# SANDBOXING CONTROLLERS FOR STOCHASTIC CYBER-PHYSICAL SYSTEMS[*]

BINGZHUO ZHONG[1], MAJID ZAMANI[3,2], AND MARCO CACCAMO[1]

ABSTRACT. Current cyber-physical systems (CPS) are expected to accomplish complex tasks. To achieve this goal, high performance, but unverified controllers (e.g. deep neural network, black-box controllers from third parties) are applied, which makes it very challenging to keep the overall CPS safe. By sandboxing these controllers, we are not only able to use them but also to enforce safety properties over the controlled physical systems at the same time. However, current available solutions for sandboxing controllers are just applicable to deterministic (a.k.a. non-stochastic) systems, possibly affected by bounded disturbances. In this paper, for the first time we propose a novel solution for sandboxing unverified complex controllers for CPS operating in noisy environments (a.k.a. stochastic CPS). Moreover, we also provide probabilistic guarantees on their safety. Here, the unverified control input is observed at each time instant and checked whether it violates the maximal tolerable probability of reaching the unsafe set. If this probability exceeds a given threshold, the unverified control input will be rejected, and the advisory input provided by the optimal safety controller will be used to maintain the probabilistic safety guarantee. The proposed approach is illustrated empirically and the results indicate that the expected safety probability is guaranteed.

## 1. INTRODUCTION

Cyber-Physical Systems (CPS) are complex systems in which physical components are interacting tightly with cyber ones. These systems are widely used in various kinds of applications, such as automotive, aviation, manufacture plants and so on. Nowadays, these systems are expected to accomplish complex missions. As a result, complex, high performance but unverified controllers (e.g., deep neural network or black-box controllers from third parties) are applied to complete these complex missions, which makes it increasingly challenging to ensure the safety of CPS. To cope with this issue, we exploit the idea of *sandbox* from the community of computer security, which is a popular security mechanism for cyber systems[RBP09]. In short, it provides a testing environment to isolate the untested and untrusted components from the critical part of a digital controller. The behaviour of the untrusted component is restricted and it can only access the critical part when it follows the rules given by the sandboxing mechanism. Hence, we designed a novel architecture that uses a **Safe**ty Advisor and a Super**visor** (Safe-visor in short). Instead of providing a testing environment and focusing on cyber security, Safe-visor architecture can be used to sandbox any types of unverified controllers in run time regarding the safety of the physical systems. The control inputs of the controller fed to the system are checked and can only be accepted when they are not disobeying the safety rule defined in the sandboxing mechanism. The architecture of safe-visor is illustrated in Figure 1. In this architecture, the safety of the physical system is characterized by the probability of fulfilling some safety specifications. In general, the Safe-visor specifies verifiable safety rules for the unverified controller to follow so that a specific level of safety probability of the physical system can be ensured. During the execution of the Safe-visor, the Safety Advisor is responsible for providing advisory input for the Supervisor based on the current state of the physical system, which seeks to maximize the safety probability. Meanwhile, the Supervisor checks the input
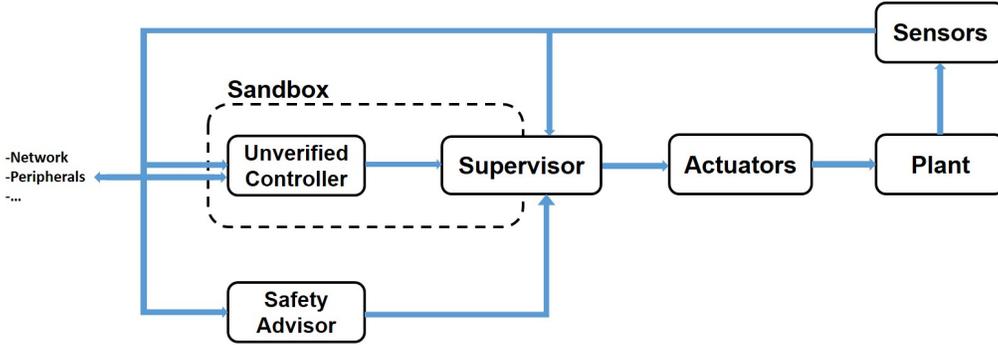
FIGURE 1. **Safe**ty Advisor - Super**visor** (Safe-visor) Architecture for sandboxing unverified controller.

given by the unverified controller according to the safety rule. Input from the unverified controller would only be accepted when it follows the rule; otherwise, the Supervisor would accept the advisory input from the Safety Advisor to maximize the safety probability of the physical system. In the rest of the paper, designing a Supervisor means designing its safety rule for checking the inputs from unverified controllers. It should be noted that inputs given by Safety Advisors only focus on the safety of the system, which should be treated as a fallback in case the unverified controllers are trying to perform some harmful actions. On the other hand, the unverified controller is designed for functionality. i.e. it is expected to realize some tasks which are much more complicated than purely keeping the system safe. By sandboxing the unverified controller, we are able to exploit its advantages for realizing complex tasks while preventing the system from being threatened by its harmful behaviour, if any.

In this paper, we deal with stochastic CPS modelled as controlled discrete-time Markov process (cdt-MP). We focus on the safety invariance specification, in which the system is expected to stay inside a pre-defined safety set. Here, we formulate the safety invariance specification as a reach-avoid problem in finite time horizon, and design the Safety Advisor based on a finite Markov Decision Process (MDP) constructed from the original cdt-MP. The inputs given by the unverified controller are checked by the Supervisor at every time instant based on an estimation of the probability of reaching the unsafe set (i.e., the complement of the pre-defined safety set).

**Related Work.** In [ABE+18][HKKT16][BKKW15], a shield is synthesized to correct erroneous output values from those unverified, complexed components in a system so that safety properties can be enforced at run time. This idea is mainly used for systems which can be modelled as automaton, e.g. reactive systems, while our method can be applied to systems with continuous state space and input space. The most relevant work to our proposed method is the one developed based on Simplex architecture [Sha01, CGR+07], in which the unverified, high-performance controller is sandboxed by an elliptic recovery region associated with a verified, high-assurance controller. Inspired by the idea of sandboxing the unverifiable controllers by using Simplex architecture, many results have been proposed for different kinds of systems and invariance specifications. In the case that bounded uncertainty exists in the system dynamic, L1-Simplex [WHS13] is applicable by using L1-adaptive controller [HCK+11] as the high-assurance controller with which the linear model uncertainty in the system dynamic is estimated and compensated. RSimplex [WHS18] uses Robust Fault-Tolerant Controller (RFTC) with the similar idea of L1-Simplex, but it is capable of dealing with non-linear model uncertainty. Net-Simplex [YLZS13] is able to cope with bounded time delay introduced by the network connection in the system. It models the system as a linear parameter-varying system and accordingly designs time-delay-related recovery region. A recent result in [ACH+18] proposes a way to sandbox an unverified controller which may suffer from undetectable cyber attacks by dynamically planning and executing high-assurance controllers so

that the physical system is not endangered. The common point of these results is that Lyapunov-function-based safety invariant sets are used as recovery regions.

The main difference between the Simplex architecture and our proposed result is that in our proposed solution, only the unverified controller is in charge of accomplishing the task, under the supervision of the Supervisor, rather than designing two parallel controllers (i.e. the high-assurance and high-performance controllers) for the given task and define a verified decision logic to decide which one to be used. The Safety Advisor is not expected to fully control the system and finish the complicated task, but it is only responsible for providing fallback to maximize the safety probability. Then, by properly designing the Supervisor, the unverified controller has more flexibility for functionality, and the safety probability can be guaranteed due to the existence of the fallback solution.

There are some other results which extend the concept of Simplex architecture using reachability analysis to cope with the aforementioned conservativeness. Results in [BMMC11] provide a backward reachability based method to generate a decision module between the mission controller and the safety one which mainly focuses on the safety of the system. Results in [BJCS14] propose a method in which real-time reachability is integrated into the Lyapunov invariance-based method. This largely increases the feasible region of the mission controller. Another idea to get rid of the conservativeness is to compute the safety invariant purely based on offline reachability analysis, as discussed in [ATR$^+$17]. It should be noted that these methods are only designed for deterministic (non-stochastic) systems. To the best of our knowledge, our result is the first work with a solution for sandboxing unverified controllers in stochastic settings.

The rest of the paper is organized as follows: we provide preliminary discussion regarding the notations, models used in this work and formulation of the problem in Section 2. Then, a scheme to design the Safety Advisor and Supervisor is proposed in Section 3, which will be empirically tested by two case studies in Section 4. Finally, Section 5 concludes the paper.

## 2. Problem Formulation

2.1. **Preliminaries.** A topological space $S$ is called a Borel space if it is homeomorphic to a Borel subset of a Polish space (i.e., a separable and completely metrizable space). One of the common examples of Borel space are the Euclidean spaces $\mathbb{R}^n$. Any Borel space $S$ is assumed to be endowed with a Borel $\sigma$-algebra denoted by $\mathcal{B}(S)$. A map $f : X \to Y$ is measurable whenever it is Borel measurable. A map $f : X \to Y$ is universally measurable if the inverse image of every Borel set is measurable with respect to every complete probability measure on $X$ that measures all Borel subsets of $X$.

For the stochastic kernel, we adopt the notation as in [TMKA13]. Given two Borel space $X$ and $Y$, the stochastic kernel on $X$ given $Y$ is the map $P : Y \times \mathcal{B}(X) \to [0,1]$ such that $P(\cdot|y)$ is a probability measure on $X$ for any point $y \in Y$ and $P(B|\cdot)$ is a measurable function on $Y$ for any set $B \in \mathcal{B}(X)$.

2.2. **Notations.** We denote by $\mathbb{R}$ the set of real numbers and by $\mathbb{N}$ the set of natural numbers. We denote by $\overline{0, n} := \{0, 1, \ldots, n\}$ an interval in $\mathbb{N}$ starting from 0 and ending at $n \in \mathbb{N}$. Set $\mathbb{R}^n$ represents the $n$-dimensional Euclidean space where $n \in \mathbb{N}$.

2.3. **Model Description and Problem Formulation.** In this paper, we focus on discrete-time stochastic control systems in the following form:

$$x(t+1) = f(x(t), u(t), w(t)), \tag{2.1}$$

in which $t, t + 1 \in \overline{0, n}$ are two successive time instants in the time domain $\overline{0, n}$ of the system, where $n \in \mathbb{N}$. Here, $x(t) \in X$ is the state of the system at time $t$, where $X \subseteq \mathbb{R}^n$ is a Borel space as the state space of the system. We denote by $(X, \mathcal{B}(X))$ the measurable space with $\mathcal{B}(X)$ being the Borel sigma-algebra on the state space. We denote by $u(t) \in U$ the input to the system at time $t$, where $U \subseteq \mathbb{R}^m$ is a Borel space as the input space of the system. We denote by $w(t)$ the uncertainty at time instant $t$ where $w : \mathbb{N} \to \mathbb{R}^d$ is a sequence of

independent and identically distributed (i.i.d.) random variables. Map $f : X \times U \times \mathbb{R}^d \to X$ is a measurable function characterizing the state dynamic of the system. In this paper, we focus on stochastic systems, which can also be formulated as controlled discrete-time Markov Processes (cdt-MP).

**Definition 2.1.** *(cdt-MP) [HLL12] A controlled discrete-time Markov process is a tuple*

$$\mathfrak{D} = (X,\, U,\, \{U(x)\}_{x \in X},\, T_{\mathfrak{D}}),$$

*where $X \subseteq \mathbb{R}^n$ is a Borel space representing the state space of the model and $U \subseteq \mathbb{R}^m$ is a Borel space referring to the input space. The set $\{U(x)\}_{x \in X}$ is a family of non-empty measurable subsets of $U$, and $U(x)$ is the set of feasible inputs when system is at state $x$. We denote by $T_{\mathfrak{D}}$ a Borel measurable stochastic kernel $T_{\mathfrak{D}} : \mathcal{B} \times X \times U \to [0,1]$, which assigns to any $x \in X$ and $u \in U(x)$ a probability measure on the Borel space $(X, \mathcal{B}(X))$ and characterizes the state transition of the Markov process.*

In the rest of the paper, we focus on systems in which $U(x) = U$, i.e. all inputs are feasible at any state in the evolution of the system. The evolution of the system is described by paths as defined below.

**Definition 2.2.** *(Path) A path of a cdt-MP $\mathfrak{D}$ is*

$$\omega = (x(0), u(0), x(1), u(1), \ldots, x(t), u(t), \ldots),$$

*where $x(t) \in X$ and $u(t) \in U$, $t \in \overline{0,n}$, $\overline{0,n} \subset \mathbb{N}$ is the time domain of the path. We denote by $\omega_x = \{x(i)\}_{i \in \overline{0,n}}$ and $\omega_u = \{u(i)\}_{i \in \overline{0,n}}$ the subsequences of states and inputs in $\omega$.*

Given a cdt-MP $\mathfrak{D}$, we are interested in Markov policies to control the system.

**Definition 2.3.** *(Markov policy) For a cdt-MP $\mathfrak{D} = (X, U, \{U(x)\}_{x \in X}, T_{\mathfrak{D}})$, a Markov policy $\mu$ is a sequence $\mu = (\mu_0, \mu_1, \mu_2, \ldots)$ of universally measurable map $\mu_t : X \to U$ at time $t \in \overline{0,n}$, where $\overline{0,n}$ is the time domain of $\mathfrak{D}$.*

With Markov policies, the input at time $t$ is only determined by the state at the same time instant, i.e. $u(t) = \mu_t(x(t))$. In this paper, we are interested in the safety specification, where the state sequences are expected to stay (with a given probability threshold) inside a safe subset of the state set. We formulate this specification as a reach-avoid problem in finite time horizon.

**Definition 2.4.** *(Reach-avoid problem) Consider a safety set $\mathcal{A} \subset \mathcal{B}(S)$, a bounded Borel set as a safe set, and $\mathcal{A}^c = \mathcal{B}(S) \backslash \mathcal{A}$, as its complement, i.e. an unsafe set. We define the reach-avoid problem under Markov policy $\mu$ over time horizon $\overline{0,N}$ as the following:*

$$p_{s_0}^{\mu}(\mathcal{A}^c) = \mathbb{P}_{s_0}^{\mu}\{s(k) \in \mathcal{A}^c | \exists k \in \overline{0,N},\, s(0) = s_0\},$$

*where $s_0 \in \mathcal{A}$. The minimal probability of reaching the unsafe set is defined as:*

$$p_{*,s_0}(\mathcal{A}^c) = \inf\{p_{s_0}^{\mu}(\mathcal{A}^c),\, \mu \in \Pi_M^N\},$$

*where $\Pi_M^N$ is the set of all Markov policies over time horizon $\overline{0,N}$. A Markov policy $\mu_*$ is optimal with respect to an initial state $s_0$ if $p_{s_0}^{\mu_*}(\mathcal{A}^c) = p_{*,s_0}(\mathcal{A}^c)$.*

## 3. Design of Safe-visor

As discussed in the introduction, designing the Safe-visor for sandboxing unverified controllers consists of designing a Safety Advisor and a Supervisor. The Safety Advisor is designed regarding the safety specification. The Supervisor, on the other hand, is designed for detecting (potential) harmful behaviours of the unverified controller and accordingly deciding the input fed to the system (either the one from the unverified controller or from the Safety Advisor).

Regarding the safety invariance specification, an optimal safety controller can be designed as discussed in [APLS08, EZS14], which provides optimal safety policy to guarantee a minimal probability of reaching the

unsafe set in a finite time horizon. We use this controller as the Safety Advisor, which is introduced in details in Section 3.1.

Since the Safety Advisor only focuses on minimizing the probability of reaching the unsafe set, we need to turn to the unverified controller for functionality. Nevertheless, we still expect a high level of safety for the system. Therefore, we denote by $\rho$ the **maximal tolerable probability of reaching the unsafe set** that we are able to accept, which quantifies the compromise between functionality and safety. Given $\rho$, the Supervisor can decide whether it should accept inputs from an unverified controller at some time instants by estimating the probability of reaching the unsafe set and compare it with $\rho$. Details for designing a Supervisor working in this way is discussed in Section 3.2. It should be noted that the design of unverified controller is not the topic of this paper. The approach proposed here can be applied to any unverified controller as long as the set of all possible inputs provided by the unverified controller is a subset of the input set of the Supervisor. Moreover, we focus on those unverified controllers whose behaviour are unpredictable, i.e. we do not know the exact action of the unverified controller in a given state unless the system actually reaches that state, and the action of unverified controller at the same state may be time dependent. Otherwise, we may be able to verify this controller and sandboxing may not be needed anymore. In the rest of the paper, we denote by $u_{uc}(x,t)$ the input provided by the unverified controller at state $x$ at time instant $t$.

3.1. **Safety Advisor.** As mentioned above, we use optimal safety controller with respect to safety invariance specification as the Safety Advisor. To synthesize the optimal safety controller, we define a value function [TMKA13]:

$$V_n^\pi(x) := P_x^\pi(\diamond^{\leq n} \mathcal{A}^c), \tag{3.1}$$

for $n \in \mathbb{N}$ to denote the probability of reaching the set $\mathcal{A}^c$ in the finite time horizon $\overline{0,n}$ from the initial state $x$, where $\pi \in \Pi_M$ is a Markov policy and $\diamond^{\leq n}\mathcal{A}^c := \{\omega \in \Omega : \omega_x(k) \in \mathcal{A}^c \text{ for some } 0 \leq k \leq n\}$, where $\Omega$ is the set of all possible paths within the time horizon $\overline{0,n}$. Since we formulate the safety invariance specification as reach-avoid problem in a finite time horizon, we should minimize the probability mentioned above and, hence, the optimal value function is given by

$$V_{*,n}(x) := \inf_{\pi \in \Pi} V_n^\pi(x), \tag{3.2}$$

initialized with $V_{*,0} = 1_{\mathcal{A}^c}(x)$, where $1_M(x) = 1$ when $x \in M$ and $1_M(x) = 0$ otherwise. In fact, this optimal value function can be recursively calculated in the following way ([TMKA13], Corollary 3):

$$V_{*,n+1}(x) = 1_{\mathcal{A}^c}(x) + 1_{\mathcal{A}}(x) \inf_{u \in U} \int_X V_{*,n}(y) T_{\mathfrak{D}}(dy|x,u), \tag{3.3}$$

and the optimal policy at time $t = k$ associated to $V_{*,n+1}(x)$ can be obtained as the following

$$\mu_{*,k} \in \arg\inf_{\mu_k} \int_X (1_{\mathcal{A}^c}(y) + 1_{\mathcal{A}}(y) V_{*,n}(y)) T_{\mathfrak{D}}(dy|x,\mu_k(x)). \tag{3.4}$$

However, analytical solution of the value function as well as the optimal policy above is very difficult to be obtained in general. Alternatively, we abstract the original cdt-MP $\mathfrak{D}$ and construct a finite Markov Decision Process (MDP) as proposed in [TMKA13], and calculate the solutions based on this finite MDP. First, we use uniform grids to partition the safety region and input set of the cdt-MP. Let $X_p = \bigcup_{i=1}^N \tilde{X}_i$ be a measurable partition of the safety set $A$ and $U_p = \bigcup_{j=1}^M \tilde{U}_j$ a measurable partition of $U$. Let $\tilde{x}_i \in \tilde{X}_i$ for $1 \leq i \leq N$ be representative points of $\tilde{X}_i$ and let $\tilde{u}_j \in \tilde{U}_j$ for $1 \leq j \leq M$ be representative points of $\tilde{U}_j$. We define the discretization parameter $\delta_x = \max_{\tilde{x}_i, \tilde{x}_i' \in X_p} \mathbf{d}_X(\tilde{x}_i, \tilde{x}_i')$ for the state set and $\delta_u = \max_{\tilde{u}_j, \tilde{u}_j' \in U_p} \mathbf{d}_U(\tilde{u}_j, \tilde{u}_j')$ for the input set where $\mathbf{d}_X$ and $\mathbf{d}_U$ are the metrics (e.g. Euclidean ones) over sets $X$ and $U$, respectively. Then, the constructed finite MDP is denoted by $\mathfrak{M} = \{\tilde{X}, \tilde{U}, \tilde{T}\}$, in which $\tilde{X} := \{\tilde{x}_i\}_{i=1}^N \cup \{\phi\}$, $\{\tilde{x}_i\}_{i=1}^N$ is the set of representative points of $X_p$, $\phi$ is a "sink" state representing the unsafe set $\mathcal{A}^c$ in the original cdt-MP, and

$\tilde{U} := \{\tilde{u}_j\}_{j=1}^{M}$ is the set of representative points of $U_p$. The stochastic kernel $\tilde{T}$ is then a matrix, which can be computed as follows:
$$\tilde{T}(\tilde{x}_m|\tilde{x}_i, \tilde{u}_j) = \begin{cases} T_{\mathfrak{D}}(\tilde{X}_m|\tilde{x}_i, \tilde{u}_j) & \text{if } \tilde{x}_i, \tilde{x}_m \in \{\tilde{x}_i\}_{i=1}^{N}, \tilde{u}_j \in \tilde{U} \\ T_{\mathfrak{D}}(\mathcal{A}^c|\tilde{x}_i, \tilde{u}_j) & \text{if } \tilde{x}_i \in \{\tilde{x}_i\}_{i=1}^{N}, \tilde{x}_m \in \{\phi\}, \tilde{u}_j \in \tilde{U} \\ 1 & \text{if } \tilde{x}_i, \tilde{x}_m \in \{\phi\}, \tilde{u}_j \in \tilde{U} \\ 0 & \text{if } \tilde{x}_i \in \{\phi\}, \tilde{x}_m \in \{\tilde{x}_i\}_{i=1}^{N}, \tilde{u}_j \in \tilde{U} \end{cases}$$

For the finite MDP $\mathfrak{M}$, we denote by $\tilde{V}_{*,n}(\tilde{x})$ the n-horizon minimal value function for the reach-avoid problem. Similar to equation (3.3), we initialize it with $\tilde{V}_{*,0} = 1_{\{\phi\}}(\tilde{x})$ and it can be calculated recursively as follows:

$$\tilde{V}_{*,n+1}(\tilde{x}) = 1_{\{\phi\}}(\tilde{x}) + 1_{\{\phi\}^c}(\tilde{x}) \min_{\tilde{u}\in\tilde{U}} \sum_{\tilde{y}\in\tilde{X}} \tilde{V}_{*,n}(\tilde{y})\tilde{T}(\tilde{y}|\tilde{x}, \tilde{u}), \tag{3.5}$$

and the optimal policy at time $t = k$ associated to $\tilde{V}_{*,n+1}(\tilde{x})$ is given by

$$\mu_{*,k}(\tilde{x}) \in \arg\min_{\tilde{\mu}_k} \sum_{\tilde{y}\in\tilde{X}} (1_{\{\phi\}^c}(\tilde{y}) + 1_{\{\phi\}}(\tilde{y})\tilde{V}_{*,n}(\tilde{y}))\tilde{T}(d\tilde{y}|\tilde{x}, \tilde{\mu}_k(\tilde{x})). \tag{3.6}$$

In principle, the optimal policy can be obtained for arbitrary long time horizon, but $\tilde{V}_{*,n}(\tilde{x})$ will keep decreasing, i.e. the probability of avoiding the unsafe set is decreasing, when $n$ increases. Therefore, the time horizon of the optimal policy cannot be arbitrarily long, but it is tunable up to some degrees by setting the maximal tolerable value of the value function, i.e. the smaller (bigger) the maximal tolerable value of the value function is, the shorter (longer) the time horizon for the optimal safety policy is. This value should not be bigger than the maximal tolerable probability of reaching the unsafe set, i.e. $\rho$, as defined in the beginning of Section 3, so that $\rho$ can be guaranteed at least by accepting advisory input from the Safety Advisor. Therefore, in our implementation, the time horizon $\overline{0, H}$ of the Safety Advisor is determined in a way such that $\forall \tilde{x} \in \tilde{X}\backslash\{\phi\}, \tilde{V}_{*,H}(\tilde{x}) \leq \rho$ and $\exists \tilde{x} \in \tilde{X}\backslash\{\phi\}, \tilde{V}_{*,H+1}(\tilde{x}) > \rho$.

3.2. **Supervisor.** As previously mentioned, the Supervisor is required to estimate the probability of reaching the unsafe set, in case it accepts inputs from the unverified controller. Since the safety guarantee given by the Safety Advisor is calculated based on the abstraction of the original stochastic system, i.e. the finite MDP, for consistency of guarantee regarding safety probability, we use the same finite MDP to design the Supervisor.

As discussed in the previous section, the probability of reaching the unsafe set of the finite MDP is quantified by value function $\tilde{V}_H^\mu(s_0)$ according to equation (3.1). When the initial state $s_0$ and the time horizon $\overline{0, H}$ are fixed, the value function is varied by different $\mu$. Meanwhile, compared with purely using the optimal safety policy $\mu_*$, sandboxing the unverified controller and accepting it at some states at some time instants intrinsically result in a new Markov policy for controlling the system, according to the architecture of Safevisor. Therefore, to ensure that the probability of reaching the unsafe set is lower than the predefined $\rho$ in a given time horizon $\overline{0, H}$, the Supervisor should be designed in a way such that the following inequality holds:

$$\tilde{V}_H^{\mu'}(s_0) \leq \rho, \tag{3.7}$$

where $\mu'$ is the Markov policy used to control the system, when the unverified controller is accepted at some states at some time instants by the Supervisor. In Section 3.1, the optimal safety policy is obtained by selecting a Markov policy minimizing the value function of each state at each time instant. In other way, when the Markov policy is fixed, we can calculate the value function in the way illustrated in the next theorem.

**Theorem 3.1.** *Given a Markov policy $\mu = (\mu_0, \mu_1, \ldots, \mu_{H-1})$ in a finite time horizon $\overline{0, H}$, the value function $\tilde{V}_n(\tilde{x})$ can be recursively calculated in the following way:*

$$\tilde{V}_{n+1}(\tilde{x}) = 1_{\{\phi\}}(\tilde{x}) + 1_{\{\phi\}^c}(\tilde{x}) \sum_{\tilde{y}\in\tilde{X}} \tilde{V}_n(\tilde{y})\tilde{T}(\tilde{y}|\tilde{x}, \mu_{H-n-1}(\tilde{x})), \tag{3.8}$$

*where $\tilde{x} \in \tilde{X}$ and $\tilde{V}_0(x) = \tilde{V}_{*,0}(x)$.*

Theorem 3.1 can be proved similar to the proof of Lemma 1 in [APLS08], since Lemma 1 in [APLS08] can be treated as a general case for Theorem 3.1. With having Theorem 3.1, the remaining question is how to determine $\mu'$ at run time. Let $\mu' = (\mu'_0, \mu'_2, \ldots, \mu'_{H-1})$. When the Supervisor is being executed, at every time instant $k \in \overline{0, H-2}$, $\mu'_t$ are unknown for all $t$ where $k < t \leq H - 1$ (i.e., the Markov policy used to control the system in the future time is unknown). To guarantee the safety threshold specified by $\rho$, at every time instant $k \in \overline{0, H-2}$, input from the unverified controller can only be accepted, when inequality (3.7) is at least fulfilled in the case that the Supervisor only accepts the advisory input from the safety advisor afterwards. This requirement is formally defined in Definition 3.2.

**Definition 3.2.** *Given current time instant $k$, where $0 \leq k \leq H - 2$, $\omega$ is the path up to $k$ and $\rho$ is the maximal tolerable probability of reaching the unsafe set, the input $u_{uc}(\omega_x(k), k)$ from the unverified controller can only be accepted, if there exists a Markov policy $\mu = \{\mu_0, \mu_1, \mu_2, \ldots, \mu_{H-1}\} \in M$ such that $\tilde{V}_H^\mu(\omega_x(0)) \leq \rho$, where $M$ denotes the set of all Markov policies, $\mu_k(\omega_x(k)) = u_{uc}(\omega_x(k), k)$, and for all $t$ where $k < t \leq H - 1$, $\mu_t = \mu_{*,t}$.*

In general, it is difficult to calculate the exact value of $\tilde{V}_H^{\mu'}(s_0)$ at run time due to the lack of adequate information from the past. At each time instant $k$ during the execution, where $k \in \overline{0, H-2}$, the only available information for the Supervisor is the path $\omega$ of the system up to $k$. In other words, the Supervisor does not have complete information about $\mu'_t$ for all $t \in \overline{0, k}$, since $\mu'_t(\tilde{x})$ is unknown when $\tilde{x} \in \tilde{X} \setminus \{\omega_x(t)\}$. To cope with this difficulty, we propose a novel Supervisor, namely *History-based Supervisor*, as defined in Definition 3.3, which is able to check the feasibility of the input provided by the unverified controller only based on the history information during the execution (i.e., path $\omega$ of the system up to the current time instant $k$ during the execution).

**Definition 3.3.** *(History-based Supervisor) For all $k \in \overline{0, H-1}$ [1], given the history of path $\omega$ up to $k$, the input $u_{uc}(\omega_x(k), k)$ from the unverified controller can only be accepted, when quantity*

$$\prod_{t=1}^{k} \sum_{\tilde{x} \in \tilde{X} \setminus \{\phi\}} \tilde{T}(\tilde{x} | \omega_x(t-1), \omega_u(t-1)) \left( 1 - \sum_{\tilde{x} \in \tilde{X}} \tilde{V}_{*, H-k-1}(\tilde{x}) \tilde{T}(\tilde{x} | \omega_x(k), u_{uc}(\omega_x(k), k)) \right)$$

*is not smaller than $1 - \rho$, where $\rho$ is the maximal tolerable probability of reaching the unsafe set.*

By using History-based Supervisor in Safe-visor architecture, it can be guaranteed that $\mu'$ fulfils inequality (3.7), as illustrated in the next theorem.

**Theorem 3.4.** *Given a finite MDP and the unsafe set $A^c$, by using History-based Supervisor at $t$ for all $t \in \overline{0, H-1}$ in Safe-visor architecture, we have*

$$p_{s_0}^{\mu'}(\diamond^{\leq H} \mathcal{A}^c) \leq \rho,$$

*where $\mu'$ is the Markov policy used to control the system when History-based Supervisor is applied.*

Proof of Theorem 3.4 is provided in the appendix.

Note that $\tilde{V}_{*,n}$ and $\tilde{T}$ are calculated offline when synthesizing the Safety Advisor. Hence, the Supervisor defined in Definition 3.3 can be readily used in real-time, since the required computation can be efficiently performed. Concretely, at every time instant $k$ during the execution:

---

[1] No input needed to be provided at $t = H$ since it is the end of the execution.

(1) The number of operations required for computing $\prod_{t=1}^{k} \sum_{\tilde{x} \in \tilde{X} \setminus \{\phi\}} \tilde{T}(\tilde{x}|\omega_x(t-1), \omega_u(t-1))$ is constant, since

$$\prod_{t=1}^{k} \sum_{\tilde{x} \in \tilde{X} \setminus \{\phi\}} \tilde{T}(\tilde{x}|\omega_x(t-1), \omega_u(t-1))$$

$$= \sum_{\tilde{x} \in \tilde{X} \setminus \{\phi\}} \tilde{T}(\tilde{x}|\omega_x(k-1), \omega_u(k-1)) \times \prod_{t=1}^{k-1} \sum_{\tilde{x} \in \tilde{X} \setminus \{\phi\}} \tilde{T}(\tilde{x}|\omega_x(t-1), \omega_u(t-1))$$

$$= (1 - \tilde{T}(\phi|\omega_x(k-1), \omega_u(k-1)) \times \prod_{t=1}^{k-1} \sum_{\tilde{x} \in \tilde{X} \setminus \{\phi\}} \tilde{T}(\tilde{x}|\omega_x(t-1), \omega_u(t-1))$$

while $\prod_{t=1}^{k-1} \sum_{\tilde{x} \in \tilde{X} \setminus \{\phi\}} \tilde{T}(\tilde{x}|\omega_x(t-1), \omega_u(t-1))$ has already been computed at the previous time instant (i.e., $k-1$), and $\tilde{T}(\phi|\omega_x(k-1), \omega_u(k-1))$ can be directly obtained from $\tilde{T}$.

(2) The number of operations required for computing

$$1 - \sum_{\tilde{x} \in \tilde{X}} \tilde{V}_{*,H-k-1}(\tilde{x}) \tilde{T}(\tilde{x}|\omega_x(k), u_{uc}(k))$$

is proportional to the number of states of the finite MDP, since $\tilde{V}_{*,H-k-1}(\tilde{x})$ and $\tilde{T}(\tilde{x}|\omega_x(k), u_{uc}(k))$ can directly be obtained in $\tilde{V}_{*,n}$ and $\tilde{T}$.

The real time applicability of the proposed Supervisor is shown in the experiments in Section 4.

## 4. CASE STUDY

In this section, we apply our approach to two case studies. The first case study is a temperature control problem and the second one is a traffic control problem. We simulate each test case $1.0 \times 10^6$ times and analyze accordingly the percentage of paths staying in the safety set in the given time horizon. For comparison, we simulate these test cases by 1) only using the unverified controller and 2) only using the proposed safety advisor. Moreover, we compute the average execution time for our Supervisor in both cases to show feasibility of running it in real-time. The simulation in this section is performed in MATLAB 2018b, on a computer equipped with Intel(R) Xeon(R) E-2186G CPU (3.8 GHz) and 32 GB of RAM running Window 10.

4.1. **Temperature Control Problem.** In the temperature control problem, a room is equipped with a heater being controlled and the temperature of the room is required to be kept between 19 and 21°C. The temperature of the room can be modelled as the following, which is adapted from [LSZ18]:

$$x(k+1) = (1 - \beta - \gamma u(k))x(k) + \gamma T_h u(k) + \beta T_e + \omega(k) \tag{4.1}$$

where $x(k)$ denotes the temperature at time $t = k$. Input $u(k)$ takes any real value between 0 to 0.6. Parameter $\beta$ is conduction factor between the external environment and the room, $\gamma$ is conduction factor between the heater and the room, $T_e$ is the temperature of the external environment and $T_h$ is the temperature of the heater. We denote by $\omega$ a Gaussian white noise. In this section, we set $\beta = 0.022$, $\gamma = 0.05$, $T_e = -1$, $T_h = 50$, the mean of $\omega$ is 0 and variance is 0.04. The sampling time interval in this example is 9 minutes.

Now, we synthesize the Safety Advisor as discussed in Section 3.1. We use the discretization parameter $\delta_x = 1.0 \times 10^{-3}$ and $\delta_u = 2.4 \times 10^{-2}$ to discretize the safety set (resulting in 2000 discrete states) and the input set (resulting in 25 discrete inputs) to construct a finite MDP. We set $\rho$ as 1% and obtain a controller for time horizon $\overline{0, 40}$ (6 hours). We set the initial state at 19.01°C. The unverified controller tries to keep the heater idle at all time, i.e. $U_{uc}(t) \equiv 0$ for all $t \in \overline{0, 39}$. This is an unacceptable input which cools down the room to an unacceptable low level. For the given $\rho$, it is expected that at least 99% of the paths stay inside the safety set in the given time horizon. The result of the simulation is shown in Table 1 and Figure 3. The temperature keeps decreasing and all paths go outside of the safety set, when the system is fully controlled
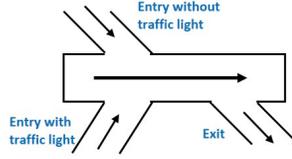
FIGURE 2. Traffic Control Problem

TABLE 1. Result of simulation for both case studies.

| | Temperature Control | Traffic Control |
|---|---|---|
| Percentage of paths in the safety set (with Safe-visor) | 99.02% | 99.958% |
| Average acceptance rate of the unverified controller | 19.12% | 8.5114% |
| Percentage of paths in the safety set (without Safe-visor) | 0% | 0% |
| Percentage of paths in the safety set (when system is fully controlled by the Safety Advisor) | 99.18% | 99.989% |
| Average execution time for the History-based Supervisor | 33.42 $\mu s$ | 31.83 $\mu s$ |

by the unverified controller. Meanwhile, more than 99% of the paths stay within the safety set when our proposed method is applied.

4.2. **Traffic Control Problem.** In the traffic control problem, we focus on a road traffic control containing a cell with 2 entries and 1 exit, as illustrated in Figure 2. One of the entry is controlled by a traffic light. The dynamic of the system can be modelled as the following, which is adapted from [LSZ19]:

$$x(k + 1) = (1 - \frac{\tau v}{l} - q)x(k) + e_1 u(k) + \sigma(k) + e_2, \tag{4.2}$$

where $x(k)$ denotes the density of traffic at time $k$, $u(k) \in \{0, 1\}$ is the input to the system (1 means the green light is on while 0 means the red light is on). Parameter $v$ is the flow speed of the vehicle on the road, $l$ is the length of the cell, $\sigma$ is a white Gaussian noise, and $\tau$ denotes the sampling time interval of the system. In one sampling interval, $e_1$ is the number of cars that pass the entry controlled by the traffic light, $e_2$ refers to the number of cars that pass the entry without traffic light, and $q$ is the percentage of cars which leave the cell through the exit. In the simulation, we set $l = 500$[m], $v = 25$[m/s], $\tau = 6$s, $e_1 = 3$, $e_2 = 6$, $q = 10\%$, the mean of $\sigma$ is 0 and variance is 2. In this case study, it is desired that the density of traffic is lower than 20.

Now, we synthesize the Safety Advisor as discussed in Section 3.1. We use the discretization parameter $\delta_x = 1.0 \times 10^{-3}$ to discretize the safety set (resulting in 20000 discrete states) to construct a finite MDP. Note that the input set is already finite. We set $\rho$ as 0.05% and obtain a controller for the time horizon $\overline{0, 8186}$ (13.64 hours). For the simulation, we set the initial state at $x = 9$, and choose the unverified controller as the following: $u_{uc}(t) = 0$ when $t \in \overline{0, 8186}$ is an odd number and $u_{uc}(t) = 1$ otherwise. For the given $\rho$, it is expected that at least 99.95% of the paths stay inside the safety set in the given time horizon. The result of the simulations is shown in Table 1 and Figure 4. All paths go outside of the safety set, when the system is fully controlled by the unverified controller. Meanwhile, more than 99.95% of paths stay within the safety set when our proposed method is applied.

According to the empirical result, by sandboxing the unverified controller with Safe-visor architecture, the probabilistic guarantees are respected while some of the inputs from the unverified controller are still accepted for functionality. The average execution time for the History-based Supervisor shows its good real-time applicability, which makes it practical to be applied in real time.
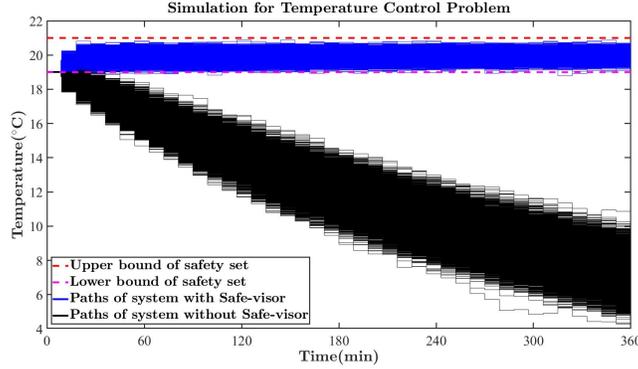
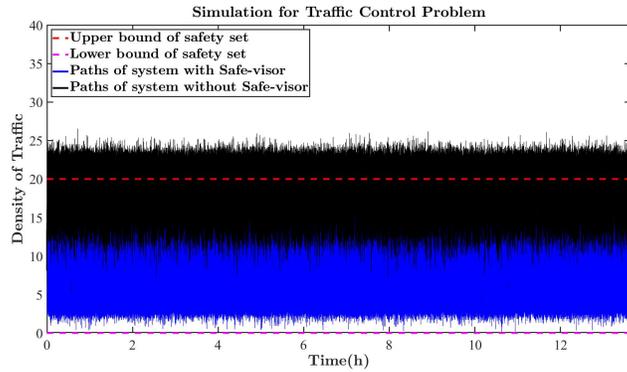FIGURE 3. Comparison between paths of system with and without Safe-visor (Temperature Control Problem).



FIGURE 4. Comparison between paths of system with and without Safe-visor (Traffic Control Problem)

## 5. CONCLUSION AND FUTURE WORK

In this paper, we developed a new framework for sandboxing unverified controllers for stochastic cyber-physical systems regarding safety invariance specification. In comparison with the Simplex architecture, our framework is applicable to stochastic systems, and provides more flexibility for the unverified controllers to accomplish complex tasks. According to the empirical results for two case studies, the pre-proposed safety probability is guaranteed by using our method. In the future, we would like to extend this method to 1) systems modelled by Partially Observable Markov Decision Processes[Mon82] 2) more general safety specifications, e.g. those expressed as co-safe linear temporal logic formulae [KV01].

## 6. ACKNOWLEDGEMENTS

The authors would like to thank Abolfazl Lavaei for the discussions on synthesizing optimal safety controllers for stochastic systems.

## 7. Appendix: Proof of Theorem 3.4

The proof of theorem 3.4 is done with the help of the following lemma.

**Lemma 7.1.** *Given a finite MDP $\mathfrak{M} = \{\tilde{X}, \tilde{U}, \tilde{T}\}$ and a Markov policy $\mu = (\mu_0, \mu_1, \ldots, \mu_{H-1})$ in a finite time horizon $\overline{0, H}$, we have*

$$1 - \tilde{V}_{n+1}(\tilde{x}) = \sum_{\tilde{y} \in \tilde{X} \setminus \{\phi\}} (1 - \tilde{V}_n(\tilde{y})) \tilde{T}(\tilde{y}|\tilde{x}, \mu_{H-n-1}(\tilde{x}))$$

*where $\tilde{V}_n(\tilde{x})$ is the value function for the reach-avoid problem and $\tilde{x} \in \tilde{X}$.*

The proof can be readily derived based on Theorem 3.1 and the definition of $\tilde{T}$. Let $\mu'$ be the Markov policy used to control the system when the unverified controller is accepted at some states at some time instants. Here, we use $\tilde{X}_s$ to represent $\tilde{X} \setminus \{\phi\}$. Let's define:

$$f(\tilde{x}(k), \mu'_k(\tilde{x}(k))) = 1 - \sum_{\tilde{x}(k+1) \in \tilde{X}_s} \tilde{V}_{*,H-k-1}(\tilde{x}(k+1)) \tilde{T}(\tilde{x}(k+1)|\tilde{x}(k), \mu'_k(\tilde{x}(k))),$$

and

$$g(\tilde{x}(k-1), \mu'_{k-1}(\tilde{x}(k-1))) = \tilde{T}(\tilde{x}(k)|\tilde{x}(k-1), \mu'_{k-1}(\tilde{x}(k-1))).$$

Given initial state $s_0 \in \tilde{X}_s$, at each time instant $t = k$ where $k \in \overline{0, H-1}$, we have

$$1 - \tilde{V}_H^{\mu'}(s_0)$$

$$= \sum_{\tilde{x}(1) \in \tilde{X}_s} \left( \sum_{\tilde{x}(2) \in \tilde{X}_s} \left( \cdots \left( \sum_{\tilde{x}(k) \in \tilde{X}_s} f(\tilde{x}(k), \mu'_k(\tilde{x}(k))) g(\tilde{x}(k-1), \mu'_{k-1}(\tilde{x}(k-1))) \right) \right. \right.$$
$$\left. \left. \cdots \right) g(\tilde{x}(1), \mu'_1(\tilde{x}(1))) \right) g(s_0, \mu'_0(s_0))$$

$$\geq \sum_{\tilde{x}(1) \in \tilde{X}_s} \left( \sum_{\tilde{x}(2) \in \tilde{X}_s} \left( \cdots \left( f(\underline{\tilde{x}(k)}, \underline{\mu'_k(\tilde{x}(k))}) \sum_{\tilde{x}(k) \in \tilde{X}_s} g(\tilde{x}(k-1), \mu'_{k-1}(\tilde{x}(k-1))) \right) \right. \right.$$
$$\left. \left. \cdots \right) g(\tilde{x}(1), \mu'_1(\tilde{x}(1))) \right) g(s_0, \mu'_0(s_0))$$

$$\geq \sum_{\tilde{x}(1) \in \tilde{X}_s} \left( \sum_{\tilde{x}(2) \in \tilde{X}_s} \left( \cdots \left( \left( \sum_{\tilde{x}(k) \in \tilde{X}_s} g(\underline{\tilde{x}(k-1)}, \underline{\mu'_{k-1}(\tilde{x}(k-1))}) \right) f(\underline{\tilde{x}(k)}, \underline{\mu'_k(\tilde{x}(k))}) \right. \right. \right.$$
$$\left. \left. \left. \sum_{\tilde{x}(k-1) \in \tilde{X}_s} g(\tilde{x}(k-2), \mu'_{k-2}(\tilde{x}(k-2))) \right) \cdots \right) g(\tilde{x}(1), \mu'_1(\tilde{x}(1))) \right) g(s_0, \mu'_0(s_0))$$

$$\geq \sum_{\tilde{x}(1) \in \tilde{X}_s} \left( \sum_{\tilde{x}(2) \in \tilde{X}_s} \left( \cdots \left( \left( \sum_{\tilde{x}(k-1) \in \tilde{X}_s} g(\underline{\tilde{x}(k-2)}, \underline{\mu'_{k-2}(\tilde{x}(k-2))}) \right) \right. \right. \right.$$
$$\left( \sum_{\tilde{x}(k) \in \tilde{X}_s} g(\underline{\tilde{x}(k-1)}, \underline{\mu'_{k-1}(\tilde{x}(k-1))}) \right) f(\underline{\tilde{x}(k)}, \underline{\mu'_k(\tilde{x}(k))})$$
$$\left. \left. \left. \sum_{\tilde{x}(k-2) \in \tilde{X}_s} g(\tilde{x}(k-3), \mu'_{k-3}(\tilde{x}(k-3))) \right) \cdots \right) g(\tilde{x}(1), \mu'_1(\tilde{x}(1))) \right) g(s_0, \mu'_0(s_0))$$
$$\cdots$$

$$\geq \prod_{t=1}^{k} \sum_{\tilde{x}(t) \in \tilde{X}_s} g(\underline{\tilde{x}(t-1)}, \underline{\mu_{t-1}(\tilde{x}(t-1))}) (f(\underline{\tilde{x}(k)}, \underline{\mu'_k(\tilde{x}(k))})$$

where

$$(\underline{\tilde{x}(t-1)}, \underline{\mu_{t-1}(\tilde{x}(t-1))}) = \mathop{\arg\min}_{\substack{\tilde{x}(t-1) \in \tilde{X}_s \\ \mu_{t-1}(\tilde{x}(t-1))}} \sum_{\tilde{x}(t) \in \tilde{X}_s} g(\tilde{x}(t-1), \mu_{t-1}(\tilde{x}(t-1)))$$

for all $t \in \overline{0,k}$, and

$$(\underline{\tilde{x}(k)}, \ \underline{\mu_k(\tilde{x}(k))}) = \underset{\substack{\tilde{x}(k) \in \tilde{X}_s \\ \mu_k(\tilde{x}(k))}}{\arg \min} f(\tilde{x}(k), \mu_k'(\tilde{x}(k))).$$

Noted that $\omega = (\underline{\tilde{x}(0)}, \underline{\mu_0(\tilde{x}(0))}, \underline{\tilde{x}(1)}, \underline{\mu_1(\tilde{x}(1))} \ldots \underline{\tilde{x}(k)})$ is one of the paths up to time instant $k$ which can be generated by the system controlled by the Markov policy $\mu'$, and the History-based Supervisor ensures that for all paths $\omega$ up to arbitrary time instant $k \in \overline{0,H}$,

$$\prod_{t=1}^{k} \sum_{\tilde{x} \in \tilde{X}_s} g(\omega_x(t-1), \omega_u(t-1)) \left( f(\omega_x(k), u_{uc}(\omega_x(k), k)) \right) \geq 1 - \rho.$$

Note that we have $1 - \tilde{V}_H^{\mu'}(s_0) \geq 1 - \rho$, i.e. $p_{s_0}^{\mu'}(\lozenge^{\leq H} \mathcal{A}^c) = \tilde{V}_H^{\mu'}(s_0) \leq \rho$.

## References

[ABE+18]  Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum, and Ufuk Topcu. Safe reinforcement learning via shielding. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

[ACH+18]  Fardin Abdi, Chien-Ying Chen, Monowar Hasan, Songran Liu, Sibin Mohan, and Marco Caccamo. Preserving physical safety under cyber attacks. *IEEE Internet of Things Journal*, 2018.

[APLS08]  Alessandro Abate, Maria Prandini, John Lygeros, and Shankar Sastry. Probabilistic reachability and safety for controlled discrete time stochastic hybrid systems. *Automatica*, 44(11):2724–2734, 2008.

[ATR+17]  Fardin Abdi, Rohan Tabish, Matthias Rungger, Majid Zamani, and Marco Caccamo. Application and system-level software fault tolerance through full system restarts. In *2017 ACM/IEEE 8th International Conference on Cyber-Physical Systems (ICCPS)*, pages 197–206. IEEE, 2017.

[BJCS14]  Stanley Bak, Taylor T Johnson, Marco Caccamo, and Lui Sha. Real-time reachability for verified simplex design. In *2014 IEEE Real-Time Systems Symposium*, pages 138–148. IEEE, 2014.

[BKKW15]  Roderick Bloem, Bettina Könighofer, Robert Könighofer, and Chao Wang. Shield synthesis. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 533–548. Springer, 2015.

[BMMC11]  Stanley Bak, Karthik Manamcheri, Sayan Mitra, and Marco Caccamo. Sandboxing controllers for cyber-physical systems. In *2011 IEEE/ACM Second International Conference on Cyber-Physical Systems*, pages 3–12. IEEE, 2011.

[CGR+07]  Tanya L Crenshaw, Elsa Gunter, Craig L Robinson, Lui Sha, and PR Kumar. The simplex reference model: Limiting fault-propagation due to unreliable components in cyber-physical system architectures. In *28th IEEE International Real-Time Systems Symposium (RTSS 2007)*, pages 400–412. IEEE, 2007.

[EZS14]  S Esmaeil Zadeh Soudjani. *Formal Abstractions for Automated Verification and Synthesis of Stochastic Systems*. PhD thesis, Technical University of Delft, 2014.

[HCK+11]  Naira Hovakimyan, Chengyu Cao, Evgeny Kharisov, Enric Xargay, and Irene M Gregory. L 1 adaptive control for safety-critical systems. *IEEE Control Systems Magazine*, 31(5):54–104, 2011.

[HKKT16]  Laura Humphrey, Bettina Könighofer, Robert Könighofer, and Ufuk Topcu. Synthesis of admissible shields. In *Haifa Verification Conference*, pages 134–151. Springer, 2016.

[HLL12]  Onésimo Hernández-Lerma and Jean B Lasserre. *Discrete-time Markov control processes: basic optimality criteria*, volume 30. Springer Science & Business Media, 2012.

[KV01]  Orna Kupferman and Moshe Y Vardi. Model checking of safety properties. *Formal Methods in System Design*, 19(3):291–314, 2001.

[LSZ18]  Abolfazl Lavaei, Sadegh Soudjani, and Majid Zamani. From dissipativity theory to compositional construction of finite markov decision processes. In *Proceedings of the 21st International Conference on Hybrid Systems: Computation and Control (part of CPS Week)*, pages 21–30. ACM, 2018.

[LSZ19]  Abolfazl Lavaei, Sadegh Soudjani, and Majid Zamani. Compositional synthesis of large-scale stochastic systems: A relaxed dissipativity approach. *arXiv preprint arXiv:1902.01223*, 2019.

[Mon82]  George E Monahan. State of the art—a survey of partially observable markov decision processes: theory, models, and algorithms. *Management Science*, 28(1):1–16, 1982.

[RBP09]  Charles Reis, Adam Barth, and Carlos Pizano. Browser security: lessons from google chrome. *Communications of the ACM*, 52(8):45–49, 2009.

[Sha01]  Lui Sha. Using simplicity to control complexity. *IEEE Software*, pages 20–28, 2001.

[TMKA13]  Ilya Tkachev, Alexandru Mereacre, Joost-Pieter Katoen, and Alessandro Abate. Quantitative automata-based controller synthesis for non-autonomous stochastic hybrid systems. In *Proceedings of the 16th international conference on Hybrid systems: computation and control*, pages 293–302. ACM, 2013.

[WHS13]  Xiaofeng Wang, Naira Hovakimyan, and Lui Sha. L1simplex: fault-tolerant control of cyber-physical systems. In *2013 ACM/IEEE International Conference on Cyber-Physical Systems (ICCPS)*, pages 41–50. IEEE, 2013.

[WHS18]     Xiaofeng Wang, Naira Hovakimyan, and Lui Sha. Rsimplex: A robust control architecture for cyber and physical
            failures. *ACM Transactions on Cyber-Physical Systems*, 2(4):27, 2018.
[YLZS13]    Jianguo Yao, Xue Liu, Guchuan Zhu, and Lui Sha. Netsimplex: Controller fault tolerance architecture in networked
            control systems. *IEEE Transactions on Industrial Informatics*, 9(1):346–356, 2013.

[2]DEPARTMENT OF MECHANICAL ENGINEERING, TECHNICAL UNIVERSITY OF MUNICH, GERMANY.

*Email address*: ingzhuo.zhong@tum.de

*Email address*: mcaccamo@tum.de

[2]DEPARTMENT OF COMPUTER SCIENCE, LMU MUNICH, GERMANY.

[3]DEPARTMENT OF COMPUTER SCIENCE, UNIVERSITY OF COLORADO BOULDER, USA.

*Email address*: majid.zamani@colorado.edu