

# ACCURATE EVALUATION OF QUICKEST CHANGEPOINT DETECTORS VIA NON-PARAMETRIC SURVIVAL ANALYSIS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We propose non-parametric estimators for the average run length (ARL) and average detection delay (ADD) in quickest changepoint detection (QCD) under finite and irregular sequence lengths. Although ARL and ADD are widely used as optimality criteria in theoretical and simulation studies, their application to real-world datasets is hindered by limited and irregular sequence lengths. To address this issue, we propose non-parametric estimators for the ARL and ADD, termed *KM-ARL* and *KM-ADD*, by drawing an analogy between QCD and survival analysis to model detection probabilities under sequence truncation. We derive estimation bias bounds and prove that they are asymptotically unbiased unless extrapolation is required. Experiments on simulated and real-world datasets demonstrate their practical utility, enhancing robustness against limited and irregular sequence lengths, improving interpretability, and facilitating empirical, intuitive model selection. Our Python code are provided in the supplementary material and will be released upon acceptance, offering ready-to-use implementations for practitioners.

## 1 INTRODUCTION

We study evaluation metrics for models in online quickest changepoint detection (QCD) with unknown pre- and post-change distributions, where their datasets with changepoint labels are available. QCD has been extensively studied theoretically, with many models proposed and their optimality proven (Tartakovsky, 2019). It also has diverse real-world applications, including statistical process control (Hawkins et al., 2003), industrial quality control (Wadinger et al., 2024), epidemiology (Johnson & Pedersen, 2025), wireless sensor networks (Hadjiliadis et al., 2009), health monitoring (Tan et al., 2023), radar target detection (Xiang et al., 2021), and seismic sensing (Li et al., 2016).

The average run length (ARL) and the average detection delay (ADD) are central to the theoretical and simulation analysis of QCD models (Tartakovsky et al., 2014; Tartakovsky, 2019). The ARL is defined as the average time a detector takes to raise a false alarm, while the ADD refers to the average delay a detector takes to identify a changepoint after it has occurred. These metrics exhibit a tradeoff: reducing ADD typically shortens ARL (making the detector trigger-happy) and vice versa.

Although the ARL and ADD are widely used as optimality criteria in theoretical and simulation studies, their application to real-world datasets is challenging because practical sequences are often limited and irregular in length. Fig. 1(a) shows that naive ARL and ADD estimators yield substantial bias and variance, hindering reliable evaluation of QCD models.

To address this issue, we draw an analogy between QCD and survival analysis (Kleinbaum & Klein, 1996) and propose non-parametric estimators for the ARL and ADD. We adapt the Kaplan-Meier estimator (KME) (Kaplan & Meier, 1958) to the QCD setting to model detection probabilities under sequence truncation. These estimators are non-parametric, requiring no assumptions on the underlying data distribution (e.g., exponential). We derive estimation bias bounds for these estimators, termed *KM-ARL* and *KM-ADD*, by decomposing the estimation bias into *finite-sample bias* and *truncation bias*. We then show that the finite-sample bias decays exponentially with increasing dataset size and that the truncation bias is smaller than that of conventional estimators. Building on these findings, we prove that the *KM-ARL* and *KM-ADD* are *asymptotically unbiased* unless extrapolation is required.

To demonstrate practical applicability, we conduct experiments on both simulated and real-world datasets with limited and irregular lengths. The results show that our estimators reduce estimation bias compared to baseline estimators and enhance robustness to limited and irregular sequence lengths, thereby improving interpretability and facilitating empirical, intuitive model selection (Fig. 1(b)). Our code is provided in the supplementary material and will be publicly released upon acceptance, offering off-the-shelf, ready-to-use implementations in Python (Van Rossum & Drake, 2009) for practitioners.

Our contributions are threefold: (1) we propose KM-ARL and KM-ADD, non-parametric estimators for the ARL and ADD, enabling their evaluation on real-world datasets with limited and irregular-length sequences; (2) we derive estimation bias bounds for these estimators and prove that they are asymptotically unbiased unless extrapolation is required; and (3) we demonstrate their practical utility through experiments and provide ready-to-use Python implementations.

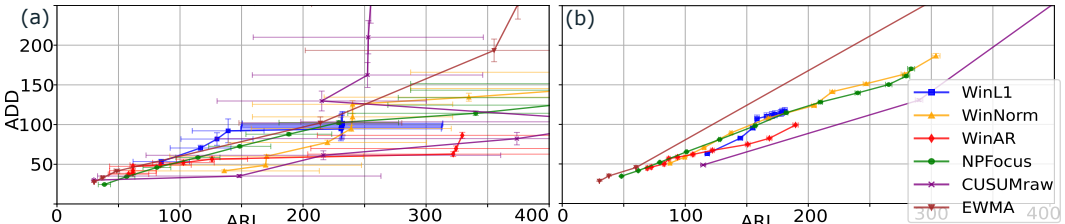


Figure 1: **Evaluation of QCD models on real-world dataset (machine-labeled subset of WISDM Actitracker).** (a) LB-ARL & LB-ADD. (b) KM-ARL & KM-ADD. Error bars represent the standard error of the mean. Our KM-ARL and KM-ADD are more robust to irregular lengths than the conventional estimators (LB-ARL and LB-ADD: see Sec. 3). The figures are shown at the same magnification scale for readability. The complete figure is provided in Fig. 14 in App. C.5. Additional evaluation on a Gaussian process dataset is provided in App. C.3.

## 2 RELATED WORK

We highlight our contribution in the context of prior research on ARL and ADD estimation using survival analysis. To our knowledge, there is no metric that explicitly address irregular sequence lengths in QCD. (Sahki et al., 2020) empirically estimate ARLs and ADDs on truncated sequences of *fixed* length, leveraging *parametric* survival analysis, assuming the survival function decays exponentially. (Bradley et al., 2023) employ a *semi-parametric* survival analysis method, the Cox proportional hazards model (Cox, 1972), which assumes an exponential response of the estimator to covariates. They also use the KME to empirically evaluate intrusion detectors under *fixed* length truncation; however, the ARL and ADD are not directly estimated, and no theoretical analysis is provided. (Lim & Lee, 2025) propose ad hoc estimators for the ARL and ADD on sequences of *fixed* length; however, their theoretical justification has yet to be established, as these estimators diverge when the sequence length  $\rightarrow \infty$ . In contrast, we build our estimators on *non-parametric* survival analysis, eliminating the exponential assumptions, derive bias bounds, and prove asymptotic unbiasedness for the proposed estimators. Moreover, our analysis includes *irregular* sequence lengths, thereby addressing more practical scenarios. We also provide a ready-to-use implementations of our estimators for practitioners. Supplementary related work is given in App. E.

## 3 KAPLAN-MEIER ARL & ADD

We first introduce our notation. For comparison, we then introduce conventional estimators under random length truncation. Finally, we present our proposed estimators.

**Definitions.**  $[n]$  denotes  $\{1, 2, \dots, n\}$  with  $n \in \mathbb{N}$ . We use  $P$  as a probability distribution.  $X^{(0,t)} = (X^{(0)}, X^{(1)}, \dots, X^{(t)})$  denote real-valued random variables representing a sequence of length  $t + 1$  ( $t \in \mathbb{Z}_{\geq 0}$ ). Each frame  $X^{(s)}$  ( $s \in \mathbb{Z}_{\geq 0}$ ) is sampled from the pre-change density  $g(X^{(s)} | X^{(0,s-1)})$  when  $s < \nu$  and from the post-change density  $f(X^{(s)} | X^{(0,s-1)})$  when  $s \geq \nu$ , where the changepoint

$\nu \in \mathbb{Z}_{\geq 0} \cup \{\infty\}$  is a random variable independent of the observations.  $\nu = \infty$  and  $\nu = 0$  indicate that all frames in the sequence are sampled from the pre-change and post-change densities, respectively. Let  $T \in \mathbb{Z}_{\geq 0}$  denote a random variable representing the sequence length, which is independent of the observations. Let  $\tau : X^{(0,T)} \mapsto \tau(X^{(0,T)}) \in \mathbb{Z}_{\geq 0} \cup \{\infty\}$  be a changepoint detector, where  $\tau(X^{(0,T)})$  is the detection point of sequence  $X^{(0,T)}$ . We set  $\tau(X^{(0,T)}) = \infty$  when no change is detected within the input sequence of finite length. The ARL and ADD are defined as

$$\begin{aligned} \mu_\infty &:= \mathbb{E}[\tau \mid \nu = \infty] \\ \text{and } M_\infty &:= \mathbb{E}[\Delta\tau \mid \Delta\tau \geq 0, \nu < \infty], \end{aligned} \quad (1)$$

respectively, where  $\Delta\tau := \tau - \nu$ .  $T = \infty$  is assumed here, and the expectation is also taken over  $\nu$ . We consider estimation of the ARL and ADD from a dataset  $\{(X_i^{(0,T_i)}, \nu_i)\}_{i=1}^N$  with predicted detection points  $\{\tau_i\}_{i=1}^N$  from a detector  $\tau$ , where  $N \in \mathbb{N}$  denotes the dataset size,  $X_i^{(0,T_i)}$  is a sequence with a random length  $T_i$ ,  $\nu_i$  is the corresponding changepoint, and  $\tau_i$  is the detection point of  $X_i^{(0,T_i)}$ . Let  $\langle \cdot \rangle_{i:\mathcal{C}} := \sum_{i:\mathcal{C}} \cdot / |\{i : \mathcal{C}\}|$  denote the empirical expectation under condition  $\mathcal{C}$ .

**Conventional estimators.** A conventional estimator of ARL under random sequence lengths, adapted from (Qiu, 2013) and referred to as the Less-Biased ARL (LB-ARL) in this paper, follows the definition of  $\mu_\infty$ :  $\hat{\mu}_T^{(\text{LB})} := \langle \tau_i \rangle_{i:\nu_i=\infty, \tau_i \leq T_i}$ , where only sequences satisfying  $\tau_i \leq T_i$  are included. Similarly, a conventional ADD estimator, the LB-ADD, can be defined as  $\hat{M}_T^{(\text{LB})} := \langle \Delta\tau_i \rangle_{i:\Delta\tau_i \geq 0, \nu_i < \infty, \Delta\tau_i \leq \Delta T_i}$ , where  $\Delta T_i := T_i - \nu_i$ .

A critical drawback of these conventional estimators is that they ignore sequences in which the QCD model fails to raise an alarm before the horizon ( $\tau > T$ ). Although the LB-ARL has been employed in Monte Carlo simulations of the ARL (Qiu, 2013; Lim & Lee, 2025) when  $T = \text{const.} < \infty$ , it exhibits a more substantial negative bias than our estimator due to truncation, as proven in Sec.4 and demonstrated in Sec.5.

**Key ideas from survival analysis.** To overcome this challenge, we model the detection probability beyond the truncation length, using a non-parametric approach inspired by *survival analysis* (Kleinbaum & Klein, 1996). A central interest in survival analysis is the *survival function*  $S(t) := P(\text{event time} > t)$ , representing the probability that a patient survives beyond time  $t$  ( $t \in \mathbb{R}_{\geq 0}$ ), where the *event time* refers to the time of death. Crucially, the estimation of  $S(t)$  on a dataset is performed under *right-censoring*; i.e., the exact event times of several patients are unknown because the patients are lost to follow-up or are still under observation at the end of the study. Thus, we only know the lower bound of the event times of these patients, called *censoring times*. The mean survival time is given by the integral of  $S(t)$  over time, i.e., the area under the survival curve, because  $\int_0^\infty S(t)dt = \int_0^\infty P(\text{event time} > t)dt = \int_0^\infty \mathbb{E}[\mathbb{1}(\text{event time} > t)]dt = \mathbb{E}[\int_0^\infty \mathbb{1}(\text{event time} > t)dt] = \mathbb{E}[\text{event time}]$ , where  $\mathbb{1}$  is the indicator function.

**KM-ARL.** To estimate the ARL under irregular sequence lengths, we draw an analogy by regarding a patient as a sequence, the event time as the detection point, the censoring time as the minimum of the changepoint and the sequence length, and the mean survival time as the ARL (see Tab. 1 in App. A for reference). Under this analogy, we estimate the *survival function of detection points*  $S^{\text{ARL}}(t) := P(\tau > t \mid \nu = \infty)$  using the Kaplan-Meier estimator (KME) (Kaplan & Meier, 1958), a non-parametric estimator of the survival function:  $\hat{S}^{\text{ARL}}(t) = \prod_{j:t_j^{\text{ARL}} \leq t} (1 - \frac{d_j^{\text{ARL}}}{n_j^{\text{ARL}}})$ , where  $0 < t_1^{\text{ARL}} < t_2^{\text{ARL}} < \dots < t_{N'}^{\text{ARL}}$  are distinct detection points, with  $N'_{\text{ARL}} \in \mathbb{N}$  ( $\leq N$ );  $d_j^{\text{ARL}} := |\{i \in [N] \mid \tau_i = t_j^{\text{ARL}}\}|$  is the number of sequences with a detection at  $t_j^{\text{ARL}}$  ( $j \in [N'_{\text{ARL}}]$ ); and  $n_j^{\text{ARL}} := |\{i \in [N] \mid \min\{\tau_i, C_i^{\text{ARL}}\} \geq t_j^{\text{ARL}}\}|$  is the number of sequences neither detected nor censored prior to  $t_j^{\text{ARL}}$  ( $j \in [N'_{\text{ARL}}]$ ), with the censoring time of sequence  $i$  defined as  $C_i^{\text{ARL}} := \min\{\nu_i, T_i\}$ . An example of  $\hat{S}^{\text{ARL}}(t)$  is shown in App. C.1. We propose a non-parametric estimator of the ARL under irregular sequence lengths, termed KM-ARL, as the integral of  $\hat{S}^{\text{ARL}}(t)$  over the range  $[0, a]$  for arbitrary  $a \in \mathbb{R}_{\geq 0}$ :

$$\hat{\mu}_T^{(\text{KM})} := \int_0^a \hat{S}^{\text{ARL}}(t)dt. \quad (2)$$

In practice, we set  $a = T_{\max} := \max_i \{\min\{\tau_i, C_i^{\text{ARL}}\}\}$ , i.e., the maximum last-observed time, following standard practice in survival analysis (Qi & Wang, 2018; Calkins et al., 2018). For a given dataset, choosing  $a > T_{\max}$  is irrelevant because it is extrapolation beyond the observed support. Theoretically ideal choices of  $a$  are given in Sec. 4.

**KM-ADD.** For the ADD, we regard a patient as a sequence with  $\nu_i < \infty$  and  $\tau_i \geq \nu_i$ , the event time as the detection delay  $\Delta\tau_i (\geq 0)$ , the censoring time as the sequence length measured from the changepoint ( $C_i^{\text{ADD}} := \Delta T_i = T_i - \nu_i$ ), and the mean survival time as the ADD (see Tab. 1 in App. A for reference). Under this analogy, we propose a non-parametric estimator of the ADD under irregular lengths, termed KM-ADD, as

$$\hat{M}_T^{(\text{KM})} := \int_0^b \hat{S}^{\text{ADD}}(t) dt, \quad (3)$$

where  $\hat{S}^{\text{ADD}}(t) := \prod_{j: t_j^{\text{ADD}} \leq t} (1 - \frac{d_j^{\text{ADD}}}{n_j^{\text{ADD}}})$  is a non-parametric estimate of the *survival function of detection delays*  $S^{\text{ADD}}(t) := P(\Delta\tau > t \mid \Delta\tau \geq 0, \nu < \infty)$ ;  $0 \leq t_1^{\text{ADD}} < t_2^{\text{ADD}} < \dots < t_{N'_{\text{ADD}}}^{\text{ADD}}$  are the distinct detection delays, i.e., the sorted unique values of  $\Delta\tau_i \geq 0$ , with  $N'_{\text{ADD}} \in \mathbb{N} (\leq N)$ ;  $d_j^{\text{ADD}} := |\{i \in [N] \mid 0 \leq \Delta\tau_i = t_j^{\text{ADD}}\}|$  is the number of sequences with positive detection delay equal to  $t_j^{\text{ADD}}$  ( $j \in [N'_{\text{ADD}}]$ ); and  $n_j^{\text{ADD}} := |\{i \in [N] \mid \min\{\Delta\tau_i, C_i^{\text{ADD}}\} \geq t_j^{\text{ADD}}, \Delta\tau_i \geq 0\}|$  is the number of sequences neither detected nor censored prior to  $t_j^{\text{ADD}}$  ( $j \in [N'_{\text{ADD}}]$ ). Again, the upper limit  $b$  is arbitrary, and we set  $b = \Delta T_{\max} := \max_i \{\min\{\Delta\tau_i, C_i^{\text{ADD}}\} \mid \Delta\tau_i \geq 0, \nu_i < \infty\}$  in experiment.

## 4 BIAS ANALYSIS

We derive bias bounds for the KM-ARL and KM-ADD and prove that they are asymptotically unbiased unless extrapolation is required. We first examine the estimation bias of the KM-ARL:

$$\mathcal{B}(\hat{\mu}_T^{(\text{KM})}) := \mathbb{E}[\hat{\mu}_T^{(\text{KM})}] - \mu_\infty, \quad (4)$$

where  $\mu_\infty := \mathbb{E}[\tau \mid \nu = \infty] < \infty$  is the true ARL under infinite sequence length (we assume all relevant finiteness hereafter). We decompose this bias into two components: the *finite-sample bias* and the *truncation bias*

$$\mathcal{B}_{\text{FS}}(\hat{\mu}_T^{(\text{KM})}) := \mathbb{E}[\hat{\mu}_T^{(\text{KM})}] - \mu_T^{(\text{KM})} \quad (5)$$

$$\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})}) := \mu_T^{(\text{KM})} - \mu_\infty^{(\text{KM})} \quad (6)$$

so that  $\mathcal{B}(\hat{\mu}_T^{(\text{KM})}) = \mathcal{B}_{\text{FS}}(\hat{\mu}_T^{(\text{KM})}) + \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})})$ , where  $\mu_T^{(\text{KM})} := \int_0^a S_{\text{ARL}}(t) dt$  is the KM-ARL with the true survival function of detection points.

**Finite-sample bias.**  $\mathcal{B}_{\text{FS}}(\hat{\mu}_T^{(\text{KM})})$  quantifies the error arising from a finite dataset of sequences with finite and irregular lengths and dominates the total bias when the dataset size is small. We derive a bound for  $\mathcal{B}_{\text{FS}}(\hat{\mu}_T^{(\text{KM})})$ , indicating that it decays exponentially as  $N \rightarrow \infty$ :

**Theorem 4.1** (Finite-sample bias bounds for KM-ARL). *We idealize the time index as continuous without loss of generality, for technical convenience. Let  $F^{\text{ARL}}, G^{\text{ARL}}$ , and  $H^{\text{ARL}}$  be the cumulative distribution functions (CDFs) of  $\tau, C^{\text{ARL}}$ , and  $\min\{\tau, C^{\text{ARL}}\}$ , respectively. Assume that (i)  $F^{\text{ARL}}$  and  $G^{\text{ARL}}$  do not have common discontinuities, (ii)  $\tau$  and  $C^{\text{ARL}}$  are independent, known as the independent censoring (or non-informative censoring) assumption (Ranganathan & Pramesh, 2012), and (iii)  $H^{\text{ARL}}$  is continuous. Then, we have for any  $a \in \mathbb{R}_{\geq 0}$*

$$\begin{aligned} & - \int_0^a t G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t) \\ & \leq \mathcal{B}_{\text{FS}}(\hat{\mu}_T^{(\text{KM})}) \leq \int_0^a a G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t). \end{aligned} \quad (7)$$

Our proof is based on (Stute & Wang, 1993; Stute, 1994) and given in App. B.1. Assumption (i) is a technical requirement and is valid unless pathological situations are considered. Assumption

(ii) holds for online QCD models that do not look ahead the input sequence, which are the focus in this paper. This is because: detection points  $\tau_i$  prior to censoring can be regarded as samples from  $P(\tau | \nu = \infty)$ ;  $P(\tau | \nu = \infty)$ ,  $P(\nu)$ , and  $P(T)$  are independent; and thus,  $P(\tau | \nu = \infty)$  and  $P(C := \min\{\nu, T\})$  are also independent. In contrast, Assumption (ii) does not hold for offline changepoint detection models because  $\tau$  depends on  $X^{(0,T)}$ , which in turn depends on  $\nu$  and  $T$ . Assumption (iii) is also a technical requirement.

According to Thm. 4.1, the finite-sample bias exhibits the following properties. First, it decays exponentially to zero as  $N \rightarrow \infty$  if  $H^{\text{ARL}}(t) < 1$  for  $t \in [0, a]$ . Second, it vanishes when  $G^{\text{ARL}} = 0$ , i.e., in the absence of censoring (no changepoints or horizons) in  $t \in [0, a]$ , which is desirable because estimation becomes more accurate as censoring decreases. Third, the bound becomes looser for larger  $a$  because  $G^{\text{ARL}}$  and  $H^{\text{ARL}}$  approach 1 monotonically as  $t$  increases. However, the truncation bias decreases for larger  $a$ , leading to a tradeoff between the finite-sample and truncation biases. Finally, we note that empirically verifying the convergence of the finite-sample bias is challenging because, for small sample sizes, estimation variance overshadows the bias.

We can derive a similar bound for the finite-sample bias of the KM-ADD (proof is given in App. B.2), ensuring that it also decays exponentially as  $N'' \rightarrow \infty$ , where  $N''$  is the number of sequences with  $\nu_i < \infty$  and  $\Delta\tau_i \geq 0$ :

**Theorem 4.2** (Finite-sample bias bounds for KM-ADD). *We idealize the time index as continuous without loss of generality, for technical convenience. Consider only sequences with  $\nu_i < \infty$  and  $\Delta\tau_i \geq 0$  in the dataset. Let  $F^{\text{ADD}}, G^{\text{ADD}}$ , and  $H^{\text{ADD}}$  be the CDFs of  $\Delta\tau$ ,  $C^{\text{ADD}}$ , and  $\min\{\Delta\tau, C^{\text{ADD}}\}$ , respectively. Assume that (i)  $F^{\text{ADD}}$  and  $G^{\text{ADD}}$  do not have common discontinuities, (ii)  $\Delta\tau$  and  $C^{\text{ADD}}$  are independent, and (iii)  $H^{\text{ADD}}$  is continuous. Then, for any  $b \in \mathbb{R}_{\geq 0}$ ,*

$$\begin{aligned} & - \int_0^b t G^{\text{ADD}}(t) H^{\text{ADD}}(t)^{N''-1} dF^{\text{ADD}}(t) \\ & \leq \mathcal{B}_{\text{FS}}(\hat{M}_T^{(\text{KM})}) \leq \int_0^b b G^{\text{ADD}}(t) H^{\text{ADD}}(t)^{N''-1} dF^{\text{ADD}}(t). \end{aligned} \quad (8)$$

Here, we defined the finite-sample bias of the KM-ADD similarly to that of the KM-ARL:  $\mathcal{B}_{\text{FS}}(\hat{M}_T^{(\text{KM})}) := \mathbb{E}[\hat{M}_T^{(\text{KM})}] - M_T^{(\text{KM})}$ , where  $M_T^{(\text{KM})} := \int_0^b S^{\text{ADD}}(t) dt$  is the KM-ADD with the true survival function of detection delays. Assumption (i) and (iii) are technical conditions, as previously noted below Thm. 4.1. Assumption (ii) is the independent censoring assumption for the KM-ADD, which can be justified in online QCD models by the approximation  $P(\Delta\tau | \Delta\tau \geq 0, \nu < \infty) \approx P(\tau | \nu = 0)$  and  $P(C := T - \nu | \Delta\tau \geq 0, \nu < \infty) \approx P(C := T | \nu = 0)$  because  $P(\tau | \nu = 0)$  and  $P(C := T | \nu = 0)$  are independent. This approximation implies that the distributions of the detection delay  $\tau - \nu$  and the censoring time  $T - \nu$  are approximately equal to their respective distributions measured from  $t = 0$  rather than from  $\nu$ .

**Truncation bias.**  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})})$  captures the error from sequence truncation, which dominates the total bias when sequence lengths are short. It vanishes as  $a$  increases because  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})}) = \int_0^a S^{\text{ARL}}(t) dt - \mathbb{E}[\tau | \nu < \infty] = - \int_a^\infty S^{\text{ARL}}(t) dt \rightarrow 0$ , where we used  $\mathbb{E}[\tau | \nu < \infty] = \int_0^\infty S^{\text{ARL}}(t) dt$ . The convergence rate depends on the underlying distributions. In Thm. 4.3 below, we derive a bound for  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})})$ , indicating that it is non-positive and minor than that of the conventional LB-ARL. In Sec. 5, we will empirically justify our result (Fig. 2). Proof is given in App. B.3. Note that given an evaluation dataset, we can easily specify  $T_{\text{max}}^*$  defined below (App. F).

**Theorem 4.3** (Truncation bias bound for KM-ARL). *Define the truncation bias of the LB-ARL as  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{LB})}) := \mu_T^{(\text{LB})} - \mu_\infty$ , where  $\mu_T^{(\text{LB})} := \mathbb{E}[\tau | \nu = \infty, \tau \leq T]$  is the true ARL under random lengths. For  $a = T_{\text{max}}^*$ , where  $T_{\text{max}}^* := \inf\{T | \text{CDF}(T) = 1\}$  is the least upper bound for the support of the CDF of  $T$ , we have  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{LB})}) \leq \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})}) \leq 0$ .*

A similar bound holds for the KM-ADD. Proof is given in App. B.4.

**Theorem 4.4** (Truncation bias bound for KM-ADD). *Define the truncation bias of the LB-ADD as  $\mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{LB})}) := M_T^{(\text{LB})} - M_\infty$ , where  $M_T^{(\text{LB})} := \mathbb{E}[\Delta\tau | \nu < \infty, 0 \leq \Delta\tau \leq \Delta T]$  is the true ADD under random sequence lengths. For  $b = \Delta T_{\text{max}}^*$ , where  $\Delta T_{\text{max}}^* := \inf\{\Delta T := T - \nu |$*

CDF( $\Delta T$ ) = 1}, we have  $\mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{LB})}) \leq \mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{KM})}) \leq 0$ , under the independent censoring assumption in Thm. 4.2.

**Total Estimation Bias.** Finally, we examine the total estimation bias of the KM-ARL:  $\mathcal{B}(\hat{\mu}_T^{(\text{KM})}) = \mathcal{B}_{\text{FS}}(\hat{\mu}_T^{(\text{KM})}) + \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})})$  (a similar discussion below holds for the KM-ADD). From Thm. 4.1 and  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})}) = -\int_a^\infty S^{\text{ARL}}(t)$ , we have for any  $a \in \mathbb{R}_{\geq 0}$

$$\begin{aligned} & -\int_0^a t G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t) - \int_a^\infty S^{\text{ARL}}(t) \\ & \leq \mathcal{B}(\hat{\mu}_T^{(\text{KM})}) \leq \int_0^a a G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t) - \int_a^\infty S^{\text{ARL}}(t). \end{aligned} \quad (9)$$

Define  $t_F := \inf\{t \in \mathbb{R}_{\geq 0} \mid F^{\text{ARL}}(t) = 1\}$  and  $t_H := \inf\{t \in \mathbb{R}_{\geq 0} \mid H^{\text{ARL}}(t) = 1\}$ , and set  $a$  to  $t_F$ . Then, since  $S^{\text{ARL}} = 1 - F^{\text{ARL}}$ , we have  $\int_a^\infty S^{\text{ARL}}(t) = 0$ ; i.e., the truncation bias vanishes. Therefore,

$$\begin{aligned} & -\int_0^{t_F} t G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t) \\ & \leq \mathcal{B}(\hat{\mu}_T^{(\text{KM})}) \leq \int_0^{t_F} t_F G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t). \end{aligned} \quad (10)$$

If  $t_F < t_H$ , the bound decays exponentially as  $N \rightarrow \infty$  because  $H^{\text{ARL}}(t) < 1$  for  $t \in [0, t_F]$ ; i.e., the KM-ARL is asymptotically unbiased if  $t_F < t_H$ . Otherwise, there are non-vanishing terms  $-\int_0^{t_F} t G^{\text{ARL}}(t) dF^{\text{ARL}}(t)$  and  $\int_0^{t_F} t_F G^{\text{ARL}}(t) dF^{\text{ARL}}(t)$  in the lower and upper bounds, respectively. This is reasonable because  $t_F \geq t_H$  implies that detection points may occur beyond censoring times with non-zero probability; i.e., not all detection points are observable. Consequently, estimating the ARL without bias or additional assumptions is impossible, necessitating extrapolation.

## 5 EXPERIMENT

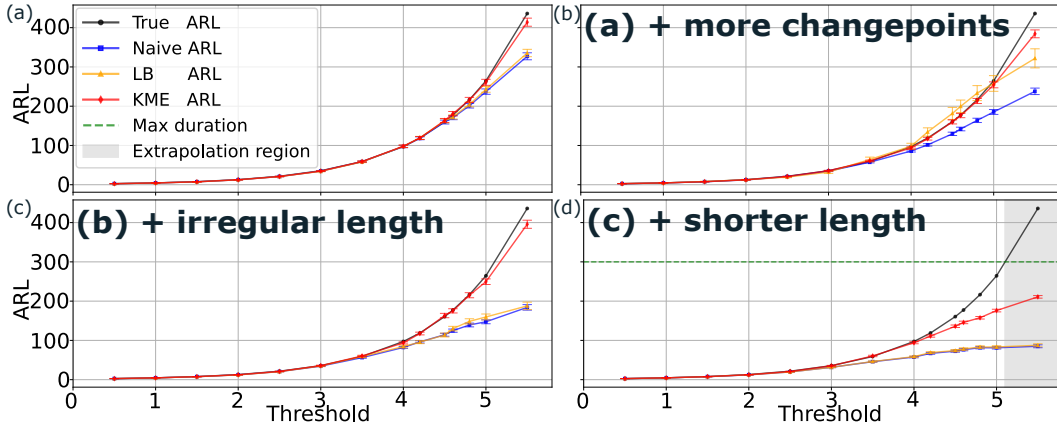


Figure 2: **Threshold of detector vs. ARL.** The KM-ARL provides more accurate estimates of the true ARL even when sequences contain changepoints and have limited and irregular lengths. The Gaussian process dataset contains 1000 sequences. The GSR with ground-truth statistics is evaluated under various thresholds. Changepoints are sampled uniformly. Error bars represent the standard error of the mean. (a) **Sequence length is 1000, with 10% of sequences containing a changepoint.** (b) **Sequence length is 1000, with 90% of sequences containing a changepoint.** (c) **Sequence lengths vary irregularly in the range [100, 1000], with 90% of sequences containing a changepoint.** (d) **Sequence lengths vary irregularly in the range [30, 300], with 90% of sequences containing a changepoint.** In (d),  $\text{ARL} > 300$  (gray area) is an extrapolation region (excluding the true ARL).

To demonstrate the practical relevance of our KM-ARL and KM-ADD for sequences with limited and irregular lengths, we use both simulated and real-world datasets. Our results show that these metrics

324 reduce estimation bias compared to baseline estimators, enhance robustness to such sequences, and  
 325 thereby improve interpretability and facilitate empirical, intuitive model selection. Our code is given  
 326 in the supplementary material and will be released upon acceptance.  
 327

328 **Datasets.** We use two simulation datasets, the Gaussian and Poisson processes, and one real-world  
 329 dataset, the WISDM Actitracker (Kwapisz et al., 2011). The pre-change and post-change Gaussian  
 330 processes have the mean of 0 and 0.1, respectively, with preserving the variance equal to 0.1. We  
 331 use two types of changepoint distributions: geometric and uniform. Several sequences have no  
 332 changepoint, and the number of with-change sequences depends on the changepoint distributions.  
 333 We also simulate irregular length by randomly truncating the sequences. The setup and experimental  
 334 results for Poisson processes are provided in App. C.4 due to page limitations. The results are similar  
 335 to those obtained for the Gaussian processes. The WISDM Actitracker is a large, real-world dataset  
 336 for smartphone-based human activity recognition (walking, jogging, stairs, sitting, standing, and  
 337 lying down) collected by the WISDM Lab at Fordham University, offering both user-labeled and  
 338 machine-labeled data. The labels specify activities and their temporal intervals. We provide the  
 339 results for the machine-labeled subset in the main text, and the results for the user-labeled subset are  
 340 provided in Fig. 15 in App. C.5. The machine-labeled subset contains 51,326 sequences, and the  
 341 sequence lengths exhibit substantial irregularity from 1 to 54,401 after our preprocesses. See App. D.3  
 342 for more statistical information of the WISDM Actitracker dataset. We remove the sequences with  
 343 length 1, as they are not informative when computing performance metrics of QCD models. The  
 344 preprocesses are detailed in our code and App. D.2.

345 5.1 ESTIMATION ON SIMULATION DATASET

346 **Detection models.** The QCD model used for the simulation datasets is the generalized  
 347 Shiryaev-Roberts (GSR) procedure with ground-truth statistics. The GSR raise an alarm when  
 348 the following statistic hits a pre-defined threshold:  $R(t) = \omega \prod_{s=0}^t \mathcal{L}(s) + \sum_{k=0}^t \prod_{s=k}^t \mathcal{L}(s)$ ,  
 349 where the likelihood ratio is denoted by  $\mathcal{L}(t) = f(X^{(t)} | X^{(0,t-1)}) / g(X^{(t)} | X^{(0,t-1)})$ , and  
 350  $\omega$  controls the strength of warm start and is set to 0 in our experiments. Additionally, we provide  
 351 the results for the cumulative sum (CUSUM) procedure in App. C.6.  
 352  
 353  
 354  
 355  
 356  
 357

358 **True ARLs & ADDs.** For the experiments on the Gaussian and Poisson process datasets,  
 359 we simulate the ground-truth ARLs and ADDs by generating sequences of effectively infinite  
 360 length. Their error bars are omitted because all errors are sufficiently small, with relative errors  
 361  $\lesssim 10^{-3}$ .  
 362  
 363  
 364  
 365

366 **Naive ARL.** We introduce another baseline metric, referred to as the Naive ARL. While  
 367 the LB-ARL use only sequences with  $\nu_i = \infty$  and  $\tau_i \leq T_i$  for computing the ARL, the Naive  
 368 ARL utilizes sequences with  $\tau_i < \nu_i$  and  $\tau_i \leq T_i$ :  $\hat{\mu}_T^{(NV)} := \langle \tau_i \rangle_{i: \tau_i < \nu_i, \tau_i \leq T_i}$ . The number  
 369 of sequences used for the Naive ARL is larger than that of the LB-ARL because the condition  
 370  $\tau_i < \nu_i \wedge \tau_i \leq T_i$  includes  $\nu_i < \infty \wedge \tau_i \leq T_i$ ; however, the Naive ARL has a non-vanishing  
 371 bias because  $\mathbb{E}[\hat{\mu}_T^{(NV)}] = \mathbb{E}[\tau | \tau < \nu, \tau < T]$  under minor assumptions, which is not equal to  
 372  $\mathbb{E}[\tau | \nu = \infty, \tau < T]$  or  $\mathbb{E}[\tau | \nu = \infty]$ .  
 373  
 374  
 375  
 376  
 377

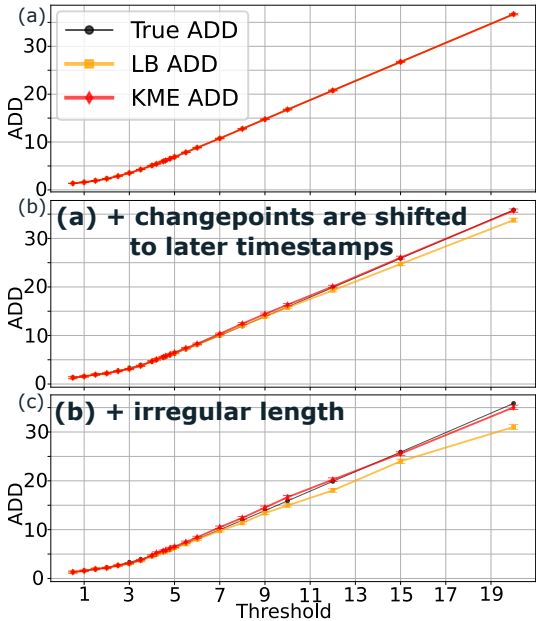


Figure 3: **Threshold of detector vs. ADD.** The KM-ADD provides more accurate estimates of the true ADD even when sequence lengths are limited and irregular. The Gaussian process dataset has 10000 sequences. Changepoints are sampled from the geometric distribution with the success probability  $p \in [0, 1]$ . A smaller  $p$  leads to more sparse and delayed changepoints. Other settings are identical to that of Fig. 2. **(a) Sequence length is 100, with  $p = 0.25$ .** **(b) Sequence length is 100, with  $p = 0.001$ .** The changepoints are shifted to later timestamps, decreasing the chance of detection due to the length limit. **(c) Sequence lengths vary irregularly in the range  $[10, 100]$ , with  $p = 0.001$ .**

378 **Result on ARL.** Fig. 2 presents the true ARL,  
 379 Naive ARL, LB-ARL, and KM-ARL across various GSR thresholds, demonstrating that the KM-ARL  
 380 provides more accurate estimates of the true ARL even when sequences contain changepoints and  
 381 have limited and irregular lengths. Fig. 2(a) simulates an ideal light-censoring scenario, where  
 382 most detection points occur before changepoints or truncation. The true ARL is less than  $T_{\max}$ , the  
 383 sequence length  $T$  is constant ( $=1000$ ), and only 10% of the sequences contain a changepoint. Fig. 2(b)  
 384 simulates a heavier-censoring scenario, where 90% of the sequences contain a changepoint. Fig. 2(c)  
 385 simulates an even more heavily censored scenario, where, in addition to Fig. 2(b), sequence lengths are  
 386 irregular and randomly sampled from  $[100, 1000]$ . Our KM-ARL remains robust under this condition,  
 387 utilizing with-change and truncated sequences, unlike the Naive ARL and the LB-ARL. Fig. 2(d)  
 388 simulates a challenging and realistic scenario, where, in addition to Fig. 2(c),  $T_{\max}$  is reduced to  
 389 300, and sequence lengths are sampled from  $[30, 300]$ . The gray area indicates  $ARL > T_{\max} = 300$ ,  
 390 where no data are observed, i.e., an extrapolation region (excluding the true ARL). Bias becomes  
 391 non-negligible in  $ARL \gtrsim T_{\max}$  due to truncation. In this region, no metric can reliably estimate the  
 392 ARL without bias or additional assumptions, as noted Sec. 4. To further mitigate truncation bias, one  
 393 may combine our estimators with parametric extrapolation methods for the survival function, such as  
 394 those proposed in (Sahki et al., 2020).

395 **Result on ADD.** Fig. 3 presents the true ADD,  
 396 LB-ADD, and KM-ADD across various GSR  
 397 thresholds, demonstrating that the KM-ADD  
 398 provides more accurate estimates of the true  
 399 ADD even when sequences have limited and  
 400 irregular lengths. Fig. 3(a) simulates an ideal  
 401 light-censoring scenario, where most detection  
 402 points occur before truncation.  $T$  is constant  
 403 and set to 100, and  $\Delta T_{\max}$  is also set to 100,  
 404 which is greater than the true ADD. Fig. 3(b)  
 405 simulates a heavier-censoring scenario, where  
 406 changepoints are shifted to later timestamps, re-  
 407 ducing the chance of detection before trunca-  
 408 tion. Fig. 3(c) simulates a challenging and re-  
 409 alistic scenario, where, in addition to Fig. 3(b),  
 410 sequence length are irregular and sampled from  
 411  $[10, 100]$ , further reducing the chance of detec-  
 412 tion before truncation. The KM-ADD still re-  
 413 mains robust, utilizing truncated sequences, un-  
 414 like the conventional LB-ADD.

415 **ARL-ADD tradeoff curve.** Fig. 4 shows  
 416 ARL-ADD tradeoff curves, which are theoret-  
 417 ically related to the optimality of QCD models  
 418 and are used for model evaluation in practice.  
 419 Fig. 4 demonstrates that the KM-ARL and KM-  
 420 ADD can more accurately estimate the ARL-  
 421 ADD curve than the others. Fig. 4(a) simulates a light-censored scenario, where most detection points  
 422 occur before censoring. The sequence length  $T$  is constant and set to  $T = T_{\max} = 500$ . Fig. 4(b)  
 423 simulates a challenging and realistic scenario, where censoring is heavier. Sequence lengths are  
 424 irregular and sampled from  $[50, 500]$ . The KM-ARL and KM-ADD remains more robust than the  
 425 baselines. Again, bias is non-negligible in the region  $ARL \approx T_{\max} = 500$ , where extrapolation  
 426 region is nearby.

## 427 5.2 EVALUATION OF MODELS ON REAL-WORLD DATASET

428  
 429 We further evaluate six QCD models on a real-world, challenging dataset: the WISDM Actitracker  
 430 (Kwapisz et al., 2011), which contains 51,326 sequences with substantial irregularity of lengths from  
 431 1 to 54,401. The computational costs for KM-ARL and KM-ADD are negligible in our experiments  
 compared with the time required to run the QCD algorithms.

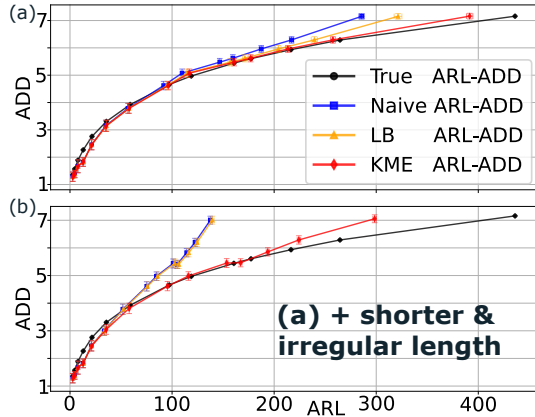


Figure 4: **ARL-ADD tradeoff curves.** The KM-ARL and KM-ADD provide more accurate estimates of the true ARL-ADD curve even when sequences contain changepoints and have limited and irregular lengths. The setup is identical to that of Fig. 2, except that the dataset size is 10000. 50% of sequences contain a changepoint. The LB-ADD is used for the Naive ARL’s ADD. **(a) Sequence length is 1000.** **(b) Sequence lengths vary irregularly in the range  $[50, 500]$ .** The complete ablation study is given in App. C.2.

**QCD models.** We use the following online QCD models, including window-based, frame-based, parametric, and non-parametric models: Window L1 (Bai, 1995), Window Normal (Lavielle, 1999; Lavielle & Teysriere, 2006), Window AR (Bai, 2000), non-parametric focused changepoint detection (NP-FOCuS) (Romano et al., 2024a), CUSUM (Page, 1954), and exponentially weighted moving average (EWMA) (Roberts, 1959). The window size and burn-in interval are fixed at 30, ensuring much smaller than the average length. Most other hyperparameters left at default values. See App. D.1 for details of the models and their hyperparameters.

**Result.** The ARL-ADD tradeoff curves are presented in Fig. 1, demonstrating that our KM-ARL and KM-ADD are more robust even when censoring is significantly heavy. Fig. 1(a) shows that LB-ARL and LB-ADD become unstable and exhibit high variance, particularly at large QCD thresholds. This instability arises because the detector often fails to raise an alarm within the sequence length, and the limited number of sequences further amplifies the empirical variance. In contrast, Fig. 1(b) shows that our KME-based estimators do not suffer from this issue. These metrics reduce variance compared to baseline estimators because they are calculated from a constant number of sequences, regardless of the threshold, which enhances their robustness against censoring. Therefore, the KM-ARL and KM-ADD improves interpretability and facilitating empirical, intuitive model selection (see App. D.4 for more details of the variance computation).

## 6 DISCUSSION, LIMITATIONS, AND FUTURE WORK

**To further reduce estimation bias.** For further bias reduction, one may use the bootstrap method (Efron & Tibshirani, 1994) to mitigate the finite-sample bias, although this effect is typically overshadowed by the truncation bias. To address the truncation bias beyond extrapolation, a parametric estimator can be combined with our estimators, as in (Sahki et al., 2020); however, note that the estimation accuracy deteriorates when the parametric assumption is invalid.

**Changepoint isolation and multiple changepoint detection.** Our work assumes a single changepoint type; however, changepoint isolation (classification) (Tartakovsky et al., 2014) is often required in practice. We conjecture that our survival analysis–based approach can be extended to this challenging setting. A promising direction is to use survival models under *competitive risks* (Morita, 2021), treating different changepoint types as distinct causes of death. Similarly, while our method does not currently support multiple changepoint detection (Niu et al., 2016), it may leverage *multi-state models under competing risks* (Therneau et al., 2020; Beyersmann et al., 2011; Mills, 2011), which are well established in survival analysis.

**Probability of false alarms.** The probability of false alarms (PFA) is another standard metric in QCD theory, commonly used to establish Bayesian optimality (Tartakovsky, 2019). However, it is also affected by right-censoring in real-world scenarios. We hypothesized that multi-state models under competitive risks, such as the Aalen–Johansen estimator (AJE) (Aalen & Johansen, 1978), could alleviate this issue. Our preliminary experiments, however, failed to accurately estimate PFAs, possibly because the AJE assumes no event ties, an assumption often violated when continuous-time sequences are discretized into frames. This finding suggests that more careful approaches, such as those from discrete-time survival analysis (Tutz et al., 2016), are required.

**Dependent censoring.** In Sec. 4, we assume independent censoring for KM-ADD, supported by the distributional approximation therein. This assumption can potentially be relaxed by leveraging extensive research on dependent censoring in survival analysis (Hsu & Taylor, 2010; Lin et al., 2023; Crommen et al., 2025). Doing so would not only eliminate the assumption but also extend our bias analysis to offline QCD, where censoring depends on detection. See App. F for more details.

**Heavy censoring.** We do not recommend using KM-ARL or KM-ADD when datasets are severely imbalanced between pre-change and post-change sequences because this leads to significantly heavy censoring and inflated finite-sample and truncation bias. While the extent of acceptable censoring depends on the dataset and underlying distributions, Malmquist (2025) reports that, for 30, 50, 100, or 150 subjects (sequences), censoring rates up to 90% are tolerable if censoring is uniform. However, for dependent censoring occurring just before event time (detection), estimation bias increases sharply

486 (see App. F). Nonetheless, several studies have proposed modifications to the KME to address heavy  
487 censoring (Shafiq et al., 2007; Zare & Mahmoodi, 2013), which can be integrated with our KM-ARL  
488 and KM-ADD.  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

## 540 REPRODUCIBILITY STATEMENT

541

542 For reproducibility, we provide the code for our experiments and proposed estimators in the Sup-  
 543 plementary Materials, which will be publicly available upon acceptance. Complete descriptions of  
 544 data processing steps for all datasets are given in Sec. 5, App. D, and in our code. Assumptions,  
 545 definitions, and full proofs of our theoretical results appear in Sec. 3, Sec. 4, and App. B.

546

## 547 REFERENCES

548

549 Odd O Aalen and Søren Johansen. An empirical transition matrix for non-homogeneous Markov  
 550 chains based on censored observations. *Scandinavian journal of statistics*, pp. 141–150, 1978.

551

552 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S.  
 553 Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew  
 554 Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath  
 555 Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah,  
 556 Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker,  
 557 Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin  
 558 Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine  
 559 learning on heterogeneous systems, 2015. License: Apache License 2.0. Software available from  
 tensorflow.org.

560

561 Michael G. Akritas. The central limit theorem under censoring. *Bernoulli*, 6(6):1109–1120, 2000.  
 562 ISSN 13507265. URL <http://www.jstor.org/stable/3318473>.

563

564 Burak Alakent and Ece C Mutlu. Application of robust estimators in Shewhart S-charts. *arXiv  
 preprint arXiv:1812.11132*, 2018.

565

566 Phipps Arabie and Scott A Boorman. Multidimensional scaling of measures of distance between  
 567 partitions. *Journal of Mathematical Psychology*, 10(2):148–203, 1973.

568

569 Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and  
 570 hierarchical image segmentation. *IEEE transactions on pattern analysis and machine intelligence*,  
 33(5):898–916, 2010.

571

572 Yupaporn Areepong and Wilasinee Peerajit. Integral equation solutions for the average run length for  
 573 monitoring shifts in the mean of a generalized seasonal ARFIMAX( $p, d, q, r$ )<sub>s</sub> process running on  
 574 a CUSUM control chart. *Plos one*, 17(2):e0264283, 2022.

575

576 Jushan Bai. Least absolute deviation estimation of a shift. *Econometric Theory*, 11(3):403–436,  
 1995.

577

578 Jushan Bai. Vector autoregressive models with structural changes in regression coefficients and in  
 579 variance-covariance matrices. Technical report, China Economics and Management Academy,  
 580 Central University of Finance and Economics, 2000.

581

582 Jan Beyersmann, Arthur Allignol, and Martin Schumacher. *Competing risks and multistate models  
 with R*. Springer Science & Business Media, 2011.

583

584 Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and  
 585 Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006. ISBN 0387310738.

586

587 Pinaki Biswas and John D Kalbfleisch. A risk-adjusted cusum in continuous time based on the cox  
 model. *Statistics in medicine*, 27(17):3382–3406, 2008.

588

589 {Connie M.} Borrer, J. {Bert Keats}, and Douglas Montgomery. Robustness of the time between  
 590 events CUSUM. *International Journal of Production Research*, 41(15):3435–3444, October 2003.  
 591 ISSN 0020-7543. doi: 10.1080/0020754031000138321.

592

593 Ikram Bouchikhi, André Ferrari, Cédric Richard, Anthony Bourrier, and Marc Bernot. Kernel based  
 online change point detection. In *2019 27th European Signal Processing Conference (EUSIPCO)*,  
 pp. 1–5. IEEE, 2019.

- 594 Taylor Bradley, Elie Alhajjar, and Nathaniel D. Bastian. Novelty detection in network traffic: Using  
595 survival analysis for feature identification. In *2023 IEEE International Conference on Assured  
596 Autonomy (ICAA)*, pp. 11–18, 2023. doi: 10.1109/ICAA58325.2023.00010.
- 597
- 598 Keri L Calkins, Chelsea E Canan, Richard D Moore, Catherine R Lesko, and Bryan Lau. An  
599 application of restricted mean survival time in a competing risks setting: comparing time to art  
600 initiation by injection drug use. *BMC medical research methodology*, 18(1):27, 2018.
- 601
- 602 D. R. Cox. Regression models and life-tables. *Journal of the Royal Statistical Society. Series B  
603 (Methodological)*, 34(2):187–220, 1972. ISSN 00359246. URL [http://www.jstor.org/  
604 stable/2985181](http://www.jstor.org/stable/2985181).
- 605 Gilles Crommen, Negera Wakgari Deresa, Myrthe D’Haen, Jie Ding, Ilias Willems, and Ingrid  
606 Van Keilegom. Recent advances in copula-based methods for dependent censoring. *SORT-Statistics  
607 and Operations Research Transactions*, pp. 3–42, 2025.
- 608
- 609 Stephen v Crowder. A simple method for studying run-length distributions of exponentially weighted  
610 moving average charts. *Technometrics*, 29(4):401–407, 1987.
- 611
- 612 Cameron Davidson-Pilon. lifelines: survival analysis in Python. *Journal of Open Source Software*, 4  
613 (40):1317, 2019.
- 614
- 615 Nathan Eagle and Alex Pentland. Reality mining: sensing complex social systems. *Personal and  
616 ubiquitous computing*, 10(4):255–268, 2006.
- 617
- 618 Nader Ebrahimi and Daniel Molefe. Survival function estimation when lifetime and censoring time  
619 are dependent. *Journal of multivariate analysis*, 87(1):101–132, 2003.
- 620
- 621 Bradley Efron and Robert J Tibshirani. *An introduction to the bootstrap*. Chapman and Hall/CRC,  
622 1994.
- 623
- 624 Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The  
625 PASCAL visual object classes (voc) challenge. *International journal of computer vision*, 88(2):  
626 303–338, 2010.
- 627
- 628 James C Fu, Fred A Spiring, and Hansheng Xie. On the average run lengths of quality control  
629 schemes using a Markov chain approach. *Statistics & Probability Letters*, 56(4):369–380, 2002.
- 630
- 631 A. GANDY, J. T. KVALØY, A. BOTTLE, and F. ZHOU. Risk-adjusted monitoring of time to event.  
632 *Biometrika*, 97(2):375–388, 2010. ISSN 00063444, 14643510. URL [http://www.jstor.  
633 org/stable/25734092](http://www.jstor.org/stable/25734092).
- 634
- 635 Kristine Gierz. *Inference and Estimation in Change Point Models for Censored Data*. PhD thesis,  
636 Old Dominion University, 2020.
- 637
- 638 Daniel Gomon, Hein Putter, Rob GHH Nelissen, and Stéphanie Van Der Pas. CGR-CUSUM: a  
639 continuous time generalized rapid response cumulative sum chart. *Biostatistics*, 25(1):253–269,  
640 2024.
- 641
- 642 Aditya Gopalan, Braghadeesh Lakshminarayanan, and Venkatesh Saligrama. Bandit quickest change-  
643 point detection. *Advances in Neural Information Processing Systems*, 34:29064–29073, 2021.
- 644
- 645 Olympia Hadjiliadis, Tobias Schaefer, and H Vincent Poor. Quickest detection in coupled systems.  
646 In *Proceedings of the 48th IEEE Conference on Decision and Control (CDC) held jointly with 2009  
647 28th Chinese Control Conference*, pp. 4723–4728. IEEE, 2009.
- 648
- 649 Abdul Haq and William H Woodall. A note on an average run length calculation for the EWMA and  
650 other charts. *Quality and Reliability Engineering International*, 38(8):4351–4355, 2022.
- 651
- 652 B C Haris, Gayadhar Pradhan, A Misra, SRM Prasanna, Rohan Kumar Das, and Rohit Sinha.  
653 Multivariability speaker recognition database in indian scenario. *International Journal of Speech  
654 Technology*, 15(4):441–453, 2012.

- 648 C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser,  
649 J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett,  
650 A. Haldane, J. F. Del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy,  
651 W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*,  
652 585(7825):357–362, 09 2020. License: BSD 3-Clause "New" or "Revised" License.
- 653 Douglas M Hawkins, Peihua Qiu, and Chang Wook Kang. The changepoint model for statistical  
654 process control. *Journal of quality technology*, 35(4):355–366, 2003.
- 655 F. C. Heilbron, V. Escorcia, B. Ghanem, and J. C. Niebles. Activitynet: A large-scale video benchmark  
656 for human activity understanding. In *2015 IEEE Conference on Computer Vision and Pattern  
657 Recognition (CVPR)*, pp. 961–970, June 2015. doi: 10.1109/CVPR.2015.7298698.
- 658 José Hernández-Orallo, Cèsar Ferri, Nicolas Lachiche, and Peter Flach. The 1st workshop on ROC  
659 analysis in artificial intelligence (ROCAI-2004). *ACM SIGKDD Explorations Newsletter*, 6(2):  
660 159–161, 2004.
- 661 Chiu-Hsieh Hsu and Jeremy MG Taylor. A robust weighted Kaplan-Meier approach for data with  
662 dependent censoring using linear combinations of prognostic covariates. *Statistics in medicine*, 29  
663 (21):2215–2223, 2010.
- 664 Yu-Han Huang and Venugopal V Veeravalli. High probability latency quickest change detection  
665 over a finite horizon. In *2024 IEEE International Symposium on Information Theory (ISIT)*, pp.  
666 1047–1052. IEEE, 2024.
- 667 Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*, 2(1):193–218,  
668 1985.
- 669 *Proceedings of the Second Workshop on ROC Analysis in Machine Learning (ROCML'05)*,  
670 <http://www.dsic.upv.es/flip/ROCML2005/>, 2005. International Conference on Machine Learning  
671 (ICML'05). URL <http://publis.icube.unistra.fr/index.php/11-LFMR05>.
- 672 *Proceedings of the Third Workshop on ROC Analysis in Machine Learning (ROCML'06)*,  
673 <http://www.dsic.upv.es/flip/ROCML2006/>, 2006. International Conference on Machine Learning  
674 (ICML'06). URL <http://publis.icube.unistra.fr/11-LFM06>.
- 675 Y.-G. Jiang, J. Liu, A. Roshan Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. THUMOS  
676 challenge: Action recognition with a large number of classes. [http://crcv.ucf.edu/  
677 THUMOS14/](http://crcv.ucf.edu/THUMOS14/), 2014.
- 678 Peter Johnson and Jesper Lund Pedersen. Bayesian changepoint detection for epidemic models.  
679 *Scientific Reports*, 15(1):20545, 2025.
- 680 Edward L Kaplan and Paul Meier. Nonparametric estimation from incomplete observations. *Journal  
681 of the American statistical association*, 53(282):457–481, 1958.
- 682 Yik Lun Kei, Jialiang Li, Hangjian Li, Yanzhen Chen, and OSCAR HERNAN MADRID PADILLA.  
683 Change point detection in dynamic graphs with decoder-only latent space model. *Transactions  
684 on Machine Learning Research*, 2025. ISSN 2835-8856. URL [https://openreview.net/  
685 forum?id=DVeFqV56Iz](https://openreview.net/forum?id=DVeFqV56Iz).
- 686 Victor Khamesi. ocpdet: A Python package for online changepoint detection in univariate and  
687 multivariate data, October 2022. URL <https://doi.org/10.5281/zenodo.7632721>.
- 688 JP Klein and ML Moeschberger. Asymptotic bias of the product limit estimator under dependent  
689 competing risks. *Indian Journal of Productivity, Reliability and Quality Control*, 9:1–7, 1984.
- 690 David G Kleinbaum and Mitchel Klein. *Survival analysis: A self-learning text*. Springer, 1996.
- 691 Bryan Klimt and Yiming Yang. Introducing the Enron Corpus. In *Proceedings of the First Conference  
692 on Email and Anti-Spam (CEAS)*, Mountain View, CA, 2004. URL [http://www.ceas.cc/  
693 papers-2004/168.pdf](http://www.ceas.cc/papers-2004/168.pdf).

- 702 Seven Knoth and Seven Knoth. Exact average run lengths of cusum schemes for erlang distributions:  
703 Exact average run lengths of cusum schemes. *Sequential Analysis*, 17(2):173–184, 1998.  
704
- 705 Jennifer R Kwapisz, Gary M Weiss, and Samuel A Moore. Activity recognition using cell phone  
706 accelerometers. *ACM SigKDD Explorations Newsletter*, 12(2):74–82, 2011.  
707
- 708 Marc Lavielle. Detection of multiple changes in a sequence of dependent variables. *Stochastic  
709 Processes and their applications*, 83(1):79–102, 1999.
- 710 Marc Lavielle and Gilles Teyssiere. Detection of multiple change-points in multivariate time series.  
711 *Lithuanian Mathematical Journal*, 46(3):287–306, 2006.  
712
- 713 Alexander Lavin and Subutai Ahmad. Evaluating real-time anomaly detection algorithms – the  
714 numenta anomaly benchmark. In *2015 IEEE 14th International Conference on Machine Learning  
715 and Applications (ICMLA)*, pp. 38–44, 2015. doi: 10.1109/ICMLA.2015.141.
- 716 Shuang Li, Yang Cao, Christina Leamon, Yao Xie, Lei Shi, and WenZhan Song. Online seismic  
717 event picking via sequential change-point detection. In *2016 54th Annual Allerton Conference on  
718 Communication, Control, and Computing (Allerton)*, pp. 774–779. IEEE, 2016.  
719
- 720 Zhonghua Li, Changliang Zou, Zhen Gong, and Zhaojun Wang. The computation of average run  
721 length and average time to signal: an overview. *Journal of Statistical Computation and Simulation*,  
722 84(8):1779–1802, 2014.
- 723 Johan Lim and Sungim Lee. Efficient ARL estimation for general control charts using censored run  
724 lengths. *Quality Engineering*, 37(3):359–368, 2025.  
725
- 726 Hung-Mo Lin, Sean TH Liu, Matthew A Levin, John Williamson, Nicole M Bouvier, Judith A  
727 Aberg, David Reich, and Natalia Egorova. Informative censoring—a cause of bias in estimating  
728 COVID-19 mortality using hospital data. *Life*, 13(1):210, 2023.
- 729 Gary Lorden. Procedures for reacting to a change in distribution. *The annals of mathematical  
730 statistics*, pp. 1897–1908, 1971.  
731
- 732 Ashutosh Makone. Twitter US Airline Sentiment, 2016. URL [https://www.kaggle.com/  
733 datasets/crowdfLOWER/twitter-airline-sentiment](https://www.kaggle.com/datasets/crowdfLOWER/twitter-airline-sentiment). Original data from Crowd-  
734 Flower’s “Data for Everyone” library; accessed 2025-11-24.
- 735 Sara Malmquist. What rates of censoring induce bias in the Kaplan–Meier estimator for small  
736 samples? Master’s thesis, Uppsala University, Department of Statistics, 2025.  
737
- 738 Yajun Mei. Is average run length to false alarm always an informative criterion? *Sequential Analysis*,  
739 27(4):354–376, 2008.
- 740 Melinda Mills. Competing risk and multi-state models. *Introducing Survival and Event History  
741 Analysis*, pp. 190–212, 2011.  
742
- 743 Jagabandhu Mishra and SR Mahadeva Prasanna. Spoken language change detection inspired by  
744 speaker change detection. *Circuits, Systems, and Signal Processing*, 43(10):6373–6398, 2024.  
745
- 746 Kojiro Morita. Introduction to survival analysis in the presence of competing risks. *Annals of Clinical  
747 Epidemiology*, 3(4):97–100, 2021.
- 748 S Mostafa Mousavi, Yixiao Sheng, Weiqiang Zhu, and Gregory C Beroza. Stanford earthquake  
749 dataset (STEAD): A global data set of seismic signals for AI. *IEEE Access*, 7:179464–179476,  
750 2019.  
751
- 752 Yue S Niu, Ning Hao, and Heping Zhang. Multiple change-point detection: A selective overview.  
753 *Statistical Science*, pp. 611–623, 2016.  
754
- 755 E. S. Page. Continuous inspection schemes. *Biometrika*, 41(1/2):100–115, 1954. ISSN 00063444.  
URL <http://www.jstor.org/stable/2333009>.

- 756 Juan Carlos Pardo-Fernández and Pablo Rodríguez Castro. International workshop on ROC anal-  
757 ysis and related topics (ROC2025): Programme and book of abstracts, January 2025. URL  
758 [https://roc2025.webs.uvigo.es/programme\\_and\\_BoA\\_ROC2025.pdf](https://roc2025.webs.uvigo.es/programme_and_BoA_ROC2025.pdf). Inter-  
759 national Workshop on ROC Analysis and Related Topics, Vigo, January 23–24, 2025.
- 760 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
761 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward  
762 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,  
763 Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep  
764 learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran  
765 Associates, Inc., 2019. License: <https://github.com/pytorch/pytorch/blob/master/LICENSE>.
- 766 Canan Pehlivan and Murat Caner Testik. Impact of model misspecification on the exponential  
767 EWMA charts: a robustness study when the time-between-events are not exponential. *Quality and*  
768 *Reliability Engineering International*, 26(2):177–190, 2010.
- 770 Ioannis Phinikettos and Axel Gandy. An omnibus CUSUM chart for monitoring time to event data.  
771 *Lifetime data analysis*, 20(3):481–494, 2014.
- 772 Moshe Pollak. Optimal detection of a change in distribution. *The Annals of Statistics*, pp. 206–227,  
773 1985.
- 775 Aleksey S. Polunchenko. Exact distribution of the generalized Shiryaev–Roberts stopping time  
776 under the minimax Brownian motion setup. *Sequential Analysis*, 35(1):108–143, 2016. doi: 10.  
777 1080/07474946.2016.1132066. URL [https://www.tandfonline.com/doi/abs/10.](https://www.tandfonline.com/doi/abs/10.1080/07474946.2016.1132066)  
778 [1080/07474946.2016.1132066](https://www.tandfonline.com/doi/abs/10.1080/07474946.2016.1132066).
- 779 H Vincent Poor. Quickest detection with exponential penalty for delay. *The Annals of Statistics*, 26  
780 (6):2179–2205, 1998.
- 782 Harsh Purohit, Ryo Tanabe, Kenji Ichige, Takashi Endo, Yuki Nikaido, Kaori Suefusa, and Yohei  
783 Kawaguchi. MIMII dataset: Sound dataset for malfunctioning industrial machine investigation  
784 and inspection. *arXiv preprint arXiv:1909.09347*, 2019.
- 785 Tony Qi and Jiuzhou Wang. Calculating restricted mean survival time. In *Phar-*  
786 *maSUG 2018*, 2018. URL [https://pharmasug.org/proceedings/2018/AA/](https://pharmasug.org/proceedings/2018/AA/PharmaSUG-2018-AA04.pdf)  
787 [PharmaSUG-2018-AA04.pdf](https://pharmasug.org/proceedings/2018/AA/PharmaSUG-2018-AA04.pdf).
- 788 Peihua Qiu. *Introduction to statistical process control*. CRC press, 2013.
- 790 R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical  
791 Computing, Vienna, Austria, 2021. URL <https://www.R-project.org/>.
- 792 Priya Ranganathan and CS Pramesh. Censoring in survival analysis: potential for bias. *Perspectives*  
793 *in clinical research*, 3(1):40, 2012.
- 795 Marion R Reynolds. Approximations to the average run length in cumulative sum control charts.  
796 *Technometrics*, 17(1):65–71, 1975.
- 797 Louis-Paul Rivest and Martin T Wells. A martingale approach to the copula-graphic estimator for the  
798 survival function under dependent censoring. *Journal of Multivariate Analysis*, 79(1):138–155,  
799 2001.
- 800 S. W. Roberts. Control chart tests based on geometric moving averages. *Technometrics*, 1(3):239–250,  
801 1959. ISSN 00401706. URL <http://www.jstor.org/stable/1266443>.
- 803 Gaetano Romano, Idris A. Eckley, and Paul Fearnhead. A log-linear nonparametric online changepoint  
804 detection algorithm based on functional pruning. *IEEE Transactions on Signal Processing*, 72:  
805 594–606, 2024a. doi: 10.1109/TSP.2023.3343550.
- 806 Gaetano Romano, Daniel Grose, Kes Ward, Austin Edward, Liudmila Pishchagina, Guillem Rigauil,  
807 Vincent Runge, Paul Fearnhead, and Idris A. Eckley. changepoint.online: A collection of methods  
808 for online changepoint detection, April 2024b. URL [https://github.com/](https://github.com/grosted/changepoint_online)  
809 [grosted/changepoint\\_online](https://github.com/grosted/changepoint_online).

- 810 Patrick Royston and Mahesh KB Parmar. Restricted mean survival time: an alternative to the hazard  
811 ratio for the design and analysis of randomized trials with a time-to-event outcome. *BMC medical*  
812 *research methodology*, 13(1):152, 2013.
- 813  
814 Neville Ryant, Kenneth Church, Christopher Cieri, Alejandrina Cristia, Jun Du, Sriram Ganapathy,  
815 and Mark Liberman. First DIHARD challenge evaluation plan. In *Technical Report*. Linguistic  
816 Data Consortium, University of Pennsylvania, 2018.
- 817 Michael S Saccucci and James M Lucas. Average run lengths for exponentially weighted moving  
818 average control schemes using the Markov chain approach. *Journal of Quality Technology*, 22(2):  
819 154–162, 1990.
- 820 Nassim Sahki, Anne Gégout-Petit, and Sophie Wantz-Mézières. Performance study of change-point  
821 detection thresholds for cumulative sum statistic in a sequential context. *Quality and Reliability*  
822 *Engineering International*, 36(8):2699–2719, 2020.
- 823  
824 Rebecca Salles, Janio Lima, Michel Reis, Rafaelli Coutinho, Esther Pacitti, Florent Masseglia, Reza  
825 Akbarinia, Chao Chen, Jonathan Garibaldi, Fabio Porto, et al. SoftED: Metrics for soft evaluation  
826 of time series event detection. *Computers & Industrial Engineering*, 198:110728, 2024.
- 827 R. Sasikumr and M. Sujatha. A comprehensive study on monitoring treatment outcomes using risk-  
828 adjusted CUSUM control charts. *Journal of Information Systems Engineering and Management*,  
829 10(35s):248–263, April 2025. ISSN 2468-4376. doi: 10.52783/jisem.v10i35s.5988. URL  
830 <https://jisem-journal.com/index.php/journal/article/view/5988>.
- 831  
832 Landon H Sego, Marion R Reynolds Jr, and William H Woodall. Risk-adjusted monitoring of survival  
833 times. *Statistics in medicine*, 28(9):1386–1401, 2009.
- 834 Mohammad Shafiq, Shuhrat Shah, and M Alamgir. Modified weighted Kaplan-Meier estimator.  
835 *Pakistan Journal of Statistics and Operation Research*, pp. 39–44, 2007.
- 836  
837 Albert N Shiryaev. On optimum methods in quickest detection problems. *Theory of Probability & Its*  
838 *Applications*, 8(1):22–46, 1963.
- 839 Salvatore Stolfo, Wei Fan, Wenke Lee, Andreas Prodromidis, and Philip Chan. KDD Cup 1999 Data.  
840 UCI Machine Learning Repository, 1999. DOI: <https://doi.org/10.24432/C51C7N>.
- 841  
842 Winfried Stute. The bias of Kaplan-Meier integrals. *Scandinavian Journal of Statistics*, pp. 475–484,  
843 1994.
- 844 Winfried Stute and J-L Wang. The strong law under random censorship. *The Annals of statistics*, pp.  
845 1591–1607, 1993.
- 846  
847 Rapin Sunthornwat, Saowanit Sukparungsee, and Yupaporn Areepong. Analytical explicit formulas  
848 of average run length of homogenously weighted moving average control chart based on a MAX  
849 process. *Symmetry*, 15(12):2112, 2023.
- 850 Eugene Tan, Shannon D Algar, Débora Corrêa, Thomas Stemler, and Michael Small. Network  
851 representations of attractors for change point detection. *Communications Physics*, 6(1):340, 2023.
- 852  
853 Alexander Tartakovsky. *Sequential change detection and hypothesis testing: General non-iid*  
854 *stochastic models and asymptotically optimal rules*. Chapman and Hall/CRC, 2019.
- 855 Alexander Tartakovsky, Igor Nikiforov, and Michele Basseville. *Sequential Analysis: Hypothesis*  
856 *Testing and Changepoint Detection*. Chapman & Hall/CRC, 1st edition, 2014.
- 857  
858 Nesime Tatbul, Tae Jun Lee, Stan Zdonik, Mejbah Alam, and Justin Gottschlich. Precision and  
859 recall for time series. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and  
860 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Asso-  
861 ciates, Inc., 2018. URL [https://proceedings.neurips.cc/paper\\_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2018/file/8f468c873a32bb0619eae2050ba45d1-Paper.pdf)  
862 [2018/file/8f468c873a32bb0619eae2050ba45d1-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/8f468c873a32bb0619eae2050ba45d1-Paper.pdf).
- 863  
864 Terry Therneau, Cynthia Crowson, and Elizabeth Atkinson. Multi-state models and competing risks.  
*CRAN-R* (<https://cran.r-project.org/web/packages/survival/vignettes/compete.pdf>), 2020.

864 Charles Truong, Laurent Oudre, and Nicolas Vayatis. Selective review of offline change point  
865 detection methods. *Signal Processing*, 167:107299, 2020.  
866

867 Gerhard Tutz, Matthias Schmid, et al. *Modeling discrete time-to-event data*, volume 3. Springer,  
868 2016.

869 G. J. J. Van den Burg and C. K. I. Williams. An evaluation of change point detection algorithms.  
870 *arXiv preprint arXiv:2003.06222*, 2020.  
871

872 Guido Van Rossum and Fred L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA,  
873 2009. ISBN 1441412697.

874 Marek Wadinger, Michal Kvasnica, and Yoshinobu Kawahara. Change-point detection in indus-  
875 trial data streams based on online dynamic mode decomposition with control. *arXiv preprint*  
876 *arXiv:2407.05976*, 2024.  
877

878 Binglu Wang, Yongqiang Zhao, Le Yang, Teng Long, and Xuelong Li. Temporal action localization in  
879 the deep learning era: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
880 46(4):2171–2190, 2023a.

881 Xiuheng Wang, Ricardo Augusto Borsoi, Cédric Richard, and Jie Chen. Change point detection with  
882 neural online density-ratio estimator. In *ICASSP 2023-2023 IEEE International Conference on*  
883 *Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023b.  
884

885 Jack Woollam, Jannes Münchmeyer, Frederik Tilmann, Andreas Rietbrock, Dietrich Lange, Thomas  
886 Bornstein, Tobias Diehl, Carlo Giunchi, Florian Haslinger, Dario Jozinović, et al. SeisBench—a  
887 toolbox for machine learning in seismology. *Seismological Society of America*, 93(3):1695–1709,  
888 2022.

889 Yijian Xiang, Murat Akcakaya, Satyabrata Sen, and Arye Nehorai. Target detection via cognitive  
890 radars using change-point detection, learning, and adaptation. *Circuits, Systems, and Signal*  
891 *Processing*, 40(1):233–261, 2021.

892 Ali Zare and Mahmood Mahmoodi. Modified Kaplan-Meier estimator based on competing risks for  
893 heavy censoring data. *Int J Statist Med Res*, 2(4):297–304, 2013.  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

APPENDIX

**Contents**

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Related Work</b>	<b>2</b>
<b>3</b>	<b>Kaplan-Meier ARL &amp; ADD</b>	<b>2</b>
<b>4</b>	<b>Bias Analysis</b>	<b>4</b>
<b>5</b>	<b>Experiment</b>	<b>6</b>
5.1	Estimation on Simulation Dataset . . . . .	7
5.2	Evaluation of Models on Real-world Dataset . . . . .	8
<b>6</b>	<b>Discussion, Limitations, and Future Work</b>	<b>9</b>
<b>A</b>	<b>QCD-Survival Analysis Correspondence</b>	<b>19</b>
A.1	Background of KME in Survival Analysis . . . . .	19
A.1.1	Introduction . . . . .	19
A.1.2	Decomposition via Conditional Survival Probabilities . . . . .	19
A.1.3	Risk Set $n_j$ and Number of Events $d_j$ . . . . .	21
A.1.4	Why $d_j/n_j$ ? Conditional Probability Viewpoint . . . . .	21
A.1.5	KME as Product-limit Estimator . . . . .	22
<b>B</b>	<b>Proofs</b>	<b>22</b>
B.1	Proof of Theorem 4.1 . . . . .	23
B.2	Proof of Theorem 4.2 . . . . .	25
B.3	Proof of Thm. 4.3 . . . . .	26
B.4	Proof of Theorem 4.4 . . . . .	29
<b>C</b>	<b>Supplementary Experimental Results</b>	<b>32</b>
C.1	Example of Survival Curve of Detection Points . . . . .	32
C.2	Complete Ablation of Fig. 4 . . . . .	32
C.3	Gaussian Process Dataset . . . . .	32
C.4	Poisson Process Dataset . . . . .	32
C.5	WISDM Actitracker Dataset . . . . .	34
C.6	ARL-ADD with CUSUM . . . . .	34
C.7	ARL-ADD with Geometric Changeoint Distribution . . . . .	34
<b>D</b>	<b>Detailed Experimental Settings</b>	<b>40</b>
D.1	Online QCD Models . . . . .	40
D.2	Preprocesses for WISDM Actitracker . . . . .	40
D.3	Statistics of WISDM Actitracker . . . . .	41
D.4	Numerical Computation of Variance . . . . .	41
<b>E</b>	<b>Supplementary Related Work</b>	<b>44</b>
<b>F</b>	<b>Supplementary Discussion</b>	<b>45</b>
F.1	More about Independent Censoring Assumption . . . . .	46
F.1.1	Background: Restricted Mean Survival Time (RMST) . . . . .	46
F.1.2	Solutions . . . . .	47
F.2	Relevance of Requiring Datasets with Multiple Sequences with Changeoint Labels . . . . .	47
<b>G</b>	<b>Use of Large Language Models</b>	<b>48</b>

## A QCD-SURVIVAL ANALYSIS CORRESPONDENCE

**Proposed Correspondence.** We summarize our key idea, an analogy between QCD and survival analysis in Tab. 1.

ARL	ADD	Survival analysis
detection time $\tau_i$	detection delay $\tau_i - \nu_i$	event time
$\min\{\nu_i, T_i\}$	$T_i - \nu_i$	censoring time $C_i$
KM-ARL $\hat{\mu}_T^{(KM)}$	KM-ADD $\hat{M}_T^{(KM)}$	restricted mean survival time

Table 1: QCD-survival analysis analogy.

**Illustration.** For clarity, we illustrate the original KME, KM-ARL, and the computation of  $d_j$  and  $n_j$  in Fig. 5.

### A.1 BACKGROUND OF KME IN SURVIVAL ANALYSIS

#### A.1.1 INTRODUCTION

We provide the motivation and intuition for the original KME in survival analysis Kaplan & Meier (1958); Kleinbaum & Klein (1996). Let  $\tau$  be a nonnegative random variable representing a lifetime, and let

$$S(t) := P(\tau > t) \quad (11)$$

denote the survival function. In the ideal (textbook) setting without censoring, and with i.i.d. samples  $\tau_1, \dots, \tau_N$ , a naive nonparametric estimator of  $S(t)$  is the following empirical survival function

$$\hat{S}_{\text{emp}}(t) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{\tau_i > t\}. \quad (12)$$

However, in practice, we often do *not* observe all  $\tau_i$  fully. Instead, some subjects drop out, are lost to follow-up, or the study ends before they fail. This leads to *right-censoring*, where we observe

$$Y_i = \min(\tau_i, C_i), \quad \Delta_i = \mathbb{1}\{\tau_i \leq C_i\}, \quad (13)$$

with  $C_i$  a censoring time. Then  $Y_i$  is either the failure time (if  $\Delta_i = 1$ ) or a censoring time (if  $\Delta_i = 0$ ), where  $\Delta_i$  is the event flag.

If we treat censored observations as if they were failures, we underestimate survival; if we drop them entirely, we discard information about the period in which we do know they survived. The KME is motivated precisely by this need:

- use all information available prior to censoring;
- do not make parametric assumptions (nonparametric);
- remain interpretable as a product of empirical conditional survival probabilities.

We derive the KME formula in the following to clarify its motivation.

#### A.1.2 DECOMPOSITION VIA CONDITIONAL SURVIVAL PROBABILITIES

A key idea is to write  $S(t)$  as a product of conditional survival probabilities at each distinct failure time (or event time, i.e., the time of death). Let

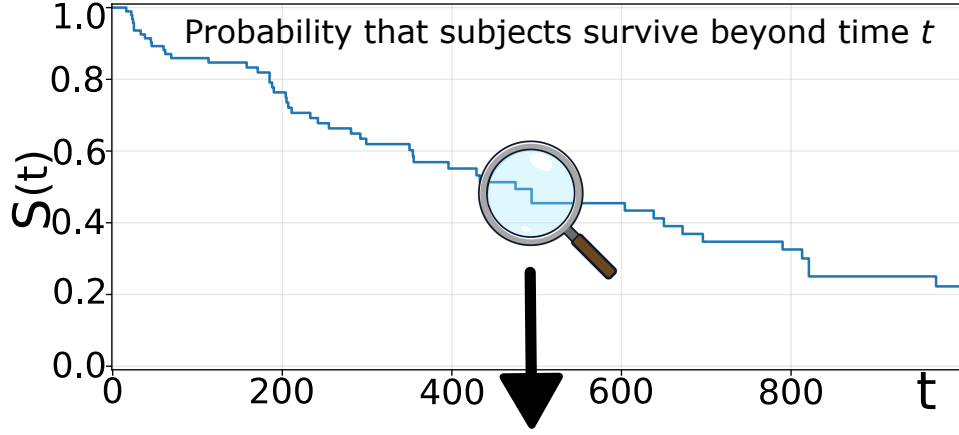
$$t_{(1)} < t_{(2)} < \dots < t_{(K)} \quad (14)$$

be the distinct *failure* times in the population (not yet the sample). For  $t$  between  $t_{(j)}$  and  $t_{(j+1)}$ , we can write

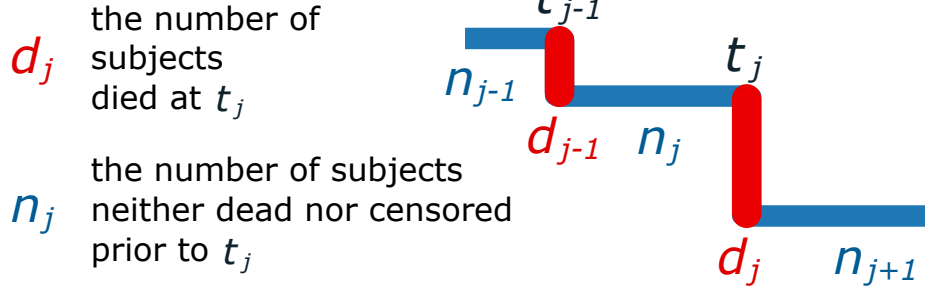
$$S(t) = P(\tau > t) = \prod_{t_{(m)} \leq t} P(\tau > t_{(m)} \mid \tau \geq t_{(m)}). \quad (15)$$

This is just repeated application of the chain rule of probability. Intuitively,

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079



## Survival Analysis



## QCD

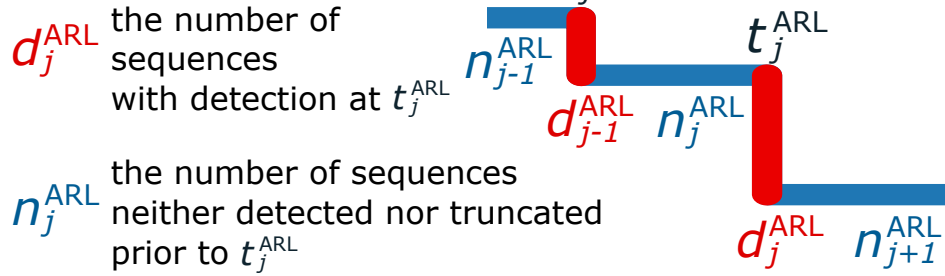


Figure 5: **Original KME, KM-ARL, and computation of  $d_j$  and  $n_j$ .** we estimate the *survival function of detection points*  $S^{\text{ARL}}(t) := P(\tau > t \mid \nu = \infty)$  using the Kaplan-Meier estimator (KME) (Kaplan & Meier, 1958), a non-parametric estimator of the survival function:  $\hat{S}^{\text{ARL}}(t) = \prod_{j: t_j^{\text{ARL}} \leq t} (1 - \frac{d_j^{\text{ARL}}}{n_j^{\text{ARL}}})$ , where  $0 < t_1^{\text{ARL}} < t_2^{\text{ARL}} < \dots < t_{N'}^{\text{ARL}}$  are distinct detection points, with  $N'_{\text{ARL}} \in \mathbb{N} (\leq N)$ ;  $d_j^{\text{ARL}} := |\{i \in [N] \mid \tau_i = t_j^{\text{ARL}}\}|$  is the number of sequences with a detection at  $t_j^{\text{ARL}}$  ( $j \in [N'_{\text{ARL}}]$ ); and  $n_j^{\text{ARL}} := |\{i \in [N] \mid \min\{\tau_i, C_i^{\text{ARL}}\} \geq t_j^{\text{ARL}}\}|$  is the number of sequences neither detected nor censored prior to  $t_j^{\text{ARL}}$  ( $j \in [N'_{\text{ARL}}]$ ), with the censoring time of sequence  $i$  defined as  $C_i^{\text{ARL}} := \min\{\nu_i, T_i\}$ . We propose a non-parametric estimator of the ARL under irregular sequence lengths, termed KM-ARL, as the integral of  $\hat{S}^{\text{ARL}}(t)$  over the range  $[0, a]$  for arbitrary  $a \in \mathbb{R}_{\geq 0}$ :  $\hat{\mu}_T^{(\text{KM})} := \int_0^a \hat{S}^{\text{ARL}}(t) dt$ .

- at each time  $t_{(m)}$ , we look at those who are still alive (have “survived so far”);
- we multiply by the probability that they survive *beyond*  $t_{(m)}$ ;
- the overall survival up to time  $t$  is the product of these step-by-step survival probabilities.

The KM estimator simply replaces each conditional probability  $P(\tau > t_{(m)} \mid \tau \geq t_{(m)})$  by a natural empirical estimate using the sample with censoring.

### A.1.3 RISK SET $n_j$ AND NUMBER OF EVENTS $d_j$

Suppose we have  $N$  subjects, and we observe  $(Y_i, \Delta_i)$  for  $i = 1, \dots, N$ . Let

- $t_{(1)} < \dots < t_{(J)}$  be the distinct observed *failure* times in the sample (i.e., times where at least one  $\Delta_i = 1$  and  $Y_i = t_{(j)}$ ),
- $d_j$  be the number of failures at time  $t_{(j)}$ ,
- $n_j$  be the size of the *risk set* at time  $t_{(j)}$ , i.e., the number of subjects who are known to be alive and uncensored just *before*  $t_{(j)}$ .

**Intuition for the risk set  $n_j$ .** The risk set  $n_j$  at time  $t_{(j)}$  collects all individuals for whom failure at  $t_{(j)}$  is still a possibility:

- subjects who failed earlier ( $Y_i < t_{(j)}$  and  $\Delta_i = 1$ ) are removed: they are already dead;
- subjects who were censored earlier ( $Y_i < t_{(j)}$  and  $\Delta_i = 0$ ) are removed: we no longer observe them, and we cannot know whether they would have failed at  $t_{(j)}$  or not;
- subjects whose observed time  $Y_i$  is at least  $t_{(j)}$  remain in the risk set: we know they survived at least up to just before  $t_{(j)}$ .

Thus,

$$n_j = \sum_{i=1}^N \mathbb{1}\{Y_i \geq t_{(j)}\}. \quad (16)$$

This definition uses all partial information: even if some  $Y_i$  correspond to future censoring events, until they are censored, they contribute to the information that the subject has survived up to that time.

**Intuition for the number of failures  $d_j$ .** At  $t_{(j)}$ , we also count the number of failures:

$$d_j = \sum_{i=1}^N \mathbb{1}\{Y_i = t_{(j)}, \Delta_i = 1\}. \quad (17)$$

Only *failures* contribute to  $d_j$ . Censoring at time  $t_{(j)}$  does not count as a failure; it simply removes subjects from future risk sets (for times after  $t_{(j)}$ ). This aligns with the goal of estimating the distribution of failure times, not censoring times.

### A.1.4 WHY $d_j/n_j$ ? CONDITIONAL PROBABILITY VIEWPOINT

Consider the underlying quantity

$$P(\tau = t_{(j)} \mid \tau \geq t_{(j)}) \quad (18)$$

as the probability that a subject who has survived up to  $t_{(j)}$  fails exactly at  $t_{(j)}$ . In a discrete-time picture, this is analogous to a “hazard probability” at  $t_{(j)}$ . Given  $n_j$  subjects in the risk set at  $t_{(j)}$ , and  $d_j$  of them failing at  $t_{(j)}$ , it is natural to estimate this conditional probability by the empirical proportion

$$\hat{P}(\tau = t_{(j)} \mid \tau \geq t_{(j)}) = \frac{d_j}{n_j}. \quad (19)$$

Consequently, the conditional survival probability at  $t_{(j)}$  is estimated by

$$\hat{P}(\tau > t_{(j)} \mid \tau \geq t_{(j)}) = 1 - \frac{d_j}{n_j}. \quad (20)$$

Note that censored individuals are *included* in  $n_j$  as long as they are not yet censored at  $t_{(j)}$ , but never in  $d_j$ . This reflects the idea that:

- up to the censoring time, their “survival experience” is fully observed and contributes to the risk set  $n_j$ ;
- at the censoring time, we stop knowing what happens afterward, so they drop from future risk sets (but are not treated as failures  $d_j$ ).

#### A.1.5 KME AS PRODUCT-LIMIT ESTIMATOR

Putting everything together, the KM estimator of  $S(t)$  is defined as

$$\hat{S}_{\text{KM}}(t) = \prod_{t_{(j)} \leq t} \left(1 - \frac{d_j}{n_j}\right). \quad (21)$$

This is exactly the product of estimated conditional survival probabilities:

$$\hat{S}_{\text{KM}}(t) = \prod_{t_{(j)} \leq t} \hat{P}(\tau > t_{(j)} \mid \tau \geq t_{(j)}). \quad (22)$$

Intuitively:

- At the beginning ( $t < t_{(1)}$ ), survival is 1: no one has failed.
- At each failure time  $t_{(j)}$ , we shrink the survival curve by multiplying by  $1 - d_j/n_j$ .
- Between failure times,  $\hat{S}_{\text{KM}}(t)$  is constant, yielding a right-continuous step function.

This construction uses only the ordering of failure and censoring times and assumes (in the standard theory) that censoring is non-informative (independent of the failure mechanism, given the observed history). Under these conditions, the KM estimator can be seen as the nonparametric maximum likelihood estimator of  $S(t)$  under right-censoring.

In summary, the roles of  $d_j$  and  $n_j$  are:

- $n_j$  is the size of the risk set just before  $t_{(j)}$ . It collects all individuals whose failure at  $t_{(j)}$  is still observable: they have not failed and have not been censored yet. This ensures that  $d_j/n_j$  is a valid empirical conditional probability.
- $d_j$  is the number of failures at  $t_{(j)}$ . It quantifies how much the survival curve should drop at that time, relative to the current risk set.
- The ratio  $d_j/n_j$  is the natural empirical estimator of the conditional failure probability at  $t_{(j)}$ . Its complement  $1 - d_j/n_j$  is the empirical conditional survival probability.
- The KM estimator is then the product over time of these conditional survival probabilities, giving a nonparametric estimator of the entire survival function that properly accounts for right-censoring.

## B PROOFS

Our proofs for the finite-sample bounds are based on (Stute, 1994). In our proofs for the finite-sample bounds, we idealize the time index as continuous without loss of generality, for technical convenience. Define  $\nu_1, \dots, \nu_N$  as i.i.d. random variables (the changepoints), independent of the observations of sequence, where  $N \in \mathbb{N}$ . Assume that  $\tau_1, \dots, \tau_N$  are independent and identically distributed (i.i.d.) random variables (the detection points of a QCD model). Assume that  $T_1, \dots, T_N$  are i.i.d. random variables (the lengths of sequences), independent of the observations of sequences.

## B.1 PROOF OF THEOREM 4.1

Let  $C_1^{\text{ARL}}, \dots, C_N^{\text{ARL}}$  are i.i.d. random variables (the censoring times for the KM-ARL), defined as  $C_i^{\text{ARL}} := \min\{\nu_i, T_i\}$  for  $i \in [N]$ . Assume that  $C^{\text{ARL}}$  is independent of  $\tau$ , called the independent censoring assumption (or non-informative censoring assumption) for the KM-ARL. This assumption holds for online QCD models that do not look ahead the input sequence, which are the focus of this paper. This is because: detection points  $\tau_i$  prior to censoring can be regarded as samples from  $P(\tau \mid \nu = \infty)$ ;  $P(\tau \mid \nu = \infty)$ ,  $P(\nu)$ , and  $P(T)$  are independent; and thus,  $P(\tau \mid \nu = \infty)$  and  $P(C^{\text{ARL}} := \min\{\nu, T\})$  are also independent. Here,  $\mathbb{P}$  denotes a probability or a probability measure, with slight abuse of notation throughout the paper. In contrast, the independent censoring assumption does not hold for offline changepoint detection models because  $\tau$  depends on  $X^{(0,T)}$ , which in turn depends on  $\nu$  and  $T$ . Let  $F^{\text{ARL}}$  and  $G^{\text{ARL}}$  denote the cumulative distribution functions (CDFs) of  $\tau$  and  $C^{\text{ARL}}$ , respectively. The survival function of detection points is defined as  $S^{\text{ARL}}(t) := 1 - F^{\text{ARL}}(t)$ , where  $t \in \mathbb{R}_{\geq 0}$ .

We adopt a different convention of the KME from that in the main text. Consider the random variables

$$Z_i^{\text{ARL}} := \min(\tau_i, C_i^{\text{ARL}}), \quad (23)$$

$$\delta_i^{\text{ARL}} := \mathbb{1}_{\{\tau_i < C_i^{\text{ARL}}\}}, \quad (24)$$

where  $\delta_i^{\text{ARL}}$  is the indicator representing whether detection  $\tau_i$  has been observed before a changepoint or the end of the sequence. Let  $H^{\text{ARL}}$  denote the CDF of  $Z^{\text{ARL}}$ , and assume that it is continuous on  $t \in \mathbb{R}_{\geq 0}$ . Let  $Z_{1:N}^{\text{ARL}} \leq Z_{2:N}^{\text{ARL}} \leq \dots \leq Z_{N:N}^{\text{ARL}}$  be the *order statistics* of  $Z^{\text{ARL}}$  (i.e.,  $Z_{i:N}^{\text{ARL}}$  are sorted in increasing order) and  $\delta_{[i:N]}^{\text{ARL}}$  be the *concomitant* of  $Z_{i:N}^{\text{ARL}}$ ; i.e.,  $\delta_{i:N}^{\text{ARL}}$  are sorted according to  $Z_{i:N}^{\text{ARL}}$ , meaning that  $\delta_{[i:N]}^{\text{ARL}} = \delta_j$  if  $Z_{i:N}^{\text{ARL}} = Z_j^{\text{ARL}}$ . The estimation of  $S^{\text{ARL}}(t)$  is given by

$$1 - \hat{F}_N^{\text{ARL}}(t) := \hat{S}^{\text{ARL}}(t) = \prod_{i=1}^N \left(1 - \frac{\delta_{[i:N]}^{\text{ARL}}}{N - i + 1}\right) \mathbb{1}_{\{Z_{i:N}^{\text{ARL}} \leq t\}}, \quad (25)$$

where  $t \in \mathbb{R}_{\geq 0}$ . This is equivalent to the convention in the main text, which can be confirmed by canceling factors telescopically. In case of ties, treat a death as if it occurs before a censoring at the same time point. Our KM-ARL is formally defined as

$$\hat{\mu}_T^{(\text{KM})} = \int_0^a t d\hat{F}_N^{\text{ARL}}, \quad (26)$$

where  $a \in \mathbb{R}_{\geq 0}$  is arbitrary.

**Theorem B.1** (Finite-sample bias bounds for KM-ARL (Thm. 4.1)). *We idealize the time index as continuous without loss of generality, for technical convenience. Let  $F^{\text{ARL}}$ ,  $G^{\text{ARL}}$ , and  $H^{\text{ARL}}$  be the CDFs of  $\tau$ ,  $C^{\text{ARL}}$ , and  $Z^{\text{ARL}}$ , respectively. Assume that (i)  $F^{\text{ARL}}$  and  $G^{\text{ARL}}$  do not have common discontinuities, (ii)  $\tau$  and  $C^{\text{ARL}}$  are independent, known as the independent censoring (or non-informative censoring), and (iii)  $H^{\text{ARL}}$  is continuous. Then, we have for any  $a \in \mathbb{R}_{\geq 0}$*

$$- \int_0^a t G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} F^{\text{ARL}}(dt) \quad (27)$$

$$\leq \mathcal{B}_{\text{FS}}(\hat{\mu}_T^{(\text{KM})}) \leq \quad (28)$$

$$\int_0^a a G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} F^{\text{ARL}}(dt), \quad (29)$$

Assumption (i) is a technical requirement and is valid unless pathological situations are considered. In most cases, it holds with probability 1 because we consider continuous time. Note that  $F^{\text{ARL}}$  and  $G^{\text{ARL}}$  may have separate discontinuities. Assumption (ii) holds for online QCD models that do not look ahead the input sequence, which are the focus of this paper. This is because: detection points  $\tau_i$  prior to censoring can be regarded as samples from  $P(\tau \mid \nu = \infty)$ ;  $P(\tau \mid \nu = \infty)$ ,  $P(\nu)$ , and  $P(T)$  are independent; and thus,  $P(\tau \mid \nu = \infty)$  and  $P(C := \min\{\nu, T\})$  are also independent. In contrast, Assumption (ii) does not hold for offline changepoint detection models because  $\tau$  depends on  $X^{(0,T)}$ , which in turn depends on  $\nu$  and  $T$ . Assumption (iii) is also a technical requirement.

1242 *Proof.*

$$1243 \int_0^a t dF^{\text{ARL}}(t)$$

$$1244 = \int_0^a \left( \int_0^a \mathbb{1}(t > s) ds \right) dF^{\text{ARL}}(t) \quad (30)$$

$$1245$$

$$1246 = \int_0^a \left( \int_0^a \mathbb{1}(t > s) dF^{\text{ARL}}(t) \right) ds \quad (31)$$

$$1247$$

$$1248 = \int_0^a \left( \int_s^a dF^{\text{ARL}}(t) \right) ds \quad (32)$$

$$1249$$

$$1250 = \int_0^a (F^{\text{ARL}}(a) - F^{\text{ARL}}(s)) ds \quad (33)$$

$$1251$$

$$1252 = aF^{\text{ARL}}(a) - a + \int_0^a S^{\text{ARL}}(t) dt. \quad (34)$$

1253 Similarly, we have

$$1254 \int_0^a t d\hat{F}^{\text{ARL}}(t) = a\hat{F}^{\text{ARL}}(a) - a + \int_0^a \hat{S}^{\text{ARL}}(t) dt. \quad (35)$$

1255 Thus, the finite-sample bias is given by

$$1256 B_{\text{FS}}(\hat{\mu}_T^{(\text{KM})})$$

$$1257 = \mathbb{E} \left[ \int_0^a \hat{S}^{\text{ARL}}(t) dt \right] - \int_0^a S^{\text{ARL}}(t) dt$$

$$1258 = \mathbb{E} \left[ \int_0^a t d\hat{F}^{\text{ARL}}(t) \right] - \int_0^a t dF^{\text{ARL}}(t)$$

$$1259 - a(\mathbb{E}[\hat{F}^{\text{ARL}}(a)] - F^{\text{ARL}}(a))$$

$$1260 =: I_1 - aI_2, \quad (36)$$

1261 where we defined

$$1262 I_1 = \mathbb{E} \left[ \int_0^a t d\hat{F}^{\text{ARL}}(t) \right] - \int_0^a t dF^{\text{ARL}}(t) \quad (37)$$

$$1263 I_2 = \mathbb{E}[\hat{F}^{\text{ARL}}(a)] - F^{\text{ARL}}(a). \quad (38)$$

1264 We bound them by invoking the following lemma from (Stute, 1994), adapted to our setting. The full proof of Lem. B.2 is lengthy but given in (Stute, 1994). Refer also to (Stute & Wang, 1993), on which Stute (1994) builds his result.

1265 **Lemma B.2** (Cor. 1.1 in (Stute, 1994)). *For any  $F^{\text{ARL}}$ ,  $G^{\text{ARL}}$ ,  $H^{\text{ARL}}$ , and  $\hat{F}^{\text{ARL}}$  that satisfy the assumptions described in this section, the following inequality holds for any  $N \in \mathbb{N}$ , where  $\varphi(t) \geq 0$  is a Borel measurable function on  $t \in \mathbb{R}_{\geq 0}$ :*

$$1266 - \int \varphi(t) G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t)$$

$$1267 \leq \mathbb{E} \left[ \int \varphi(t) d\hat{F}^{\text{ARL}}(t) \right] - \int \varphi(t) dF^{\text{ARL}}$$

$$1268 \leq 0. \quad (39)$$

1269 According to Lem. B.2, for  $\varphi(t) = \mathbb{1}(t \in [0, a])$ , we have

$$1270 - \int_0^a G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t)$$

$$1271 \leq \mathbb{E} \left[ \int_0^a d\hat{F}^{\text{ARL}}(t) \right] - \int_0^a dF^{\text{ARL}} \leq 0. \quad (40)$$

Since  $\mathbb{E}[\int_0^a d\hat{F}^{\text{ARL}}(t)] - \int_0^a dF^{\text{ARL}} = \mathbb{E}[\hat{F}^{\text{ARL}}(a)] - F^{\text{ARL}}(a) = I_2$ , we have

$$- \int_0^a G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t) \leq I_2 \leq 0. \quad (41)$$

Additionally, for  $\varphi(t) = t \mathbf{1}(t \in [0, 1])$ , we have

$$\begin{aligned} & - \int_0^a t G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t) \\ & \leq \mathbb{E}[\int_0^a t d\hat{F}^{\text{ARL}}(t)] - \int_0^a t dF^{\text{ARL}}(t) (= I_1) \leq 0. \end{aligned} \quad (42)$$

From Eqs. (36), (37), (38), (41), and (42), we have proved

$$\begin{aligned} & - \int_0^a t G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t) \\ & \leq I_1 - aI_2 (= B_{\text{FS}}(\hat{\mu}_T^{(\text{KM})})) \leq \\ & \int_0^a a G^{\text{ARL}}(t) H^{\text{ARL}}(t)^{N-1} dF^{\text{ARL}}(t) \end{aligned} \quad (43)$$

for arbitrary  $a \in \mathbb{R}_{\geq 0}$ .  $\square$

## B.2 PROOF OF THEOREM 4.2

Consider only sequences with  $\nu < \infty$  and  $\Delta\tau \geq 0$  only, as others do not contribute to the computation of the ADD, defined as  $\mathbb{E}[\tau - \nu \mid \tau \geq \nu, \nu < \infty]$ . In this section, we use  $N \in \mathbb{N}$  as the number of sequences with  $\nu_i < \infty$  and  $\Delta\tau_i \geq 0$  in the dataset. Assume that  $\Delta\tau_1, \dots, \Delta\tau_N$  are independent and identically distributed (i.i.d.) random variables (the detection delays of a QCD model), defined as  $\Delta\tau_i := \tau_i - \nu_i$ . Let  $C_1^{\text{ADD}}, \dots, C_N^{\text{ADD}}$  are i.i.d. random variables (the censoring times for the KM-ADD), defined as  $C_i^{\text{ADD}} := \Delta T_i = T_i - \nu_i$  for  $i \in [N]$ . Assume that  $C^{\text{ADD}}$  is independent of  $\Delta\tau$ , called the independent censoring assumption (or non-informative censoring assumption) for the KM-ADD. The independent censoring assumption is justified in online QCD models by the approximate  $P(\Delta\tau \mid \Delta\tau \geq 0, \nu < \infty) \approx P(\tau \mid \nu = 0)$  and  $P(C := \Delta T \mid \Delta\tau \geq 0, \nu < \infty) \approx P(C := T \mid \nu = 0)$  because  $P(\tau \mid \nu = 0)$  and  $P(C := T \mid \nu = 0)$  are independent. This indicates that the distributions of detection delay  $\tau - \nu$  and censoring time  $T - \nu$  are approximately equal to the distributions of them measured from  $t = 0$ , not  $\nu$ . Under this approximation, the definition of the ADD becomes  $\mathbb{E}[\tau \mid \nu = 0]$ , which is called the steady-state ARL (Saccucci & Lucas, 1990; Sasikumr & Sujatha, 2025; Lim & Lee, 2025) and often used in the theory of control charts. For simplicity, we assume the independent censoring in the proof of Thm. 4.2, and define  $F^{\text{ADD}}, G^{\text{ADD}}$ , and  $H^{\text{ADD}}$  accordingly. The survival function of detection delays is defined as  $S^{\text{ADD}}(t) := 1 - F^{\text{ADD}}(t)$ , where  $t \in \mathbb{R}_{\geq 0}$ .

Consider the random variables

$$Z_i^{\text{ADD}} := \min(\Delta\tau_i, C_i^{\text{ADD}}), \quad (44)$$

$$\delta_i^{\text{ADD}} := \mathbf{1}_{\{\Delta\tau_i < C_i^{\text{ADD}}\}}, \quad (45)$$

where  $\delta_i^{\text{ADD}}$  is the indicator representing whether the detection  $\tau_i$  has been observed before the end of the sequence. Let  $H^{\text{ADD}}$  denote the CDF of  $Z^{\text{ADD}}$ , and assume that it is continuous on  $t \in \mathbb{R}_{\geq 0}$ . Let  $Z_{1:N}^{\text{ADD}} \leq Z_{2:N}^{\text{ADD}} \leq \dots \leq Z_{N:N}^{\text{ADD}}$  be the *order statistics* of  $Z^{\text{ADD}}$  (i.e.,  $Z_i^{\text{ADD}}$  are sorted in increasing order) and  $\delta_{[i:N]}^{\text{ADD}}$  be the *concomitant* of  $Z_{i:N}^{\text{ADD}}$ ; i.e.,  $\delta_i^{\text{ADD}}$  are sorted according to  $Z_{i:N}^{\text{ADD}}$ , meaning that  $\delta_{[i:N]}^{\text{ADD}} = \delta_j$  if  $Z_{i:N}^{\text{ADD}} = Z_j^{\text{ADD}}$ . The *survival function of detection points* is given by

$$1 - \hat{F}_N^{\text{ADD}}(t) := \hat{S}^{\text{ADD}}(t) = \prod_{i=1}^N \left(1 - \frac{\delta_{[i:N]}^{\text{ADD}}}{N - i + 1}\right) \mathbf{1}_{(Z_{i:N}^{\text{ADD}} \leq t)}, \quad (46)$$

where  $t \in \mathbb{R}_{\geq 0}$ . This is equivalent to the convention in the main text, which can be confirmed by canceling factors telescopically. In case of ties, treat a death as if it occurs before a censoring at the same time point. Our KM-ARL is formally defined as

$$\hat{M}_T^{(\text{KM})} = \int_0^b t d\hat{F}_N^{\text{ADD}}, \quad (47)$$

where  $a \in \mathbb{R}_{\geq 0}$  is arbitrary.

**Theorem B.3** (Finite-sample bias bounds for KM-ADD (Thm. 4.2)). *We idealize the time index as continuous without loss of generality, for technical convenience. Consider only sequences with  $\nu_i < \infty$  and  $\Delta\tau_i \geq 0$  in the dataset. Let  $F^{\text{ADD}}$ ,  $G^{\text{ADD}}$ , and  $H^{\text{ADD}}$  be the CDFs of  $\Delta\tau$ ,  $C^{\text{ADD}}$ , and  $Z^{\text{ADD}}$ , respectively. Assume that (i)  $F^{\text{ADD}}$  and  $G^{\text{ADD}}$  do not have common discontinuities, (ii)  $\Delta\tau$  and  $C^{\text{ADD}}$  are independent, known as the independent censoring (or non-informative censoring), and (iii)  $H^{\text{ADD}}$  is continuous. Then, we have for any  $b \in \mathbb{R}_{\geq 0}$*

$$- \int_0^b t G^{\text{ADD}}(t) H^{\text{ADD}}(t)^{N-1} F^{\text{ADD}}(dt) \quad (48)$$

$$\leq \mathcal{B}_{\text{FS}}(\hat{M}_T^{(\text{KM})}) \leq \quad (49)$$

$$\int_0^b b G^{\text{ADD}}(t) H^{\text{ADD}}(t)^{N-1} F^{\text{ADD}}(dt), \quad (50)$$

Assumption (i) is a technical requirement and is valid unless pathological situations are considered. In most cases, it holds with probability 1 because we consider continuous time. Note that  $F^{\text{ADD}}$  and  $G^{\text{ADD}}$  may have separate discontinuities. Assumption (ii) is the independent censoring assumption discussed above. Assumption (iii) is also a technical requirement.

*Proof.* Thanks to the setup detailed in this section, the proof is identical to that of Thm. 4.1 with relevant replacements, such as  $(\cdot)^{\text{ARL}}$  with  $(\cdot)^{\text{ADD}}$ .  $\square$

### B.3 PROOF OF THM. 4.3

**Theorem B.4** (Truncation bias bounds for KM-ARL (Thm. 4.3)). *Define the truncation bias of the LB-ARL as  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{LB})}) := \mu_T^{(\text{LB})} - \mu_\infty$ , where  $\mu_T^{(\text{LB})} := \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T]$  is the true ARL under random sequence lengths. For  $a = T_{\text{max}}^*$ , we have  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{LB})}) \leq \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})}) \leq 0$ .*

*Proof.* Recall that

$$\begin{aligned} \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})}) &= \mu_T^{(\text{KM})} - \mu_\infty \\ &= \int_0^{T_{\text{max}}^*} S^{\text{ARL}}(t) dt - \mathbb{E}[\tau \mid \nu = \infty] \end{aligned} \quad (51)$$

$$\begin{aligned} \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{LB})}) &= \mu_T^{(\text{LB})} - \mu_\infty \\ &= \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T] - \mathbb{E}[\tau \mid \nu = \infty]. \end{aligned} \quad (52)$$

Since

$$\int_0^\infty S^{\text{ARL}}(t) dt \quad (53)$$

$$= \int_0^\infty P(\tau > t \mid \nu = \infty) dt \quad (54)$$

$$= \int_0^\infty \mathbb{E}[\mathbf{1}(\tau > t) \mid \nu = \infty] dt \quad (55)$$

$$= \mathbb{E}\left[\int_0^\infty \mathbf{1}(\tau > t) dt \mid \nu = \infty\right] \quad (56)$$

$$= \mathbb{E}[\tau \mid \nu = \infty], \quad (57)$$

we have  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})}) = -\int_{T_{\max}^*}^{\infty} S^{\text{ARL}}(t) dt \leq 0$ . Below, we show that  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{LB})}) \leq \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})})$ .  
Since

$$\begin{aligned} & \mu_T^{(\text{KM})} \\ &= \int_0^{T_{\max}^*} S^{\text{ARL}}(t) dt \end{aligned} \quad (58)$$

$$= \int_0^{T_{\max}^*} P(\tau > t \mid \nu = \infty) dt \quad (59)$$

$$= \int_0^{T_{\max}^*} \mathbb{E}[\mathbb{1}(\tau > t) \mid \nu = \infty] dt \quad (60)$$

$$= \mathbb{E}\left[\int_0^{T_{\max}^*} dt \mid \nu = \infty\right] \quad (61)$$

$$\begin{aligned} &= P(\tau \leq T_{\max}^* \mid \nu = \infty) \\ &\quad \times \mathbb{E}\left[\int_0^{T_{\max}^*} \mathbb{1}(\tau > t) dt \mid \nu = \infty, \tau \leq T_{\max}^*\right] \\ &\quad + P(\tau > T_{\max}^* \mid \nu = \infty) \\ &\quad \times \mathbb{E}\left[\int_0^{T_{\max}^*} \mathbb{1}(\tau > t) dt \mid \nu = \infty, \tau > T_{\max}^*\right] \end{aligned} \quad (62)$$

$$\begin{aligned} &= P(\tau \leq T_{\max}^* \mid \nu = \infty) \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*] \\ &\quad + P(\tau > T_{\max}^* \mid \nu = \infty) \mathbb{E}[T_{\max}^* \mid \nu = \infty, \tau > T_{\max}^*] \end{aligned} \quad (63)$$

$$\begin{aligned} &= P(\tau \leq T_{\max}^* \mid \nu = \infty) \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*] \\ &\quad + P(\tau > T_{\max}^* \mid \nu = \infty) T_{\max}^* \end{aligned} \quad (64)$$

$$\begin{aligned} &= (1 - P(\tau > T_{\max}^* \mid \nu = \infty)) \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*] \\ &\quad + P(\tau > T_{\max}^* \mid \nu = \infty) T_{\max}^* \end{aligned} \quad (65)$$

$$\begin{aligned} &= \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*] \\ &\quad + P(\tau > T_{\max}^* \mid \nu = \infty) \\ &\quad \times (T_{\max}^* - \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*]) \end{aligned} \quad (66)$$

we have

$$\begin{aligned} & \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})}) - \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{LB})}) \\ &= \mu_T^{(\text{KM})} - \mu_T^{(\text{LB})} \end{aligned} \quad (67)$$

$$\begin{aligned} &= \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*] \\ &\quad + P(\tau > T_{\max}^* \mid \nu = \infty) \\ &\quad \times (T_{\max}^* - \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*]) \\ &\quad - \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T] \end{aligned} \quad (68)$$

$$\begin{aligned} &= \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*] - \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T] \\ &\quad + P(\tau > T_{\max}^* \mid \nu = \infty) \\ &\quad \times (T_{\max}^* - \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*]). \end{aligned} \quad (69)$$

Finally, we can see that  $\mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*] - \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T] \geq 0$  (Lem. B.5), and  $T_{\max}^* - \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*] \geq 0$ ; hence,  $\mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{KM})}) - \mathcal{B}_{\text{TR}}(\hat{\mu}_T^{(\text{LB})}) \geq 0$ .  $\square$

**Lemma B.5.**  $\mathbb{E}[\tau \mid \nu = \infty, \tau \leq T_{\max}^*] - \mathbb{E}[\tau \mid \nu = \infty, \tau \leq T] \geq 0$ .

*Proof.* Without loss of generality, this inequality is equivalent to  $\mathbb{E}[\tau \mid 0 \leq \tau \leq T_{\max}^*] - \mathbb{E}[\tau \mid 0 \leq \tau \leq T] \geq 0$ , where  $T \leq T_{\max}^*$  ( $\in \mathbb{R}$ ) holds a.s., and  $\tau$  and  $T$  are independent. For notational simplicity, we will show that

$$\mathbb{E}_{X,Y}[X \mid X \leq y_{\max}] - \mathbb{E}_{X,Y}[X \mid X \leq Y] \geq 0, \quad (70)$$

where  $X$  is a non-negative random variable,  $Y$  is a non-negative random variable with  $0 \leq Y \leq y_{\max}$  a.s.,  $y_{\max} (\geq 0)$  is a real constant, and  $X$  and  $Y$  are independent.

Write  $F_X$  for the CDF of  $X$ , and let

$$g(x) := P(Y \geq x), \quad 0 \leq x \leq y_{\max}. \quad (71)$$

Because  $Y \in [0, y_{\max}]$ ,  $g$  is non-increasing in  $x$ , while  $g(x) = 0$  for  $x > y_{\max}$ . Thus, we have

$$P(X \leq Y) = \mathbb{E}[\mathbb{1}(X \leq Y)] = \int_0^\infty P(Y \geq x) dF_X(x) = \int_0^{y_{\max}} g(x) dF_X(x), \quad (72)$$

where we used the independence of  $X$  and  $Y$ . Additionally, we have

$$\mathbb{E}[X \mathbb{1}(X \leq Y)] = \int_0^\infty x P(Y \geq x) dF_X(x) = \int_0^{y_{\max}} x g(x) dF_X(x). \quad (73)$$

Therefore,

$$\mathbb{E}[X | X \leq Y] = \frac{\mathbb{E}[X \mathbb{1}(X \leq Y)]}{P(X \leq Y)} = \frac{\int_0^{y_{\max}} x g(x) dF_X(x)}{\int_0^{y_{\max}} g(x) dF_X(x)}. \quad (74)$$

On the other hand, we can similarly derive

$$\mathbb{E}[X | X \leq y_{\max}] = \frac{\int_0^{y_{\max}} x dF_X(x)}{\int_0^{y_{\max}} dF_X(x)}. \quad (75)$$

Next, let  $\nu$  be the finite measure on  $[0, y_{\max}]$  given by  $d\nu(x) = dF_X(x)$ , and define another measure  $\mu$  by

$$d\mu(x) = g(x) d\nu(x), \quad 0 \leq x \leq y_{\max}. \quad (76)$$

Note that  $g(x)$  is non-increasing in  $x$  and that  $x$  is increasing in  $x$ . Thus, Eqs. (74) and (75) can be rewritten as

$$\mathbb{E}[X | X \leq Y] = \frac{\int x d\mu(x)}{\int d\mu(x)}, \quad \mathbb{E}[X | X \leq y_{\max}] = \frac{\int x d\nu(x)}{\int d\nu(x)}. \quad (77)$$

Consider the covariance under the normalized measure proportional to  $\nu$ :

$$\text{Cov}_\nu(x, g(x)) = \mathbb{E}_\nu[Xg(X)] - \mathbb{E}_\nu[X] \mathbb{E}_\nu[g(X)] \leq 0, \quad (78)$$

because  $x$  is increasing and  $g(x)$  is decreasing (Lem. B.6). Therefore,

$$\left( \frac{\int x g(x) d\nu(x)}{\int d\nu(x)} \right) \leq \left( \frac{\int x d\nu(x)}{\int d\nu(x)} \right) \left( \frac{\int g(x) d\nu(x)}{\int d\nu(x)} \right). \quad (79)$$

This is equivalent to

$$\frac{\int x g(x) d\nu(x)}{\int g(x) d\nu(x)} \leq \frac{\int x d\nu(x)}{\int d\nu(x)}. \quad (80)$$

In other words,

$$\mathbb{E}[X | X \leq Y] \leq \mathbb{E}[X | X \leq y_{\max}]. \quad (81)$$

This concludes the proof.  $\square$

**Lemma B.6.** *Let  $X$  be a real random variable,  $f$  an increasing functions, and  $h$  and non-increasing function. Then, it holds that  $\text{Cov}(f(X), h(X)) \leq 0$ .*

*Proof.* Take two i.i.d. copies,  $X$  and  $X'$ , with the same distribution. Let us consider the random quantity

$$(f(X) - f(X'))(h(X) - h(X')). \quad (82)$$

Expanding its expectation, we have

$$\begin{aligned} \mathbb{E}[(f(X) - f(X'))(h(X) - h(X'))] &= \mathbb{E}[f(X)h(X)] + \mathbb{E}[f(X')h(X')] \\ &\quad - \mathbb{E}[f(X)h(X')] - \mathbb{E}[f(X')h(X)] \end{aligned} \quad (83)$$

$$= 2\mathbb{E}[f(X)h(X)] - 2\mathbb{E}[f(X)]\mathbb{E}[h(X)] \quad (84)$$

$$= 2\text{Cov}(f(X), h(X)). \quad (85)$$

Thus, we have the identity

$$\text{Cov}(f(X), h(X)) = \frac{1}{2} \mathbb{E}[(f(X) - f(X'))(h(X) - h(X'))]. \quad (86)$$

Now use the fact that:  $f$  is increasing, and  $h$  is non-increasing. Fix any  $x, y$  in the support. Then:

- 1512 • If  $x > y$ , then  $f(x) > f(y)$  and  $h(x) \leq h(y)$ , so  $(f(x) - f(y)) > 0$  and  $(h(x) - h(y)) \leq 0$ ,  
 1513 hence  $(f(x) - f(y))(h(x) - h(y)) \leq 0$ .  
 1514  
 1515 • If  $x < y$ , then  $f(x) < f(y)$  and  $h(x) \geq h(y)$ , so  $(f(x) - f(y)) < 0$  and  $(h(x) - h(y)) \geq 0$ ,  
 1516 again  $(f(x) - f(y))(h(x) - h(y)) \leq 0$ .  
 1517  
 1518 • If  $x = y$ , the product is 0.  
 1519

1520 Thus, we have

$$1521 (f(x) - f(y))(h(x) - h(y)) \leq 0 \text{ for all } x, y. \quad (87)$$

1522 Therefore, when  $X$  and  $X'$  are i.i.d.,

$$1523 (f(X) - f(X'))(h(X) - h(X')) \leq 0 \text{ a.s.} \quad (88)$$

1524 Taking expectations,

$$1525 \mathbb{E}[(f(X) - f(X'))(h(X) - h(X'))] \leq 0. \quad (89)$$

1526 Plug this back into Eq. (86), we have

$$1527 \text{Cov}(f(X), h(X)) = \frac{1}{2} \mathbb{E}[(f(X) - f(X'))(h(X) - h(X'))] \leq 0. \quad (90)$$

1528 □

#### 1529 B.4 PROOF OF THEOREM 4.4

1530 **Theorem B.7** (Truncation bias bound for KM-ADD). *Define the truncation bias of the LB-ADD*  
 1531 *as  $\mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{LB})}) := M_T^{(\text{LB})} - M_\infty$ , where  $M_T^{(\text{LB})} := \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T]$  is the true*  
 1532 *ADD under random sequence lengths. For  $b = \Delta T_{\text{max}}^*$ , we have  $\mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{LB})}) \leq \mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{KM})}) \leq 0$ ,*  
 1533 *under the independent censoring assumption in Thm. 4.2.*  
 1534

1535 *Proof.* Recall that

$$1536 \mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{KM})}) = M_T^{(\text{KM})} - M_\infty$$

$$1537 = \int_0^{\Delta T_{\text{max}}^*} S^{\text{ADD}}(t) dt - \mathbb{E}[\Delta\tau \mid \nu < \infty, \Delta\tau \geq 0] \quad (91)$$

$$1538 \mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{LB})}) = M_T^{(\text{LB})} - M_\infty$$

$$1539 = \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T] - \mathbb{E}[\Delta\tau \mid \nu < \infty, \Delta\tau \geq 0]. \quad (92)$$

1540 First, since

$$1541 \int_0^\infty S^{\text{ADD}}(t) dt \quad (93)$$

$$1542 = \int_0^\infty P(\Delta\tau > t \mid \nu < \infty, \Delta\tau \geq 0) dt \quad (94)$$

$$1543 = \int_0^\infty \mathbb{E}[\mathbf{1}(\Delta\tau > t) \mid \nu < \infty, \Delta\tau \geq 0] dt \quad (95)$$

$$1544 = \mathbb{E}\left[\int_0^\infty \mathbf{1}(\Delta\tau > t) dt \mid \nu < \infty, \Delta\tau \geq 0\right] \quad (96)$$

$$1545 = \mathbb{E}[\Delta\tau \mid \nu < \infty, \Delta\tau \geq 0], \quad (97)$$

1546 we have  $\mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{KM})}) = - \int_{\Delta T_{\text{max}}^*}^\infty S^{\text{ADD}}(t) dt \leq 0$ .

Second, we will show that  $\mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{LB})}) \leq \mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{KM})})$ . Since

$$M_T^{(\text{KM})} = \int_0^{\Delta T_{\max}^*} S^{\text{ADD}}(t) dt \quad (98)$$

$$= \int_0^{\Delta T_{\max}^*} P(\Delta\tau > t \mid \nu < \infty, \Delta\tau \geq 0) dt \quad (99)$$

$$= \int_0^{\Delta T_{\max}^*} \mathbb{E}[\mathbf{1}(\Delta\tau > t) \mid \nu < \infty, \Delta\tau \geq 0] dt \quad (100)$$

$$= \mathbb{E}[\int_0^{\Delta T_{\max}^*} \mathbf{1}(\Delta\tau > t) dt \mid \nu < \infty, \Delta\tau \geq 0] \quad (101)$$

$$\begin{aligned} &= P(\Delta\tau \leq \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \\ &\quad \times \mathbb{E}[\int_0^{\Delta T_{\max}^*} \mathbf{1}(\Delta\tau > t) dt \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*] \\ &\quad + P(\Delta\tau > \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \\ &\quad \times \mathbb{E}[\int_0^{\Delta T_{\max}^*} \mathbf{1}(\Delta\tau > t) dt \mid \nu < \infty, 0 \leq \Delta\tau > \Delta T_{\max}^*] \end{aligned} \quad (102)$$

$$\begin{aligned} &= P(\Delta\tau \leq \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \\ &\quad \times \mathbb{E}[\Delta\tau \mid \nu < \infty, \Delta\tau \leq \Delta T_{\max}^*] \\ &\quad + P(\Delta\tau > \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \\ &\quad \times \mathbb{E}[\Delta T_{\max}^* \mid \nu < \infty, 0 \leq \Delta\tau > \Delta T_{\max}^*] \end{aligned} \quad (103)$$

$$\begin{aligned} &= P(\Delta\tau \leq \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \\ &\quad \times \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*] \\ &\quad + P(\Delta\tau > \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \Delta T_{\max}^* \end{aligned} \quad (104)$$

$$\begin{aligned} &= (1 - P(\Delta\tau > \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0)) \\ &\quad \times \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*] \\ &\quad + P(\Delta\tau > \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \Delta T_{\max}^* \end{aligned} \quad (105)$$

$$\begin{aligned} &= \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*] \\ &\quad + P(\Delta\tau > \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \\ &\quad \times (\Delta T_{\max}^* - \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*]) \end{aligned} \quad (106)$$

we have

$$\begin{aligned} &\mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{KM})}) - \mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{LB})}) \\ &= M_T^{(\text{KM})} - M_T^{(\text{LB})} \end{aligned} \quad (107)$$

$$\begin{aligned} &= \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*] \\ &\quad + P(\Delta\tau > \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \\ &\quad \times (\Delta T_{\max}^* - \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*]) \\ &\quad - \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T] \end{aligned} \quad (108)$$

$$\begin{aligned} &= \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*] \\ &\quad - \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T] \\ &\quad + P(\Delta\tau > \Delta T_{\max}^* \mid \nu < \infty, \Delta\tau \geq 0) \\ &\quad \times (\Delta T_{\max}^* - \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*]). \end{aligned} \quad (109)$$

Finally, we can see that  $\mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*] - \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T] \geq 0$  (Lem. B.8), and  $\Delta T_{\max}^* - \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*] \geq 0$ ; hence,  $\mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{KM})}) - \mathcal{B}_{\text{TR}}(\hat{M}_T^{(\text{LB})}) \geq 0$ .

□

1620 **Lemma B.8.**  $\mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T_{\max}^*] - \mathbb{E}[\Delta\tau \mid \nu < \infty, 0 \leq \Delta\tau \leq \Delta T] \geq 0.$   
 1621

1622 *Proof.* Without loss of generality, this inequality is equivalent to  $\mathbb{E}[\Delta\tau \mid 0 \leq \Delta\tau \leq \Delta T_{\max}^*] - \mathbb{E}[\Delta\tau \mid$   
 1623  $0 \leq \Delta\tau \leq \Delta T] \geq 0$ , where  $\Delta T \leq \Delta T_{\max}^* (\in \mathbb{R})$  holds a.s.  $\Delta\tau$  and  $\Delta T (= C^{\text{ADD}})$  are independent  
 1624 because of the independent censoring assumption. This inequality can be rewritten as

$$1625 \mathbb{E}_{X,Y}[X \mid X \leq y_{\max}] - \mathbb{E}_{X,Y}[X \mid X \leq Y] \geq 0, \quad (110)$$

1626 where  $X$  is a non-negative random variable,  $Y$  is a non-negative random variable with  $0 \leq Y \leq y_{\max}$   
 1627 a.s.,  $y_{\max} (\geq 0)$  is a real constant, and  $X$  and  $Y$  are independent. Because Ineq. (110) coincides with  
 1628 Ineq. (70), the proof follows the same steps as Lem. B.5.  $\square$   
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673

## C SUPPLEMENTARY EXPERIMENTAL RESULTS

### C.1 EXAMPLE OF SURVIVAL CURVE OF DETECTION POINTS

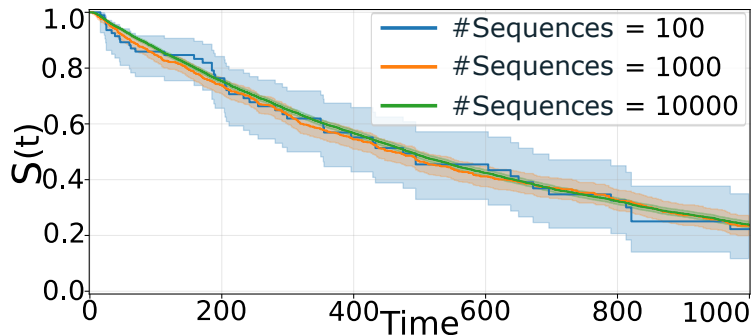


Figure 6: **Estimated survival functions of detection points.** The shaded area represents the standard error of the mean.

Survival functions of detection points are shown in Fig. 6. The experiment uses the Gaussian process dataset at three dataset sizes. The QCD model used is the GSR procedure with ground-truth statistics, evaluated under a given threshold. Changepoint locations are sampled from a geometric distribution with the success probability of  $p = 0.001$ . Error bars represent the standard error of the mean.

### C.2 COMPLETE ABLATION OF FIG. 4

Fig. 7 presents the complete ablation study corresponding to Fig. 4, demonstrating that KM-ARL and KM-ADD yield more accurate estimates of the true ARL-ADD curve, even under limited and irregular sequence lengths. The gray areas indicate regions of extrapolation (excluding the true ARL), where ARLs cannot be estimated due to the absence of data, unless additional assumptions are imposed on the underlying distribution.

### C.3 GAUSSIAN PROCESS DATASET

We provide additional experimental results on the Gaussian process dataset. Fig. 8 is an example sequence. Fig. 9 shows temporal statistics of a sequence in our Gaussian process dataset. Fig. 10 presents the ARL-ADD tradeoff curves of the QCD models evaluated on a Gaussian process dataset, while we use the WISDM Actitracker dataset in Fig. 1 in the main text. It demonstrates that our KM-ARL and KM-ADD are more robust to irregular lengths than the conventional LB-ARL and LB-ADD. Note that some curves in Fig. 10(a) are non-monotonic because the sequences used to compute the LB-ARL and LB-ADD differ drastically across different thresholds, causing unstable estimates. For example, short sequences are excluded from the computation of the LB metrics when the threshold is high. This issue does not arise for our KME-based metrics, contributing to their robustness to finite and irregular sequence lengths.

### C.4 POISSON PROCESS DATASET

We provide experimental results on the Poisson process dataset. The pre-change and post-change Poisson processes have the mean of 1 and 4, respectively. We use two types of changepoint distributions: geometric and uniform. Fig. 8 is an example sequence from our Poisson process dataset. Fig. 12 shows statistics of a sequence in our Poisson process dataset. Fig. 13 presents the ARL-ADD curve evaluated on the Poisson process dataset, demonstrating that the KM-ARL and KM-ADD provide more accurate estimates of the true ARL-ADD curve, even when sequence lengths are limited and irregular.

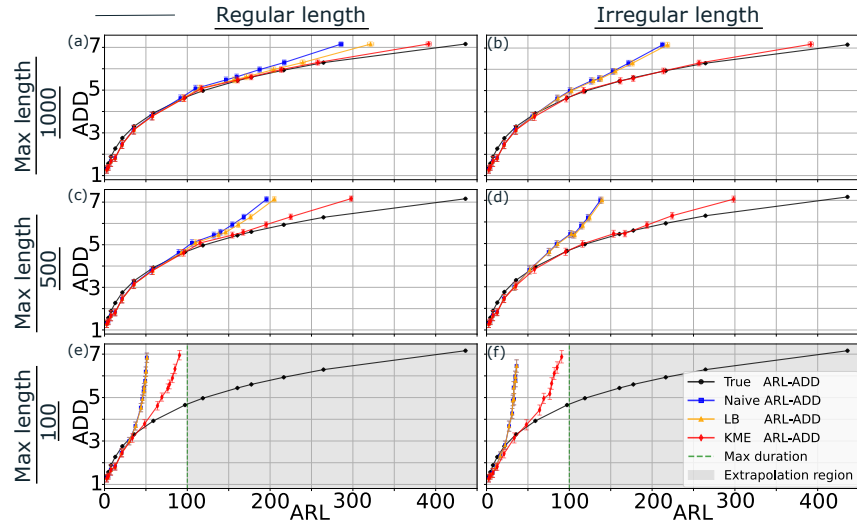


Figure 7: **Complete ablation study of Fig. 4.** The KM-ARL and KM-ADD provide more accurate estimates of the true ARL-ADD curve, even when sequence lengths are limited and irregular. The experiments use the Gaussian process dataset, comprising 10000 sequences. The employed QCD algorithm is the GSR procedure using ground-truth statistics. Change-point locations are sampled uniformly. 50% of sequences contain a change-point. Error bars represent the standard error of the mean. The gray areas indicate regions of extrapolation (excluding the true ARL), where ARLs cannot be estimated due to the absence of data, unless additional assumptions are imposed on the underlying distribution. (a) Sequence length is 1000. (b) Sequence lengths vary irregularly in the range  $[100, 1000]$ . (c) Sequence length is 500. (d) Sequence lengths vary irregularly in the range  $[50, 500]$ . (e) Sequence length is 100. (f) Sequence lengths vary irregularly in the range  $[10, 100]$ .

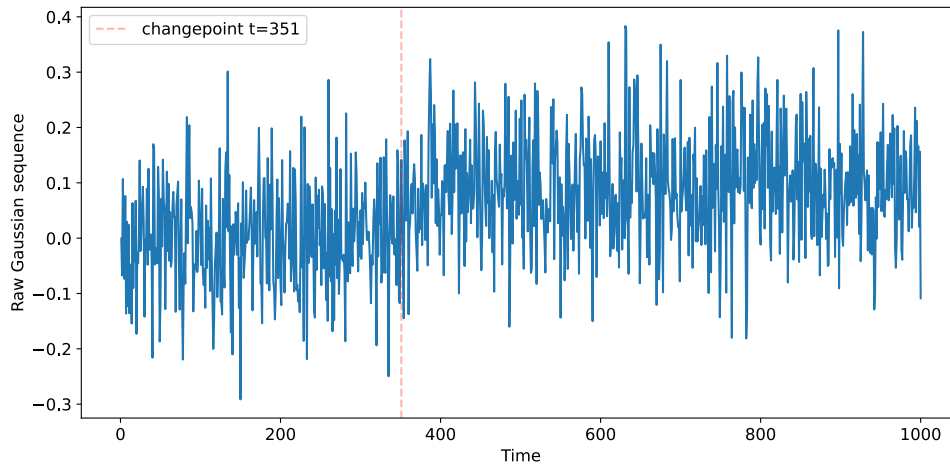


Figure 8: **Example sequence in Gaussian process dataset.**

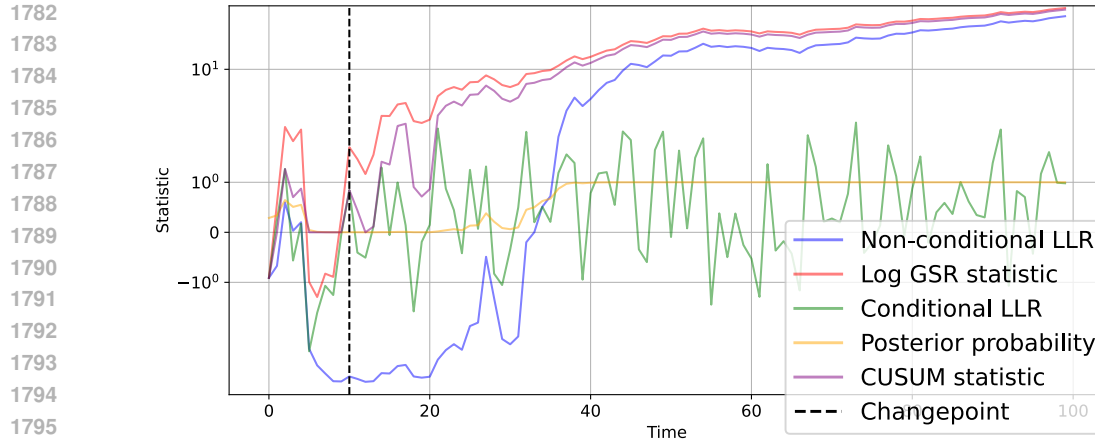


Figure 9: **Example statistics in Gaussian process dataset.** The observation interval is  $\Delta t = 1$ . LLR is short for log-likelihood ratio  $\log(f(X^{(0,t)})/g(X^{(0,t)}))$ , where  $f$  and  $g$  are the pre- and post-change density, respectively. Non-conditional LLR is given by  $\log(f(X^{(0,t)})/g(X^{(0,t)}))$ , while Conditional LLR is given by  $\log(f(X^{(t)} | X^{(0,t-1)})/g(X^{(t)} | X^{(0,t-1)}))$ . Posterior probability means  $g(X^{(0,t)})$ .

### C.5 WISDM ACTITRACKER DATASET

We provide additional experimental results on the WISDM Actitracker dataset. Fig. 14 is the original figure of Fig. 1.

While we use the machine-labeled subset in the main text, Fig. 15 presents the ARL-ADD curves of QCD models evaluated on the user-labeled subset (“labeled” subset) of the WISDM Actitracker dataset, which is about  $10^3 (= 51326/83)$  times smaller than the machine-labeled subset (“unlabeled” subset) of the WISDM Actitracker dataset (see Tab. 2 for statistics). Evaluation is challenging on this subset because the number of sequences is limited to only 83. For such small datasets, we recommend using min-max-based metrics, such as the box plot, rather than average-based metrics, such as the ARL and ADD.

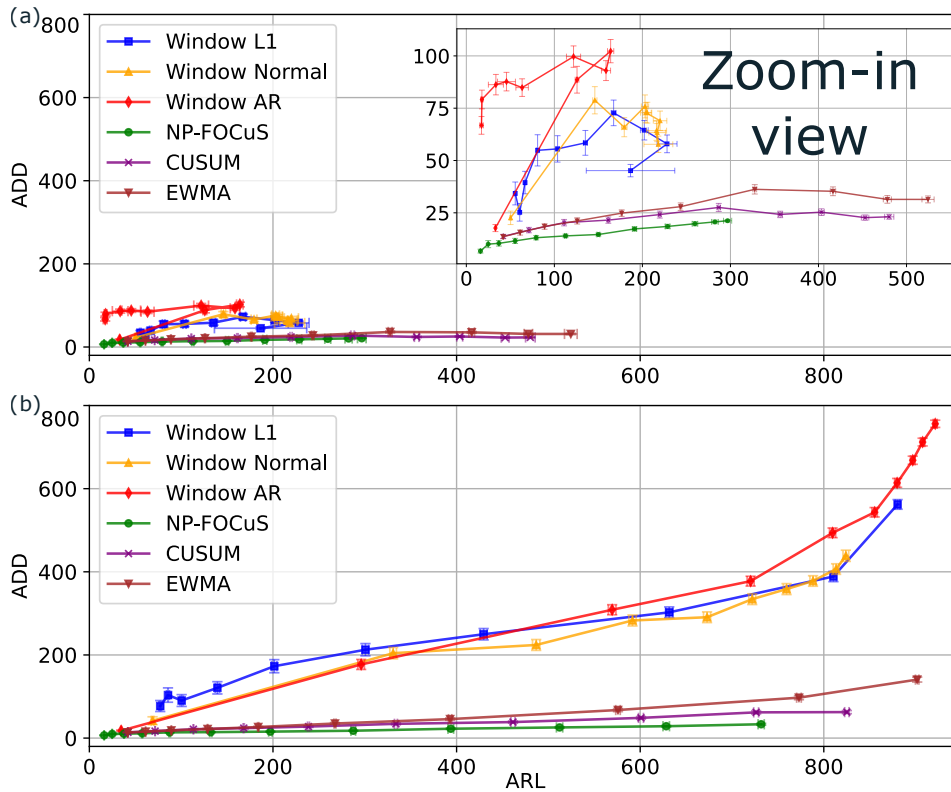
### C.6 ARL-ADD WITH CUSUM

While we use the GSR procedure in the main text, we provide additional experimental results obtained with the CUSUM procedure (Page, 1954) with ground-truth statistics. The ARL-ADD curves are provided in Fig. 16 and support our findings in the main text. The CUSUM procedure is evaluated under various thresholds.

### C.7 ARL-ADD WITH GEOMETRIC CHANGEPOINT DISTRIBUTION

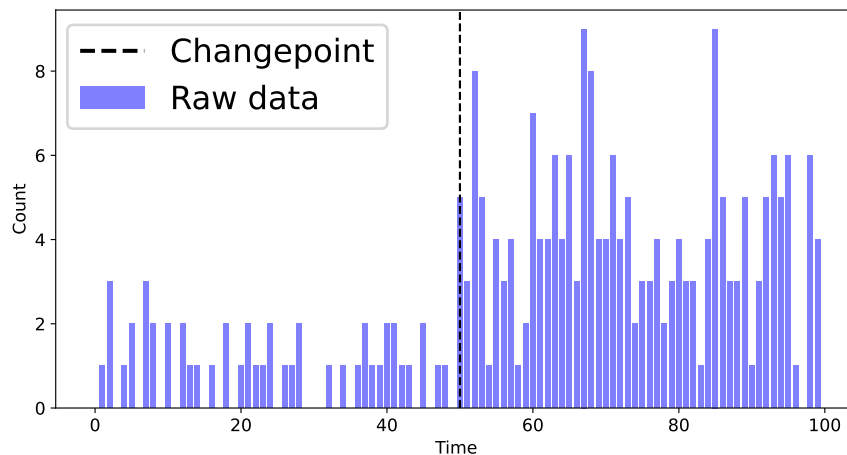
We provide additional experimental results obtained with the geometric changepoint distribution, instead of the uniform distribution. We evaluate the ARL-ADD tradeoff curve on a Gaussian process dataset. The results are provided in Fig. 17 and support our findings in the main text.

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861



1862 **Figure 10: Evaluation of QCD models on Gaussian process dataset. (a) LB-ARL and LB-ADD.**  
1863 **(b) KM-ARL and KM-ADD.** Our KM-ARL and KM-ADD are more robust to irregular lengths  
1864 than the conventional LB-ARL and LB-ADD. The experiments use the Gaussian process dataset,  
1865 comprising 10000 sequences. Changepoint locations are sampled uniformly. 50% of the sequences  
1866 contain a changepoint. Sequence lengths vary irregularly in the range  $[100, 1000]$ . Note that some  
1867 curves in (a) are non-monotonic because the sequences used to compute the LB-ARL and LB-ADD  
1868 differ drastically across different thresholds, causing unstable estimates. For example, short sequences  
1869 are excluded from the computation of the LB metrics when the threshold is high. This issue does not  
1870 arise for our KME-based metrics, contributing to their robustness to finite and irregular sequence  
1871 lengths.

1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889



**Figure 11: Example sequence in Poisson process dataset.**

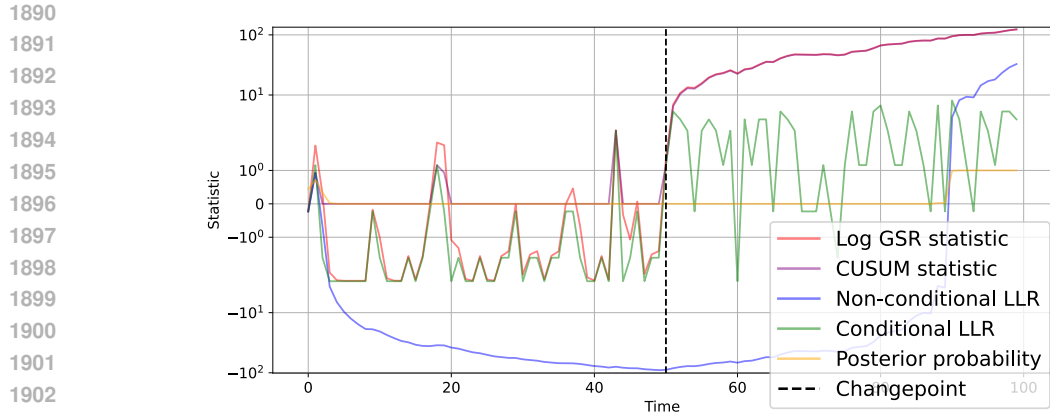


Figure 12: **Example statistics in Gaussian process dataset.** Observation interval is  $\Delta t = 1$ . LLR is short for log-likelihood ratio  $\log(f(X^{(0,t)})/g(X^{(0,t)}))$ , where  $f$  and  $g$  are the pre- and post-change density, respectively. Non-conditional LLR is given by  $\log(f(X^{(0,t)})/g(X^{(0,t)}))$ , while Conditional LLR is given by  $\log(f(X^{(t)} | X^{(0,t-1)})/g(X^{(t)} | X^{(0,t-1)}))$ . Posterior probability means  $g(X^{(0,t)})$ .

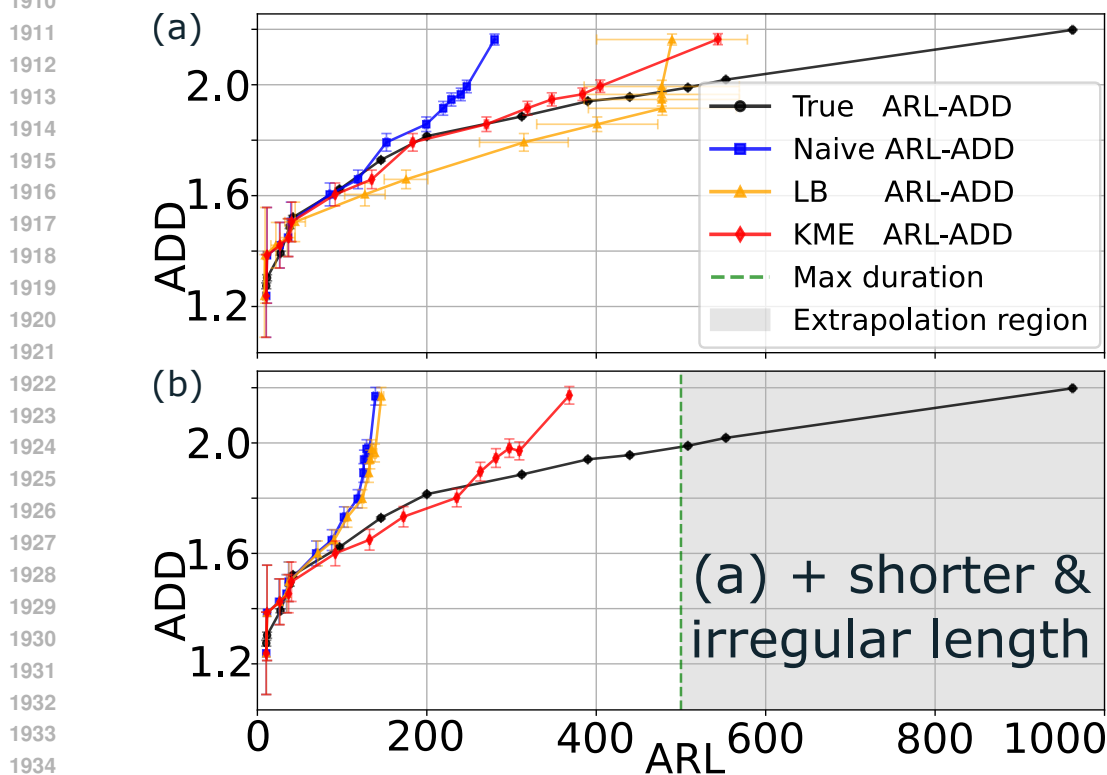


Figure 13: **ARL-ADD tradeoff curves on Poisson process dataset.** The KM-ARL and KM-ADD provide more accurate estimates of the true ARL-ADD curve, even when sequence lengths are limited and irregular. (a) **Sequence length is 1000.** (b) **Sequence lengths vary irregularly in the range [50, 500].** The experiments use the Gaussian process dataset, comprising 10000 sequences. The employed QCD algorithm is the GSR procedure using ground-truth statistics. Changepoint locations are sampled uniformly. 50% of sequences contain a changepoint. Error bars represent the standard error of the mean. In (b),  $ARL > 500$  (gray area) indicates a region of extrapolation (excluding the true ARL), where ARLs cannot be estimated due to the absence of data, unless additional assumptions are imposed on the underlying distribution.

1944  
 1945  
 1946  
 1947  
 1948  
 1949  
 1950  
 1951  
 1952  
 1953  
 1954  
 1955  
 1956  
 1957  
 1958  
 1959  
 1960  
 1961  
 1962  
 1963  
 1964  
 1965  
 1966  
 1967  
 1968  
 1969  
 1970  
 1971  
 1972  
 1973  
 1974  
 1975  
 1976  
 1977  
 1978  
 1979  
 1980  
 1981  
 1982  
 1983  
 1984  
 1985  
 1986  
 1987  
 1988  
 1989  
 1990  
 1991  
 1992  
 1993  
 1994  
 1995  
 1996  
 1997

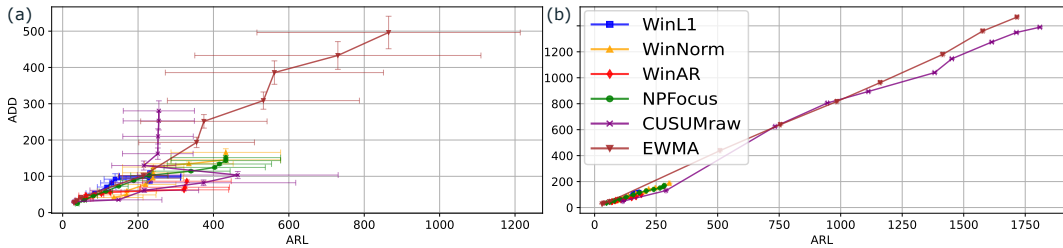


Figure 14: **Original figure of Fig. 1.** Our KM-ARL and KM-ADD are more interpretable than the conventional LB-ARL and LB-ADD.

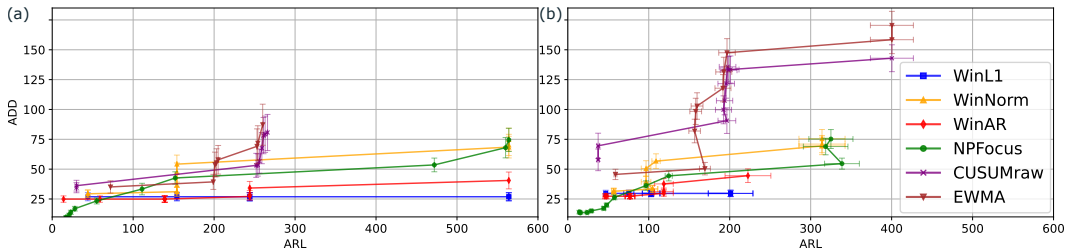


Figure 15: **Evaluation of QCD models on real-world dataset (user-labeled subset of WISDM Actitracker).** The user-labeled subset (“labeled” subset) of the WISDM Actitracker is used, which is about  $10^3 (= 51326/83)$  times smaller than the machine-labeled subset (“unlabeled” subset) of the WISDM Actitracker dataset (see Tab. 2 for statistics). Evaluation is challenging on this subset because the number of sequences is limited to only 83. For such small datasets, we recommend using min-max-based metrics, such as the box plot, rather than average-based metrics, such as the ARL and ADD.

1998  
 1999  
 2000  
 2001  
 2002  
 2003  
 2004  
 2005  
 2006  
 2007  
 2008  
 2009  
 2010  
 2011  
 2012  
 2013  
 2014  
 2015  
 2016  
 2017  
 2018  
 2019  
 2020  
 2021  
 2022  
 2023  
 2024  
 2025  
 2026  
 2027  
 2028  
 2029  
 2030  
 2031  
 2032  
 2033  
 2034  
 2035  
 2036  
 2037  
 2038  
 2039  
 2040  
 2041  
 2042  
 2043  
 2044  
 2045  
 2046  
 2047  
 2048  
 2049  
 2050  
 2051

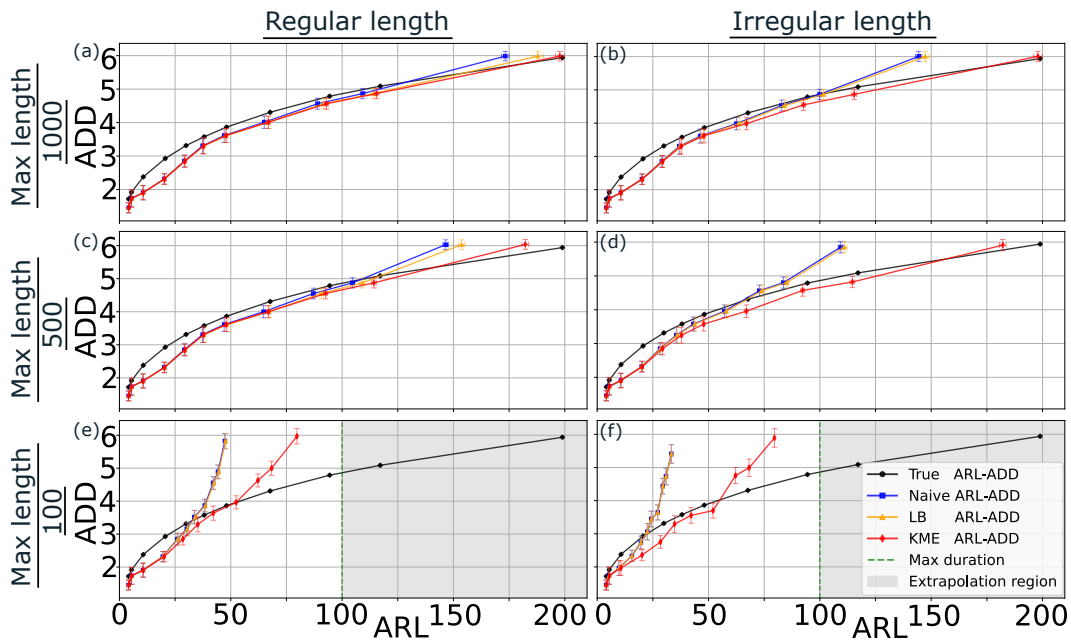


Figure 16: **ARL-ADD curve of CUSUM procedure with ground-truth statistics.** The results closely parallel those obtained with the GSR procedure in the main text: The KM-ARL and KM-ADD provide more accurate estimates of the true ARL-ADD curve, even when sequence lengths are limited and irregular. CUSUM is evaluated under various thresholds. The experiments use the Gaussian process dataset, comprising 10000 sequences. Change point locations are sampled uniformly. 50% of sequences contain a change point. Error bars represent the standard error of the mean. The gray areas indicate regions of extrapolation (excluding the true ARL), where ARLs cannot be estimated due to the absence of data, unless additional assumptions are imposed on the underlying distribution. (a) Sequence length is 1000. (b) Sequence lengths vary irregularly in the range  $[100, 1000]$ . (c) Sequence length is 500. (d) Sequence lengths vary irregularly in the range  $[50, 500]$ . (e) Sequence length is 100. (f) Sequence lengths vary irregularly in the range  $[10, 100]$ .

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

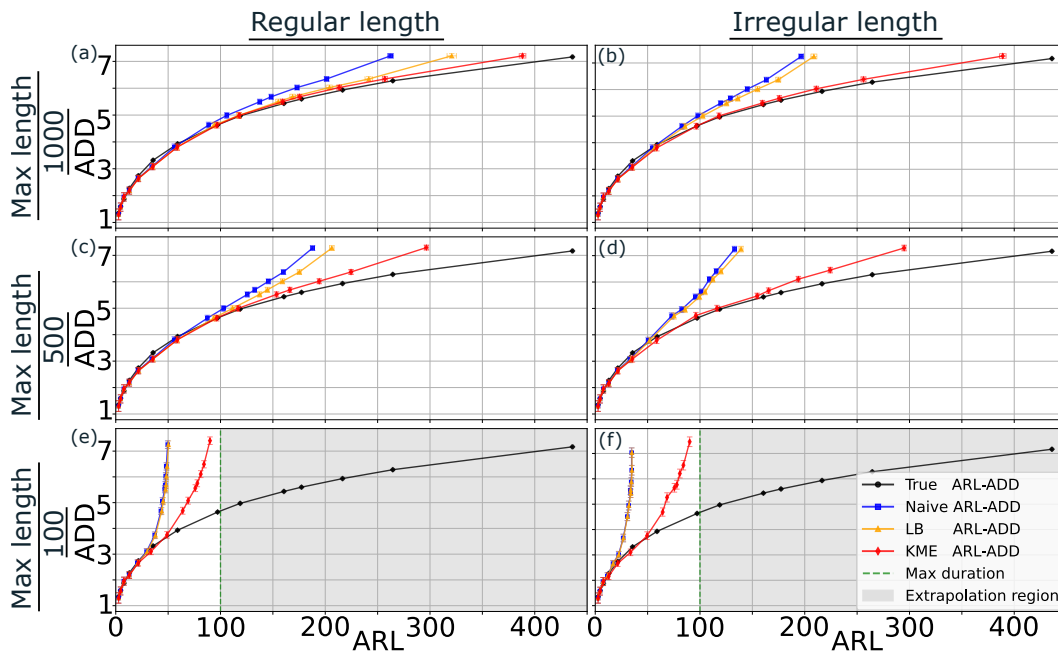


Figure 17: **ARL-ADD curve on Gaussian process dataset with geometric changepoint distribution.** The results closely parallel those obtained with the uniform changepoint distribution in the main text: The KM-ARL and KM-ADD provide more accurate estimates of the true ARL-ADD curve, even when sequence lengths are limited and irregular. The experiments use the Gaussian process dataset, comprising 10000 sequences. The employed QCD algorithm is the GSR procedure using ground-truth statistics. Changepoint locations are sampled from a geometric distribution with the success probability  $p = 0.001$ . Error bars represent the standard error of the mean. The gray areas indicate regions of extrapolation (excluding the true ARL), where ARLs cannot be estimated due to the absence of data, unless additional assumptions are imposed on the underlying distribution. (a) Sequence length is 1000. (b) Sequence lengths vary irregularly in the range  $[100, 1000]$ . (c) Sequence length is 500. (d) Sequence lengths vary irregularly in the range  $[50, 500]$ . (e) Sequence length is 100. (f) Sequence lengths vary irregularly in the range  $[10, 100]$ .

## D DETAILED EXPERIMENTAL SETTINGS

**Runtime.** The computational costs of metric evaluation, including KM-ARL and KM-ADD, are negligible in our experiments compared with the time required to run the QCD algorithms. Computing a single data point in the figures takes from a few seconds to several hours, depending on the sequence length, the number of sequences, and the QCD algorithm used.

**Computing Infrastructure.** We use Python 3.11.9 (Van Rossum & Drake, 2009), Numpy 1.26.4 (Harris et al., 2020), lifelines 0.30.0 (Davidson-Pilon, 2019), changepoint-online 1.2.1 (Romano et al., 2024b), ruptures 1.1.9 (Truong et al., 2020), ocpdet 0.0.6 (Khamesi, 2022) (we pick up relevant implementations only (`ocpdet/CUSUM.py` and `ocpdet/EWMA.py`), and our implementation does not require TensorFlow (Abadi et al., 2015), which is required when installing `ocpdet`). PyTorch 2.6.0 (Paszke et al., 2019) (used only for saving and loading the Gaussian and Poisson datasets and can be replaced with Numpy etc.), The operating system (OS) is Ubuntu 20.04. Intel(R) Core(TM) i9-7980XE CPU @ 2.60 GHz (18 cores) are used. The total random access memory (RAM) size of our server is 125 GBs. We do not use GPUs.

### D.1 ONLINE QCD MODELS

To see a practical relevance of our KME-based metrics, we evaluate six QCD models on the simulated Gaussian process dataset (App. C.3) and real-world WISDM Actitracker dataset (main text and App. C.5). We use the following online QCD models: Window L1 (Bai, 1995), Window Normal (Lavielle, 1999; Lavielle & Teyssiere, 2006), Window AR (Bai, 2000), non-parametric focused changepoint detection (NP-FOCuS) (Romano et al., 2024a), CUSUM (Page, 1954), and exponentially weighted moving average (EWMA) (Roberts, 1959). For Window L1, Window Normal, Window AR, the discrepancy measures for detection are derived from *cost functions* as specified in (Truong et al., 2020).

Window L1 detects changes in the median and uses the  $L^1$  norm for the cost function. Window Normal detects changes in the mean and covariance matrix of a sequence of multivariate Gaussian random variables. It uses the log likelihood of empirical Gaussian distribution, for the computation of the cost function. Window AR estimates the least-squares estimates of the break dates obtained from a piecewise autoregressive (AR) model. NP-FOCuS is an online, non-parametric changepoint detection algorithm designed to efficiently detect distributional changes in real-time data streams by leveraging functional pruning techniques. For multivariate time-series, we instantiate  $n$  detectors and take the minimum detection time, where  $n$  is the number of features. CUSUM is a well-known QCD algorithm, which detects persistent shifts in the mean of a sequence. For multivariate time-series, we use the norm of the feature vector for the CUSUM chart. EWMA computes a running average of datapoints, placing exponentially decreasing weights on older observations.

Window L1, Window Normal, Window AR are implemented with `ruptures` (Truong et al., 2020), NP-FOCuS is implemented with `changepoint-online` (Romano et al., 2024b), and CUSUM and EWMA are implemented with `ocpdet` (Khamesi, 2022). We adapt window-based models from `ruptures`, originally designed for offline changepoint detection, to online QCD.

We fix both the window size and the burn-in interval at 30 frames, which is about three times smaller than the minimum maximum sequence length in our simulation experiments (100). It is also much smaller than the average lengths in the WISDM Actitracker. Most other hyperparameters remain at their default values; otherwise, they are specified in our code.

### D.2 PREPROCESSES FOR WISDM ACTITRACKER

Our preprocesses for the WISDM Actitracker dataset for both “labeled” and “unlabeled” subsets are comprised of the following procedures:

1. Extract user IDs (“user”), features (“X0”, “X1”, “X2”, ..., “RESULTANT”), and frame labels (class”).
2. Convert the string labels (“Jobbing”, “Walking”, “Stairs”, “Sitting”, “Standing”, “Lying-Down”) to binary numeric labels (0 for pre-change and 1 for post-change). In the “labeled”

	User-labeled set	Machine-labeled set
#Sequences	83	51,326
#Frames	5,435	1,369,349
#Seqs. w/ 100% positive labels	37	76
#Seqs. w/ 0% positive labels	17	61
#Seqs. w/ mixed labels	29	51,189
Positive frame ratio	0.741	0.684
Mean length	65.5	26.7
Min length	1	1
Max length	565	54,401

Table 2: **Statistics of WISDM Actitracker after preprocesses.** A positive label here means that a frame (timestamp) is sampled from the post-change distribution. The user-labeled set and the machine-labeled set corresponds to the “labeled” set and the “unlabeled” set in the original WISDM Actitracker dataset. Note that the “unlabeled” set is actually labeled with a system developed in the WISDM Lab.

subset, “Sitting” is mapped to 0 and all other activities to 1. In the “unlabeled” subset, “Walking” and “Sitting” are mapped to 1, while the remaining activities are mapped to 0.

3. Split sequences that have a label transition from 1 to 0 to remove turn-back sequences.
4. Zero-pad outlier features (value  $> 10^{12}$ ) (although some huge values such as  $10^{10}$  still remain).
5. Normalize each feature to the range  $[-1, 1]$ .
6. Save feature vectors, labels, and changepoints to a file.

All the preprocesses are detailed in our code (`WISDMactitracker.ipynb`). Finally, we manually exclude the sequences of length 1 when evaluating QCD algorithms (in `calc_esARL_WISDM_cpmodels.py` and `calc_esADD_WISDM_cpmodels.py` of our code).

### D.3 STATISTICS OF WISDM ACTITRACKER

Tab. 2 summarizes the statistics of the WISDM Actitracker dataset used in our experiments. The user-labeled and the machine-labeled subsets correspond to the “labeled” and the “unlabeled” subsets, respectively, in the original WISDM Actitracker dataset. Note that the “unlabeled” subset is, in fact, labeled with a system developed in the WISDM Lab (Kwapisz et al., 2011).

Fig. 18 shows histograms of sequence lengths of the user-labeled and machine-labeled subsets of the WISDM Actitracker dataset after the preprocesses, which are detailed in App. D.2. The sequence lengths exhibit substantial irregularity, and estimating the ARL and ADD is challenging.

Fig. 19 shows histograms of changepoint indices (timestamps) of the user-labeled and machine-labeled subsets of the WISDM Actitracker dataset after the preprocesses, which are detailed in App. D.2. They contain a variety of pre-change lengths, but the number of without-change sequences are limited, and estimating the ARL is challenging.

### D.4 NUMERICAL COMPUTATION OF VARIANCE

We describe the computation methods of the restricted variance of the detection time,  $\hat{V} := \widehat{\text{Var}}[\min\{\tau, a\}]$ , used in our experiments.

The restricted variance of LB-ARL is defined as the naive empirical variance:  $\hat{V} = \widehat{\text{Var}}[\min\{\tau, T\}]$ , which is consistent with the mean of LB-ARL ( $\mathbb{E}[\min\{\tau, T\}]$ ). Given a dataset with size  $N$ , it is computed as  $\frac{1}{N_{\text{LB}}} \sum_{i=1}^{N_{\text{LB}}} (\tau_i - \bar{\tau})^2$ .  $\bar{\tau}$  is the empirical LB-ARL, i.e., the detection time averaged over the subset in which  $\tau_i \leq$  sequence length. The subset size is denoted by  $N_{\text{LB}} (\leq N)$ .

There are several methods to numerically compute the variance of RMST (KM-ARL). We adopt `lifelines` (Davidson-Pilon, 2019), a standard Python library for survival analysis, to compute the KME and its variance. Specifically, we use

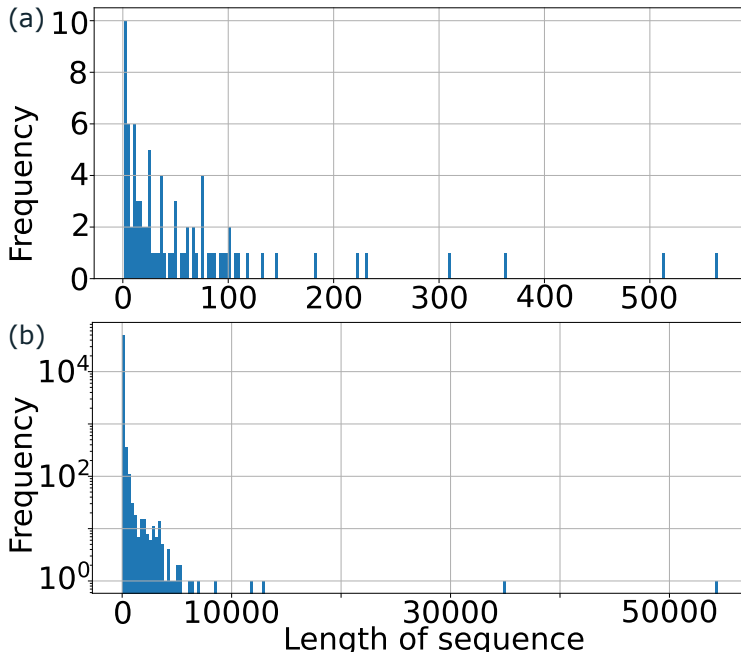


Figure 18: **Lengths of WISDM Actitracker sequences.** The sequence lengths exhibit substantial irregularity, and estimating the ARL and ADD is challenging. **(a) User-labeled subset ( $y$ -linear scale).** **(b) Machine-labeled subset ( $y$ -log scale).** The user-labeled and the machine-labeled subsets correspond to the “labeled” and the “unlabeled” subsets, respectively, in the original WISDM Actitracker dataset. Note that the “unlabeled” subset is, in fact, labeled with a system developed in the WISDM Lab (Kwapisz et al., 2011).

`lifelines.utils.restricted_mean_survival_time` with `return_variance=True`, in which  $\hat{V} = 2 \int_0^a t \hat{S}(t) dt - (\int_0^a \hat{S}(t) dt)^2$ , where  $\hat{S}(t)$  denotes the empirical survival function, computed from the given dataset with size  $N$ . The derivation is given in the appendix in (Royston & Parmar, 2013), as is mentioned in the documentation<sup>1</sup>.

**Why do KM-ARL/ADD reduce variance compared to LB-ARL/ADD?** We elaborate on the discussion of our experimental results on the real-world dataset. In short,  $\hat{V}$  of LB-ARL tends to larger than that of KM-ARL because the former can use only  $N_{LB} \leq N$ , while the latter can exploit all  $N$  sequences. LB-ARL/ADD become unstable and exhibit high variance, particularly at large QCD thresholds, where  $N_{LB}$  tends to be small because the detector often fails to raise an alarm within the sequence length. In contrast, KM-ARL/ADD do not suffer from this issue because they are calculated from a constant number of sequences  $N$ , regardless of the threshold, which enhances their robustness against censoring.

<sup>1</sup>[https://lifelines.readthedocs.io/en/latest/lifelines.utils.html#lifelines.utils.restricted\\_mean\\_survival\\_time](https://lifelines.readthedocs.io/en/latest/lifelines.utils.html#lifelines.utils.restricted_mean_survival_time).

2268  
 2269  
 2270  
 2271  
 2272  
 2273  
 2274  
 2275  
 2276  
 2277  
 2278  
 2279  
 2280  
 2281  
 2282  
 2283  
 2284  
 2285  
 2286  
 2287  
 2288  
 2289  
 2290  
 2291  
 2292  
 2293  
 2294  
 2295  
 2296  
 2297  
 2298  
 2299  
 2300  
 2301  
 2302  
 2303  
 2304  
 2305  
 2306  
 2307  
 2308  
 2309  
 2310  
 2311  
 2312  
 2313  
 2314  
 2315  
 2316  
 2317  
 2318  
 2319  
 2320  
 2321

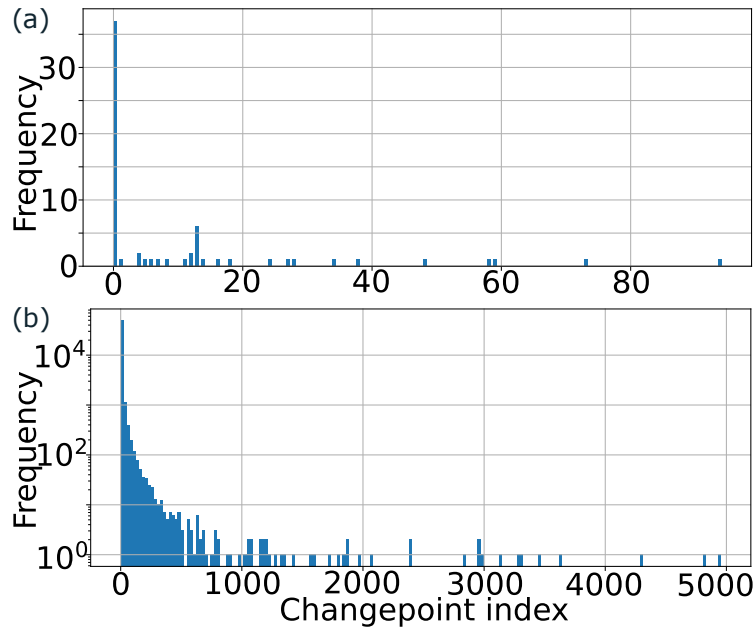


Figure 19: **Changepoint indices of WISDM Actitracker sequences.** Changepoint index here means the timestamp at which a change occurs, i.e., the pre-change length. They contain a variety of pre-change lengths, but the number of without-change sequences are limited, and estimating the ARL is challenging. **(a) User-labeled subset ( $y$ -linear scale).** The number of sequences without a changepoint is 17 out of 83. **(b) Machine-labeled subset ( $y$ -log scale).** The number of sequences without a changepoint is 61 out of 51326. The user-labeled and the machine-labeled subsets correspond to the “labeled” and the “unlabeled” subsets, respectively, in the original WISDM Actitracker dataset. Note that the “unlabeled” subset is, in fact, labeled with a system developed in the WISDM Lab (Kwapisz et al., 2011).

## E SUPPLEMENTARY RELATED WORK

**QCD metrics.** General QCD metrics include the Shiryaev’s ARL (Shiryaev, 1963), Lorden’s worst-case ADD (Lorden, 1971), Pollak’s supremum conditional ADD (Pollak, 1985), exponential penalty detection delay (Poor, 1998), probability of false alarm (PFA) (Tartakovsky, 2019), precision-recall (Van den Burg & Williams, 2020), Jaccard index (Van den Burg & Williams, 2020),  $F_\beta$  score (Van den Burg & Williams, 2020), range-based precision-recall (Tatbul et al., 2018), NAB (Lavin & Ahmad, 2015), SoftED (Salles et al., 2024), Hausdorff distance (Van den Burg & Williams, 2020), adjusted Rand index (Hubert & Arabie, 1985), variation of information (Arabie & Boorman, 1973), and segmentation covering metric, (Everingham et al., 2010; Arbelaez et al., 2010). In this paper, we focus on the standard ARL and Bayesian ADD (expectation is taken over the changepoint) because they are intuitive for practical evaluation and are widely used in proving the optimality of many QCD algorithms (Tartakovsky, 2019).

To our knowledge, no existing metric explicitly addresses random irregular sequence lengths in QCD. Some metrics consider sequence truncation, such as the inverse of the false alarm rate over a finite time interval  $\Delta t$  (Alakent & Mutlu, 2018), which approximates the ARL for large  $\Delta t$ . The PFA in horizon (Huang & Veeravalli, 2024) is defined as the PFA measured on truncated sequences. However, these metrics are designed for fixed, finite-length sequences and introduce significant bias when applied to irregular-length sequences.

**ARL and ADD estimation.** There are many studies that derive the ARL and ADD of specific control charts or distributions: (Reynolds, 1975; Crowder, 1987; Saccucci & Lucas, 1990; Knoth & Knoth, 1998; Fu et al., 2002; Areepong & Peerajit, 2022; Haq & Woodall, 2022; Sunthornwat et al., 2023). Li et al. (2014) review simulation methods for computing the ARL and ADD of CUSUM (Page, 1954) and EWMA (Roberts, 1959), on simulation data. The references therein analyze Markov chain methods and integral equation methods. Alakent & Mutlu (2018) define the ARL of the Shewhart chart as the inverse of the alarm probability. In contrast, we propose the estimators of the ARL and ADD for arbitrary control charts, QCD models, and underlying distributions. Qiu (2013) defines a maximum length of the sequence for the ARL computation and excludes sequences that fail to detect changepoints by the maximum length. This estimator is similar to the LB-ARL and LB-ADD and exhibits a more substantial negative bias than our ARL estimator due to truncation, as proven in Sec.4 and demonstrated in Sec.5.

**Is ARL always informative?** (Mei, 2008) is the first paper in which the author systematically questions the appropriateness of the ARL. He claims that (i) a detection scheme with finite detection delay can have infinite ARL, where the detector can raise false alarms with probability 1, and (ii) under the standard minimax formulation with the ARL, we are in danger of finding a detection scheme that focuses on detecting larger changes instead of smaller changes. Claim (i) is theoretically relevant but practically irrelevant. It points out, in theory, that if the ARL is not assumed to be finite, one can construct a detector that raise false alarms. In practice, however, a detector with a small finite detection delay and an (approximately or numerically) infinite ARL is considered to perform well. This is not problematic in real-world applications. We agree with Claim (ii) because the minimax formulation of QCD, or more generally, the min and max operations are sensitive to outliers and fail to capture the average behavior of the model. For this reason, when used, they are typically complemented by average-based metrics such as the ARL and ADD. In summary, we encourage researchers to choose evaluation metrics with a clear understanding of their strengths and limitations.

**QCD model for finite sequence length.** To our knowledge, (Huang & Veeravalli, 2024) is the first and only work on the QCD problem on sequences of finite length. They propose a CUSUM variant for a fixed finite sequence length. If the model fails to detect the change before the horizon, they set the delay to  $T - \nu$ , causing a downward bias in the ADD.

**QCD for survival analysis.** There are many studies on the application of changepoint detection to survival analysis. In (Biswas & Kalbfleisch, 2008; Sego et al., 2009; GANDY et al., 2010; Phinikettos & Gandy, 2014; Gomon et al., 2024; Sasikumr & Sujatha, 2025), they propose variants of CUSUM to estimate survival time on right-censored data, which is orthogonal to our focus: accurately estimating ARLs and ADDs of arbitrary QCD algorithms on right-censored sequential data with irregular lengths. Gierz (2020) consider changepoint problems concerning a shift in a distribution for a set of

time-ordered observations under censoring or truncation. Polunchenko (2016) derive a closed-form formula for the survival function of the generalized Shiryayev-Roberts (GSR) detection time under Brownian motion.

**Parametric and closed-form approaches to ARL and ADD.** The exponential approximation of the underlying distribution provides accurate estimation of ARL and ADD in some scenarios. Also, closed-form expressions for ARL are known for some limited cases (e.g., (Fu et al., 2002; Areepong & Peerajit, 2022; Sunthornwat et al., 2023)). However, it has been also known that: (1) the exponential approximation does not work well if the distribution deviates from exponential (Borror et al., 2003; Pehlivan & Testik, 2010); and (2) closed-form expressions are not always available—it depends on distributions and detectors. This is why we emphasize that our estimators are applicable to *arbitrary* underlying distributions (=non-parametric) and *arbitrary* QCD models, a key contributions that clearly distinguishes our work from prior approaches.

## F SUPPLEMENTARY DISCUSSION

**Terminology.** In general, the terms “average”, “mean”, and “expected” are used interchangeably when referring to the ARL and ADD. The ARL is also referred to as the ARL to false alarm (Tartakovsky, 2019) or the average time to signal (Li et al., 2014). The ADD with a fixed changepoint is also referred to as the conditional expected delay to detection (CEDD) (Tartakovsky, 2019). The ADD with taking expectation over the changepoint is referred to as the Bayesian ADD or, simply, ADD (Tartakovsky, 2019). In our paper, we refer the Bayesian ADD as the ADD (with  $T = \infty$ ). In the statistical theory of control charts, including change detection in survival monitoring (or monitoring time-to-event data), the ARL and ADD are sometimes referred to as the in-control or zero-state ARL and the out-of-control or steady-state ARL, respectively (Saccucci & Lucas, 1990; Sasikumr & Sujatha, 2025; Lim & Lee, 2025). Strictly speaking, the steady-state ARL is defined as  $E[\tau \mid \nu = 0]$ , not  $E[\tau - \nu \mid \tau \geq \nu, \nu < \infty]$ .

**How to define  $T$  and  $T_{\max}^*$ .** We clarify how to define the distribution of the truncation length (denoted by  $F_T$  here), focusing on ARL,  $T$  (truncation length), and  $T_{\max}^*$  (the least upper bound for the support of  $F_T$ , i.e., the maximum possible sequence length). A parallel discussion applies to ADD. In brief, given an evaluation dataset, we can arbitrarily define  $F_T$ . In practice, once a test dataset is collected for model evaluation,  $F_T$  may be taken as the empirical CDF of the test dataset (or a mollified (smoothed) version of it). Accordingly,  $T_{\max}^*$  can be defined as the maximum sequence length in the test dataset, denoted by  $T_{\max}$  in the main text.

**Tighter bound.** In Cor. 1.1 in (Stute, 1994), an additional tighter *upper* bound is available, which potentially can be used to derive tighter bounds for the KM-ARL and KM-ADD.

**Theoretical computational efficiency.** The computational complexity of KM-ARL is  $\mathcal{O}(N) + \mathcal{O}(T_{\max})$  (plus event time sorting if necessary), where  $N$  is the number of sequences in the dataset and  $T_{\max}$  is the max sequence length. Bucketing the sequences into  $T_{\max}$  bins takes  $\mathcal{O}(N)$ , counting  $n_j^{\text{ARL}}$  and  $d_j^{\text{ARL}}$  takes  $\mathcal{O}(T_{\max})$ , computing  $\hat{S}^{\text{ARL}}$  takes  $\mathcal{O}(T_{\max})$ , and integrating the step-like function  $\hat{S}^{\text{ARL}}$  takes  $\mathcal{O}(T_{\max})$ .

**Empirical computational efficiency.** As mentioned in App. D, the computational costs of metric evaluation, including KM-ARL and KM-ADD, are negligible in our experiments compared with the time required to run the QCD algorithms. For significantly longer sequences (e.g.,  $\gtrsim 10^3$ , which is our maximum length) and larger datasets (e.g.,  $\gtrsim 10^4$ , which is our maximum size), we recommend splitting the dataset and aggregating the ARLs and ADDs, although it can degrade evaluation accuracy.

**Evaluation on significantly small datasets.** Evaluation of ARLs and ADDs on significantly small datasets such as the user-labeled subset of the WISDM Actitracker is challenging, as shown in Fig. 15 in App. C.5. In such cases, consider applying finite-sample bias correction methods (e.g., bootstrapping) or increasing samples.

2430 **Implementation in Python.** We implement our estimators based on Python (Van Rossum & Drake,  
2431 2009) because it is more compatible with machine learning development than R (R Core Team, 2021).  
2432

2433 **Variance estimation.** While our primary focus is on bias estimation, we also provide a rough  
2434 estimate of variance. The variance of LB-ARL (naive, standard definition of variance) is given by  
2435

$$2436 \int_0^a (t - \bar{t})^2 dF(t), \quad (111)$$

2437 where  $a$  denotes the cut-off time,  $F$  the CDF of the event time, and  $\bar{t}$  the mean event time. On the  
2438 other hand, Akritas (2000) provides the asymptotic variance of the RMST—analogue to the variance  
2439 of KM-ARL in our setting:  
2440

$$2441 \int_0^a \frac{S(t)}{1 - H(t)} (t - \bar{t})^2 dF(t) \quad (112)$$

2442 (Eq. (4) in (Akritas, 2000)), where  $S$  denotes the survival function and  $H$  the CDF of  
2443  $\min\{\text{event time, censoring}\}$ . Compared with LB-ARL (111), the variance of KM-ARL (112) is  
2444 computed as a weighted expectation of  $(t - \bar{t})^2$ , with its scale determined by the weight  $\frac{S(t)}{1 - H(t)}$ .  
2445 Because  $S(t) \in [0, 1]$ , this weight can grow large when  $1 - H(t)$  is small unless  $S(t)$  decreases  
2446 sufficiently in the same region. Consequently, the variance integral in (112) can be dominated by  
2447 the interval where  $t$  is large and close to the least upper bound of  $H(t)$ . Therefore, the variance of  
2448 KM-ARL can exceed that of LB-ARL when both event time and censoring have heavy tails and  $a$  is  
2449 large; otherwise, it may be smaller.  
2450

2451 Finally, the assumptions and background theory in Akritas (2000) are intricate and require careful  
2452 examination to validate this expression in QCD—one of the main challenges in developing our bias  
2453 theory. A full treatment would warrant a separate paper.  
2454

2455 **More about Heavy Censoring.** Demonstrating the numerical effect of bias under violations of  
2456 the independent censoring assumption would be helpful for practitioners. According to Fig. 2 in  
2457 (Malmquist, 2025), under a constant hazard rate and with the number of subjects (sequences) = 150,  
2458 it is reported that the mean squared error of the survival function over time is approximately given by:

- 2459 • For independent censoring (uniform censoring):
  - 2460 – < 0.01 for 10% censoring.
  - 2461 – < 0.01 for 50% censoring.
  - 2462 – < 0.01 for 90% censoring.
- 2463 • For dependent censoring (occurring just before event time (detection)):
  - 2464 – < 0.01 for 10% censoring.
  - 2465 –  $\approx 0.4$  for 50% censoring.
  - 2466 –  $\approx 2.1$  for 90% censoring.

2467 This result shows that heavy, dependent censoring inflates the bias. It aligns with our statement in  
2468 Sec. 6, where we clarify the appropriate use cases for KM-ARL and KM-ADD.  
2469

## 2470 F.1 MORE ABOUT INDEPENDENT CENSORING ASSUMPTION

2471 We extend our discussion in Sec. 6 about potential alleviation of the independent censoring assump-  
2472 tion.  
2473

### 2474 F.1.1 BACKGROUND: RESTRICTED MEAN SURVIVAL TIME (RMST)

2475 To mitigate the bias arising from dependent censoring and finite time interval, we must quantify the  
2476 effect of the bias on the survival function or its integral within a finite time interval, known as the  
2477 restricted mean survival time (RMST) in survival analysis. It is the counterpart of KM-ARL and  
2478 KM-ADD in survival analysis and is defined as the expected survival time up to a specified horizon  
2479 (denoted by  $h$  here), equivalently the area under the survival curve from 0 to  $h$ . It has been known  
2480 that the RMST is systematically biased under dependent censoring (Klein & Moeschberger, 1984;  
2481 Rivest & Wells, 2001; Ebrahimi & Molefe, 2003).  
2482  
2483

2484 F.1.2 SOLUTIONS  
2485

2486 Many studies have developed mitigation techniques for bias arising from dependent censoring. The  
2487 following works, among others, suggest promising directions for deriving KM-ARL/ADD under  
2488 dependent (informative) censoring or alleviating the resulting bias.

- 2489 • (Ebrahimi & Molefe, 2003): Develops a consistent estimator of the survival function under  
2490 dependent censoring, which may extend to estimating ARL and ADD in QCD.
- 2491 • (Hsu & Taylor, 2010): Proposes a robust weighted KME that uses baseline prognostic  
2492 covariates to correct bias in the marginal survival function when censoring is dependent.
- 2493 • (Lin et al., 2023): Applies inverse probability of censoring weighting (IPCW); at each time  
2494  $t$ , each subject (patient) still under observation is weighted by the inverse of their estimated  
2495 probability of remaining uncensored up to  $t$ .
- 2496 • (Crommen et al., 2025): Summarizes several bias-correction strategies:
  - 2497 – Nonparametric marginal estimation under a known copula: Given a specified copula  
2498 for  $(\tau, C)$ , observable probabilities uniquely determine the marginals. Step-function  
2499 estimators are constructed that reduce to the KME under independence and yield a  
2500 consistent estimator under dependent censoring.
  - 2501 – Likelihood estimation in copula models: Corrects bias by modeling the joint distribution  
2502 of  $(\tau, C)$ ; parametric copulas enable joint likelihood maximization with identifiable  
2503 parameters and consistent, asymptotically normal estimates.
  - 2504 – Correction of regression functionals: Embeds dependent censoring in copula-based  
2505 Cox and related semi-parametric models to adjust hazard ratios, covariate effects, and  
2506 causal effects.
  - 2507 – Machine-learning and partial-identification approaches: Addresses bias when flexible  
2508 modeling or weaker assumptions are required, using both deep and non-deep learning  
2509 methods.
  - 2510

2511 **Conclusion.** To avoid technical complications, we adopt the independent censoring assumption,  
2512 a standard assumption in survival analysis, examine the conditions under which it holds in QCD,  
2513 and specify them in the main text and this appendix. A full analysis of the above studies and their  
2514 integration into QCD would warrant a separate paper and is indeed an ongoing research direction for  
2515 the authors. Nonetheless, we are happy to present our current preliminary insights here, as we believe  
2516 they offer useful starting points for future work and provide essential background for developing the  
2517 emerging link between survival analysis and QCD.

2518  
2519 F.2 RELEVANCE OF REQUIRING DATASETS WITH MULTIPLE SEQUENCES WITH  
2520 CHANGEPOINT LABELS

2521 In this paper, we focus on the situations where datasets of multiple sequences with changepoint labels  
2522 are available. This problem setting is common in machine learning. Although there are many studies  
2523 of QCD on unlabeled datasets or pure simulations, there are also many studies on labeled datasets.

2524  
2525 There are a wide variety of labeled datasets with multiple sequences and labeled changepoints (or  
2526 annotations that can be seen as changepoints).

- 2527 • Human sensing: WISDM Actitracker (Kwapisz et al., 2011) includes both human- and  
2528 machine-labeled changepoints. We use this real-world sensing data in our experiments.
- 2529 • Twitter sentiment-change detection: Twitter Data Stream (US Airline Sentiment) (Makone,  
2530 2016), a collection of tweets annotated with the tags of positive/negative/neutral, is used to  
2531 detect sentiment changes in social media data stream (Bouchikhi et al., 2019): e.g., a sudden  
2532 increase in negative tweets after an event.
- 2533 • Speaker change detection: Many speech datasets, such as DIHARD dataset (Ryant et al.,  
2534 2018) and IITG-MV phase 3 dataset (Haris et al., 2012), have speaker or language annota-  
2535 tions, which can be used for speaker change detection (Mishra & Prasanna, 2024).
- 2536 • Temporal action localization: THUMOS14 Jiang et al. (2014) and ActivityNet (Heilbron  
2537 et al., 2015) are standard benchmark datasets in temporal action localization (Wang et al.,

2538 2023a). They contain a number of videos with temporal annotations of human actions (e.g.,  
2539 Horseback riding, Diving, Long jump, etc.), where timestamps with action changes can be  
2540 seen as changepoints.

2541 • Earthquake detection: STEAD (Mousavi et al., 2019) and several datasets in SeisBench  
2542 (Woollam et al., 2022) offer datasets of seismic signals with annotation of event times.  
2543

2544 Furthermore, we can create labeled datasets from unlabeled real-world dataset by injecting change-  
2545 points:

2546 • Sound anomaly detection: Normal and abnormal audio data recorded from industry machine  
2547 are concatenated to synthesize changepoints (Gopalan et al., 2021). The base dataset is the  
2548 MIMI Dataset (Purohit et al., 2019).

2549 • Social event detection: Based on the unlabeled MIT Cellphone Data (Eagle & Pentland,  
2550 2006), annotations are defined by linking changepoints to calendar events (e.g., sponsor  
2551 meeting, Presidents Day, spring break, etc.) (Kei et al., 2025).

2552 • Social event detection: Similarly to (Kei et al., 2025), annotations are provided for Enron  
2553 Email Data (Klimt & Yang, 2004) based on actual corporate events (e.g., Federal Energy  
2554 Regulatory Commission’s decisions, the CEO’s public protests, the bankruptcy filing, etc.)  
2555 (Wang et al., 2023b).  
2556

2557 Of course, we can create labeled datasets via simulation (e.g., KDD Cup 1999 intrusion detection  
2558 dataset (Stolfo et al., 1999)).  
2559

2560 Thus, real-world applications include: human sensing and action recognition (Bishop, 2006; Wang  
2561 et al., 2023a), sentiment-change detection on social media data stream (Bouchikhi et al., 2019),  
2562 speaker/language change detection (Mishra & Prasanna, 2024), earthquake detection (Woollam et al.,  
2563 2022), sound anomaly detection (Gopalan et al., 2021), event detection from cellphone data or email  
2564 data (Kei et al., 2025; Wang et al., 2023b), and intrusion detection (Stolfo et al., 1999), among others.

2565 In summary, although our estimators cannot be used for unlabeled datasets, as noted in the main text,  
2566 they apply to many real-world tasks listed above, among others.

2567 **Our contributions to machine learning community.** Finally, we would like to note that analyzing  
2568 how we evaluate performance metrics has been of general interest in machine learning (Hernández-  
2569 Orallo et al., 2004; Int, 2005; 2006; Pardo-Fernández & Castro, 2025). Recent top-tier conferences  
2570 themselves emphasize this by requiring careful reporting of metrics in their Author Instructions and  
2571 Reviewer Instructions, underscoring that reproducible and interpretable evaluation is a first-class  
2572 concern, not an afterthought. Our work advances more robust, interpretable evaluation of QCD  
2573 models, enabling empirical, intuitive model selection, as stated in the Abstract and Experiment  
2574 sections.  
2575

## 2576 G USE OF LARGE LANGUAGE MODELS

2577 During our research, we use LLMs to polish writing, finding related work, and verify our proof.  
2578  
2579  
2580  
2581  
2582  
2583  
2584  
2585  
2586  
2587  
2588  
2589  
2590  
2591