# Quickest Anomaly Detection in Sensor Networks With Unlabeled Samples

Zhongchang Sun     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 large-scale sensor networks, samples may lack of 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 large-scale 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 as the number of sensors grows [3]. However, these battery-powered IoT devices are usually expected to survive for years without battery change. In this case, message delivered to the fusion center may be constrained not to contain the identity information. For the same reason, the data transmitted to the fusion center is usually quantized to be in a finite alphabet. 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]–[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 (air/water quality) detection, rumor detection, and seismic wave detection.

### A. Contributions and Major Challenges

The first part of this paper focuses 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 is 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.

The second part of this paper focuses 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 Bayesian approach to address the challenge raised by the unknown trajectory of the anomaly, 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.

We also conduct extensive numerical results to demonstrate the performance of our proposed algorithms. The numerical results show that for the static setting, our GM-CuSum algorithm outperforms a heuristic Bayesian 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 arbitrarily 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]–[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 has been studied in the literature [4]–[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 are studied, where two practical algorithms, the unlabeled log-likelihood ratio test and the generalized likelihood ratio test are proposed. A more specific problem is studied in [7] where samples follow 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 more 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 system.

Existing studies of quickly detecting a dynamic change mostly focus on the labeled setting, e.g., [27], [32], [33]. Our problem is similar to the one in [33] but we focus on unlabeled samples. Our major technical 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]–[38]. Our work is different from the existing literature. Due to unlabeled samples and the dynamic nature of the anomaly, the unknown parameter, i.e., the identify and the label of the affected sensor, is time-varying. Therefore, the generalized likelihood approach which estimates the unknown parameters using their MLEs may not perform well. Moreover, unlike studies in [36]–[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 $\nu$, an anomaly emerges in the network, and changes the data-generating distributions of the sensors. 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.

After the anomaly emerges, one sensor of an unknown type is affected by the anomaly. 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 samples are generated by $p_{1,k-K}$.

Denote by $X^n[t] = \{X_1[t], \ldots, X_n[t]\}$ the $n$ arbitrarily permuted samples at time $t$ received by the fusion center. We assume that $X_1[t], \ldots, X_n[t]$ 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.

Let $\mathcal{K} = \{1, 2, \cdots, K\}$. 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 $\boldsymbol{S} \triangleq \{S[t]\}_{t=1}^{\infty}$ denote the trajectory of the anomaly. Here $\boldsymbol{S}$ is *unknown* to the decision maker. Even if the trajectory of the anomaly $\boldsymbol{S}$ 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 $\sigma_t^{S[t]} : \{1, \ldots, n\} \to \{1, \ldots, K, S[t] + K\}$. This function associates sample $X_i[t]$, $1 \leq i \leq n$, to group $j$ for some $j \in \{1, 2, \ldots, K, K + S[t]\}$, i.e., specifies the probability distribution of $X_i[t]$. Specifically, if $\sigma_t^{S[t]}(i) = j$, then

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

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, and changes with time.

Let $\Omega_{\boldsymbol{S}} = \{\sigma_1^{S[1]}, ..., \sigma_\infty^{S[\infty]}\}$ be the labels when the trajectory of the anomaly is $\boldsymbol{S}$, which is unknown. Let $\mathbb{P}_{\Omega_S}^{\boldsymbol{S},\nu}$ and $\mathbb{E}_{\Omega_S}^{\boldsymbol{S},\nu}$ denote the probability measure and the corresponding expectation when the change point is at $\nu$ and the samples received by the fusion center is permuted according to the label $\Omega_{\boldsymbol{S}}$ (see Appendix A for more details). We further let $\mathbb{P}_\Omega^\infty$ and $\mathbb{E}_\Omega^\infty$ denote the probability measure and the corresponding expectation when there is no change, i.e., $\nu = \infty$, where $\Omega = \Omega_{\boldsymbol{S}}$ 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 $\tau$:

$$\text{WADD}(\tau) = \sup_{\nu \geq 1} \sup_{\boldsymbol{S}} \sup_{\Omega_S} \operatorname*{esssup} \mathbb{E}_{\Omega_S}^{\boldsymbol{S},\nu}[(\tau - \nu)^+ | \mathbf{X}^n[1, \nu - 1]],$$
$$\text{WARL}(\tau) = \inf_\Omega \mathbb{E}_\Omega^\infty[\tau], \tag{2}$$

where $\mathbf{X}^n[t_1, t_2] = \{X^n[t_1], \cdots, X^n[t_2]\}$, for any $t_1 \leq t_2$. The goal is to design a stopping rule that minimizes the WADD subject to a constraint on the WARL:

$$\inf_{\tau : \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau), \tag{3}$$

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 $\gamma$, and $1/\gamma$ can be interpreted as the false alarm rate.

## 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 $j \in \{1, 2, \cdots, K, k + K\}$, there are $\binom{n}{n_1, \ldots, n_k - 1, \ldots, n_K, 1}$ possible $\sigma_t^k$ 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., $t < \nu$, the samples $X^n[t]$ follows the distribution

$$\mathbb{P}_{0,\sigma_t^0}(X^n[t]) \triangleq \prod_{i=1}^{n} p_{0,\sigma_t^0(i)}(X_i[t]), \tag{4}$$

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

$$\mathbb{P}_{\sigma_t^k}^k(X^n[t]) \triangleq \prod_{i:\sigma_t^k(i)\leq K} p_{0,\sigma_t^k(i)}(X_i[t]) \times \prod_{i:\sigma_t^k(i)>K} p_{1,\sigma_t^k(i)-K}(X_i[t]), \tag{5}$$

for some unknown $\sigma_t^k \in \mathcal{S}_{n,k}$. Let $\Omega_k = \{\sigma_1^0, \ldots, \sigma_{\nu-1}^0, \sigma_\nu^k, \ldots, \sigma_\infty^k\}$ be the labels over time, when the anomaly emerges at $\nu$ (similarly defined as $\Omega_S$). Let $\mathbb{P}_{\Omega_k}^{k,\nu}$ denote the probability measure when the change point is at $\nu$ and the samples are generated according to (4), (5) and $\Omega_k$. We further let $\mathbb{E}_{\Omega_k}^{k,\nu}$ denotes the corresponding expectation.

Then, the WADD for any stopping time $\tau$ can be written as

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

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:

$$\inf_{\tau:\text{WARL}(\tau)\geq\gamma} \text{WADD}(\tau). \tag{6}$$

### A. Universal Lower Bound on WADD

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

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

**Theorem 1.** *As* $\gamma \to \infty$,

$$\inf_{\tau:WARL(\tau)\geq\gamma} WADD(\tau) \geq \frac{\log \gamma}{I^*} + O(1). \tag{7}$$

The proof of Theorem 1 can be found in Appendix B. The main challenges 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 (6) is lower bounded by $\frac{\log \gamma}{I^*} + O(1)$ for any stopping rule that satisfy 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.

A first idea is to use the MLE to estimate the unknown label $\sigma_t^k$ and the unknown affected sensor $k$. In the static setting, $k$ does not change with time, however, $\sigma_t^k$ changes with time, thus a direct MLE for $\sigma_t^k$ at each time $t$ may not work well. Therefore, we take a mixture approach w.r.t. the unknown label, and then take a MLE approach w.r.t. the unknown affected sensor. Our algorithm is constructed as follows.

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

$$T_G = \inf\{t : W[t] \geq b\}, \tag{8}$$

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 $W_k[t] = \max_{1 \leq j \leq t} \sum_{i=j}^{t} \log \frac{\widetilde{\mathbb{P}}^k(X^n[i])}{\widetilde{\mathbb{P}}_0(X^n[i])}$. The test statistic $W[t]$ has the following recursion:

$$W[t+1] = \max_{k \in \mathcal{K}} \left\{ (W_k[t])^+ + \log \frac{\widetilde{\mathbb{P}}^k(X^n[t+1])}{\widetilde{\mathbb{P}}_0(X^n[t+1])} \right\}, \tag{9}$$

where $W_k[0] = 0$ for any $k \in \mathcal{K}$. We then take their maximum over $k$ as $W[t]$.

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

**Theorem 2.** *1) Let $b = \log(K\gamma)$ in (8). Then $WARL(T_G) \geq \gamma$; and 2) As $\gamma \to \infty$, $WADD(T_G) \leq \frac{\log \gamma}{I^*} + O(1)$.*

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 [40] and the optional sampling theorem [40]. The major challenge lies in that we consider the worst-case label. A key property we develop and use in the proof of the WARL lower bound 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 $\mathbb{E}_{0,\sigma^0} \left[ \log \frac{\widetilde{\mathbb{P}}^k(X^n)}{\widetilde{\mathbb{P}}_0(X^n)} \right]$ is invariant for different $\sigma^0$'s.

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

Based on Theorem 1 and Theorem 2, we then establish the second-order asymptotic optimality of $T_G$ in the following theorem.

**Theorem 3.** *$T_G$ is second-order asymptotically optimal for the problem in (6).*

*Proof.* By Theorem 1 and Theorem 2, we establish the second-order asymptotic optimality of $T_G$. $\square$

## 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: $\widetilde{\mathbb{P}}^{\boldsymbol{\beta}}(X^n) = \sum_{k=1}^{K} \beta_k \widetilde{\mathbb{P}}^k(X^n)$, where $\boldsymbol{\beta} = \{\beta_k\}_{k=1}^{K}$, $0 \leq \beta_k \leq 1$ and $\sum_{k=1}^{K} \beta_k = 1$. Denote by $I_{\boldsymbol{\beta}}$ the KL divergence between $\widetilde{\mathbb{P}}^{\boldsymbol{\beta}}$ and $\widetilde{\mathbb{P}}_0$. Let $\boldsymbol{\beta}^* = \arg\min_{\boldsymbol{\beta}} I_{\boldsymbol{\beta}}$.

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

**Theorem 4.** *As $\gamma \to \infty$, we have that*

$$\inf_{\tau:WARL(\tau)\geq\gamma} WADD(\tau) \geq \frac{\log \gamma}{I_{\boldsymbol{\beta}*}}(1 + o(1)). \tag{10}$$

The proof of Theorem 4 can be found in Appendix D.

It can be seen from Theorem 4 that the WADD for the problem in (3) is lower bounded by $\frac{\log \gamma}{I_{\boldsymbol{\beta}*}}(1+o(1))$ for large $\gamma$. This motivates us to apply the optimal weight $\boldsymbol{\beta}^*$ to design an algorithm that can achieves the WADD lower bound asymptotically. Moreover, we have that $I^* \geq I_{\boldsymbol{\beta}*}$ 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 Bayesian approach where the probability that the $k$-th group is affected by the anomaly is $\beta_k^*$. We then construct our weighted mixture CuSum algorithm as follows. Define the log of weighted mixture likelihood ratio using $\boldsymbol{\beta}^*$:

$$\ell_{\boldsymbol{\beta}*}(X^n) = \log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)}. \tag{11}$$

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

We then construct the following weighted mixture CuSum algorithm:

$$T_{\boldsymbol{\beta}*}(b) = \inf \left\{ t : \max_{1 \leq j \leq t+1} \sum_{i=j}^{t} \ell_{\boldsymbol{\beta}*}(X^n[i]) \geq b \right\}. \tag{12}$$

Let $\widehat{W}[t] = \max\limits_{1 \leq j \leq t+1} \sum_{i=j}^{t} \ell_{\boldsymbol{\beta}*}(X^n[i])$. The test statistic $\widehat{W}[t]$ has the following recursion: $\widehat{W}[t+1] = (\widehat{W}[t])^+ + \ell_{\boldsymbol{\beta}*}(X^n[t+1]), \widehat{W}[0] = 0$.

Note that different from the way that we handle the unknown and time-varying label $\sigma$, here, for the unknown type of the affected sensor, we take the mixture according to $\boldsymbol{\beta}^*$ instead of a uniform distribution over $\mathcal{K}$. As will be shown later both theoretically in Theorem 6 and numerically in Section V, taking a uniform mixture over $\mathcal{K}$ may not lead to the optimal performance.

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

**Lemma 1.** *For any $k \in \mathcal{K}$,*

$$\widetilde{\mathbb{E}}^k \left[ \log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)} \right] \geq I_{\boldsymbol{\beta}*}. \tag{13}$$

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 $\mathbb{E}_{\sigma^k}^k[\ell_{\boldsymbol{\beta}*}(X^n)]$ is invariant for different $\sigma^k$'s. Specifically, let $\mathbb{E}_{\sigma^k}^k$ denote the

expectation under $\mathbb{P}^k_{\sigma^k}$, where a sensor of type $k$ is affected, and the data received is labeled according to $\sigma^k$. For any $\pi$, let $\hat{\sigma}^k = \sigma^k \circ \pi$. Then $\mathbb{E}^k_{\sigma^k}[\ell_{\boldsymbol{\beta}^*}(\pi(X^n))] = \mathbb{E}^k_{\sigma^k \circ \pi}[\ell_{\boldsymbol{\beta}^*}(X^n)] = \mathbb{E}^k_{\hat{\sigma}^k}[\ell_{\boldsymbol{\beta}^*}(X^n)]$. For any $\hat{\sigma}^k \in \mathcal{S}_{n,k}$, a $\pi$ can always be found so that $\sigma^k \circ \pi = \hat{\sigma}^k$. Thus, for any $\sigma^k, \hat{\sigma}^k \in \mathcal{S}_{n,k}$, $\mathbb{E}^k_{\hat{\sigma}^k}[\ell_{\boldsymbol{\beta}^*}(X^n)] = \mathbb{E}^k_{\sigma^k}[\ell_{\boldsymbol{\beta}^*}(X^n)]$. Therefore, $\mathbb{E}^k_{\sigma^k}[\ell_{\boldsymbol{\beta}^*}(X^n)]$ is invariant for different $\sigma^k$'s. Then, under the pre-change distribution $\mathbb{P}_{0,\sigma^0_t}$, the expectation of the weighted mixture likelihood ratio $\mathbb{E}_{0,\sigma^0_t}[\ell_{\boldsymbol{\beta}^*}(X^n)]$ is invariant for different $\sigma^0_t$'s, we have that

$$
\begin{aligned}
\mathbb{E}_{0,\sigma^0_t}&\left[\log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)}\right] \\
&= \frac{1}{|\mathcal{S}_{n,0}|} \sum_{\sigma^0_t \in \mathcal{S}_{n,0}} \mathbb{E}_{0,\sigma^0_t}\left[\log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)}\right] \\
&= \widetilde{\mathbb{E}}_0\left[\log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)}\right] \\
&= -D(\widetilde{\mathbb{P}}_0 || \widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}) \leq 0.
\end{aligned}
\tag{14}
$$

Therefore, before the change time $\nu$, $\widehat{W}[t]$ has a negative drift. Similarly, from Lemma 1, after the change time $\nu$, under any group assignment $\Omega_{\boldsymbol{S}}$ and trajectory $\boldsymbol{S}$, $\widehat{W}[t]$ has a positive drift whose expectation is no less than $I_{\boldsymbol{\beta}^*}$, and evolves towards $\infty$.

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

**Theorem 5.** *1) For $T_{\boldsymbol{\beta}^*}$ defined in* (12)*, let $b = \log \gamma$, then $WARL(T_{\boldsymbol{\beta}^*}) \geq \gamma$.*

*2) As $\gamma \to \infty$, we have that*

$$
WADD(T_{\boldsymbol{\beta}^*}) \leq \frac{\log \gamma}{I_{\boldsymbol{\beta}^*}}(1 + o(1)).
\tag{15}
$$

The proof of Theorem 5 can be found in Appendix F. The proof of Theorem 5 is based on the Weak Law of Large Numbers for the weighted mixture likelihood ratio similarly to [37]. The major challenge lies in that here we are interested in the worst-case label and the worst-case anomaly trajectory. Note that 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 $T_{\boldsymbol{\beta}^*}$. To show the asymptotically optimal performance of $T_{\boldsymbol{\beta}^*}$, instead of finding the worst-case label and anomaly trajectory, we apply the symmetric property of $T_{\boldsymbol{\beta}^*}$ and Lemma 1 to show that the WADD and WARL of $T_{\boldsymbol{\beta}^*}$ are bounded under all possible labels and trajectories.

We then establish the first-order asymptotic optimality of $T_{\boldsymbol{\beta}^*}$ in the following theorem.

**Theorem 6.** *$T_{\boldsymbol{\beta}^*}$ is first-order asymptotically optimal for problem* (3)*.*

*Proof.* Combining Theorem 4 and Theorem 5, we establish the first-order asymptotic optimality of $T_{\boldsymbol{\beta}^*}$.
$\square$

If we apply $T_{\boldsymbol{\beta}^*}$ (designed for the dyanmic setting) to the static setting, the WADD of $T_{\boldsymbol{\beta}^*}$ can also be upper bounded by $\frac{\log \gamma}{I_{\boldsymbol{\beta}^*}}(1 + o(1))$. However, $T_{\boldsymbol{\beta}^*}$ 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 weighted mixture CuSum algorithm works better than the GM-CuSum.
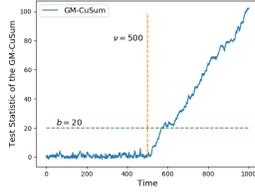
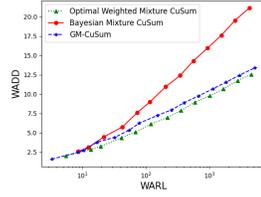Fig. 1. Evolution path of the GM-CuSum algorithm.



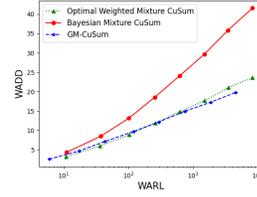Fig. 2. Comparison of the three algorithms in static setting: $n = 4, K = 2$.



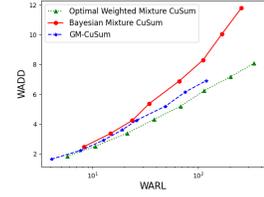Fig. 3. Comparison of the three algorithms in static setting: $n = 8, K = 2$.



Fig. 4. Comparison of the three algorithms in static setting: $n = 4, K = 4$.
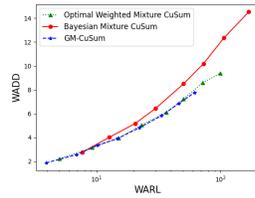


Fig. 5. Comparison of the three algorithms in static setting: $n = 8, K = 4$.
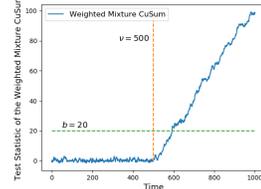


Fig. 6. Evolution path of the weighted mixture CuSum algorithm.
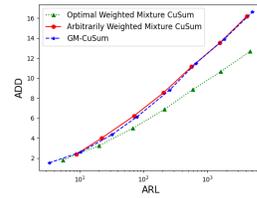


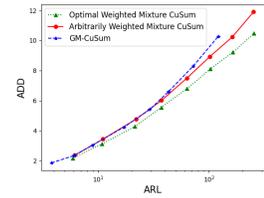Fig. 7. Comparison of the three algorithms in dynamic setting: $n = 4, K = 2$.



Fig. 8. Comparison of the three algorithms in dynamic setting: $n = 8, K = 2$.

## V. SIMULATION RESULTS

We first consider the static setting. We show an example evolution path of the GM-CuSum algorithm. We set $n = 2$ and $K = 2$. For type I sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.3)$ and $\mathcal{B}(10, 0.4)$, respectively, where $\mathcal{B}$ denotes binomial distribution. For type II sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.8)$ and $\mathcal{B}(10, 0.6)$, respectively. We set the change point to be 500 and $b = 20$. We plot one sample evolution path of the GM-CuSum algorithm when one sensor of type one is affected. It can be seen from Fig. 1 that before the change point, the test statistic fluctuates around zero, and after the change point, it starts to increase with a positive drift.

We then compare our GM-CuSum algorithm with a Bayesian mixture CuSum algorithm $T_B = \inf \Big\{ t :$ $\max_{1 \leq j \leq t} \sum_{i=j}^{t} \log \frac{\frac{1}{|\mathcal{K}|} \sum_{k \in \mathcal{K}} \widetilde{\mathbb{P}}^k(X^n[i])}{\widetilde{\mathbb{P}}_0(X^n[i])} \geq b \Big\}$ and the optimal weighted mixture CuSum algorithm. We plot the WADD as a function of the WARL under the worst-case static trajectory.

We consider four cases with different number of sensors and types. For the cases 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. We plot the figures for the cases where each type has two sensors and each type has four sensors in Fig. 2 and Fig. 3, respectively. For the cases 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 distributions are $\mathcal{B}(10, 0.4)$ and $\mathcal{B}(10, 0.7)$ respectively. We plot the figures for the cases where each type has one sensor and each type has two sensors in Fig. 4 and Fig. 5, respectively. We apply the Monte-Carlo approximation idea to obtain the optimal weight for our optimal weighted mixture CuSum algorithm. We repeat the experiment for 5000 times.
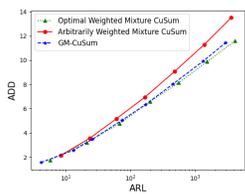
Fig. 9. Comparison of the three algorithms in dynamic setting: $n = 4, K = 4$.
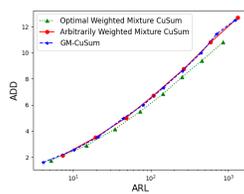
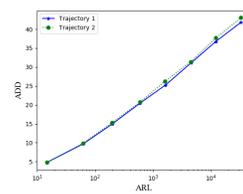Fig. 10. Comparison of the three algorithms in dynamic setting: $n = 8, K = 4$.

Fig. 11. Optimal Weighted CuSum Algorithm under Different Trajectories.

It can be seen from Fig. 2, Fig. 3, Fig. 4 and Fig. 5 that our GM-CuSum outperforms the Bayesian algorithm and the performance of the optimal weighted mixture CuSum algorithm are close to the GM-CuSum algorithm. The simulation results show that our optimal weighted mixture CuSum algorithm are also robust under the static setting. Moreover, the relationship between the WADD and log of the WARL is linear, which validates our theoretical results.

We then consider the dynamic anomaly. We use the same parameters of distributions as in the static setting. We first show an evolution path of the weighted mixture CuSum algorithm under a random trajectory $S$ in Fig. 6. To generate this random trajectory $S$, at each time step, let the probability that one sensor of type I is affected be $0.8$ and the probability that one sensor of type II is affected be $0.2$. Similar to the GM-CuSum, before the change point, the test statistic fluctuates around zero, and after the change point, it starts to increase with a positive drift.

We then compare our optimal weighted mixture CuSum algorithm with an arbitrary weighted mixture CuSum, i.e., replace $\beta^{*}$ in (12) with some arbitrarily $\beta$, e.g., $\beta = (\frac{1}{2}, \frac{1}{2})$ for the case with two types and the GM-CuSum. Here, we plot the average detection delay (ADD) and the average run length (ARL) for some randomly generated trajectories. It can be seen from Fig. 7, Fig. 8, Fig. 9 and Fig. 10 that our optimal weighted mixture CuSum algorithm outperforms the Bayesian weighted mixture CuSum algorithm and the GM-CuSum. The relationship between the WADD and log of the WARL is linear. It can also be observed that the GM-CuSum algorithm does not perform well under the dynamic setting.

We then compare the performance of our weighted mixture CuSum algorithm under two different trajectories. We choose $n = 2$ and $K = 2$. For type I sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.3)$ and $\mathcal{B}(10, 0.4)$, respectively. For type II sensors, the pre- and post-change distributions are $\mathcal{B}(10, 0.8)$ and $\mathcal{B}(10, 0.6)$, respectively. For trajectory 1, at each time, let the probability that one sensor of type one is affected be $0.8$ and the probability that one sensor of type two is affected be $0.2$. For trajectory 2, at each time, let the probability that one sensor of type one is affected be $0.2$ and the probability that one sensor of type two is affected be $0.8$. We plot the ADD as function of the ARL. It can be seen from Fig. 11 that for two different trajectories, our optimal weighted mixture CuSum algorithm have the same performance, which demonstrates the robustness of our optimal weighted mixture CuSum algorithm under different trajectories.

## VI. CONCLUSION

In this paper, we investigated the problem of quickest detection of an anomaly in networks with unlabeled samples. We first investigated the case with a static anomaly. We used the MLE to estimate the type of the affected sensor. A GM-CuSum algorithm was proposed. We showed that it is 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 a 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 $\binom{n}{n_1,\ldots,n_K}$ possible $\sigma_t^{S[t]}: \{1,\ldots,n\} \to \{1,\ldots,K\}$ satisfying $|\{i : \sigma_t^{S[t]}(i) = k\}| = n_k$, for any $k = 1,\ldots,K$. 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 $\binom{n}{n_1,\ldots,n_{S[t]}-1,\ldots,n_K,1}$ possible $\sigma_t^{S[t]}: \{1,\ldots,n\} \to \{1,\ldots,K,S[t]+K\}$ satisfying

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

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

$$
\mathbb{P}_{0,\sigma_t^0}(X^n[t]) = \prod_{i=1}^{n} p_{0,\sigma_t^0(i)}(X_i[t]), \tag{16}
$$

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

$$
\mathbb{P}_{\sigma_t^{S[t]}}^{S[t]}(X^n[t]) \triangleq \prod_{i:\sigma_t^{S[t]}(i) \leq K} p_{0,\sigma_t^{S[t]}(i)}(X_i[t]) \times \prod_{i:\sigma_t^{S[t]}(i) > K} p_{1,\sigma_t^{S[t]}(i)-K}(X_i[t]), \tag{17}
$$

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 $\widetilde{\mathbb{P}}_0$ and a post-change distribution $\widetilde{\mathbb{P}}^k$, respectively. Define the $\widetilde{\mathrm{WADD}}_k$ and $\widetilde{\mathrm{ARL}}$ for any stopping rule $\tau$ as follows:

$$
\widetilde{\mathrm{WADD}}_k(\tau) = \sup_{\nu \geq 1} \mathrm{esssup}\,\widetilde{\mathbb{E}}^{k,\nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1,\nu-1]],
$$

$$
\widetilde{\mathrm{ARL}}(\tau) = \widetilde{\mathbb{E}}^{\infty}[\tau], \tag{18}
$$

where $\widetilde{\mathbb{E}}^{k,\nu}$ denotes the expectation when the change is at $\nu$, the pre- and post-change distributions are $\widetilde{\mathbb{P}}_0$ and $\widetilde{\mathbb{P}}^k$, and $\widetilde{\mathbf{X}}^n[t]$ for $1 \leq t \leq \nu - 1$ are i.i.d. from $\widetilde{\mathbb{P}}_0$, $\widetilde{\mathbb{E}}^{\infty}$ denotes the expectation when there is no change and samples are generated according to $\widetilde{\mathbb{P}}_0$.

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 the $\mathrm{WADD}_k$ and WARL for any stopping rule $\tau$ as follows:

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

$$\text{WARL}(\tau) = \inf_{\Omega} \mathbb{E}_{\Omega}^{\infty}[\tau]. \tag{19}$$

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

$$\begin{aligned}
\text{WADD}(\tau) &= \sup_{k \in \mathcal{K}} \text{WADD}_k(\tau) \\
&\geq \sup_{\nu \geq 1} \sup_{\Omega_k} \text{esssup} \mathbb{E}_{\Omega_k}^{k,\nu}[(\tau - \nu)^+ | \mathbf{X}^n[1, \nu - 1]] \\
&\geq \sup_{\nu \geq 1} \text{esssup} \widetilde{\mathbb{E}}^{k,\nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]] \\
&= \widetilde{\text{WADD}}_k(\tau).
\end{aligned} \tag{20}$$

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

$$\begin{aligned}
\inf_{\tau : \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau) &\geq \inf_{\tau : \widetilde{\text{ARL}}(\tau) \geq \gamma} \widetilde{\text{WADD}}_k(\tau) \\
&\geq \frac{\log \gamma}{I_k} + O(1), \text{ as } \gamma \to \infty.
\end{aligned} \tag{21}$$

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

$$\inf_{\tau : \text{WARL}(\tau) \geq \gamma} \text{WADD}(\tau) \geq \frac{\log \gamma}{I^*} + O(1), \text{ as } \gamma \to \infty. \tag{22}$$

## APPENDIX C
## PROOF OF THEOREM 2

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

$$r_{m+1} = \inf \left\{ t > r_m : \sup_k \sum_{i=r_m+1}^{t} \log \frac{\widetilde{\mathbb{P}}^k(X_i^n)}{\widetilde{\mathbb{P}}_0(X_i^n)} \leq 0 \right\}. \tag{23}$$

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

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

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

$$= \int \frac{\widetilde{\mathbb{P}}^k(x^n)}{\widetilde{\mathbb{P}}_0(x^n)} \cdot \widetilde{\mathbb{P}}_0(x^n)\mathrm{d}x^n$$

$$= \int \widetilde{\mathbb{P}}^k(x^n)\mathrm{d}x^n = 1. \tag{24}$$

Therefore, for any $\Omega$ and $t > r_m$,

$$\mathbb{E}_\Omega^\infty\left[ \prod_{i=r_m+1}^{t+1} \frac{\widetilde{\mathbb{P}}^k(X_i^n)}{\widetilde{\mathbb{P}}_0(X_i^n)} \middle| \mathcal{F}_t \right]$$

$$= \mathbb{E}_\Omega^\infty\left[ \prod_{i=r_m+1}^{t} \frac{\widetilde{\mathbb{P}}^k(X_i^n)}{\widetilde{\mathbb{P}}_0(X_i^n)} \cdot \frac{\widetilde{\mathbb{P}}^k(X_{t+1}^n)}{\widetilde{\mathbb{P}}_0(X_{t+1}^n)} \middle| \mathcal{F}_t \right]$$

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

$$= \prod_{i=r_m+1}^{t} \frac{\widetilde{\mathbb{P}}^k(X_i^n)}{\widetilde{\mathbb{P}}_0(X_i^n)} \cdot \mathbb{E}_{0,\sigma^0}\left[ \frac{\widetilde{\mathbb{P}}^k(X_{t+1}^n)}{\widetilde{\mathbb{P}}_0(X_{t+1}^n)} \right]$$

$$= \prod_{i=r_m+1}^{t} \frac{\widetilde{\mathbb{P}}^k(X_i^n)}{\widetilde{\mathbb{P}}_0(X_i^n)}. \tag{25}$$

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

We then have that for any $\Omega$,

$$\mathbb{P}_\Omega^\infty\left\{ \sup_k \sum_{i=r_m+1}^{t} \log \frac{\widetilde{\mathbb{P}}^k(X_i^n)}{\widetilde{\mathbb{P}}_0(X_i^n)} \geq b \text{ for some } t > r_m \middle| \mathcal{F}_{r_m} \right\}$$

$$\leq \sum_{k=1}^{K} \mathbb{P}_\Omega^\infty\left\{ \sum_{i=r_m+1}^{t} \log \frac{\widetilde{\mathbb{P}}^k(X_i^n)}{\widetilde{\mathbb{P}}_0(X_i^n)} \geq b \text{ for some } t > r_m \middle| \mathcal{F}_{r_m} \right\}$$

$$= \sum_{k=1}^{K} \mathbb{P}_\Omega^\infty\left\{ \prod_{i=r_m+1}^{t} \frac{\widetilde{\mathbb{P}}^k(X_i^n)}{\widetilde{\mathbb{P}}_0(X_i^n)} \geq e^b \text{ for some } t > r_m \middle| \mathcal{F}_{r_m} \right\}$$

$$\leq K \frac{\mathbb{E}_{0,\sigma^0}\left[ \frac{\widetilde{\mathbb{P}}^k(X_{r_m+1}^n)}{\widetilde{\mathbb{P}}_0(X_{r_m+1}^n)} \right]}{e^b} = Ke^{-b}, \tag{26}$$

where the last inequality is due to Doob's submartingale inequality [40] and the optional sampling theorem [40].

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

$$\mathbb{P}_\Omega^\infty\left( M \geq m+1 \middle| \mathcal{F}_{r_m} \right)$$

$$\geq \mathbb{P}_\Omega^\infty\left\{ \sup_k \sum_{i=r_m+1}^{t} \log \frac{\widetilde{\mathbb{P}}^k(X_i^n)}{\widetilde{\mathbb{P}}_0(X_i^n)} < b \text{ for all } t > r_m \middle| \mathcal{F}_{r_m} \right\}$$

$$\geq 1 - Ke^{-b}. \tag{27}$$

We then have that for any $\Omega$,

$$\begin{aligned}
\mathbb{P}_\Omega^\infty(M > m) &= \mathbb{E}_\Omega^\infty\left[\mathbb{P}_\Omega^\infty\left(M \geq m + 1 | \mathcal{F}_{r_m}\right) \cdot \mathbb{1}_{\{M \geq m\}}\right] \\
&\geq (1 - Ke^{-b})\mathbb{P}_\Omega^\infty(M > m - 1) \\
&\geq (1 - Ke^{-b})^2\mathbb{P}_\Omega^\infty(M > m - 2) \\
&\geq (1 - Ke^{-b})^m\mathbb{P}_\Omega^\infty(M > 0) \\
&= (1 - Ke^{-b})^m.
\end{aligned} \tag{28}$$

It then follows that

$$\begin{aligned}
\mathrm{WARL}(T_G) = \inf_\Omega \mathbb{E}_\Omega^\infty[T_G] &\geq \inf_\Omega \mathbb{E}_\Omega^\infty[M] \\
&\geq \inf_\Omega \sum_{m=0}^\infty \mathbb{P}_\Omega^\infty(M > m) \\
&\geq \sum_{m=0}^\infty (1 - Ke^{-b})^m = \frac{e^b}{K}.
\end{aligned} \tag{29}$$

Let $b = \log K\gamma$, we have that $\mathrm{WARL}(T_G) \geq \gamma$. Let $T_k$ be the mixture CuSum algorithm for problem in (19):

$$T_k = \inf\left\{t : \max_{1 \leq j \leq t} \sum_{i=j}^t \log \frac{\widetilde{\mathbb{P}}^k(X^n[i])}{\widetilde{\mathbb{P}}_0(X^n[i])} \geq b\right\}. \tag{30}$$

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

$$\begin{aligned}
\mathrm{WADD}_k(T_G) = \sup_{\nu \geq 1} \sup_{\Omega_k} \mathrm{esssup}\mathbb{E}_{\Omega_k}^{k,\nu}[(T_G - \nu)^+ | \mathbf{X}^n[1, \nu - 1]] \\
\leq \sup_{\nu \geq 1} \sup_{\Omega_k} \mathrm{esssup}\mathbb{E}_{\Omega_k}^{k,\nu}[(T_k - \nu)^+ | \mathbf{X}^n[1, \nu - 1]] \\
\leq \frac{\log b}{I_k} + O(1),
\end{aligned} \tag{31}$$

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, choose $b = \log K\gamma$, we then have that

$$\begin{aligned}
\mathrm{WADD}(T_G) = \sup_{k \in \mathcal{K}} \mathrm{WADD}_k(T_G) &\leq \sup_{k \in \mathcal{K}} \mathrm{WADD}_k(T_k) \\
&= \sup_{k \in \mathcal{K}} \frac{\log K\gamma}{I_k} + O(1) \\
&= \frac{\log \gamma}{I^*} + \frac{\log K}{I^*} + O(1), \text{ as } \gamma \to \infty.
\end{aligned} \tag{32}$$

## APPENDIX D
## PROOF OF THEOREM 4

For any trajectory $\mathbf{S}$ and stopping time $\tau$, define the WADD and WARL

$$\begin{aligned}
\mathrm{WADD}_{\mathbf{S}}(\tau) &= \sup_{\nu \geq 1} \sup_{\Omega_{\mathbf{S}}} \mathrm{esssup}\mathbb{E}_{\Omega_{\mathbf{S}}}^{\mathbf{S},\nu}[(\tau - \nu)^+ | \mathbf{X}^n[1, \nu - 1]], \\
\mathrm{ARL}(\tau) &= \inf_\Omega \mathbb{E}_\Omega^\infty[\tau].
\end{aligned} \tag{33}$$

Consider QCD problem with a pre-change distribution $\widetilde{\mathbb{P}}_0 = \frac{1}{|\mathcal{S}_{n,0}|} \sum_{\sigma^0 \in \mathcal{S}_{n,0}} \mathbb{P}_{0,\sigma^0}$ and a post-change distribution $\widetilde{\mathbb{P}}^{S[t]} = \frac{1}{|\mathcal{S}_{n,S[t]}|} \sum_{\sigma^{S[t]} \in \mathcal{S}_{n,S[t]}} \mathbb{P}_{\sigma^{S[t]}}^{S[t]}$, respectively. For this pair of pre- and post-change distributions and any trajectory $\boldsymbol{S}$, define the $\widetilde{\mathrm{WADD}}_{\boldsymbol{S}}$ and $\widetilde{\mathrm{ARL}}_{\boldsymbol{S}}$ for any stopping rule $\tau$:

$$\widetilde{\mathrm{WADD}}_{\boldsymbol{S}}(\tau) = \sup_{\nu \geq 1} \mathrm{esssup} \widetilde{\mathbb{E}}^{\boldsymbol{S},\nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]],$$
$$\widetilde{\mathrm{ARL}}(\tau) = \widetilde{\mathbb{E}}^\infty[\tau]. \tag{34}$$

where $\widetilde{\mathbb{E}}^{\boldsymbol{S},\nu}$ denotes the expectation when change point is $\nu$, before the change point, the data follows distribution $\widetilde{\mathbb{P}}_0$ and after the change point, at time $t$, the data follows the distribution $\widetilde{\mathbb{P}}^{S[t]}$, and $\widetilde{\mathbf{X}}^n[1, \nu - 1]$ are i.i.d. from $\widetilde{\mathbb{P}}_0$; and $\widetilde{\mathbb{E}}^\infty$ denote the expectation when for any $t \geq 0$, the data follows distribution $\widetilde{\mathbb{P}}_0$, i.e., $\nu = \infty$.

Consider another QCD problem with pre-change distribution $\widetilde{\mathbb{P}}_0$ and post-change distribution $\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}$. Under this pair of pre- and post-change distributions, for any stopping time $\tau$, define worst-case average detection delay and average running length as follows:

$$\widetilde{\mathrm{WADD}}_{\boldsymbol{\beta}^*}(\tau) = \sup_{\nu \geq 1} \mathrm{esssup} \widetilde{\mathbb{E}}^{\boldsymbol{\beta}^*,\nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]],$$
$$\widetilde{\mathrm{ARL}}(\tau) = \widetilde{\mathbb{E}}^\infty[\tau]. \tag{35}$$

In QCD problems, ARL only depends on the pre-change distribution. Therefore, for any stopping time $\tau$, problems in (2) and (33) have the same ARL, problems in (34) and (35) have the same ARL. Let $\mathcal{C}_\gamma$ denotes the collection of all stopping times $\tau$ that satisfy $\mathrm{ARL}(\tau) \geq \gamma$ and $\widetilde{\mathcal{C}}_\gamma$ denotes the collection of all stopping times $\tau$ that satisfy $\widetilde{\mathrm{ARL}}(\tau) \geq \gamma$. Our goal is to prove that

$$\inf_{\tau \in \mathcal{C}_\gamma} \mathrm{WADD}(\tau) \geq \inf_{\tau \in \widetilde{\mathcal{C}}_\gamma} \widetilde{\mathrm{WADD}}_{\boldsymbol{\beta}^*}(\tau) \sim \frac{\log \gamma}{I_{\boldsymbol{\beta}^*}}(1 + o(1)). \tag{36}$$

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

$$\widehat{X}^n[t] \sim \begin{cases} \widetilde{\mathbb{P}}_0, & \text{if } t < \nu, \\ \mathbb{P}_{\sigma_t^{S[t]}}^{S[t]}, & \text{if } t \geq \nu. \end{cases} \tag{37}$$

For any stopping time $\tau$ and any $\boldsymbol{S}$, define the worst-case average detection delay for the model in (37) as follows:

$$\widehat{\mathrm{WADD}}_{\boldsymbol{S}}(\tau) = \sup_{\nu \geq 1} \sup_{\sigma_\nu^{S[\nu]}, \dots, \sigma_\infty^{S[\infty]}} \mathrm{esssup} \widehat{\mathbb{E}}^{\boldsymbol{S},\nu}_{\sigma_\nu^{S[\nu]}, \dots, \sigma_\infty^{S[\infty]}}[(\tau - \nu)^+ | \widehat{\mathbf{X}}^n[1, \nu - 1]], \tag{38}$$

where $\widehat{\mathbb{E}}^{\boldsymbol{S},\nu}_{\sigma_\nu^{S[\nu]}, \dots, \sigma_\infty^{S[\infty]}}$ denotes the expectation when the data is distributed according to (37).

Let $\widehat{\mathrm{WADD}}(\tau) = \sup_{\boldsymbol{S}} \widehat{\mathrm{WADD}}_{\boldsymbol{S}}(\tau)$. To prove (36), we will first show that for any $\boldsymbol{S}$, $\mathrm{WADD}_{\boldsymbol{S}}(\tau) = \widetilde{\mathrm{WADD}}_{\boldsymbol{S}}(\tau)$, and then show that $\widehat{\mathrm{WADD}}_{\boldsymbol{S}}(\tau) \geq \widetilde{\mathrm{WADD}}_{\boldsymbol{S}}(\tau)$. We will then complete our proof by showing that for any $\tau$ and $\boldsymbol{\beta}$, $\widehat{\mathrm{WADD}}(\tau) \geq \widetilde{\mathrm{WADD}}_{\boldsymbol{\beta}}(\tau)$.

**Step 1.** Denote by $\mathcal{M}$ the collection of all $\{\sigma_1^0, ..., \sigma_{\nu-1}^0\}$, and $\mu$ is an element in $\mathcal{M}$. When the trajectory is $\boldsymbol{S}$, denote by $\mathcal{N}_{\boldsymbol{S}}$ the collection of all $\{\sigma_\nu^{S[\nu]}, ..., \sigma_\infty^{S[\infty]}\}$, and $\omega$ is an element in $\mathcal{N}_{\boldsymbol{S}}$. Then, the WADD$_{\boldsymbol{S}}$ can be written as

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

where $\mathbb{E}_\omega^{\boldsymbol{S},\nu}$ denotes the expectation when change point is $\nu$, the trajectory is $\boldsymbol{S}$, and after the change point, the data follows distribution $\prod_{t=\nu}^\infty \mathbb{P}_{\sigma_t^{S[t]}}^{S[t]}$. We note that $\widehat{X}^n[t]$ and $X^n[t]$, for $t \geq \nu$, have the same distribution $\mathbb{P}_{\sigma_t^{S[t]}}^{S[t]}$. Therefore, the difference between WADD$_{\boldsymbol{S}}$ and $\widehat{\text{WADD}}_{\boldsymbol{S}}$ 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 $\mathbb{E}_\omega^{\boldsymbol{S},\nu}[(\tau-\nu)^+|\mathbf{X}^n[1,\nu-1]]$. Then, WADD$_{\boldsymbol{S}}$ and $\widehat{\text{WADD}}_{\boldsymbol{S}}$ can be written as

$$\text{WADD}_{\boldsymbol{S}}(\tau) = \sup_{\nu \geq 1} \sup_{\omega \in \mathcal{N}_{\boldsymbol{S}}} \sup_{\mu \in \mathcal{M}} \text{esssup} f_\omega(\mathbf{X}^n[1,\nu-1]),$$
$$\widehat{\text{WADD}}_{\boldsymbol{S}}(\tau) = \sup_{\nu \geq 1} \sup_{\omega \in \mathcal{N}_{\boldsymbol{S}}} \text{esssup} f_\omega(\widehat{\mathbf{X}}^n[1,\nu-1]). \tag{39}$$

It then suffices to show that for any $\omega \in \mathcal{N}_S$, $\sup_{\mu \in \mathcal{M}} \text{esssup} f_\omega(\mathbf{X}^n[1,\nu-1]) = \text{ess sup} f_\omega(\widehat{\mathbf{X}}^n[1,\nu-1])$.

For any $\omega \in \mathcal{N}_{\boldsymbol{S}}$ and $\mu \in \mathcal{M}$, let

$$b_{\omega,\mu} = \text{esssup} f_\omega(\mathbf{X}^n[1,\nu-1])$$
$$= \inf\{b : \mathbb{P}_\mu(f_\omega(\mathbf{X}^n[1,\nu-1]) > b) = 0\}, \tag{40}$$

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

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

$$b_\omega^* = \inf\left\{b : \int_{\mathbf{x}^n[1,\nu-1]} \mathbb{1}_{\{f_\omega(\mathbf{x}^n[1,\nu-1])>b\}} \times \mathrm{d} \prod_{t=1}^{\nu-1} \widetilde{\mathbb{P}}_0(x^n(t)) = 0\right\}$$
$$= \inf\left\{b : \int_{\mathbf{x}^n[1,\nu-1]} \mathbb{1}_{\{f_\omega(\mathbf{x}^n[1,\nu-1])>b\}} \times \mathrm{d} \prod_{t=1}^{\nu-1} \frac{1}{|\mathcal{S}_{n,0}|} \sum_{\sigma_t^0 \in \mathcal{S}_{n,0}} \mathbb{P}_{0,\sigma_t^0}(x^n(t)) = 0\right\}$$
$$= \inf\left\{b : \int_{\mathbf{x}^n[1,\nu-1]} \mathbb{1}_{\{f_\omega(\mathbf{x}^n[1,\nu-1])>b\}} \times \mathrm{d} \frac{1}{|\mathcal{M}|} \sum_{\mu \in \mathcal{M}} \mathbb{P}_\mu(\mathbf{x}^n[1,\nu-1]) = 0\right\}$$
$$= \inf\left\{b : \frac{1}{|\mathcal{M}|} \sum_{\mu \in \mathcal{M}} \mathbb{P}_\mu(f_\omega(\mathbf{X}^n[1,\nu-1]) > b) = 0\right\}.$$

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

$$\sup_{\mu \in \mathcal{M}} b_{\omega,\mu} \leq b_\omega^*. \tag{41}$$

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

$$b_\omega^* \leq \sup_{\mu \in \mathcal{M}} b_{\omega,\mu}. \tag{42}$$

Combining (41) and (42), we have that $\sup_{\mu \in \mathcal{M}} b_{\omega,\mu} = b_\omega^*$, and thus $\sup_{\mu \in \mathcal{M}} \mathrm{esssup} f_\omega(\mathbf{X}^n[1,\nu-1]) = \mathrm{esssup} f_\omega(\widehat{\mathbf{X}}^n[1,\nu-1])$. This implies that for any $\tau$,

$$\mathrm{WADD}_{\boldsymbol{S}}(\tau) = \widehat{\mathrm{WADD}}_{\boldsymbol{S}}(\tau). \tag{43}$$

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

$$\widetilde{\mathbb{P}}^\nu\Big(\sup_{\omega \in \mathcal{N}_{\boldsymbol{S}}} f_\omega(\widehat{\mathbf{X}}^n[1,\nu-1]) > \hat{b}\Big)$$
$$\leq \widetilde{\mathbb{P}}^\nu\Big(\cup_{\omega \in \mathcal{N}_{\boldsymbol{S}}}\big\{f_\omega(\widehat{\mathbf{X}}^n[1,\nu-1]) > \hat{b}\big\}\Big) = 0. \tag{44}$$

Therefore,

$$\hat{b} = \sup_{\omega \in \mathcal{N}_{\boldsymbol{S}}} \mathrm{esssup} f_\omega(\widehat{\mathbf{X}}^n[1,\nu-1])$$
$$\geq \mathrm{esssup} \sup_{\omega \in \mathcal{N}_{\boldsymbol{S}}} f_\omega(\widehat{\mathbf{X}}^n[1,\nu-1]). \tag{45}$$

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

$$\sup_{\substack{\{\sigma_\nu^{S[\nu]},\cdots,\sigma_T^{S[T]}\} \\ \in \mathcal{S}_{n,S[\nu]} \times,\cdots,\times \mathcal{S}_{n,S[T]}}} \sum_{t=\nu+1}^{T} (t-\nu)\mathbb{P}_{\sigma_\nu^{S[\nu]},\cdots,\sigma_T^{S[T]}}^{\boldsymbol{S},\nu}(\tau = t|\widehat{\mathbf{X}}^n[1,\nu-1])$$

$$\geq \sum_{t=\nu+1}^{T} (t-\nu)\frac{1}{|\mathcal{S}_{n,S[\nu]}| \times \cdots \times |\mathcal{S}_{n,S[T]}|} \sum_{\substack{\{\sigma_\nu^{S[\nu]},\cdots,\sigma_T^{S[T]}\} \\ \in \mathcal{S}_{n,S[\nu]} \times,\cdots,\times \mathcal{S}_{n,S[T]}}} \mathbb{P}_{\sigma_\nu^{S[\nu]},\cdots,\sigma_T^{S[T]}}^{\boldsymbol{S},\nu}(\tau = t|\widehat{\mathbf{X}}^n[1,\nu-1])$$

$$= \sum_{t=\nu+1}^{T} (t-\nu)\widetilde{\mathbb{P}}^{\boldsymbol{S},\nu}(\tau = t|\widetilde{\mathbf{X}}^n[1,\nu-1]), \tag{46}$$

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

$$f_\omega(\widehat{\mathbf{X}}^n[1,\nu-1]) \geq \widetilde{\mathbb{E}}^{\boldsymbol{S},\nu}[(\tau-\nu)^+|\widetilde{\mathbf{X}}^n[1,\nu-1]], \tag{47}$$

From (45) and (47), we have that

$$
\begin{aligned}
\widehat{\mathrm{WADD}}_{\boldsymbol{S}}(\tau) &= \sup_{\omega \in \mathcal{N}_{\boldsymbol{S}}} \mathrm{esssup} f_{\omega}(\widehat{\mathbf{X}}^n[1, \nu - 1]) \\
&\geq \mathrm{esssup} \widetilde{\mathbb{E}}^{\boldsymbol{S}, \nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]] \\
&= \widetilde{\mathrm{WADD}}_{\boldsymbol{S}}(\tau).
\end{aligned}
\tag{48}
$$

Combining (43) and (48), it follows that

$$
\mathrm{WADD}_{\boldsymbol{S}}(\tau) = \widehat{\mathrm{WADD}}_{\boldsymbol{S}}(\tau) \geq \widetilde{\mathrm{WADD}}_{\boldsymbol{S}}(\tau).
\tag{49}
$$

This holds for any trajectory $\boldsymbol{S}$. It then follows that

$$
\begin{aligned}
\mathrm{WADD}(\tau) &= \sup_{\nu \geq 0} \sup_{\Omega_{\boldsymbol{S}}} \sup_{\boldsymbol{S}} \mathrm{esssup} \mathbb{E}_{\Omega_{\boldsymbol{S}}}^{\boldsymbol{S}, \nu}[(\tau - \nu)^+ | \mathbf{X}^n[1, \nu - 1]] \\
&\geq \sup_{\nu \geq 0} \sup_{\boldsymbol{S}} \mathrm{esssup} \widetilde{\mathbb{E}}^{\boldsymbol{S}, \nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]] \\
&= \widetilde{\mathrm{WADD}}(\tau).
\end{aligned}
\tag{50}
$$

**Step 3.** The last step is to show that for any $\tau$ and any $\boldsymbol{\beta}$, $\widetilde{\mathrm{WADD}}(\tau) \geq \widetilde{\mathrm{WADD}}_{\boldsymbol{\beta}}(\tau)$. Firstly, we will show that

$$
\begin{aligned}
&\sup_{\boldsymbol{S}} \mathrm{esssup} \widetilde{\mathbb{E}}^{\boldsymbol{S}, \nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]] \\
&\geq \mathrm{esssup} \sup_{\boldsymbol{S}} \widetilde{\mathbb{E}}^{\boldsymbol{S}, \nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]].
\end{aligned}
\tag{51}
$$

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

$$
\widetilde{\mathbb{P}}^{\nu} \left( \widetilde{\mathbb{E}}^{\boldsymbol{S}, \nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]] > c \right) = 0.
\tag{52}
$$

Since $\Lambda_{\boldsymbol{S}}$ is countable, it then follows that

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

Therefore,

$$
\begin{aligned}
c &= \sup_{\boldsymbol{S}} \mathrm{esssup} \widetilde{\mathbb{E}}^{\boldsymbol{S}, \nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]] \\
&\geq \mathrm{esssup} \sup_{\boldsymbol{S}} \widetilde{\mathbb{E}}^{\boldsymbol{S}, \nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]].
\end{aligned}
\tag{53}
$$

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

$$
\begin{aligned}
&\sup_{\boldsymbol{S}} \sum_{t=\nu+1}^{T} (t - \nu) \widetilde{\mathbb{P}}^{S[\nu], \dots, S[T]}(\tau = t | \widetilde{\mathbf{X}}^n[1, \nu - 1]) \\
&\geq \sum_{t=\nu+1}^{T} (t - \nu) \sum_{\{S[\nu], \dots, S[T]\} \in \Lambda_{\boldsymbol{S}}^{\otimes(T-\nu+1)}} \beta_{S[\nu]} \times \cdots \times \beta_{S[T]} \widetilde{\mathbb{P}}^{S[\nu], \dots, S[T]}(\tau = t | \widetilde{\mathbf{X}}^n[1, \nu - 1]) \\
&= \sum_{t=\nu+1}^{T} (t - \nu) \widetilde{\mathbb{P}}^{\boldsymbol{\beta}, \nu}(\tau = t | \widetilde{\mathbf{X}}^n[1, \nu - 1]),
\end{aligned}
\tag{54}
$$

where $\widetilde{\mathbb{P}}^{S[\nu],...,S[T]}$ denotes the probability measure when the trajectory is $\boldsymbol{S}$, the observations from time $\nu$ to time $T$ are generated according to $\widetilde{\mathbb{P}}^{S[\nu]},...,\widetilde{\mathbb{P}}^{S[T]}$. As $T \to \infty$, we have

$$\sup_{\boldsymbol{S}} \widetilde{\mathbb{E}}^{\boldsymbol{S},\nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]]$$
$$\geq \widetilde{\mathbb{E}}^{\boldsymbol{\beta},\nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]]. \tag{55}$$

From (53) and (55), we have that

$$\widetilde{\text{WADD}}(\tau) = \sup_{\boldsymbol{S}} \text{esssup} \widetilde{\mathbb{E}}^{\boldsymbol{S},\nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]]$$
$$\geq \text{esssup} \widetilde{\mathbb{E}}^{\boldsymbol{\beta},\nu}[(\tau - \nu)^+ | \widetilde{\mathbf{X}}^n[1, \nu - 1]]$$
$$= \widetilde{\text{WADD}}_{\boldsymbol{\beta}}(\tau). \tag{56}$$

Combining (50) and (56), we have that for any $\tau$ and $\boldsymbol{\beta}$

$$\text{WADD}(\tau) \geq \widetilde{\text{WADD}}(\tau) \geq \widetilde{\text{WADD}}_{\boldsymbol{\beta}}(\tau). \tag{57}$$

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

$$\inf_{\substack{\{\sigma_1^0,...,\sigma_T^0\} \\ \in \mathcal{S}_{n,0}^{\otimes T}}} \sum_{t=1}^{T} t \mathbb{P}_{\sigma_1^0,...,\sigma_T^0}^{\infty}(\tau = t)$$

$$\leq \sum_{t=1}^{T} t \frac{1}{|\mathcal{S}_{n,0}|^T} \sum_{\substack{\{\sigma_1^0,...,\sigma_T^0\} \\ \in \mathcal{S}_{n,0}^{\otimes T}}} \mathbb{P}_{\sigma_1^0,...,\sigma_T^0}^{\infty}(\tau = t)$$

$$= \sum_{t=1}^{T} t \widetilde{\mathbb{P}}^{\infty}(\tau = t). \tag{58}$$

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

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

Since (57) holds for any $\boldsymbol{\beta}$, it holds for $\boldsymbol{\beta}^*$. Problem (35) is a classical QCD problem. From the asymptotic lower bound analysis in [37], we have that for large $\gamma$,

$$\inf_{\tau \in \mathcal{C}_\gamma} \text{WADD}(\tau) \geq \inf_{\tau \in \widetilde{\mathcal{C}}_\gamma} \widetilde{\text{WADD}}_{\boldsymbol{\beta}^*}(\tau) \sim \frac{\log \gamma}{I_{\boldsymbol{\beta}^*}}(1 + o(1)). \tag{59}$$

## APPENDIX E
## PROOF OF LEMMA 1

The minimization of $I_{\boldsymbol{\beta}}$ is to solve the following problem:

$$\inf_{\boldsymbol{\beta}} \quad I_{\boldsymbol{\beta}}$$
$$\text{s.t.} \quad -\beta_k \leq 0, \text{ for } k \in [1, K] \tag{60}$$
$$\sum_{k=1}^{K} \beta_k - 1 = 0.$$

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

$$L(\boldsymbol{\beta}, \eta, \boldsymbol{\mu}) = I_{\boldsymbol{\beta}} + \eta\Big(\sum_{k=1}^{K} \beta_k - 1\Big) - \sum_{k=1}^{K} \mu_k \beta_k. \tag{61}$$

The minimizer $\boldsymbol{\beta}^*$ 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

$$\frac{\partial L}{\partial \beta_k}|_{\boldsymbol{\beta}^*} = \widetilde{\mathbb{E}}^k\Big[\log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)}\Big] + 1 + \eta - \mu_k = 0, \tag{62}$$

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

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

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

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

$$
\begin{aligned}
I_{\boldsymbol{\beta}^*} &= \sum_{k=1}^{K} \beta_k^* \widetilde{\mathbb{E}}^k\Big[\log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)}\Big] \\
&= \sum_{k \in K_1} \beta_k^* \widetilde{\mathbb{E}}^k\Big[\log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)}\Big] \\
&= \widetilde{\mathbb{E}}^k\Big[\log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)}\Big], \text{ for all } k \in \mathcal{K}_1.
\end{aligned} \tag{63}
$$

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

$$\widetilde{\mathbb{E}}^k\Big[\log \frac{\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}(X^n)}{\widetilde{\mathbb{P}}_0(X^n)}\Big] \geq I_{\boldsymbol{\beta}^*}. \tag{64}$$

## APPENDIX F
## PROOF OF THEOREM 5

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

$$
\begin{aligned}
\text{WADD}_{\boldsymbol{S}}(T_{\boldsymbol{\beta}^*}) &= \sup_{\nu \geq 0} \sup_{\Omega_{\boldsymbol{S}}} \text{esssup} \mathbb{E}_{\Omega_{\boldsymbol{S}}}^{\boldsymbol{S}, \nu}\big[(T_{\boldsymbol{\beta}^*} - \nu)^+ \mid \mathbf{X}^n[1, \nu-1]\big] \\
&= \sup_{\Omega_{\boldsymbol{S}}} \mathbb{E}_{\Omega_{\boldsymbol{S}}}^{\boldsymbol{S}, 0}[T_{\boldsymbol{\beta}^*}].
\end{aligned} \tag{65}
$$

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

$$\sup_{\substack{\{\sigma_1^{S[1]}, \cdots, \sigma_T^{S[T]}\} \\ \in \mathcal{S}_{n, S[1]} \times, \cdots, \times \mathcal{S}_{n, S[T]}}} \sum_{t=1}^{T} t \mathbb{P}_{\sigma_1^{S[1]}, \cdots, \sigma_T^{S[T]}}^{\boldsymbol{S}, 0}(T_{\boldsymbol{\beta}^*} = t)$$

$$= \sum_{t=1}^{T} t \frac{1}{\mid \mathcal{S}_{n,S[1]} \mid \times \cdots \times \mid \mathcal{S}_{n,S[T]} \mid} \sum_{\substack{\{\sigma_1^{S[1]}, \cdots, \sigma_T^{S[T]}\} \\ \in \mathcal{S}_{n,S[1]} \times, \cdots, \times \mathcal{S}_{n,S[T]}}} \mathbb{P}_{\sigma_1^{S[1]}, \cdots, \sigma_T^{S[T]}}^{\boldsymbol{S},0}(T_{\boldsymbol{\beta}^*} = t)$$

$$= \sum_{t=1}^{T} t \widetilde{\mathbb{P}}^{\boldsymbol{S},0}(T_{\boldsymbol{\beta}^*} = t). \tag{66}$$

As $T \to \infty$, we have that

$$\sup_{\Omega_{\mathcal{S}}} \mathbb{E}_{\Omega_{\mathcal{S}}}^{\boldsymbol{S},0}[T_{\boldsymbol{\beta}^*}] = \widetilde{\mathbb{E}}^{\boldsymbol{S},0}[T_{\boldsymbol{\beta}^*}] = \widetilde{\mathrm{WADD}}_{\boldsymbol{S}}(T_{\boldsymbol{\beta}^*}). \tag{67}$$

For any $\boldsymbol{S}$, we have $\mathrm{WADD}_{\boldsymbol{S}}(T_{\boldsymbol{\beta}^*}) = \widetilde{\mathrm{WADD}}_{\boldsymbol{S}}(T_{\boldsymbol{\beta}^*})$. Therefore, $\mathrm{WADD}(T_{\boldsymbol{\beta}^*}) = \widetilde{\mathrm{WADD}}(T_{\boldsymbol{\beta}^*})$ by taking sup over $\boldsymbol{S}$ on both sides. It then follows that

$$\mathrm{WADD}(T_{\boldsymbol{\beta}^*}) = \widetilde{\mathrm{WADD}}(T_{\boldsymbol{\beta}^*}) = \sup_{\boldsymbol{S}} \widetilde{\mathbb{E}}^{\boldsymbol{S},0}[T_{\boldsymbol{\beta}^*}]. \tag{68}$$

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

$$\widetilde{\mathbb{E}}^{\boldsymbol{S},0}\left[\frac{T_{\boldsymbol{\beta}^*}}{n_b}\right] = \int_0^{\infty} \widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\frac{T_{\boldsymbol{\beta}^*}}{n_b} > x\right) \mathrm{d}x$$

$$\leq \sum_{t=1}^{\infty} \widetilde{\mathbb{P}}^{\boldsymbol{S},0}(T_{\boldsymbol{\beta}^*} > tn_b) + 1. \tag{69}$$

For any $\boldsymbol{S}$, we have that

$$\widetilde{\mathbb{P}}^{\boldsymbol{S},0}(T_{\boldsymbol{\beta}^*} > tn_b)$$

$$= \widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\max_{1 \leq k \leq tn_b} \max_{1 \leq i \leq k} \sum_{j=i}^{k} \ell_{\boldsymbol{\beta}^*}(X_j^n) < b\right)$$

$$\leq \widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\max_{1 \leq i \leq mn_b} \sum_{j=i}^{mn_b} \ell_{\boldsymbol{\beta}^*}(X_j^n) < b, \forall m \in [t]\right)$$

$$\leq \widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\boldsymbol{\beta}^*}(X_j^n) < b, \forall m \in [t]\right)$$

$$= \widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b} < I_{\boldsymbol{\beta}^*} - \epsilon, \forall m \in [t]\right)$$

$$= \prod_{m=1}^{t} \widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b} < I_{\boldsymbol{\beta}^*} - \epsilon\right). \tag{70}$$

It then follows that

$$\sup_{\boldsymbol{S}} \sum_{t=1}^{\infty} \widetilde{\mathbb{P}}^{\boldsymbol{S},0}(T_{\boldsymbol{\beta}^*} > tn_b)$$

$$\leq \sup_{\boldsymbol{S}} \sum_{t=1}^{\infty} \prod_{m=1}^{t} \widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b} \ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b} < I_{\boldsymbol{\beta}^*} - \epsilon\right).$$

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

Let $I_{\boldsymbol{S}_m} = \widetilde{\mathbb{E}}^{\boldsymbol{S},0}\left[\dfrac{\sum_{j=(m-1)n_b+1}^{mn_b}\ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b}\right]$. From (63) and (64), we have that

$$
\begin{aligned}
I_{\boldsymbol{S}_m} &= \widetilde{\mathbb{E}}^{\boldsymbol{S},0}\left[\frac{\sum_{j=(m-1)n_b+1}^{mn_b}\ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b}\right]\\
&= \sum_{j=(m-1)n_b+1}^{mn_b}\widetilde{\mathbb{E}}^{S[j]}\left[\frac{\ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b}\right]\\
&= \frac{1}{n_b}\sum_{j=(m-1)n_b+1}^{mn_b}\widetilde{\mathbb{E}}^{S[j]}\left[\ell_{\boldsymbol{\beta}^*}(X_j^n)\right] \geq I_{\boldsymbol{\beta}^*}.
\end{aligned}
\tag{71}
$$

It then follows that for any $\boldsymbol{S}$ and $m$

$$
\begin{aligned}
&\widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b}\ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b} < I_{\boldsymbol{\beta}^*} - \epsilon\right)\\
&\leq \widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b}\ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b} < I_{\boldsymbol{S}_m} - \epsilon\right)\\
&\leq \widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\left|\frac{\sum_{j=(m-1)n_b+1}^{mn_b}\ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b} - I_{\boldsymbol{S}_m}\right| > \epsilon\right).
\end{aligned}
\tag{72}
$$

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

$$
\begin{aligned}
&\widetilde{\mathbb{P}}^{\boldsymbol{S},0}\left(\left|\frac{\sum_{j=(m-1)n_b+1}^{mn_b}\ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b} - I_{\boldsymbol{S}_m}\right| > \epsilon\right)\\
&\leq \mathrm{Var}_{\widetilde{\mathbb{P}}^{\boldsymbol{S}}}\left(\frac{\sum_{j=(m-1)n_b+1}^{mn_b}\ell_{\boldsymbol{\beta}^*}(X_j^n)}{n_b}\right)\frac{1}{\epsilon^2}\\
&= \frac{1}{\epsilon^2 n_b^2}\sum_{j=(m-1)n_b+1}^{mn_b}\mathrm{Var}_{\widetilde{\mathbb{P}}^{S[j]}}(\ell_{\boldsymbol{\beta}^*}(X_j^n))\\
&\leq \frac{\sum_{j=(m-1)n_b+1}^{mn_b}\sigma^2}{n_b^2\epsilon^2} = \frac{\sigma^2}{n_b\epsilon^2}.
\end{aligned}
\tag{73}
$$

Let $\delta = \frac{\sigma^2}{n_b \epsilon^2}$. From (69) and (73), we have that

$$\sup_{\boldsymbol{S}} \widetilde{\mathbb{E}}^{\boldsymbol{S},0}\left[\frac{T_{\boldsymbol{\beta}^*}}{n_b}\right] \leq 1 + \sup_{\boldsymbol{S}} \sum_{t=1}^{\infty} \widetilde{\mathbb{P}}^{\boldsymbol{S},0}(T_{\boldsymbol{\beta}^*} > tn_b)$$

$$\leq 1 + \sum_{t=1}^{\infty} (\frac{\sigma^2}{n_b \epsilon^2})^t$$

$$= 1 + \sum_{t=1}^{\infty} \delta^t = \frac{1}{1-\delta}. \tag{74}$$

Therefore, we have

$$\sup_{\boldsymbol{S}} \widetilde{\mathbb{E}}^{\boldsymbol{S},0}\left[T_{\boldsymbol{\beta}^*}\right] \leq \frac{b}{(I_{\boldsymbol{\beta}^*} - \epsilon)(1-\delta)}. \tag{75}$$

(75) holds for all $\epsilon$. It then follows that as $b \to \infty$,

$$\text{WADD}(T_{\boldsymbol{\beta}^*}) = \sup_{\boldsymbol{S}} \widetilde{\mathbb{E}}^{\boldsymbol{S},0}\left[T_{\boldsymbol{\beta}^*}\right] \leq \frac{b}{I_{\boldsymbol{\beta}^*}}(1 + o(1)). \tag{76}$$

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

$$\inf_{\substack{\{\sigma_1^0, \ldots, \sigma_T^0\} \\ \in \mathcal{S}_{n,0}\otimes^T}} \sum_{t=1}^{T} t\mathbb{P}_{\sigma_1^0, \ldots, \sigma_T^0}^{\infty}(T_{\boldsymbol{\beta}^*} = t)$$

$$= \sum_{t=1}^{T} t\frac{1}{\mid \mathcal{S}_{n,0}\mid^T} \sum_{\substack{\{\sigma_1^0, \ldots, \sigma_T^0\} \\ \in \mathcal{S}_{n,0}\otimes^T}} \mathbb{P}_{\sigma_1^0, \ldots, \sigma_T^0}^{\infty}(T_{\boldsymbol{\beta}^*} = t)$$

$$= \sum_{t=1}^{T} t\widetilde{\mathbb{P}}^{\infty}(T_{\boldsymbol{\beta}^*} = t). \tag{77}$$

As $T \to \infty$, we have that $\text{WARL}(T_{\boldsymbol{\beta}^*}) = \widetilde{\text{ARL}}(T_{\boldsymbol{\beta}^*})$. $T_{\boldsymbol{\beta}^*}$ is the CuSum algorithm for a simple QCD problem with pre-change distribution $\widetilde{\mathbb{P}}_0$ and post-change distribution $\widetilde{\mathbb{P}}^{\boldsymbol{\beta}^*}$. From the optimal property of CuSum algorithm in [39] and [41], we have that when $b = \log \gamma$, $\text{WARL}(T_{\boldsymbol{\beta}^*}) = \widetilde{\text{ARL}}(T_{\boldsymbol{\beta}^*}) \geq \gamma$.

## REFERENCES

[1] Z. Sun and S. Zou, "Quickest dynamic anomaly detection in anonymous heterogeneous sensor networks," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, pp. 106–111, 2021.

[2] T. E. Humphreys, B. M. Ledvina, M. L. Psiaki, B. W. O'Hanlon, P. M. Kintner, *et al.*, "Assessing the spoofing threat: Development of a portable gps civilian spoofer," in *Proceedings of the 21st International Technical Meeting of the Satellite Division of The Institute of Navigation (ION GNSS 2008)*, pp. 2314–2325, 2008.

[3] L. Keller, M. J. Siavoshani, C. Fragouli, K. Argyraki, and S. Diggavi, "Identity aware sensor networks," in *Proc. Int. Conf. Commun., Computing Control Appl.*, pp. 2177–2185, IEEE, 2009.

[4] W. N. Chen and I. H. Wang, "Anonymous heterogeneous distributed detection: Optimal decision rules, error exponents, and the price of anonymity," *IEEE Trans. Inform. Theory*, vol. 65, no. 11, pp. 7390–7406, 2019.

[5] W. Li and Y. Huang, "Bandwidth-constrained distributed quickest change detection in heterogeneous sensor networks: Anonymous vs non-anonymous settings," *arXiv preprint arXiv: 2202.02697*, 2022.

[6] S. Marano and P. K. Willett, "Algorithms and fundamental limits for unlabeled detection using types," *IEEE Trans. Signal Proc.*, vol. 67, no. 8, pp. 2022–2035, 2019.

[7] S. Marano and P. Willett, "Making decisions by unlabeled bits," *IEEE Trans. Signal Proc.*, vol. 68, pp. 2935–2947, 2020.

[8] J. Unnikrishnan, S. Haghighatshoar, and M. Vetterli, "Unlabeled sensing with random linear measurements," *IEEE Trans. Inform. Theory*, vol. 64, no. 5, pp. 3237–3253, 2018.

[9] S. Haghighatshoar and G. Caire, "Signal recovery from unlabeled samples," *IEEE Trans. Signal Proc.*, vol. 66, no. 5, pp. 1242–1257, 2017.

[10] A. Abid, A. Poon, and J. Zou, "Linear regression with shuffled labels," *arXiv preprint arXiv:1705.01342*, 2017.

[11] V. Emiya, A. Bonnefoy, L. Daudet, and R. Gribonval, "Compressed sensing with unknown sensor permutation," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, pp. 1040–1044, IEEE, 2014.

[12] Z. Liu and J. Zhu, "Signal detection from unlabeled ordered samples," *IEEE Commun. Lett.*, vol. 22, no. 12, pp. 2431–2434, 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. Inform. Theory*, vol. 64, no. 5, pp. 3286–3300, 2017.

[14] G. Elhami, A. Scholefield, B. B. Haro, and M. Vetterli, "Unlabeled sensing: Reconstruction algorithm and theoretical guarantees," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, pp. 4566–4570, Ieee, 2017.

[15] Y. M. Lu and M. N. Do, "A theory for sampling signals from a union of subspaces," *IEEE Trans. Signal Proc.*, vol. 56, no. 6, pp. 2334–2345, 2008.

[16] G. Wang, J. Zhu, R. S. Blum, P. Willett, S. Marano, V. Matta, and P. Braca, "Signal amplitude estimation and detection from unlabeled binary quantized samples," *IEEE Trans. Signal Proc.*, vol. 66, no. 16, pp. 4291–4303, 2018.

[17] Z. Sun, S. Zou, R. Zhang, and Q. Li, "Quickest change detection in anonymous heterogeneous sensor networks," *IEEE Trans. Signal Proc.*, 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, 2021.

[19] A. G. Tartakovsky and V. V. Veeravalli, "Change-point detection in multichannel and distributed systems," *Applied Sequential Methodologies: Real-World Examples with Data Analysis*, 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 Proc.*, vol. 54, no. 9, pp. 3372–3382, 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.*, pp. 670–692, 2013.

[23] G. Fellouris and G. Sokolov, "Second-order asymptotic optimality in multisensor sequential change detection," *IEEE Trans. Inform. Theory*, vol. 62, no. 6, pp. 3662–3675, 2016.

[24] V. Raghavan and V. V. Veeravalli, "Quickest change detection of a Markov process across a sensor array," *IEEE Trans. Inform. Theory*, vol. 56, no. 4, pp. 1961–1981, 2010.

[25] O. Hadjiliadis, H. Zhang, and H. V. Poor, "One shot schemes for decentralized quickest change detection," *IEEE Trans. Inform. Theory*, vol. 55, no. 7, pp. 3346–3359, 2009.

[26] M. Ludkovski, "Bayesian quickest detection in sensor arrays," *Seq. Anal.*, vol. 31, no. 4, pp. 481–504, 2012.

[27] S. Zou, V. V. Veeravalli, J. Li, and D. Towsley, "Quickest detection of dynamic events in networks," *IEEE Trans. Inform. Theory*, vol. 66, no. 4, pp. 2280–2295, 2020.

[28] V. V. Veeravalli, "Decentralized quickest change detection," *IEEE Trans. Inform. Theory*, vol. 47, no. 4, pp. 1657–1665, 2001.

[29] A. G. Tartakovsky and V. V. Veeravalli, "Asymptotically optimal quickest change detection in distributed sensor systems," *Seq. Anal.*, vol. 27, no. 4, pp. 441–475, 2008.

[30] 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. (ICASSP)*, pp. 8454–8458, 2019.

[31] 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, 2021.

[32] G. Rovatsos, G. Moustakides, and V. Veeravalli, "Quickest detection of a dynamic anomaly in a sensor network," in *Proc. Asilomar Conf. Signals, Systems and Computers*, pp. 98–102, 2019.

[33] G. Rovatsos, S. Zou, and V. V. Veeravalli, "Sequential algorithms for moving anomaly detection in networks," *Seq. 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. Inform. Theory*, vol. 65, no. 3, pp. 1397–1412, 2018.

[35] R. Zhang, R. Yao, Y. Xie, and F. Qiu, "Quickest detection of cascading failure," *arXiv preprint arXiv:1911.05610*, 2019.

[36] D. Siegmund and E. S. Venkatraman, "Using the generalized likelihood ratio statistic for sequential detection of a change-point," *Ann. Statist.*, pp. 255–271, 1995.

[37] T. L. Lai, "Information bounds and quick detection of parameter changes in stochastic systems," *IEEE Trans. Inform. Theory*, vol. 44, no. 7, pp. 2917–2929, 1998.

[38] T. Banerjee and V. V. Veeravalli, "Data-efficient minimax quickest change detection with composite post-change distribution," *IEEE Trans. Inform. Theory*, vol. 61, no. 9, pp. 5172–5184, 2015.

[39] G. Lorden, "Procedures for reacting to a change in distribution," *Ann. Math. Statist.*, vol. 42, no. 6, pp. 1897–1908, 1971.

[40] D. Williams, *Probability with Martingales*. Cambridge University Press, 1991.

[41] G. V. Moustakides, "Optimal stopping times for detecting changes in distributions," *Ann. Statist.*, vol. 14, no. 4, pp. 1379–1387, 1986.