# Provably Correct Design of Observations for Fault Detection with Privacy Preservation

Zhe Xu, Sayan Saha and Agung Julius

*Abstract*— During the operation of complex cyber-physical systems, detection of faults needs to be performed using limited state information for practicality and privacy concerns. While a well-designed observation can distinguish a faulty behavior from the normal behavior, it can also represent the action of hiding some of the state information or discrete mode transitions. In this paper, we present a framework for constructing the observation maps in the form of metric temporal logic (MTL) formulae that can be formally proven to detect fault in a switched system while preserving certain privacy conditions. We simulate finitely many nominal trajectories and use the robustness tubes around the simulated trajectories to cover the infinite trajectories that constitute the system behavior. Thus the inferred MTL formulae from the simulated trajectories can be used for classifying the system behaviors in a provably correct fashion. We implement our approach on the simulation model of a smart building testbed to detect the open window fault while preserving the privacy of the room occupancy.

## I. Introduction

Consider the task of designing a monitoring system that can detect faults during the operation of a safety critical cyber-physical system. To do so, the monitor needs to collect enough information from the cyber-physical system that can distinguish potential faulty operation from normal operation. On the other hand, practicality and privacy limit the amount of information that can be collected. For example, the number of sensors that can be deployed may be limited, or certain information is withheld by the system owner for privacy concerns. It is imperative to determine how information can be extracted from the operation of a cyber-physical system to enable fault detection, while respecting limitations such as privacy.

Presently, there are mainly two categories of approaches for fault detection (identifying whether a fault has occurred). The first category relies on the pattern recognition of sensor readings using reasoning or learning based techniques, such as neural networks [1], support vector machine [2], etc. The second category relies on the mathematical model of systems as they compare available measurements with information analytically derived from the system model. As modeled, whenever the difference between the actual system's output and the estimated output (the residual) [3] is above a certain threshold value, one can assume that a fault has occurred. For hybrid systems and switched systems, the main challenge of the model-based fault detection is due to the difficulty in capturing the combined continuous and discrete measurements.

In this paper, we propose an approach that utilizes both the pattern recognition techniques and the model-based methods for switched system fault detection and privacy preservation. We first approximate the switched system behavior using finitely many simulated execution trajectories of the systems. With the notion of trajectory robustness, we provide a guarantee on how far the system's state trajectories can deviate as a result of initial state variations [4]. Then we design an observation map that projects the different behaviors (normal, faulty, presence and absence of the privacy conditions) into an observation space where the images of the normal behavior and the faulty behavior are separate (fault detection) while the images of the behavior with the presence and absence of the privacy conditions are close (privacy preservation). As the fault detection is essentially a classification between the normal behavior and the faulty behavior, we modify the methods in [5], [6], [7], [8], [9] to automatically infer metric temporal logic (MTL) formulae directly from the simulated trajectories to classify the faulty and normal trajectories while considering the initial state variations and preserving the privacy conditions. We extend the previous works above in the following aspects: (i) instead of classifying trajectories in different sets, we classify robustness tubes of trajectories with initial state variations and disturbances; (ii) in performing the classification for fault detection (such as detecting the open window fault in a smart building), we also consider the preservation of privacy conditions (such as the room occupancy).

## II. Preliminaries

### A. Switched Systems

*Definition 1 (Switched System):* A switched system is a 5-tuple $\mathcal{S} = (\mathcal{Q}, \mathcal{X}, \mathcal{X}^0, \mathcal{F}, \mathcal{E})$ where

- $\mathcal{Q} = \{1, 2, \ldots, M\}$ is the set of indices for the modes (or subsystems).
- $\mathcal{X}$ is the domain of the continuous state, $x \in \mathcal{X}$ is the continuous state of the system, $\mathcal{X}^0 \subset \mathcal{X}$ is the initial set of states.
- $\mathcal{F} = \{f_q | q \in \mathcal{Q}\}$ where $f_q$ describes the continuous time-invariant dynamics for the mode $\dot{x} = f_q(x)$, which is assumed to admit a unique global solution $\xi_q(t, x_q^0)$, where $\xi_q$ satisfies $\frac{\partial \xi_q(t, x_q^0)}{\partial t} = f_q(\xi_q(t, x_q^0))$, and $\xi_q(0, x_q^0) = x_q^0$ is an initial condition in mode $q$.
- $\mathcal{E}$ is a subset of $\mathcal{Q} \times \mathcal{Q}$ which contains the valid transitions. If a transition $e = (q, q') \in \mathcal{E}$ takes place, the system switches from mode $q$ to $q'$.

Zhe Xu, Sayan Saha and Agung Julius are with the Department of Electrical, Computer, and Systems Engineering, Rensselaer Polytechnic Institute, Troy, NY 12180, USA e-mail: xuz8, sahas3, juliua2@rpi.edu.

*Remark 1:* Definition 1 is more restricted than that in [10], as we only consider state-independent transitions in this paper.

We assume that there is a minimal dwell time $\delta_{\min}$ between any two transitions to avoid Zeno behaviors.

*Definition 2 (Trajectory):* A trajectory of a switched system $\mathcal{S}$ is denoted as a sequence $\rho = \{(q^i, x^i, T^i)\}_{i=0}^{N_q}$ ($N_q \in \mathbb{N}$), where

- $\forall i \geq 0$, $q^i \in \mathcal{Q}$, $x^i \in \mathcal{X}$, $x^0 \in \mathcal{X}^0$ is the initial state, $x^{i+1} = \xi_{q^i}(T^i, x^i)$ is the initial state at mode $q^{i+1}$;
- $\forall i \geq 0$, $T^i \geq \delta_{\min} > 0$ is the dwell time at mode $q^i$, while the transition times are $T^0, T^0 + T^1, \ldots, T^0 + T^1 + \cdots + T^{N_q-1}$;
- $\forall i \geq 0$, $(q^i, q^{i+1}) \in \mathcal{E}$.

*Definition 3 (Output Trajectory):* For a trajectory $\rho = \{(q^i, x^i, T^i)\}_{i=0}^{N_q}$, we define the output trajectory $s_\rho$ as follows:

$$s_\rho(t) = \begin{cases} \xi_{q^0}(t, x^0), & \text{if } t < T^0, \\ \xi_{q^i}(t - \sum_{k=0}^{i-1} T^k, x^i), \\ \qquad \text{if } \sum_{k=0}^{i-1} T^k \leq t < \sum_{k=0}^{i} T^k, 1 \leq i \leq N_q. \end{cases}$$

### B. Metric Temporal Logic

In this section, we briefly review the metric temporal logic [11]. We focus on MTL formulae that are interpreted over an output trajectory of a switched system. The continuous state of the system we are studying is described by a set of $n$ variables that can be written as a vector $x = \{x_1, x_2, \ldots, x_n\}$. The domain of $x$ is denoted by $\mathcal{X}$. The domain $\mathbb{B} = \{true, false\}$ is the Boolean domain and the time set is $\mathbb{T} = \mathbb{R}_{\geq 0}$. The output trajectory $s_\rho$ of a switched system is defined in Sec. II-A. A set $AP = \{\pi_1, \pi_2, \ldots \pi_n\}$ is a set of atomic propositions, each mapping $\mathcal{X}$ to $\mathbb{B}$. The syntax of MTL is defined recursively as follows:

$$\phi := \top \mid \pi \mid \neg\phi \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2 \mid \phi_1 \mathcal{U}_\mathcal{I} \phi_2$$

where $\top$ stands for the Boolean constant True, $\pi$ is an atomic proposition, $\neg$ (negation), $\wedge$ (conjunction), $\vee$ (disjunction) are standard Boolean connectives, $\mathcal{U}$ is a temporal operator representing "until", $\mathcal{I}$ is a time interval of the form $I = [i_1, i_2)$. We can also derive two useful temporal operators from "until"($\mathcal{U}$), which are "eventually"$\Diamond_\mathcal{I}\phi = \top\mathcal{U}_\mathcal{I}\phi$ and "always"$\Box_\mathcal{I}\phi = \neg\Diamond_\mathcal{I}\neg\phi$. We define the set of states that satisfy the atomic proposition $\pi$ as $\mathcal{O}(\pi) \in \mathcal{X}$.

For a set $S \subseteq \mathcal{X}$, we define the signed distance from $x$ to $S$ as

$$\mathbf{Dist}_d(x, S) \triangleq \begin{cases} -\inf\{d(x,y) | y \in cl(S)\}, & \text{if } x \notin S, \\ \inf\{d(x,y) | y \in \mathcal{X} \setminus S\}, & \text{if } x \in S. \end{cases} \quad (1)$$

where $d$ is a metric on $\mathcal{X}$ and $cl(S)$ denotes the closure of the set $S$. In this paper, we use the metric $d(x,y) = \|x - y\|$, where $\|\cdot\|$ denotes the 2-norm.

If the output trajectory $s_\rho$ is only defined on a finite time interval ($\triangleq \mathbb{T}_o$), then the time domain of the formula $\phi$ is

also finite. It is defined recursively as follows [12]:

$$D(\pi, s_\rho) = \mathbb{T}_o,$$
$$D(\neg\phi, s_\rho) = D(\phi, s_\rho),$$
$$D(\phi_1 \wedge \phi_2, s_\rho) = D(\phi_1, s_\rho) \cap D(\phi_2, s_\rho),$$
$$D(\phi_1 \mathcal{U}_\mathcal{I} \phi_2, s_\rho) = \{t | t + \mathcal{I} \subset (D(\phi_1, s_\rho) \cap D(\phi_2, s_\rho))\}.$$

The necessary length of a formula $\phi$, denoted as $\|\phi\|$, is defined as follows [12]:

$$\|\phi\| := \min\{T \mid \text{if } \mathbb{T}_o = [0, T], D(\phi, s_\rho) \neq \emptyset\}$$

We use $[[\phi]](s_\rho, t)$ to denote the robustness degree of the output trajectory $s_\rho$ with respect to the formula $\phi$ at time $t$. The robust semantics of a formula $\phi$ with respect to $s_\rho$ are defined recursively as follows [11]:

$$[[\top]](s_\rho, t) := +\infty,$$
$$[[\pi]](s_\rho, t) := \mathbf{Dist}_d(s_\rho(t), \mathcal{O}(\pi)),$$
$$[[\neg\phi]](s_\rho, t) := -[[\phi]](s_\rho, t),$$
$$[[\phi_1 \wedge \phi_2]](s_\rho, t) := \min\big([[\phi_1]](s_\rho, t), [[\phi_2]](s_\rho, t)\big),$$
$$[[\phi_1 \mathcal{U}_\mathcal{I} \phi_2]](s_\rho, t) := \max_{t' \in (t+\mathcal{I})} \Big(\min\big([[\phi_2]](s_\rho, t'),$$
$$\min_{t \leq t'' < t'} [[\phi_1]](s_\rho, t'')\big)\Big).$$

## III. FAULT DETECTION AND PRIVACY PRESERVATION BY APPROXIMATED BEHAVIOR AND OBSERVATION MAPS

We define the behavior $\mathfrak{B}$ of a switched system $\mathcal{S} = (\mathcal{Q}, \mathcal{X}, \mathcal{X}^0, \mathcal{F}, \mathcal{E})$ as the collection of all possible execution trajectories of $\mathcal{S}$ starting from the initial set $\mathcal{X}^0$. An external observation of the system can be formulated as a mapping $\Pi$ from the behavior $\mathfrak{B}$ to an observation space $\mathfrak{D}$.

We denote $q_{\text{fault}} \in \mathcal{Q}$ and $q_{\text{nom}} \in \mathcal{Q}$ as the faulty modes and normal modes, i.e. the modes that represent presence and absence of fault respectively. For a certain trajectory of $\mathfrak{B}$, we assume that the fault occurs at most once and would not disappear by itself after the occurrence. We denote $q_\sigma \in \mathcal{Q}$ and $q_{\bar\sigma} \in \mathcal{Q}$ as the modes that represent presence and absence of privacy condition $\sigma$ respectively.

*Definition 4:* Assume that $\mathfrak{B} = \mathfrak{B}_{\text{nom}} \cup \mathfrak{B}_{\text{fault}}$, $\mathfrak{B}_{\text{nom}} \cap \mathfrak{B}_{\text{fault}} = \emptyset$, where $\mathfrak{B}_{\text{nom}}$ is the set of all normal trajectories and $\mathfrak{B}_{\text{fault}}$ is the set of all faulty trajectories. We assume that all faulty trajectories extend at least until $\delta_d$ time units after a fault occurrence. By the observation map $\Pi : \mathfrak{B} \to \mathfrak{D}$, a fault is $(\delta_{\text{d}}, \epsilon)$-detectable if $\forall \rho \in \mathfrak{B}_{\text{nom}}, \forall \rho' \in \mathfrak{B}_{\text{fault}}$, $d_O(\Pi(\rho), \Pi(\rho')) > \epsilon$, where $d_O$ is a metric in the observation space $\mathfrak{D}$.

*Remark 2:* From [13], a system is $(\delta_d, \delta_m)$-detectable if any fault can be detected $\delta_d$ time units after its occurrence and $\delta_m$ represents the observation accuracy of the time intervals with respect to a specific metric (called label sequence metric $d_{\Sigma^*}$). In this paper, we generalize $(\delta_d, \delta_m)$-detectability in [13] to $(\delta_d, \epsilon)$-detectability where $\epsilon$ is the observation accuracy with respect to any metric $d_O$ in the observation space $\mathfrak{D}$.

*Definition 5:* Assume that $\mathfrak{B} = \mathfrak{B}_\sigma \cup \mathfrak{B}_{\bar\sigma}$, $\mathfrak{B}_\sigma \cap \mathfrak{B}_{\bar\sigma} = \emptyset$, where $\mathfrak{B}_\sigma$ is the set of all trajectories where the privacy

condition $\sigma$ is present for some time during the execution and $\mathfrak{B}_{\bar{\sigma}}$ is the set of all trajectories where the privacy condition $\sigma$ is absent during the execution. By the observation map $\Pi : \mathfrak{B} \to \mathfrak{D}$, the privacy condition $\sigma$ is $\epsilon$-preservable if $\forall \rho \in \mathfrak{B}_{\sigma}, \exists \rho' \in \mathfrak{B}_{\bar{\sigma}}, d_O(\Pi(\rho), \Pi(\rho')) \leq \epsilon$ and vice versa, where $d_O$ is a metric in the observation space $\mathfrak{D}$.

*Proposition 1:* If the fault is $(\delta_d, \epsilon)$-detectable, then it is $(\tilde{\delta}_d, \epsilon)$-detectable for all $\tilde{\delta}_d \geq \delta_d$. If the fault is $(\delta_d, \epsilon)$-detectable, then it is $(\delta_d, \tilde{\epsilon})$-detectable for all $\tilde{\epsilon} \leq \epsilon$. If the privacy condition $\sigma$ is $\epsilon$-preservable, then it is $\tilde{\epsilon}$-preservable for all $\tilde{\epsilon} \geq \epsilon$.

*Proof:* Straightforward from Definition 4 and 5. ∎

*Problem 1:* For $\mathfrak{B} = \mathfrak{B}_{\text{nom}} \cup \mathfrak{B}_{\text{fault}} = \mathfrak{B}_{\sigma} \cup \mathfrak{B}_{\bar{\sigma}}$, find an observation map $\Pi : \mathfrak{B} \to \mathfrak{D}$ and a positive number $\epsilon$ such that the fault is $(\delta_d, \epsilon)$-detectable and the privacy condition $\sigma$ are $\epsilon$-preservable in the observation space $\mathfrak{D}$.

Using the idea of approximate bisimulation, we can approximate the behavior $\mathfrak{B}$ (with infinitely many trajectories) with a set $\hat{\mathfrak{B}}$ of finitely many simulated trajectories [14], [15]. For each mode $q$, suppose that we can construct a continuously differentiable function $\Phi_q : \mathcal{X} \times \mathcal{X} \to \mathbb{R}_+$ as an autobisimulation function [16], which requires that for any $x, x' \in \mathcal{X}$,

$$\begin{aligned} \Phi_q(x, x') &\geq \|x - x'\|, \\ \nabla_x \Phi_q(x, x') f_q(x) + \nabla_{x'} \Phi_q(x, x') f_q(x') &\leq 0. \end{aligned} \quad (2)$$

We define $B_q(x, \gamma) \triangleq \{x' | \Phi_q(x, x') \leq \gamma\}$ as the robust neighborhood of $x$, $\gamma$ is referred to as the robustness radius. As the autobisimulation function is non-increasing through time, it is guaranteed that for any initial states $x^0, x'^0 \in \mathcal{X}$, if the initial distance $\Phi_q(x^0, x'^0) = \gamma^0$, then for any time $t$, $\xi_q(t, x'^0) \in B_q(\xi_q(t, x^0), \gamma^0)$.

*Definition 6 (Robustness Tube):* For a switched system $\mathcal{S}$, the robustness tube with the initial robustness radius $\gamma^0$ around a nominal (simulated) trajectory $\rho = \{(q^i, x^i, T^i)\}_{i=0}^{N_q}$, denoted as $Tube(\rho, \gamma^0)$, is defined as the set of trajectories that initiate from $B_{q^0}(x^0, \gamma^0)$ and have the same mode transition sequences and transition times with $\rho$.

The robustness tube can be over-approximated by the robust neighborhoods around the nominal (simulated) trajectory. We denote $\hat{\mathfrak{B}} = \{\rho^1, \cdots, \rho^{\hat{N}}\}$ as the set of nominal (simulated) trajectories, where for $k \in \{1, 2, \ldots, \hat{N}\}$, $\rho^k = \{(q^{i,k}, x^{i,k}, T^{i,k})\}_{i=0}^{N_q^k}$.

If for each mode $q$, the continuous dynamics is affine and stable, then there exists a quadratic autobisimulation function $\Phi_q(\xi_q(t, x_q^0), \xi_q(t, x)) = [(\xi_q(t, x_q^0) - \xi_q(t, x))^T M_q(\xi_q(t, x_q^0) - \xi_q(t, x))]^{\frac{1}{2}}$, where $M_q$ is positive definite. We denote $d_{\text{s}}(\rho', \rho) \triangleq \sup_{t \geq 0} d(s_{\rho'}(t), s_\rho(t))$ and $\vec{h}(\mathfrak{B}, \hat{\mathfrak{B}}) \triangleq \sup_{\rho' \in \mathfrak{B}} \inf_{\rho \in \hat{\mathfrak{B}}} d_{\text{s}}(\rho', \rho)$ as the directed Hausdorff distance from $\mathfrak{B}$ to $\hat{\mathfrak{B}}$. For a matrix $A$, we denote $\|A\|$ as the largest singular value of $A$.

*Proposition 2:* If $\mathfrak{B} \subset \bigcup_{1 \leq k \leq \hat{N}} Tube(\rho^k, \gamma^{0,k})$ for some positive numbers $\gamma^{0,k}$, and there exist sequences $\Gamma^k =$

$\{\gamma^{i,k}\}_{i=0}^{N_q^k}$ such that $\forall i \geq 1$, $\forall k \geq 1$,

$$B_{q^{i-1,k}}(\xi_{q^{i-1,k}}(T^{i-1,k}, x^{i-1,k}), \gamma^{i-1,k}) \subset B_{q^{i,k}}(x^{i,k}, \gamma^{i,k}),$$

where $B_{q^{i,k}}(\xi_{q^{i,k}}(t, x^{i,k}), \gamma^{i,k}) =$

$$\{z | [(z - \xi_{q^{i,k}}(t, x^{i,k}))^T M_{q^{i,k}} (z - \xi_{q^{i,k}}(t, x^{i,k}))]^{\frac{1}{2}} \leq \gamma^{i,k}\},$$

then $\vec{h}(\mathfrak{B}, \hat{\mathfrak{B}}) \leq \tilde{\gamma}_{\max}$, where $\tilde{\gamma}_{\max} = \max_{0 \leq i \leq N_q^k, 1 \leq k \leq \hat{N}}$ $\gamma^{i,k} \left\| M_{q^{i,k}}^{-\frac{1}{2}} \right\|$.

We denote $\hat{\mathfrak{B}}_{\text{nom}}$, $\hat{\mathfrak{B}}_{\text{fault}}$, $\hat{\mathfrak{B}}_{\sigma}$ and $\hat{\mathfrak{B}}_{\bar{\sigma}}$ as the sets of nominal (simulated) trajectories that approximate $\mathfrak{B}_{\text{nom}}$, $\mathfrak{B}_{\text{fault}}$, $\mathfrak{B}_{\sigma}$ and $\mathfrak{B}_{\bar{\sigma}}$ respectively.

*Proposition 3:* If $\vec{h}(\mathfrak{B}, \hat{\mathfrak{B}}) \leq \tilde{\gamma}_{\max}$, then for the observation map $\Pi(x) = x$ and the metric $d_{\text{s}}$ in the observation space $\mathfrak{D}$, the following statements are true:
1) If the fault is $(\delta_d, \epsilon)$-detectable for the approximated behavior $\hat{\mathfrak{B}} = \hat{\mathfrak{B}}_{\text{nom}} \cup \hat{\mathfrak{B}}_{\text{fault}}$ and $\epsilon \geq 2\tilde{\gamma}_{\max}$, then the fault is $(\delta_d, \epsilon - 2\tilde{\gamma}_{\max})$-detectable for the original behavior $\mathfrak{B} = \mathfrak{B}_{\text{nom}} \cup \mathfrak{B}_{\text{fault}}$.
2) If the condition $\sigma$ is $\epsilon$-preservable for the approximated behavior $\hat{\mathfrak{B}} = \hat{\mathfrak{B}}_{\sigma} \cup \hat{\mathfrak{B}}_{\bar{\sigma}}$, then the condition $\sigma$ is $(\epsilon + 2\tilde{\gamma}_{\max})$-preservable for the original behavior $\mathfrak{B} = \mathfrak{B}_{\sigma} \cup \mathfrak{B}_{\bar{\sigma}}$.

Next, we formulate the observation maps in the form of metric temporal logic.

*Theorem 1:* Given a trajectory $\rho = \{(q^i, x^i, T^i)\}_{i=0}^{N_q}$ of a switched system $\mathcal{S}$, assume that for any trajectory $\rho' = \{(q^i, x'^i, T^i)\}_{i=0}^{N_q} \in Tube(\rho, \gamma^0)$ and any $t \in [0, T^i]$,

$$\begin{aligned} \xi_{q^i}(t, x'^i) &\in B_{q^i}(\xi_{q^i}(t, x^i), \gamma^i) \\ &= \{z | [(z - \xi_{q^i}(t, x^i))^T M_{q^i} (z - \xi_{q^i}(t, x^i))]^{\frac{1}{2}} \leq \gamma^i\}. \end{aligned} \quad (3)$$

Then for any $\rho' \in Tube(\rho, \gamma^0)$ and any $t \in D(\phi, s_\rho)$,

$$[[\phi]](s_\rho, t) - \hat{\gamma}_{\max} \leq [[\phi]](s_{\rho'}, t) \leq [[\phi]](s_\rho, t) + \hat{\gamma}_{\max},$$

where $\phi$ is any MTL formula, $\hat{\gamma}_{\max} \triangleq \max_{0 \leq i \leq N_q} \hat{\gamma}^i$, $\hat{\gamma}^i = \gamma^i \left\| M_{q^i}^{-\frac{1}{2}} \right\|$.

*Definition 7:* The satisfaction signature of an output trajectory $s_\rho$ with respect to the MTL formula $\phi$ at time $t$, denoted as $\mu(s_\rho, \phi, t)$, is defined as follows: $\mu(s_\rho, \phi, t) = 1$ if $s_\rho$ satisfies $\phi$ at time $t$; $\mu(s_\rho, \phi, t) = 0$ if $s_\rho$ violates $\phi$ at time $t$.

*Definition 8:* The robustness signature of an output trajectory $s_\rho$ with respect to the MTL formula $\phi$ and a positive number $\gamma$ at time $t$, denoted as $\chi(s_\rho, \phi, \gamma, t)$, is defined as follows:

$$\chi(s_\rho, \phi, \gamma, t) = \begin{cases} \text{RS}, & \text{if } [[\phi]](s_\rho, t) > \gamma, \\ \text{NS}, & \text{if } 0 \leq [[\phi]](s_\rho, t) \leq \gamma, \\ \text{NV}, & \text{if } -\gamma \leq [[\phi]](s_\rho, t) < 0, \\ \text{RV}, & \text{if } [[\phi]](s_\rho, t) < -\gamma. \end{cases}$$

*Remark 3:* RS, NS, NV, RV are short for Robust Satisfaction, Nominal Satisfaction, Nominal Violation, Robust Violation respectively.

*Corollary 1:* With the same assumptions of Theorem 1, the following statements are true ($\hat{\gamma}_{\max} \triangleq \max_{0 \leq i \leq N_q} \hat{\gamma}^i$):

1) If $\chi(s_\rho, \phi, \hat{\gamma}_{\max}, t) = \text{RS}$, then $\mu(s_{\rho'}, \phi, t) = 1$ for any $\rho' \in Tube(\rho, \gamma^0)$;

2) If $\chi(s_\rho, \phi, \hat{\gamma}_{\max}, t) = \text{RV}$, then $\mu(s_{\rho'}, \phi, t) = 0$ for any $\rho' \in Tube(\rho, \gamma^0)$.

We design the MTL observation map as a monitor and the satisfaction (violation) of an output trajectory $s_\rho$ with respect to the MTL formula $\phi$ is checked at every time point in the time domain $D(\phi, s_\rho)$. We first simulate the sets of trajectories $\hat{\mathfrak{B}}_{\text{nom},\sigma} = \hat{\mathfrak{B}}_{\text{nom}} \cap \hat{\mathfrak{B}}_\sigma = \{\rho^1_{\text{nom},\sigma}, \cdots, \rho^{\hat{N}_{\text{nom},\sigma}}_{\text{nom},\sigma}\}$, $\hat{\mathfrak{B}}_{\text{nom},\bar{\sigma}} = \hat{\mathfrak{B}}_{\text{nom}} \cap \hat{\mathfrak{B}}_{\bar{\sigma}} = \{\rho^1_{\text{nom},\bar{\sigma}}, \cdots, \rho^{\hat{N}_{\text{nom},\sigma}}_{\text{nom},\bar{\sigma}}\}$, $\hat{\mathfrak{B}}_{\text{fault},\sigma} = \hat{\mathfrak{B}}_{\text{fault}} \cap \hat{\mathfrak{B}}_\sigma = \{\rho^1_{\text{fault},\sigma}, \cdots, \rho^{\hat{N}_{\text{fault},\sigma}}_{\text{fault},\sigma}\}$ and $\hat{\mathfrak{B}}_{\text{fault},\bar{\sigma}} = \hat{\mathfrak{B}}_{\text{fault}} \cap \hat{\mathfrak{B}}_{\bar{\sigma}} = \{\rho^1_{\text{fault},\bar{\sigma}}, \cdots, \rho^{\hat{N}_{\text{fault},\bar{\sigma}}}_{\text{fault},\bar{\sigma}}\}$. The trajectory $\rho^k_{\text{nom},\sigma} = \{(q^{i,k}, x^{i,k}, T^{i,k})\}_{i=0}^{N^k_{\text{nom},\sigma}}$ for $k \in \{1, 2, \ldots, \hat{N}_{\text{nom},\sigma}\}$ have robustness tubes around them, which can be over-approximated by robust neighborhoods with sequences of robustness radii $\Gamma^k_{\text{nom},\sigma} = \{\gamma^{i,k}_{\text{nom},\sigma}\}_{i=0}^{N^k_{\text{nom},\sigma}}$. Similarly, the sequences of robustness radii corresponding to $\rho^k_{\text{nom},\bar{\sigma}}$, $\rho^k_{\text{fault},\sigma}$ and $\rho^k_{\text{fault},\bar{\sigma}}$ are $\Gamma^k_{\text{nom},\bar{\sigma}}$ $(k = 1, 2, \ldots, \hat{N}_{\text{nom},\bar{\sigma}})$, $\Gamma^k_{\text{fault},\sigma}$ $(k = 1, 2, \ldots, \hat{N}_{\text{fault},\sigma})$ and $\Gamma^k_{\text{fault},\bar{\sigma}}$ $(k = 1, 2, \ldots, \hat{N}_{\text{fault},\bar{\sigma}})$, respectively. The MTL formula $\phi(\alpha)$ is chosen from a set of templates and $\alpha$ denote all the parameters that define the formula. The search starts from a basis of candidate formulae in the form of $\square_{[\tau_1, \tau_2]}\pi$ or $\lozenge_{[\tau_3, \tau_4]}\pi$ and adding Boolean connectives until a satisfactory formula is found. For each template of the formula $\phi$, we find the parameters $\alpha$ that minimize the following cost function:

$$
\begin{aligned}
J(\alpha) = & \sum_{k=1}^{\hat{N}_{\text{nom},\sigma}} g\big(-\max_{t[j] \in \mathbb{T}^k_{\text{nom},\sigma}} [[\phi(\alpha)]](s_{\rho^k_{\text{nom},\sigma}}, t[j]), \hat{\gamma}^k_{\text{nom},\sigma,\max}\big) \\
& + \sum_{k=1}^{\hat{N}_{\text{nom},\bar{\sigma}}} g\big(-\max_{t[j] \in \mathbb{T}^k_{\text{nom},\bar{\sigma}}} [[\phi(\alpha)]](s_{\rho^k_{\text{nom},\bar{\sigma}}}, t[j]), \hat{\gamma}^k_{\text{nom},\bar{\sigma},\max}\big) \\
& + \sum_{k=1}^{\hat{N}_{\text{fault},\sigma}} g\big(\max_{t[j] \in \mathbb{T}^k_{\text{fault},\sigma}} [[\phi(\alpha)]](s_{\rho^k_{\text{fault},\sigma}}, t[j]), \hat{\gamma}^k_{\text{fault},\sigma,\max}\big) \\
& + \sum_{k=1}^{\hat{N}_{\text{fault},\bar{\sigma}}} g\big(\max_{t[j] \in \mathbb{T}^k_{\text{fault},\bar{\sigma}}} [[\phi(\alpha)]](s_{\rho^k_{\text{fault},\bar{\sigma}}}, t[j]), \hat{\gamma}^k_{\text{fault},\bar{\sigma},\max}\big),
\end{aligned}
$$
$$(4)$$

where

$$\hat{\gamma}^k_{\text{nom},\sigma,\max} \triangleq \max_{0 \le i \le N^k_{\text{nom},\sigma}} \hat{\gamma}^{i,k}_{\text{nom},\sigma}, \quad \mathbb{T}^k_{\text{nom},\sigma} \triangleq D(\phi(\alpha), s_{\rho^k_{\text{nom},\sigma}}),$$

$$\hat{\gamma}^k_{\text{nom},\bar{\sigma},\max} \triangleq \max_{0 \le i \le N^k_{\text{nom},\bar{\sigma}}} \hat{\gamma}^{i,k}_{\text{nom},\bar{\sigma}}, \quad \mathbb{T}^k_{\text{nom},\bar{\sigma}} \triangleq D(\phi(\alpha), s_{\rho^k_{\text{nom},\bar{\sigma}}}),$$

$$\hat{\gamma}^k_{\text{fault},\sigma,\max} \triangleq \max_{0 \le i \le N^k_{\text{fault},\sigma}} \hat{\gamma}^{i,k}_{\text{fault},\sigma}, \quad \mathbb{T}^k_{\text{fault},\sigma} \triangleq D(\phi(\alpha), s_{\rho^k_{\text{fault},\sigma}}),$$

$$\hat{\gamma}^k_{\text{fault},\bar{\sigma},\max} \triangleq \max_{0 \le i \le N^k_{\text{fault},\bar{\sigma}}} \hat{\gamma}^{i,k}_{\text{fault},\bar{\sigma}}, \quad \mathbb{T}^k_{\text{fault},\bar{\sigma}} \triangleq D(\phi(\alpha), s_{\rho^k_{\text{fault},\bar{\sigma}}}).$$

The function $g(\cdot)$ is a penalty function defined as

$$g(x, \gamma) \triangleq \begin{cases} \gamma - x & \text{if } x \le \gamma, \\ 0 & \text{if } x > \gamma. \end{cases}$$
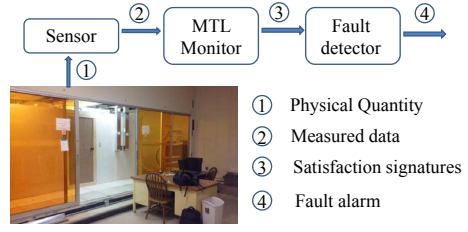


Fig. 1. Block diagram of the fault detector and the observation maps.

① Physical Quantity
② Measured data
③ Satisfaction signatures
④ Fault alarm

The cost function $J(\alpha)$ awards the robust satisfaction of $\phi(\alpha)$ by $\hat{\mathfrak{B}}_{\text{fault},\sigma}$ and $\hat{\mathfrak{B}}_{\text{fault},\bar{\sigma}}$ at some time points in the time domain and awards the robust violation of $\phi(\alpha)$ by $\hat{\mathfrak{B}}_{\text{nom},\sigma}$ and $\hat{\mathfrak{B}}_{\text{nom},\bar{\sigma}}$ at any time point in the time domain.

We use Particle Swarm Optimization (PSO)[17] to optimize $J(\alpha)$. We denote the parameters corresponding to the obtained optimal solution as $\alpha^*$. If $J(\alpha^*) = 0$, then we can design the fault detector as $\mu(s_\rho, \phi(\alpha^*), t)$, as whenever $\mu(s_\rho, \phi(\alpha^*), t)$ turns 1, the fault alarm is raised (see Fig. 1). Then with the observation map $\Pi(\rho) = \omega_\rho$, where $\omega_\rho = \sup_{t \ge 0} \mu(s_\rho, \phi(\alpha^*), t)$, $\rho$ is mapped to a binary value of 0 for the normal trajectories and a binary value of 1 for the faulty trajectories. For the metric $d_O(\omega_{\rho'}, \omega_\rho) \triangleq d(\omega_{\rho'}, \omega_\rho)$, the fault is $(t_\text{r} - t_\text{f} + \|\phi\|, \epsilon)$-detectable and the privacy condition is $\epsilon$-preservable for any $\epsilon \in (0, 1)$, where $t_\text{f}$ is the fault occurring time and $t_\text{r}$ is the first time instant $\mu(s_\rho, \phi(\alpha^*), t)$ turns 1. If $J(\alpha^*) > 0$, then we need to add another layer of observation map for the MTL monitor (see the example in Sec. IV).

## IV. IMPLEMENTATION

In this section, we implement our fault detection and privacy preservation method on detecting an open window fault of a room while preserving the privacy of the room occupancy in the simulation model of a smart building testbed [18]. The system is modeled as a switched system $\mathcal{S}$ with 4 different modes (we assume that closed and open window correspond to normal and faulty modes respectively, occupied room and unoccupied room correspond to presence and absence of the privacy condition $\sigma$ respectively), as shown in Fig. 2. We assume that at most one person will enter the room, the case for room occupation of more than one person can be done in a similar manner but with more modes. The state $x = [T, w]$ represents the temperature and humidity ratio of the room respectively. The continuous dynamics in the 4 modes are given as follows:

For mode $q_{\text{nom},\bar{\sigma}}$ (normal, unoccupied):

$$\begin{cases} C\dot{x}_1 = \dot{m}C_\text{p}(T_\text{s} - x_1) + \beta G(x_2 - w_\infty) - K(x_1 - T_\infty) + \varepsilon_1; \\ M\dot{x}_2 = \dot{m}(w_\text{s} - x_2) - G(x_2 - w_\infty) + \varepsilon_2. \end{cases}$$

For mode $q_{\text{nom},\sigma}$ (normal, occupied):

$$\begin{cases} C\dot{x}_1 = \dot{m}C_{\text{p}}(T_{\text{s}} - x_1) + \beta G(x_2 - w_\infty) - K(x_1 - \\ \qquad T_\infty) + \dot{Q}_{\text{gen}} - \beta\dot{W}_{\text{gen}} + \varepsilon_1; \\ M\dot{x}_2 = \dot{m}(w_{\text{s}} - x_2) - G(x_2 - w_\infty) + \dot{W}_{\text{gen}} + \varepsilon_2. \end{cases}$$

For mode $q_{\text{fault},\bar{\sigma}}$ (fault, unoccupied):

$$\begin{cases} C\dot{x}_1 = \dot{m}C_{\text{p}}(T_{\text{s}} - x_1) + \beta G_{\text{f}}(x_2 - w_\infty) - K_{\text{f}}(x_1 \\ \qquad - T_\infty) + \varepsilon_1; \\ M\dot{x}_2 = \dot{m}(w_{\text{s}} - x_2) - G_{\text{f}}(x_2 - w_\infty) + \varepsilon_2. \end{cases}$$

For mode $q_{\text{fault},\sigma}$ (fault, occupied):

$$\begin{cases} C\dot{x}_1 = \dot{m}C_{\text{p}}(T_{\text{s}} - x_1) + \beta G_{\text{f}}(x_2 - w_\infty) - K_{\text{f}}(x_1 \\ \qquad - T_\infty) + \dot{Q}_{\text{gen}} - \beta\dot{W}_{\text{gen}} + \varepsilon_1; \\ M\dot{x}_2 = \dot{m}(w_{\text{s}} - x_2) - G_{\text{f}}(x_2 - w_\infty) + \dot{W}_{\text{gen}} + \varepsilon_2. \end{cases}$$

where $\dot{m}$ is the mass flow rate, $C$ is the thermal capacitance of the room, $M$ is mass of air in the room, $G$ and $G_{\text{f}}$ are the mass transfer conductance between the room and the ambient for the normal and faulty situations respectively, $w_{\text{s}}$, $T_{\text{s}}$ are the supply air humidity ratio and temperature respectively, $w_\infty$, $T_\infty$ are the ambient humidity ratio and temperature respectively, $\dot{W}_{\text{gen}}$ and $\dot{Q}_{\text{gen}}$ are humidity and heat generation within the room (i.e. from human), $C_{\text{p}}$ is specific heat of air at constant pressure, $\beta$ is latent heat of vaporization of water, $K$ and $K_{\text{f}}$ is the wall thermal conductance for the normal and faulty situations respectively. The disturbances $\varepsilon_1$ and $\varepsilon_2$ are added to account for the modeling uncertainties.

We set $T_\infty = 303\text{K}$ (29.85°C), $T_{\text{s}} = 290\text{K}$ (16.85°C), $w_\infty = 0.0105$, $w_{\text{s}} = 0.01$. As the ambient temperature and humidity ratio are both higher than the room temperature and humidity ratio, when there is a fault (the window is open), the room temperature and humidity ratio will increase towards the new equilibrium. When the room becomes occupied, the room temperature and humidity ratio will also increase. We use $x_{\text{nom},\bar{\sigma}}, x_{\text{nom},\sigma}, x_{\text{fault},\bar{\sigma}}, x_{\text{fault},\sigma}$ to denote the equilibrium states of the 4 modes. It can be seen from Fig. 3 (a) that $x_{\text{fault},\bar{\sigma}}$ and $x_{\text{nom},\sigma}$ are almost the same, therefore a temporal logic formula is a proper observation to distinguish their temporal patterns in the transient period.
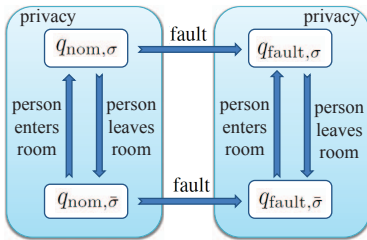


Fig. 2. The switched system $\mathcal{S}$ for the smart building model.

We set the minimal dwell time $\delta_{\min} = 480\text{s}$. We simulate the following nominal trajectories in the time window $\mathbb{T}_o = [t_{\text{s}} - 240\text{s}, t_{\text{s}} + 480\text{s}]$ (we assume that there is a transition at

$t_{\text{s}}$ and $t_{\text{s}} = 240\text{s}$):

$$\begin{aligned} \rho^1_{\text{nom},\bar{\sigma}} &= (q_{\text{nom},\bar{\sigma}}, x_{\text{nom},\bar{\sigma}}, T_{\text{end}}), \\ \rho^1_{\text{nom},\sigma} &= (q_{\text{nom},\sigma}, x_{\text{nom},\sigma}, T_{\text{end}}), \\ \rho^2_{\text{nom},\sigma} &= (q_{\text{nom},\bar{\sigma}}, x_{\text{nom},\bar{\sigma}}, t_{\text{s}}), (q_{\text{nom},\sigma}, x_{\text{nom},\bar{\sigma}}, T_{\text{end}} - t_{\text{s}}), \\ \rho^3_{\text{nom},\sigma} &= (q_{\text{nom},\sigma}, x_{\text{nom},\sigma}, t_{\text{s}}), (q_{\text{nom},\bar{\sigma}}, x_{\text{nom},\sigma}, T_{\text{end}} - t_{\text{s}}), \\ \rho^1_{\text{fault},\bar{\sigma}} &= (q_{\text{nom},\bar{\sigma}}, x_{\text{nom},\bar{\sigma}}, t_{\text{s}}), (q_{\text{fault},\bar{\sigma}}, x_{\text{nom},\bar{\sigma}}, T_{\text{end}} - t_{\text{s}}), \\ \rho^1_{\text{fault},\sigma} &= (q_{\text{nom},\sigma}, x_{\text{nom},\sigma}, t_{\text{s}}), (q_{\text{fault},\sigma}, x_{\text{nom},\sigma}, T_{\text{end}} - t_{\text{s}}). \end{aligned}$$

We set $T_{\text{end}} = 720\text{s}$. Note that the transitions may occur at different times than $t_{\text{s}}$, but then they can always be shifted to the trajectories in the robustness tube around the simulated trajectories due to time invariance. As the states may not be at the equilibrium at time $(t_{\text{s}} + 240\text{s})$ (see points $x'_{\text{nom},\bar{\sigma}}$ and $x'_{\text{nom},\sigma}$ in Fig. 3 (a)), there exist some initial state variations as well at time $(t_{\text{s}} - 240\text{s})$ due to a possible transition at time $(t_{\text{s}} - 480\text{s})$ or earlier. Another reason for the state variations come from the disturbances. Therefore, we equip the nominal trajectories with robust neighborhoods that cover the initial state variations and disturbances, as shown in Fig. 3 (b) (using Theorem 1 of [19] and applying zero input, we can find an autobisimulation function such that both the initial state variations and the disturbances can be covered with the robustness radii $\gamma^{i,k}_{\text{nom},\sigma} = \gamma^{i,k}_{\text{nom},\bar{\sigma}} = \gamma^{i,k}_{\text{fault},\sigma} = \gamma^{i,k}_{\text{fault},\bar{\sigma}} = 0.002$). Although we only simulate the trajectories for 720s, the robustness tubes around the normal trajectories starting from time $(t_{\text{s}} + 240\text{s})$ are repeating the robustness tubes starting from time $(t_{\text{s}} - 240\text{s})$ until another transition occurs. In this way, if we constrain $\|\phi\| < 240$ (which makes $[t_{\text{s}} - \|\phi\|, t_{\text{s}} + \|\phi\|] \subset D(\phi, s_\rho)$), then the robust neighborhoods around the simulated trajectories approximate all the possible behaviors with time invariance for the MTL monitor. We use the data of the simulated room temperature to design the MTL monitor and the case for the room humidity ratio can be done in a similar manner. The optimal formula is obtained as follows (the time unit is second):

$$\begin{aligned} \phi(\alpha^*) = (\square_{[160,180]}T > 290.75) \vee ((\Diamond_{[0,20)}T < 290.516) \\ \wedge (\Diamond_{[39,65)}T > 290.525) \wedge (\Diamond_{[122,161)}T < 290.625)). \end{aligned}$$

The robust signatures of the simulated trajectories with respect to $\phi(\alpha^*)$ and $\hat{\gamma}^k_{\text{nom},\sigma,\max} = \hat{\gamma}^k_{\text{nom},\bar{\sigma},\max} = \hat{\gamma}^k_{\text{fault},\sigma,\max} = \hat{\gamma}^k_{\text{fault},\sigma,\bar{\max}} = 0.002$ are shown in Fig. 4. The robust signatures of all the simulated faulty output trajectories are RS for at least 19s (which means the satisfaction signatures of all the faulty output trajectories in the robustness tube are 1 for at least 19s) while the robust signatures of the simulated normal output trajectories are not RV for at most 11s (which means the satisfaction signatures of all the normal output trajectories in the robustness tube are 1 for at most 11s). Therefore, we design the fault detector $\mu(s_\rho, \phi', t)$, where $\phi' = \square_{[0,t_c]}\phi(\alpha^*)$ for any $t_c \in (11s, 19s]$. Whenever the fault detector detects 1, the fault alarm is raised. Then with the observation map $\Pi(\rho) = \Omega_\rho$, where $\Omega_\rho = \sup_{t \geq 0} \mu(s_\rho, \phi', t)$, the fault is $(148s + t_c, \epsilon)$-detectable and $\sigma$ is $\epsilon$-preservable for any $\epsilon \in (0, 1)$.

The inference takes about 27 minutes on a Dell desktop computer with a 3.20 GHz Intel Xeon CPU and 8 GB RAM. So the inference process can be run every 30 minutes using the updated ambient temperature and humidity values.
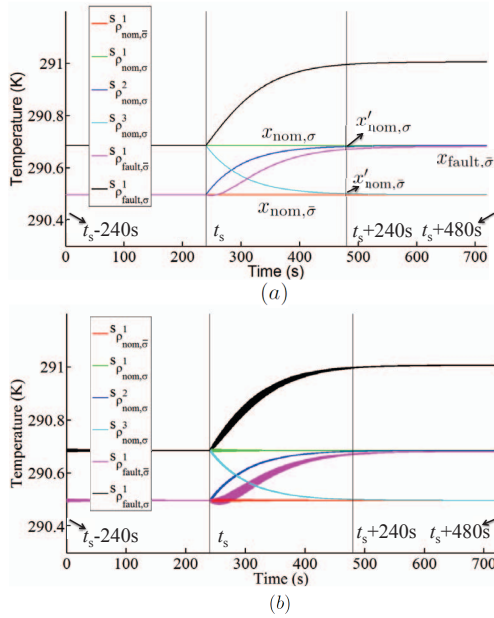


Fig. 3. The simulated output trajectories of the room temperature when $\delta_{\min}=480s$ in the time window [0, 720s] (a) and the robustness tubes around the simulated trajectories (b).
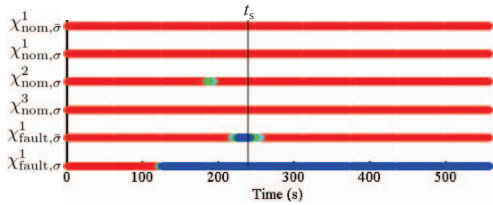


Fig. 4. Robust signatures $\chi(s_\rho, \phi(\alpha^*), 0.002, t)$ (denoted as $\chi_\rho$ for brevity) of the simulated output trajectories when $\delta_{\min}=480s$ (the colors blue, green, cyan and red represent RS, NS, NV and RV, respectively).

## V. CONCLUSION

In this paper, we present a method for constructing the observation maps in the form of metric temporal logic (MTL) formulae for fault detection and privacy preservation in a provably correct fashion. The results from a smart building testbed show that the open window fault can be detected while the privacy information about the room occupancy is preserved. In the future work, we can extend the theory of privacy preservation to include multiple privacy conditions.

## REFERENCES

[1] G. Cardoso, J. Rolim, and H. Zurn, "Application of neural-network modules to electric power system fault section estimation," *IEEE Trans. Power Delivery*, vol. 19, no. 3, pp. 1034–1041, July 2004.

[2] D. Thukaram, H. Khincha, and B. Ravikumar, "An intelligent approach using support vector machines for monitoring and identification of faults on transmission systems," in *IEEE Power India Conference*, 2006, pp. 184–190.

[3] F. Chowdhury and J. Aravena, "A modular methodology for fast fault detection and classification in power systems," *IEEE Trans. Control Syst. Technology*, vol. 6, no. 5, pp. 623–634, Sep 1998.

[4] Z. Xu, A. A. Julius, and J. H. Chow, "Robust testing of cascading failure mitigations based on power dispatch and quick-start storage," *IEEE Systems Journal*, Early access on IEEE Xplore.

[5] E. Asarin, A. Donzé, O. Maler, and D. Nickovic, "Parametric identification of temporal properties," in *Proc. of the Second Int. Conf. on Runtime Verification*, ser. RV'11.  Berlin, Heidelberg: Springer-Verlag, 2012, pp. 147–160.

[6] Z. Kong, A. Jones, A. Medina Ayala, E. Aydin Gol, and C. Belta, "Temporal logic inference for classification and prediction from data," in *Proc. of the 17th Int. Conf. on Hybrid Systems: Computation and Control*, ser. HSCC '14.  New York, NY, USA: ACM, 2014, pp. 273–282.

[7] B. Hoxha, A. Dokhanchi, and G. Fainekos, "Mining parametric temporal logic properties in model-based design for cyber-physical systems," *International Journal on Software Tools for Technology Transfer*, Feb 2017. [Online]. Available: http://dx.doi.org/10.1007/s10009-017-0447-4

[8] Z. Xu, C. Belta, and A. A. Julius, "Temporal logic inference with prior information: An application to robot arm movements," *IFAC-PapersOnLine*, vol. 48, no. 27, pp. 141 – 146, 2015.

[9] Z. Xu, M. Birtwistle, C. Belta, and A. Julius, "A temporal logic inference approach for model discrimination," *IEEE Life Sciences Letters*, vol. 2, no. 3, pp. 19–22, Sept 2016.

[10] X. Xu and P. J. Antsaklis, "Results and perspectives on computational methods for optimal control of switched systems," in *Proc. of the 6th Int. Conf. on Hybrid Systems: Computation and Control*, ser. HSCC'03.  Berlin, Heidelberg: Springer-Verlag, 2003, pp. 540–555.

[11] G. E. Fainekos and G. J. Pappas, "Robustness of temporal logic specifications for continuous-time signals," *Theoretical Computer Science*, vol. 410, no. 42, pp. 4262 – 4291, 2009.

[12] Z. Xu and A. A. Julius, "Census signal temporal logic inference for multiagent group behavior analysis," *IEEE Trans. Autom. Science and Eng.*, Early access on IEEE Xplore.

[13] Y. Deng, A. DInnocenzo, M. D. Di Benedetto, S. Di Gennaro, and A. A. Julius, "Verification of hybrid automata diagnosability with measurement uncertainty," *IEEE Trans. Autom. Control*, vol. 61, no. 4, pp. 982–993, April 2016.

[14] A. Girard and G. J. Pappas, "Approximate bisimulation relations for constrained linear systems," *Automatica*, vol. 43, no. 8, pp. 1307 – 1317, 2007.

[15] A. A. Julius, G. E. Fainekos, M. Anand, I. Lee, and G. J. Pappas, "Robust test generation and coverage for hybrid systems," in *Proc. Hybrid Syst.: Computat. Control*.  Springer, 2007, pp. 329–342.

[16] A. A. Julius, "Trajectory-based controller design for hybrid systems with affine continuous dynamics," in *2010 IEEE Int. Conf. on Autom. Science and Eng.*, Aug 2010, pp. 1007–1012.

[17] F. Zhang, J. Cao, and Z. Xu, "An improved particle swarm optimization particle filtering algorithm," in *Int. Conf. on Communications, Circuits and Systems (ICCCAS)*, vol. 2, Nov 2013, pp. 173–177.

[18] C. C. Okaeme, S. Mishra, and J. T. Wen, "A comfort zone set-based approach for coupled temperature and humidity control in buildings," in *IEEE Int. Conf. on Autom. Science and Eng. (CASE)*, Aug 2016, pp. 456–461.

[19] Z. Xu, A. Julius, and J. H. Chow, "Optimal energy storage control for frequency regulation under temporal logic specifications," in *Proc. American Control Conference (ACC)*, May 2017, pp. 1874–1879.