

Quickest Anomaly Detection in Sensor Networks With Unlabeled Samples

Zhongchang Sun  and Shaofeng Zou 

Abstract—The problem of quickest anomaly detection in networks with unlabeled samples is studied. At some unknown time, an anomaly emerges in the network and changes the data-generating distribution of some unknown sensor. The data vector received by the fusion center at each time step undergoes some unknown and arbitrary permutation of its entries (unlabeled samples). The goal of the fusion center is to detect the anomaly with minimal detection delay subject to false alarm constraints. With unlabeled samples, existing approaches that combines local cumulative sum (CuSum) statistics cannot be used anymore. Several major questions include whether detection is still possible without the label information, if so, what is the fundamental limit and how to achieve that. Two cases with static and dynamic anomaly are investigated, where the sensor affected by the anomaly may or may not change with time. For the two cases, practical algorithms based on the ideas of mixture likelihood ratio and/or maximum likelihood estimate are constructed. Their average detection delays and false alarm rates are theoretically characterized. Universal lower bounds on the average detection delay for a given false alarm rate are also derived, which further demonstrate the asymptotic optimality of the two algorithms.

Index Terms—Quickest change detection, unlabeled samples, permuted samples, asymptotically optimal, fundamental limits.

I. INTRODUCTION

IN SENSOR networks, samples may lack label information such as identity due to, e.g., malicious attacks and limited communication resources. For example, wireless ad-hoc sensor networks are usually vulnerable to spoofing attacks [2], and samples received by the fusion center may then lose their label information. Furthermore, in Internet-of-things (IoT) networks, where devices are commonly small and low-cost sensing devices powered by battery with limited communication bandwidth, and are usually deployed in a massive scale, the communication overhead of identifying individual sensors increases drastically with the number of sensors [3]. However, these battery-powered IoT devices are usually expected to survive for years without

battery change. In this case, message that is delivered to the fusion center may be constrained not to contain the identity information. Furthermore, in social sensing applications, participants may choose to be anonymous in order to protect privacy, i.e., sharing the data without including identity information. Motivated by these applications, there is a recent surge of interest in the problem of signal processing with unlabeled data (see e.g., [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17]), which refers to various signal processing problems where the data vector undergoes an unknown permutation of its entries, and the original position of each datum in the vector is unknown.

In this paper, we investigate the problem of quickest anomaly detection in sensor networks with unlabeled samples. Specifically, at some unknown time, an anomaly emerges in the network and leads to a change in the data-generating distribution of some unknown sensor. The fusion center sequentially receives unlabeled (arbitrarily permuted) samples from all the sensors at each time step. The goal of the fusion center is to detect the anomaly as quickly as possible, subject to false alarm constraints. This problem is of particular relevance to applications where an anomaly affects some sensor in the network, and the affected sensor may change over time [18], e.g., surveillance system, intrusion detection, environmental change detection, rumor detection, and seismic wave detection.

A. Contributions and Major Challenges

We first focus on the static anomaly, where the sensor affected by the anomaly does not change with time, but which sensor is affected is still unknown. We consider the detection delay under the worst-case affected sensor. The goal here is to minimize the detection delay subject to false alarm constraints. The major challenges here are two-fold. First of all, the labels of the samples are unknown and time-varying. Second, even if the labels are known, i.e., each sample is associated with its sensor, the sensor the anomaly affects is still unknown. For this problem, we construct a generalized mixture CuSum (GM-CuSum) algorithm. The basic idea is to estimate the unknown identity of the affected sensor using the maximum likelihood estimate (MLE), and further employ a mixture likelihood w.r.t. all possible labels. We prove that the GM-CuSum is second-order asymptotically optimal.¹

¹An algorithm is second-order asymptotically optimal if as the false alarm rate goes to zero, its detection delay is within an $O(1)$ term of the best possible detection delay.

Manuscript received 3 September 2022; revised 4 January 2023; accepted 6 March 2023. Date of publication 17 March 2023; date of current version 31 March 2023. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Ali Tajer. The work of Zhongchang Sun and Shaofeng Zou was supported by the National Science Foundation under Grants 1948165 and 2112693. This paper was presented in part at the 2021 IEEE International Symposium on Information Theory [DOI: 10.1109/ISIT45174.2021.9517771]. (Corresponding author: Shaofeng Zou.)

The authors are with the Department of Electrical Engineering, University at Buffalo, Buffalo, NY 14228 USA (e-mail: zhongcha@buffalo.edu; szou3@buffalo.edu).

Digital Object Identifier 10.1109/TSP.2023.3256275

We then focus on a general and more challenging setting with dynamic anomaly, where the sensor affected by the anomaly changes with time. Here, we refer to the sequence of sensors affected by the anomaly over time as the trajectory of the anomaly. We consider the detection delay under the worst-case trajectory. Compared to the static setting, the additional challenge is that the affected sensor changes with time, and thus the change is not persistent at any particular sensor. Therefore, estimating the identity of the affected sensor over time is not applicable. We then propose a weighted approach to address this challenge, and find the optimal weight to construct a weighted mixture CuSum algorithm. We prove that the weighted mixture CuSum algorithm is first-order asymptotically optimal.² We also discuss two computationally efficient approximations for large-scale networks.

We also conduct extensive numerical experiments to demonstrate the performance of our proposed algorithms. The numerical results show that for the static setting, our GM-CuSum algorithm outperforms a heuristic uniformly weighted mixture CuSum algorithm; the optimal weighted mixture CuSum algorithm also performs well for the static setting; and for the dynamic setting, our optimal weighted mixture CuSum algorithm outperforms an uniformly weighted one and the GM-CuSum algorithm. These numerical results validate our theoretical optimality results.

B. Related Work

The quickest change detection (QCD) problem in sensor networks with labeled samples was extensively studied in the literature, e.g., [19], [20], [21], [22], [23], [24], [25], [26], [27], [28], [29], [30], [31], where the fusion center knows the identity of each sample, i.e., knows which sensor that each sample is from. Therefore, one CuSum algorithm can be implemented at each sensor and then be combined to make the decision. This type of algorithms were shown to be asymptotically optimal for various settings. In this paper, we investigate the setting with unlabeled samples, where at each time step samples are arbitrarily permuted, and the permutation is time-varying. The fusion center does not know which sensor each sample comes from, and then cannot implement a CuSum algorithm for each sensor.

Various learning and inference problems with unlabeled data have been studied in the literature [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], which mainly focus on the offline setting with non-sequential data. Here we only review several closely related ones on detection problems. In [6], hypothesis testing with unlabeled samples was studied, where two practical algorithms, the unlabeled log-likelihood ratio test and the generalized likelihood ratio test were proposed. A more specific problem was studied in [7] where samples follow the Bernoulli distribution and an approximated log-likelihood test based on the central limit theorem was proposed. In [4], the binary hypothesis testing problem with unlabeled samples was

studied, and an optimal mixture likelihood ratio test (MLRT) was developed. In [5], the bandwidth-constrained QCD problem with unlabeled samples was investigated, where each sensor sends 1-bit quantized feedback to the fusion center. In [17], the QCD problem with unlabeled samples was studied where the change affects all the sensors simultaneously. In this paper, we investigate a practical scenario where an anomaly may not affect all the sensors, which is of particular interest in the distributed setting, and the anomaly may also be dynamic and affect different sensors at different times, e.g., a moving target in surveillance systems.

Existing studies of quickly detecting a dynamic change mostly focus on the labeled setting, e.g., [31], [32], [33]. Our problem is similar to the one in [33] but we focus on the unlabeled setting. Our major challenge is due to the additional ambiguity of unknown labels. The QCD problem with a slowly changing post-change distribution was studied in [34], [35], whereas in this paper, the anomaly can move arbitrarily fast.

With unlabeled samples, our problem is also related to the composite QCD problem with unknown pre- and post-change parameters e.g., [22], [36], [37], [38]. Our work is different from the existing literature since the unknown parameters, i.e., the identity and the label of the affected sensor, are time-varying. Therefore, the generalized likelihood approach which estimates the unknown parameters using their MLEs may not perform well. Moreover, unlike studies in [36], [37], [38] where the distributions are assumed to belong to the exponential family, we do not have any assumptions on the distributions.

II. PROBLEM FORMULATION

Consider a network monitored in real time by a set of n heterogeneous sensors. These sensors can be clustered into K types, and each type k has n_k sensors, $1 \leq k \leq K$. The data-generating distributions of samples from type k sensors are denoted by $p_{\theta,k}$, $\theta \in \{0, 1\}$, which are known to the fusion center. At some unknown time ν , an anomaly emerges in the network, and changes the data-generating distribution of one sensor. The fusion center does not know which type of sensor is affected. If a sensor of type k is affected by the anomaly, then its samples are generated by $p_{1,k}$, otherwise, by $p_{0,k}$. The goal is to detect the anomaly as quickly as possible subject to false alarm constraints. We focus on the case with unlabeled samples, where the data vector at each time step undergoes an unknown permutation of its entries, and the original position of each datum in the vector is unknown to the fusion center. In other words, the fusion center does not know which type of sensors that each sample comes from, and therefore, does not know the sample's exact data-generating distribution.

Based on whether the sensor is affected by the anomaly and the type of the sensor, we rearrange the sensors into $2K$ groups. The first K groups consists of sensors that are not affected by the anomaly; and the second K groups consists of sensors that are affected by the anomaly. Specifically, for sensors in group $1 \leq k \leq K$, their samples are generated by $p_{0,k}$, and for sensors in group $K < k \leq 2K$, their samples are generated by $p_{1,k-K}$.

²An algorithm is first-order asymptotically optimal if as the false alarm rate goes to zero, the ratio between its detection delay and the best possible detection delay goes to 1.

In this paper, we use capital letters to denote random variables and lower case letters to denote their realizations. Denote by $X^n[t] = \{X_1[t], \dots, X_n[t]\}$ the n arbitrarily permuted samples at time t received by the fusion center. We assume that $X_1[t], \dots, X_n[t]$ are independent, and $X^n[t_1]$ is independent from $X^n[t_2]$ for any $t_1 \neq t_2$. Note that $X_i[t]$ is not necessarily the sample from sensor i since samples are permuted/unlabeled. The results in our paper hold for both continuous and discrete random variables $X_i[t]$.

Let $\mathcal{K} = \{1, 2, \dots, K\}$. Denote by $s[t] \in \mathcal{K} \cup \{0\}$ the type of the affected sensor at t . For notational convenience, we use $s[t] = 0$ to denote the case when there is no anomaly, i.e., $t < \nu$. Let $\mathbf{s} \triangleq \{s[t]\}_{t=1}^\infty$ denote the trajectory of the anomaly. Here \mathbf{s} is *unknown* to the decision maker. Even if the trajectory of the anomaly \mathbf{s} is given, the distribution of $X^n[t]$ still cannot be fully specified due to lack of label information. To characterize the distribution of $X^n[t]$, we define a label function $\sigma_t^{s[t]} : \{1, \dots, n\} \rightarrow \{1, \dots, K, s[t] + K\}$. This function associates sample $X_i[t]$, $1 \leq i \leq n$, to group j for some $j \in \{1, 2, \dots, K, K + s[t]\}$, i.e., specifies the probability distribution of $X_i[t]$. Specifically, if $\sigma_t^{s[t]}(i) = j$, then

$$X_i[t] \sim \begin{cases} p_{0,j}, & \text{if } 1 \leq j \leq K, \\ p_{1,j-K}, & \text{if } K < j \leq 2K. \end{cases} \quad (1)$$

Here $\sigma_t^{s[t]}$ can be interpreted as the inverse of the permutation applied to the data vector. We further note that $\sigma_t^{s[t]}$ is *unknown* to the decision maker.

Let $\Omega_s = \{\sigma_1^{s[1]}, \dots, \sigma_\infty^{s[\infty]}\}$ be the labels when the trajectory of the anomaly is \mathbf{s} , which is unknown. Let $\mathbb{P}_{\Omega_s}^{s,\nu}$ and $\mathbb{E}_{\Omega_s}^{s,\nu}$ denote the probability measure and the corresponding expectation when the change point is at ν and the samples received by the fusion center are permuted according to the label Ω_s (see Appendix A for more details). We further let \mathbb{P}_Ω^∞ and \mathbb{E}_Ω^∞ denote the probability measure and the corresponding expectation when there is no change, i.e., $\nu = \infty$, where $\Omega = \Omega_s$ with $s[t] = 0, \forall t \geq 1$.

We extend Lorden's criterion [39], and define the worst-case average detection delay (WADD) and the worst-case average running length (WARL) for any stopping time τ :

$$\text{WADD}(\tau) = \sup_{\nu \geq 1} \sup_s \sup_{\Omega_s} \text{esssup} \mathbb{E}_{\Omega_s}^{s,\nu}[(\tau - \nu)^+ | \mathbf{X}^n[1, \nu - 1]],$$

$$\text{WARL}(\tau) = \inf_{\Omega} \mathbb{E}_{\Omega}^\infty[\tau], \quad (2)$$

where $\mathbf{X}^n[t_1, t_2] = \{X^n[t_1], \dots, X^n[t_2]\}$, for any $t_1 \leq t_2$ and $(\tau - \nu)^+ = \max\{\tau - \nu, 0\}$. Let $f : X \rightarrow \mathbb{R}$ be a real-valued function and (X, \mathcal{X}, μ) be a probability space. The essential supremum is then defined as

$$\text{esssup} f = \inf \{a \in \mathbb{R} : \mu(\{x : f(x) > a\}) = 0\}. \quad (3)$$

The goal is to design a stopping rule that minimizes the WADD subject to a constraint on the WARL:

$$\inf_{\tau: \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau), \quad (4)$$

TABLE I
SUMMARY OF NOTATIONS

Symbol	Definition
n	Number of sensors
K	Number of different types
n_k	Number of sensors of type k
$s[t]$	Type of the affected sensor at time t
$\mathbf{s} \triangleq \{s[t]\}_{t=1}^\infty$	Trajectory of the anomaly
$\sigma_t^{s[t]}$	Label function that associates sample to group
Ω_s	Labels when the trajectory of the anomaly is \mathbf{s}
$\mathbb{P}_{\Omega_s}^{s,\nu}(\mathbb{E}_{\Omega_s}^{s,\nu})$	Probability measure (expectation) when the change point is at ν and the samples are permuted according to the label Ω_s
$\mathbb{P}_\Omega^\infty(\mathbb{E}_\Omega^\infty)$	Probability measure (expectation) when there is no change
$D(P Q)$	Kullback-Leibler (KL) divergence between two distributions P and Q .

where $\gamma > 0$ is a pre-specified threshold. Here the false alarm constraint is to guarantee that under all possible sample permutations, the average running length to a false alarm is always lower bounded by γ .

A stopping time T is second-order asymptotically optimal if $\text{WARL}(T) \geq \gamma$ and for large γ

$$\text{WADD}(T) = \inf_{\tau: \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau) + O(1). \quad (5)$$

A stopping time T is first-order asymptotically optimal if $\text{WARL}(T) \geq \gamma$ and for large γ

$$\text{WADD}(T) = \inf_{\tau: \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau) (1 + o(1)). \quad (6)$$

In Table I, we summarize important notations in this paper.

III. STATIC ANOMALY

We first investigate the case with static anomaly, i.e., the sensor affected by the anomaly does not change with time. In this case, for any $t \geq \nu$, $s[t] = k$ for some unknown type k . Then, for all $j \in \{1, 2, \dots, K, k + K\}$, there are $(n_1, \dots, n_{k-1}, \dots, n_K, 1)$ possible σ_t^k 's to associate each sample with a data-generating distribution, and we denote the collection of all possible labels by $\mathcal{S}_{n,k}$ (see Appendix A for more details). Before the anomaly emerges, i.e., $t < \nu$, $X^n[t]$ follows the distribution

$$\mathbb{P}_{0,\sigma_t^0}(X^n[t]) \triangleq \prod_{i=1}^n p_{0,\sigma_t^0(i)}(X_i[t]), \quad (7)$$

for some unknown $\sigma_t^0 \in \mathcal{S}_{n,0}$. At time $t \geq \nu$, $s[t] = k$, $X^n[t]$ follows the distribution

$$\begin{aligned} \mathbb{P}_{\sigma_t^k}^k(X^n[t]) &\triangleq \prod_{i: \sigma_t^k(i) \leq K} p_{0,\sigma_t^k(i)}(X_i[t]) \\ &\times \prod_{i: \sigma_t^k(i) > K} p_{1,\sigma_t^k(i)-K}(X_i[t]), \end{aligned} \quad (8)$$

for some unknown $\sigma_t^k \in \mathcal{S}_{n,k}$. Let $\Omega_k = \{\sigma_1^k, \dots, \sigma_{\nu-1}^k, \sigma_\nu^k, \dots, \sigma_\infty^k\}$ be the labels over time, when the anomaly emerges at ν (similarly defined as Ω_s). Let $\mathbb{P}_{\Omega_k}^{k,\nu}$ denote the probability measure when the change point is at ν and the samples are generated according to (7), (8) and Ω_k . Further let $\mathbb{E}_{\Omega_k}^{k,\nu}$ denote the corresponding expectation.

The WADD for a stopping time τ can be written as

$$\text{WADD}(\tau) = \sup_{\nu \geq 1} \sup_k \sup_{\Omega_k} \text{esssup} \mathbb{E}_{\Omega_k}^{k, \nu} [(\tau - \nu)^+ | \mathbf{X}^n[1, \nu - 1]].$$

The WARL is defined in the same way as in (2).

The goal is to design a stopping rule that minimizes the WADD subject to a constraint on the WARL:

$$\inf_{\tau: \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau). \quad (9)$$

A. Universal Lower Bound on WADD

We first derive a universal lower bound on WADD for any τ satisfying the false alarm constraint: $\inf_{\Omega} \mathbb{E}_{\Omega}^{\infty}[\tau] \geq \gamma$.

Let $I_k = D(\tilde{\mathbb{P}}_k || \tilde{\mathbb{P}}_0)$ denote the Kullback-Leibler (KL) divergence between two mixture distributions $\tilde{\mathbb{P}}^k = \frac{1}{|\mathcal{S}_{n,k}|} \sum_{\sigma \in \mathcal{S}_{n,k}} \mathbb{P}_{\sigma}^k$ and $\tilde{\mathbb{P}}_0 = \frac{1}{|\mathcal{S}_{n,0}|} \sum_{\sigma \in \mathcal{S}_{n,0}} \mathbb{P}_{0,\sigma}$. Here, $\tilde{\mathbb{P}}^k$ is the uniform mixture of all possible labels when the affected sensor is type k . Let $I^* = \min_{1 \leq k \leq K} I_k$. We then have the following theorem.

Theorem 1: As $\gamma \rightarrow \infty$,

$$\inf_{\tau: \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau) \geq \frac{\log \gamma}{I^*} + O(1). \quad (10)$$

The proof of Theorem 1 can be found in Appendix B. The main challenge in the proof of Theorem 1 is due to the worst-case over all labels and affected sensors in WADD and WARL. From Theorem 1, it can be seen that the WADD for problem (9) is lower bounded by $\frac{\log \gamma}{I^*} + O(1)$ for any stopping rule that satisfies the constraint on WARL. Theorem 1 motivates us to find the k that minimizes I_k , i.e., achieves I^* , and design an algorithm to achieve this universal lower bound.

B. Generalized Mixture CuSum Algorithm

In this section, we construct an algorithm that achieves the universal lower bound asymptotically.

When there are unknown parameters, MLE is commonly used to estimate the unknown parameters. In the static setting, k does not change with time, however, σ_t^k changes with time, thus a direct MLE for σ_t^k at each time t may not work well.

If k is known, then our problem is invariant under the group of transformations of all possible labels (permutations), that is, our problem is independent of the order of collected samples at each time. Therefore, our problem is related to the invariant theory in [40, Section 6]. This motivates us to take a mixture approach w.r.t. the unknown labels, and then take a MLE approach w.r.t. the unknown affected sensor.

Let $W[t] = \max_{k \in \mathcal{K}} \max_{1 \leq j \leq t} \sum_{i=j}^t \log \frac{\tilde{\mathbb{P}}^k(X^n[i])}{\tilde{\mathbb{P}}_0(X^n[i])}$. We then define the GM-CuSum stopping time as follows:

$$T_G = \inf\{t : W[t] \geq b\}, \quad (11)$$

where $b > 0$ is the threshold. Here $W[t]$ can be updated efficiently. We keep K CuSums in parallel. Note that this can be done recursively. Let $W_k[t] = \max_{1 \leq j \leq t} \sum_{i=j}^t \log \frac{\tilde{\mathbb{P}}^k(X^n[i])}{\tilde{\mathbb{P}}_0(X^n[i])}$.

The test statistic $W[t]$ has the following recursion:

$$W[t+1] = \max_{k \in \mathcal{K}} \left\{ (W_k[t])^+ + \log \frac{\tilde{\mathbb{P}}^k(X^n[t+1])}{\tilde{\mathbb{P}}_0(X^n[t+1])} \right\}, \quad (12)$$

where $W_k[0] = 0, \forall k$. We then take their maximum as $W[t]$.

In the following theorem, we show 1) the WARL lower bound of T_G and 2) the WADD upper bound of T_G .

Theorem 2: 1) Let $b = \log(K\gamma)$ in (11). Then $\text{WARL}(T_G) \geq \gamma$; and 2) As $\gamma \rightarrow \infty$, $\text{WADD}(T_G) \leq \frac{\log \gamma}{I^*} + O(1)$.

The proof of Theorem 2 can be found in Appendix C. The proof of the lower bound on WARL is based on Doob's submartingale inequality [41] and the optional sampling theorem [41]. The major challenge lies in that we consider the worst-case label. A key property we develop and use in the proof is that under the pre-change distribution $\mathbb{P}_{0,\sigma_t^0}$, for any $k \in \mathcal{K}$, the expectation of the mixture likelihood ratio $\mathbb{E}_{0,\sigma^0} \left[\log \frac{\tilde{\mathbb{P}}^k(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right]$ is invariant for different σ^0 's.

Theorem 2 suggests that to meet the WARL constraint, b should be chosen such that $b = \log K\gamma$.

Based on Theorems 1 and 2, we then establish the second-order asymptotic optimality of T_G .

Theorem 3: T_G is second-order asymptotically optimal for the problem in (9).

The asymptotic optimality of T_G can be derived similarly under the Pollak's criterion [42]. We omit the details here.

IV. QUICKEST DYNAMIC ANOMALY DETECTION

In this section, we consider the general problem with a dynamic anomaly, where the sensor affected by the anomaly changes with time. The GM-CuSum algorithm designed for static anomaly may not work well anymore since the sensor affected by the anomaly changes with time.

A. Universal Lower Bound on WADD

Define the following weighted mixture distribution: $\tilde{\mathbb{P}}^{\beta}(X^n) = \sum_{k=1}^K \beta_k \tilde{\mathbb{P}}^k(X^n)$, where $\beta = \{\beta_k\}_{k=1}^K$, $0 \leq \beta_k \leq 1$ and $\sum_{k=1}^K \beta_k = 1$. Denote by I_{β} the KL divergence between $\tilde{\mathbb{P}}^{\beta}$ and $\tilde{\mathbb{P}}_0$. Let $\beta^* = \arg \min_{\beta} I_{\beta}$.

For the universal lower bound on WADD, we have the following theorem.

Theorem 4: As $\gamma \rightarrow \infty$, we have that

$$\inf_{\tau: \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau) \geq \frac{\log \gamma}{I_{\beta^*}} (1 + o(1)). \quad (13)$$

The proof of Theorem 4 can be found in Appendix D. From Theorem 4, the WADD for the problem in (4) is lower bounded by $\frac{\log \gamma}{I_{\beta^*}} (1 + o(1))$ for large γ . This motivates us to apply the optimal weight β^* to design an algorithm that can achieve the WADD lower bound asymptotically. Moreover, we have that $I^* \geq I_{\beta^*}$ which implies that a dynamic anomaly is more difficult to detect than a static anomaly.

B. Weighted Mixture CuSum

In the static setting, the unknown affected sensor can be estimated by its MLE. However, in the dynamic setting, the affected sensor changes with time, and the MLE approach may not work well. Theorem 4 motivates us to tackle the unknown anomaly trajectory using a weighted approach where the probability that the k -th group is affected by the anomaly is β_k^* . We then construct our optimal weighted mixture CuSum algorithm as follows. Define the log of weighted mixture likelihood ratio using β^* :

$$\ell_{\beta^*}(X^n) = \log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)}. \quad (14)$$

It can be easily shown that $\ell_{\beta^*}(X^n)$ is invariant to any permutations on X^n , i.e., for any permutation $\pi(X^n) = (X_{\pi(1)}, X_{\pi(2)}, \dots, X_{\pi(n)})$, $\ell_{\beta^*}(X^n) = \ell_{\beta^*}(\pi(X^n))$. This is due to the fact that $\ell_{\beta^*}(X^n)$ takes the sum over all possible labels thus is invariant to the actual permutation.

We then construct the following optimal weighted mixture CuSum algorithm:

$$T_{\beta^*}(b) = \inf \left\{ t : \max_{1 \leq j \leq t+1} \sum_{i=j}^t \ell_{\beta^*}(X^n[i]) \geq b \right\}. \quad (15)$$

Let $\widehat{W}[t] = \max_{1 \leq j \leq t+1} \sum_{i=j}^t \ell_{\beta^*}(X^n[i])$, then $\widehat{W}[t+1] = (\widehat{W}[t])^+ + \ell_{\beta^*}(X^n[t+1])$, and $\widehat{W}[0] = 0$.

Different from the way that we handle the unknown label σ , here, for the unknown type of the affected sensor, we take the mixture according to β^* instead of a uniform distribution over \mathcal{K} . As will be shown later both theoretically in Theorem 6 and numerically in Section VI, taking a uniform mixture over \mathcal{K} may not lead to the optimal performance.

Let $\tilde{\mathbb{E}}^k$ and $\tilde{\mathbb{E}}_0$ denote the expectation under the probability $\tilde{\mathbb{P}}^k$ and $\tilde{\mathbb{P}}_0$ respectively. The following property of β^* plays an important role in developing the asymptotic optimality of the weighted mixture CuSum algorithm.

Lemma 1: For any $k \in \mathcal{K}$, $\tilde{\mathbb{E}}^k \left[\log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] \geq I_{\beta^*}$.

The proof of Lemma 1 can be found in Appendix E.

In the following, we provide a heuristic explanation of how $\widehat{W}[t]$ evolves in the pre- and post-change regimes. We first argue that $\mathbb{E}_{\sigma^k}^k[\ell_{\beta^*}(X^n)]$ is invariant for different σ^k 's. Specifically, let $\mathbb{E}_{\sigma^k}^k$ denote the expectation under $\mathbb{P}_{\sigma^k}^k$, where a sensor of type k is affected, and the data received is labeled according to σ^k . For any π , let $\hat{\sigma}^k = \sigma^k \circ \pi$. Then $\mathbb{E}_{\sigma^k}^k[\ell_{\beta^*}(\pi(X^n))] = \mathbb{E}_{\sigma^k \circ \pi}^k[\ell_{\beta^*}(X^n)] = \mathbb{E}_{\hat{\sigma}^k}^k[\ell_{\beta^*}(X^n)]$. For any $\hat{\sigma}^k \in \mathcal{S}_{n,k}$, a π can always be found so that $\sigma^k \circ \pi = \hat{\sigma}^k$. Thus, for any $\sigma^k, \hat{\sigma}^k \in \mathcal{S}_{n,k}$, $\mathbb{E}_{\sigma^k}^k[\ell_{\beta^*}(X^n)] = \mathbb{E}_{\hat{\sigma}^k}^k[\ell_{\beta^*}(X^n)]$. Therefore, $\mathbb{E}_{\sigma^k}^k[\ell_{\beta^*}(X^n)]$ is invariant for different σ^k 's. Then, under the pre-change distribution $\mathbb{P}_{0,\sigma_t^0}$, the expectation of the weighted mixture likelihood ratio $\mathbb{E}_{0,\sigma_t^0}[\ell_{\beta^*}(X^n)]$ is invariant for different σ_t^0 's. This implies that

$$\mathbb{E}_{0,\sigma_t^0} \left[\log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] = -D(\tilde{\mathbb{P}}_0 || \tilde{\mathbb{P}}^{\beta^*}) \leq 0. \quad (16)$$

Therefore, before the change time ν , $\widehat{W}[t]$ has a negative drift. Similarly, from Lemma 1, after the change time ν , under any group assignment Ω_s and trajectory s , $\widehat{W}[t]$ has a positive drift whose expectation is no less than I_{β^*} , and evolves towards ∞ .

The following theorem establishes 1) the WARL lower bound of T_{β^*} , and 2) the WADD upper bound of T_{β^*} .

Theorem 5: 1) For T_{β^*} defined in (15), let $b = \log \gamma$, then $\text{WARL}(T_{\beta^*}) \geq \gamma$. 2) As $\gamma \rightarrow \infty$, we have that $\text{WADD}(T_{\beta^*}) \leq \frac{\log \gamma}{I_{\beta^*}}(1 + o(1))$.

The proof of Theorem 5 can be found in Appendix F. The proof is based on the Weak Law of Large Numbers for the weighted mixture likelihood ratio (similar to [37]). The major challenge lies in that here we are interested in the worst-case label and the worst-case anomaly trajectory. In our problem, the label and the affected sensor change with time. Therefore, it's challenging to explicitly characterize the worst-case label and anomaly trajectory for T_{β^*} . To show the asymptotically optimal performance of T_{β^*} , instead of finding the worst-case label and anomaly trajectory, we apply the symmetric property of T_{β^*} and Lemma 1 to show that the WADD and WARL of T_{β^*} are bounded under all possible labels and trajectories.

We then establish the first-order asymptotic optimality of T_{β^*} in the following theorem.

Theorem 6: T_{β^*} is first-order asymptotically optimal for problem (4).

Proof: Combining Theorems 4 and 5, we establish the first-order asymptotic optimality of T_{β^*} . ■

The asymptotic optimality of T_{β^*} can be derived similarly under the Pollak's criterion [42]. We omit the details here.

If we apply T_{β^*} (designed for the dynamic setting) to the static setting, the WADD of T_{β^*} can also be upper bounded by $\frac{\log \gamma}{I_{\beta^*}}(1 + o(1))$. However, T_{β^*} may not be asymptotically optimal. On the other hand, in the dynamic setting, the sensor affected by the anomaly changes with time, and thus the MLE may not work well. Therefore, the optimal weighted mixture CuSum algorithm works better than the GM-CuSum.

V. COMPUTATIONAL COMPLEXITY AND EFFICIENT APPROXIMATION

In the previous sections, we proved that the GM-CuSum algorithm and the optimal weighted mixture CuSum algorithm are asymptotically optimal. However, the test statistic involves computing the mixture likelihood over all possible $\sigma_t^{s[t]}$, which is expensive when n is large.

At each time t , we have $\binom{n_1, \dots, n_{s[t]}-1, \dots, n_K, 1}{n}$ possible $\sigma_t^{s[t]}$'s. Consider the case with large n , and assume that $\lim_{n \rightarrow \infty} \frac{n_k}{n} = \alpha_k$, which is a constant. Let $\alpha = [\alpha_1, \dots, \alpha_K]^T$. Consider discrete distributions,³ and denote by \mathcal{X} the support set of the samples. From the exponential bound on the size of the type

³Samples in sensor networks are usually quantized before transmitting to fusion center to reduce the power consumption.

class [43], we have that

$$\frac{2^{nH\left(\left[\frac{n_1}{n}, \dots, \frac{n_{s[t]-1}}{n}, \dots, \frac{n_K}{n}, \frac{1}{n}\right]\right)}}{(n+1)^{|\mathcal{X}|}} \leq (n_1, \dots, n_{s[t]-1}, \dots, n_K, 1) \\ \leq 2^{nH\left(\left[\frac{n_1}{n}, \dots, \frac{n_{s[t]-1}}{n}, \dots, \frac{n_K}{n}, \frac{1}{n}\right]\right)},$$

where H denotes the Shannon entropy. We then have that

$$\lim_{n \rightarrow \infty} H\left(\left[\frac{n_1}{n}, \dots, \frac{n_{s[t]-1}}{n}, \dots, \frac{n_K}{n}, \frac{1}{n}\right]\right) = H(\alpha).$$

Therefore, the computational complexity of the GM-CuSum and the optimal weighted mixture CuSum increase exponentially with n , which is expensive for large n .

In this following, we discuss two computationally efficient methods to approximate the test statistics of the GM-CuSum and the optimal weighted mixture CuSum when n is large, and then evaluate their performance in Section VI-D.

The first method is based on the method of types [43]. Let Π_{X^n} denote the empirical distribution of X^n and let $\mathcal{T}(\Pi_{X^n})$ denote the type class of Π_{X^n} . For any k , we have that [4]

$$\log \frac{\tilde{\mathbb{P}}^k(X^n)}{\tilde{\mathbb{P}}_0(X^n)} = \log \frac{\tilde{\mathbb{P}}^k(\mathcal{T}(\Pi_{X^n}))}{\tilde{\mathbb{P}}_0(\mathcal{T}(\Pi_{X^n}))}. \quad (17)$$

From the generalized Sanov's theorem [4], [17], the probability of types $\log \tilde{\mathbb{P}}^k(\mathcal{T}(\Pi_{X^n}))$ can further be approximated by the following optimization problem

$$- \inf_{\substack{U=[U_1, \dots, U_{K+1}]^T \in (\mathcal{P}_{\mathcal{X}})^{K+1} \\ \phi^T U = \Pi_{X^n}}} \sum_{j=1, j \neq k}^K n_j D(U_j \| p_{0,j}) \\ + (n_k - 1) D(U_k \| p_{0,k}) + D(U_{K+1} \| p_{1,k}), \quad (18)$$

where $\phi = \left[\frac{n_1}{n}, \dots, \frac{n_{k-1}}{n}, \dots, \frac{n_K}{n}, \frac{1}{n}\right]^T$ and $\mathcal{P}_{\mathcal{X}}$ denotes the probability simplex on \mathcal{X} . Problem (18) is a convex optimization problem whose computational complexity is independent of n . Therefore, the computation of test statistics of the GM-CuSum algorithm and the optimal weighted mixture CuSum algorithm can be converted to solving convex optimization problems and the overall complexity at each time step is only linear in the number of sensors.

The second method is to estimate the unknown labels using the MLE and to use the generalized likelihood ratio test (GLRT) to approximate the mixture likelihood ratio [6]. Computing the GLRT is a special case of the assignment problem and efficient algorithms have been developed. In [6], two efficient greedy algorithms were proposed to solve the assignment problem approximately with complexity $O(n^2)$. Therefore, at each time t , the test statistics of the GM-CuSum and the optimal weighted mixture CuSum algorithm can be approximated with computational complexity $O(Kn^2)$.

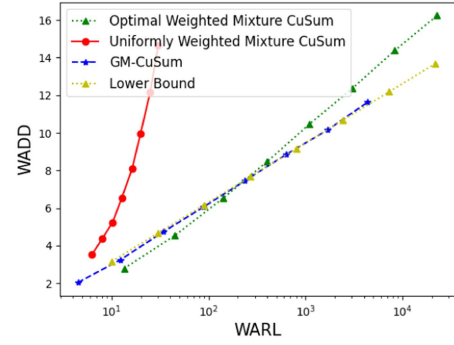


Fig. 1. Static setting: $n = 4, K = 4$.

VI. NUMERICAL RESULTS

A. Static Anomaly Detection

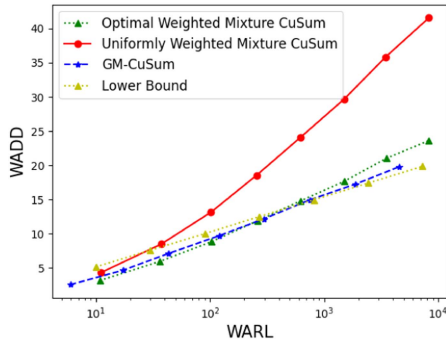
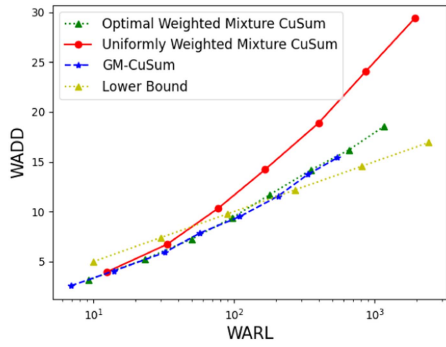
We first consider the static setting, and compare our GM-CuSum with a uniformly weighted mixture CuSum

$$T_B = \inf \left\{ t : \max_{1 \leq j \leq t} \sum_{i=j}^t \log \frac{\frac{1}{|\mathcal{K}|} \sum_{k \in \mathcal{K}} \tilde{\mathbb{P}}^k(X^n[i])}{\tilde{\mathbb{P}}_0(X^n[i])} \geq b \right\}$$

and the optimal weighted mixture CuSum. We plot the WADD as a function of the WARL under the worst-case static trajectory. We also plot the asymptotic lower bound $\frac{\log \gamma}{I^*}$.

We first compare the three algorithms under the Gaussian distributions. There are four types of sensors and each type contains one sensor. For the type I sensor, the pre- and post-change distributions are $\mathcal{N}(-1, 1)$ and $\mathcal{N}(2, 1)$, for the type II sensor, the pre- and post-change distributions are $\mathcal{N}(1, 1)$ and $\mathcal{N}(3, 1)$, for the type III sensor, the pre- and post-change distributions are $\mathcal{N}(-1, 1)$ and $\mathcal{N}(3, 1)$, for the type IV sensor, the pre- and post-change distributions are $\mathcal{N}(1, 1)$ and $\mathcal{N}(-1, 1)$ respectively. The optimal weight for our weighted mixture CuSum algorithm is solved by Monte-Carlo. It can be seen from Fig. 1 that with the same false alarm rate, our GM-CuSum has the lowest WADD, which implies that it detects the anomaly with the smallest detection delay. Moreover, the slope of the GM-CuSum matches the lower bound, which validates that the GM-CuSum is asymptotically optimal. The optimal weighted mixture CuSum algorithm also has a good performance in the static setting.

We then compare the three algorithms under the binomial distributions. We consider two cases with different K . For the case where there are two types of sensors, for type I sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.2)$ and $\mathcal{B}(10, 0.5)$, for type II sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.8)$ and $\mathcal{B}(10, 0.6)$, respectively. Here \mathcal{B} denotes binomial distribution, the first parameter denotes the number of trials and the second parameter denotes the success probability of each trial. We plot the results for the case where each type has four sensors in Fig. 2. For the case where there are four types of sensors, for type I sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.2)$ and $\mathcal{B}(10, 0.8)$, for type II sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.3)$ and $\mathcal{B}(10, 0.6)$, for type III sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.5)$ and $\mathcal{B}(10, 0.9)$, for type IV sensors, the pre- and post-change

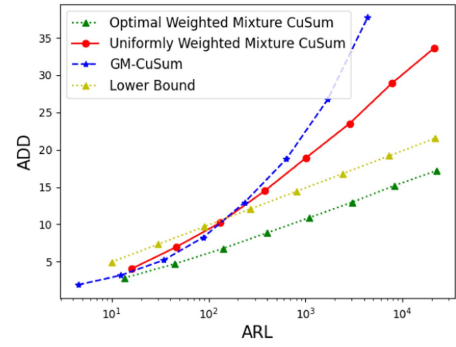
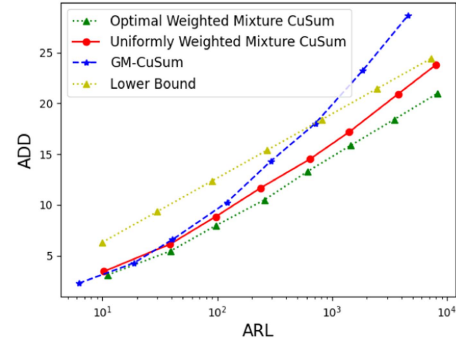
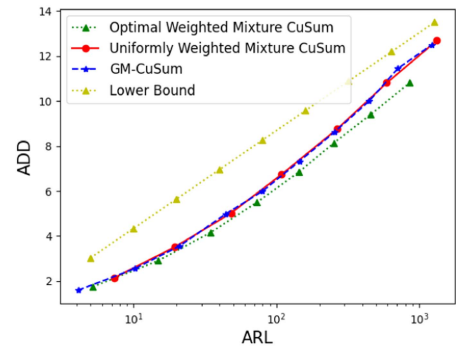
Fig. 2. Static setting: $n = 8, K = 2$.Fig. 3. Static setting: $n = 8, K = 4$.

distributions are $\mathcal{B}(10, 0.4)$ and $\mathcal{B}(10, 0.7)$ respectively. We plot the results for the case where each type has two sensors in Fig. 3. We use Monte-Carlo to obtain the optimal weight for our optimal weighted mixture CuSum algorithm. It can be seen from Figs. 2 and 3 that with the same false alarm rate, our GM-CuSum has the lowest WADD, which implies that it detects the anomaly with the smallest detection delay. Moreover, the relationship between the WADD and log of the WARL is linear, and the slope of the GM-CuSum matches with the one of the lower bound, validating its asymptotic optimality.

B. Dynamic Anomaly Detection

For detecting the dynamic anomaly, we use the same parameters of distributions as in the static setting.

We compare our optimal weighted mixture CuSum algorithm with a uniformly weighted mixture CuSum, i.e., replace β^* in (15) with $\beta = (\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$ and the GM-CuSum under the Gaussian distribution. We plot the average detection delay (ADD) and the average run length (ARL) for some randomly generated trajectories since the worst-case trajectory is difficult to simulate. It can be seen from Fig. 4 that the weighted mixture CuSum algorithm outperforms the uniformly weighted mixture CuSum algorithm and the GM-CuSum. Therefore, in the dynamic setting, the optimal weighted mixture CuSum algorithm detects the presence of the anomaly with the lowest detection delay. Moreover, the slope of our optimal weighted mixture CuSum algorithm matches with the one of the lower bound, which validates its asymptotic optimality.

Fig. 4. Dynamic setting: $n = 4, K = 4$.Fig. 5. Dynamic setting: $n = 8, K = 2$.Fig. 6. Dynamic setting: $n = 8, K = 4$.

We then compare the three algorithms using the binomial distributions as in the static setting. It can be seen from Figs. 5 and 6 that our optimal weighted mixture CuSum algorithm outperforms the uniformly weighted mixture CuSum algorithm and the GM-CuSum since with the same false alarm rate, the optimal weighted mixture CuSum algorithm has the smallest detection delay. The relationship between the ADD and log of the ARL is linear. Moreover, the slope of our optimal weighted mixture CuSum algorithm matches the theoretical lower bound, which demonstrates its asymptotic optimality. It can also be observed that the GM-CuSum algorithm does not perform well under the dynamic setting.

C. Moving Target Detection With Unlabeled Samples

In this section, we consider a practical application of target detection [7], [44]. For simplicity, consider a 3×3 grid. One

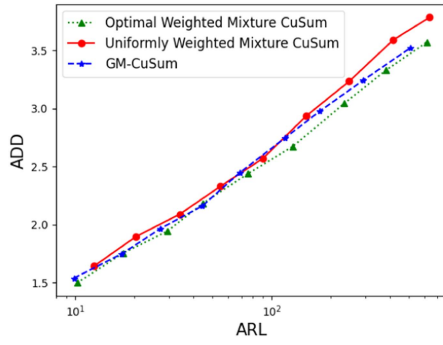


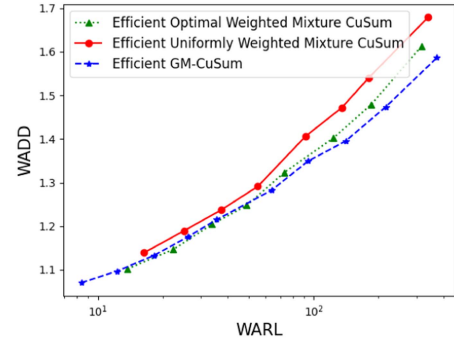
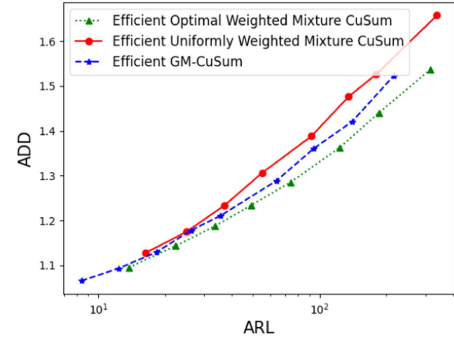
Fig. 7. Moving target detection.

sensor is deployed at the center of each cell. The fusion center can only collect unlabeled data. Before the target presents, at each time, sensors will send exponentially distributed samples with mean 1 to the fusion center. When the target appears, the distribution of each sample changes to another exponential distribution with an elevated mean, depending on the distance between the cell and the target. The target moves across different cells at different time. In this case, the mixture likelihood $\tilde{\mathbb{P}}^k(X^n)$ can be computed when the target lies in a specific cell and the GM-CuSum algorithm and the optimal weighted mixture CuSum algorithm can further be designed. We compare the GM-CuSum algorithm, the uniformly weighted mixture CuSum algorithm and the optimal weighted mixture CuSum algorithm. From Fig. 7, with the same false alarm rate, our optimal weighted mixture CuSum has the smallest detection delay which implies that it detects the presence of the target quickly, e.g., with an ARL of 10^3 , it only takes about 3.5 samples to detect the target.

D. Computationally Efficient Approximation

In this section, we implement the computationally efficient approximation in Section V by replacing the test statistics in uniformly weighted mixture CuSum algorithm, GM-CuSum algorithm and optimal weighted mixture CuSum with the value of the optimization problem in (18) and demonstrate their performance when n is large.

Let $n = 60, K = 3$ and $n_1 = n_2 = n_3 = 20$. After the anomaly emerges, the distribution of all the sensors of an unknown type changes. Our GM-CuSum algorithm and optimal weighted mixture CuSum algorithm can be designed with a slight modification of the mixture likelihood $\tilde{\mathbb{P}}^k(X^n)$. For type I sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.45)$ and $\mathcal{B}(10, 0.6)$, for type II sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.4)$ and $\mathcal{B}(10, 0.6)$, and for type III sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.7)$ and $\mathcal{B}(10, 0.55)$ respectively. The optimal weight for our optimal weighted mixture CuSum algorithm is obtained by Monte-Carlo. We consider the static setting, and plot the WADD as a function of the WARL. It can be seen from Fig. 8 that the GM-CuSum has the smallest WADD given the same WARL. Also the WADD is small (it only takes around 1.6 samples to detect the target for a given a WARL at 4×10^2). Moreover, our GM-CuSum has

Fig. 8. Static setting: $n = 60, K = 3$.Fig. 9. Dynamic setting: $n = 60, K = 3$.

the best performance for detecting the static anomaly when n is large.

We then consider the dynamic setting, we plot the ADD as a function of ARL for some random trajectories. It can be seen from Fig. 9 that with the same ARL, our optimal weighted mixture CuSum algorithm has the smallest detection delay and thus has the best performance for detecting the dynamic anomaly when n is large.

VII. CONCLUSION

In this paper, we studied the problem of quickest detection of an anomaly in networks with unlabeled samples. We first investigated the case with a static anomaly. A GM-CuSum algorithm was proposed and shown to be second-order asymptotically optimal. We then extended our study to the case with a dynamic anomaly, that is, the affected sensor changes with time. We proposed an optimal weighted mixture CuSum algorithm, and proved that it is first-order asymptotically optimal. Our approaches provide useful insights for general (sequential) statistical inference problems with unlabeled samples.

APPENDIX A

Before the anomaly emerges, i.e., $t < \nu$, there are n_k sensors in group k , $\forall 1 \leq k \leq K$, and 0 sensors in group k , $\forall K < k \leq 2K$. Then, there are in total $\binom{n}{n_1, \dots, n_K}$ possible $\sigma_t^{s[t]}$: $\{1, \dots, n\} \rightarrow \{1, \dots, K\}$ satisfying $|\{i : \sigma_t^{s[t]}(i) = k\}| = n_k$, for any $k = 1, \dots, K$. We denote the collection of all such labels by $\mathcal{S}_{n,0}$. After the anomaly emerges, i.e., $t \geq \nu$, one sensor of type $s[t] \neq 0$ is affected by anomaly. Therefore, the

number of sensors in group $s[t]$ and $s[t] + K$ are $n_{s[t]} - 1$ and 1 respectively. Then, there are $(n_1, \dots, n_{s[t]} - 1, \dots, n_K, 1)$ possible $\sigma_t^{s[t]}: \{1, \dots, n\} \rightarrow \{1, \dots, K, s[t] + K\}$ satisfying

$$\left| \left\{ i : \sigma_t^{s[t]}(i) = k \right\} \right| = \begin{cases} n_k, & \text{if } 1 \leq k \leq K \text{ and } k \neq s[t], \\ n_k - 1, & \text{if } k = s[t], \\ 1, & \text{if } k = s[t] + K, \\ 0, & \text{otherwise.} \end{cases}$$

We then denote the collection of all such labels by $\mathcal{S}_{n,s[t]}$.

Before the anomaly emerges, i.e., $t < \nu$, the samples $X^n[t]$ follows the distribution

$$\mathbb{P}_{0,\sigma_t^0}(X^n[t]) = \prod_{i=1}^n p_{0,\sigma_t^0(i)}(X_i[t]), \quad (19)$$

for some unknown $\sigma_t^0 \in \mathcal{S}_{n,0}$. At time $t \geq \nu$, $X^n[t]$ follows the distribution

$$\begin{aligned} \mathbb{P}_{\sigma_t^{s[t]}}^{s[t]}(X^n[t]) &\triangleq \prod_{i:\sigma_t^{s[t]}(i) \leq K} p_{0,\sigma_t^{s[t]}(i)}(X_i[t]) \\ &\times \prod_{i:\sigma_t^{s[t]}(i) > K} p_{1,\sigma_t^{s[t]}(i)-K}(X_i[t]), \end{aligned} \quad (20)$$

for some unknown $\sigma_t^{s[t]} \in \mathcal{S}_{n,s[t]}$.

APPENDIX B PROOF OF THEOREM 1

Consider a simple QCD problem with a pre-change distribution $\tilde{\mathbb{P}}_0$ and a post-change distribution $\tilde{\mathbb{P}}^k$, respectively. Let

$$\begin{aligned} \widetilde{\text{WADD}}_k(\tau) &= \sup_{\nu \geq 1} \text{esssup}_{\tilde{\mathbb{E}}^{k,\nu}} \left[(\tau - \nu)^+ |\tilde{\mathbf{X}}^n[1, \nu - 1]| \right], \\ \widetilde{\text{ARL}}(\tau) &= \tilde{\mathbb{E}}^\infty[\tau], \end{aligned} \quad (21)$$

where $\tilde{\mathbb{E}}^{k,\nu}$ denotes the expectation when the change is at ν , the pre- and post-change distributions are $\tilde{\mathbb{P}}_0$ and $\tilde{\mathbb{P}}^k$, and $\tilde{\mathbf{X}}^n[t]$ for $1 \leq t \leq \nu - 1$ are i.i.d. from $\tilde{\mathbb{P}}_0$, $\tilde{\mathbb{E}}^\infty$ denotes the expectation when samples are generated according to $\tilde{\mathbb{P}}_0$.

For any $1 \leq k \leq K$, consider another QCD problem with a pre-change distribution $\mathbb{P}_{0,\sigma_t^0}$ and a post-change distribution $\mathbb{P}_{\sigma_t^k}^k$, respectively. For this pair of pre- and post-change distributions, define

$$\begin{aligned} \text{WADD}_k(\tau) &= \sup_{\nu \geq 1} \sup_{\Omega_k} \text{esssup}_{\mathbb{E}_{\Omega_k}^{k,\nu}} \left[(\tau - \nu)^+ |\mathbf{X}^n[1, \nu - 1]| \right], \\ \text{WARL}(\tau) &= \inf_{\Omega} \mathbb{E}_{\Omega}^\infty[\tau]. \end{aligned} \quad (22)$$

For any $1 \leq k \leq K$ and any τ satisfying $\text{WARL}(\tau) \geq \gamma$, it can be shown that

$$\begin{aligned} \text{WADD}(\tau) &= \sup_{k \in \mathcal{K}} \text{WADD}_k(\tau) \\ &\geq \sup_{\nu \geq 1} \sup_{\Omega_k} \text{esssup}_{\mathbb{E}_{\Omega_k}^{k,\nu}} \left[(\tau - \nu)^+ |\mathbf{X}^n[1, \nu - 1]| \right] \\ &\geq \sup_{\nu \geq 1} \text{esssup}_{\tilde{\mathbb{E}}^{k,\nu}} \left[(\tau - \nu)^+ |\tilde{\mathbf{X}}^n[1, \nu - 1]| \right] \end{aligned}$$

$$= \widetilde{\text{WADD}}_k(\tau). \quad (23)$$

The second inequality is due to the fact that for any τ , $\text{WADD}_k(\tau) \geq \widetilde{\text{WADD}}_k(\tau)$ [17, eq. (18)]. Similarly, we have that for any τ , $\text{WARL}(\tau) \leq \widetilde{\text{ARL}}(\tau)$ [17, eq. (18)]. It then follows that for any $k \in \mathcal{K}$,

$$\begin{aligned} \inf_{\tau: \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau) &\geq \inf_{\tau: \widetilde{\text{ARL}}(\tau) \geq \gamma} \widetilde{\text{WADD}}_k(\tau) \\ &\geq \frac{\log \gamma}{I_k} + O(1), \text{ as } \gamma \rightarrow \infty. \end{aligned} \quad (24)$$

The last inequality is due to the universal lower bound on WADD for a simple QCD problem [37]. We then have that

$$\inf_{\tau: \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau) \geq \frac{\log \gamma}{I^*} + O(1), \text{ as } \gamma \rightarrow \infty. \quad (25)$$

APPENDIX C PROOF OF THEOREM 2

For any $m \geq 0$, let $r_0 = 0$ and define the stopping time

$$r_{m+1} = \inf \left\{ t > r_m : \sup_k \sum_{i=r_m+1}^t \log \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\tilde{\mathbb{P}}_0(X_i^n)} \leq 0 \right\}. \quad (26)$$

For any permutation $\pi(X^n) = (X_{\pi(1)}, X_{\pi(2)}, \dots, X_{\pi(n)})$, we have that $\log \frac{\tilde{\mathbb{P}}^k(X^n)}{\tilde{\mathbb{P}}_0(X^n)} = \log \frac{\tilde{\mathbb{P}}^k(\pi(X^n))}{\tilde{\mathbb{P}}_0(\pi(X^n))}$. For any π , let $\hat{\sigma}^0 = \sigma^0 \circ \pi$, where “ \circ ” denotes the composition of two functions. Then $\mathbb{E}_{0,\sigma^0} \left[\log \frac{\tilde{\mathbb{P}}^k(\pi(X^n))}{\tilde{\mathbb{P}}_0(\pi(X^n))} \right] = \mathbb{E}_{0,\sigma^0 \circ \pi} \left[\log \frac{\tilde{\mathbb{P}}^k(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] = \mathbb{E}_{0,\hat{\sigma}^0} \left[\log \frac{\tilde{\mathbb{P}}^k(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right]$. For any $\hat{\sigma}^0 \in \mathcal{S}_{n,0}$, a π can always be found so that $\sigma^0 \circ \pi = \hat{\sigma}^0$. Thus, for any $\sigma^0, \hat{\sigma}^0 \in \mathcal{S}_{n,0}$, $\mathbb{E}_{0,\hat{\sigma}^0} \left[\log \frac{\tilde{\mathbb{P}}^k(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] = \mathbb{E}_{0,\sigma^0} \left[\log \frac{\tilde{\mathbb{P}}^k(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right]$.

We then have that for any $\sigma^0 \in \mathcal{S}_{n,0}$,

$$\begin{aligned} \mathbb{E}_{0,\sigma^0} \left[\frac{\tilde{\mathbb{P}}^k(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] &= \frac{1}{|\mathcal{S}_{n,0}|} \sum_{\sigma^0 \in \mathcal{S}_{n,0}} \mathbb{E}_{0,\sigma^0} \left[\frac{\tilde{\mathbb{P}}^k(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] \\ &= \int \frac{\tilde{\mathbb{P}}^k(x^n)}{\tilde{\mathbb{P}}_0(x^n)} \cdot \tilde{\mathbb{P}}_0(x^n) dx^n = 1. \end{aligned} \quad (27)$$

Therefore, for any Ω and $t > r_m$,

$$\begin{aligned} \mathbb{E}_{\Omega}^\infty \left[\prod_{i=r_m+1}^{t+1} \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\tilde{\mathbb{P}}_0(X_i^n)} \middle| \mathcal{F}_t \right] &= \mathbb{E}_{\Omega}^\infty \left[\prod_{i=r_m+1}^t \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\tilde{\mathbb{P}}_0(X_i^n)} \middle| \mathcal{F}_t \right] \cdot \mathbb{E}_{0,\sigma^0} \left[\frac{\tilde{\mathbb{P}}^k(X_{t+1}^n)}{\tilde{\mathbb{P}}_0(X_{t+1}^n)} \middle| \mathcal{F}_t \right] \\ &= \prod_{i=r_m+1}^t \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\tilde{\mathbb{P}}_0(X_i^n)} \cdot \mathbb{E}_{0,\sigma^0} \left[\frac{\tilde{\mathbb{P}}^k(X_{t+1}^n)}{\tilde{\mathbb{P}}_0(X_{t+1}^n)} \right] \\ &= \prod_{i=r_m+1}^t \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\tilde{\mathbb{P}}_0(X_i^n)}. \end{aligned} \quad (28)$$

Therefore, $\left\{ \prod_{i=r_m+1}^t \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\mathbb{P}_0(X_i^n)}, \mathcal{F}_t, t > r_m \right\}$ is a martingale under \mathbb{P}_Ω^∞ for any Ω with mean 1.

We then have that for any Ω ,

$$\begin{aligned} & \mathbb{P}_\Omega^\infty \left\{ \sup_k \sum_{i=r_m+1}^t \log \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\mathbb{P}_0(X_i^n)} \geq b \text{ for some } t > r_m \middle| \mathcal{F}_{r_m} \right\} \\ & \leq \sum_{k=1}^K \mathbb{P}_\Omega^\infty \left\{ \prod_{i=r_m+1}^t \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\mathbb{P}_0(X_i^n)} \geq e^b \text{ for some } t > r_m \middle| \mathcal{F}_{r_m} \right\} \\ & \leq K \frac{\mathbb{E}_{0,\sigma^0} \left[\frac{\tilde{\mathbb{P}}^k(X_{r_m+1}^n)}{\mathbb{P}_0(X_{r_m+1}^n)} \right]}{e^b} = K e^{-b}, \end{aligned} \quad (29)$$

where the last inequality is due to Doob's submartingale inequality [41] and the optional sampling theorem [41].

Let $M = \inf \left\{ m \geq 0 : r_m < \infty, \sup_k \sum_{i=r_m+1}^t \log \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\mathbb{P}_0(X_i^n)} \geq b \text{ for some } t > r_m \right\}$. We have that for any Ω ,

$$\begin{aligned} & \mathbb{P}_\Omega^\infty (M \geq m+1 | \mathcal{F}_{r_m}) \\ & \geq \mathbb{P}_\Omega^\infty \left\{ \sup_k \sum_{i=r_m+1}^t \log \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\mathbb{P}_0(X_i^n)} < b \text{ for all } t > r_m \middle| \mathcal{F}_{r_m} \right\} \\ & \geq 1 - K e^{-b}. \end{aligned} \quad (30)$$

We then have that for any Ω ,

$$\begin{aligned} \mathbb{P}_\Omega^\infty (M > m) &= \mathbb{E}_\Omega^\infty [\mathbb{P}_\Omega^\infty (M \geq m+1 | \mathcal{F}_{r_m}) \cdot \mathbb{1}_{\{M \geq m\}}] \\ &\geq (1 - K e^{-b}) \mathbb{P}_\Omega^\infty (M > m-1) \\ &\geq (1 - K e^{-b})^m. \end{aligned} \quad (31)$$

It then follows that

$$\begin{aligned} \text{WARL}(T_G) &= \inf_\Omega \mathbb{E}_\Omega^\infty [T_G] \geq \inf_\Omega \mathbb{E}_\Omega^\infty [M] \\ &\geq \inf_\Omega \sum_{m=0}^{\infty} \mathbb{P}_\Omega^\infty (M > m) \geq \sum_{m=0}^{\infty} (1 - K e^{-b})^m = \frac{e^b}{K}. \end{aligned} \quad (32)$$

Let $b = \log K\gamma$, we have that $\text{WARL}(T_G) \geq \gamma$. Let T_k be the mixture CuSum algorithm for problem in (22):

$$T_k = \inf \left\{ t : \max_{1 \leq j \leq t} \sum_{i=j}^t \log \frac{\tilde{\mathbb{P}}^k(X_i^n)}{\mathbb{P}_0(X_i^n)} \geq b \right\}. \quad (33)$$

It then follows that for any $1 \leq k \leq K$,

$$\begin{aligned} \text{WADD}_k(T_G) &= \sup_{\nu \geq 1} \sup_{\Omega_k} \text{esssup} \mathbb{E}_{\Omega_k}^{k,\nu} [(T_G - \nu)^+ | \mathbf{X}^n[1, \nu-1]] \\ &\leq \sup_{\nu \geq 1} \sup_{\Omega_k} \text{esssup} \mathbb{E}_{\Omega_k}^{k,\nu} [(T_k - \nu)^+ | \mathbf{X}^n[1, \nu-1]] \\ &\leq \frac{\log b}{I_k} + O(1), \end{aligned} \quad (34)$$

where the last equality is because of the exact optimality of the mixture CuSum algorithm (see Theorem 1 in [17]).

To satisfy the WARL constraint, set $b = \log K\gamma$, we have

$$\begin{aligned} \text{WADD}(T_G) &= \sup_{k \in \mathcal{K}} \text{WADD}_k(T_G) \leq \sup_{k \in \mathcal{K}} \text{WADD}_k(T_k) \\ &= \sup_{k \in \mathcal{K}} \frac{\log K\gamma}{I_k} + O(1) = \frac{\log \gamma}{I^*} + O(1), \text{ as } \gamma \rightarrow \infty. \end{aligned} \quad (35)$$

APPENDIX D

PROOF OF THEOREM 4

For any trajectory \mathbf{s} and stopping time τ , define the WADD and WARL

$$\begin{aligned} \text{WADD}_{\mathbf{s}}(\tau) &= \sup_{\nu \geq 1} \sup_{\Omega_{\mathbf{s}}} \text{esssup} \mathbb{E}_{\Omega_{\mathbf{s}}}^{s,\nu} [(\tau - \nu)^+ | \mathbf{X}^n[1, \nu-1]], \\ \text{ARL}(\tau) &= \inf_{\Omega} \mathbb{E}_{\Omega}^{\infty} [\tau]. \end{aligned} \quad (36)$$

Consider QCD problem with a pre-change distribution $\tilde{\mathbb{P}}_0 = \frac{1}{|\mathcal{S}_{n,0}|} \sum_{\sigma^0 \in \mathcal{S}_{n,0}} \mathbb{P}_{0,\sigma^0}$ and a post-change distribution $\tilde{\mathbb{P}}^{s[t]} = \frac{1}{|\mathcal{S}_{n,s[t]}|} \sum_{\sigma^{s[t]} \in \mathcal{S}_{n,s[t]}} \mathbb{P}_{\sigma^{s[t]}}^{s[t]}$, respectively. For this pair of pre- and post-change distributions and any trajectory \mathbf{s} , define

$$\begin{aligned} \widetilde{\text{WADD}}_{\mathbf{s}}(\tau) &= \sup_{\nu \geq 1} \text{esssup} \tilde{\mathbb{E}}^{s,\nu} [(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1]], \\ \widetilde{\text{ARL}}(\tau) &= \tilde{\mathbb{E}}^{\infty} [\tau]. \end{aligned} \quad (37)$$

where $\tilde{\mathbb{E}}^{s,\nu}$ denotes the expectation when change point is ν , before the change point, the data follows distribution $\tilde{\mathbb{P}}_0$ and after the change point, at time t , the data follows the distribution $\tilde{\mathbb{P}}^{s[t]}$, and $\tilde{\mathbf{X}}^n[1, \nu-1]$ are i.i.d. from $\tilde{\mathbb{P}}_0$; and $\tilde{\mathbb{E}}^{\infty}$ denote the expectation when the data follows distribution $\tilde{\mathbb{P}}_0$.

Consider another QCD problem with pre-change distribution $\tilde{\mathbb{P}}_0$ and post-change distribution $\tilde{\mathbb{P}}^{\beta^*}$. Under this pair of pre- and post-change distributions, define

$$\begin{aligned} \widetilde{\text{WADD}}_{\beta^*}(\tau) &= \sup_{\nu \geq 1} \text{esssup} \tilde{\mathbb{E}}^{\beta^*,\nu} [(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1]], \\ \widetilde{\text{ARL}}(\tau) &= \tilde{\mathbb{E}}^{\infty} [\tau]. \end{aligned} \quad (38)$$

In QCD problems, ARL only depends on the pre-change distribution. Therefore, for any stopping time τ , problems in (2) and (36) have the same ARL, problems in (37) and (38) have the same ARL. Let \mathcal{C}_γ denotes the collection of all stopping times τ that satisfy $\text{ARL}(\tau) \geq \gamma$ and $\tilde{\mathcal{C}}_\gamma$ denotes the collection of all stopping times τ that satisfy $\widetilde{\text{ARL}}(\tau) \geq \gamma$. Our goal is to prove that

$$\inf_{\tau \in \mathcal{C}_\gamma} \text{WADD}(\tau) \geq \inf_{\tau \in \tilde{\mathcal{C}}_\gamma} \widetilde{\text{WADD}}_{\beta^*}(\tau) \sim \frac{\log \gamma}{I_{\beta^*}} (1 + o(1)). \quad (39)$$

Construct a new sequence of random variables $\{\hat{X}^n[t]\}_{t=1}^{\infty}$. Before the change point, $\hat{X}^n[t]$ are i.i.d. according to the mixture distribution $\tilde{\mathbb{P}}_0 = \frac{1}{|\mathcal{S}_{n,0}|} \sum_{\sigma^0 \in \mathcal{S}_{n,0}} \mathbb{P}_{0,\sigma^0}$. After the change point, i.e., $t \geq \nu$, $\hat{X}^n[t]$ follows the distribution $\mathbb{P}_{\sigma_t^{s[t]}}^{s[t]}$ for some $\sigma_t^{s[t]} \in \mathcal{S}_{n,s[t]}$. Specifically,

$$\hat{X}^n[t] \sim \begin{cases} \tilde{\mathbb{P}}_0, & \text{if } t < \nu, \\ \mathbb{P}_{\sigma_t^{s[t]}}^{s[t]}, & \text{if } t \geq \nu. \end{cases} \quad (40)$$

For any stopping time τ and any \mathbf{s} , let

$$\widehat{\text{WADD}}_{\mathbf{s}}(\tau) = \sup_{\nu \geq 1} \sup_{\sigma_{\nu}^{s[\nu]}, \dots, \sigma_{\infty}^{s[\infty]}} \text{esssup} \mathbb{E}_{\sigma_{\nu}^{s[\nu]}, \dots, \sigma_{\infty}^{s[\infty]}}^{s, \nu} \left[(\tau - \nu)^+ |\widehat{\mathbf{X}}^n[1, \nu - 1] \right], \quad (41)$$

where $\mathbb{E}_{\sigma_{\nu}^{s[\nu]}, \dots, \sigma_{\infty}^{s[\infty]}}^{s, \nu}$ denotes the expectation when the data is distributed according to (40).

Let $\widehat{\text{WADD}}(\tau) = \sup_{\mathbf{s}} \widehat{\text{WADD}}_{\mathbf{s}}(\tau)$. To prove (39), we first show that for any \mathbf{s} , $\text{WADD}_{\mathbf{s}}(\tau) = \widehat{\text{WADD}}_{\mathbf{s}}(\tau)$, and then show that $\widehat{\text{WADD}}_{\mathbf{s}}(\tau) \geq \text{WADD}_{\mathbf{s}}(\tau)$. We complete our proof by showing that for any τ and β , $\widehat{\text{WADD}}(\tau) \geq \widehat{\text{WADD}}_{\beta}(\tau)$.

Step 1: Denote by \mathcal{M} the collection of all $\{\sigma_1^0, \dots, \sigma_{\nu-1}^0\}$, and μ is an element in \mathcal{M} . When the trajectory is \mathbf{s} , denote by $\mathcal{N}_{\mathbf{s}}$ the collection of all $\{\sigma_{\nu}^{s[\nu]}, \dots, \sigma_{\infty}^{s[\infty]}\}$, and ω is an element in $\mathcal{N}_{\mathbf{s}}$. Then, the $\text{WADD}_{\mathbf{s}}$ can be written as

$$\begin{aligned} \text{WADD}_{\mathbf{s}}(\tau) &= \sup_{\nu \geq 1} \sup_{\Omega_{\mathbf{s}}} \text{esssup} \mathbb{E}_{\Omega_{\mathbf{s}}}^{s, \nu} [(\tau - \nu)^+ |\mathbf{X}^n[1, \nu - 1]] \\ &= \sup_{\nu \geq 1} \sup_{\omega \in \mathcal{N}_{\mathbf{s}}} \sup_{\mu \in \mathcal{M}} \text{esssup} \mathbb{E}_{\omega}^{s, \nu} [(\tau - \nu)^+ |\mathbf{X}^n[1, \nu - 1]], \end{aligned}$$

where $\mathbb{E}_{\omega}^{s, \nu}$ denotes the expectation when change point is ν , the trajectory is \mathbf{s} , and after the change point, the data follows distribution $\prod_{t=\nu}^{\infty} \mathbb{P}_{\sigma_t^{s[t]}}^{s[t]}$. We note that $\widehat{\mathbf{X}}^n[t]$ and $\mathbf{X}^n[t]$, for $t \geq \nu$, have the same distribution $\mathbb{P}_{\sigma_t^{s[t]}}^{s[t]}$. Therefore, the difference between $\text{WADD}_{\mathbf{s}}$ and $\widehat{\text{WADD}}_{\mathbf{s}}$ lies in that they take esssup with respect to different distributions, i.e., the distributions of $\mathbf{X}^n[1, \nu - 1]$ and $\widehat{\mathbf{X}}^n[1, \nu - 1]$ are different. Let $f_{\omega}(\mathbf{X}^n[1, \nu - 1])$ denote $\mathbb{E}_{\omega}^{s, \nu} [(\tau - \nu)^+ |\mathbf{X}^n[1, \nu - 1]]$. Then, $\text{WADD}_{\mathbf{s}}$ and $\widehat{\text{WADD}}_{\mathbf{s}}$ can be written as

$$\begin{aligned} \text{WADD}_{\mathbf{s}}(\tau) &= \sup_{\nu \geq 1} \sup_{\omega \in \mathcal{N}_{\mathbf{s}}} \sup_{\mu \in \mathcal{M}} \text{esssup} f_{\omega}(\mathbf{X}^n[1, \nu - 1]), \\ \widehat{\text{WADD}}_{\mathbf{s}}(\tau) &= \sup_{\nu \geq 1} \sup_{\omega \in \mathcal{N}_{\mathbf{s}}} \text{esssup} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1]). \end{aligned} \quad (42)$$

It then suffices to show that for any $\omega \in \mathcal{N}_{\mathbf{s}}$,

$$\begin{aligned} \sup_{\mu \in \mathcal{M}} \text{esssup} f_{\omega}(\mathbf{X}^n[1, \nu - 1]) &= \text{esssup} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1]). \\ \text{For any } \omega \in \mathcal{N}_{\mathbf{s}} \text{ and } \mu \in \mathcal{M}, \text{ let} \\ b_{\omega, \mu} &= \text{esssup} f_{\omega}(\mathbf{X}^n[1, \nu - 1]) \\ &= \inf \{b : \mathbb{P}_{\mu}(f_{\omega}(\mathbf{X}^n[1, \nu - 1]) > b) = 0\}, \end{aligned} \quad (43)$$

where \mathbb{P}_{μ} denotes the probability measure when the data is generated from $\mathbb{P}_{0, \sigma_1^0}, \dots, \mathbb{P}_{0, \sigma_{\nu-1}^0}$ before change point ν .

Let $b_{\omega}^* = \text{esssup} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1])$. It can be shown that

$$\begin{aligned} b_{\omega}^* &= \inf \left\{ b : \int_{\mathbf{x}^n[1, \nu-1]} \mathbb{1}_{\{f_{\omega}(\mathbf{x}^n[1, \nu-1]) > b\}} \right. \\ &\quad \left. \times d \prod_{t=1}^{\nu-1} \widetilde{\mathbb{P}}_0(x^n(t)) = 0 \right\} \end{aligned}$$

$$\begin{aligned} &= \inf \left\{ b : \int_{\mathbf{x}^n[1, \nu-1]} \mathbb{1}_{\{f_{\omega}(\mathbf{x}^n[1, \nu-1]) > b\}} \right. \\ &\quad \left. \times d \frac{1}{|\mathcal{M}|} \sum_{\mu \in \mathcal{M}} \mathbb{P}_{\mu}(\mathbf{x}^n[1, \nu - 1]) = 0 \right\} \\ &= \inf \left\{ b : \frac{1}{|\mathcal{M}|} \sum_{\mu \in \mathcal{M}} \mathbb{P}_{\mu}(f_{\omega}(\mathbf{X}^n[1, \nu - 1]) > b) = 0 \right\}. \end{aligned}$$

It then follows that for any $\mu \in \mathcal{M}$, and $\omega \in \mathcal{N}_{\mathbf{s}}$, $\mathbb{P}_{\mu}(f_{\omega}(\mathbf{X}^n[1, \nu - 1]) > b_{\omega}^*) = 0$. Therefore, for any $\mu \in \mathcal{M}$, we have that $b_{\omega, \mu} \leq b_{\omega}^*$. Then

$$\sup_{\mu \in \mathcal{M}} b_{\omega, \mu} \leq b_{\omega}^*. \quad (44)$$

Conversely, for any $\mu \in \mathcal{M}$, we have that $\mathbb{P}_{\mu}(f_{\omega}(\mathbf{X}^n[1, \nu - 1]) > \sup_{\mu \in \mathcal{M}} b_{\omega, \mu}) = 0$. Then, $\frac{1}{|\mathcal{M}|} \sum_{\mu \in \mathcal{M}} \mathbb{P}_{\mu}(f_{\omega}(\mathbf{X}^n[1, \nu - 1]) > \sup_{\mu \in \mathcal{M}} b_{\omega, \mu}) = 0$. This further implies that

$$b_{\omega}^* \leq \sup_{\mu \in \mathcal{M}} b_{\omega, \mu}. \quad (45)$$

Combining (44) and (45), we have that $\sup_{\mu \in \mathcal{M}} b_{\omega, \mu} = b_{\omega}^*$, and thus $\sup_{\mu \in \mathcal{M}} \text{esssup} f_{\omega}(\mathbf{X}^n[1, \nu - 1]) = \text{esssup} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1])$. This implies that for any τ ,

$$\text{WADD}_{\mathbf{s}}(\tau) = \widehat{\text{WADD}}_{\mathbf{s}}(\tau). \quad (46)$$

Step 2: The next step is to show $\widehat{\text{WADD}}_{\mathbf{s}}(\tau) \geq \widehat{\text{WADD}}_{\beta}(\tau)$. We first show that $\sup_{\omega \in \mathcal{N}_{\mathbf{s}}} \text{esssup} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1]) \geq \text{esssup} \sup_{\omega \in \mathcal{N}_{\beta}} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1])$. Denote by $\widetilde{\mathbb{P}}^{\nu}$ the probability measure when the change is at ν , the pre- and post-change distributions are $\widetilde{\mathbb{P}}_0$ and $\widetilde{\mathbb{P}}^{s[t]}$ at time t , respectively. Let $\hat{b} = \sup_{\omega \in \mathcal{N}_{\beta}} \text{esssup} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1])$. For any $\omega \in \mathcal{N}_{\mathbf{s}}$, we have that $\widetilde{\mathbb{P}}^{\nu}(f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1]) > \hat{b}) = 0$. Since $\mathcal{N}_{\mathbf{s}}$ is countable, and a countable union of sets of measure zero has measure zero, we then have that

$$\begin{aligned} &\widetilde{\mathbb{P}}^{\nu} \left(\sup_{\omega \in \mathcal{N}_{\mathbf{s}}} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1]) > \hat{b} \right) \\ &\leq \widetilde{\mathbb{P}}^{\nu} \left(\cup_{\omega \in \mathcal{N}_{\mathbf{s}}} \{f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1]) > \hat{b}\} \right) = 0. \end{aligned} \quad (47)$$

Therefore,

$$\hat{b} \geq \text{esssup} \sup_{\omega \in \mathcal{N}_{\beta}} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1]). \quad (48)$$

Before the change point ν , $\widehat{\mathbf{X}}^n[t]$ and $\widetilde{\mathbf{X}}^n[t]$ follow the same distribution. For any $T \geq \nu + 1$, we have that

$$\begin{aligned} &\sup_{\{\sigma_{\nu}^{s[\nu]}, \dots, \sigma_T^{s[T]}\} \in \mathcal{S}_{n, s[\nu]} \times \dots \times \mathcal{S}_{n, s[T]}} \sum_{t=\nu+1}^T (t - \nu) \\ &\quad \times \mathbb{P}_{\sigma_{\nu}^{s[\nu]}, \dots, \sigma_T^{s[T]}}^{s, \nu}(\tau = t | \widehat{\mathbf{X}}^n[1, \nu - 1]) \end{aligned}$$

$$\begin{aligned}
&\geq \sum_{t=\nu+1}^T (t-\nu) \frac{1}{|\mathcal{S}_{n,s[\nu]}| \times \cdots \times |\mathcal{S}_{n,s[T]}|} \\
&\quad \times \sum_{\substack{\{\sigma_\nu^{s[\nu]}, \dots, \sigma_T^{s[T]}\} \\ \in \mathcal{S}_{n,s[\nu]} \times \cdots \times \mathcal{S}_{n,s[T]}}} \mathbb{P}_{\sigma_\nu^{s[\nu]}, \dots, \sigma_T^{s[T}}}^{s, \nu} \left(\tau = t | \tilde{\mathbf{X}}^n[1, \nu-1] \right) \\
&= \sum_{t=\nu+1}^T (t-\nu) \tilde{\mathbb{P}}^{s, \nu} \left(\tau = t | \tilde{\mathbf{X}}^n[1, \nu-1] \right), \quad (49)
\end{aligned}$$

where $\mathbb{P}_{\sigma_\nu^{s[\nu]}, \dots, \sigma_T^{s[T]}}^{s, \nu}$ denotes the probability measure when change point is ν , the trajectory is \mathbf{s} , the observations from time ν to time T are generated according to $\mathbb{P}_{\sigma_\nu^{s[\nu]}, \dots, \sigma_T^{s[T]}}^{s[\nu]}, \dots, \mathbb{P}_{\sigma_T^{s[T]}}^{s[T]}$. As $T \rightarrow \infty$, we have that

$$f_\omega \left(\tilde{\mathbf{X}}^n[1, \nu-1] \right) \geq \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right], \quad (50)$$

From (48) and (50), we have that

$$\begin{aligned}
\widehat{\text{WADD}}_{\mathbf{s}}(\tau) &\geq \text{esssup} \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right] \\
&= \widehat{\text{WADD}}_{\mathbf{s}}(\tau). \quad (51)
\end{aligned}$$

Combining (46) and (51), it follows that

$$\text{WADD}_{\mathbf{s}}(\tau) = \widehat{\text{WADD}}_{\mathbf{s}}(\tau) \geq \widehat{\text{WADD}}_{\mathbf{s}}(\tau). \quad (52)$$

This holds for any trajectory \mathbf{s} . It then follows that

$$\begin{aligned}
\text{WADD}(\tau) &\geq \sup_{\nu \geq 0} \sup_{\mathbf{s}} \text{esssup} \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right] \\
&= \widehat{\text{WADD}}(\tau). \quad (53)
\end{aligned}$$

Step 3: The last step is to show that for any τ and any β , $\widehat{\text{WADD}}(\tau) \geq \widehat{\text{WADD}}_{\beta}(\tau)$. Firstly, we will show that

$$\begin{aligned}
&\sup_{\mathbf{s}} \text{esssup} \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right] \\
&\geq \text{esssup} \sup_{\mathbf{s}} \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right]. \quad (54)
\end{aligned}$$

Let $c = \sup_{\mathbf{s}} \text{esssup} \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right]$. Denote by $\Lambda_{\mathbf{s}}$ the collection of all trajectory \mathbf{s} . For any \mathbf{s} , we have that

$$\tilde{\mathbb{P}}^{\nu} \left(\tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right] > c \right) = 0. \quad (55)$$

Since $\Lambda_{\mathbf{s}}$ is countable, it then follows that

$$\begin{aligned}
&\tilde{\mathbb{P}}^{\nu} \left(\sup_{\mathbf{s}} \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right] > c \right) \\
&\leq \tilde{\mathbb{P}}^{\nu} \left(\bigcup_{\mathbf{s} \in \Lambda_{\mathbf{s}}} \left\{ \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right] > c \right\} \right) = 0.
\end{aligned}$$

Therefore,

$$\begin{aligned}
c &= \sup_{\mathbf{s}} \text{esssup} \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right] \\
&\geq \text{esssup} \sup_{\mathbf{s}} \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right]. \quad (56)
\end{aligned}$$

For any $T \geq \nu + 1$, we have that

$$\begin{aligned}
&\sup_{\mathbf{s}} \sum_{t=\nu+1}^T (t-\nu) \tilde{\mathbb{P}}^{s[\nu], \dots, s[T]} \left(\tau = t | \tilde{\mathbf{X}}^n[1, \nu-1] \right) \\
&\geq \sum_{t=\nu+1}^T (t-\nu) \sum_{\{s[\nu], \dots, s[T]\} \in \Lambda_{\mathbf{s}}^{\otimes (T-\nu+1)}} \beta_{s[\nu]} \times \cdots \times \beta_{s[T]} \\
&\quad \times \tilde{\mathbb{P}}^{s[\nu], \dots, s[T]} \left(\tau = t | \tilde{\mathbf{X}}^n[1, \nu-1] \right) \\
&= \sum_{t=\nu+1}^T (t-\nu) \tilde{\mathbb{P}}^{\beta, \nu} \left(\tau = t | \tilde{\mathbf{X}}^n[1, \nu-1] \right), \quad (57)
\end{aligned}$$

where $\tilde{\mathbb{P}}^{s[\nu], \dots, s[T]}$ denotes the probability measure when the trajectory is \mathbf{s} , the observations from time ν to time T are generated according to $\tilde{\mathbb{P}}^{s[\nu]}, \dots, \tilde{\mathbb{P}}^{s[T]}$. As $T \rightarrow \infty$, we have

$$\begin{aligned}
&\sup_{\mathbf{s}} \tilde{\mathbb{E}}^{s, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right] \\
&\geq \tilde{\mathbb{E}}^{\beta, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right]. \quad (58)
\end{aligned}$$

From (56) and (58), we have that

$$\begin{aligned}
\widehat{\text{WADD}}(\tau) &\geq \text{esssup} \tilde{\mathbb{E}}^{\beta, \nu} \left[(\tau - \nu)^+ | \tilde{\mathbf{X}}^n[1, \nu-1] \right] \\
&= \widehat{\text{WADD}}_{\beta}(\tau). \quad (59)
\end{aligned}$$

Combining (53) and (59), we have that for any τ and β

$$\text{WADD}(\tau) \geq \widehat{\text{WADD}}(\tau) \geq \widehat{\text{WADD}}_{\beta}(\tau). \quad (60)$$

For any $T \geq 1$, we have that

$$\begin{aligned}
&\inf_{\{\sigma_1^0, \dots, \sigma_T^0\} \in \mathcal{S}_{n,0}^{\otimes T}} \sum_{t=1}^T t \mathbb{P}_{\sigma_1^0, \dots, \sigma_T^0}^{\infty}(\tau = t) \\
&\leq \sum_{t=1}^T t \frac{1}{|\mathcal{S}_{n,0}|^T} \sum_{\{\sigma_1^0, \dots, \sigma_T^0\} \in \mathcal{S}_{n,0}^{\otimes T}} \mathbb{P}_{\sigma_1^0, \dots, \sigma_T^0}^{\infty}(\tau = t) \\
&= \sum_{t=1}^T t \tilde{\mathbb{P}}^{\infty}(\tau = t). \quad (61)
\end{aligned}$$

As $T \rightarrow \infty$, we have that $\text{ARL}(\tau) \leq \widehat{\text{ARL}}(\tau)$.

Therefore, for any stopping time τ satisfying $\text{ARL}(\tau) \geq \gamma$, it will also satisfy $\widehat{\text{ARL}}(\tau) \geq \gamma$. We then have that $\mathcal{C}_{\gamma} \subseteq \widehat{\mathcal{C}}_{\gamma}$.

Since (60) holds for any β , it holds for β^* . Problem (38) is a classical QCD problem. We have that for large γ [37],

$$\inf_{\tau \in \mathcal{C}_{\gamma}} \text{WADD}(\tau) \geq \inf_{\tau \in \widehat{\mathcal{C}}_{\gamma}} \widehat{\text{WADD}}_{\beta^*}(\tau) \sim \frac{\log \gamma}{I_{\beta^*}} (1 + o(1)).$$

APPENDIX E PROOF OF LEMMA 1

The minimization of I_{β} is to solve the following problem:

$$\inf_{\beta} I_{\beta}$$

$$\text{s.t. } -\beta_k \leq 0, \text{ for } k \in [1, K], \sum_{k=1}^K \beta_k - 1 = 0. \quad (62)$$

This is a convex optimization problem with linear constraints. Define the Lagrange function $L(\beta, \eta, \mu)$:

$$L(\beta, \eta, \mu) = I_\beta + \eta \left(\sum_{k=1}^K \beta_k - 1 \right) - \sum_{k=1}^K \mu_k \beta_k. \quad (63)$$

The minimizer β^* satisfies the Karush–Kuhn–Tucker(KKT) conditions: $\mu_k, \beta_k^* \geq 0, \mu_k \beta_k^* = 0, \sum_{k=1}^K \beta_k^* - 1 = 0$ and

$$\frac{\partial L}{\partial \beta_k} |_{\beta^*} = \tilde{\mathbb{E}}^k \left[\log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] + 1 + \eta - \mu_k = 0, \quad (64)$$

where $\tilde{\mathbb{E}}^k$ denotes the expectation under the distribution $\tilde{\mathbb{P}}^k$.

When $\beta_k^* > 0$, we have $\mu_k = 0$. Therefore, for any $k, k' \in \mathcal{K}$ with $\beta_k^*, \beta_{k'}^* > 0$, we have

$$\tilde{\mathbb{E}}^k \left[\log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] = \tilde{\mathbb{E}}^{k'} \left[\log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] = -(1 + \eta).$$

The set \mathcal{K} can be divided into two disjoint parts \mathcal{K}_1 and \mathcal{K}_2 : $\beta_k^* > 0$ if $k \in \mathcal{K}_1$ and $\beta_k^* = 0$ if $k \in \mathcal{K}_2$. We have that

$$\begin{aligned} I_{\beta^*} &= \sum_{k \in \mathcal{K}_1} \beta_k^* \tilde{\mathbb{E}}^k \left[\log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] \\ &= \tilde{\mathbb{E}}^k \left[\log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right], \text{ for all } k \in \mathcal{K}_1. \end{aligned} \quad (65)$$

For all $k \in \mathcal{K}_2$, $\beta_k^* = 0$. By the KKT conditions, we have that $\mu_k \geq 0$. Therefore, for any $k \in \mathcal{K}_2$, $\tilde{\mathbb{E}}^k \left[\log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] + 1 + \eta = \mu_k \geq 0$. We then have that for any $k \in \mathcal{K}_2$,

$$\tilde{\mathbb{E}}^k \left[\log \frac{\tilde{\mathbb{P}}^{\beta^*}(X^n)}{\tilde{\mathbb{P}}_0(X^n)} \right] \geq I_{\beta^*}. \quad (66)$$

APPENDIX F PROOF OF THEOREM 5

Due to the fact that $\max_{1 \leq k \leq t+1} \sum_{i=k}^t \ell_{\beta^*}(X_i^n)$ has initial value 0 and remains non-negative, the delay is the largest when the change happens at $\nu = 0$. Therefore, for any s , we have

$$\text{WADD}_s(T_{\beta^*}) = \sup_{\Omega_s} \mathbb{E}_{\Omega_s}^{s,0}[T_{\beta^*}]. \quad (67)$$

For any $T \geq \nu + 1$, we have that

$$\begin{aligned} &\sup_{\{\sigma_1^{s[1]}, \dots, \sigma_T^{s[T]}\} \in \mathcal{S}_{n,s[1]} \times \dots \times \mathcal{S}_{n,s[T]}} \sum_{t=1}^T t \tilde{\mathbb{P}}^{s,0}_{\sigma_1^{s[1]}, \dots, \sigma_T^{s[T]}}(T_{\beta^*} = t) \\ &= \sum_{t=1}^T t \frac{1}{|\mathcal{S}_{n,s[1]}| \times \dots \times |\mathcal{S}_{n,s[T]}|} \\ &\quad \times \sum_{\{\sigma_1^{s[1]}, \dots, \sigma_T^{s[T]}\} \in \mathcal{S}_{n,s[1]} \times \dots \times \mathcal{S}_{n,s[T]}} \mathbb{P}_{\sigma_1^{s[1]}, \dots, \sigma_T^{s[T]}}^{s,0}(T_{\beta^*} = t) \end{aligned}$$

$$= \sum_{t=1}^T t \tilde{\mathbb{P}}^{s,0}(T_{\beta^*} = t). \quad (68)$$

As $T \rightarrow \infty$, we have that $\sup_{\Omega_s} \mathbb{E}_{\Omega_s}^{s,0}[T_{\beta^*}] = \tilde{\mathbb{E}}^{s,0}[T_{\beta^*}] = \widetilde{\text{WADD}}_s(T_{\beta^*})$. For any s , we have $\text{WADD}_s(T_{\beta^*}) = \widetilde{\text{WADD}}_s(T_{\beta^*})$. Therefore, $\text{WADD}(T_{\beta^*}) = \widetilde{\text{WADD}}(T_{\beta^*})$ by taking sup over s on both sides. It then follows that

$$\text{WADD}(T_{\beta^*}) = \widetilde{\text{WADD}}(T_{\beta^*}) = \sup_s \tilde{\mathbb{E}}^{s,0}[T_{\beta^*}]. \quad (69)$$

Let $0 < \epsilon < I_{\beta^*}$ and $n_b = \frac{b}{I_{\beta^*} - \epsilon}$. For any trajectory s , from the sum-integral inequality, we have that

$$\begin{aligned} \tilde{\mathbb{E}}^{s,0} \left[\frac{T_{\beta^*}}{n_b} \right] &= \int_0^\infty \tilde{\mathbb{P}}^{s,0} \left(\frac{T_{\beta^*}}{n_b} > x \right) dx \\ &\leq \sum_{t=1}^\infty \tilde{\mathbb{P}}^{s,0}(T_{\beta^*} > tn_b) + 1. \end{aligned} \quad (70)$$

For any s , we have that

$$\begin{aligned} \tilde{\mathbb{P}}^{s,0}(T_{\beta^*} > tn_b) &= \tilde{\mathbb{P}}^{s,0} \left(\max_{1 \leq k \leq tn_b} \max_{1 \leq i \leq k} \sum_{j=i}^k \ell_{\beta^*}(X_j^n) < b \right) \\ &\leq \tilde{\mathbb{P}}^{s,0} \left(\max_{1 \leq i \leq mn_b} \sum_{j=i}^{mn_b} \ell_{\beta^*}(X_j^n) < b, \forall m \in [t] \right) \\ &\leq \tilde{\mathbb{P}}^{s,0} \left(\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n) < b, \forall m \in [t] \right) \\ &= \tilde{\mathbb{P}}^{s,0} \left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} < I_{\beta^*} - \epsilon, \forall m \in [t] \right) \\ &= \prod_{m=1}^t \tilde{\mathbb{P}}^{s,0} \left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} < I_{\beta^*} - \epsilon \right). \end{aligned} \quad (71)$$

It then follows that

$$\begin{aligned} \sup_s \sum_{t=1}^\infty \tilde{\mathbb{P}}^{s,0}(T_{\beta^*} > tn_b) \\ &\leq \sup_s \sum_{t=1}^\infty \prod_{m=1}^t \tilde{\mathbb{P}}^{s,0} \left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} < I_{\beta^*} - \epsilon \right). \end{aligned}$$

Then we will bound $\tilde{\mathbb{P}}^{s,0} \left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} < I_{\beta^*} - \epsilon \right)$.

Let $I_{s_m} = \tilde{\mathbb{E}}^{s,0} \left[\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} \right]$. From (65) and (66),

$$I_{s_m} = \frac{1}{n_b} \sum_{j=(m-1)n_b+1}^{mn_b} \tilde{\mathbb{E}}^{s[j]} [\ell_{\beta^*}(X_j^n)] \geq I_{\beta^*}. \quad (72)$$

It then follows that for any s and m

$$\tilde{\mathbb{P}}^{s,0} \left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} < I_{\beta^*} - \epsilon \right)$$

$$\begin{aligned} &\leq \tilde{\mathbb{P}}^{s,0} \left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} < I_{s_m} - \epsilon \right) \\ &\leq \tilde{\mathbb{P}}^{s,0} \left(\left| \frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} - I_{s_m} \right| > \epsilon \right). \end{aligned} \quad (73)$$

Assume that $\max_{k \in [1, K]} \mathbb{E}^k[\ell_{\beta^*}(X^n)^2] < \infty$. Let $\sigma^2 = \max_{k \in [1, K]} \text{Var}_{\tilde{\mathbb{P}}^k}(\ell_{\beta^*}(X^n))$ where $\text{Var}_{\tilde{\mathbb{P}}^k}$ denotes the variance under the distribution $\tilde{\mathbb{P}}^k$. By Chebychev's inequality,

$$\begin{aligned} &\tilde{\mathbb{P}}^{s,0} \left(\left| \frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} - I_{s_m} \right| > \epsilon \right) \\ &\leq \text{Var}_{\tilde{\mathbb{P}}^s} \left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\beta^*}(X_j^n)}{n_b} \right) \frac{1}{\epsilon^2} \\ &= \frac{1}{\epsilon^2 n_b^2} \sum_{j=(m-1)n_b+1}^{mn_b} \text{Var}_{\tilde{\mathbb{P}}^{s[j]}}(\ell_{\beta^*}(X_j^n)) \\ &\leq \frac{\sum_{j=(m-1)n_b+1}^{mn_b} \sigma^2}{n_b^2 \epsilon^2} = \frac{\sigma^2}{n_b \epsilon^2}. \end{aligned} \quad (74)$$

Let $\delta = \frac{\sigma^2}{n_b \epsilon^2}$. From (70) and (74), we have that

$$\begin{aligned} \sup_s \tilde{\mathbb{E}}^{s,0} \left[\frac{T_{\beta^*}}{n_b} \right] &\leq 1 + \sup_s \sum_{t=1}^{\infty} \tilde{\mathbb{P}}^{s,0}(T_{\beta^*} > tn_b) \\ &\leq 1 + \sum_{t=1}^{\infty} \left(\frac{\sigma^2}{n_b \epsilon^2} \right)^t = 1 + \sum_{t=1}^{\infty} \delta^t = \frac{1}{1-\delta}. \end{aligned} \quad (75)$$

Therefore, we have

$$\sup_s \tilde{\mathbb{E}}^{s,0} [T_{\beta^*}] \leq \frac{b}{(I_{\beta^*} - \epsilon)(1 - \delta)}. \quad (76)$$

(76) holds for all ϵ . It then follows that as $b \rightarrow \infty$,

$$\text{WADD}(T_{\beta^*}) = \sup_s \tilde{\mathbb{E}}^{s,0} [T_{\beta^*}] \leq \frac{b}{I_{\beta^*}} (1 + o(1)). \quad (77)$$

For the ARL lower bound, for any $T \geq 1$, we have that

$$\begin{aligned} &\inf_{\{\sigma_1^0, \dots, \sigma_T^0\} \in \mathcal{S}_{n,0}^{\otimes T}} \sum_{t=1}^T t \mathbb{P}_{\sigma_1^0, \dots, \sigma_T^0}^{\infty}(T_{\beta^*} = t) \\ &= \sum_{t=1}^T t \frac{1}{|\mathcal{S}_{n,0}|^T} \sum_{\{\sigma_1^0, \dots, \sigma_T^0\} \in \mathcal{S}_{n,0}^{\otimes T}} \mathbb{P}_{\sigma_1^0, \dots, \sigma_T^0}^{\infty}(T_{\beta^*} = t) \\ &= \sum_{t=1}^T t \tilde{\mathbb{P}}^{\infty}(T_{\beta^*} = t). \end{aligned} \quad (78)$$

As $T \rightarrow \infty$, we have that $\text{WARL}(T_{\beta^*}) = \widetilde{\text{ARL}}(T_{\beta^*})$. T_{β^*} is the CuSum algorithm for a simple QCD problem with pre-change distribution $\tilde{\mathbb{P}}_0$ and post-change distribution $\tilde{\mathbb{P}}^{\beta^*}$. From the optimal property of CuSum algorithm in [39] and [45], we have that when $b = \log \gamma$, $\text{WARL}(T_{\beta^*}) = \widetilde{\text{ARL}}(T_{\beta^*}) \geq \gamma$.

REFERENCES

- [1] Z. Sun and S. Zou, "Quickest dynamic anomaly detection in anonymous heterogeneous sensor networks," in *Proc. IEEE Int. Symp. Inf. Theory*, 2021, pp. 106–111.
- [2] T. E. Humphreys et al., "Assessing the spoofing threat: Development of a portable GPS civilian spoofer," in *Proc. 21st Int. Tech. Meeting Satell. Division Inst. Navigation*, 2008, pp. 2314–2325.
- [3] L. Keller, M. Jafari Siavoshani, C. Fragouli, K. Argyraki, and S. Diggavi, "Identity aware sensor networks," in *Proc. IEEE Int. Conf. Commun. Comput. Control Appl.*, 2009, pp. 2177–2185.
- [4] W. N. Chen and I. H. Wang, "Anonymous heterogeneous distributed detection: Optimal decision rules, error exponents, and the price of anonymity," *IEEE Trans. Inf. Theory*, vol. 65, no. 11, pp. 7390–7406, Nov. 2019.
- [5] W. Li and Y. Huang, "Bandwidth-constrained distributed quickest change detection in heterogeneous sensor networks: Anonymous vs non-anonymous settings," 2022, *arXiv:2202.02697*.
- [6] S. Marano and P. K. Willett, "Algorithms and fundamental limits for unlabeled detection using types," *IEEE Trans. Signal Process.*, vol. 67, no. 8, pp. 2022–2035, Apr. 2019.
- [7] S. Marano and P. Willett, "Making decisions by unlabeled bits," *IEEE Trans. Signal Process.*, vol. 68, pp. 2935–2947, 2020.
- [8] J. Unnikrishnan, S. Haghighatshoar, and M. Vetterli, "Unlabeled sensing with random linear measurements," *IEEE Trans. Inf. Theory*, vol. 64, no. 5, pp. 3237–3253, May 2018.
- [9] S. Haghighatshoar and G. Caire, "Signal recovery from unlabeled samples," *IEEE Trans. Signal Process.*, vol. 66, no. 5, pp. 1242–1257, May 2018.
- [10] A. Abid, A. Poon, and J. Zou, "Linear regression with shuffled labels," 2017, *arXiv:1705.01342*.
- [11] V. Emiya, A. Bonnefof, L. Daudet, and R. Gribonval, "Compressed sensing with unknown sensor permutation," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, 2014, pp. 1040–1044.
- [12] Z. Liu and J. Zhu, "Signal detection from unlabeled ordered samples," *IEEE Commun. Lett.*, vol. 22, no. 12, pp. 2431–2434, Dec. 2018.
- [13] A. Pananjady, M. J. Wainwright, and T. A. Courtade, "Linear regression with shuffled data: Statistical and computational limits of permutation recovery," *IEEE Trans. Inf. Theory*, vol. 64, no. 5, pp. 3286–3300, May 2018.
- [14] G. Elhami, A. Scholfield, B. Bejar Haro, and M. Vetterli, "Unlabeled sensing: Reconstruction algorithm and theoretical guarantees," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, 2017, pp. 4566–4570.
- [15] Y. M. Lu and M. N. Do, "A theory for sampling signals from a union of subspaces," *IEEE Trans. Signal Process.*, vol. 56, no. 6, pp. 2334–2345, Jun. 2008.
- [16] G. Wang et al., "Signal amplitude estimation and detection from unlabeled binary quantized samples," *IEEE Trans. Signal Process.*, vol. 66, no. 16, pp. 4291–4303, Aug. 2018.
- [17] Z. Sun, S. Zou, R. Zhang, and Q. Li, "Quickest change detection in anonymous heterogeneous sensor networks," *IEEE Trans. Signal Process.*, vol. 70, pp. 1041–1055, 2022.
- [18] G. Rovatsos, G. V. Moustakides, and V. V. Veeravalli, "Quickest detection of moving anomalies in sensor networks," *IEEE J. Sel. Areas Inf. Theory*, vol. 2, no. 2, pp. 762–773, Jun. 2021.
- [19] A. G. Tartakovsky and V. V. Veeravalli, "Change-point detection in multi-channel and distributed systems," *Appl. Sequential Methodol.: Real-World Examples Data Anal.*, vol. 173, pp. 339–370, 2004.
- [20] A. G. Tartakovsky, B. L. Rozovskii, R. B. Blazek, and H. Kim, "A novel approach to detection of intrusions in computer networks via adaptive sequential and batch-sequential change-point detection methods," *IEEE Trans. Signal Process.*, vol. 54, no. 9, pp. 3372–3382, Sep. 2006.
- [21] Y. Mei, "Efficient scalable schemes for monitoring a large number of data streams," *Biometrika*, vol. 97, no. 2, pp. 419–433, 2010.
- [22] Y. Xie and D. Siegmund, "Sequential multi-sensor change-point detection," *Ann. Statist.*, vol. 41, pp. 670–692, 2013.
- [23] G. Fellouris and G. Sokolov, "Second-order asymptotic optimality in multisensor sequential change detection," *IEEE Trans. Inf. Theory*, vol. 62, no. 6, pp. 3662–3675, Jun. 2016.
- [24] V. Raghavan and V. V. Veeravalli, "Quickest change detection of a Markov process across a sensor array," *IEEE Trans. Inf. Theory*, vol. 56, no. 4, pp. 1961–1981, Apr. 2010.
- [25] O. Hadjiladis, H. Zhang, and H. V. Poor, "One shot schemes for decentralized quickest change detection," *IEEE Trans. Inf. Theory*, vol. 55, no. 7, pp. 3346–3359, Jul. 2009.
- [26] M. Ludkovski, "Bayesian quickest detection in sensor arrays," *Sequential Anal.*, vol. 31, no. 4, pp. 481–504, 2012.

- [27] V. V. Veeravalli, "Decentralized quickest change detection," *IEEE Trans. Inf. Theory*, vol. 47, no. 4, pp. 1657–1665, May 2001.
- [28] A. G. Tartakovsky and V. V. Veeravalli, "Asymptotically optimal quickest change detection in distributed sensor systems," *Sequential Anal.*, vol. 27, no. 4, pp. 441–475, 2008.
- [29] S. Zou, V. V. Veeravalli, J. Li, D. Towsley, and A. Swami, "Distributed quickest detection of significant events in networks," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, 2019, pp. 8454–8458.
- [30] L. Xie, S. Zou, Y. Xie, and V. V. Veeravalli, "Sequential (quickest) change detection: Classical results and new directions," *IEEE J. Sel. Areas Inf. Theory*, vol. 2, no. 2, pp. 494–514, Jun. 2021.
- [31] S. Zou, V. V. Veeravalli, J. Li, and D. Towsley, "Quickest detection of dynamic events in networks," *IEEE Trans. Inf. Theory*, vol. 66, no. 4, pp. 2280–2295, Apr. 2020.
- [32] G. Rovatsos, G. V. Moustakides, and V. V. Veeravalli, "Quickest detection of a dynamic anomaly in a sensor network," in *Proc. IEEE 53rd Asilomar Conf. Signals Syst. Comput.*, 2019, pp. 98–102.
- [33] G. Rovatsos, S. Zou, and V. V. Veeravalli, "Sequential algorithms for moving anomaly detection in networks," *Sequential Anal.*, vol. 39, no. 1, pp. 6–31, 2020.
- [34] S. Zou, G. Fellouris, and V. V. Veeravalli, "Quickest change detection under transient dynamics: Theory and asymptotic analysis," *IEEE Trans. Inf. Theory*, vol. 65, no. 3, pp. 1397–1412, Mar. 2019.
- [35] R. Zhang, Y. Xie, R. Yao, and F. Qiu, "Online detection of cascading change-points," in *Proc. 58th Annu. Allerton Conf. Commun., Control, Comput.*, 2022, pp. 1–6, doi: [10.1109/Allerton49937.2022.9929381](https://doi.org/10.1109/Allerton49937.2022.9929381).
- [36] D. Siegmund and E. S. Venkatraman, "Using the generalized likelihood ratio statistic for sequential detection of a change-point," *Ann. Statist.*, vol. 23, pp. 255–271, 1995.
- [37] T. Leung Lai, "Information bounds and quick detection of parameter changes in stochastic systems," *IEEE Trans. Inf. Theory*, vol. 44, no. 7, pp. 2917–2929, Nov. 1998.
- [38] T. Banerjee and V. V. Veeravalli, "Data-efficient minimax quickest change detection with composite post-change distribution," *IEEE Trans. Inf. Theory*, vol. 61, no. 9, pp. 5172–5184, Sep. 2015.
- [39] G. Lorden, "Procedures for reacting to a change in distribution," *Ann. Math. Statist.*, vol. 42, no. 6, pp. 1897–1908, 1971.
- [40] E. L. Lehmann, J. P. Romano, and G. Casella, *Testing Statistical Hypotheses*. Berlin, Germany: Springer, 2005.
- [41] D. Williams, *Probability With Martingales*. New York, NY, USA: Cambridge Univ. Press, 1991.
- [42] M. Pollak et al., "Optimal detection of a change in distribution," *Ann. Statist.*, vol. 13, no. 1, pp. 206–227, 1985.
- [43] T. M. Cover, *Elements of Information Theory*. New York, NY, USA: Wiley, 2006.
- [44] Z. Sutton, P. Willet, and S. Marano, "Sensor network target detection with unlabeled observations," in *Proc. IEEE Aerosp. Conf.*, 2020, pp. 1–10.
- [45] G. V. Moustakides, "Optimal stopping times for detecting changes in distributions," *Ann. Statist.*, vol. 14, no. 4, pp. 1379–1387, 1986.



Zhongchang Sun (Student, IEEE) received the B.S. degree from the Beijing Institute of Technology, Beijing, China in 2019. He is working toward the Ph.D. degree with the Department of Electrical Engineering, University at Buffalo, the State University of New York, Buffalo, NY, USA. His research interests include hypothesis testing, quickest change detection, and distributionally robust optimization.



Shaofeng Zou received the B.E. degree (with honors) from Shanghai Jiao Tong University, Shanghai, China, in 2011, the Ph.D. degree in electrical and computer engineering from Syracuse University, Syracuse, NY, USA, in 2016. He is currently an Assistant Professor with the Department of Electrical Engineering, University at Buffalo, the State University of New York, Buffalo, NY, USA. During 2016–2018, he was a Postdoctoral Research Associate with the Coordinated Science Lab, University of Illinois at Urbana-Champaign, Champaign, IL, USA. His research interests include reinforcement learning, machine learning, statistical signal processing and information theory. He was the recipient of the National Science Foundation CRII Award in 2019 and the 2023 AAAI Distinguished Paper Award.