proof.* We prove the two implications from the lemma in turn.

"$\Longrightarrow$": Let $\alpha$ be a strategy with $\mu_i = subDist_{\mathcal{D},\alpha}(a_1 a_2 \cdots a_i)$ for all $i \in \{0, 1, \ldots, n\}$. By Lemma 2.1 we can assume that $\alpha$ is trace-based. For all $i \in \{0, 1, \ldots, n-1\}$ define a local strategy $\alpha_i$ with $\alpha_i = \alpha[a_1 a_2 \cdots a_i]$. Then we have for all $i \in \{0, 1, \ldots, n-1\}$:

$$\begin{aligned}
\mu_{i+1} &= subDist_{\mathcal{D},\alpha}(a_1 a_2 \cdots a_{i+1}) && \text{definition of } \alpha \\
&= subDist_{\mathcal{D},\alpha}(a_1 a_2 \cdots a_i) \cdot \Delta_{\alpha[a_1 a_2 \cdots a_i]}(a_{i+1}) && \text{Lemma 5.1} \\
&= \mu_i \cdot \Delta_{\alpha_i}(a_{i+1}) && \text{definitions of } \alpha, \alpha_i \\
&= \mathsf{Succ}(\mu_i, \alpha_i, a_{i+1}) && \text{by (5.1)}
\end{aligned}$$

"$\Longleftarrow$": Let $\alpha_0, \alpha_1, \ldots, \alpha_{n-1}$ be local strategies with $\mu_{i+1} = \mathsf{Succ}(\mu_i, \alpha_i, a_{i+1})$ for all $i \in \{0, 1, \ldots, n-1\}$. Define a trace-based strategy $\alpha$ such that $\alpha[a_1 a_2 \cdots a_i] = \alpha_i$ for all $i \in \{0, 1, \ldots, n-1\}$. (This condition need not completely determine $\alpha$.) We prove by induction on $i$ that $\mu_i = subDist_{\mathcal{D},\alpha}(a_1 a_2 \cdots a_i)$ for all $i \in \{0, 1, \ldots, n\}$. For $i = 0$ this is trivial. For the step, we have:

$$\begin{aligned}
\mu_{i+1} &= \mathsf{Succ}(\mu_i, \alpha_i, a_{i+1}) && \text{definition of } \alpha_i \\
&= \mu_i \cdot \Delta_{\alpha_i}(a_{i+1}) && \text{by (5.1)} \\
&= subDist_{\mathcal{D},\alpha}(a_1 a_2 \cdots a_i) \cdot \Delta_{\alpha_i}(a_{i+1}) && \text{induction hypothesis} \\
&= subDist_{\mathcal{D},\alpha}(a_1 a_2 \cdots a_i) \cdot \Delta_{\alpha[a_1 a_2 \cdots a_i]}(a_{i+1}) && \text{definition of } \alpha \\
&= subDist_{\mathcal{D},\alpha}(a_1 a_2 \cdots a_{i+1}) && \text{Lemma 5.1} \qquad \square
\end{aligned}$$

Let $\mathcal{D} = \langle Q_{\mathcal{D}}, \mu_0^{\mathcal{D}}, \mathsf{L}, \delta_{\mathcal{D}} \rangle$ and $\mathcal{E} = \langle Q_{\mathcal{E}}, \mu_0^{\mathcal{E}}, \mathsf{L}, \delta_{\mathcal{E}} \rangle$ be two MDPs over the same set $\mathsf{L}$ of labels. A *bisimulation* is a relation $\mathcal{R} \subseteq \mathsf{subDist}(Q_{\mathcal{D}}) \times \mathsf{subDist}(Q_{\mathcal{E}})$ such that whenever $\mu_{\mathcal{D}} \ \mathcal{R} \ \mu_{\mathcal{E}}$ then

- $\|\mu_{\mathcal{D}}\| = \|\mu_{\mathcal{E}}\|$;
- for all local strategies $\alpha_{\mathcal{D}}$ there exists a local strategy $\alpha_{\mathcal{E}}$ such that for all $a \in \mathsf{L}$ we have $\mathsf{Succ}(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}}, a) \ \mathcal{R} \ \mathsf{Succ}(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}}, a)$;
- for all local strategies $\alpha_{\mathcal{E}}$ there exists a local strategy $\alpha_{\mathcal{D}}$ such that for all $a \in \mathsf{L}$ we have $\mathsf{Succ}(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}}, a) \ \mathcal{R} \ \mathsf{Succ}(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}}, a)$.

As usual, a union of bisimulations is a bisimulation. Denote by $\sim$ the union of all bisimulations, i.e., $\sim$ is the largest bisimulation. We write $\mathcal{D} \sim \mathcal{E}$ if $\mu_0^{\mathcal{D}} \sim \mu_0^{\mathcal{E}}$. In general, the set $\sim$ is uncountably infinite, so methods for computing state-based bisimulation (e.g., partition refinement) are not applicable.

Proposition 5.3 below establishes a link between trace refinement and bisimulation. An intuitive interpretation of the proposition is that if $\mathcal{D}$ is an MDP and $\mathcal{C}$ an MC, then the best way of disproving bisimilarity between $\mathcal{D}$ and $\mathcal{C}$ is to exhibit a sequence of local strategies in $\mathcal{D}$ so that the resulting behaviour of $\mathcal{D}$ cannot be matched by $\mathcal{C}$. Using Lemma 5.2 this sequence of local strategies can be assembled to a strategy for $\mathcal{D}$, which then witnesses that $\mathcal{D} \not\sqsubseteq \mathcal{C}$.

**Proposition 5.3.** *Let $\mathcal{D}$ be an MDP and $\mathcal{C}$ be an MC. Then $\mathcal{D} \sim \mathcal{C}$ if and only if $\mathcal{D} \sqsubseteq \mathcal{C}$.*

*Proof.* Let $\mathcal{D} = \langle Q_{\mathcal{D}}, \mu_0^{\mathcal{D}}, \mathsf{L}, \delta^{\mathcal{D}} \rangle$ and $\mathcal{C} = \langle Q_{\mathcal{C}}, \mu_0^{\mathcal{C}}, \mathsf{L}, \delta^{\mathcal{C}} \rangle$.

"$\Longrightarrow$": Let $\mathcal{D} \sim \mathcal{C}$. Hence $\mu_0^{\mathcal{D}} \sim \mu_0^{\mathcal{C}}$. We show that $\mathcal{D} \sqsubseteq \mathcal{C}$. Let $\alpha^{\mathcal{D}}$ be a strategy for $\mathcal{D}$. Let $w = a_1 a_2 \cdots a_n \in \mathsf{L}^*$. Let $\mu_0^{\mathcal{D}}, \mu_1^{\mathcal{D}}, \ldots, \mu_n^{\mathcal{D}}$ be the subdistributions with $\mu_i^{\mathcal{D}} = subDist_{\mathcal{D},\alpha^{\mathcal{D}}}(a_1 a_2 \cdots a_i)$ for all $i$. By Lemma 5.2 there exist local strategies $\alpha_0^{\mathcal{D}}, \alpha_1^{\mathcal{D}}, \ldots, \alpha_{n-1}^{\mathcal{D}}$ with $\mu_{i+1}^{\mathcal{D}} = \mathsf{Succ}(\mu_i^{\mathcal{D}}, \alpha_i^{\mathcal{D}}, a_{i+1})$ for all $i$. Since $\mu_0^{\mathcal{D}} \sim \mu_0^{\mathcal{C}}$, there exist local strategies $\alpha_0^{\mathcal{C}}, \alpha_1^{\mathcal{C}}, \ldots, \alpha_{n-1}^{\mathcal{C}}$ for $\mathcal{C}$ and subdistributions $\mu_1^{\mathcal{C}}, \mu_2^{\mathcal{C}}, \ldots, \mu_n^{\mathcal{C}}$ with $\mu_{i+1}^{\mathcal{C}} = \mathsf{Succ}(\mu_i^{\mathcal{C}}, \alpha_i^{\mathcal{C}}, a_{i+1})$ for all $i$ and $\mu_i^{\mathcal{D}} \sim \mu_i^{\mathcal{C}}$ for all $i$. Since $\mathcal{C}$ is an MC, the local strategies $\alpha_i^{\mathcal{C}}$ are, in fact, irrelevant. By Lemma 5.2 we have $\mu_i^{\mathcal{C}} = subDist_{\mathcal{C}}(a_1 a_2 \cdots a_i)$ for all $i$. So we have:

$$
\begin{aligned}
Tr_{\mathcal{D},\alpha^{\mathcal{D}}}(w) &= \|subDist_{\mathcal{D},\alpha^{\mathcal{D}}}(w)\| && \text{by (2.1)} \\
&= \|\mu_n^{\mathcal{D}}\| && \mu_n^{\mathcal{D}} = subDist_{\mathcal{D},\alpha^{\mathcal{D}}}(w) \\
&= \|\mu_n^{\mathcal{C}}\| && \mu_n^{\mathcal{D}} \sim \mu_n^{\mathcal{C}} \\
&= \|subDist_{\mathcal{C}}(w)\| && \mu_n^{\mathcal{C}} = subDist_{\mathcal{C}}(w) \\
&= Tr_{\mathcal{C}}(w) && \text{by (2.1)}
\end{aligned}
$$

Since $\alpha^{\mathcal{D}}$ and $w$ were chosen arbitrarily, we conclude that $\mathcal{D} \sqsubseteq \mathcal{C}$.

"$\Longleftarrow$": Let $\mathcal{D} \sqsubseteq \mathcal{C}$. We show $\mu_0^{\mathcal{D}} \sim \mu_0^{\mathcal{C}}$. Define a relation $\mathcal{R} \subseteq \mathsf{subDist}(Q_{\mathcal{D}}) \times \mathsf{subDist}(Q_{\mathcal{C}})$ such that $\mu^{\mathcal{D}} \mathcal{R} \mu^{\mathcal{C}}$ if and only if there exist a strategy $\alpha^{\mathcal{D}}$ for $\mathcal{D}$ and a trace $w$ with $\mu^{\mathcal{D}} = subDist_{\mathcal{D},\alpha^{\mathcal{D}}}(w)$ and $\mu^{\mathcal{C}} = subDist_{\mathcal{C}}(w)$. We claim that $\mathcal{R}$ is a bisimulation. To prove the claim, consider any $\mu^{\mathcal{D}}, \mu^{\mathcal{C}}$ with $\mu^{\mathcal{D}} \mathcal{R} \mu^{\mathcal{C}}$. Then there exist a strategy $\alpha^{\mathcal{D}}$ for $\mathcal{D}$ and a trace $w$ with $\mu^{\mathcal{D}} = subDist_{\mathcal{D},\alpha^{\mathcal{D}}}(w)$ and $\mu^{\mathcal{C}} = subDist_{\mathcal{C}}(w)$. Since $\mathcal{D} \sqsubseteq \mathcal{C}$, we have $Tr_{\mathcal{D},\alpha^{\mathcal{D}}}(w) = Tr_{\mathcal{C}}(w)$. So we have:

$$
\begin{aligned}
\|\mu^{\mathcal{D}}\| &= \|subDist_{\mathcal{D},\alpha^{\mathcal{D}}}(w)\| && \mu^{\mathcal{D}} = subDist_{\mathcal{D},\alpha^{\mathcal{D}}}(w) \\
&= Tr_{\mathcal{D},\alpha^{\mathcal{D}}}(w) && \text{by (2.1)} \\
&= Tr_{\mathcal{C}}(w) && \text{as argued above} \\
&= \|subDist_{\mathcal{C}}(w)\| && \text{by (2.1)} \\
&= \|\mu^{\mathcal{C}}\| && \mu^{\mathcal{C}} = subDist_{\mathcal{C}}(w)
\end{aligned}
$$

This proves the first condition for $\mathcal{R}$ being a bisimulation.

For the rest of the proof assume $w = a_1 a_2 \cdots a_n$. Write $\mu_n^{\mathcal{D}} = \mu^{\mathcal{D}}$ and $\mu_n^{\mathcal{C}} = \mu^{\mathcal{C}}$. Let $\alpha_n^{\mathcal{D}}$ be a local strategy for $\mathcal{D}$. Let $a_{n+1} \in \mathsf{L}$. Define $\mu_{n+1}^{\mathcal{D}} = \mathsf{Succ}(\mu_n^{\mathcal{D}}, \alpha_n^{\mathcal{D}}, a_{n+1})$, and $\mu_{n+1}^{\mathcal{C}} = \mathsf{Succ}(\mu_n^{\mathcal{C}}, \alpha_n^{\mathcal{C}}, a_{n+1})$ for an arbitrary (and unimportant as $\mathcal{C}$ is an MC) local strategy $\alpha_n^{\mathcal{C}}$ for $\mathcal{C}$. For the second and the third condition of $\mathcal{R}$ being a bisimulation we need to prove $\mu_{n+1}^{\mathcal{D}} \mathcal{R} \mu_{n+1}^{\mathcal{C}}$. Define $\mu_1^{\mathcal{D}}, \mu_2^{\mathcal{D}}, \ldots, \mu_{n-1}^{\mathcal{D}}$ such that $\mu_i^{\mathcal{D}} = subDist_{\mathcal{D},\alpha^{\mathcal{D}}}(a_1 a_2 \cdots a_i)$ for all $i \in \{0, 1, \ldots, n\}$. By Lemma 5.2 there are local strategies $\alpha_0^{\mathcal{D}}, \alpha_1^{\mathcal{D}}, \ldots, \alpha_{n-1}^{\mathcal{D}}$ such that $\mu_{i+1}^{\mathcal{D}} = \mathsf{Succ}(\mu_i^{\mathcal{D}}, \alpha_i^{\mathcal{D}}, a_{i+1})$ for all $i \in \{0, 1, \ldots, n-1\}$. We also have $\mu_{n+1}^{\mathcal{D}} = \mathsf{Succ}(\mu_n^{\mathcal{D}}, \alpha_n^{\mathcal{D}}, a_{n+1})$, so again by Lemma 5.2 there is a strategy $\beta^{\mathcal{D}}$ with $\mu_i^{\mathcal{D}} = subDist_{\mathcal{D},\beta^{\mathcal{D}}}(a_1 a_2 \cdots a_i)$ for all $i \in \{0, 1, \ldots, n+1\}$. In particular, $\mu_{n+1}^{\mathcal{D}} = subDist_{\mathcal{D},\beta^{\mathcal{D}}}(w a_{n+1})$. Similarly, we have $\mu_{n+1}^{\mathcal{C}} = subDist_{\mathcal{C}}(w a_{n+1})$. Thus, $\mu_{n+1}^{\mathcal{D}} \mathcal{R} \mu_{n+1}^{\mathcal{C}}$. Hence we have proved that $\mathcal{R}$ is a bisimulation.

Considering the empty trace, we see that $\mu_0^{\mathcal{D}} \mathcal{R} \mu_0^{\mathcal{C}}$. Since $\mathcal{R} \subseteq \sim$, we also have $\mu_0^{\mathcal{D}} \sim \mu_0^{\mathcal{C}}$, as desired. $\qquad\square$

5.2. **A Necessary and Sufficient Condition for Bisimilarity.** In the following we consider MDPs $\mathcal{D} = \langle Q, \mu_0^{\mathcal{D}}, \mathsf{L}, \delta \rangle$ and $\mathcal{E} = \langle Q, \mu_0^{\mathcal{E}}, \mathsf{L}, \delta \rangle$ over the same state space. This is without loss of generality, since we might take the disjoint union of the state spaces. Since $\mathcal{D}$ and $\mathcal{E}$ differ only in the initial distribution, we will focus on $\mathcal{D}$.

Let $B \in \mathbb{R}^{Q \times k}$ with $k \geq 1$. Assume the label set is $\mathsf{L} = \{a_1, \ldots, a_{|\mathsf{L}|}\}$. For $\mu \in \mathsf{subDist}(Q)$ and a local strategy $\alpha$ we define a point $p(\mu, \alpha) \in \mathbb{R}^{|\mathsf{L}| \cdot k}$ such that

$$p(\mu, \alpha) = \begin{pmatrix} \mu \Delta_\alpha(a_1)B & \mu \Delta_\alpha(a_2)B & \cdots & \mu \Delta_\alpha(a_{|\mathsf{L}|})B \end{pmatrix}.$$

For the reader's intuition, we remark that we will choose matrices $B \in \mathbb{R}^{Q \times k}$ so that if two subdistributions $\mu_{\mathcal{D}}, \mu_{\mathcal{E}}$ are bisimilar then $\mu_{\mathcal{D}}B = \mu_{\mathcal{E}}B$. (In fact, one can compute $B$ so that the converse holds as well, i.e., $\mu_{\mathcal{D}} \sim \mu_{\mathcal{E}}$ if and only if $\mu_{\mathcal{D}}B = \mu_{\mathcal{E}}B$.) It follows that, for subdistributions $\mu_{\mathcal{D}}, \mu_{\mathcal{E}}$ and local strategies $\alpha_{\mathcal{D}}, \alpha_{\mathcal{E}}$, if $\mathsf{Succ}(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}}, a) \sim \mathsf{Succ}(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}}, a)$ holds for all $a \in \mathsf{L}$ then $p(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}}) = p(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}})$. Let us also remark that for fixed $\mu \in \mathsf{subDist}(Q)$, the set $P_\mu = \{p(\mu, \alpha) \mid \alpha \text{ is a local strategy}\} \subseteq \mathbb{R}^{|\mathsf{L}| \cdot k}$ is a (bounded and convex) polytope. As a consequence, if $\mu_{\mathcal{D}} \sim \mu_{\mathcal{E}}$ then the polytopes $P_{\mu_{\mathcal{D}}}$ and $P_{\mu_{\mathcal{E}}}$ must be equal. In the next paragraph we define *"extremal"* strategies $\widehat{\alpha}$, which intuitively are local strategies such that $p(\mu, \widehat{\alpha})$ is a vertex of the polytope $P_\mu$.

Let $\boldsymbol{v} \in \mathbb{R}^{|\mathsf{L}| \cdot k}$ be a *column* vector; we denote column vectors in boldface. We view $\boldsymbol{v}$ as a "direction". Recall that $d_q$ is the Dirac distribution on the state $q$. A pure local strategy $\widehat{\alpha}$ is *extremal in direction $\boldsymbol{v}$ with respect to $B$* if

$$p(d_q, \alpha)\boldsymbol{v} \leq p(d_q, \widehat{\alpha})\boldsymbol{v} \tag{5.2}$$

$$p(d_q, \alpha)\boldsymbol{v} = p(d_q, \widehat{\alpha})\boldsymbol{v} \quad \text{implies} \quad p(d_q, \alpha) = p(d_q, \widehat{\alpha}) \tag{5.3}$$

for all states $q \in Q$ and all pure local strategies $\alpha$.

By linearity, if (5.2) and (5.3) hold for all pure local strategies $\alpha$ then (5.2) and (5.3) hold for all local strategies $\alpha$. We say a local strategy $\widehat{\alpha}$ is *extremal with respect to $B$* if there is a direction $\boldsymbol{v}$ such that $\widehat{\alpha}$ is extremal in direction $\boldsymbol{v}$ with respect to $B$.

In the following we prove some facts about extremal local strategies that will be needed later.

**Lemma 5.4.** *Let $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ be an MDP. Let $B \in \mathbb{R}^{Q \times k}$ with $k \geq 1$. Let $\mu \in \mathsf{subDist}(Q)$. Let $\alpha, \widehat{\alpha}$ be local strategies. Suppose $\boldsymbol{v} \in \mathbb{R}^{|\mathsf{L}| \cdot k}$ is a direction in which $\widehat{\alpha}$ is extremal and $p(\mu, \alpha)\boldsymbol{v} = p(\mu, \widehat{\alpha})\boldsymbol{v}$. Then $p(\mu, \alpha) = p(\mu, \widehat{\alpha})$.*

*Proof.* We have:

$$\sum_{q \in Q} \mu(q) \cdot p(d_q, \alpha)\boldsymbol{v} = p(\mu, \alpha)\boldsymbol{v} \qquad\qquad \text{definition of } p$$

$$= p(\mu, \widehat{\alpha})\boldsymbol{v} \qquad\qquad \text{assumption on } \widehat{\alpha}$$

$$= \sum_{q \in Q} \mu(q) \cdot p(d_q, \widehat{\alpha})\boldsymbol{v} \qquad\qquad \text{definition of } p$$

With (5.2) it follows that for all $q \in \mathsf{Supp}(\mu)$ we have $p(d_q, \alpha)\boldsymbol{v} = p(d_q, \widehat{\alpha})\boldsymbol{v}$. Hence by (5.3) we obtain $p(d_q, \alpha) = p(d_q, \widehat{\alpha})$ for all $q \in \mathsf{Supp}(\mu)$. Thus:

$$p(\mu, \alpha) = \sum_{q \in Q} \mu(q) \cdot p(d_q, \alpha) = \sum_{q \in Q} \mu(q) \cdot p(d_q, \widehat{\alpha}) = p(\mu, \widehat{\alpha}) \qquad\qquad \square$$

For a subdistribution $\mu$ define the bounded, convex polytope $P_\mu \subseteq \mathbb{R}^{|\mathsf{L}| \cdot k}$ with

$$P_\mu = \{p(\mu, \alpha) \mid \alpha : Q \to \mathsf{Dist}(\mathsf{moves})\}.$$

Comparing two polytopes $P_{\mu_\mathcal{D}}$ and $P_{\mu_\mathcal{E}}$ for subdistributions $\mu_\mathcal{D}, \mu_\mathcal{E}$ will play a key role for deciding bisimulation. First we prove the following lemma, which states that any vertex of the polytope $P_\mu$ can be obtained by applying an extremal local strategy. Although this is intuitive, the proof is not very easy.

**Lemma 5.5.** *Let $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ be an MDP. Let $B \in \mathbb{R}^{Q \times k}$ with $k \geq 1$. Let $\mu \in \mathsf{subDist}(Q)$. If $x \in P_\mu$ is a vertex of $P_\mu$ then there is an extremal local strategy $\widehat{\alpha}$ with $x = p(\mu, \widehat{\alpha})$.*

*Proof.* Let $x \in P_\mu$ be a vertex of $P_\mu$. Let $\alpha_1 : Q \to \mathsf{Dist}(\mathsf{moves})$ be a local strategy so that $x = p(\mu, \alpha_1)$. Since $x$ is a vertex, we can assume that $\alpha_1$ is pure. Since $x$ is a vertex of $P_\mu$, there is a hyperplane $H \subseteq \mathbb{R}^{|\mathsf{L}| \cdot k}$ such that $\{x\} = P_\mu \cap H$. Let $\boldsymbol{v}_1 \in \mathbb{R}^{|\mathsf{L}| \cdot k}$ be a normal vector of $H$. Since $\{x\} = P_\mu \cap H$, we have $x\boldsymbol{v}_1 = \max_{y \in P_\mu} y\boldsymbol{v}_1$ or $x\boldsymbol{v}_1 = \min_{y \in P_\mu} y\boldsymbol{v}_1$; without loss of generality, say $x\boldsymbol{v}_1 = \max_{y \in P_\mu} y\boldsymbol{v}_1$. Since $\{x\} = P_\mu \cap H$, we have for all $q \in \mathsf{Supp}(\mu)$ and all $\alpha$:

$$p(d_q, \alpha)\boldsymbol{v}_1 = p(d_q, \alpha_1)\boldsymbol{v}_1 \quad \text{implies} \quad p(d_q, \alpha) = p(d_q, \alpha_1). \tag{5.4}$$

For all $q \in Q \setminus \mathsf{Supp}(\mu)$, redefine the pure local strategy $\alpha_1(q)$ so that all $q \in Q$ and all local strategies $\alpha$ satisfy $p(d_q, \alpha)\boldsymbol{v}_1 \leq p(d_q, \alpha_1)\boldsymbol{v}_1$. Since $Q$ and $\mathsf{moves}$ are finite, there is $\varepsilon > 0$ such that all $q \in Q$ and all *pure* local strategies $\alpha$ either satisfy $p(d_q, \alpha)\boldsymbol{v}_1 = p(d_q, \alpha_1)\boldsymbol{v}_1$ or $p(d_q, \alpha)\boldsymbol{v}_1 \leq p(d_q, \alpha_1)\boldsymbol{v}_1 - \varepsilon$.

Define

$$\Sigma = \{\alpha : Q \to \mathsf{Dist}(\mathsf{moves}) \mid \forall\, q \in Q : p(d_q, \alpha)\boldsymbol{v}_1 = p(d_q, \alpha_1)\boldsymbol{v}_1\}.$$

Consider the bounded, convex polytope $P_2 \subseteq \mathbb{R}^{|\mathsf{L}| \cdot k}$ defined by

$$P_2 = \left\{ \sum_{q \in Q} p(d_q, \alpha) \; \middle| \; \alpha \in \Sigma \right\}.$$

By an argument similar to the one above, there are a pure local strategy $\widehat{\alpha} \in \Sigma$, a vertex $x_2 = \sum_{q \in Q} p(d_q, \widehat{\alpha})$ of $P_2$, and a vector $\boldsymbol{v}_2 \in \mathbb{R}^{|\mathsf{L}| \cdot k}$ such that for all $q \in Q$ and all $\alpha \in \Sigma$, we have $p(d_q, \alpha)\boldsymbol{v}_2 \leq p(d_q, \widehat{\alpha})\boldsymbol{v}_2$, and if $p(d_q, \alpha)\boldsymbol{v}_2 = p(d_q, \widehat{\alpha})\boldsymbol{v}_2$ then $p(d_q, \alpha) = p(d_q, \widehat{\alpha})$. By scaling down $\boldsymbol{v}_2$ by a small positive scalar, we can assume that all $q \in Q$ and all local strategies $\alpha$ satisfy

$$|p(d_q, \alpha)\boldsymbol{v}_2| \leq \frac{\varepsilon}{3}. \tag{5.5}$$

Since $\widehat{\alpha} \in \Sigma$, all $q \in Q$ satisfy $p(d_q, \widehat{\alpha})\boldsymbol{v}_1 = p(d_q, \alpha_1)\boldsymbol{v}_1$. By (5.4) all $q \in \mathsf{Supp}(\mu)$ satisfy $p(d_q, \widehat{\alpha}) = p(d_q, \alpha_1)$. Hence:

$$p(\mu, \widehat{\alpha}) = \sum_{q \in \mathsf{Supp}(\mu)} \mu(q)p(d_q, \widehat{\alpha}) = \sum_{q \in \mathsf{Supp}(\mu)} \mu(q)p(d_q, \alpha_1) = p(\mu, \alpha_1) = x$$

It remains to show that there is a direction $\boldsymbol{v}$ in which $\widehat{\alpha}$ is extremal. Take $\boldsymbol{v} = \boldsymbol{v}_1 + \boldsymbol{v}_2$. Let $q \in Q$ and let $\alpha$ be a pure local strategy. We consider two cases:

- Assume $p(d_q, \alpha)\boldsymbol{v}_1 = p(d_q, \alpha_1)\boldsymbol{v}_1$. Then there is $\beta \in \Sigma$ with $\alpha(q) = \beta(q)$, hence $p(d_q, \alpha) = p(d_q, \beta)$. We have:

$$
\begin{aligned}
p(d_q, \alpha)\boldsymbol{v} &= p(d_q, \beta)\boldsymbol{v} && p(d_q, \alpha) = p(d_q, \beta) \\
&= p(d_q, \beta)\boldsymbol{v}_1 + p(d_q, \beta)\boldsymbol{v}_2 && \text{definition of } \boldsymbol{v} \\
&= p(d_q, \alpha_1)\boldsymbol{v}_1 + p(d_q, \beta)\boldsymbol{v}_2 && \beta \in \Sigma \\
&= p(d_q, \widehat{\alpha})\boldsymbol{v}_1 + p(d_q, \beta)\boldsymbol{v}_2 && \widehat{\alpha} \in \Sigma \\
&\le p(d_q, \widehat{\alpha})\boldsymbol{v}_1 + p(d_q, \widehat{\alpha})\boldsymbol{v}_2 && \text{definition of } \widehat{\alpha} \\
&= p(d_q, \widehat{\alpha})\boldsymbol{v} && \text{definition of } \boldsymbol{v}
\end{aligned}
$$

  Hence (5.2) holds for $\widehat{\alpha}$. To show (5.3), assume $p(d_q, \alpha)\boldsymbol{v} = p(d_q, \widehat{\alpha})\boldsymbol{v}$. Then all terms in the computation above are equal, and $p(d_q, \beta)\boldsymbol{v}_2 = p(d_q, \widehat{\alpha})\boldsymbol{v}_2$. By the definition of $\widehat{\alpha}$, this implies $p(d_q, \beta) = p(d_q, \widehat{\alpha})$. Hence $p(d_q, \alpha) = p(d_q, \beta) = p(d_q, \widehat{\alpha})$. Hence (5.3) holds for $\widehat{\alpha}$.
- Assume $p(d_q, \alpha)\boldsymbol{v}_1 \ne p(d_q, \alpha_1)\boldsymbol{v}_1$. By the definition of $\varepsilon$ it follows $p(d_q, \alpha)\boldsymbol{v}_1 \le p(d_q, \alpha_1)\boldsymbol{v}_1 - \varepsilon$. We have:

$$
\begin{aligned}
p(d_q, \alpha)\boldsymbol{v} &= p(d_q, \alpha)\boldsymbol{v}_1 + p(d_q, \alpha)\boldsymbol{v}_2 && \text{definition of } \boldsymbol{v} \\
&\le p(d_q, \alpha_1)\boldsymbol{v}_1 - \varepsilon + p(d_q, \alpha)\boldsymbol{v}_2 && \text{as argued above} \\
&= p(d_q, \widehat{\alpha})\boldsymbol{v}_1 - \varepsilon + p(d_q, \alpha)\boldsymbol{v}_2 && \widehat{\alpha} \in \Sigma \\
&\le p(d_q, \widehat{\alpha})\boldsymbol{v}_1 - \varepsilon + \frac{\varepsilon}{3} && \text{by (5.5)} \\
&\le p(d_q, \widehat{\alpha})\boldsymbol{v}_1 + p(d_q, \widehat{\alpha})\boldsymbol{v}_2 - \varepsilon + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} && \text{by (5.5)} \\
&< p(d_q, \widehat{\alpha})\boldsymbol{v}_1 + p(d_q, \widehat{\alpha})\boldsymbol{v}_2 && \varepsilon > 0 \\
&= p(d_q, \widehat{\alpha})\boldsymbol{v} && \text{definition of } \boldsymbol{v}
\end{aligned}
$$

  This implies (5.2) and (5.3) for $\widehat{\alpha}$.

Hence, $\widehat{\alpha}$ is extremal in direction $\boldsymbol{v}$. $\qquad\square$

The following lemma states the intuitive fact that in order to compare the polytopes $P_{\mu_{\mathcal{D}}}$ and $P_{\mu_{\mathcal{E}}}$, it suffices to compare the vertices obtained by applying extremal local strategies:

**Lemma 5.6.** Let $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ be an MDP. Let $B \in \mathbb{R}^{Q \times k}$ with $k \ge 1$. Then for all $\mu_{\mathcal{D}}, \mu_{\mathcal{E}} \in \mathsf{subDist}(Q)$ we have $P_{\mu_{\mathcal{D}}} = P_{\mu_{\mathcal{E}}}$ if and only if for all extremal local strategies $\widehat{\alpha}$ we have $p(\mu_{\mathcal{D}}, \widehat{\alpha}) = p(\mu_{\mathcal{E}}, \widehat{\alpha})$.

*Proof.* We prove the two implications from the lemma in turn.

"$\Longrightarrow$": Suppose $P_{\mu_{\mathcal{D}}} = P_{\mu_{\mathcal{E}}}$. Let $\widehat{\alpha}$ be a local strategy that is extremal in direction $\boldsymbol{v}$. Since $P_{\mu_{\mathcal{D}}} = P_{\mu_{\mathcal{E}}}$, there are $\alpha_{\mathcal{E}}$ and $\alpha_{\mathcal{D}}$ such that $p(\mu_{\mathcal{D}}, \widehat{\alpha}) = p(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}})$ and $p(\mu_{\mathcal{E}}, \widehat{\alpha}) = p(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}})$. We have:

$$
\begin{aligned}
p(\mu_{\mathcal{D}}, \widehat{\alpha})\boldsymbol{v} &= p(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}})\boldsymbol{v} && p(\mu_{\mathcal{D}}, \widehat{\alpha}) = p(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}}) \\
&\le p(\mu_{\mathcal{E}}, \widehat{\alpha})\boldsymbol{v} && \widehat{\alpha} \text{ is extremal in direction } \boldsymbol{v} \\
&= p(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}})\boldsymbol{v} && p(\mu_{\mathcal{E}}, \widehat{\alpha}) = p(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}}) \\
&\le p(\mu_{\mathcal{D}}, \widehat{\alpha})\boldsymbol{v} && \widehat{\alpha} \text{ is extremal in direction } \boldsymbol{v}
\end{aligned}
$$

So all inequalities are in fact equalities. In particular, we have $p(\mu_{\mathcal{D}}, \widehat{\alpha})\boldsymbol{v} = p(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}})\boldsymbol{v}$. It follows:

$$\begin{aligned} p(\mu_{\mathcal{D}}, \widehat{\alpha}) &= p(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}}) && \text{Lemma 5.4} \\ &= p(\mu_{\mathcal{E}}, \widehat{\alpha}) && \text{definition of } \alpha_{\mathcal{D}} \end{aligned}$$

"$\Longleftarrow$": Let $x$ be a vertex of $P_{\mu_{\mathcal{D}}}$. By Lemma 5.5 there exists an extremal local strategy $\widehat{\alpha}$ with $x = p(\mu_{\mathcal{D}}, \widehat{\alpha})$. By the assumption we have $p(\mu_{\mathcal{D}}, \widehat{\alpha}) = p(\mu_{\mathcal{E}}, \widehat{\alpha})$. Hence $x = p(\mu_{\mathcal{D}}, \widehat{\alpha}) = p(\mu_{\mathcal{E}}, \widehat{\alpha}) \in P_{\mu_{\mathcal{E}}}$. Since $x$ is an arbitrary vertex of $P_{\mu_{\mathcal{D}}}$, and $P_{\mu_{\mathcal{D}}}, P_{\mu_{\mathcal{E}}}$ are bounded, convex polytopes, it follows $P_{\mu_{\mathcal{D}}} \subseteq P_{\mu_{\mathcal{E}}}$. The reverse inclusion is shown similarly. $\square$

The following lemma shows that the alternation of quantifiers over local strategies can be replaced by quantifying over extremal local strategies:

**Lemma 5.7.** *Let $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ be an MDP. Let $B \in \mathbb{R}^{Q \times k}$ with $k \geq 1$. Let $\mu_{\mathcal{D}}, \mu_{\mathcal{E}} \in \mathsf{subDist}(Q)$. In the following let $\alpha_{\mathcal{D}}, \alpha_{\mathcal{E}}$ range over local strategies, $\widehat{\alpha}$ over extremal local strategies, and $a$ over $\mathsf{L}$. Then*

$$\begin{aligned} \forall \alpha_{\mathcal{D}} \exists \alpha_{\mathcal{E}} \forall a &: \mu_{\mathcal{D}} \Delta_{\alpha_{\mathcal{D}}}(a)B = \mu_{\mathcal{E}} \Delta_{\alpha_{\mathcal{E}}}(a)B \\ \wedge \quad \forall \alpha_{\mathcal{E}} \exists \alpha_{\mathcal{D}} \forall a &: \mu_{\mathcal{D}} \Delta_{\alpha_{\mathcal{D}}}(a)B = \mu_{\mathcal{E}} \Delta_{\alpha_{\mathcal{E}}}(a)B \end{aligned} \qquad (5.6)$$

*holds if and only the following holds:*

$$\forall \widehat{\alpha} \forall a : \mu_{\mathcal{D}} \Delta_{\widehat{\alpha}}(a)B = \mu_{\mathcal{E}} \Delta_{\widehat{\alpha}}(a)B$$

*Proof.* We have:

$$\begin{aligned} & \forall \alpha_{\mathcal{D}} \exists \alpha_{\mathcal{E}} \forall a : \mu_{\mathcal{D}} \Delta_{\alpha_{\mathcal{D}}}(a)B = \mu_{\mathcal{E}} \Delta_{\alpha_{\mathcal{E}}}(a)B \\ \Longleftrightarrow \quad & \forall \alpha_{\mathcal{D}} \exists \alpha_{\mathcal{E}} : p(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}}) = p(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}}) && \text{definition of } p \\ \Longleftrightarrow \quad & P_{\mu_{\mathcal{D}}} \subseteq P_{\mu_{\mathcal{E}}} \end{aligned}$$

It follows:

$$\begin{aligned} (5.6) \quad \Longleftrightarrow \quad & P_{\mu_{\mathcal{D}}} = P_{\mu_{\mathcal{E}}} \\ \Longleftrightarrow \quad & \forall \widehat{\alpha} \forall a : \mu_{\mathcal{D}} \Delta_{\widehat{\alpha}}(a)B = \mu_{\mathcal{E}} \Delta_{\widehat{\alpha}}(a)B && \text{Lemma 5.6} \qquad \square \end{aligned}$$

The following proposition provides necessary and sufficient conditions for bisimilarity, which—as we will see—can be effectively checked.

**Proposition 5.8.** *Let $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$ be an MDP. Let $B \in \mathbb{R}^{Q \times k}$ with $k \geq 1$.*

(1) *Suppose that for all $\mu_{\mathcal{D}}, \mu_{\mathcal{E}} \in \mathsf{subDist}(Q)$ with $\mu_{\mathcal{D}} \sim \mu_{\mathcal{E}}$ we have $\mu_{\mathcal{D}} B = \mu_{\mathcal{E}} B$. Then for all $\mu_{\mathcal{D}}, \mu_{\mathcal{E}} \in \mathsf{subDist}(Q)$ with $\mu_{\mathcal{D}} \sim \mu_{\mathcal{E}}$ we have $\mu_{\mathcal{D}} \Delta_{\widehat{\alpha}}(a)B = \mu_{\mathcal{E}} \Delta_{\widehat{\alpha}}(a)B$ for all extremal local strategies $\widehat{\alpha}$ and all $a \in \mathsf{L}$.*

(2) *Suppose that $B$ includes the column vector $\mathbf{1} = (1\ 1 \cdots 1)^T$ (where the superscript $T$ denotes transpose) and that for all extremal local strategies $\widehat{\alpha}$ and all $a \in \mathsf{L}$ the columns of $\Delta_{\widehat{\alpha}}(a)B$ are in the linear span of the columns of $B$. Then for all $\mu_{\mathcal{D}}, \mu_{\mathcal{E}} \in \mathsf{subDist}(Q)$ with $\mu_{\mathcal{D}} B = \mu_{\mathcal{E}} B$ we have $\mu_{\mathcal{D}} \sim \mu_{\mathcal{E}}$.*

*Proof.*

(1) Let $\mu_{\mathcal{D}}, \mu_{\mathcal{E}} \in \mathsf{subDist}(Q)$ with $\mu_{\mathcal{D}} \sim \mu_{\mathcal{E}}$. By the definition of bisimulation and using (5.1), we obtain:

– for all local strategies $\alpha_{\mathcal{D}}$ there exists a local strategy $\alpha_{\mathcal{E}}$ such that for all $a \in \mathsf{L}$ we have $\mu_{\mathcal{D}} \Delta_{\alpha_{\mathcal{D}}}(a) \sim \mu_{\mathcal{E}} \Delta_{\alpha_{\mathcal{E}}}(a)$;

– for all local strategies $\alpha_{\mathcal{E}}$ there exists a local strategy $\alpha_{\mathcal{D}}$ such that for all $a \in \mathsf{L}$ we have $\mu_{\mathcal{D}} \Delta_{\alpha_{\mathcal{D}}}(a) \sim \mu_{\mathcal{E}} \Delta_{\alpha_{\mathcal{E}}}(a)$.

Using our assumption on $B$, we see that (5.6) from Lemma 5.7 holds for $\mu_{\mathcal{D}}, \mu_{\mathcal{E}}$. By Lemma 5.7 we have $\mu_{\mathcal{D}} \Delta_{\widehat{\alpha}}(a)B = \mu_{\mathcal{E}} \Delta_{\widehat{\alpha}}(a)B$ for all extremal local strategies $\widehat{\alpha}$ and all $a \in \mathsf{L}$.

(2) It suffices to show that the relation $\sim_B \subseteq \mathsf{subDist}(Q) \times \mathsf{subDist}(Q)$ defined by

$$\mu_{\mathcal{D}} \sim_B \mu_{\mathcal{E}} \quad \Longleftrightarrow \quad \mu_{\mathcal{D}}B = \mu_{\mathcal{E}}B$$

is a bisimulation. Let $\mu_{\mathcal{D}} \sim_B \mu_{\mathcal{E}}$, i.e., $\mu_{\mathcal{D}}B = \mu_{\mathcal{E}}B$. Since $B$ includes the column $\mathbf{1}$, we have $\|\mu_{\mathcal{D}}\| = \|\mu_{\mathcal{E}}\|$. Since for all extremal local strategies $\widehat{\alpha}$ and all $a \in \mathsf{L}$ the columns of $\Delta_{\widehat{\alpha}}(a)B$ are in the linear span of the columns of $B$, we have $\mathbf{0}^T = (\mu_{\mathcal{D}} - \mu_{\mathcal{E}})B = (\mu_{\mathcal{D}} - \mu_{\mathcal{E}})\Delta_{\widehat{\alpha}}(a)B$ for all extremal local strategies $\widehat{\alpha}$ and all $a \in \mathsf{L}$. Lemma 5.7 implies that (5.6) holds for $\mu_{\mathcal{D}}, \mu_{\mathcal{E}}$. Using (5.1) and the definition of $\sim_B$, we obtain:
– for all local strategies $\alpha_{\mathcal{D}}$ there exists a local strategy $\alpha_{\mathcal{E}}$ such that for all $a \in \mathsf{L}$ we have $\mathsf{Succ}(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}}, a) \sim_B \mathsf{Succ}(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}}, a)$;
– for all local strategies $\alpha_{\mathcal{E}}$ there exists a local strategy $\alpha_{\mathcal{D}}$ such that for all $a \in \mathsf{L}$ we have $\mathsf{Succ}(\mu_{\mathcal{D}}, \alpha_{\mathcal{D}}, a) \sim_B \mathsf{Succ}(\mu_{\mathcal{E}}, \alpha_{\mathcal{E}}, a)$.

Thus the relation $\sim_B$ is a bisimulation. $\qquad\square$

5.3. **A coNP Algorithm for Checking Bisimilarity of two MDPs.** Proposition 5.8 suggests an algorithm for determining bisimilarity in a given MDP $\mathcal{D} = \langle Q, \mu_0, \mathsf{L}, \delta \rangle$. More concretely, we can compute a matrix $B$ such that for all subdistributions $\mu_{\mathcal{D}}, \mu_{\mathcal{E}}$ we have $\mu_{\mathcal{D}} \sim \mu_{\mathcal{E}}$ if and only if $\mu_{\mathcal{D}}B = \mu_{\mathcal{E}}B$. The algorithm initializes $B$ with the column vector $\mathbf{1}$ and henceforth maintains the invariant that $\mu_{\mathcal{D}} \sim \mu_{\mathcal{E}}$ implies $\mu_{\mathcal{D}}B = \mu_{\mathcal{E}}B$.

- If there exists an extremal local strategy $\widehat{\alpha}$ and a label $a \in \mathsf{L}$ such that a column of $\Delta_{\widehat{\alpha}}(a)B$ is linearly independent of the columns of $B$, then add this column to $B$. This maintains the invariant by Proposition 5.8 (1). Repeat.
- Otherwise (i.e., such $\widehat{\alpha}$ and $a$ do not exist) terminate. Then by Proposition 5.8 (2) we have $\mu_{\mathcal{D}}B = \mu_{\mathcal{E}}B \implies \mu_{\mathcal{D}} \sim \mu_{\mathcal{E}}$. Together with the invariant we get $\mu_{\mathcal{D}} \sim \mu_{\mathcal{E}} \iff \mu_{\mathcal{D}}B = \mu_{\mathcal{E}}B$, as claimed.

This algorithm terminates because $B$ can have at most $|Q|$ linearly independent columns.

Along these lines we can also prove the following theorem.

**Theorem 5.9.** *The problem that, given two MDPs $\mathcal{D}$ and $\mathcal{E}$, asks whether $\mathcal{D} \sim \mathcal{E}$ is in* coNP.

*Proof.* Without loss of generality we assume $\mathcal{D} = \langle Q, \mu_{\mathcal{D}}, \mathsf{L}, \delta \rangle$ and $\mathcal{E} = \langle Q, \mu_{\mathcal{E}}, \mathsf{L}, \delta \rangle$. Hence we wish to decide in NP whether $\mu_{\mathcal{D}} \not\sim \mu_{\mathcal{E}}$.

We proceed along the lines of the algorithm above. Specifically, it follows from the arguments there that $\mu_{\mathcal{D}} \not\sim \mu_{\mathcal{E}}$ holds if and only if the following condition *Cond* holds:

There are $k \in \{1, 2, \ldots, |Q|\}$ and $\boldsymbol{b}_0 = \mathbf{1}$ and $\boldsymbol{b}_1, \ldots, \boldsymbol{b}_{k-1} \in \mathbb{R}^Q$ and $i_0, i_1, \ldots, i_{k-1} \in \{0, 1, \ldots, k - 2\}$ and $a_1, a_2, \ldots, a_{k-1} \in \mathsf{L}$ and pure local strategies $\widehat{\alpha}_1, \widehat{\alpha}_2, \ldots, \widehat{\alpha}_{k-1}$ such that for all $j \in \{1, 2, \ldots, k - 1\}$
- $\widehat{\alpha}_j$ is extremal with respect to the matrix formed by the column vectors $\boldsymbol{b}_0, \boldsymbol{b}_1, \ldots, \boldsymbol{b}_{j-1}$ and
- $i_j < j$ and
- $\boldsymbol{b}_j = \Delta_{\widehat{\alpha}_j}(a_j)\boldsymbol{b}_{i_j}$,

and $\mu_{\mathcal{D}} \boldsymbol{b}_{k-1} \neq \mu_{\mathcal{E}} \boldsymbol{b}_{k-1}$.

It remains to argue that $Cond$ can be checked in NP. We can nondeterministically guess $k \leq |Q|$ and $i_0, i_1, \ldots, i_{k-1} \leq k-2$ and $a_1, a_2, \ldots, a_{k-1} \in \mathsf{L}$ and pure local strategies $\widehat{\alpha}_1, \widehat{\alpha}_2, \ldots, \widehat{\alpha}_{k-1}$. This determines $\boldsymbol{b}_1, \ldots, \boldsymbol{b}_{k-1}$. All conditions in $Cond$ are straightforward to check in polynomial time, except the condition that for all $j \in \{1, 2, \ldots, k-1\}$ we have that $\widehat{\alpha}_j$ is extremal with respect to $\boldsymbol{b}_0, \boldsymbol{b}_1, \ldots, \boldsymbol{b}_{j-1}$. In the remainder of the proof, we argue that this can also be checked in polynomial time.

Let $j \in \{1, 2, \ldots, k-1\}$. Let $B \in \mathbb{R}^{Q \times j}$ be the matrix with columns $\boldsymbol{b}_0, \boldsymbol{b}_1, \ldots, \boldsymbol{b}_{j-1}$. We want to check that $\widehat{\alpha}_j$ is extremal with respect to $B$. For all $q \in Q$, compute in polynomial time the set $\mathsf{eqmoves}(q) \subseteq \mathsf{moves}(q)$ defined by

$$\mathsf{eqmoves}(q) = \{\mathsf{m} \in \mathsf{moves}(q) \mid p(d_q, \alpha_{q,\mathsf{m}}) = p(d_q, \widehat{\alpha}_j)\},$$

where $\alpha_{q,\mathsf{m}}$ is a pure local strategy with $\alpha_{q,\mathsf{m}}(q)(\mathsf{m}) = 1$ (it does not matter how $\alpha_{q,\mathsf{m}}(q')$ is defined for $q' \neq q$). We want to verify that (5.2) and (5.3) holds for $\widehat{\alpha}_j$. Hence we need to find $\boldsymbol{v} \in \mathbb{R}^{|\mathsf{L}| \cdot j}$ so that for all $q \in Q$ and all $\mathsf{m} \in \mathsf{moves}(q) \setminus \mathsf{eqmoves}(q)$ we have $p(d_q, \alpha_{q,\mathsf{m}})\boldsymbol{v} < p(d_q, \widehat{\alpha}_j)\boldsymbol{v}$. If such a vector $\boldsymbol{v}$ exists, it can be scaled up by a large positive scalar so that we have:

$$p(d_q, \alpha_{q,\mathsf{m}})\boldsymbol{v} + 1 \leq p(d_q, \widehat{\alpha}_j)\boldsymbol{v} \quad \forall q \in Q \quad \forall \mathsf{m} \in \mathsf{moves}(q) \setminus \mathsf{eqmoves}(q) \tag{5.7}$$

Hence it suffices to check if there exists a vector $\boldsymbol{v}$ that satisfies (5.7). This amounts to a feasibility check of a linear program of polynomial size. Such a check can be carried out in polynomial time [12]. □

5.4. **An NC Algorithm for Trace Refinement.** In the following we consider an MDP $\mathcal{D} = \langle Q, \mu_0^{\mathcal{D}}, \mathsf{L}, \delta \rangle$ and an MC $\mathcal{C} = \langle Q_{\mathcal{C}}, \mu_0^{\mathcal{C}}, \mathsf{L}, \delta_{\mathcal{C}} \rangle$. Without loss of generality, we assume $Q_{\mathcal{C}} \subseteq Q$. Similarly, we also assume that $\delta_{\mathcal{C}}$ is a restriction of $\delta$, and hence we write $\mathcal{C} = \langle Q_{\mathcal{C}}, \mu_0^{\mathcal{C}}, \mathsf{L}, \delta \rangle$. We may view subdistributions $\mu_{\mathcal{C}} \in \mathsf{subDist}(Q_{\mathcal{C}})$ as $\mu_{\mathcal{C}} \in \mathsf{subDist}(Q)$ in the natural way. The following proposition is analogous to Proposition 5.8. The key difference is that the need for considering *extremal* strategies has disappeared. This is due to the fact that only one of the two models is nondeterministic.

**Proposition 5.10.** *Let $\mathcal{D} = \langle Q, \mu_0^{\mathcal{D}}, \mathsf{L}, \delta \rangle$ be an MDP and $\mathcal{C} = \langle Q_{\mathcal{C}}, \mu_0^{\mathcal{C}}, \mathsf{L}, \delta \rangle$ be an MC with $Q_{\mathcal{C}} \subseteq Q$. Let $B \in \mathbb{R}^{Q \times k}$ with $k \geq 1$. In the following let $\mu_{\mathcal{D}}$ range over $\mathsf{subDist}(Q)$ and $\mu_{\mathcal{C}}$ over $\mathsf{subDist}(Q_{\mathcal{C}})$.*

(1) *Suppose that for all $\mu_{\mathcal{D}}, \mu_{\mathcal{C}}$ with $\mu_{\mathcal{D}} \sim \mu_{\mathcal{C}}$ we have $\mu_{\mathcal{D}} B = \mu_{\mathcal{C}} B$. Then for all $\mu_{\mathcal{D}}, \mu_{\mathcal{C}}$ with $\mu_{\mathcal{D}} \sim \mu_{\mathcal{C}}$ we have $\mu_{\mathcal{D}} \Delta_\alpha(a) B = \mu_{\mathcal{C}} \Delta_\alpha(a) B$ for all local strategies $\alpha$ and all $a \in \mathsf{L}$.*
(2) *Suppose that $B$ includes the column vector $\mathbf{1} = (1 \ 1 \cdots 1)^T$ and that for all local strategies $\alpha$ and all $a \in \mathsf{L}$ the columns of $\Delta_\alpha(a) B$ are in the linear span of the columns of $B$. Then for all $\mu_{\mathcal{D}}, \mu_{\mathcal{C}}$ with $\mu_{\mathcal{D}} B = \mu_{\mathcal{C}} B$ we have $\mu_{\mathcal{D}} \sim \mu_{\mathcal{C}}$.*

*Proof.* The proof is similar to but simpler than the proof of Proposition 5.8. For completeness, we give it explicitly.

(1) Let $\mu_{\mathcal{D}} \sim \mu_{\mathcal{C}}$. By the definition of bisimulation and using (5.1) we have $\mu_{\mathcal{D}} \Delta_\alpha(a) \sim \mu_{\mathcal{C}} \Delta_\alpha(a)$ for all local strategies $\alpha$ and all $a \in \mathsf{L}$. By our assumption on $B$, we have $\mu_{\mathcal{D}} \Delta_\alpha(a) B = \mu_{\mathcal{C}} \Delta_\alpha(a) B$ for all local strategies $\alpha$ and all $a \in \mathsf{L}$.

(2) It suffices to show that the relation $\sim_B \subseteq \mathsf{subDist}(Q) \times \mathsf{subDist}(Q_\mathcal{C})$ defined by

$$\mu_\mathcal{D} \sim_B \mu_\mathcal{C} \quad \Longleftrightarrow \quad \mu_\mathcal{D} B = \mu_\mathcal{C} B$$

is a bisimulation. Let $\mu_\mathcal{D} \sim_B \mu_\mathcal{C}$, i.e., $\mu_\mathcal{D} B = \mu_\mathcal{C} B$. Since $B$ includes the column $\mathbf{1}$, we have $\|\mu_\mathcal{D}\| = \|\mu_\mathcal{C}\|$. Since for all local strategies $\alpha$ and all $a \in \mathsf{L}$ the columns of $\Delta_\alpha(a)B$ are in the linear span of the columns of $B$, we have $\mathbf{0}^T = (\mu_\mathcal{D} - \mu_\mathcal{C})B = (\mu_\mathcal{D} - \mu_\mathcal{C})\Delta_\alpha(a)B$ for all local strategies $\alpha$ and all $a \in \mathsf{L}$. Using (5.1) and the definition of $\sim_B$, we see that for all local strategies $\alpha$ and all $a \in \mathsf{L}$ we have $\mathsf{Succ}(\mu_\mathcal{D}, \alpha, a) \sim_B \mathsf{Succ}(\mu_\mathcal{C}, \alpha, a)$. Thus the relation $\sim_B$ is a bisimulation. $\qquad\square$

**Corollary 5.11.** *Let $\mathcal{V} \subseteq \mathbb{R}^Q$ be the smallest column-vector space with $\mathbf{1} \in \mathcal{V}$ and $\Delta_\alpha(a)\boldsymbol{u} \in \mathcal{V}$ for all $\boldsymbol{u} \in \mathcal{V}$, all labels $a \in \mathsf{L}$, and all local strategies $\alpha$. Then for all $\mu_\mathcal{D} \in \mathsf{subDist}(Q)$ and all $\mu_\mathcal{C} \in \mathsf{subDist}(Q_\mathcal{C})$ we have:*

$$\mu_\mathcal{D} \sim \mu_\mathcal{C} \quad \Longleftrightarrow \quad \mu_\mathcal{D}\boldsymbol{u} = \mu_\mathcal{C}\boldsymbol{u} \ \text{ for all } \boldsymbol{u} \in \mathcal{V}$$

Notice the differences to Proposition 5.8: there we considered all extremal local strategies (potentially exponentially many), here we consider all local strategies (in general infinitely many). However, we show that one can efficiently find few local strategies that span all local strategies. This allows us to reduce (in logarithmic space) the bisimulation problem between an MDP and an MC to the bisimulation problem between two MCs, which is equivalent to the trace-equivalence problem in MCs (by Proposition 5.3). The latter problem is known to be in $\mathsf{NC}$ [24]. Theorem 5.12 then follows with Proposition 5.3.

**Theorem 5.12.** *The problem $\mathsf{MDP} \sqsubseteq \mathsf{MC}$ is in $\mathsf{NC}$, hence in $\mathsf{P}$.*

*Proof.* Let $\mathcal{D} = \langle Q, \mu_0^\mathcal{D}, \mathsf{L}, \delta \rangle$ be an MDP and $\mathcal{C} = \langle Q_\mathcal{C}, \mu_0^\mathcal{C}, \mathsf{L}, \delta \rangle$ be an MC with $Q_\mathcal{C} \subseteq Q$.

Let $\alpha_0$ denote an arbitrary pure local strategy. For each $q \in Q$ and each $\mathsf{m} \in \mathsf{moves}(q)$ denote by $\alpha_{q,\mathsf{m}}$ the pure local strategy such that $\alpha_{q,\mathsf{m}}(q)(\mathsf{m}) = 1$ and $\alpha_{q,\mathsf{m}}(q') = \alpha_0(q')$ for all $q' \in Q \setminus \{q\}$. Define

$$\Sigma = \{\alpha_0\} \cup \{\alpha_{q,\mathsf{m}} \mid q \in Q, \ \mathsf{m} \in \mathsf{moves}(q)\} \qquad\qquad \text{and}$$

$$\mathcal{M} = \{\Delta_\alpha(a) \in \mathbb{R}^{Q \times Q} \mid \alpha \in \Sigma, \ a \in \mathsf{L}\} \qquad\qquad \text{and}$$

$$\mathcal{M}_\infty = \{\Delta_\alpha(a) \in \mathbb{R}^{Q \times Q} \mid \alpha \text{ is a local strategy}, \ a \in \mathsf{L}\}.$$

The vector space $\mathcal{V} \subseteq \mathbb{R}^Q$ from Corollary 5.11 is the smallest vector space with
- $\mathbf{1} = (1\ 1 \cdots 1)^T \in \mathcal{V}$ and
- $M\boldsymbol{u} \in \mathcal{V}$, for all $\boldsymbol{u} \in \mathcal{V}$ and all $M \in \mathcal{M}_\infty$.

We have $\mathcal{M} \subseteq \mathcal{M}_\infty$, where $|\mathcal{M}|$ is finite and $|\mathcal{M}_\infty|$ is infinite. Every matrix in $\mathcal{M}_\infty$ can be expressed as a linear combination of matrices from $\mathcal{M}$: Indeed, let $\alpha$ be a local strategy. Then for all $a \in \mathsf{L}$ we have:

$$\Delta_\alpha(a) = \Delta_{\alpha_0}(a) + \sum_{q \in Q}\left(-\Delta_{\alpha_0}(a) + \sum_{\mathsf{m} \in \mathsf{moves}(q)} \alpha(q)(\mathsf{m}) \cdot \Delta_{\alpha_{q,\mathsf{m}(q)}}(a)\right)$$

So by linearity, the vector space $\mathcal{V}$ is the smallest column-vector space such that
- $\mathbf{1} = (1\ 1 \cdots 1)^T \in \mathcal{V}$ and
- $M\boldsymbol{u} \in \mathcal{V}$, for all $\boldsymbol{u} \in \mathcal{V}$ and all $M \in \mathcal{M}$.

Define a finite set of labels $\mathsf{L}' = \{b_{\alpha,a} \mid \alpha \in \Sigma,\ a \in \mathsf{L}\}$, and for each $\alpha \in \Sigma$ and each $a \in \mathsf{L}$ a matrix

$$\Delta'(b_{\alpha,a}) = \frac{1}{|\Sigma|}\Delta_\alpha(a).$$

The matrix $\sum_{b \in \mathsf{L}'} \Delta'(b)$ is stochastic. Define the MCs $\mathcal{D}' = \langle Q, \mu_0^{\mathcal{D}}, \mathsf{L}', \delta' \rangle$ and $\mathcal{C}' = \langle Q, \mu_0^{\mathcal{C}}, \mathsf{L}', \delta' \rangle$ such that $\delta'$ induces the transition matrices $\Delta'(b)$ for all $b \in \mathsf{L}'$. The MCs $\mathcal{D}'$ and $\mathcal{C}'$ are computable in logarithmic space. Let $\mathcal{V}' \subseteq \mathbb{R}^Q$ be the smallest column-vector space such that

- $\mathbf{1} = (1\ 1 \cdots 1)^T \in \mathcal{V}$ and
- $\Delta'(b)\boldsymbol{u} \in \mathcal{V}$, for all $\boldsymbol{u} \in \mathcal{V}$ and all $b \in \mathsf{L}'$.

Since the matrices in $\mathcal{M}$ and the matrices $\Delta'(b)$ are scalar multiples of each other, we have $\mathcal{V} = \mathcal{V}'$. It holds:

$$
\begin{aligned}
\mathcal{D} \sqsubseteq \mathcal{C} \quad &\Longleftrightarrow \quad \mathcal{D} \sim \mathcal{C} \text{ in } \mathcal{D} && \text{Proposition 5.3} \\
&\Longleftrightarrow \quad \mu_0^{\mathcal{D}} \sim \mu_0^{\mathcal{C}} \text{ in } \mathcal{D} && \text{definition} \\
&\Longleftrightarrow \quad \forall\, \boldsymbol{u} \in \mathcal{V} : \mu_0^{\mathcal{D}}\boldsymbol{u} = \mu_0^{\mathcal{C}}\boldsymbol{u} && \text{Corollary 5.11} \\
&\Longleftrightarrow \quad \forall\, \boldsymbol{u} \in \mathcal{V}' : \mu_0^{\mathcal{D}}\boldsymbol{u} = \mu_0^{\mathcal{C}}\boldsymbol{u} && \mathcal{V} = \mathcal{V}' \\
&\Longleftrightarrow \quad \mu_0^{\mathcal{D}} \sim \mu_0^{\mathcal{C}} \text{ in } \mathcal{D}' && \text{Corollary 5.11} \\
&\Longleftrightarrow \quad \mathcal{D}' \sim \mathcal{C}' \text{ in } \mathcal{D}' && \text{definition} \\
&\Longleftrightarrow \quad \mathcal{D}' \sqsubseteq \mathcal{C}' && \text{Proposition 5.3}
\end{aligned}
$$

As mentioned in Section 2.2, deciding whether $\mathcal{D}' \sqsubseteq \mathcal{C}'$ holds amounts to the trace-equivalence problem for MCs. It follows from Tzeng [24] that the latter is decidable in $\mathsf{NC}$, hence in $\mathsf{P}$. $\quad\square$

## 6. Conclusions

We have settled the decidability and complexity status of most subproblems of trace refinement between two MDPs. Key technical ingredients were links to a certain notion of bisimulation, linear-algebra arguments, and comparisons of polytopes.

As an open problem, we highlight the complexity of the distribution-based notion of bisimulation, which we have shown to be in $\mathsf{coNP}$. Is this notion of bisimulation $\mathsf{coNP}$-complete or in $\mathsf{P}$?

## References

[1] S. Arora, R. Ge, R. Kannan, and A. Moitra. Computing a nonnegative matrix factorization - provably. In *STOC*, pages 145–162. ACM, 2012.

[2] V. D. Blondel and V. Canterini. Undecidable problems for probabilistic automata of fixed dimension. *Theoretical Computer Science*, 36 (3):231–245, 2003.

[3] J. Canny. Some algebraic and geometric computations in PSPACE. In *STOC*, pages 460–467, 1988.

[4] J. Cohen and U. Rothblum. Nonnegative ranks, decompositions, and factorizations of nonnegative matrices. *Linear Algebra and its Applications*, 190:149–168, 1993.

[5] A. Condon and R. Lipton. On the complexity of space bounded interactive proofs (extended abstract). In *FOCS*, pages 462–467, 1989.

[6] T. Cormen, C. Stein, R. Rivest, and C. E. Leiserson. *Introduction to algorithms*. McGraw-Hill Higher Education, 2nd edition, 2001.

[7] L. Doyen, T.A. Henzinger, and J.-F. Raskin. Equivalence of labeled Markov chains. *International Journal on Foundations of Computer Science*, 19(3):549–563, 2008.

[8] J. Fearnley and M. Jurdzinski. Reachability in two-clock timed automata is PSPACE-complete. *Information and Computation*, 243:26–36, 2015.

[9] N. Fijalkow. Undecidability results for probabilistic automata. *SIGLOG News*, 4(4):10–17, 2017.

[10] R. Greenlaw, H.J. Hoover, and W.L. Ruzzo. *Limits to parallel computation: P-completeness theory*. Oxford University Press, 1995.

[11] H. Hermanns, J. Krčál, and J. Křetínský. Probabilistic bisimulation: Naturally on distributions. In *CONCUR*, volume 8704 of *LNCS*, pages 249–265. Springer, 2014. Technical report at http://arxiv.org/abs/1404.5084.

[12] L. Khachiyan. A polynomial algorithm for linear programming. *Doklady Akademii Nauk SSSR.*, 224:1093–1096, 1979.

[13] S. Kiefer, A.S. Murawski, J. Ouaknine, B. Wachter, and J. Worrell. Language equivalence for probabilistic automata. In *CAV*, volume 6806 of *LNCS*, pages 526–540. Springer, 2011.

[14] S. Kiefer, A.S. Murawski, J. Ouaknine, B. Wachter, and J. Worrell. APEX: An analyzer for open probabilistic programs. In *CAV*, volume 7358 of *LNCS*, pages 693–698. Springer, 2012.

[15] S. Kiefer and B. Wachter. Stability and complexity of minimising probabilistic automata. In *ICALP*, volume 8573 of *LNCS*, pages 268–279, 2014.

[16] L. Li and Y. Feng. Quantum Markov chains: Description of hybrid systems, decidability of equivalence, and model checking linear-time properties. *Information and Computation*, 244:229–244, 2015.

[17] T.M. Ngo, M. Stoelinga, and M. Huisman. Confidentiality for probabilistic multi-threaded programs and its verification. In *Engineering Secure Software and Systems*, volume 7781 of *LNCS*, pages 107–122. Springer, 2013.

[18] A. Paz. *Introduction to Probabilistic Automata*. Academic Press, 1971.

[19] S. Peyronnet, M. de Rougemont, and Y. Strozecki. Approximate verification and enumeration problems. In *ICTAC*, volume 7521 of *LNCS*, pages 228–242. Springer, 2012.

[20] J. Renegar. On the computational complexity and geometry of the first-order theory of the reals. Parts I–III. *Journal of Symbolic Computation*, 13(3):255–352, 1992.

[21] M.-P. Schützenberger. On the definition of a family of automata. *Information and Control*, 4:245–270, 1961.

[22] Y. Shitov. A universality theorem for nonnegative matrix factorizations. 2016. Report at https://arxiv.org/abs/1606.09068.

[23] W. Tzeng. A polynomial-time algorithm for the equivalence of probabilistic automata. *SIAM Journal on Computing*, 21(2):216–227, 1992.

[24] W. Tzeng. On path equivalence of nondeterministic finite automata. *Information Processing Letters*, 58(1):43–46, 1996.

[25] S. Vavasis. On the complexity of nonnegative matrix factorization. *SIAM Journal on Optimization*, 20(3):1364–1377, 2009.