# Monitoring Algorithmic Fairness
# under Partial Observations

Thomas A. Henzinger, Konstantin Kueffner, and Kaushik Mallik

Institute of Science and Technology Austria (ISTA)

**Abstract.** As AI and machine-learned software are used increasingly for making decisions that affect humans, it is imperative that they remain fair and unbiased in their decisions. To complement design-time bias mitigation measures, runtime verification techniques have been introduced recently to monitor the algorithmic fairness of deployed systems. Previous monitoring techniques assume full observability of the states of the (unknown) monitored system. Moreover, they can monitor only fairness properties that are specified as arithmetic expressions over the probabilities of different events. In this work, we extend fairness monitoring to systems modeled as partially observed Markov chains (POMC), and to specifications containing arithmetic expressions over the expected values of numerical functions on event sequences. The only assumptions we make are that the underlying POMC is aperiodic and starts in the stationary distribution, with a bound on its mixing time being known. These assumptions enable us to estimate a given property for the entire distribution of possible executions of the monitored POMC, by observing only a single execution. Our monitors observe a long run of the system and, after each new observation, output updated PAC-estimates of how fair or biased the system is. The monitors are computationally lightweight and, using a prototype implementation, we demonstrate their effectiveness on several real-world examples.

## 1 Introduction

Runtime verification complements traditional static verification techniques, by offering lightweight approaches for verifying properties of systems from a single long observed execution trace [9]. Recently, runtime verification was used to monitor biases in machine-learned decision-making softwares [3,32,31]. Decision-making softwares are being increasingly used for making critical decisions affecting humans; example areas include judiciary [14,19], policing [21,46], and banking [45]. It is important that these softwares are unbiased towards the protected attributes of humans, like gender and ethnicity. However, they were shown to be biased on many occasions in the past [19,43,50,54,55]. While many offline approaches were proposed for mitigating such biases [12,61,13,63,35,40], runtime verification introduces a new complementary tool to oversee *algorithmic fairness* of deployed decision-making systems [3,32,31]. In this work, we extend runtime verification to monitor algorithmic fairness for a broader class of system models and a more expressive specification language.

Prior works on monitoring algorithmic fairness assumed that the given system is modeled as a Markov chain with unknown transition probabilities but with fully observable states [3,32]. A sequence of states visited by the Markov chain represents a (randomized) sequence of events generated from the interaction of the decision-making agent and its environment. The goal is to design a monitor that will observe one such long sequence of states, and, after observing every new state in the sequence, will compute an updated PAC-estimate of how fair or biased the system is.

In the prior works, the PAC guarantee on the output hinges on the full observability and the Markovian structure of the system [3,32,31]. While this setup is foundational, it is also very basic, and is not fulfilled by many real-world examples. Consider a lending scenario where at every step a bank (the decision-maker) receives the features (e.g., the age, gender, and ethnicity) of a loan applicant, and decides whether to grant or reject the loan. To model this system using the existing setup, we would need to assume that the monitor can observe the full state of the system which includes all the features of every applicant. In reality, the monitor will often be a third-party system, having only partial view of the system's states.

We address the problem of designing monitors when the systems are modeled using partially observed Markov chains (POMC) with unknown transition probabilities. The difficulty comes from the fact that a random observation sequence that is visible to the monitor may not follow a Markovian pattern, even though the underlying state sequence is Markovian. We overcome this by making the assumption that the POMC starts in the stationary distribution, which in turn guarantees a certain uniformity in how the observations follow each other. We argue that the stationarity assumption is fulfilled whenever the system has been running for a long time, which is suitable for long term monitoring of fairness properties. With the help of a few additional standard assumptions on the POMC, like aperiodicity and the knowledge of a bound on the mixing time, we can compute PAC estimates on the degree of algorithmic fairness over the distribution of all runs of the system from a single monitored observation sequence.

Besides the new system model, we also introduce a richer specification language— called bounded specification expressions (BSE). BSE-s can express many common algorithmic fairness properties from the literature, such as demographic parity [20], equal opportunity [30], and disparate impact [24]. Furthermore, BSE-s can express new fairness properties which were not expressible using the previous formalism [3,32]. In particular, BSE-s can express quantitative fairness properties, including fair distribution of expected credit scores and fair distribution of expected wages across different demographic groups of the population; details can be found, respectively, in Ex. 5 and 6 in Sec. 3.2.

The building block of a BSE is an atomic function, which is a function that assigns bounded numerical values to observation sequences of a particular length. Using an atomic function, we can express weighted star-free regular expressions (every word satisfying the given regular expression has a numerical weight), average response time-like properties, etc. A BSE can contain many different atomic

functions combined together through a restricted set of arithmetic, relational, and logical operations. We define two fragments of BSE-s: The first one is called QuantBSE, which contains only arithmetic expressions over atomic functions, and whose semantic value is the expected value of the given expression over the distribution of runs of the POMC. The second one is called QualBSE, which turns the QuantBSE expressions into boolean expressions through relational (e.g., whether a QuantBSE expression is greater than zero) and logical operators (e.g., conjunction of two relational sentences), and whose semantic value is the expected truth or falsehood of the given expression over the distribution of runs of the POMC.

For any given BSE, we show how to construct a monitor that observes a single long observation sequence generated by the given POMC with unknown transition probabilities, and after each observation outputs an updated numerical estimate of the actual semantic value of the BSE for the observed system. The heart of our approach is a PAC estimation algorithm for the semantic values of the atomic functions. The main difficulty stems from the statistical dependence between any two consecutive observations, which is a side-effect of the partial observability of the states of the Markov chain, and prevents us from using the common PAC bounds that were used in the prior works that assumed full observability of the POMC states [3,32]. We show how the problem can be cast as the statistical estimation problem of estimating the expected value of a function over the states of a POMC which satisfies a certain bounded difference property. This estimation problem can be solved using a version of McDiarmid's concentration inequality [52], for which we need the additional assumptions that the given POMC is aperiodic and that a bound on its mixing time is known. We use McDiarmid's inequality to find the PAC estimate of every individual atomic function of the given BSE. The individual PAC estimates can then be combined using known methods to obtain the overall PAC estimate of the given BSE [3].

Our monitors are computationally lightweight, and produce reasonably tight PAC bounds of the monitored properties. Using a prototype implementation, we present the effectiveness of our monitors on two different examples. On a real-world example, we showed how our monitors can check if a bank has been fair in giving loans to individuals from two different demographic groups in the population, and on an academic example, we showed how our monitors' outputs improve as the known bound on the mixing time gets tighter.

The proofs of the technical claims can be found in the appendices.

## 1.1 Related Work

There are many works in AI and machine-learning which address how to eliminate or minimize decision biases in learned models through improved design principles [48,20,30,42,39,56,12,61,13,63,35,40]. In formal methods, too, there are some works which statically verify absence of biases of learned models [2,11,59,28,49,37,7,29]. All of these works are static interventions and rely on the availability of the system model, which may not be always true.

Runtime verification of algorithmic fairness, through continuous monitoring of decision events, is a relatively new area pioneered by the work of Albarghouthi et al. [3]. We further advanced their idea in our other works which appeared recently [32,31]. In those works, on one hand, we generalized the class of supported system models to Markov chains and presented the new Bayesian statistical view of the problem [32]. On the other hand, we relaxed the time-invariance assumption on the system [31]. In this current paper, we limit ourselves to time-invariant systems but extend the system models to partially observed Markov chains and consider the broader class of BSE properties, which enables us to additionally express properties whose values depend on observation sequences.

Traditional runtime verification techniques support mainly temporal properties and employ finite automata-based monitors [57,38,23,47,18,8,5]. In contrast, runtime verification of algorithmic fairness requires checking statistical properties, which is beyond the limit of what automata-based monitors can accomplish. Although there are some works on quantitative runtime verification using richer types of monitors (with counters/registers like us) [27,33,51,34], the considered specifications usually do not extend to statistical properties such as algorithmic fairness.

Among the few works on monitoring statistical properties of systems, a majority of them only provides asymptotic correctness guarantees [25,60], whereas we provide anytime guarantees. On the other hand, works on monitoring statistical properties with finite-sample (nonasymptotic) guarantees are rare and are restricted to simple properties, such as probabilities of occurrences of certain events [10] and properties specified using certain fragments of LTL [53]. Monitoring POMCs (the same modeling formalism as us) were studied before by Stoller et al. [58], though the setting was a bit different from ours. Firstly, they only consider LTL properties, and, secondly, they assume the system model to be known by the monitor. This way the task of the monitor effectively reduces to a state estimation problem from a given observation sequence.

Technique-wise, there are some similarities between our work and the works on statistical model-checking [4,62,15,17,1] in that both compute PAC-guarantees on satisfaction or violation of a given specification. However, to the best of our knowledge, the existing statistical model-checking approaches do not consider algorithmic fairness properties.

## 2 Preliminaries

### 2.1 Notation

We write $\mathbb{R}$, $\mathbb{R}^+$, $\mathbb{N}$, and $\mathbb{N}^+$ to denote the sets of real numbers, positive real numbers, natural numbers (including zero), and positive integers, respectively.

Let $\Sigma$ be a countable alphabet. We write $\Sigma^*$ and $\Sigma^\omega$ to denote, respectively, the set of every finite and infinite word over $\Sigma$. Moreover, $\Sigma^\infty$ denotes the set of finite and infinite words, i.e., $\Sigma^\infty := \Sigma^* \cup \Sigma^\omega$. We use the convention that symbols with arrow on top will denote words, whereas symbols without arrow

will denote alphabet elements. Let $\vec{s} = s_1 s_2 \ldots$ be a word. We write $\vec{s}_i$ to denote the $i$-th symbol $s_i$, and write $\vec{s}_{i..j}$ to denote the subword $s_i \ldots s_j$, for $i < j$. We use the convention that the indices of a word begin at 1, so that the length of a word matches the index of the last symbol.

Let $\vec{s} \in \Sigma^*$ and any $\vec{t} \in \Sigma^\infty$ be two words. We denote the concatenation of $\vec{s}$ and $\vec{t}$ as $\vec{s}\,\vec{t}$. We generalize this to sets of words: For $S \subseteq \Sigma^*$ and $T \subseteq \Sigma^\infty$, we define the concatenation $ST := \{\vec{s}\,\vec{t} \mid \vec{s} \in S, \vec{t} \in T\}$. We say $\vec{s}$ is a prefix of $\vec{r}$, written $\vec{s} \prec \vec{r}$, if there exists a word $\vec{t} \in \Sigma^\infty$ such that $\vec{s}\,\vec{t} = \vec{r}$.

Suppose $\mathbb{T} \subseteq \mathbb{R}$ is a subset of real numbers, $v \in \mathbb{T}^n$ is a vector of length $n$ over $\mathbb{T}$, and $M \in \mathbb{T}^{n \times m}$ is a matrix of dimension $n \times m$ over $\mathbb{T}$; here $m, n$ can be infinity. We use $v_i$ to denote the $i$-th element of $v$, and $M_{ij}$ to denote the element at the intersection of the $i$-th row and the $j$-th column of $M$. A probability distribution over a set $S$ is a vector $v \in [0,1]^{|S|}$, such that $\sum_{i \in [1;|S|]} v_i = 1$.

## 2.2 Randomized Event Generators: Partially Observed Markov Chains

We use partially observed Markov chains (POMC) as sequential randomized generators of events. A POMC is a tuple $(Q, M, \lambda, \Sigma, \ell)$, where $Q = \mathbb{N}^+$ is a countable set of states, $M$ is a stochastic matrix of dimension $|Q| \times |Q|$, called the *transition probability matrix*, $\lambda$ is a probability distribution over $Q$ representing the *initial state distribution*, $\Sigma$ is a countable set of *observations*, and $\ell \colon Q \to \Sigma$ is a function mapping every state to an *observation*. All POMCs in this paper are time-homogeneous, i.e., their transition probabilities do not vary over time.

Semantically, every POMC $\mathcal{M}$ induces a probability measure $\mathbb{P}_{\mathcal{M}}(\cdot)$ over the generated state and observation sequences. For every finite state sequence $\vec{q} = q_1 q_2 \ldots q_t \in Q^*$, the probability that $\vec{q}$ is generated by $\mathcal{M}$ is given by $\mathbb{P}_{\mathcal{M}}(\vec{q}) = \lambda_{q_1} \cdot \prod_{i=1}^{t-1} M_{q_i q_{i+1}}$. Every finite state sequence $\vec{q} \in Q^*$ for which $\mathbb{P}_{\mathcal{M}}(\vec{q}) > 0$ is called a finite *internal path* of $\mathcal{M}$; we omit $\mathcal{M}$ if it is clear from the context. The set of every internal path of length $n$ is denoted as $Q^n(\mathcal{M})$, and the set of every finite internal path is denoted as $Q^*(\mathcal{M})$.

Every finite internal path $\vec{q}$ can be extended to a set of infinite internal paths, which is called the cylinder set induced by $\vec{q}$, and is defined as $Cyl(\vec{q}) := \{\vec{r} \in Q^\omega \mid \vec{q} \prec \vec{r}\}$. The probability measure $\mathbb{P}_{\mathcal{M}}(\cdot)$ on finite internal paths induces a pre-measure on the respective cylinder sets, which can be extended to a unique measure on the infinite internal paths by means of the Carathéodory's extension theorem [6, pp. 757]. The probability measure on the set of infinite internal paths is also denoted using $\mathbb{P}_{\mathcal{M}}(\cdot)$.

An external observer can only observe the observable part of an internal path of a POMC. Given an internal path $\vec{q} = q_1 q_2 \ldots \in Q^\infty$, we write $\ell(\vec{q})$ to denote the observation sequence $\ell(q_1)\ell(q_2)\ldots \in \Sigma^\infty$. For a set of internal paths $S \subseteq Q^\infty$, we write $\ell(S)$ to denote the respective set of observation sequences $\{\vec{w} \in \Sigma^\infty \mid \exists \vec{q} \,.\, \vec{w} = \ell(\vec{q})\}$. An observation sequence $\vec{w} \in \Sigma^\infty$ is called an *observed* path (of $\mathcal{M}$) if there exists an internal path $\vec{q}$ for which $\ell(\vec{q}) = \vec{w}$. As before, we write $\Sigma^n(\mathcal{M})$ for the set of every observed path of length $n$, and $\Sigma^*(\mathcal{M})$ for the set of every finite observed path.

We also use the inverse operator of $\ell$ to map every observed path $\vec{w}$ to the set of possible internal paths: $\ell^{-1}(\vec{w}) := \{\vec{q} \in Q^\infty \mid \ell(\vec{q}) = \vec{w}\}$. Furthermore, we extend $\ell^{-1}(\cdot)$ to operate over sets of observation sequences in the following way: For any given $S \subseteq \Sigma^\infty$, define $\ell^{-1}(S) := \{\vec{q} \mid \exists \vec{w} \in S . \ell(\vec{q}) = \vec{w}\}$.

We abuse the notation and use $\mathbb{P}_{\mathcal{M}}(\cdot)$ to denote the induced probability measure on the set of observed paths, defined in the following way. Given every set of finite observed paths $S \subseteq \Sigma^*$, we define $\mathbb{P}_{\mathcal{M}}(S) := \sum_{\vec{q} \in \ell^{-1}(S)} \mathbb{P}_{\mathcal{M}}(\vec{q})$. When the paths in a given set are infinite, the sum is replaced by integral. We write $\vec{W} \sim \mathcal{M}$ to denote the random variable that represents the distribution over finite sample observed paths generated by the POMC $\mathcal{M}$.

*Example 1.* As a running example, we introduce a POMC that models the sequential interaction between a bank and loan applicants. Suppose there is a population of loan applicants, where each applicant has a credit score between 1 and 4, and belongs to either an advantaged group $A$ or a disadvantaged group $B$. At every step, the bank receives loan application from one applicant, and, based on some unknown (but non-time-varying) criteria, decides whether to grant loan or reject the application. We want to monitor, for example, the difference between loan acceptance probabilities for people belonging to the two groups.

The underlying POMC $\mathcal{M}$ that models the sequence of loan application events is shown in Fig. 1. A possible internal path is $S(A,1)NS(A,4)YSB(A,3)N\ldots$, whose corresponding observed path is $SANSAYSBN\ldots$. In our experiments, we use a more realistic model of the POMC with way more richer set of features for the individuals.
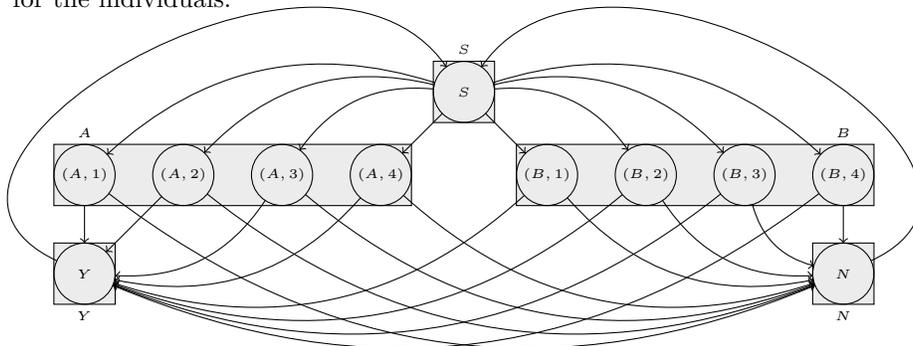


Fig. 1: The POMC modeling the sequential interaction between the bank and the loan applicants. The states $S$, $Y$, and $N$ respectively denote the start state, the event that the loan was granted ("$Y$" stands for "Yes"), and the event that the loan was rejected ("$N$" stands for "No"). Every middle state $(X,i)$, for $X \in \{A,B\}$ and $i \in \{1,2,3,4\}$, represents the group ($A$ or $B$) and the credit score $i$ of the current applicant. The states $S, Y, N$ are fully observable, i.e., their observation symbols coincide with their state symbols. The middle states are partially observable, with every $(A,i)$ being assigned the observation $A$ and every $(B,i)$ being assigned the observation $B$. The states with the same observation belong to the same shaded box.

### 2.3 Register Monitors

Our register monitors are adapted from the polynomial monitors of Ferrère et al. [26], and were also used in our previous work (in a more general randomized form) [32]. Let $R$ be a finite set of integer variables called registers. A function $v \colon R \to \mathbb{N}$ assigning concrete value to every register in $R$ is called a valuation of $R$. Let $\mathbb{N}^R$ denote the set of all valuations of $R$. Registers can be read and written according to relations in the signature $S = \langle 0, 1, +, -, \times, \div, \leq \rangle$. We consider two basic operations on registers:

  – A *test* is a conjunction of atomic formulas over $S$ and their negation;
  – An *update* is a mapping from variables to terms over $S$.

We use $\Phi(R)$ and $\Gamma(R)$ to respectively denote the set of tests and updates over $R$. *Counters* are special registers with a restricted signature $S = \langle 0, 1, +, -, \leq \rangle$.

**Definition 1 (Register monitor).** *A register monitor is a tuple $(\Sigma, \Lambda, R, v_{\mathsf{in}}, f, T)$ where $\Sigma$ is a finite input alphabet, $\Lambda$ is an output alphabet, $R$ is a finite set of registers, $v_{\mathsf{in}} \in \mathbb{N}^R$ is the initial valuation of the registers, $f \colon \mathbb{N}^R \to \Lambda$ is an output function, and $T \colon \Sigma \times \Phi(R) \to \Gamma(R)$ is the transition function such that for every $\sigma \in \Sigma$ and for every valuation $v \in \mathbb{N}^R$, there exists a unique $\phi \in \Phi(R)$ with $v \models \phi$ and $T(\sigma, \phi) \in \Gamma(R)$.*

We refer to register monitors simply as monitors, and we fix the output alphabet $\Gamma$ as the set of every real interval.

A *state* of a monitor $\mathcal{A}$ is a valuation of its registers $v \in \mathbb{N}^R$; the initial valuation $v_{\mathsf{in}}$ is the initial state. The monitor $\mathcal{A}$ *transitions* from state $v$ to another state $v'$ on input $\sigma \in \Sigma$ if there exists $\phi$ such that $v \models \phi$, there exists an update $\gamma = T(\sigma, \phi)$, and if $v'$ maps every register $x$ to $v'(x) = v(\gamma(x))$. The transition from $v$ to $v'$ on input $\sigma$ is written as $v \xrightarrow{\sigma} v'$. A *run* of $\mathcal{A}$ on a word $w_1 \ldots w_t \in \Sigma^*$ is a sequence of transitions $v_1 = v_{\mathsf{in}} \xrightarrow{w_1} v_2 \xrightarrow{w_2} \ldots \xrightarrow{w_t} v_{t+1}$. The *semantics* of the monitor is the function $[\![\mathcal{A}]\!] \colon \Sigma^* \to \Lambda$ that maps every finite input word to the last output of the monitor on the respective run. For instance, the semantics of $\mathcal{A}$ on the word $\vec{w}$ is $[\![\mathcal{A}]\!](\vec{w}) = f(v_{t+1})$. An illustrative example of register monitors can be found in our earlier work [32, Sec. 2.2].

## 3 Monitoring Quantitative Algorithmic Fairness Properties

In our prior work on monitoring algorithmic fairness for *fully* observable Markov chains [32], we formalized (quantitative) algorithmic fairness properties using the so-called Probabilistic Specification Expressions (PSE). A PSE $\varphi$ is an arithmetic expression over the variables of the form $v_{ij}$, for $i, j \in Q$ for a finite set $Q$. The semantics of $\varphi$ is interpreted statically over a given Markov chain $M$ with state space $Q$, by replacing every $v_{ij}$ with the transition probability from the state $i$ to the state $j$ in $M$. The algorithmic question we considered is that given a PSE $\varphi$, how to construct a monitor that will observe one long path of an unknown

Markov chain, and after each observation will output a PAC estimate of the value of $\varphi$ with a pre-specified confidence level.

An exact representation of the above problem formulation is not obvious for POMCs. In particular, while it is reasonable to generalize the semantics of PSEs to be over the probabilities between *observations* instead of probabilities between states, it is unclear how these probabilities will be defined. In the following, we use simple examples to illustrate several cruxes of formalizing algorithmic fairness on POMCs, and motivate the use of the assumptions of stationary distribution, irreducibility, and positive recurrence (formally stated in Assump. 1) to mitigate the difficulties. These assumptions will later be used to formalize the algorithmic fairness properties in Sec. 3.2.

In the following, we will write $\pi$ to denote the *stationary distribution* of Markov chains with transition matrix $M$, i.e., $\pi = M\pi$.

## 3.1   Role of the Stationary Distribution

First, we demonstrate in the following example that POMCs made up of unfair sub-components may have overall fair behavior in the stationary distribution, which does not happen for fully observable Markov chains.

*Example 2.* Suppose there are two coins $A$ and $B$, where $A$ comes up with heads with probability 0.9 and $B$ comes up with tails with probability 0.9. We observe a sequence of coin tosses (i.e., the observations are heads and tails), without knowing which of the two coins (the state) was tossed. If the choice of the coin at each step is made uniformly at random, then, intuitively, the system will produce fair outcomes in the long run, with equal proportions of heads and tails being observed in expectation. Thus, although each coin was unfair, we can still observe overall fair outcome, provided the fraction of times each coin was chosen in the stationary distribution balances out the unfairness in the coins themselves.

To make the above situation more concrete, imagine that the underlying POMC has two states $a, b$ (e.g., $a, b$ represent the states when $A, B$ are selected for tossing, respectively) with the same observation (which coin is selected is unknown to the observer), where the measures of the given fairness condition (e.g., the biases of the coins $A, B$) are given by $f_a, f_b$. We argue that, intuitively, the overall fairness of the POMC is given by $\pi_a f_a + \pi_b f_b$. This type of analysis is unique to POMCs, whereas for fully observable Markov chains, computation of fairness is simpler and can be done without involving the stationary distribution.

In the next example, we demonstrate some challenges of monitoring fairness when we express fairness by weighing in the stationary distribution as above.

*Example 3.* Consider the setting of Ex. 2, and suppose now only the initial selection of the coin happens uniformly at random but subsequently the same coin is used forever. If we consider the underlying POMC, both $\pi_a, \pi_b$ will be 0.5, because the initial selection of the coin happens uniformly at random. However, the monitor will observe the toss outcomes of only one of the two coins on a given trace. It is unclear how the monitor can extrapolate its estimate to the overall fairness property $\pi_a f_a + \pi_b f_b$ in this case.

To deal with the situations described in Ex. 2 and Ex. 3, we will make the following assumption.

**Assumption 1** *We assume that the POMCs are irreducible, positively recurrent, and are initialized in their stationary distributions.*

The irreducibility and positive recurrence guarantees existence of the stationary distribution. Assump. 1 ensures that, firstly, we will see every state infinitely many times (ruling out the above corner-case), and, secondly, the proportion of times the POMC will spend in all the states will be the same (given by the stationary distribution) all the time. While Assump. 1 makes it easier to formulate and analyze the algorithmic fairness properties over POMCs, monitoring these properties over POMCs still remains a challenging problem due to the non-Markovian nature of the observed path.

### 3.2 Bounded Specification Expressions

We introduce bounded specification expressions (BSE) to formalize the fairness properties that we want to monitor. A BSE assigns values to finite word patterns of a given alphabet. The main components of a BSE are *atomic functions*, where an atomic function $f_n$ assigns bounded real values to observation sequences of length $n$, for a given $n \in \mathbb{N}^+$. An atomic function $f_n$ can express quantitative star-free regular expressions, assigning real values to words of length $n$.

Following are some examples. Let $\Sigma = \{r, g\}$ be an observation alphabet, where $r$ stands for "request" and $g$ stands for "grant." A boolean atomic function $f_2$, with $f_2(rr) = 0$ and $f_2(\vec{w}) = 1$ for every $\vec{w} \in \Sigma^2 \setminus \{rr\}$, can express the property that two requests should not appear consecutively. An integer-valued atomic function $f_{10}$, with $f_{10}(rr^i g \vec{w}) = i$ when $i \in [0; 8]$ and $\vec{w} \in \Sigma^{8-i}$, and with $f_{10}(\vec{z}) = 8$ when $\vec{z} \in \Sigma^{10} \setminus rr^i g \Sigma^{8-i}$, assigns to any sub-sequence the total waiting time between a request and the subsequent grant, while saturating the waiting time to 8 when it is above 8. The specified word-length $n$ for any atomic function $f_n$ is called the *arity* of $f_n$. Let $P$ be the set of all atomic functions over a given observation alphabet.

A BSE may also contain arithmetic and/or logical connectives and relational operators to express complex value-based properties of an underlying probabilistic generator, like the POMCs. We consider two fragments of BSE-s, expressing qualitative and quantitative properties, and called, respectively, QualBSE and QuantBSE in short. The syntaxes of the two types of BSE-s are given as:

$$\text{(QuantBSE)} \qquad \varphi ::= \kappa \in \mathbb{R} \mid f \in P \mid \varphi + \varphi \mid \varphi \cdot \varphi \mid 1 \div \varphi \mid (\varphi), \qquad \text{(1a)}$$

$$\text{(QualBSE)} \qquad \psi ::= \textit{true} \mid \varphi \geq 0 \mid \neg \psi \mid \psi \wedge \psi. \qquad \text{(1b)}$$

The semantics of a QuantBSE $\varphi$ over the alphabet $\Sigma$ is interpreted over POMCs satisfying Assump. 1 and with observations $\Sigma$. When $\varphi$ is an atomic function $f \colon \Sigma^n \to [a, b]$ for some $n \in \mathbb{N}^+$, $a, b \in \mathbb{R}$, then, for a given POMC $\mathcal{M}$,

the semantics of $\varphi$ is defined as follows. For every time $t \in \mathbb{N}^+$,

$$\varphi(\mathcal{M}) = f(\mathcal{M}) := \int_{\Sigma^\omega} f(\vec{w}_{t:t+n-1}) d\mathbb{P}_\mathcal{M}(\vec{w}). \tag{2}$$

The definition of $f(\mathcal{M})$ is well-defined, because $f(\mathcal{M})$ will be the same for every $t$, since the POMC will remain in the stationary distribution forever (by Assump. 1 and by the property of stationary distributions). Intuitively, the semantics $f(\mathcal{M})$ represents the expected value of the function $f$ on any sub-word of length $n$ on any observed path of the POMC, when it is known that the POMC is in the stationary distribution (Assump. 1).

The arithmetic operators in QuantBSE-s have the usual semantics ("+" for addition, "−" for difference, "·" for multiplication, and "÷" for division).

On the other hand, the semantics of a QualBSE $\psi$ is boolean, which inductively uses the semantics of the constituent $\varphi$ expressions. For a QualBSE $\psi = \varphi \geq 0$, the semantics of $\psi$ is given by:

$$\psi(\mathcal{M}) := \begin{cases} true & \text{if } \varphi(\mathcal{M}) \geq 0, \\ false & \text{otherwise.} \end{cases}$$

The semantics of the boolean operators in $\psi$ is the usual semantics of boolean operators in propositional logic. The following can be added as syntactic sugar: "$\varphi \geq c$" for a constant $c$ denotes "$\varphi' \geq 0$" with $\varphi' := \varphi - c$, "$\varphi \leq c$" denotes "$-\varphi \geq -c$," "$\varphi = c$" denotes "$(\varphi \geq c) \wedge (\varphi \leq c)$," "$\varphi > c$" denotes "$\neg(\varphi \leq c)$," "$\varphi < c$" denotes "$\neg(\varphi \geq c)$," and "$\psi \vee \psi$" denotes "$\neg(\neg\psi \wedge \neg\psi)$."

**Fragment of BSE: Probabilistic Specification Expressions (PSEs):** In our prior work [32], we introduced PSEs to model algorithmic fairness properties of Markov chains with fully observable state space. PSEs are arithmetic expressions over atomic variables of the form $v_{ij}$, where $i, j$ are the states of the given Markov chain, and whose semantic value equals the transition probability from $i$ to $j$. The semantics of a PSE is then the valuation of the expression obtained by plugging in the respective transition probabilities. We can express PSEs using QuantBSE-s as below. For every variable $v_{ij}$ appearing in a given PSE, we use the atomic function $f$ that assigns to every finite word $\vec{w} \in \Sigma^*$ the ratio of the number of $(i, j)$ transitions to the number of occurrences of $i$ in $\vec{w}$. We will denote the function $f$ as $P(j \mid i)$ in this case, and, in general, $i, j$ can be observation labels for the case of QuantBSE-s. It is straightforward to show that semantically the two expressions will be the same. On the other hand, QuantBSE-s are strictly more expressive than PSEs. For instance, unlike PSEs, QuantBSE-s can specify probability of transitioning from one observation label to another, the average number of times a given state is visited on any finite path of a Markov chain, etc.

**Fragment of BSE: Probabilities of Sequences:** We consider a useful fragment that expresses the probability that a sequence from a given set $S \subseteq \Sigma^*$ of finite observation sequences will be observed at any point in time on any observed

path. We assume that the length of every sequence in $S$ is uniformly bounded by some integer $n$. Let $\overline{S} \subseteq \Sigma^n$ denote the set of extensions of sequences in $S$ up to length $n$, i.e., $\overline{S} := \{\vec{w} \in \Sigma^n \mid \exists \vec{u} \in S \, . \, \vec{u} \prec \vec{w}\}$. Then the desired property will be expressed simply using an atomic function with $f \colon \Sigma^n \to \{0,1\}$ being the indicator function of the set $\overline{S}$, i.e., $f(\vec{w}) = 1$ iff $\vec{w} \in \overline{S}$. It is straightforward to show that, for a given POMC $\mathcal{M}$, the semantics $f(\mathcal{M})$ expresses the desired property. For a set of finite words $S \subseteq \Sigma^*$, we introduce the shorthand notation $P(S)$ to denote the probability of seeing an observation from the set $S$ at any given point in time. Furthermore, for a pair of sets of finite words $S, T \subseteq \Sigma^*$, we use the shorthand notation $P(S \mid T)$ to denote $P(TS)/P(T)$, which represents the conditional probability of seeing a word in $S$ after we have seen a word in $T$.

*Example 4 (Group fairness.).* Consider the setting in Ex. 1. We show how we can represent various group fairness properties using QuantBSE-s. Demographic parity [20] quantifies bias as the difference between the probabilities of individuals from the two demographic groups getting the loan, which can be expressed as $P(Y \mid A) - P(Y \mid B)$. Disparate impact [24] quantifies bias as the ratio between the probabilities of getting the loan across the two demographic groups, which can be expressed as $P(Y \mid A) \div P(Y \mid B)$.

In prior works [3,32], group fairness properties could be expressed on strictly less richer class of fully observed Markov chain models, where the features of each individual were required to contain only their group information. An extension to the model of Ex. 1 is not straightforward as the confidence interval used in these works would not be applicable.

*Example 5 (Social fairness.).* Consider the setting in Ex. 1, except that now the credit score of each individual will be observable along with their group memberships, i.e., each observation is a pair of the form $(X, i)$ with $X \in \{A, B\}$ and $i \in \{1, 2, 3, 4\}$. There may be other non-sensitive features, such as age, which may be hidden. We use the social fairness property [31] quantified as the difference between the expected credit scores of the groups $A$ and $B$. To express this property, we use the unary atomic functions $f_1^X \colon \Sigma \to \mathbb{N}$, for $X \in \{A, B\}$, such that $f_1^X \colon (Y, i) \mapsto i$ if $Y = X$ and is 0 otherwise. The semantics of $f_1^X$ is the expected credit score of group $X$ scaled by the probability of seeing an individual from group $X$. Then social fairness is given by the QuantBSE $\varphi = \frac{f_1^A}{P(A)} - \frac{f_1^B}{P(B)}$.

*Example 6 (Quantitative group fairness.).* Consider a sequential hiring scenario where at each step the salary and a sensitive feature (like gender) of a new recruit are observed. We denote the pair of observations as $(X, i)$, where $X \in \{A, B\}$ represents the group information based on the sensitive feature and $i$ represents the salary. We can express the disparity in expected salary of the two groups in a similar manner as in Ex. 5. Define the unary functions $f_1^X \colon \Sigma \to \mathbb{N}$, for $X \in \{A, B\}$, such that $f_1^X \colon (Y, i) \mapsto i$ if $Y = X$ and is 0 otherwise. The semantics of $f_1^X$ is the expected salary of group $X$ scaled by the probability of seeing an individual from group $X$. Then the group fairness property is given by the QuantBSE $\varphi = \frac{f_1^A}{P(A)} - \frac{f_1^B}{P(B)}$.

### 3.3 Problem Statement

Informally, our goal is to build monitors that will observe randomly generated observed paths of increasing length from a given unknown POMC, and, after each observation, will generate an updated estimate of how fair or biased the system was until the current time. Since the monitor's estimate is based on statistics collected from a finite path, the output may be incorrect with some probability. That is, the source of randomness is from the fact that the prefix is a finite sample of the fixed but unknown POMC.

For a given $\delta \in (0,1)$, and a given BSE $\varphi$, we define a *problem instance* as the tuple $(\varphi, \delta)$.

**Problem 1 (Monitoring QuantBSE-s)** *Suppose $(\varphi, \delta)$ is a problem instance where $\varphi$ is a QuantBSE. Design a monitor $\mathcal{A}$, with output alphabet $\{[l, u] \mid l, u \in \mathbb{R} \ . \ l < u\}$, such that for every POMC $\mathcal{M}$ satisfying Assump. 1, we have:*

$$\mathbb{P}_{\overrightarrow{W} \sim \mathcal{M}} \left( \varphi(\mathcal{M}) \in [\![\mathcal{A}]\!](\overrightarrow{W}) \right) \geq 1 - \delta. \tag{3}$$

The estimate $[l, u] = [\![\mathcal{A}]\!](\overrightarrow{w})$ is called the $(1 - \delta) \cdot 100\%$ *confidence interval* for $\varphi(M)$. The radius, given by $\varepsilon = 0.5 \cdot (u - l)$, is called the *estimation error*, the quantity $\delta$ is called the *failure probability*, and the quantity $1 - \delta$ is called the *confidence*. Intuitively, the monitor outputs the estimated confidence interval that contains the range of values within which the true semantic value of $\varphi$ falls with $(1 - \delta) \cdot 100\%$ probability. The estimate gets more precise as the error gets smaller, and the confidence gets higher. We will prefer the monitor with the maximum possible precision, i.e., having the least estimation error for a given $\delta$.

**Problem 2 (Monitoring QualBSE-s)** *Suppose $(\varphi, \delta)$ is a problem instance where $\varphi$ is a QualBSE. Design a monitor $\mathcal{A}$, with output alphabet $\{true, false\}$, such that for every POMC $\mathcal{M}$ satisfying Assump. 1, we have:*

$$\mathbb{P}_{\overrightarrow{W} \sim \mathcal{M}} \left( \psi(\mathcal{M}) \mid [\![\mathcal{A}]\!](\overrightarrow{W}) = true \right) \geq 1 - \delta, \tag{4}$$

$$\mathbb{P}_{\overrightarrow{W} \sim \mathcal{M}} \left( \neg\psi(\mathcal{M}) \mid [\![\mathcal{A}]\!](\overrightarrow{W}) = false \right) \geq 1 - \delta. \tag{5}$$

Unlike Prob. 1, the monitors addressing Prob. 2 do not output an interval but output a boolean verdict. Intuitively, the output of the monitor for Prob. 2 is either *true* or *false*, and it is required that the semantic value of the property $\psi$ is, respectively, *true* or *false* with $(1 - \delta) \cdot 100\%$ probability.

## 4 Construction of the Monitor

Our overall approach in this work is similar to the prior works [3,32,31]: We first compute a point estimate of the given BSE from the finite observation sequence of the POMC, and then compute an interval estimate through known concentration inequalities. However, the same concentration inequalities as the

prior works cannot be applied, because they required two successive observed events be independent, which is not true for POMCs. For instance, in Ex. 3, if we start the sequence of tosses by first tossing coin $A$, then we know that the subsequent tosses are going to be done using $A$ only, thereby implying that the outcomes of the future tosses will be statistically dependent on the initial random process that chooses between the two coins at the first step.

We present a novel theory of monitors for BSE-s on POMCs satisfying Assump. 1, using McDiarmid-style concentration inequalities for hidden Markov chains. In Sec. 4.1 and 4.2, we first present, respectively, the point estimator and the monitor for an individual atom. In Sec. 4.3, we build the overall monitor by combining the interval estimates of the individual atoms through interval arithmetic and union bound.

## 4.1  A Point Estimator for the Atoms

Consider a BSE atom $f$. We present a point estimator for $f$, which computes an estimated value of $f(\mathcal{M})$ from a finite observed path $\vec{w} \in \Sigma^t$, of an arbitrary length $t$, of the unknown POMC $\mathcal{M}$. The point estimator $\hat{f}(\cdot)$ is given as:

$$\hat{f}(\vec{w}) \coloneqq \frac{1}{t-n+1} \sum_{i=1}^{t-n+1} f(\vec{w}_{i..i+n-1}). \tag{6}$$

In the following proposition, we establish the unbiasedness of the estimator $\hat{f}(\cdot)$, a desirable property that says that the expected value of the estimator's output will coincide with the true value of the property that is being estimated.

**Proposition 1.** *Let $\mathcal{M}$ be a POMC satisfying Assump. 1, $f \colon \Sigma^n \to [a,b]$ be a function for fixed $n$, $a$, and $b$, and $\vec{W} \sim \mathcal{M}$ be a random observed path of an arbitrary length $|\vec{W}| = t > n$. Then $\mathbb{E}(\hat{f}(\vec{W})) = f(\mathcal{M})$.*

The following corollary establishes the counterpart of Prop. 1 for the fragment of BSE with probabilities of sequences.

**Corollary 1.** *Let $\mathcal{M}$ be a POMC satisfying Assump. 1, $\Lambda \subset \Sigma^*$ be a set of bounded length observation sequences with bound $n$, $f \colon \Sigma^n \to \{0,1\}$ be the indicator function of the set $\overline{\Lambda}$, and $\vec{W} \sim \mathcal{M}$ be a random observed path of an arbitrary length $|\vec{W}| > n$. Then $\mathbb{E}(\hat{f}(\vec{W})) = P(\Lambda)$.*

## 4.2  The Atomic Monitor

A monitor for each individual atom is called an atomic monitor, which serves as the building block for the overall monitor. Each atomic monitor is constructed by computing an interval estimate of the semantic value $f(\mathcal{M})$ for the respective atom $f$ on the unknown POMC $\mathcal{M}$. For computing the interval estimate, we use the McDiarmid-style inequality (see Thm. 4 in appendix) to find a bound on the width of the interval around the point estimate $\hat{f}(\cdot)$.

**Algorithm 1** $Monitor_{(f,\delta)}$: Monitor for $(f,\delta)$ where $f\colon \Sigma^n \to [a,b]$ is an atomic function of a BSE

<div>

1: **function** $Init()$
2:    $t \leftarrow 0$   ▷current time
3:    $y \leftarrow 0$ ▷current point estimate
4:    $\vec{w} \leftarrow \underbrace{\bot \ldots \bot}_{n \text{ times}}$ ▷a dummy word of length $n$, where $\bot$ is the dummy symbol
5: **end function**

</div>

<div>

1: **function** $Next(\sigma)$
2:    $t \leftarrow t + 1$         ▷progress time
3:    **if** $t < n$ **then**   ▷too short observation sequence
4:       $\vec{w}_t \leftarrow \sigma$
5:       **return** $\bot$         ▷inconclusive
6:    **else**
7:       $\vec{w}_{1..n-1} \leftarrow \vec{w}_{2..n}$     ▷shift window
8:       $\vec{w}_n \leftarrow \sigma$   ▷add the new observation
9:       $x \leftarrow f(\vec{w})$     ▷latest evaluation of $f$
10:      $y \leftarrow (y * (t - n) + x)/(t - n + 1)$   ▷running av. impl. of Eq. 6
11:      $\varepsilon \leftarrow \sqrt{-\ln(\delta/2) \cdot \frac{t \cdot \min(t-n+1,n)\cdot 9 \cdot \tau_{mix}}{2(t-n+1)^2}}$ ▷ PAC bound, see Thm. 4 in the appendix
12:       **return** $[y - \varepsilon, y + \varepsilon]$   ▷confidence interval
13:    **end if**
14: **end function**

</div>

McDiarmid's inequality is a concentration inequality bounding the distance between the sample value and the expected value of a function satisfying the bounded difference property when evaluated on independent random variables. There are several works extending this result to functions evaluated over a sequence of dependent random variables, including Markov chains [52,22,41]. In order to use McDiarmid's inequality, we will need the following standard [44] additional assumption on the underlying POMC.

**Assumption 2** *We assume that the POMCs are aperiodic, and that the mixing time of the POMC is bounded by a known constant $\tau_{\mathsf{mix}}$.*

We summarize the algorithmic computation of the atomic monitor in Alg. 1, and establish its correctness in the following theorem.

**Theorem 1 (Solution of Prob. 1 for atomic formulas).** *Let $(f,\delta)$ be a problem instance where $f\colon \Sigma^n \to [a,b]$ is an atomic formula for some fixed $n$, $a$, and $b$. Moreover, suppose the given unknown POMC satisfies Assump. 2. Then Algorithm 1 implements a monitor solving Problem 1 for the given problem instance. The monitor requires $\mathcal{O}(n)$-space, and, after arrival of each new observation, computes the updated output in $\mathcal{O}(n)$-time.*

The confidence intervals generated by McDiarmid-style inequalities for Markov chains tighten in relation to the mixing time of the Markov chain. This means the slower a POMC mixes, the longer the monitor needs to watch to be able to obtain an output interval of the same quality.

### 4.3 The Complete Monitor

The final monitors for QuantBSE-s and QualBSE-s are presented in Alg. 3 and Alg. 2, respectively, where we recursively combine the interval estimates of the

constituent sub-expressions using interval arithmetic and the union bound. Similar idea was used by Albarghouthi et al. [3]. The correctness and computational complexities of the monitors are formally stated below.

**Theorem 2 (Solution of Prob. 1).** *Let $(\varphi_1 \odot \varphi_2, \delta_1 + \delta_2)$ be a problem instance where $\varphi_1, \varphi_2$ are a pair of QuantBSE-s and $\odot \in \{+, \cdot, \div\}$. Moreover, suppose the given unknown POMC satisfies Assump. 2. Then Alg. 3 implements the monitor $\mathcal{A}$ solving Problem 1 for the given problem instance. If the total number of atoms in $\varphi_1 \odot \varphi_2$ is $k$ and if the arity of the largest atom in $\varphi_1 \odot \varphi_2$ is $n$, then $\mathcal{A}$ requires $\mathcal{O}(k+n)$-space, and, after arrival of each new observation, computes the updated output in $\mathcal{O}(k \cdot n)$-time.*

**Theorem 3 (Solution of Prob. 2).** *Let $(\psi, \delta)$ be a problem instance where $\psi$ is a QualBSE. Moreover, suppose the given unknown POMC satisfies Assump. 2. Then Alg. 2 implements the monitor $\mathcal{A}$ solving Problem 2 for the given problem instance. If the total number of atoms in $\psi$ is $k$ and if the arity of the largest atom in $\psi$ is $n$, then $\mathcal{A}$ requires $\mathcal{O}(k+n)$-space, and, after arrival of each new observation, computes the updated output in $\mathcal{O}(k \cdot n)$-time.*

---

**Algorithm 2** $Monitor_{(\psi,\delta)}$

---

1: **function** $Init()$
2:      **if** $\psi \equiv \varphi \geq 0$ **then**
3:          $\mathcal{A} \leftarrow Monitor_{(\varphi,\delta)}$
4:          $\mathcal{A}.Init()$
5:      **else if** $\psi \equiv \neg\psi_1$ **then**
6:          $\mathcal{A} \leftarrow Monitor_{(\psi_1,\delta)}$
7:          $\mathcal{A}.Init()$
8:      **else if** $\psi \equiv \psi_1 \wedge \psi_2$ **then**
9:          Choose $\delta_1, \delta_2$ s.t. $\delta = \delta_1 + \delta_2$
10:        $\mathcal{A}_1 \leftarrow Monitor_{(\psi_1,\delta_1)}$
11:        $\mathcal{A}_2 \leftarrow Monitor_{(\psi_2,\delta_2)}$
12:        $\mathcal{A}_1.Init()$
13:        $\mathcal{A}_2.Init()$
14:      **end if**
15: **end function**

1: **function** $Next(\sigma)$
2:      **if** $\psi \equiv \varphi \geq 0$ **then**
3:          $[l, u] \leftarrow \mathcal{A}.Next(\sigma)$
4:          **if** $l \geq 0$ **then return** *true*
5:          **else if** $u \leq 0$ **then return** *false*
6:          **else return** $\bot$   ▷don't know, we assume $\neg\bot = \bot \wedge true = \bot \wedge false = \bot$.
7:          **end if**
8:      **else if** $\psi \equiv \neg\psi_1$ **then**
9:          **return** $\neg(\mathcal{A}.Next(\sigma))$
10:      **else if** $\psi \equiv \psi_1 \wedge \psi_2$ **then**
11:        **return** $\mathcal{A}_1.Next(\sigma) \wedge \mathcal{A}_2.Next(\sigma)$
12:      **end if**
13: **end function**

---

## 5 Experiments

We implemented our monitoring algorithm in Python, and applied it to the real-world lending example [16] described in Ex. 1 and to an academic example called hypercube. We ran the experiments on a MacBook Pro (2023) with Apple M2 Pro processor and 16GB of RAM.

**The Lending Example.** The underlying POMC model (unknown to the monitor) of the system is approximately as shown in Fig. 1 with a few differences.

Firstly, we added a low-probability self-loop on the state $S$ to ensure aperiodicity. Secondly, we considered only two credit score levels. Thirdly, there are more hidden states (in total 171 states) in the system, like the action of the individual (repaying the loan or defaulting), etc. We monitor demographic parity, defined as $\varphi_{\mathsf{DP}} := P(Y \mid A) - P(Y \mid B)$, and an absolute version of it, defined as $\varphi_{\mathsf{TDP}} := P(AY) - P(BY)$. While $\varphi_{\mathsf{DP}}$ represents the difference in probabilities of giving loans to individuals from the two groups ($A$ and $B$), $\varphi_{\mathsf{TDP}}$ represents the difference in joint probabilities of selecting and then giving loans to individuals from the two groups. None of the two properties can be expressed using the previous formalism [3,32], because $\varphi_{\mathsf{DP}}$ requires conditioning on observations, and $\varphi_{\mathsf{TDP}}$ requires expressing absolute probabilities, which were not considered before.

---

**Algorithm 3** $Monitor_{(\varphi_1 \odot \varphi_2, \delta_1 + \delta_2)}$

---

1: **function** $Init()$
2: $\quad \mathcal{A}_1 \leftarrow Monitor_{(\varphi_1, \delta_1)}$
3: $\quad \mathcal{A}_2 \leftarrow Monitor_{(\varphi_2, \delta_2)}$
4: $\quad \mathcal{A}_1.Init()$
5: $\quad \mathcal{A}_2.Init()$
6: **end function**

1: **function** $Next(\sigma)$
2: $\quad [l_1, u_1] \leftarrow \mathcal{A}_1.Next(\sigma)$
3: $\quad [l_2, u_2] \leftarrow \mathcal{A}_2.Next(\sigma)$
4: $\quad$ **return** $[l_1, u_1] \odot [l_2, u_2]$ ▷interval arithmetic
5: **end function**

---

After receiving new observations, the monitors for $\varphi_{\mathsf{DP}}$ and $\varphi_{\mathsf{TDP}}$ took, respectively, $47\,\mu s$ and $18\,\mu s$ on an average (overall $43\,\mu s$–$0.2\,s$ and $12\,\mu s$–$3.2\,s$) to update their outputs, showing that our monitors are fast in practice.

Fig. 2 shows the outputs of the monitors for $\delta = 0.05$ (i.e., 95% confidence interval). For the POMC of the lending example, we used a pessimistic bound $\tau_{\mathsf{mix}} = 170589.78$ steps on the mixing time (computation as in [36]), with which the estimation error $\varepsilon$ shrinks rather slowly in both cases. For example, for $\varphi_{\mathsf{TDP}}$, in order to get from the trivial value $\varepsilon = 1$ (the confidence interval spans the entire range of possible values) down to $\varepsilon = 0.1$, the monitor requires about $4 \cdot 10^9$ observations. For $\varphi_{\mathsf{DP}}$, the monitor requires even more number of observations ($\sim 10^{12}$) to reach the same error level. This is because $\varphi_{\mathsf{DP}}$ involves conditional probabilities requiring divisions, which amplify the error when composed using interval arithmetics. We conclude that a direct division-free estimation of the conditional probabilities, together with tighter bounds on the mixing time will significantly improve the long-run accuracy of the monitor.

**The Hypercube Example.** We considered a second example [44, pp. 63], whose purpose is to demonstrate that the tightness of our monitors' outputs is sensitive to the choice of the bound on the mixing time. The POMC models a random walk along the edges of a hypercube $\{0,1\}^n$, where each vertex of the hypercube represents a state in the POMC and states starting with 0 and 1 are mapped to the observations $a$ and $b$, respectively. We fix $n$ to 3 in our experiments. At every step, the current vertex is chosen with probability $1/2$, and every neighbor is chosen with probability $1/2n$. A tight bound on the mixing time of this POMC is given by $\tau_{\mathsf{true\,mix}} = n(\log n + \log 4)$ steps [44, pp. 63]. We consider the properties $\psi_{\mathsf{DP}} := P(a \mid a) - P(b \mid b)$ and $\psi_{\mathsf{TDP}} := P(aa) - P(bb)$.
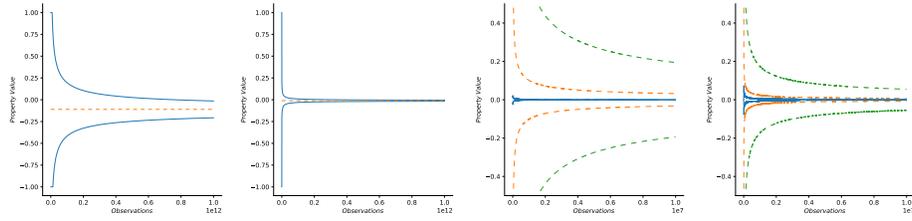
Fig. 2: Monitoring $\varphi_{\mathsf{DP}}$ (first, third) and $\varphi_{\mathsf{TDP}}$ (second, fourth) on the lending (first, second) and the hypercube (third, fourth) examples. The <u>first and second</u> plots show the computed 95%-confidence interval (solid) and the true value of the property (dashed) for the lending POMC. In reality, the monitor was run for about $7 \times 10^8$ steps until the point estimate nearly converged, though the confidence interval was trivial at this point (the whole interval $[-1, 1]$), owing to the pessimistic bound $\tau_{\mathsf{mix}}$. In the figure, we have plotted a projection of how the confidence interval would taper over time, had we kept the monitor running. The <u>third and fourth</u> plots summarize the monitors' outputs over 100 executions of the hypercube POMC. The solid lines are the max and min values of the point estimates, the dashed lines are the boundaries of all the 95%-confidence intervals (among the 100 executions) with the conservative bound $\tau_{\mathsf{mix}}$ (green) and the sharper bound $\tau_{\mathsf{true\,mix}}$ (orange) on the mixing time.

We empirically evaluated the quality of the confidence intervals computed by our monitor (for $\psi_{\mathsf{DP}}$ and $\psi_{\mathsf{TDP}}$) over a set of 100 sample runs, and summarize the findings in the third and fourth plots of Fig. 2. We used $\tau_{\mathsf{mix}} = 204.94$ steps and $\tau_{\mathsf{true\,mix}} = 7.45$ steps, and we can observe that in both cases, the output with $\tau_{\mathsf{true\,mix}}$ is significantly tighter than with $\tau_{\mathsf{mix}}$. Compared to the lending example, we obtain reasonably tight estimate with significantly smaller number of observations, which is due to the smaller bounds on the mixing time.

## 6   Conclusion

We generalized runtime verification of algorithmic fairness properties to systems modeled using POMCs and a specification language (BSE) with arithmetic expressions over numerical functions assigning values to observation sequences. Under the assumptions of stationary initial distribution, aperiodicity, and the knowledge of a bound on the mixing time, we presented a runtime monitor, which monitors a long sequence of observations generated by the POMC, and after each observation outputs an updated PAC estimate of the value of the given BSE.

While the new stationarity assumption is important for defining the semantics of the BSE expressions, the aperiodicity and the knowledge of the bound on the mixing time allow us to use the known McDiarmid's inequality for computing

the PAC estimate. In future, we intend to eliminate the latter two assumptions, enabling us to use our approach for a broader class of systems. Additionally, eliminating the time-homogeneity assumption would also be an important step for monitoring algorithmic fairness of the real-world systems with time-varying probability distributions [31].

# References

1. Agha, G., Palmskog, K.: A survey of statistical model checking. ACM Transactions on Modeling and Computer Simulation (TOMACS) **28**(1), 1–39 (2018)
2. Albarghouthi, A., D'Antoni, L., Drews, S., Nori, A.V.: Fairsquare: probabilistic verification of program fairness. Proceedings of the ACM on Programming Languages **1**(OOPSLA), 1–30 (2017)
3. Albarghouthi, A., Vinitsky, S.: Fairness-aware programming. In: Proceedings of the Conference on Fairness, Accountability, and Transparency. pp. 211–219 (2019)
4. Ashok, P., Křetínský, J., Weininger, M.: Pac statistical model checking for markov decision processes and stochastic games. In: International Conference on Computer Aided Verification. pp. 497–519. Springer (2019)
5. Baier, C., Haverkort, B., Hermanns, H., Katoen, J.P.: Model-checking algorithms for continuous-time markov chains. IEEE Transactions on Software Engineering **29**(6), 524–541 (2003). `https://doi.org/10.1109/TSE.2003.1205180`
6. Baier, C., Katoen, J.P.: Principles of model checking. MIT press (2008)
7. Balunovic, M., Ruoss, A., Vechev, M.: Fair normalizing flows. In: International Conference on Learning Representations (2021)
8. Bartocci, E., Deshmukh, J., Donzé, A., Fainekos, G., Maler, O., Ničković, D., Sankaranarayanan, S.: Specification-based monitoring of cyber-physical systems: a survey on theory, tools and applications. In: Lectures on Runtime Verification, pp. 135–175. Springer (2018)
9. Bartocci, E., Falcone, Y.: Lectures on Runtime Verification. Springer (2018)
10. Bartolo Burlò, C., Francalanza, A., Scalas, A., Trubiani, C., Tuosto, E.: Towards probabilistic session-type monitoring. In: International Conference on Coordination Languages and Models. pp. 106–120. Springer (2021)
11. Bastani, O., Zhang, X., Solar-Lezama, A.: Probabilistic verification of fairness properties via concentration. Proceedings of the ACM on Programming Languages **3**(OOPSLA), 1–27 (2019)
12. Bellamy, R.K., Dey, K., Hind, M., Hoffman, S.C., Houde, S., Kannan, K., Lohia, P., Martino, J., Mehta, S., Mojsilović, A., et al.: Ai fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. IBM Journal of Research and Development **63**(4/5), 4–1 (2019)
13. Bird, S., Dudik, M., Edgar, R., Horn, B., Lutz, R., Milan, V., Sameki, M., Wallach, H., Walker, K.: Fairlearn: A toolkit for assessing and improving fairness in ai. Microsoft, Tech. Rep. MSR-TR-2020-32 (2020)
14. Chouldechova, A.: Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. Big data **5**(2), 153–163 (2017)

15. Clarke, E.M., Zuliani, P.: Statistical model checking for cyber-physical systems. In: International symposium on automated technology for verification and analysis. pp. 1–12. Springer (2011)

16. D'Amour, A., Srinivasan, H., Atwood, J., Baljekar, P., Sculley, D., Halpern, Y.: Fairness is not static: Deeper understanding of long term fairness via simulation studies. In: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency. p. 525–534. FAT* '20 (2020)

17. David, A., Du, D., Guldstrand Larsen, K., Legay, A., Mikučionis, M.: Optimizing control strategy using statistical model checking. In: NASA Formal Methods Symposium. pp. 352–367. Springer (2013)

18. Donzé, A., Maler, O.: Robust satisfaction of temporal logic over real-valued signals. In: International Conference on Formal Modeling and Analysis of Timed Systems. pp. 92–106. Springer (2010)

19. Dressel, J., Farid, H.: The accuracy, fairness, and limits of predicting recidivism. Science advances $4$(1), eaao5580 (2018)

20. Dwork, C., Hardt, M., Pitassi, T., Reingold, O., Zemel, R.: Fairness through awareness. In: Proceedings of the 3rd innovations in theoretical computer science conference. pp. 214–226 (2012)

21. Ensign, D., Friedler, S.A., Neville, S., Scheidegger, C., Venkatasubramanian, S.: Runaway feedback loops in predictive policing. In: Conference on Fairness, Accountability and Transparency. pp. 160–171. PMLR (2018)

22. Esposito, A.R., Mondelli, M.: Concentration without independence via information measures. arXiv preprint arXiv:2303.07245 (2023)

23. Faymonville, P., Finkbeiner, B., Schwenger, M., Torfah, H.: Real-time stream-based monitoring. arXiv preprint arXiv:1711.03829 (2017)

24. Feldman, M., Friedler, S.A., Moeller, J., Scheidegger, C., Venkatasubramanian, S.: Certifying and removing disparate impact. In: proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining. pp. 259–268 (2015)

25. Ferrere, T., Henzinger, T.A., Kragl, B.: Monitoring event frequencies. In: 28th EACSL Annual Conference on Computer Science Logic. vol. 152 (2020)

26. Ferrère, T., Henzinger, T.A., Saraç, N.E.: A theory of register monitors. In: Proceedings of the 33rd Annual ACM/IEEE Symposium on Logic in Computer Science. pp. 394–403 (2018)

27. Finkbeiner, B., Sankaranarayanan, S., Sipma, H.: Collecting statistics over runtime executions. Electronic Notes in Theoretical Computer Science $70$(4), 36–54 (2002)

28. Ghosh, B., Basu, D., Meel, K.S.: Justicia: A stochastic sat approach to formally verify fairness. arXiv preprint arXiv:2009.06516 (2020)

29. Ghosh, B., Basu, D., Meel, K.S.: Algorithmic fairness verification with graphical models. arXiv preprint arXiv:2109.09447 (2021)

30. Hardt, M., Price, E., Srebro, N.: Equality of opportunity in supervised learning. Advances in neural information processing systems $29$ (2016)

31. Henzinger, T., Karimi, M., Kueffner, K., Mallik, K.: Runtime monitoring of dynamic fairness properties. In: Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency. pp. 604–614 (2023)

32. Henzinger, T.A., Karimi, M., Kueffner, K., Mallik, K.: Monitoring algorithmic fairness. In: Enea, C., Lal, A. (eds.) Computer Aided Verification. pp. 358–382. Springer (2023)

33. Henzinger, T.A., Saraç, N.E.: Monitorability under assumptions. In: International Conference on Runtime Verification. pp. 3–18. Springer (2020)

34. Henzinger, T.A., Saraç, N.E.: Quantitative and approximate monitoring. In: 2021 36th Annual ACM/IEEE Symposium on Logic in Computer Science (LICS). pp. 1–14. IEEE (2021)

35. Jagielski, M., Kearns, M., Mao, J., Oprea, A., Roth, A., Sharifi-Malvajerdi, S., Ullman, J.: Differentially private fair learning. In: International Conference on Machine Learning. pp. 3000–3008. PMLR (2019)

36. Jerison, D.: General mixing time bounds for finite markov chains via the absolute spectral gap. arXiv preprint arXiv:1310.8021 (2013)

37. John, P.G., Vijaykeerthy, D., Saha, D.: Verifying individual fairness in machine learning models. In: Conference on Uncertainty in Artificial Intelligence. pp. 749–758. PMLR (2020)

38. Junges, S., Torfah, H., Seshia, S.A.: Runtime monitors for markov decision processes. In: International Conference on Computer Aided Verification. pp. 553–576. Springer (2021)

39. Kearns, M., Neel, S., Roth, A., Wu, Z.S.: Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. In: International Conference on Machine Learning. pp. 2564–2572. PMLR (2018)

40. Konstantinov, N.H., Lampert, C.: Fairness-aware pac learning from corrupted data. Journal of Machine Learning Research **23** (2022)

41. Kontorovich, A., Raginsky, M.: Concentration of measure without independence: a unified approach via the martingale method. In: Convexity and Concentration. pp. 183–210. Springer (2017)

42. Kusner, M.J., Loftus, J., Russell, C., Silva, R.: Counterfactual fairness. Advances in neural information processing systems **30** (2017)

43. Lahoti, P., Gummadi, K.P., Weikum, G.: ifair: Learning individually fair data representations for algorithmic decision making. In: 2019 ieee 35th international conference on data engineering (icde). pp. 1334–1345. IEEE (2019)

44. Levin, D.A., Peres, Y.: Markov chains and mixing times, vol. 107. American Mathematical Soc. (2017)

45. Liu, L.T., Dean, S., Rolf, E., Simchowitz, M., Hardt, M.: Delayed impact of fair machine learning. In: International Conference on Machine Learning. pp. 3150–3158. PMLR (2018)

46. Lum, K., Isaac, W.: To predict and serve? Significance **13**(5), 14–19 (2016)

47. Maler, O., Nickovic, D.: Monitoring temporal properties of continuous signals. In: Formal Techniques, Modelling and Analysis of Timed and Fault-Tolerant Systems, pp. 152–166. Springer (2004)

48. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., Galstyan, A.: A survey on bias and fairness in machine learning. ACM Computing Surveys (CSUR) **54**(6), 1–35 (2021)

49. Meyer, A., Albarghouthi, A., D'Antoni, L.: Certifying robustness to programmable data bias in decision trees. Advances in Neural Information Processing Systems **34**, 26276–26288 (2021)

50. Obermeyer, Z., Powers, B., Vogeli, C., Mullainathan, S.: Dissecting racial bias in an algorithm used to manage the health of populations. Science **366**(6464), 447–453 (2019)

51. Otop, J., Henzinger, T.A., Chatterjee, K.: Quantitative automata under probabilistic semantics. Logical Methods in Computer Science **15** (2019)

52. Paulin, D.: Concentration inequalities for markov chains by marton couplings and spectral methods (2015)

53. Ruchkin, I., Sokolsky, O., Weimer, J., Hedaoo, T., Lee, I.: Compositional probabilistic analysis of temporal properties over stochastic detectors. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems **39**(11), 3288–3299 (2020)

54. Scheuerman, M.K., Paul, J.M., Brubaker, J.R.: How computers see gender: An evaluation of gender classification in commercial facial analysis services. Proceedings of the ACM on Human-Computer Interaction **3**(CSCW), 1–33 (2019)

55. Seyyed-Kalantari, L., Liu, G., McDermott, M., Chen, I.Y., Ghassemi, M.: Chexclusion: Fairness gaps in deep chest x-ray classifiers. In: BIOCOMPUTING 2021: proceedings of the Pacific symposium. pp. 232–243. World Scientific (2020)

56. Sharifi-Malvajerdi, S., Kearns, M., Roth, A.: Average individual fairness: Algorithms, generalization and experiments. Advances in Neural Information Processing Systems **32** (2019)

57. Stoller, S.D., Bartocci, E., Seyster, J., Grosu, R., Havelund, K., Smolka, S.A., Zadok, E.: Runtime verification with state estimation. In: International conference on runtime verification. pp. 193–207. Springer (2011)

58. Stoller, S.D., Bartocci, E., Seyster, J., Grosu, R., Havelund, K., Smolka, S.A., Zadok, E.: Runtime verification with state estimation. In: Runtime Verification: Second International Conference, RV 2011, San Francisco, CA, USA, September 27-30, 2011, Revised Selected Papers 2. pp. 193–207. Springer (2012)

59. Sun, B., Sun, J., Dai, T., Zhang, L.: Probabilistic verification of neural networks against group fairness. In: International Symposium on Formal Methods. pp. 83–102. Springer (2021)

60. Waudby-Smith, I., Arbour, D., Sinha, R., Kennedy, E.H., Ramdas, A.: Time-uniform central limit theory, asymptotic confidence sequences, and anytime-valid causal inference. arXiv preprint arXiv:2103.06476 (2021)

61. Wexler, J., Pushkarna, M., Bolukbasi, T., Wattenberg, M., Viégas, F., Wilson, J.: The what-if tool: Interactive probing of machine learning models. IEEE transactions on visualization and computer graphics **26**(1), 56–65 (2019)

62. Younes, H.L., Simmons, R.G.: Probabilistic verification of discrete event systems using acceptance sampling. In: International Conference on Computer Aided Verification. pp. 223–235. Springer (2002)

63. Zemel, R., Wu, Y., Swersky, K., Pitassi, T., Dwork, C.: Learning fair representations. In: International conference on machine learning. pp. 325–333. PMLR (2013)

# A   Proof of Claims in Sec. 4.1

*Proof (Proof of Prop. 1).* Let $N = t - n + 1$. By definition,

$$\mathbb{E}(\hat{f}(\overrightarrow{W})) = \sum_{\overrightarrow{w} \in \Sigma^t} \left( \frac{1}{N} \sum_{i=1}^{N} f(\overrightarrow{w}_{i..i+n-1}) \right) \cdot \mathbb{P}(\overrightarrow{w})$$

where $\mathbb{P}(\overrightarrow{u})$ is defined with respect to the Markov chain. We now express this with respect to the joint probability measure over paths of length $t$ as defined

by $\mathcal{M}$.

$$\mathbb{E}(\hat{f}(\overrightarrow{W})) = \sum_{\overrightarrow{q} \in Q^t} \left( \frac{1}{N} \sum_{i=1}^{N} f(\ell(\overrightarrow{q}_{i..i+n-1})) \right) \cdot \mathbb{P}(\overrightarrow{q})$$

$$= \sum_{\overrightarrow{q} \in Q^t} \frac{1}{N} \sum_{i=1}^{N} f(\ell(\overrightarrow{q}_{i..i+n-1})) \cdot \mathbb{P}(\overrightarrow{q})$$

$$= \frac{1}{N} \sum_{i=1}^{N} \sum_{\overrightarrow{q} \in Q^t} f(\ell(\overrightarrow{q}_{i..i+n-1})) \cdot \mathbb{P}(\overrightarrow{q}).$$

Fix a particular $i$, i.e. let $A_i := \sum_{\overrightarrow{q} \in Q^t} f(\ell(\overrightarrow{q}_{i..i+n-1})) \cdot \mathbb{P}(\overrightarrow{q})$. Split up the internal sum to obtain

$$A_i = \sum_{\overrightarrow{q}^A \in Q^i} \sum_{\overrightarrow{q}^B \in Q^n} \sum_{\overrightarrow{q}^C \in Q^{t-n-i}} f(\ell(\overrightarrow{q}^B)) \cdot \mathbb{P}(\overrightarrow{q}^A) \cdot \mathbb{P}(\overrightarrow{q}^B \mid q_i^A) \cdot \mathbb{P}(\overrightarrow{q}^C \mid q_n^B).$$

Now notice that $\sum_{\overrightarrow{q}^C \in Q^{t-n-i}} \mathbb{P}(\overrightarrow{q}^C \mid \overrightarrow{q}_n^B) = 1$. Hence, by rearranging the sums we obtain

$$A_i = \sum_{\overrightarrow{q}^B \in Q^n} f(\ell(\overrightarrow{q}^B)) \cdot \sum_{\overrightarrow{q}^A \in Q^i} \mathbb{P}(\overrightarrow{q}^A) \cdot \mathbb{P}(\overrightarrow{q}^B \mid \overrightarrow{q}_i^A)$$

$$= \sum_{\overrightarrow{q}^B \in Q^n} f(\ell(\overrightarrow{q}^B)) \cdot \prod_{j=1}^{n-1} M_{q_j^B q_{j+1}^B} \cdot \sum_{\overrightarrow{q}^A \in Q^i} \pi_{q_1^A} \cdot \prod_{j=1}^{i-1} M_{q_j^A q_{j+1}^A} \cdot M_{q_i^A q_1^B}.$$

We can use the fact that $\mathcal{M}$ is in its stationary distribution to express the internal sum as $\pi_{q_1^B}$. That is,

$$\sum_{\overrightarrow{q}^A \in Q^i} \pi_{q_1^A} \cdot \prod_{j=1}^{i-1} M_{q_j^A q_{j+1}^A} \cdot M_{q_i^A q_1^B} = \sum_{q_i^A} \cdots \sum_{q_2^A} \sum_{q_1^A} \pi_{q_1^A} \cdot \prod_{j=1}^{i-1} M_{q_j^A q_{j+1}^A} \cdot M_{q_i^A q_1^B}$$

$$= \sum_{q_i^A} M_{q_i^A q_1^B} \cdots \sum_{q_2^A} M_{q_2^A q_3^A} \cdot \sum_{q_1^A} \pi_{q_1^A} \cdot M_{q_1^A q_2^A} = \sum_{q_i^A} M_{q_i^A q_1^B} \cdots \sum_{q_2^A} M_{q_2^A q_3^A} \pi_{q_2^A}$$

$$= \sum_{q_i^A} M_{q_i^A q_1^B} \cdot \pi_{q_i^A} = \pi_{q_1^B}.$$

Hence, we obtain

$$A_i = \sum_{\overrightarrow{q}^B \in Q^n} f(\ell(\overrightarrow{q}^B)) \cdot \prod_{j=1}^{n-1} M_{q_j^B q_{j+1}^B} \cdot \pi_{q_1^B} = \mathbb{E}(f(\overrightarrow{U}))$$

where $\overrightarrow{U}$ is a random word of length $n$ generated by $\mathcal{M}$. Therefore, we obtain

$$\mathbb{E}(\hat{f}(\overrightarrow{W})) = \frac{1}{N} \sum_{i=1}^{N} \mathbb{E}(f(\overrightarrow{U})) = \mathbb{E}(f(\overrightarrow{U})).$$

Moreover, this demonstrates that due to stationarity the expected value of $f$ evaluated on any infix of length $n$ is the same and thus $\mathbb{E}(\hat{f}(\overrightarrow{W})) = f(\mathcal{M})$.

*Proof (Proof of Cor. 1).* This follows directly from Prop. 1 and the observation that $f(\ell(\overrightarrow{q^B}))$ removes the probability mass of all the paths of length $n$ whose corresponding words do not belong to $\Lambda$. Therefore, $\mathbb{E}(f(\overrightarrow{U})) = P(\Lambda)$ where $\overrightarrow{U}$ is as defined in Prop. 1.

## B  Proof of Thm. 1

We use a McDiarmid-style inequality to compute the finite-sample confidence bounds. The version below is a restricted version of the Corollary 2.19 found in [36].

**Theorem 4 ([36]).** *Let $\overrightarrow{X} \coloneqq X_1, \ldots, X_n$ be an ergodic Markov chain with countable state space $Q$, unique stationary distribution $\pi$, and finite mixing time bounded by $\tau_{mix}$. Suppose that some function $f : Q^n \to \mathbb{R}$ with artiy $n$ satisfies*

$$f(x) - f(y) \leq \sum_{i=1}^{g} nc_i \mathbb{1}(x_i \neq y_i)$$

*for some $c \in \mathbb{R}^n$ with positive entries. Then for any $\varepsilon > 0$*

$$\mathbb{P}\left( \left| f(\overrightarrow{X}) - \mathbb{E}(f(\overrightarrow{X})) \right| \geq \varepsilon \right) \leq 2 \exp\left( -\frac{2\varepsilon^2}{\sqrt{\sum_{i=1}^{n} c_i^2} \cdot 9 \cdot \tau_{mix}} \right)$$

To apply Theorem 4 it is required that the observation labels should not interfere with the so-called bounded difference property of the function. Below we establish that this requirement is fulfilled by the atoms of BSE.

**Lemma 1.** *Let $f \colon \Sigma^n \to [a, b]$ be a function with fixed $n$, $a$, and $b$, $t \geq n$ be a constant, $\overrightarrow{w}, \overrightarrow{w}' \in \Sigma^t$ be a pair of observation sequences such that the Hamming distance $|\overrightarrow{w} - \overrightarrow{w}'|_H$ is 1. Then*

$$\hat{f}(\overrightarrow{w}) - \hat{f}(\overrightarrow{w}') \leq \frac{\min(t - n + 1, n)}{t - n + 1} \cdot (b - a).$$

*Proof.* Since the Hamming distance is 1, $\overrightarrow{w}$ and $\overrightarrow{w}'$ differ only in one symbol. We know that $f$ is evaluated on a substring of length $n$. Hence, if the string is sufficiently long only $n$ evaluations of $f$ in $\hat{f}$ (i.e., only $n$ terms in the sum in Eq. 6) will be affected, while if the string is short, then only $t - n + 1$ evaluations of $f$ will be affected. Therefore, the evaluation of $f$ can differ in the worst case by at most $\frac{\min(t-n+1,n)}{t-n+1} \cdot (b - a)$.

**Lemma 2.** *Let $\mathcal{M}$ be a POMC satisfying Assump. 1 and 2, $f : \Sigma^n \to [a, b]$ be a function for a fixed $n$, $a$, and $b$, $t \geq n$ be a constant, and $\overrightarrow{W} \sim \mathcal{M}$ be a random observed $\mathcal{M}$-path of length $|\overrightarrow{W}| = t$. Then*

$$\mathbb{P}\left( |f(\overrightarrow{\mathcal{M}}) - \hat{f}(\overrightarrow{W})| \geq \varepsilon \right) \leq 2 \exp\left( -\frac{2 \cdot \varepsilon^2 (t - n + 1)^2}{t \cdot \min(t - n + 1, n) \cdot 9 \cdot \tau_{mix}} \right).$$

*Proof.* Notice that $\vec{w}, \vec{w}' \in \Sigma^t$

$$|f(w) - f(\vec{w}')| \leq \sum_{i=1}^{t} \frac{\min(t-n+1, n)}{t-n+1} \cdot (b-a) \cdot \mathbb{1}(w_i \neq w_i')$$

Therefore, we conclude that

$$\left( \sqrt{\sum_{i=1}^{t} \left( \frac{\min(t-n+1, n)}{t-n+1} \cdot (b-a) \right)^2} \right)^2 = \frac{t \cdot \min(t-n+1, n)^2}{(t-n+1)^2} \cdot (b-a)^2$$

as required by Theorem 4.

*Proof (Proof of Thm. 1).* The soundness claim follows as a consequence of Lem. 1 and Prop. 1. By combining Theorem 4 and Corollary 2.17 from [36] we obtain the result for POMC. The computational complexity is dominated by the use of the set of $n$ registers $\vec{w}$ to store the last $n$ sub-sequence of the observed path: allocation of memory for $\vec{w}$ takes $n$ space, and, after every new observation, the update of $\vec{w}$ takes $n$ write operations (Line 7).