# WISER: Segmenting watermarked region- an epidemic change-point perspective.

**Soham Bonnerjee**
Department of Statistics
University of Chicago
sohambonnerjee@uchicago.edu

**Sayar Karmakar**
Department of Statistics
University of Florida
sayarkarmakar@ufl.edu

**Subhrajyoty Roy**
Department of Statistics and Data Science
Washington University at St. Louis
roy.s@wustl.edu

## ABSTRACT

With the increasing popularity of large language models, concerns over content authenticity have led to the development of myriad watermarking schemes. These schemes can be used to detect a machine-generated text via an appropriate key, while being imperceptible to readers with no such keys. The corresponding detection mechanisms usually take the form of statistical hypothesis testing for the existence of watermarks, spurring extensive research in this direction. However, the finer-grained problem of identifying which segments of a mixed-source text are actually watermarked, is much less explored; the existing approaches either lack scalability or theoretical guarantees robust to paraphrase and post-editing. In this work, we introduce a unique perspective to such watermark segmentation problems through the lens of *epidemic change-points*. By highlighting the similarities as well as differences of these two problems, we motivate and propose WISER: a novel, computationally efficient, watermark segmentation algorithm. We theoretically validate our algorithm by deriving finite sample error-bounds, and establishing its consistency in detecting multiple watermarked segments in a single text. Complementing these theoretical results, our extensive numerical experiments show that WISER outperforms state-of-the-art baseline methods, both in terms of computational speed as well as accuracy, on various benchmark datasets embedded with diverse watermarking schemes. Our theoretical and empirical findings establish WISER as an effective tool for watermark localization in most settings. It also shows how insights from a classical statistical problem can lead to a theoretically valid and computationally efficient solution of a modern and pertinent problem.

## 1 INTRODUCTION

An unfortunate consequence of the exponential ascent of the Large Language Models (LLM), influencing all aspects of content creation, has been an increased propagation of synthetic texts across the internet. This has raised significant doubts for content authenticity and copyright infringement over multiple domains (Megías et al., 2022; Bender et al., 2021; Crothers et al., 2023; Liang et al., 2024; Milano et al., 2023; Radford et al., 2023; Chen & Shu, 2023; Woodcock, 2023), indicating an urgent need to distinguish human authorship from machine generation. "Watermarking methods" have been proposed (Christ et al., 2024; Aaronson, 2023), and widely adopted (Biden, 2023; Bartz & Hu, 2023) as a detection mechanism, embedding statistical signals into LLM-generated tokens that remain largely un-noticeable without additional information. The key insight into the watermark-based detection schemes is the use of the underlying randomness of LLM-generated outputs by incorporating pseudo-randomness into the text-generation process. When a third-party user publishes text potentially containing LLM-generated outputs with watermarks, the coupling between the LLM-generated text and the pseudo-random numbers serves as a signal that can be used for detecting the watermark. The knowledge of these pseudo-random numbers is imperative for the detection

arXiv:2509.21160v1 [stat.ML] 25 Sep 2025

mechanism to work, making the effect of watermarking un-traceable for general users, who usually do not have access to such "keys".

Such usefulness has stimulated a plethora of research proposing myriad watermarking schemes (Kirchenbauer et al., 2024; Fernandez et al., 2023; Golowich & Moitra, 2024; Hu et al., 2024; Wu et al., 2024; Zhao et al., 2025; 2024a; Liu & Bu, 2024; Zhu et al., 2024). Concurrently, much attention has landed on the pursuit of efficient, statistically valid detection schemes (Li et al., 2025a; Kuditipudi et al., 2024; Cai et al., 2024; Huang et al., 2023; Li et al., 2024a; Cai et al., 2025). These detection schemes usually employ the knowledge of the pseudo-random keys or deterministic hash functions to perform a composite-vs-composite test of hypotheses: $H_0$ : the entire text $\omega_{1:n}$ is unwatermarked (i.e. human generated), vs $H_1$ : the entire text $\omega_{1:n}$ is watermarked or $H_1'$ : the text $\omega_{1:n}$ contains watermarked segments. Interestingly, the literature on the more fine-grained problem of identifying/localizing the said watermarked segments, is relatively sparse; some of the available algorithms are painstakingly slow, unsuited to large texts. Moreover, to the best of our knowledge, no such algorithm to efficiently identify multiple watermarked segments has sufficient theoretical validity. This gap in the literature is also pointed out by Li et al. (2025b). In this paper, we propose WISER (**W**atermark **I**dentification via **S**egmenting **E**pidemic **R**egions): a *first-of-its-kind* computationally efficient and provably consistent algorithm to locate multiple watermarked segments from mixed-source input texts. Our method is inspired from the classical notion of *epidemic* change-points; this perspective is instrumental to both the theoretical validity and computational efficiency of our algorithm. We summarize our main contributions as follows.

## 1.1 Main Contributions

Our key contributions are as follows.

***Novel Perspective.*** In §2, we introduce a novel, *epidemic change-point* perspective to the watermark segmentation problem by exploiting an inherent property of the watermarking schemes; see Figure 1 below. Since the segments can occur anywhere, the interpretation as an epidemic change-point enables us to re-purpose some of the classical insights into a state-of-the-art algorithm to solve a modern problem. At the same time, as discussed in §2.2.2 and 2.2.3, the particular setting of watermark segmentation introduces new challenges, and makes our analysis significantly different from the usual change-point theory.
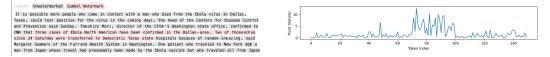


Figure 1: (Left) A mixed source text with watermarked tokens 70-100. (Right) The corresponding plot of pivot statistics vs. token.

**WISER algorithm.** The theoretical validity of the WISER segmentation algorithm arises as an automatic consequence of our perspective. In principle, our algorithm is simple to describe; the *epidemic* interpretation produces a natural estimate for the case of only one watermarked segment; the general case of multiple watermarked segments can then be dealt with by appropriately restricting the search spaces for each of these segments. The number of such segments is estimated by a series of carefully orchestrated steps, before further restriction on the search space is ensured to lessen the computational burden. The ingenuity of our algorithm is not only in its amalgamation of different ideas from statistics, but also in its practicability. We describe the algorithm in detail in Figure 2.

***Theoretical Contribution.*** Our proposed algorithm WISER is backed by the following key result.

**Theorem 1.1** (Informal version of Theorem 3.2). *Let $\hat{I}_j, j \in [\hat{K}]$ be the output of the WISER algorithm. With explicitly mentioned choices of the tuning parameters, under standard regularity conditions, it holds that $\liminf_{n\to\infty} \mathbb{P}\big(\hat{K} = K, \max_{k\in[K]} |\hat{I}_k \Delta I_k| \approx small\big) \approx 1$, where, $I_j, j \in [K]$ are the true watermarked segments; $\Delta$ is the symmetric difference operator, and $K$ and $\hat{K}$ are true and estimated number of segments, respectively.*

All the theoretical results are rigorously proved in Appendix §D. Additionally, we motivate the local estimate used in the last stage of WISER by proving in Theorem 3.1 that it is consistent in the single watermarked-segment case. To the best of our knowledge, WISER is the *first watermark segmentation algorithm with complete theoretical guarantees in the most general case*.

***Computational efficiency.*** In the numerical experiments performed in §4 and Appendix §C, the theoretical guarantee shines through in WISER's superiority over the other competitive methods across different watermarking schemes and different language models. Another key aspect of its enhanced performance is its speed. WISER is specifically designed with many localized steps that reduces its run-time, thereby making it the *only $O(n)$ watermark segmentation algorithm with provable theoretical guarantees*. We empirically also verify its speed-up in Figure 3. For a more comprehensive set of experiments and additional insights, we direct the readers to Appendix §C.

***Other contributions.*** We make some additional contributions that might be of independent interest. In terms of theory, the arbitrary dependence between the pivot statistics (introduced in §2) from the watermarked tokens, poses a significant hindrance to using the standard proof techniques from the change-point literature. Instead, we develop novel proof techniques based on moment and cumulant generating functions as well as Danskin (1967)'s results to conclude our proofs. On the application front, we address the inherent asymmetry of the watermark segmentation problem (see §2.2.3) by introducing a **M**odified **R**and **I**ndex (MRI). We argue that this provides a more accurate description of the performance of various algorithms. Due to space constraints, we have relegated both these discussions in the Appendix, in Sections D and C.1.1 respectively.

## 1.2 RELATED LITERATURE

There has been an abundance of literature on testing for existence of watermarks and the more general problems of machine-generated text detection or model equality testing (Lavergne et al., 2008; Solaiman et al., 2019; Gehrmann et al., 2019; Su et al., 2023; Mitchell et al., 2023; Huang et al., 2023; Vasilatos et al., 2023; Hans et al., 2024; Li et al., 2025a; Kuditipudi et al., 2024; Cai et al., 2024; Gao et al., 2025; Song et al., 2025; Radvand et al., 2025). However, the relatively harder problem of precisely localizing the watermarked segments from an input text has received only sparse attention. Apart from WinMax (Kirchenbauer et al., 2024), which focuses only on Red-Green watermarking, to the best of our knowledge, the only algorithms tackling the segmentation problem in its generality are Li et al. (2024b); Pan et al. (2025) and Zhao et al. (2024b). Most of these algorithms are prohibitively slow to be useful for long texts, while having little theoretical validity. In Appendix §C.1.3, we discuss the crucial limitations of each of these algorithms in contrast to WISER.

## 1.3 NOTATIONS

In this paper, we denote the set $\{1, \ldots, n\}$ by $[n]$. The $d$-dimensional Euclidean space is $\mathbb{R}^d$. We also denote in-probability convergence, and stochastic boundedness by $o_{\mathbb{P}}$ and $O_{\mathbb{P}}$, respectively. $\mathcal{L}(X)$ denotes the law of $X$. For any interval $I$, $I_L$ and $I_R$ denote its left and right end-point respectively.

## 2 WATERMARK LOCALIZATION: AN EPIDEMIC CHANGE-POINT PERSPECTIVE

Before we introduce our novel perspective in the context of locating watermarked segments, it is instrumental to establish a consistent framework of watermarking in LLM-generated texts. Let $\mathcal{W}$ denotes the dictionary, enumerated as $1, 2, \ldots, |\mathcal{W}|$. Given a text input in a tokenized form $\omega_1 \ldots \omega_{t-1}$, a watermarked LLM generates the next token $\omega_t$ in an autoregressive manner as $\omega_t = S(P_t, \zeta_t)$, where $P_t = (P_{t,w})_{w=1}^{|\mathcal{W}|}$ is the next token probability (NTP) distribution at step $t$; $S$ is a deterministic decoder function, and $\zeta_t$ is the pseudo-random variable at $t$. We grant Assumption 2.1 for the $\zeta_t$'s.

**Assumption 2.1.** *For any text $\omega_{1:n}$, there exists corresponding pseudo-random variables $\zeta_{1:n}$ available to the verifier, such that if the token $\omega_t$ at step $t$ is un-watermarked, then $\omega_t$ and $\zeta_t$ are independent conditional on $\omega_{1:(t-1)}$.*

It may seem that this assumption invalidates human edits after LLM generates a text. However, in Appendix §A, we discuss how Assumption 2.1 applies even to the mixed-source texts.

## 2.1 PIVOT STATISTICS AND ELEVATED ALTERNATIVES

Note that, a text $\omega_{1:n}$ with $K$ disjoint watermarked intervals $I_1, \ldots, I_K$, $I_j \subset [n]$ for $j \in [K]$, can be modeled as

$$w_t \sim \begin{cases} P_t, t \notin I_0 := \cup_{l=1}^K I_k \\ S(P_t, \zeta_t), \text{ otherwise,} \end{cases} \quad t = 1, 2, \ldots, n. \tag{2.1}$$

We are interested in the statistical problem of estimating the individual intervals $I_1, \ldots, I_K$ as well as $K$. Before proceeding further, it is appropriate to formally introduce the notion of pivot statistics.

**Definition 2.1.** *$Y(\omega, \zeta)$ is called a pivot statistic if $\mathcal{L}(Y)$ is same for all $\omega \in \mathcal{W}$.*

Pivot statistic has been extremely effective in providing statistically valid testing strategies for the existence of watermarks in mixed-source texts (Li et al., 2025a; 2024a; Cai et al., 2024), however, in what follows, we will demonstrate their effectiveness in aiding a localization algorithm. This effectiveness is a result of a simple property of the pivot statistics; they metamorphose the conditional independence of $\omega_t$ and $\zeta_t$ for un-watermarked tokens into $P_t$-independent distributions. Formally, this property is described in the following result.

**Lemma 2.2.** *If $S$ denotes the set of un-watermarked tokens, then $\{Y_t\}_{t \in S}$ are i.i.d.*

This ancillarity is heavily used in all the available statistical analysis of watermarked schemes; nevertheless, for the sake of completion we provide a proof in Appendix §D.3. Lemma 2.2 enables us to use the notation $\mu_0 := \mathbb{E}_0[Y(\omega, \zeta)]$ as the expectation of the pivot statistic $Y$ when the token $\omega \sim P$ is not watermarked; on the other hand, $\mathbb{E}_{1,P}[Y(\omega, \zeta)]$ will denote expectation with respect to the randomness of $\zeta$ (i.e. conditional on $P$) when $\omega$ is watermarked according to $(S, \zeta)$-mechanism. Finally, we denote $Y_t := Y_t(\omega_t, \zeta_t)$. Note that since $Y_t$ is a pivot statistic, so is $h(Y_t)$ for any *score function* $h : \mathbb{R} \to \mathbb{R}$. Usual tests for watermark detection look at $\sum_{t=1}^n h(Y_t)$ as a statistic for a one-sided test, and put considerable effort into constructing an effective score function $h$ (Kirchenbauer et al., 2024; Zhao et al., 2024b; Li et al., 2025a; Cai et al., 2025). Intrinsic to this construction, even though never explicitly stated, is the assumption that $\mathbb{E}_{1,P}[h(Y)]$ is usually larger than $\mu_0$ for any possible NTP distribution $P$. This hypothesis of "elevated alternatives" can also be empirically viewed in Figure 1.

We formalize this observation with the following hypothesis.

**Assumption 2.2** (Elevated Alternatives Hypothesis). *Assume that the next token distribution (NTP) $P$ belongs to a distribution class $\mathcal{P}$. Then, there exists $d > 0$ such that $\inf_{P \in \mathcal{P}} \mathbb{E}_{1,P}[h(Y)] \geq \mu_0 + d$, where $\mathbb{E}_{1,P}(\cdot) = \mathbb{E}_1[\cdot|P]$ denotes the unknown distribution of $h(Y)$ when watermarking is implemented on the NTP $P \in \mathcal{P}$.*

This assumption entails that the pivot statistics is effective conditional on any possible NTP from the class $\mathcal{P}$, ruling out trivial cases such as $Y(\omega, \zeta) \equiv \zeta$. Most standard watermarking schemes satisfy Assumption 2.2; see §B for some concrete examples. To summarize, the pivot statistics $Y_t$ has a mean level $\mu_0$ when the token $\omega_t$ is un-watermarked; on the other hand, we expect the pivot statistics to take comparatively larger values inside the watermarked segments. Interestingly, this observation establishes a ready-made connection to the notion of "epidemic change-points", sporadically explored in the classical time-series literature for the past few decades. We discuss this novel perspective in the following section.

## 2.2 WATERMARK AND EPIDEMIC CHANGE-POINT

We first provide some background on epidemic change-points, given their relative obscurity, for the convenience of readers who may be unfamiliar with the concept.

### 2.2.1 WHAT IS AN EPIDEMIC CHANGE-POINT?

An epidemic change-point refers to a situation where a stochastic process deviates in one of its features in an interval and returns to the baseline. The simplest and yet the most popular formulation of a 'mean-shift' epidemic model is as follows. Consider the time-series $X_i = \mu_i + Z_i$, where $Z_i$ is mean-zero stationary process and

$$\mu_i = \mu \text{ if } i \in \{1, \cdots, p\} \cup \{q+1, \cdots, n\} \text{ and } \mu_i = \mu + \delta \text{ if } i \in \{p+1, \cdots, q\} \tag{2.2}$$

The epidemic change-point framework originated with Levin & Kline (1985), who studied the testing for existence of such epidemic patches for epidemiology applications, with a more comprehensive discussion in Yao (1993); Inclán & Tiao (1994). Later on, Hušková (1995); Csörgö & Horváth (1997); Chen et al. (2016) have discussed consistency, asymptotic theory as well as statistical powers of these epidemic estimators and accompanying tests. Other related papers discussing inference tailored to epidemic alternatives can be found in Račkauskas & Suquet (2004; 2006); Ning et al. (2012). Compared to the vast literature for usual change-point analysis, the epidemic change-point literature has been quite sparse, and even then, the focus has remained mostly on testing for the existence of such temporary departure rather than on locating these patches with provable statistical guarantees.

### 2.2.2 Epidemic change-point with irregular signals

Note that, in the watermarked patches, it is unrealistic to assume a fixed mean of the pivot statistics, since the next token probability distribution usually changes at each step. Therefore, results pertaining to model (2.2) are not directly applicable here. However, invoking Assumption 2.2, we can assume that the means of the pivot statistics are separated from the null by at least some margin. This puts us in a position to solve an epidemic mean-shift problem of a new kind. Very recently Kley et al. (2024) proposed usual change-point detection under the presence of such irregular signals. Concretely, for noisy data of the form $X_t = \mu_t + Z_t$, $t = 1, \ldots, n$ where $\mu_t$ are means or signals and $(Z_t)_{t \in \mathbb{Z}}$ is a stationary mean-zero noise, they considered the following hypothesis testing problem with irregular 'non-constant-mean' alternative:

$$H_0 : \mu_1 = \cdots = \mu_n \text{ vs. } H_1 : \exists \tau \in \{2, \ldots, n\}, d > 0 : \mu_1 = \cdots = \mu_{\tau-1}, \quad \mu_\tau, \ldots, \mu_n \geq \mu_1 + d.$$

They also proposed an estimation procedure for the location parameter $\tau$. In this work, in the light of the mean pattern of the pivot statistic corresponding to the watermarked region, we extend their estimators to the epidemic alternative. Moreover, the intrinsic dependence introduced by the context of how an LLM token sequence is generated also makes our premise for the error specification quite novel and thus brings out significant technical challenges.

### 2.2.3 A subtle difference with change-point problem

Although watermark segmentation closely resembles epidemic change-point detection, a crucial difference arises in algorithm evaluation. Standard change-point problems are symmetric; under model (2.2), the edge cases $p = 1, q = n$ and $p = q$ are equivalent. On the other hand, watermarking problems exhibit asymmetry; the edge cases (i) "the entire sequence is un-watermarked" and (ii) "the entire sequence is watermarked", differ due to irregular means of the pivot statistics under watermarking. In fact, the widely popular Rand Index (RI) - being borrowed from clustering literature, and used in watermark segmentation (Li et al., 2024b; Pan et al., 2025) - fails to capture this distinction. For the interested readers, we address this by introducing a **M**odified **R**and **I**ndex (MRI), and demonstrate its advantages over RI in Appendix §C.1.1.

## 3 Theory for Watermark Localization

In this section, we develop our algorithm by proceeding step-by-step. In the §3.1, we propose an estimator to localize a single watermarked segment inside a text, and establish its theoretical consistency with finite sample results. Building on this estimate, in §3.2 we formally propose the WISER algorithm. Subsequently, we theoretically establish its consistency in segmenting multiple watermarked patches, while also discussing its linear-time computational complexity.

### 3.1 Segmenting single watermarked patch

Let us denote $X_t = h(Y_t)$. Recall Lemma 2.2, the notation $\mu_0 = \mathbb{E}_0 X_t$, and Assumption 2.2. Let $\tilde{d}$ be such that there exists $\rho \in (0, 1)$ satisfying $d > 2\rho\tilde{d}$. Based on our discussion in §2.2.2, we adapt the estimator from Kley et al. (2024) for our particular setting.

$$\hat{I} = \underset{s,t \in [n]}{\arg \min} \sum_{k \notin [s,t]} (X_k - \mu_0 - \rho\tilde{d}). \tag{3.1}$$

The following theorem analyzes its convergence properties for the case of a single, un-interrupted watermarked region. Subsequently, we discuss some of its connotations in successive remarks.

**Theorem 3.1.** *Let $\{X_t\}_{t=1}^n := \{h(Y_t)\}_{t=1}^n$ be the pivot statistics based on the given input text, and assume that $I_0 \subset \{1, \ldots, n\}$ is the only watermarked interval. Grant Assumption 2.2. Denote*

$$\varepsilon_t = \begin{cases} X_t - \mu_0, t \notin I_0, \\ X_t - \mu_t, \ \mu_t := \mathbb{E}_{1,P_t}[X_t], t \in I_0. \end{cases}$$

*Suppose the class of distributions $\mathcal{P}$ is closed and compact, and there exists $\eta > 0$ such that $\sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(\eta|\varepsilon|)] < \infty$. Moreover, assume that $\min\{\mathrm{Var}_0(\varepsilon), \sup_P \mathrm{Var}_{1,P}(\varepsilon)\} > 0$. If there exists a constant $c > 0$ such that $\tilde{d} \geq c$, then $|\hat{I} \Delta I_0| = O_{\mathbb{P}}(\tilde{d}^{-1})$. Here $\Delta$ is the symmetric difference operator and $O_{\mathbb{P}}$ hides constants independent of $n, \tilde{d}, \rho$, and $\mu_0$.*

The $O(\tilde{d}^{-1})$ rate can further be sharpened to $O(\tilde{d}^{-2})$ under a local sub-Gaussianity condition (see Proposition 1 in the Appendix §D ). In fact, under very mild conditions, Theorem 3.1 already tackles a more general scenario compared to the only other theoretical result available in a similar context (Li et al., 2024b). In contrast to a general watermarked patch, Li et al. (2024b) considered a specialized scenario, where only the first half of the text till an arbitrary point is watermarked, reducing the problem to a classical change-point setting.

The parameter $\tilde{d}$ serves as the *signal strength* in the convergence diagnostics of $\hat{I}$. It allows $\hat{I}$ to look for intervals such that the $\tilde{d}$-biased mean outside that interval is minimized. To ensure accuracy, $\tilde{d}$ has to be large, but $\tilde{d} \gg d$ might lead to overestimation. Since the minimum separation $d$ in Assumption 3.1 is typically unknown, it cannot be used directly. In most cases (see examples in Appendix §B), a distribution-dependent lower bound $d_L \leq d$ may be available, but relying on $\tilde{d} = d_L$ often sacrifices power, as $\inf_{t \in [n]} \mathbb{E}_{1,P_t}[X_t - \mu_0]$ is usually much larger. Thus, a key step in practice is a data-driven yet valid choice of $\tilde{d}$, which we discuss in §3.2. The tuning parameter $\rho$ adjusts the impact of $\tilde{d}$ and mitigates small errors in its selection. Choosing $\rho \approx 0$ is undesirable, as it causes $\hat{I}$ to overestimate $I$ due to fluctuations above $\mu_0$ under the null. Conversely, setting $\rho \approx 1$ can violate the requirement $d > 2\rho\tilde{d}$ when $\tilde{d}$ is large. Empirically, $\rho \in [0.1, 0.5]$ provides robust performance, and we revisit these choices in our discussion of WISER as well as the ablation studies in Appendix §C.3.

*Remark* 3.1 (Connection with other performance metric). Even though Theorem 3.1 controls the estimation error in terms of symmetric difference between estimated and true watermarked patches $\hat{I}$ and $I$ respectively, it is straightforward to transform this result in terms of the more familiar Intersection-Over-Union metric $\mathrm{IOU}(I, \hat{I}) = |I \cap \hat{I}|/|I \cup \hat{I}|$ as $1 - \mathrm{IOU}(I, \hat{I}) = \frac{|I \Delta \hat{I}|}{|I \cup \hat{I}|} = O_{\mathbb{P}}\left(\frac{1}{|I|\tilde{d}}\right)$. As the text size increases ($n \to \infty$), if $|I| = O(1)$, then the number of un-watermarked tokens is too large, overpowering the signal from the watermarked tokens. Under this "heavy-edit" regime, no non-trivial test statistic can differentiate between $H_0$ : the entire text $\omega_{1:n}$ is un-watermarked (i.e. human-generated) and $H_1$ : the entire text $\omega_{1:n}$ is watermarked, with reasonable power (Li et al., 2025b). The estimation being a harder problem than testing, it is therefore reasonable to assume $|I| \to \infty$ as $n \to \infty$. Therefore, Theorem 3.1 essentially entails that $\mathrm{IOU}(I, \hat{I}) \to 1$ as $n \to \infty$.

Despite the attractive theoretical properties of $\hat{I}$ given in (3.1), notwithstanding the yet unclear choice of $\tilde{d}$, there are a couple of practical roadblocks to deploying $\hat{I}$. Firstly, $\hat{I}$ has a computational complexity of $O(n^2)$, which is quite prohibitive for a large body of text one usually encounters. Secondly, it is not straightforward as to how $\hat{I}$ can be generalized to localize multiple watermarked segments. We answer these questions by proposing our WISER algorithm.

### 3.2 WISER: SEGMENTING MULTIPLE WATERMARKED PATCHES

The main motivation behind our proposed algorithm WISER is to use the estimator $\hat{I}$ on localized disjoint intervals that are more-or-less guaranteed to contain the true watermarked segments. Such intervals with guarantees are usually recovered as a consequence of some first-stage screening. For the convenience of readers, the detailed algorithm, along with a schematic diagram of WISER containing the key steps, is illustrated in Figure 2.

Subsequently, we make a mild assumption that the true watermarked segments have a minimum length, and are also well-separated to be considered as distinct segments. Formally, for two disjoint intervals $I_1 = (I_{1,L}, I_{1,R})$ and $I_2 = (I_{2,L}, I_{2,R})$, let $d(I_1, I_2) := \min\{|I_{1,L} - I_{2,R}|, |I_{1,R} - I_{2,L}|\}$.

**Assumption 3.1** (Minimum separation). *Let $K$ be the number of true watermarked segments, with the segments themselves denoted by $I_j, j \in [K]$. Then there exists a constant $C_0 > 0$, such that $\min_{k \in [K]} \{|I_k| \wedge d(I_k, I_{k-1})\} \geq C_0 n^{1/2+\gamma'} \log n$ for some $\gamma' > 0$.*

In what follows, we explain the step-by-step rationale behind the algorithm. For clarity, we ignore the niceties of $\lfloor \cdot \rfloor$'s and $\lceil \cdot \rceil$'s.

- **Blocking stage.** Let $b = \sqrt{n}$ and the threshold $\mathcal{Q}$ be given. In the first stage, we partition the data into $\sqrt{n}$ consecutive blocks, each of size $\sqrt{n}$. Among these, we retain only those blocks for which the corresponding sum of pivot statistics exceeds $\mathcal{Q}$. Typically, to avoid multiple testing issues, $\mathcal{Q}$ is chosen as the $(1 - \alpha)$-quantile of the *null* (i.e. when there is no watermarking in the entire text) distribution of the maximum block sum over all $\sqrt{n}$ blocks.
- **Discarding stage.** Under Assumption 3.1, by definition of $\mathcal{Q}$, $O(\sqrt{\log n})$ successive unwatermarked blocks will have sum exceeding $\mathcal{Q}$ *only* with vanishing probability. Therefore, any connected interval of selected blocks from the first stage, with length at most $c\sqrt{n \log n}$, must necessarily be spurious. Hence, at this stage, we join consecutive selected blocks, and discard any connected intervals smaller than $c\sqrt{n \log n}$.
- **Enlargement stage.** The above two steps ensure $\hat{K} = K$ with probability approaching 1. Also, the intervals from the previous stage are almost accurate estimates of the true segments, except for some additional watermarked regions that were part of discarded blocks. Because of Assumption 3.1 and the size of the discarded blocks, such regions have size at most $O(\sqrt{n})$. Therefore, we enlarge each interval by $\asymp n^{1/2}$ for a small $0 < \gamma \ll 1/2$. These enlarged intervals $D_j$'s remain disjoint with high probability due to Assumption 3.1, and are therefore each amenable to (3.1) to yield $\hat{I}_j$'s.
- **Estimating $\tilde{d}$.** To estimate $\tilde{d}$, we take the sample mean of $(X_t - \mu_0)$ over $\cup_{j=1}^{\hat{K}} D_j$. This serves as a proxy for the oracle average of $(X_t - \mu_0)$ over $\cup_{j \in [K]} I_j$, which may overestimate $d$. We choose $\rho$ to calibrate it so that $d > 2\rho\tilde{d}$.
- **Reducing computational cost.** We alleviate the increased computational aspect of a naive implementation of (3.1) by leveraging additional information from the screening stage to reduce the search space. Indeed, due to our blocking and discarding steps, it can be guaranteed with high probability that, for each $j \in [K]$, $D_{j,L}$ is at most $\asymp \sqrt{n}$ distance apart from $I_{j,L}$; similarly $D_{j,R}$ is also at most $\asymp \sqrt{n}$ distance apart from $I_{j,R}$. Therefore, from $D_j$ we can produce search intervals $L_j, R_j$ of lengths $\asymp n^{1/2}$ such that $I_{j,L} \in L_j$ ad $I_{j,R} \in R_j$ with high probability, and restrict the search to $s \in L_j, t \in R_j$. Consequently, now each implementation of this modified (3.1) (see Figure 2) takes $O((n^{1/2})^2) = O(n)$ amount of computational time, leading to a speed-up while maintaining theoretical validity.

The following result summarizes these insights into a formal consistency guarantee.

**Theorem 3.2.** *Assume that the null distribution of the pivot statistics is absolutely continuous with respect to the Lebesgue measure. Fix $\alpha \in (0, 1)$, and recall the quantities defined in* WISER *described in Figure 2. Let the block length $b = b_n$ satisfy $b_n \asymp \sqrt{n}$, and suppose the threshold $\mathcal{Q} = \mathcal{Q}_n$ is selected so that $\mathbb{P}_0(\max_{1 \leq k \leq \lceil n/b \rceil} S_k > \mathcal{Q}) = \alpha$. Also assume that $\mathbb{E}_0[|X - \mu_0|^{3+\delta}] < \infty$ for some $\delta > 0$. Let $\sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[X] < \infty$, and assume there exists $\tau > 0$ such that*

$$\kappa := \inf_{\theta \geq 0} \theta(\mu_0 + \tau d) + \log \sup_P \mathbb{E}_{1,p}[\exp(-\theta X)] < 0. \tag{3.2}$$

*Additionally, let the number of watermarked intervals $K$ be bounded, and Assumption 3.1 be granted for the watermarked intervals $I_k, k \in [K]$. Then, given $\varepsilon > 0$ and $d \geq c$ for some constant $c > 0$, under the assumptions of Theorem 3.1, there exists $M_\varepsilon \in \mathbb{R}_+$, independent of $n, K$, and $d$, such that,*

$$\liminf_{n \to \infty} \mathbb{P}\big(\hat{K} = K, \max_{k \in [K]} |\hat{I}_k \Delta I_k| < M_\varepsilon d^{-1}\big) \geq 1 - \varepsilon. \tag{3.3}$$

*Remark* 3.2. We briefly discuss arguably the only technical condition (3.2) in Theorem 3.2. This can be construed as a Donsker-Varadhan strengthened version of Assumption 2.2. For an appropriate choice of the score function $h$ and some NTP distribution $P^\star$ depending on $\mathcal{P}$, the Donsker Varadhan representation (Donsker & Varadhan, 1983) entails

$$\inf_{\theta \geq 0} \theta\mu_0 + \log \sup_{P \in \mathcal{P}} \mathbb{E}_{1,p}[\exp(-\theta X)] = -D_{\mathrm{KL}}(\mathcal{L}_0(X), \mathcal{L}_{1,P^\star}(X)),$$

where $D_{\mathrm{KL}}$ denotes the Kullback-Leibler divergence, $\mathcal{L}_0$ denotes the law of un-watermarked pivot statistics, and $\mathcal{L}_{1,P^\star}$ denotes the law of watermarked pivot statistics when the NTP is $P^\star$. In light of
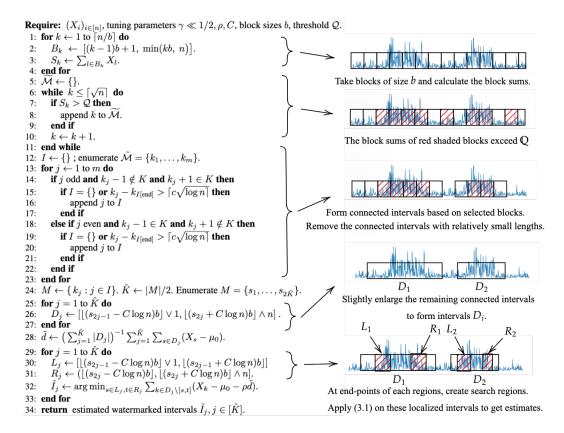
**Require:** $(X_i)_{i\in[n]}$, tuning parameters $\gamma \ll 1/2, \rho, C$, block sizes $b$, threshold $\mathcal{Q}$.
1: **for** $k \leftarrow 1$ to $\lceil n/b \rceil$ **do**
2:   $B_k \leftarrow [(k-1)b+1, \min(kb, n)]$.
3:   $S_k \leftarrow \sum_{l\in B_k} X_l$.
4: **end for**
5: $\widetilde{\mathcal{M}} \leftarrow \{\}$.
6: **while** $k \le \lceil\sqrt{n}\rceil$ **do**
7:   **if** $S_k > \mathcal{Q}$ **then**
8:     append $k$ to $\widetilde{\mathcal{M}}$.
9:   **end if**
10:   $k \leftarrow k+1$.
11: **end while**
12: $I \leftarrow \{\}$ ; enumerate $\tilde{\mathcal{M}} = \{k_1, \ldots, k_m\}$.
13: **for** $j \leftarrow 1$ to $m$ **do**
14:   **if** $j$ odd and $k_j - 1 \notin K$ and $k_j + 1 \in K$ **then**
15:     **if** $I = \{\}$ or $k_j - k_{I[\text{end}]} > \lceil c\sqrt{\log n}\rceil$ **then**
16:       append $j$ to $I$
17:     **end if**
18:   **else if** $j$ even and $k_j - 1 \in K$ and $k_j + 1 \notin K$ **then**
19:     **if** $I = \{\}$ or $k_j - k_{I[\text{end}]} > \lceil c\sqrt{\log n}\rceil$ **then**
20:       append $j$ to $I$
21:     **end if**
22:   **end if**
23: **end for**
24: $M \leftarrow \{k_j : j \in I\}$. $\hat{K} \leftarrow |M|/2$. Enumerate $M = \{s_1, \ldots, s_{2\hat{K}}\}$.
25: **for** $j \leftarrow 1$ to $\hat{K}$ **do**
26:   $D_j \leftarrow [\lfloor(s_{2j-1} - C\log n)b\rfloor \vee 1, \lfloor(s_{2j} + C\log n)b\rfloor \wedge n]$.
27: **end for**
28: $\tilde{d} \leftarrow \left(\sum_{j=1}^{\hat{K}} |D_j|\right)^{-1} \sum_{j=1}^{\hat{K}} \sum_{s\in D_j}(X_s - \mu_0)$.
29: **for** $j \leftarrow 1$ to $\hat{K}$ **do**
30:   $L_j \leftarrow [\lfloor(s_{2j-1} - C\log n)b\rfloor \vee 1, \lfloor(s_{2j-1} + C\log n)b\rfloor]$
31:   $R_j \leftarrow (\lfloor(s_{2j} - C\log n)b\rfloor, \lfloor(s_{2j} + C\log n)b\rfloor \wedge n]$.
32:   $\hat{I}_j \leftarrow \arg\min_{s\in L_j, t\in R_j} \sum_{k\in D_j\setminus[s,t]}(X_k - \mu_0 - \rho\tilde{d})$.
33: **end for**
34: **return** estimated watermarked intervals $\hat{I}_j, j \in [\hat{K}]$.

Take blocks of size $b$ and calculate the block sums.

The block sums of red shaded blocks exceed $\mathbb{Q}$

Form connected intervals based on selected blocks. Remove the connected intervals with relatively small lengths.

$D_1$   $D_2$

Slightly enlarge the remaining connected intervals to form intervals $D_i$.

$L_1$      $R_1$ $L_2$      $R_2$

$D_1$            $D_2$

At end-points of each regions, create search regions. Apply (3.1) on these localized intervals to get estimates.

Figure 2: (Left): The Algorithm WISER; (Right) WISER in action with key steps.

this, $\kappa$ lifts the minimum separation between the un-watermarked and watermarked distributions into a gap between the cumulant functions, and can therefore be understood to be mild. Equation (3.2) establishes a weak uniform control over the behavior of pivot statistics under watermarked segments. This allows us to rigorously bypass the possibly arbitrary and strong dependence across the pivot statistics corresponding to watermarked tokens while deriving Theorem 3.2.

We reiterate that with $b \asymp \sqrt{n}$, WISER has a run-time only of $O(n)$ ignoring log factors. This, to the best of our knowledge, is among the *least computationally expensive* algorithms available in the literature. In view of its theoretical validity under very general conditions, this makes it a useful tool for practical applications. We showcase it through a series of extensive numerical experiments.

## 4 NUMERICAL EXPERIMENTS

Building on the theoretical validation established in the previous sections, in this section we undertake an empirical evaluation of the proposed WISER method, demonstrating its superiority over existing state-of-the-art (SOTA) algorithms. In §4.1, we compare its accuracy against competitive methods on a benchmark dataset across multiple watermarking schemes, and in §4.2, we assess its computational efficiency. Due to space constraints, we provide additional numerical experiments in Appendix §C. We encourage the readers to check it out for more practical insights, including, **(i)** a detailed explanation of the benefits of WISER over other SOTA algorithms (§C.1.3), **(ii)** experiments quantifying the effect of watermark intensity and length across different algorithms (§C.2), and **(iii)** an ablation study (§C.3) highlighting the stability of our method across tuning parameter choices. The datasets and the large language models were acquired from the open-source Huggingface library. All the relevant reproducible codes and figures, as well as the generated datasets can be found in the Github repository.

8

## 4.1 COMPARATIVE PERFORMANCE OF WISER

Within the relatively limited body of literature on the identification of watermarked segments from mixed-source texts, `Aligator` (Zhao et al., 2024b), `SeedBS-NOT` (Li et al., 2024b) and `Waterseeker` (Pan et al., 2025) algorithms have emerged as the leading methods, producing the most accurate results so far. For an extensive comparison, our experimental setup involves completion of randomly selected 200 prompts from the Google C4 news dataset[1]. We include language models spanning a wide range of scales: parameter sizes varying from 125 million to 8 billion, and vocabulary sizes ranging in 32-262 thousands; for watermarking schemes, we consider Gumbel-max trick (Aaronson, 2023), Inverse transform (Kuditipudi et al., 2024), Red-green watermark (Kirchenbauer et al., 2023) and Permute-and-Flip watermark (Zhao et al., 2025). In each scenario, the first 50 tokens of a news article have been provided as inputs to the language models, and $n = 500$ output tokens are recorded. Among these 500 output tokens, there are two watermarked segments: 100-200 and 325-400. The specific tuning parameter choices for WISER are provided in §C. Table 1 showcases the results for the Gumbel watermarking scheme. It is evident that WISER outperforms all the other algorithms across all the metrics for each model. The detailed discussion, including the specific metrics used and additional results and insights, are provided in Appendix §C.1.

| Model Name | Vocab Size | Method | IOU | Precision | Recall | F1 | RI | MRI |
|---|---|---|---|---|---|---|---|---|
| facebook/opt-125m | 50272 | WISER | **0.944** | **1.000** | **0.995** | **0.997** | **0.984** | **0.979** |
| | | Aligator | 0.734 | 0.382 | 0.988 | 0.551 | 0.939 | 0.931 |
| | | Waterseeker | 0.672 | 1.000 | 0.802 | 0.890 | 0.864 | 0.850 |
| | | SeedBS-NOT | 0.479 | 0.730 | 0.625 | 0.673 | 0.844 | 0.823 |
| google/gemma-3-270m | 262144 | WISER | **0.896** | **0.965** | **0.960** | **0.962** | **0.953** | **0.950** |
| | | Aligator | 0.506 | 0.234 | 0.912 | 0.373 | 0.881 | 0.861 |
| | | Waterseeker | 0.645 | 0.968 | 0.775 | 0.861 | 0.851 | 0.836 |
| | | SeedBS-NOT | 0.362 | 0.610 | 0.478 | 0.536 | 0.753 | 0.704 |
| facebook/opt-1.3b | 50272 | WISER | **0.934** | **1.000** | **0.995** | **0.997** | **0.981** | **0.974** |
| | | Aligator | 0.497 | 0.235 | 0.920 | 0.375 | 0.892 | 0.871 |
| | | Waterseeker | 0.657 | 1.000 | 0.808 | 0.893 | 0.860 | 0.846 |
| | | SeedBS-NOT | 0.360 | 0.618 | 0.465 | 0.531 | 0.766 | 0.731 |
| princeton-nlp/Sheared-LLaMA-1.3B | 32000 | WISER | **0.939** | **1.000** | **0.998** | **0.999** | **0.983** | **0.978** |
| | | Aligator | 0.459 | 0.236 | 0.912 | 0.376 | 0.886 | 0.862 |
| | | Waterseeker | 0.659 | 1.000 | 0.812 | 0.897 | 0.862 | 0.847 |
| | | SeedBS-NOT | 0.278 | 0.520 | 0.388 | 0.444 | 0.731 | 0.699 |
| mistralai/Mistral-7B-v0.1 | 32000 | WISER | **0.909** | **1.000** | **0.998** | **0.999** | **0.975** | **0.961** |
| | | Aligator | 0.292 | 0.215 | 0.745 | 0.334 | 0.811 | 0.774 |
| | | Waterseeker | 0.621 | 1.000 | 0.765 | 0.867 | 0.840 | 0.824 |
| | | SeedBS-NOT | 0.240 | 0.442 | 0.320 | 0.371 | 0.657 | 0.593 |
| meta-llama/Meta-Llama-3-8B | 128256 | WISER | **0.926** | **1.000** | **0.988** | **0.994** | **0.977** | **0.975** |
| | | Aligator | 0.546 | 0.367 | 0.925 | 0.525 | 0.911 | 0.891 |
| | | Waterseeker | 0.570 | 1.000 | 0.720 | 0.837 | 0.814 | 0.791 |
| | | SeedBS-NOT | 0.379 | 0.620 | 0.515 | 0.563 | 0.778 | 0.741 |

Table 1: Results for Gumbel Watermarking

## 4.2 TIME COMPARISON

As established in §3.2, the proposed WISER algorithm achieves a computational complexity of $\approx O(n)$. Figure 3 provides empirical evidence supporting this theoretical claim and, in addition, compares the runtime behavior of WISER with other state-of-the-art methods. For this experiment, we randomly create an $n/6$-length watermarked segment using the Gumbel-max trick with NTP generated by Google's Gemma-3 model; block size was taken as $\lceil \sqrt{n} \rceil$ and $\rho = 0.1$. The results clearly indicate that WISER consistently outperforms competing approaches in terms of computational efficiency, emerging as the fastest among all methods considered in this study.

## 5 CONCLUDING REMARKS

In this paper, we introduced WISER, a first-of-its-kind algorithm for efficient and theoretically valid segmentation of watermarked intervals in mixed-source texts. By framing watermark localization as

---

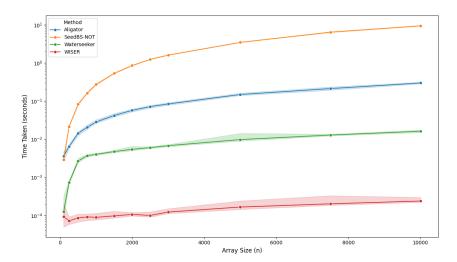[1]https://www.tensorflow.org/datasets/catalog/c4

Figure 3: Time complexity (seconds) for various algorithms as a function of completion lengths ($n$). Y-axis is in log-scale, with $95\%$ confidence interval shown in shades.

an epidemic change-point problem, we bridged a novel connection between classical statistical theory and a modern challenge in generative AI, and also designed a linear time algorithm with provable consistency guarantees, which were further confirmed by our extensive numerical experiments. Beyond the findings of this paper, it is also crucial to theoretically investigate the robustness of the proposed algorithm under human edits (Li et al., 2024a); as a roadmap, we have already included some relevant discussion in Appendix §A. Its applicability to multimodal (e.g. audio, image, video) settings (Qiu et al., 2024) also presents opportunities for future research.

## AUTHOR CONTRIBUTIONS

All the authors contributed equally to this research.

## REFERENCES

Scott Aaronson. Watermarking of large language models. https://simons.berkeley.edu/talks/scottaaronson-ut-austin-openai-2023-08-17, August 2023. Talk at the Simons Institute for the Theory of Computing.

Jushan Bai. Least squares estimation of a shift in linear processes. *J. Time Ser. Anal.*, 15(5): 453–472, 1994. ISSN 0143-9782,1467-9892. doi: 10.1111/j.1467-9892.1994.tb00204.x. URL https://doi.org/10.1111/j.1467-9892.1994.tb00204.x.

Diane Bartz and Krystal Hu. Openai, google, others pledge to watermark ai content for safety, white house says. https://www.reuters.com/technology/openai-google-others-pledge-watermark-ai-content-safety-white-house-2023-07-21/, 2023. Accessed: 2023-10-03.

Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 610–623, 2021.

Joseph R. Biden. Fact sheet: President biden issues executive order on safe, secure, and trustworthy artificial intelligence. https://bidenwhitehouse.archives.gov/briefing-room/statements-releases/2023/10/30/fact-sheet-president-biden-issues-executive-order-on-safe-secure-and-trustworthy-a October 2023. The White House, October 30, 2023.

Soham Bonnerjee, Sayar Karmakar, Maggie Cheng, and Wei Biao Wu. Testing synchronization of change-points for multiple time series. *Preprint*, 2025. URL https://sohamb01.github.io/drafts/test-of-synchronization.pdf.

Yinpeng Cai, Lexin Li, and Linjun Zhang. A statistical hypothesis testing framework for data misappropriation detection in large language models. *arXiv preprint arXiv:2501.02441*, 2025.

Zhongze Cai, Shang Liu, Hanzhao Wang, Huaiyang Zhong, and Xiaocheng Li. Towards better statistical understanding of watermarking llms. *arXiv preprint arXiv:2403.13027*, 2024.

Canyu Chen and Kai Shu. Can llm-generated misinformation be detected? *arXiv preprint arXiv:2309.13788*, 2023.

Zhenmin Chen, Zihao Li, and Min Zhou. Detecting change-points in epidemic models. *Journal of advanced statistics*, 1(4):181, 2016.

Miranda Christ, Sam Gunn, and Or Zamir. Undetectable watermarks for language models. In *The Thirty Seventh Annual Conference on Learning Theory*, pp. 1125–1139. PMLR, 2024.

Evan N Crothers, Nathalie Japkowicz, and Herna L Viktor. Machine-generated text: A comprehensive survey of threat models and detection methods. *IEEE Access*, 11:70977–71002, 2023.

Miklós Csörgö and Lajos Horváth. *Limit theorems in change-point analysis*. 1997.

John M. Danskin. *The theory of max-min and its application to weapons allocation problems*, volume V of *Econometrics and Operations Research*. Springer-Verlag New York, Inc., New York, 1967.

Nabarun Deb, Promit Ghosal, and Bodhisattva Sen. Measuring association on topological spaces using kernels and geometric graphs. *arXiv preprint arXiv:2010.01768*, 2020.

Monroe D. Donsker and S. R. S. Varadhan. Asymptotic evaluation of certain markov process expectations for large time. IV. *Comm. Pure Appl. Math.*, 36(2):183–212, 1983. ISSN 0010-3640. doi: 10.1002/cpa.3160360204. URL https://doi.org/10.1002/cpa.3160360204.

Pierre Fernandez, Antoine Chaffin, Karim Tit, Vivien Chappelier, and Teddy Furon. Three bricks to consolidate watermarks for large language models. In *2023 IEEE international workshop on information forensics and security (WIFS)*, pp. 1–6. IEEE, 2023.

D Kh Fuk and Sergey V Nagaev. Probability inequalities for sums of independent random variables. *Theory of Probability & Its Applications*, 16(4):643–660, 1971.

Irena Gao, Percy Liang, and Carlos Guestrin. Model equality testing: Which model is this API serving? In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=QCDdI7X3f9.

Sebastian Gehrmann, Hendrik Strobelt, and Alexander M Rush. Gltr: Statistical detection and visualization of generated text. *arXiv preprint arXiv:1906.04043*, 2019.

Noah Golowich and Ankur Moitra. Edit distance robust watermarks via indexing pseudorandom codes. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=FZ45kf5pIA.

J. Hájek and A. Rényi. Generalization of an inequality of Kolmogorov. *Acta Math. Acad. Sci. Hungar.*, 6:281–283, 1955. ISSN 0001-5954,1588-2632. doi: 10.1007/BF02024392. URL https://doi.org/10.1007/BF02024392.

P. Hall and C. C. Heyde. *Martingale limit theory and its application*. Probability and Mathematical Statistics. Academic Press, Inc. [Harcourt Brace Jovanovich, Publishers], New York-London, 1980. ISBN 0-12-319350-8.

Abhimanyu Hans, Avi Schwarzschild, Valeriia Cherepanova, Hamid Kazemi, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Spotting llms with binoculars: Zero-shot detection of machine-generated text. *arXiv preprint arXiv:2401.12070*, 2024.

Zhengmian Hu, Lichang Chen, Xidong Wu, Yihan Wu, Hongyang Zhang, and Heng Huang. Unbiased watermark for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=uWVC5FVidc.

Baihe Huang, Banghua Zhu, Hanlin Zhu, Jason Lee, Jiantao Jiao, and Michael Jordan. Towards optimal statistical watermarking. In *Socially Responsible Language Modelling Research*, 2023. URL https://openreview.net/forum?id=Fc2FaS9mYJ.

Marie Hušková. Estimators for epidemic alternatives. *Commentationes Mathematicae Universitatis Carolinae*, 36(2):279–291, 1995.

Carmen Inclán and George C Tiao. Use of cumulative sums of squares for retrospective detection of changes of variance. *Journal of the American Statistical Association*, 89(427):913–923, 1994. doi: 10.1080/01621459.1994.10476824.

John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A watermark for large language models. In *International Conference on Machine Learning*, pp. 17061–17084. PMLR, 2023.

John Kirchenbauer, Jonas Geiping, Yuxin Wen, Manli Shu, Khalid Saifullah, Kezhi Kong, Kasun Fernando, Aniruddha Saha, Micah Goldblum, and Tom Goldstein. On the reliability of watermarks for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=DEJIDCmWOz.

Tobias Kley, Yuhan Philip Liu, Hongyuan Cao, and Wei Biao Wu. Change-point analysis with irregular signals. *The Annals of Statistics*, 52(6):2913–2930, 2024.

Rohith Kuditipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. Robust distortion-free watermarks for language models. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=FpaCL1MO2C.

Thomas Lavergne, Tanguy Urvoy, and François Yvon. Detecting fake content with relative entropy scoring. *Pan*, 8(27-31):4, 2008.

Bruce Levin and Jennie Kline. The cusum test of homogeneity with an application in spontaneous abortion epidemiology. *Statistics in Medicine*, 4(4):469–488, 1985.

Xiang Li, Feng Ruan, Huiyuan Wang, Qi Long, and Weijie J Su. Robust detection of watermarks for large language models under human edits. *arXiv preprint arXiv:2411.13868*, 2024a.

Xiang Li, Feng Ruan, Huiyuan Wang, Qi Long, and Weijie J Su. A statistical framework of watermarks for large language models: Pivot, detection efficiency and optimal rules. *The Annals of Statistics*, 53(1):322–351, 2025a.

Xiang Li, Garrett Wen, Weiqing He, Jiayuan Wu, Qi Long, and Weijie J Su. Optimal estimation of watermark proportions in hybrid ai-human texts. *arXiv preprint arXiv:2506.22343*, 2025b.

Xingchi Li, Guanxun Li, and Xianyang Zhang. Segmenting watermarked texts from language models. *Advances in Neural Information Processing Systems*, 37:14634–14665, 2024b.

Weixin Liang, Zachary Izzo, Yaohui Zhang, Haley Lepp, Hancheng Cao, Xuandong Zhao, Lingjiao Chen, Haotian Ye, Sheng Liu, Zhi Huang, et al. Monitoring ai-modified content at scale: A case study on the impact of chatgpt on ai conference peer reviews. *arXiv preprint arXiv:2403.07183*, 2024.

Yepeng Liu and Yuheng Bu. Adaptive text watermark for large language models. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024.

David Megías, Minoru Kuribayashi, Andrea Rosales, Krzysztof Cabaj, and Wojciech Mazurczyk. Architecture of a fake news detection system combining digital watermarking, signal processing, and machine learning. *Journal of Wireless Mobile Networks, Ubiquitous Computing, and Dependable Applications (JoWUA), 2022, 13 (1): 33-55,*, 2022.

Silvia Milano, Joshua A McGrane, and Sabina Leonelli. Large language models challenge the future of higher education. *Nature Machine Intelligence*, 5(4):333–334, 2023.

Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. Detectgpt: Zero-shot machine-generated text detection using probability curvature. In *International conference on machine learning*, pp. 24950–24962. PMLR, 2023.

Wei Ning, Junvie Pailden, and Arjun Gupta. Empirical likelihood ratio test for the epidemic change model. *Journal of Data science*, 10(1):107–127, 2012.

Leyi Pan, Aiwei Liu, Yijian LU, Zitian Gao, Yichen Di, Shiyu Huang, Lijie Wen, Irwin King, and Philip S. Yu. Waterseeker: Pioneering efficient detection of watermarked segments in large documents. In *AAAI 2025 Workshop on Preventing and Detecting LLM Misinformation (PDLM)*, 2025. URL https://openreview.net/forum?id=3dslkUEgJb.

Lucas de Oliveira Prates. A more efficient algorithm to compute the rand index for change-point problems. *arXiv preprint arXiv:2112.03738*, 2021.

Jielin Qiu, William Han, Xuandong Zhao, Shangbang Long, Christos Faloutsos, and Lei Li. Evaluating durability: Benchmark insights into multimodal watermarking. *CoRR*, abs/2406.03728, 2024. URL https://doi.org/10.48550/arXiv.2406.03728.

Alfredas Račkauskas and Charles Suquet. Hölder norm test statistics for epidemic change. *Journal of statistical planning and inference*, 126(2):495–520, 2004.

Alfredas Račkauskas and Charles Suquet. Testing epidemic changes of infinite dimensional parameters. *Statistical Inference for Stochastic Processes*, 9(2):111–134, 2006.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pp. 28492–28518. PMLR, 2023.

Tara Radvand, Mojtaba Abdolmaleki, Mohamed Mostagir, and Ambuj Tewari. Zero-shot statistical tests for llm-generated text detection using finite sample concentration inequalities. *arXiv preprint arXiv:2501.02406*, 2025.

Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, et al. Release strategies and the social impacts of language models. *arXiv preprint arXiv:1908.09203*, 2019.

Yiliao Song, Zhenqiao Yuan, Shuhai Zhang, Zhen Fang, Jun Yu, and Feng Liu. Deep kernel relative test for machine-generated text detection. In *The Thirteenth International Conference on Learning Representations*, 2025.

Jinyan Su, Terry Yue Zhuo, Di Wang, and Preslav Nakov. DetectLLM: Leveraging log rank information for zero-shot detection of machine-generated text. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL https://openreview.net/forum?id=Dy2mbQIdMz.

Christoforos Vasilatos, Manaar Alam, Talal Rahwan, Yasir Zaki, and Michail Maniatakos. Howkgpt: Investigating the detection of chatgpt-generated university student homework through context-aware perplexity analysis. *arXiv preprint arXiv:2305.18226*, 2023.

Claire Woodcock. Ai is tearing wikipedia apart, May 2023. URL https://www.vice.com/en/article/ai-is-tearing-wikipedia-apart/. Accessed: 2025-09-14.

Yihan Wu, Zhengmian Hu, Junfeng Guo, Hongyang Zhang, and Heng Huang. A resilient and accessible distribution-preserving watermark for large language models. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024.

Qiwei Yao. Tests for change-points with epidemic alternatives. *Biometrika*, 80(1):179–191, 1993.

Xuandong Zhao, Prabhanjan Vijendra Ananth, Lei Li, and Yu-Xiang Wang. Provable robust watermarking for AI-generated text. In *The Twelfth International Conference on Learning Representations*, 2024a. URL https://openreview.net/forum?id=SsmT8aO45L.

Xuandong Zhao, Chenwen Liao, Yu-Xiang Wang, and Lei Li. Efficiently identifying watermarked segments in mixed-source texts. In *Neurips Safe Generative AI Workshop 2024*, 2024b.

Xuandong Zhao, Lei Li, and Yu-Xiang Wang. Permute-and-flip: An optimally stable and watermarkable decoder for LLMs. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=YyVVicZ32M.

Chaoyi Zhu, Jeroen Galjaard, Pin-Yu Chen, and Lydia Chen. Duwak: Dual watermarks in large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 11416–11436, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.678. URL https://aclanthology.org/2024.findings-acl.678/.

# Appendix

This appendix is devoted to further elaboration on the key ideas of the texts, detailed proofs of our theoretical statements, and additional experimental evidence justifying WISER. In §A, we start with discussing how Assumption 2.1 can be implemented even in presence of human-edits. Then, in §B, we provide some practical examples of watermarking schemes satisfying Assumption 2.2.§C complements the short experimental section in §4 by providing extensive numerical studies concerning the empirical behavior of WISER. Finally, in §D, we provide the detailed mathematical arguments behind WISER.

## A  DEALING WITH MIXED-SOURCE TEXTS

The assumption of knowledge of $\zeta_t$ can be too restrictive in most realistic scenarios where human edits are possible. In such cases, one assumes that the pseudo-random numbers $\zeta_t$ can also be reconstructed based on the available text and a Key with the help of a *Hash function* $\mathcal{A}$:

$$\zeta_t = \mathcal{A}(\omega_{(t-m):(t-1)}, \text{Key}). \tag{A.1}$$

Suppose $\tilde{\omega}_1 \ldots \tilde{\omega}_n$ be a mixed-source text, with segments of un-interrupted watermarked texts punctuated by human-generated texts through substitution, insertion or deletion of LLM generated texts. As a reference, we refer the readers to Procedure 1 of human edits in Li et al. (2024a). Note that it is impossible for any verifier to retrieve the exact pseudo-random numbers corresponding to each token in a mixed-source texts. Nevertheless, with the knowledge of the hash function and Key, one can construct $\tilde{\zeta}_t = \mathcal{A}(\omega_{(t-m):(t-1)}, \text{Key})$. Once there is a stretch of un-interrupted watermarked interval with length at least $m \geq 1$, the pseudo-random numbers $\zeta_{t+m}, \zeta_{t+m+1}, \ldots$ can be reliably re-constructed through (A.1) as the corresponding $\tilde{\zeta}$'s. On the other hand, if $\zeta_t$ is not the correct pseudo-random variable associated with $\omega_t$, then either

1. $\tilde{\omega}_t$ is human generated, in which case Working Hypothesis 2.2 of Li et al. (2025a) applies to yield $\omega_t$ and $\tilde{\zeta}_t$ are independent conditional on $P_t$;

2. $\tilde{\omega}_t$ is watermarked, which must mean if $\tilde{\omega}_t = \omega_{\tilde{t}}$ in the original watermarked text, then $\omega_{\tilde{t}} = S(P_{\tilde{t}}, \zeta_{\tilde{t}})$ for some true, unknown, pseudo-random number $\zeta_{\tilde{t}}$. In this case we invoke the sensitive nature of the hash function to conclude that $\omega_t$ and $\tilde{\zeta}_t$ are independent.

This argument appears in more detail in Section A.1 of Li et al. (2024a). In conclusion, the verifier can always obtain access to a sequence $\zeta_{m:n}$ corresponding to a given text $\omega_{1:n}$ such that (i) if $\omega_{(t-m):t}$ is NOT watermarked then $\omega_t$ and $\zeta_t$ are independent conditional on $P_t$; (ii), otherwise, $\zeta_t$ and $\omega_t$ may be intricately dependent on each other. This latter observation is crucial to our subsequent analysis and proposals, for it allows us to construct valid pivotal statistics. In light of the above discussion, we can be excused in making the Assumption 2.1.

Assumption 2.1 can be seen through the lens of constructing the $\tilde{\zeta}_t$ with $m = 1$. We make this slightly simplistic assumption to avoid the un-necessary measure theoretical niceties which might potentially cloud the novelty of our approach. Even with this assumption, proposing a computationally efficient algorithm and establishing its theoretical validity in a setting with multiple watermarked intervals, is an arguably non-trivial task in itself, and to the best of our knowledge, our paper is the first one to deal with this problem with full mathematical rigor. Finally, for a general mixed-source text, we remark that the WISER algorithm can be trivially extended to the setting with general $m \in \mathbb{N}$ by padding an interval of length $m$ to the left of the watermarked segments located by WISER. Its theoretical analysis will none-the-less remain non-trivial, and requires specific attention.

## B  EXAMPLES TO ASSUMPTION 2.2

In this section, we justify the elevated alternative hypothesis Assumption 2.2 by illustrating its occurrence through two popular watermarking schemes.

*Example* (Gumbel Watermark, Aaronson (2023)). Let $\zeta = (U_w)_{w \in \mathcal{W}}$ consist of $|\mathcal{W}|$ i.i.d. copies of $U(0, 1)$. The Gumbel watermark is implemented as:

$$S^{\text{gum}}(\zeta, P) := \arg \max_{w \in \mathcal{W}} \frac{\log U_w}{P_w}, \tag{B.1}$$

The pivot statistic is taken as $Y_t = U_{t,\omega_t}$, $t \in [n]$. From Lemma 2, when $\Delta = 1/2$, $\inf_{P \in \mathcal{P}_\Delta} \mathbb{E}_{1,P}[h(Y)] \geq \sum_{n=1}^{\infty}(\frac{1}{n} - \frac{1}{n+2})$, which, in light of $h(Y) \sim \mathrm{Exp}(1)$ entails that $d \geq 1/2$.

*Example* (Inverse Transform Watermark, Kuditipudi et al. (2024)). Consider an NTP distribution $P$ and a permutation $\pi : \mathcal{W} \mapsto S_{|\mathcal{W}|}$, where $S_{|\mathcal{W}|}$ is the group of permutations of $\{1, 2, \ldots, |\mathcal{W}|\}$. Further consider the multinomial distribution $\{P_{\pi^{-1}(w)}\}_{w=1}^{|\mathcal{W}|}$. The CDF of this distribution takes the form

$$F(x; \pi) = \sum_{w' \in \mathcal{W}} P_{w'} \cdot \mathbf{1}_{\{\pi(w') \leq x\}}.$$

Taking as input $U \sim U(0,1)$, the generalized inverse of this CDF is defined as

$$F^{-1}(U; \pi) = \min\Big\{ i : \sum_{w' \in \mathcal{W}} P_{w'} \cdot \mathbf{1}_{\{\pi(w') \leq i\}} \geq U \Big\},$$

which, under the $H_0$ of no watermark, follows the multinomial distribution $P$ after applying the permutation $\pi$. The inverse transform watermark is defined as the decoder:

$$\mathcal{S}^{\mathrm{inv}}(P, \zeta) := \pi^{-1}\big(F^{-1}(U; \pi)\big).$$

Lemma 4.1 of Li et al. (2025a) indicates that under alternate, the distribution of $S^{\mathrm{inv}}$ is intricately inter-related with the NTP $P$. To make the verification of Assumption 2.2 tractable, we impose a few assumption. Assume $|\mathcal{W}| \to \infty$, and with $P_{t,(i)}$ denoting the $i$-th largest co-ordinate of the probability vector $P_{t,(i)}$ for every token $t$ and $i \in [|\mathcal{W}|]$, we also assume

$$\lim_{|\mathcal{W}| \to \infty} P_{t,(1)} = 1 - \Delta \quad \text{and} \quad \lim_{|\mathcal{W}| \to \infty} \log |\mathcal{W}| \cdot P_{t,(2)} = 0.$$

Consider the pivot statistic

$$Y_t = \big|U_t - \eta(\pi_t(w_t))\big|, \quad \eta(i) := \frac{i-1}{|\mathcal{W}| - 1}.$$

Under Theorem 4.1 of Li et al. (2025a), $\mathbb{E}_1[1 - Y] = \frac{2+\Delta}{3}$, and $\mathbb{E}_0[1 - Y] = \frac{2}{3}$. Therefore, here $d = \frac{\Delta}{3}$.

## C  EXTENDED NUMERICAL EXPERIMENTS

In this section we provide additional numerical experiments complementing those in §4. In §C.1, we compare the accuracy of WISER with other competitive methods in the literature, on various benchmark datasets on myriad standard watermarking schemes. Moving on to §C.2, we investigate the effect of watermark intensity as well as the watermarked length on the performance of the algorithms . Finally, in §C.3, we provide some ablation studies corresponding to the hyper-parameters in WISER.

### C.1  COMPARATIVE PERFORMANCE OF WISER

From §4.1, recall the experimental set-up, the SOTA benchmark algorithms as well as the considered watermarking schemes. For each of the experiments, we implement WISER with block size equal to $b = 65$, $\rho = 0.5$, $\alpha = 0.05$ and $\gamma = 0.1$. Before we provide detailed comparison studies, we elaborate on the performance metrics used.

### C.1.1  PERFORMANCE METRIC

To ensure consistency with the prior works, we primarily treat the intersection-over-union (IOU) as a performance measure. Let, $\boldsymbol{I} := (I_1, \ldots, I_K)$ denotes the true watermarked intervals and $\widehat{\boldsymbol{I}} := (\widehat{I}_1, \ldots, \widehat{I}_{\hat{K}})$ be the estimated watermarked segments. Then, the intersection-over-union metric is given by

$$\mathrm{IOU}(\boldsymbol{I}, \widehat{\boldsymbol{I}}) = \frac{|(\cup_{i=1}^K I_i) \cap (\cup_{j=1}^{\hat{K}} \hat{I}_j)|}{|(\cup_{i=1}^K I_i) \cup (\cup_{j=1}^{\hat{K}} \hat{I}_j)|}.$$

Owing to Theorem 3.1, it is obvious that the IOU measure is expected to be close to 1 for the `WISER` method. Following the definition of Pan et al. (2025), we also compute the precision, recall and F1-score based on whether any of the estimated intervals have a nonempty intersection with any of the true intervals, i.e.,

$$\text{Precision} = \frac{|\{i : 1 \leq i \leq \hat{K}, \hat{I}_i \cap (\cup_{j=1}^{K} I_j) \neq \phi\}|}{\hat{K}}, \ \text{Recall} = \frac{|\{i : 1 \leq i \leq \hat{K}, \hat{I}_i \cap (\cup_{j=1}^{K} I_j) \neq \phi\}|}{K}.$$

**Rand Index and asymmetry of the watermark segmentation.**

In addition to these metrics, the Rand Index (RI) is also usually used to measure coherence between the estimated and true watermarked segments, using the algorithm illustrated in Prates (2021). For the standard definition of Rand Index, see Equation (2) of Prates (2021). However, as briefly discussed in Section 2.2.3, the Rand Index may depict a wrong picture of the performance of a watermarked segment identification algorithm. Before proceeding, we briefly deliberate on these issues.

As an illustration, consider the situation where most of the tokens (say 90%) are watermarked, while the watermark detection algorithm fails to detect any watermarked segment. While the performance of such an algorithm should reflect poorly, the standard Rand Index fails to capture this due to the exchangeability of the watermarked segment and the non-watermarked segment: any pair of indices $(i, j)$ that is truly watermarked trivially is also part of the estimated non-watermarked segment and considered as a concordant pair.

To circumvent these limitations described there, we consider a Modified Rand Index (MRI) given as

$$\text{MRI}(\boldsymbol{I}, \widehat{\boldsymbol{I}}) := \text{RI}(\boldsymbol{I}, \widehat{\boldsymbol{I}})$$
$$- \frac{\sum_{i \neq j} \left( \sum_{k=1}^{K} \mathbf{1}\{\{i,j\} \subseteq I_k \cap (\cup_{l=1}^{\hat{K}} \hat{I}_l)^c\} + \sum_{l=1}^{\hat{K}} \mathbf{1}\{\{i,j\} \subseteq \hat{I}_l \cap (\cup_{k=1}^{K} I_k)^c\} \right)}{\binom{n}{2}},$$

where $\mathbf{1}\{\cdot\}$ is the indicator function, and $n$ is the number of tokens. The MRI simply adjusts the RI by restricting its exchangeability only within each of the watermarked or non-watermarked intervals, but not in between. Intuitively, the MRI removes the specific pairs of indices $(i, j)$ from the calculation of RI for which both the indices $i$ and $j$ lie either in a true watermarked interval but are estimated to be in the non-watermarked region, or are estimated to be in a watermarked interval but actually lie in a non-watermarked region.

### C.1.2 EXPERIMENTAL RESULTS AND EXPLANATION

The comparison results are summarized in Tables 1, 2-4, corresponding to each of the watermarking schemes considered. Across all watermarking settings, `WISER` consistently delivers the strongest performance across every model and metric. In the Gumbel case, it achieves near-perfect results with IOU scores above 0.90, precision of 1.0, and recall above 0.98 across both small and large models. Competing methods like `Aligator` and `SeedBS-NOT` often fail to balance recall and precision, either collapsing to very low precision (`Aligator`) or producing weaker recall (`SeedBS-NOT`), while `Waterseeker` attains moderate balance but still lags well behind `WISER`.

The trend is even more pronounced in the cases of Inverse and Red-Green setups, where the pivot statistics remain uniformly bounded. In these cases, `Aligator` fail to detect any watermarked intervals, while both `SeedBS-NOT` and `Waterseeker` suffer a significant decline in performance. In contrast, `WISER` maintains F1-scores in the range of 0.95 - 0.99 with stable IOU values across model sizes, showing robustness to different architectures and vocabulary sizes. `Waterseeker` provides the next best alternative, but with noticeable drops in IOU and F1, especially for larger models. These findings clearly demonstrate that `WISER` not only generalises across watermarking schemes but also offers substantial gains in both detection accuracy and reliability, marking a clear benefit over existing baselines.

| Model Name | Vocab Size | Method | IOU | Precision | Recall | F1 | RI | MRI |
|---|---|---|---|---|---|---|---|---|
| facebook/opt-125m | 50272 | WISER | **0.906** | **0.995** | **0.980** | **0.987** | **0.968** | **0.931** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.558 | 0.988 | 0.710 | 0.826 | 0.804 | 0.783 |
| | | SeedBS-NOT | 0.178 | 0.282 | 0.228 | 0.252 | 0.529 | 0.428 |
| google/gemma-3-270m | 262144 | WISER | **0.874** | **0.965** | **0.958** | **0.961** | **0.945** | **0.934** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.547 | 0.983 | 0.695 | 0.814 | 0.797 | 0.775 |
| | | SeedBS-NOT | 0.221 | 0.316 | 0.272 | 0.293 | 0.575 | 0.544 |
| facebook/opt-1.3b | 50272 | WISER | **0.846** | **0.980** | **0.928** | **0.953** | **0.934** | **0.904** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.555 | 0.985 | 0.698 | 0.817 | 0.802 | 0.781 |
| | | SeedBS-NOT | 0.189 | 0.322 | 0.250 | 0.282 | 0.526 | 0.437 |
| princeton-nlp/Sheared-LLaMA-1.3B | 32000 | WISER | **0.656** | **0.990** | **0.962** | **0.976** | **0.871** | **0.850** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.582 | 0.992 | 0.750 | 0.854 | 0.826 | 0.807 |
| | | SeedBS-NOT | 0.181 | 0.286 | 0.235 | 0.258 | 0.541 | 0.515 |
| mistralai/Mistral-7B-v0.1 | 32000 | WISER | **0.718** | **0.935** | **0.822** | **0.875** | **0.859** | **0.847** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.590 | 0.996 | 0.760 | 0.862 | 0.830 | 0.813 |
| | | SeedBS-NOT | 0.154 | 0.265 | 0.192 | 0.223 | 0.502 | 0.489 |
| meta-llama/Meta-Llama-3-8B | 128256 | WISER | **0.878** | **0.995** | **0.965** | **0.980** | **0.955** | **0.913** |
| | | Aligator | 0.000 | 0.005 | 0.002 | 0.003 | 0.205 | 0.144 |
| | | Waterseeker | 0.510 | 0.980 | 0.652 | 0.783 | 0.774 | 0.748 |
| | | SeedBS-NOT | 0.143 | 0.242 | 0.185 | 0.210 | 0.511 | 0.480 |

Table 2: Results for Inverse Watermarking

| Model Name | Vocab Size | Method | IOU | Precision | Recall | F1 | RI | MRI |
|---|---|---|---|---|---|---|---|---|
| facebook/opt-125m | 50272 | WISER | **0.853** | **1.000** | **0.975** | **0.987** | **0.914** | **0.903** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.259 | 0.209 |
| | | Waterseeker | 0.730 | 0.998 | 0.815 | 0.897 | 0.889 | 0.882 |
| | | SeedBS-NOT | 0.570 | 0.665 | 0.615 | 0.639 | 0.897 | 0.870 |
| google/gemma-3-270m | 262144 | WISER | **0.838** | **0.973** | **0.970** | **0.972** | **0.908** | **0.896** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.643 | 0.982 | 0.820 | 0.894 | 0.864 | 0.850 |
| | | SeedBS-NOT | 0.600 | 0.749 | 0.738 | 0.743 | 0.900 | 0.872 |
| facebook/opt-1.3b | 50272 | WISER | **0.846** | **0.993** | **0.990** | **0.992** | **0.923** | **0.913** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.623 | 0.990 | 0.815 | 0.894 | 0.851 | 0.836 |
| | | SeedBS-NOT | 0.597 | 0.764 | 0.735 | 0.749 | 0.901 | 0.874 |
| princeton-nlp/Sheared-LLaMA-1.3B | 32000 | WISER | **0.850** | **1.000** | **0.990** | **0.995** | **0.919** | **0.908** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.619 | 0.995 | 0.810 | 0.893 | 0.851 | 0.836 |
| | | SeedBS-NOT | 0.570 | 0.775 | 0.738 | 0.756 | 0.898 | 0.860 |
| mistralai/Mistral-7B-v0.1 | 32000 | WISER | **0.814** | **0.995** | **0.955** | **0.975** | **0.909** | **0.898** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.559 | 0.993 | 0.742 | 0.850 | 0.818 | 0.799 |
| | | SeedBS-NOT | 0.507 | 0.718 | 0.672 | 0.695 | 0.877 | 0.843 |
| meta-llama/Meta-Llama-3-8B | 128256 | WISER | **0.864** | **1.000** | **0.995** | **0.997** | **0.929** | **0.919** |
| | | Aligator | 0.000 | 0.000 | 0.000 | 0.000 | 0.203 | 0.141 |
| | | Waterseeker | 0.647 | 1.000 | 0.838 | 0.912 | 0.866 | 0.851 |
| | | SeedBS-NOT | 0.590 | 0.778 | 0.770 | 0.774 | 0.919 | 0.883 |

Table 3: Results for Red-Green Watermarking

| Model Name | Vocab Size | Method | IOU | Precision | Recall | F1 | RI | MRI |
|---|---|---|---|---|---|---|---|---|
| facebook/opt-125m | 50272 | WISER | **0.925** | **0.998** | **0.998** | **0.998** | **0.980** | **0.979** |
| | | Aligator | 0.665 | 0.345 | 0.978 | 0.510 | 0.935 | 0.927 |
| | | Waterseeker | 0.712 | 1.000 | 0.905 | 0.950 | 0.891 | 0.884 |
| | | SeedBS-NOT | 0.469 | 0.725 | 0.560 | 0.632 | 0.867 | 0.817 |
| google/gemma-3-270m | 262144 | WISER | **0.935** | **1.000** | **1.000** | **1.000** | **0.982** | **0.973** |
| | | Aligator | 0.558 | 0.252 | 0.952 | 0.399 | 0.906 | 0.889 |
| | | Waterseeker | 0.614 | 1.000 | 0.782 | 0.878 | 0.841 | 0.824 |
| | | SeedBS-NOT | 0.334 | 0.610 | 0.440 | 0.511 | 0.766 | 0.686 |
| facebook/opt-1.3b | 50272 | WISER | **0.904** | **1.000** | **0.990** | **0.995** | **0.972** | **0.969** |
| | | Aligator | 0.446 | 0.216 | 0.928 | 0.350 | 0.887 | 0.863 |
| | | Waterseeker | 0.677 | 1.000 | 0.840 | 0.913 | 0.873 | 0.861 |
| | | SeedBS-NOT | 0.350 | 0.573 | 0.430 | 0.491 | 0.753 | 0.717 |
| princeton-nlp/Sheared-LLaMA-1.3B | 32000 | WISER | **0.919** | **1.000** | **1.000** | **1.000** | **0.979** | **0.977** |
| | | Aligator | 0.397 | 0.202 | 0.870 | 0.328 | 0.870 | 0.842 |
| | | Waterseeker | 0.653 | 1.000 | 0.778 | 0.875 | 0.851 | 0.837 |
| | | SeedBS-NOT | 0.264 | 0.486 | 0.350 | 0.407 | 0.688 | 0.666 |
| mistralai/Mistral-7B-v0.1 | 32000 | WISER | **0.896** | **1.000** | **0.998** | **0.999** | **0.973** | **0.972** |
| | | Aligator | 0.215 | 0.164 | 0.672 | 0.263 | 0.817 | 0.774 |
| | | Waterseeker | 0.646 | 1.000 | 0.795 | 0.886 | 0.853 | 0.838 |
| | | SeedBS-NOT | 0.238 | 0.468 | 0.315 | 0.376 | 0.650 | 0.575 |
| meta-llama/Meta-Llama-3-8B | 128256 | WISER | **0.908** | **1.000** | **0.998** | **0.999** | **0.976** | **0.976** |
| | | Aligator | 0.551 | 0.351 | 0.950 | 0.513 | 0.911 | 0.891 |
| | | Waterseeker | 0.535 | 1.000 | 0.712 | 0.832 | 0.799 | 0.773 |
| | | SeedBS-NOT | 0.413 | 0.658 | 0.545 | 0.596 | 0.775 | 0.730 |

Table 4: Results for Permute and Flip Watermarking

### C.1.3 WHY WISER OUTPERFORMS OTHER METHODS

The enhanced performance of WISER does not come out-of-the-blue, rather we argue that it is a byproduct of our unique, epidemic change-point perspective that marries theoretical validity with practical insights. While these methods—SeedBS-NOT, Aligator, and Waterseeker— each contribute useful perspectives, they also exhibit important limitations that the generality of our method usually overcomes.

**Limitations of SeedBS-NOT:** The limitations of SeedBS-NOT primarily arise from its reliance on a permutation-based change-point detection framework, which is inherently computationally expensive. Moreover, nowhere they restrict their attention to the specific scenario of watermarked segments, which consigns the change-points to occur in pairs, corresponding to the start and end of a watermarked segment. This is automatically alleviated by WISER through its adoption of a natural epidemic change-point formulation. This structural assumption substantially reduces the search space, yielding both computational efficiency and improved statistical stability. Additionally, SeedBS-NOT works with the sequence of p-values that are computed from a single observation of the pivot statistic at that location. Due to the complicated nature of the dependence between these p-values, they are difficult to combine to increase the statistical power. Our approach circumvents this by aggregating the pivot statistic at the block level (Step 7 in Figure 2), enhancing the effective sample size and increasing the power of the detection.

**Limitations of Aligator:** The Aligator algorithm frames the task as a reinforcement learning problem, producing a smoothed estimate of the underlying generative process and subsequently applying token-level hypothesis tests with a p-value threshold. While this strategy can capture localized deviations, it often results in a large number of short and fragmented detections, many of whom might be spurious due to possible multiple testing. Consequently, the method tends to produce many disjoint intervals, which severely diminishes its precision. By contrast, the discarding stage of WISER enforces structural coherence at the segment level, before returning fine-grained estimate through applying (3.1). This ensures that localized intervals correspond more closely to contiguous watermark insertions.

**Limitations of Waterseeker:** The Waterseeker algorithm may seem structurally similar to the proposed WISER method, in that it also employs a two-stage detection framework. However, Waterseeker considers a sliding window-based testing mechanism in its first stage, which has a crucial limitation. Consider a very realistic scenario when one of the pivot statistics corresponding

to an un-watermarked token is high simply due to random chance. In `Waterseeker`, this will push the score up for $W$ consecutive windows, usually resulting in a false positive in the first stage. On the other hand, for `WISER` such anomalous pivot statistics will affect only one block, which, being usually part of a connected interval with small length, can potentially be discarded with a very high probability in our *discarding stage*. For larger model, this scenario is extremely likely, making this reduction in precision much more pronounced (see Models google/gemma-3-270m and meta-Ilama/Meta-Llama-3-8B in Tables 1, 2 - 4). Moreover, Pan et al. (2025) provide only limited theoretical validation of their approach, making the optimal tuning of hyper-parameters difficult to justify. This lack of statistical guarantees limits its reliability across watermarking schemes and model sizes, in contrast to the rigorous and general guarantees underlying `WISER`.

## C.2 EFFECT OF WATERMARK INTENSITY

| Type | Method | IOU | F1 | RI | MRI |
|---|---|---|---|---|---|
| Strong but short | WISER | 0.794 | 0.984 | 0.933 | 0.925 |
| | SeedBS-NOT | 0.639 | 0.785 | 0.919 | 0.900 |
| | Waterseeker | 0.878 | 0.997 | 0.969 | 0.967 |
| Weak but long | WISER | 0.745 | 0.779 | 0.628 | 0.551 |
| | SeedBS-NOT | 0.172 | 0.321 | 0.675 | 0.187 |
| | Waterseeker | 0.268 | 0.847 | 0.519 | 0.172 |

Table 5: Effect on watermarking signal strength

Following the experimental design of Pan et al. (2025), we evaluate the comparative performance of the proposed `WISER` algorithm under varying levels of watermark intensity. As a demonstration, we choose Google's Gemma-3 series model (270 million) to generate a completion of 500 tokens for each input prompt. The watermark strength is modulated through the bias parameter $\delta$ of the Red-Green watermarking scheme (Kirchenbauer et al., 2023), while another parameter $m$ specifies the length of the watermarked region by applying the decoding strategy to the middlemost $m$ tokens within the 500-token output.

In the "strong but short" configuration ($\delta = 2.0, m = 100$), as shown in Table 5, all methods perform well, achieving a Rand Index exceeding 0.9. Although `WISER` is not the best-performing method in this particular case, it remains competitive with `Waterseeker`, which achieves the highest score. By contrast, in the "weak but long" configuration ($\delta = 1.0, m = 400$), only `WISER` maintains robust performance. While `SeedBS-NOT` appears to achieve a higher Rand Index, this outcome is primarily attributed to the issues described in §2.2.3 and §C.1.1. The Modified Rand Index (MRI) offers a more reliable assessment, highlighting the superiority of `WISER` in this setting.

## C.3 ABLATION STUDIES

We also perform an ablation study to understand the effectiveness of the hyper-parameters (e.g. -block size and $\rho$) of `WISER`. Our results are arguably quite stable across wide choices of the tuning parameters; nevertheless we provide more informed choices along with additional insights.

For this study, we consider a single watermarked segment from token index 100 to 200, fix $\rho = 0.25$ and vary the tuning parameter $b$ of the `WISER` algorithm. As one would have hoped, increasing the block size too much decreases the performance, as the smaller watermarked segments gets subsumed in the noise of unwatermarked segments when block sizes are too large. On the other hand, decreasing the block size would reduce the statistical power of the detection algorithm in the first stage itself. Therefore, one requires a judicious choice of the block size to optimally balance these two aspects, which is empirically observed through the upper plot of Figure 4. Based on empirical evidence, we recommend the choice $b \in (\lceil \sqrt{n} \rceil, 3\lceil \sqrt{n} \rceil)$, which works quite well in various settings that we have experimented with, while being also theoretically supported.

A similar conclusion also holds for the choice of $\rho$, for which we fix the block size as $b = 25$ and vary the tuning parameter $\rho$. As the choice of $d$ in Assumption 2.2 is exogeneously determined based on the language model and watermarking scheme, a large value of $\rho$ would imply a smaller $\tilde{d}$ and
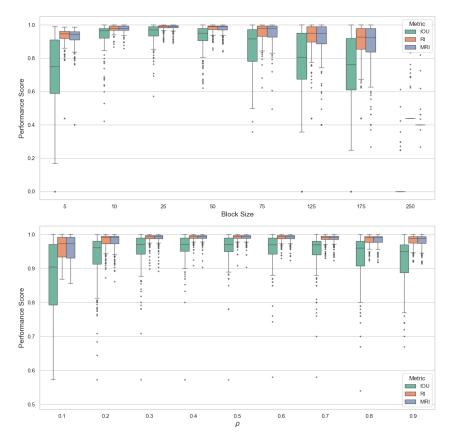
Figure 4: Effect on performance metrics (IOU and Rand Index) due to modification of the hyper-parameters of the `WISER` algorithm, namely block size (Top) and $\rho$ (Bottom).

by virtue of Theorem 3.1 would imply a larger error. The lower plot of Figure 4 demonstrates this empirically. However, any value of $\rho$ between 0.1 and 0.5 provides reasonable and relatively stable estimates.

# D  PROOF OF THEORETICAL RESULTS

In this section, we collect the proofs of theoretical results in the §3. Before we proceed further, we establish some notations. In the following, we write $a_n \lesssim b_n$ if $a_n \leq Cb_n$ for some constant $C > 0$, and $a_n \asymp b_n$ if $C_1 b_n \leq a_n \leq C_2 b_n$ for some constants $C_1, C_2 > 0$. Often we denote $a_n \lesssim b_n$ by $a_n = O(b_n)$. Additionally, if $a_n/b_n \to 0$, we write $a_n = o(b_n)$. For a function $f : \mathbb{R}^n \otimes \mathbb{R}^m \to \mathbb{R}$, let $f^{(1)}(\theta, w) = \frac{\partial}{\partial \theta} f(\theta, w)$, $\theta \in \mathbb{R}^n$, $w \in \mathbb{R}^m$, $n, m \geq 1$, be the partial derivative function with respect to $\theta$.

## D.1  PROOF OF THEOREM 3.1

In the following, we first state and prove a more generalized version of Theorem 3.1.

**Theorem D.1.** *Let $\{X_t\}_{t=1}^n := \{h(Y_t)\}_{t=1}^n$ be the pivot statistics based on the given input text, and assume that $I_0 \subset \{1, \ldots, n\}$ be the watermarked interval. Grant Assumption 2.2. Let us also denote*

$$\varepsilon_t = \begin{cases} X_t - \mu_0, t \notin I_0, \\ X_t - \mu_t, \ \mu_t := \mathbb{E}_{1,P_t}[X_t], t \in I_0. \end{cases}$$

*Suppose the class of distributions $\mathcal{P}$ is closed and compact, and there exists $\eta > 0$ such that $\sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(\eta|\varepsilon|)] < \infty$. Moreover, assume that $\min\{\mathrm{Var}_0(\varepsilon), \sup_P \mathrm{Var}_{1,P}(\varepsilon)\} > 0$. Then*

*it holds that*

$$|\hat{I} \Delta I_0| = O_{\mathbb{P}}\big((\sup_{\theta \geq 0}\{\theta \rho \tilde{d} - \Psi(\theta)\})^{-1}\big),$$

*where $\Delta$ denotes the symmetric difference operator, $O_{\mathbb{P}}$ hides constants independent of $n$ and $\tilde{d}$, and*

$$\Psi(\theta) = \log \mathbb{E}_0[\exp(\theta \varepsilon)] + 2^{-1} \log \sup_P \mathbb{E}_{1,P}[\exp(2\theta \varepsilon)] + 2^{-1} \log \sup_P \mathbb{E}_{1,P}[\exp(-2\theta \varepsilon)].$$

Theorem D.1 is proved by showing that the probability $\mathbb{P}(|\hat{I} \Delta I_0| > M)$ is small for all sufficiently large $M$. This probability is controlled by considering the objective function $V_I = S_{I^c} - (\mu_0 + \rho \tilde{d})|I^c|$, where $S_I = \sum_{k \in I} X_k$ and $S_{I^c} = \sum_{k=1}^n X_k - S_I$, and noting that, by construction of $\hat{I}$, $\mathbb{P}(|\hat{I} \Delta I_0| > M) \leq \mathbb{P}(\inf_{I:|I \Delta I_0| > M} V_I - V_{I_0} \leq 0)$. Usually, in change-point literature, one controls terms such as $\inf_{I:|I \Delta I_0| > M} V_I - V_{I_0}$ through Hàjek-Rényi type inequality Hájek & Rényi (1955); see Bai (1994); Bonnerjee et al. (2025). Such inequalities are usually derived by dividing the domain, on which infimum is taken, into smaller intervals, and applying Doob's inequality or Rosenthal's inequality piece-meal. However, the main bottleneck in a this particular setting is the potentially strong dependence between the pivot statistics in watermarked patches. We develop novel arguments that exploits $\sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(\eta|\varepsilon|)] < \infty$ to provide an extended version of the Hajek-Renyi theory through the lens of cumulant generating function. The proof is provided below.

*Proof of Theorem D.1.* For a candidate watermarked interval $I$, let $A_1(I) = I \cap I_0^c$, $A_2(I) = I \cap I_0$, $A_3(I) = I^c \cap I_0$, $A_4(I) = (I \cup I_0)^c$, and correspondingly $x_i(I) = |A_i(I)|$, $i = 1(1)4$. Subsequently, we omit the argument $I$ when it is clear from the context. Note that $|I_0| = x_2 + x_3$, $|I| = x_1 + x_2$, and $|\hat{I} \Delta I_0| = x_1 + x_3$. Note that, by definition of $\hat{I}$ it follows that $V_{\hat{I}} \leq V_{I_0}$. Finally, denote $S_i = \sum_{k \in A_i} X_k$, and $S_i^\varepsilon = \sum_{k \in A_i} \varepsilon_k$. With these notations established, we proceed through the following series of implications.

$$
\begin{aligned}
V_I - V_{I_0} &= (S_{I^c} - S_{I_0^c}) - (|I^c| - |I_0^c|)(\mu_0 + \rho \tilde{d}) \\
&= S_3 - S_1 - (x_3 - x_1)(\mu_0 + \rho \tilde{d}) \\
&= (S_3^\varepsilon + \sum_{t \in A_3} \mu_t) - (S_1^\varepsilon + x_1 \mu_0) - (x_3 - x_1)(\mu_0 + \rho \tilde{d}) \\
&= S_3^\varepsilon - S_1^\varepsilon + \sum_{t \in A_3}(\mu_t - \mu_0) + (x_1 - x_3)\rho \tilde{d} \\
&\geq S_3^\varepsilon - S_1^\varepsilon + x_3(d - \rho \tilde{d}) + x_1 \rho \tilde{d} \tag{D.1} \\
&\geq S_3^\varepsilon - S_1^\varepsilon + (x_1 + x_3)\rho \tilde{d}, \tag{D.2}
\end{aligned}
$$

where (D.1) follows from Assumption 2.2 and (D.2) uses $d \geq 2\rho \tilde{d}$. For some $M > 0$, let $D_M := \{I : |I \Delta I_0| > M\}$. Let $I_0 = [L, R]$. Note that, a candidate interval $I$ can belong to any of the following five sub-classes:

- $\mathcal{P}_1 := \{I : I \subseteq I_0, I \in D_M\}$.

- $\mathcal{P}_2 := \{I : I \supseteq I_0, I \in D_M\}$.

- $\mathcal{P}_3 := \{I : I \cap I_0 = \phi, I \in D_M\}$.

- $\mathcal{P}_4 := \{(a, b) : a < L < b < R, I = (a, b) \in D_M\}$.

- $\mathcal{P}_5 := \{(a, b) : L < a < R < b, I = (a, b) \in D_M\}$.

22

Subsequently, we detail the analysis for the relatively harder case $\mathcal{P}_4$. The arguments for the other cases are similar. Observe that:

$$\mathbb{P}(\hat{I} \in \mathcal{P}_4) \leq \mathbb{P}(\min_{I:I\in\mathcal{P}_4} V_I - V_{I_0} \leq 0)$$

$$\leq \mathbb{P}\big(\max_{I:I\in\mathcal{P}_4,x_1+x_3>M} \frac{S_1^\varepsilon - S_3^\varepsilon}{x_1+x_3} \geq \rho\tilde{d}\big)$$

$$\leq \sum_{j=M+1}^\infty \inf_{\theta\geq 0} \mathbb{P}\big(\max_{a,b:a<L<b<R:L-a+R-b=j} \exp(\theta(S_{[a,L]}^\varepsilon - S_{[b,R]}^\varepsilon)) \geq \exp(\theta\rho\tilde{d}j)\big)$$

$$\leq \sum_{j=M+1}^\infty \inf_{\theta\geq 0} \exp(-\theta\rho\tilde{d}j)\mathbb{E}\big[\max_{a,b:a<L<b<R:L-a+R-b=j} \exp(\theta(S_{[a,L]}^\varepsilon - S_{[b,R]}^\varepsilon))\big]$$

$$\leq \sum_{j=M+1}^\infty \inf_{\theta\geq 0} \exp(-\theta\rho\tilde{d}j)\mathbb{E}\big[\max_{a,b:a\in\{L-j+1,\cdots,L\},b\in\{(R-j+1)\vee L,\cdots,R\}} \exp(\theta(S_{[a,L]}^\varepsilon - S_{[b,R]}^\varepsilon))\big]$$

$$\text{(D.3)}$$

For $j \in [n]$, let $\mathcal{F}_j := \sigma(\{(\omega_{s-1},\zeta_s) : s < j\})$. Write

$$\mathbb{E}\big[\max_{a,b:a\in\{L-j+1,\cdots,L\},b\in\{(R-j+1)\vee L,\cdots,R\}} \exp(\theta(S_{[a,L]}^\varepsilon - S_{[b,R]}^\varepsilon))\big]$$

$$= \mathbb{E}\big[\max_{a:a\in\{L-j+1,\cdots,L\}} \exp(\theta S_{[a,L]}^\varepsilon)\mathbb{E}\big[\max_{b\in\{(R-j+1)\vee L,\cdots,R\}} \exp(-\theta S_{[b,R]}^\varepsilon) \mid \mathcal{F}_{(R-j)\vee L}\big]\big]$$

$$\leq \mathbb{E}\big[\max_{a:a\in\{L-j+1,\cdots,L\}} \exp(\theta S_{[a,L]}^\varepsilon)\sqrt{\mathbb{E}[\exp(-2\theta S_{[(R-j+1)\vee L,R]}^\varepsilon) \mid \mathcal{F}_{(R-j)\vee L}]}$$

$$\sqrt{\mathbb{E}\big[\max_{b\in\{(R-j+1)\vee L,\cdots,R\}} \exp(2\theta S_{[(R-j+1)\vee L,b]}^\varepsilon) \mid \mathcal{F}_{(R-j)\vee L}\big]}\big], \qquad \text{(D.4)}$$

where, (D.4) follows from Cauchy-Schwartz inequality. Now, note that, by construction of $\varepsilon_t$, conditional on $\mathcal{F}_{(R-j)\vee L}$, $\varepsilon_t$ is a martingale difference sequence adapted to $\sigma(\{(\omega_{s-1},\zeta_s) : (R-j+1)\vee L \leq s \leq t\})$. Since $x \mapsto \exp(2\theta x)$ is convex, hence $\exp(2\theta S_{[(R-j+1)\vee L,b]}^\varepsilon), b \in \{(R-j+1)\vee L,\ldots,R\}$ is a sub-martingale sequence. Consequently, Doob's maximal inequality (Hall & Heyde, 1980) applies. Further sequential conditioning yields the following series of inequalities.

$$\mathbb{E}\big[\max_{b\in\{(R-j+1)\vee L,\cdots,R\}} \exp(2\theta S_{[(R-j+1)\vee L,b]}^\varepsilon) \mid \mathcal{F}_{(R-j)\vee L}\big]$$

$$\leq 4\mathbb{E}[\exp(2\theta S_{[(R-j+1)\vee L,R]}^\varepsilon) \mid \mathcal{F}_{(R-j)\vee L}]$$

$$\leq 4\mathbb{E}[\exp(2\theta S_{[(R-j+1)\vee L,R-1]}^\varepsilon)\mathbb{E}[\exp(2\theta\varepsilon_R)|\mathcal{F}_{R-1}] \mid \mathcal{F}_{(R-j)\vee L}]$$

$$\leq 4\sup_P \mathbb{E}_{1,P}[\exp(2\theta\varepsilon)]\mathbb{E}[\exp(2\theta S_{[(R-j+1)\vee L,R-1]}^\varepsilon) \mid \mathcal{F}_{(R-j)\vee L}]$$

$$\leq 4\big(\sup_P \mathbb{E}_{1,P}[\exp(2\theta\varepsilon)]\big)^j. \qquad \text{(D.5)}$$

Proceeding along similar lines, we obtain

$$\mathbb{E}[\exp(-2\theta S_{[(R-j+1)\vee L,R]}^\varepsilon) \mid \mathcal{F}_{(R-j)\vee L}] \leq 4\big(\sup_P \mathbb{E}_{1,P}[\exp(-2\theta\varepsilon)]\big)^j, \qquad \text{(D.6)}$$

and

$$\mathbb{E}\big[\max_{a:a\in\{L-j+1,\cdots,L\}} \exp(\theta S_{[a,L]}^\varepsilon)\big] \leq 4\big(\mathbb{E}_0[\exp(\theta\varepsilon)]\big)^j. \qquad \text{(D.7)}$$

Combining (D.5)-(D.7) and plugging them in (D.4) and (D.3), one obtains

$$\mathbb{P}(I \in \mathcal{P}_4) \leq 16 \sum_{j=M+1}^\infty \inf_{\theta\geq 0} \bigg(\exp(-\theta\rho\tilde{d})\mathbb{E}_0[\exp(\theta\varepsilon)]\sqrt{\sup_P \mathbb{E}_{1,P}[\exp(2\theta\varepsilon)]\sup_P \mathbb{E}_{1,P}[\exp(-2\theta\varepsilon)]}\bigg)^j.$$

$$\text{(D.8)}$$

To deliver the coup de grâce of our argument, we are required to bound (D.8). To that end, define $\phi : \mathbb{R}_+ \to \mathbb{R}$ as

$$\phi(\theta) = -\theta\rho\tilde{d} + \log\mathbb{E}_0[\exp(\theta\varepsilon)] + 2^{-1}\log\sup_P \mathbb{E}_{1,P}[\exp(2\theta\varepsilon)] + 2^{-1}\log\sup_P \mathbb{E}_{1,P}[\exp(-2\theta\varepsilon)].$$

By definition of $\phi$, $\mathbb{P}(I \in \mathcal{P}_4) \le \sum_{j=M+1}^{\infty} \inf_{\theta \ge 0} \exp(j\phi(\theta))$. Moreover, for $\lambda \in (0,1)$, $\theta_1, \theta_2 \in \mathbb{R}_+$, Hölder's inequality produces

$$\log \sup_{P} \mathbb{E}_{1,P}[\exp(2(\lambda\theta_1 + (1-\lambda)\theta_2)\varepsilon)] \le \sup_{P} \Big( \lambda \log \mathbb{E}_{1,P}[\exp(2\theta_1\varepsilon)] + (1-\lambda)\log \mathbb{E}_{1,P}[\exp(2\theta_2\varepsilon)] \Big). \tag{D.9}$$

Similar arguments for $\log \mathbb{E}_0[\exp(\theta\varepsilon)]$ and $\log\sup_P \mathbb{E}_{1,P}[\exp(-2\theta\varepsilon)]$ show that $\phi$, being a linear combination of convex functions with non-negative weights (note that $-\theta\rho\tilde{d}$ is linear) , is itself convex.

Let $f : \mathbb{R} \otimes \mathbb{R}^{|W|} \mapsto \mathbb{R}$ be given by $f(\theta, P) = \log \mathbb{E}_0[\exp(2\theta\varepsilon)] + \log \mathbb{E}_{1,P}[\exp(2\theta\varepsilon)]$. Recalling that $f^{(1)}(\theta, w) = \frac{\partial}{\partial\theta} f(\theta, w)$, observe that

$$f^{(1)}(0, P) = 0 \text{ for any } P \in \mathcal{P}, \tag{D.10}$$

since $\mathbb{E}_0[\varepsilon] = \mathbb{E}_{1,P}[\varepsilon] = 0$. Therefore, noting that $\mathcal{P}$ is a compact subset of the $|W|$-dimensional simplex, in light of $\sup_{P \in \mathcal{P}} \mathbb{E}[\exp(-\eta|\varepsilon|)] \le \sup_{P \in \mathcal{P}} \mathbb{E}[\exp(\eta|\varepsilon|)] < \infty$, Danskin's Theorem (Danskin, 1967) entails

$$\frac{\partial}{\partial\theta} \sup_{P \in \mathcal{P}} f(\theta, P)\Big|_{\theta\downarrow 0} = \sup_{P \in \mathcal{P}} f^{(1)}(0, P) = 0,$$

where in the second inequality we use that $f(0, P) = 0$ for any $P \in \mathcal{P}$, and the third equality follows from (D.10). Similarly, $\frac{\partial}{\partial\theta} \sup_{P \in \mathcal{P}} f(\theta, P)\Big|_{\theta\uparrow 0} = -\inf_{P \in \mathcal{P}} f^{(1)}(0, P) = 0$. Therefore, $\phi'(0) = -\rho\tilde{d} < 0$. On the other hand, since $\min\{\text{Var}_0(\varepsilon), \sup_P \text{Var}_{1,P}(\varepsilon)\} > 0$, hence $\phi_\theta \to \infty$ as $\theta \uparrow \infty$. In conjunction with $\phi$ being convex, there must exist $\kappa \in (0,1)$ such that $\log \kappa := \inf_{\theta \ge 0} \phi(\theta)$. Consequently, from (D.8), one obtains,

$$\mathbb{P}(I \in \mathcal{P}_4) \le 16 \sum_{j=M+1}^{\infty} \kappa^j = O(\kappa^M).$$

Suppose $\delta \in (0,1)$ be given. A choice of $M > \frac{\log 1/\varepsilon}{\log 1/\kappa}$ ensures that $\mathbb{P}(I \in \mathcal{P}_4) < \delta$. This completes the proof. $\square$

Finally, Theorem 3.1 is proved by invoking Theorem D.1 and Proposition 3.

We can further sharpen the $O(\tilde{d}^{-1})$ rate in Theorem 3.1 to $O(\tilde{d}^{-2})$ by assuming a mild condition: local sub-Gaussianity of the pivot statistics. The following proposition also follows from Theorem D.1 and Proposition 3.

**Proposition 1.** *Grant the assumptions of Theorem 3.1. If*

$$\max\{\mathbb{E}_0[\exp(r|\varepsilon|)], \sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(r|\varepsilon|)]\} \le \exp(r^2/2),$$

*for all $r \in [0, \eta]$, then choosing $\rho > 0$ such that $\rho\tilde{d} < \frac{5}{2}\eta$, then $|\hat{I} \Delta I_0| = O_{\mathbb{P}}(\tilde{d}^{-2})$.*

## D.2 PROOF OF THEOREM 3.2

For convenience, we first re-state the theorem.

**Theorem D.2.** *Assume that the null distribution of the pivot statistics is absolutely continuous with respect to the Lebesgue measure. Fix $\alpha \in (0,1)$, and recall the quantities defined in* `WISER` *described in Figure 2. Let the block length $b = b_n$ satisfy $b_n \asymp \sqrt{n}$, and suppose the threshold $\mathcal{Q} = \mathcal{Q}_n$ is selected so that*

$$\mathbb{P}_0\big( \max_{1 \le k \le \lceil n/b \rceil} S_k > \mathcal{Q} \big) = \alpha.$$

*Also assume that that $\mathbb{E}_0[|X - \mu_0|^{3+\delta}] < \infty$ for some $\delta > 0$, $d \ge c$ for some constant $c > 0$,*

$$\sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[X] < \infty, \tag{D.11}$$

24

*and there exists $\tau > 0$ such that*

$$\kappa := \inf_{\theta \geq 0} \theta(\mu_0 + \tau d) + \log \sup_P \mathbb{E}_{1,p}[\exp(-\theta X)] < 0. \tag{D.12}$$

*Additionally, let the number of watermarked intervals $K$ be bounded, and Assumption 3.1 is granted for the watermarked intervals $I_k, k \in [K]$. Then, given $\varepsilon > 0$, under the assumptions of Theorem 3.1, there exists $M_\varepsilon \in \mathbb{R}_+$, independent of $n, K$, and $d$, such that,*

$$\liminf_{n \to \infty} \mathbb{P}\big(\hat{K} = K, \max_{k \in [K]} |\hat{I}_k \Delta I_k| < M_\varepsilon d^{-1}\big) \geq 1 - \varepsilon. \tag{D.13}$$

Let $\widetilde{\mathcal{B}} = \{1 \leq k \leq \lceil n/b \rceil : B_k \subseteq I_j \text{ for some } j \in [K]\}$. Our proof proceeds through a series of arguments, each carefully orchestrated to establish the validity of the corresponding steps of our algorithm. We comment that subsequently, all statements involving $n$ but without a limit attached to it are meant to be considered for all sufficiently large values of $n$.

### Step 1: Validity of first stage thresholding.

In this step, we show that

$$\mathbb{P}(\min_{k \in \widetilde{\mathcal{B}}} S_k > \mathcal{Q}) \to 1, \text{ as } n \to \infty. \tag{D.14}$$

To begin with, note that

$$\limsup_{n \to \infty} \max_{k \in \widetilde{\mathcal{B}}} \mathbb{P}(S_k \leq \mathcal{Q}_n)^{1/b} \leq \limsup_{n \to \infty} \inf_{\theta \geq 0} \exp(\theta \mathcal{Q}_n b_n^{-1} + \log \sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(-\theta X)])$$

$$\leq \inf_{\theta \geq 0} \exp(\theta \mu_0 + \log \sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(-\theta X)]) \tag{D.15}$$

$$\leq \exp(\kappa) < 1, \tag{D.16}$$

where (D.15) is obtained through an application of Proposition 4, and (D.16) follows from (D.12). Since $\kappa < 0$, one has $\frac{n}{b} \exp(\kappa b) \to 0$ as $n \to \infty$, and consequently

$$\mathbb{P}(\min_{k \in \widetilde{\mathcal{B}}} S_k \leq \mathcal{Q}) \leq \frac{n}{b} \max_{k \in \widetilde{\mathcal{B}}} \mathbb{P}(S_k \leq \mathcal{Q}_n) \to 0, \text{ as } n \to \infty,$$

thereby establishing (D.14).

### Step 2. Estimation of the number of watermarked regions through the set $M$.

Recall $M$ from the Step 24 of WISER. In this step of our proof, we will prove $\mathbb{P}(\hat{K} = K) \to 1$, which will also imply that $|M|$ is even with probability approaching 1. Therefore, we may be excused for assuming that $|M|$ is even.

Let $C_1, \ldots, C_{\hat{K}}$ be the disjoint set of intervals in $M$, with $C_j = [(s_{2j-1} - 1)b + 1, s_{2j}b]$. Note that for each $k \in \mathcal{B}$ such that $S_k > \mathcal{Q}$, $B_k \subseteq C_j$ for some $j$. Let $\widetilde{\mathcal{B}}_j = \{k \in \widetilde{\mathcal{B}} : B_k \subseteq I_j\}, j \in [K]$. Clearly, $\widetilde{\mathcal{B}} = \cup_{j=1}^K \widetilde{\mathcal{B}}_j$, and $\widetilde{\mathcal{B}}_j$ are disjoint. Therefore, in light of the construction of $M$ from blocks surpassing the threshold $\mathcal{Q}$, it follows,

$$\mathbb{P}(\min_{k \in \widetilde{\mathcal{B}}} S_k > \mathcal{Q}) = \mathbb{P}(\min_{j \in [K]} \min_{k \in \widetilde{\mathcal{B}}_j} S_k > \mathcal{Q}) \leq \mathbb{P}(\text{for each } j \in [K], \text{ there exists } i_j \in [\hat{K}] \text{ such that } \widetilde{\mathcal{B}}_j \subseteq C_{i_j}),$$

which implies, in light of (D.14),

$$\mathbb{P}(A_n) \to 1, \text{ as } n \to \infty, \text{ where}, A_n := \{\text{for each } j \in [K], \text{ there exists } i_j \in [\hat{K}] \text{ such that } \widetilde{\mathcal{B}}_j \subseteq C_{i_j}\}. \tag{D.17}$$

It is crucial to note that since both $\widetilde{\mathcal{B}}_j$ and $C_j$'s are defined to occur from left-to-right and since $C_j$'s are connected intervals, under the event $A_n$ it also holds that $i_1 \leq i_2 \leq \ldots \leq i_K$. At this stage, the relationship between $\hat{K}$ and $K$ is still not entirely clear. Subsequently, we will show that under the event $A_n$, the mapping $j \mapsto i_j$ is injective, establishing that $\hat{K} \geq K$ with high probability. To that end, suppose there exists $k_1 < k_2 \in [K]$ such that $i_{k_1} = i_{k_2}$. Since $C_{i_{k_1}}$ is a connected interval,

25

$i_{k_1} = i_{k_2}$ implies that that $i_{k_1} = i_{k_1+1}$. Let $\mathbb{P}_{E,F}(\cdot) = \mathbb{P}(\cdot \cap E \cap F)$ for any events $E$, $F$. Consider the following series of inequalities.

$$\mathbb{P}_{A_n}(\text{There exists } k \in [K-1] \text{ such that } C_{i_k} = C_{i_{k+1}})$$
$$\leq \mathbb{P}_{A_n}(\text{There exists } k \in [K-1] \text{ such that } (I_{k,R}, I_{k+1,L}) \subseteq C_{i_k})$$
$$\leq \mathbb{P}_{A_n}(\text{There exists } k \text{ such that } \min_{l \in (\lceil I_{k,R}/b \rceil, \lfloor I_{k+1,L}/b \rfloor)} S_l > \mathcal{Q})$$
$$\leq \mathbb{P}_0(\sum_{k=1}^{n/b} I\{S_k > \mathcal{Q}\} \geq C_0\sqrt{\log n}), \tag{D.18}$$

where the $\mathbb{P}_0$ in final inequality appears since for $l \in (\lceil I_{k,R}/b \rceil, \lfloor I_{k+1,L}/b \rfloor)$, the region $B_l$ is unwatermarked; the $\sqrt{\log n}$ appears by invoking Assumption 3.1 and noting that $b^{-1}(I_{k+1,L} - I_{k,R}) \geq C_0\sqrt{\log n}$. An application of Proposition 5 to (D.18) entails, in view of (D.17), that,

$$\mathbb{P}_{A_n}(\bar{B}_n) \to 1, \text{ as } n \to \infty, \text{ where } \bar{B}_n = \{C_{i_k} \text{ and } C_{i_s} \text{ are disjoint if } i_k \neq i_s\}.$$

Clearly, this implies that $\mathbb{P}_{A_n}(\hat{K} \geq K) \to 1$ as $n \to \infty$, which also produces $\mathbb{P}(\hat{K} \geq K) \to 1$ as $n \to \infty$. On the other hand, if $\hat{K} > K$, then under the event $A_n \cap \bar{B}_n$, there exists $j \in [\hat{K}]$ such that $C_j$ and $\cup_{s \in \mathcal{B}} B_s$ are disjoint. Consequently, it must be true that $|C_j \cap (\cup_{k=1}^K I_j)| \leq 2\sqrt{n}$. Note that by construction of $C_j$'s in Steps 13-22 of WISER $|C_j| \geq c\sqrt{n \log n}$. Therefore it must be true that there are at least $2^{-1}c\sqrt{\log n}$ many $s$'s such that $bs \notin C_j \cap (\cup_{k=1}^K I_j)$, and $S_s > \mathcal{Q}$. Hence it follows from Proposition 5 that

$$\mathbb{P}_{A_n, \bar{B}_n}(\hat{K} > K) \to 0, \text{ as } n \to \infty,$$

which immediately implies that

$$\mathbb{P}(\hat{K} = K) \to 1 \text{ as } n \to \infty. \tag{D.19}$$

**Step 3. Choice of $\tilde{d}$ and $\rho$.**

Recall $\tilde{d}$ from Step 28 of WISER in Figure 2. In this step, we establish that there exists $\rho > 0$, such that $d > 2\rho\tilde{d}$ with high probability. In conjunction to $\tilde{d}$, also define

$$d^\dagger = \frac{\sum_{j=1}^K \sum_{s \in I_j}(X_s - \mu_0)}{\sum_{j=1}^K |I_j|}.$$

Let the event $\{\hat{K} = K\}$ be denoted as $E_n$. Under $E_n$, by construction of $D_j$, $\mathbb{P}_{A_n, B_n, E_n}(I_j \subseteq D_j \text{ for all } j \in [K]) \to 1$ as $n \to \infty$. Call the latter event as $F_n$. Observe that under $E_n \cap F_n$, it holds

$$\sum_{j=1}^K |I_j| + 2Cn^{1/2}\log n \geq \sum_{j=1}^{\hat{K}} |D_j| \geq \sum_{j=1}^K |I_j| + Cn^{1/2}\log n \tag{D.20}$$

for some $C > 0$. Therefore, under the same event, it follows

$$\tilde{d} \leq d^\dagger \frac{\sum_{j=1}^K |I_j|}{\sum_{j=1}^K |I_j| + Cn^{1/2}\log n} + \frac{\sum_{s \in \cup_j(I_j^c \cap D_j)}(X_s - \mu_0)}{\sum_{j=1}^K |I_j| + Cn^{1/2}\log n}. \tag{D.21}$$

We first tackle the second term in the upper-bound in (D.21). Let $D_j^\dagger = [(I_{j,L} - \lfloor Cn^{1/2}\log^{3/2} n \rfloor) \vee 1, (I_{j,R} + \lfloor Cn^{1/2}\log^{3/2} n \rfloor) \wedge n]$. Again, by construction of $D_j$ as well as from Assumption 3.1, for all sufficiently large $n$ it follows

$$\mathbb{P}_{A_n, \bar{B}_n, E_n}(D_j \subseteq D_j^\dagger, D_i^\dagger \cap D_j^\dagger = \phi \text{ for } i \neq j) \to 1.$$

Call the above event as $G_n$. Fix $\varepsilon > 0$, and consider the following implications.

$$\mathbb{P}_{A_n, \bar{B}_n, E_n} \left( \frac{\sum_{s \in \cup_j (I_j^c \cap D_j)} (X_s - \mu_0)}{\sum_{j=1}^K |I_j| + Cn^{1/2} \log n} > \varepsilon \right)$$

$$\leq \mathbb{P}_{A_n, \bar{B}_n, E_n, G_n} \left( \frac{\sum_{s \in \cup_j (I_j^c \cap D_j^\dagger)} |X_s - \mu_0|}{\sum_{j=1}^K |I_j| + Cn^{1/2} \log n} > \varepsilon \right) + o(1)$$

$$\leq \mathbb{P} \left( \frac{\sum_{s \in \cup_j (I_j^c \cap D_j^\dagger)} |X_s - \mu_0|}{\sum_{j=1}^K |I_j| + Cn^{1/2} \log n} > \varepsilon \right) + o(1)$$

$$\leq \frac{O(n^{1/2} \log^{3/2} n)}{\varepsilon^2 n \log^2 n} + o(1) = o(1), \tag{D.22}$$

where the inequality in the final assertion follows from $| \cup_{j=1}^K (I_j^c \cap D_j^\dagger) | \lesssim n^{1/2} \log^{3/2} n$. Therefore, (D.21) and (D.22) jointly yields

$$\mathbb{P}_{A_n, \bar{B}_n, E_n, F_n} (\tilde{d} \leq 2d^\dagger) \to 1, \text{ as } n \to \infty. \tag{D.23}$$

Next, we focus on controlling $d^\dagger$ by $d$. To that end, we resort to an argument through moment generating functions. On one hand, (D.12) entails

$$\mathbb{P}(d^\dagger \leq \tau d) \leq \inf_{\theta \geq 0} \left( \exp(\theta(\mu_0 + \tau d) + \log \sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(-\theta X)]) \right)^{\sum_{j=1}^K |I_j|} \leq \exp(\kappa \sum_{j=1}^K |I_j|) \to 0. \tag{D.24}$$

On the other hand, in light of (D.11) and $d \geq c$, choose

$$\nu > \frac{\sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[X] - \mu_0}{c} \vee \frac{\tau}{4},$$

and write:

$$\mathbb{P}(d^\dagger \geq 2\nu d) \leq \inf_{\theta \geq 0} \left( \exp(-2\theta \nu c + \log \sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(\theta(X - \mu_0))]) \right)^{\sum_{j=1}^K |I_j|}. \tag{D.25}$$

Echoing the argument in the proof of Theorem 3.1, define

$$g(\theta, P \,;\, c) = -2\theta \nu c + \log \mathbb{E}_{1,P}[\exp(\theta(X - \mu_0))], \quad \widetilde{g}(\theta \,;\, c) = \sup_{P \in \mathcal{P}} g(\theta, P).$$

Since $\mathcal{P}$ is compact and $\sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(\eta|X - \mu_0|] < \infty$, Danskin's Theorem (Danskin, 1967) applies and produces

$$\widetilde{g}_+^{(1)}(0, P \,;\, c) = \frac{\partial}{\partial \theta} \sup_{P \in \mathcal{P}} g(\theta, P \,;\, c) \Big|_{\theta \downarrow 0} = \sup_{P \in \mathcal{P}} g^{(1)}(0, P \,;\, c) = -2\nu d + \sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[X - \mu_0] \leq -\nu d < 0, \tag{D.26}$$

where the final inequality is derived via (D.11). Moreover, similar to (D.9) it can be argued that $\widetilde{g}(\theta)$ is convex in $\theta$. Finally, since $\widetilde{g}(0 \,;\, c) = 0$, (D.26) coupled with its convexity implies that $\varphi(c) := \inf_{\theta \geq 0} \widetilde{g}(\theta \,;\, c) < 0$. In view of this, (D.25) results in

$$\mathbb{P}(d^\dagger \geq 2\nu d) \leq \exp(-\varphi(c) \sum_{j=1}^K |I_j|) \to 0 \text{ as } n \to \infty, \tag{D.27}$$

where the limiting assertion is due to $\sum_{j=1}^K |I_j| \geq c\sqrt{n}$. Finally, (D.23) and (D.27) jointly indicates that

$$\mathbb{P}_{A_n, B_n, E_n, F_n}(\tilde{d} \leq 4\nu d) \to 1, \text{ as } n \to \infty. \tag{D.28}$$

Subsequently, we choose $\rho = (8\nu)^{-1}$. In conclusion to this step, (D.22) along with (D.24) establishes

$$\mathbb{P}_{A_n, B_n, E_n, F_n}(G_n) \to 1 \text{ as } n \to \infty, \ G_n := \{\tau d \leq \tilde{d} \leq 4\nu d\}.$$

27

**Step 4. Localization of watermarked intervals.**

In this step, we establish the validity of our localized estimates $\hat{I}_j$. In Step 3, we argued that

$$\mathbb{P}_{A_n,B_n,E_n}(I_j \subseteq D_j \subseteq D_j^\dagger \text{ for each } j \in [K]) \to 1 \text{ as } n \to \infty.$$

Call the above event as $\tilde{F}_n$. Under $\tilde{F}_n$, it is immediate that

$$\hat{I}_j(\tilde{d}) = \underset{s\in L_j, t\in R_j}{\arg\min} \sum_{k\in D_j\backslash[s,t]} (X_k - \mu_0 - \rho\tilde{d}) = \underset{s\in L_j, t\in R_j}{\arg\min} \sum_{k\in D_j^\dagger\backslash[s,t]} (X_k - \mu_0 - \rho\tilde{d}),$$

since the operator $\sum_{k\in D_j^\dagger\backslash[s,t]}$ can be decomposed into $\sum_{k\in D_j\backslash[s,t]} + \sum_{k\in D_j^\dagger\backslash D_j}$.

We proceed towards applying Theorem 3.1 to $\hat{I}_j(\tilde{d})$. However, note that $\tilde{d}$ is a random quantity, so special care must be accorded to its treatment. To that end, define

$$\hat{I}_j(\sigma) = \underset{s\in L_j, t\in R_j}{\arg\min} \sum_{k\in D_j^\dagger\backslash[s,t]} (X_k - \mu_0 - \rho\sigma), \ \sigma \in [\tau d, 4\nu d].$$

Fix $j \in [K]$. For $M > 0$, let $D_M := \{I : |I\Delta I_j| > M\}$. For a candidate interval $I = [s,t]$, let

$$\widetilde{V}_I(\sigma) = \sum_{k\in D_j^\dagger\backslash[s,t]} (X_k - \mu_0 - \rho\sigma).$$

Clearly, by definition of $G_n$,

$$\mathbb{P}_{A_n,B_n,E_n,\tilde{F}_n,G_n}(|\hat{I}_j(\tilde{d}) \, \Delta \, I_n| > M)$$

$$\leq \mathbb{P}_{A_n,B_n,E_n,\tilde{F}_n,G_n}\left(\sup_{\sigma\in[\tau d, 4\nu d]} |\hat{I}_j(\sigma) \, \Delta \, I_n| > M\right)$$

$$\leq \mathbb{P}_{A_n,B_n,E_n,\tilde{F}_n,G_n}\left(\text{There exists } \sigma\in[\tau d, 4\nu d] \text{ such that } \inf_{s\in L_j, t\in R_j, I\in D_M} \widetilde{V}_I(\sigma) < \widetilde{V}_{I_j}(\sigma)\right)$$

$$\leq \mathbb{P}_{A_n,B_n,E_n,\tilde{F}_n,G_n}\left(\text{There exists } \sigma\in[\tau d, 4\nu d] \text{ such that } \inf_{I\in D_M} \widetilde{V}_I(\sigma) < \widetilde{V}_{I_j}(\sigma)\right)$$

$$\leq \mathbb{P}_{A_n,B_n,E_n,\tilde{F}_n,G_n}\left(\max_{I:x_1+x_3>M} \frac{S_1^\varepsilon - S_3^\varepsilon}{x_1+x_3} > (\frac{1}{2}\wedge\frac{\tau}{8\nu})d\right) \tag{D.29}$$

$$\leq \mathbb{P}_{A_n,B_n,E_n,\tilde{F}_n,G_n}\left(\max_{I:x_1+x_3>M} \frac{S_1^\varepsilon - S_3^\varepsilon}{x_1+x_3} > \frac{\tau}{8\nu}d\right) \tag{D.30}$$

$$\leq \mathbb{P}\left(\max_{I:x_1+x_3>M} \frac{S_1^\varepsilon - S_3^\varepsilon}{x_1+x_3} > \frac{\tau}{8\nu}d\right). \tag{D.31}$$

Here, (D.29) follows by recalling the notations in the proof of Theorem 3.1 and following the arguments (D.1)-(D.2) after observing $\sigma\in[\tau d, 4\nu d]$ implies $d - (8\nu)^{-1}\sigma \geq \frac{d}{2}$. Moreover, (D.30) also follows from $4\nu d \geq \sigma \geq \tau d$. Finally, (D.31) is derived from $\mathbb{P}(A \cap B) \leq \mathbb{P}(A)$; in particular, arguments of Theorem 3.1 can be followed verbatim to obtain that

$$\mathbb{P}\left(\max_{I:x_1+x_3>M} \frac{S_1^\varepsilon - S_3^\varepsilon}{x_1+x_3} > \frac{\tau}{8\nu}d\right) \leq \xi^M \text{ for some } \xi < 1.$$

Note that in the above assertion we have used the fact that $d \geq c$ to decouple $\xi$ from $d$. Given arbitrary $\varepsilon > 0$, $M_\varepsilon$ can be chosen to ensure $\xi^{M_\varepsilon} < \varepsilon$, and through $\kappa$, this choice of $M_\varepsilon$ solely depends on the constants $\nu$, $\tau$, $c$, and $\mu_0$, apart from the quantity $\varepsilon$. Therefore, in view of the number of watermarked intervals $K = O(1)$, we obtain that there exists $M_\varepsilon$ independent of $n$, $K$ and $d$ such that

$$\mathbb{P}_{A_n,B_n,E_n,\tilde{F}_n,G_n}(|\hat{I}_j(\tilde{d}) \, \Delta \, I_n| > M_\varepsilon \text{ for } j \in [K]) \leq \varepsilon$$

$$\implies \liminf_{n\to\infty} \mathbb{P}_{A_n,B_n,E_n,\tilde{F}_n,G_n}(|\hat{I}_j(\tilde{d}) \, \Delta \, I_n| \leq M_\varepsilon \text{ for } j \in [K]) \geq 1 - \varepsilon, \tag{D.32}$$

where in (D.32) we invoke

$$\lim_{n\to\infty} \mathbb{P}(A_n \cap B_n \cap E_n \cap \tilde{F}_n \cap G_n) = 1.$$

Recalling that $E_n = \{\hat{K} = K\}$ completes the proof.

### D.3 ADDITIONAL PROPOSITIONS

Firstly we provide a formal proof that the pivot statistics corresponding to un-watermarked tokens are i.i.d., a fundamental fact behind the construction and validity of our algorithm.

*Proof of Lemma 2.2.* Let $t \in S$. Since $\omega_t$ and $\zeta_t$ are independent conditional on $\omega_{1:(t-1)}$, hence $\mathcal{L}(Y_t)|\omega_{1:t} \stackrel{d}{=} \mathcal{L}(Y)$. Hence, $\{Y_t\}_{t \in S}$ are identically distributed since given $\omega_{1:t}$, the distribution of $Y_t$ is solely a function of the key $\zeta_t$ that are i.i.d.. Moreover, for $s < t \in S$, even if there is a watermarked region $I_k \subset (s, t)$, $\zeta_s$ and $\omega_l$ are independent for all $l \in (s, t]$. In view of the fact that conditional on $\omega_{1:s}$, $Y_s$ and $Y_t$ are completely determined by $\zeta_s$ and $(w_{s+1:t}, \zeta_{s+1:t})$ respectively, we deduce that $Y_s$ and $Y_t$ are independent conditional on $\omega_{1:s}$. Hence, for two Borel sets $A$ and $B$,

$$\mathbb{P}(Y_s \in A, Y_t \in B) = \mathbb{E}[\mathbb{P}(Y_s \in A|\omega_{1:s})\mathbb{E}[\mathbb{P}(Y_t \in B|\omega_{1:t})|\omega_{1:s}]]$$
$$= \mathbb{P}(Y_s \in A)\mathbb{P}(Y_t \in B) \text{ (from Definition 2.1)}.$$

This completes the proof. $\square$

Next, we collect the additional results that we have used in our theoretical arguments. The proofs are provided subsequently.

**Proposition 2.** *Let $h(x) = -\log(1-x)$, and suppose $\mathcal{P}_\Delta := \{\max_{w \in \mathcal{W}} P_w \leq 1 - \Delta\}$ for some fixed $\Delta > 0$. Then it follows that*

$$\inf_{P \in \mathcal{P}_\Delta} \mathbb{E}_{1,P}[h(Y)] \geq \sum_{n=1}^{\infty} \left( \frac{1}{n} - \left\lfloor \frac{1}{1-\Delta} \right\rfloor \frac{(1-\Delta)^2}{1+n(1-\Delta)} - \frac{1-(1-\Delta)\lfloor \frac{1}{1-\Delta} \rfloor}{1+n(1-(1-\Delta)\lfloor \frac{1}{1-\Delta} \rfloor)} \right). \quad \text{(D.33)}$$

**Proposition 3.** *Consider $\tilde{d}$ and $\Psi(\cdot)$ from Theorem 3.1. If there exists a constant $c > 0$ such that $\tilde{d} \geq c$, then*

$$(\sup_{\theta \geq 0} \{\theta \rho \tilde{d} - \Psi(\theta)\})^{-1} = O(\tilde{d}^{-1}). \quad \text{(D.34)}$$

*Recall $\varepsilon$ from Theorem 3.1. Suppose we additionally have that*

$$\max\{\mathbb{E}_0[\exp(r|\varepsilon|)], \sup_{P \in \mathcal{P}} \mathbb{E}_{1,P}[\exp(r|\varepsilon|)]\} \leq \exp(r^2/2) \text{ for all } r \in [0, \eta],$$

*$\eta$ being the same as in Theorem 3.1. Then, choosing $\rho > 0$ such that $\rho\tilde{d} < \frac{5}{2}\eta$, it holds that*

$$(\sup_{\theta \geq 0} \{\theta \rho \tilde{d} - \Psi(\theta)\})^{-1} = O(\tilde{d}^{-2}). \quad \text{(D.35)}$$

**Proposition 4.** *Let $\mathbb{E}_0[|X - \mu_0|^{3+\delta}] < \infty$ for some $\delta > 0$. Let $\mathcal{Q}$ be selected as in Theorem 3.2. Then it follows that $\mathcal{Q}/b \to \mu_0$ as $n \to \infty$.*

**Proposition 5.** *Let $X_i$ be i.i.d. with mean $\mu_0$, and let $B_k$ and $S_k$ be defined as in Steps 2 and 3 of* WISER *in Figure 2. Then it follows that*

$$\mathbb{P}_0 \left( \sum_{k=1}^{\lceil n/b \rceil} I\{S_k > \mathcal{Q}\} \geq C_0\sqrt{\log n} \right) \to 0, \text{ as } n \to \infty,$$

*where $\mathcal{Q}$ is defined as in Theorem 3.2.*

*Proof of Proposition 2.* From Lemma 3.1 of Li et al. (2025a), it follows

$$
\begin{aligned}
\mathbb{E}_{1,P}[h(X)] &= \sum_{w=1}^{|\mathcal{W}|} \int_0^1 x^{1/P_w-1}\big(-\log(1-x)\big)\,\mathrm{d}x \\
&= \sum_{w=1}^{|\mathcal{W}|} \sum_{n=1}^{\infty} \int_0^1 \frac{x^{1/P_w-1+n}}{n}\,\mathrm{d}x \\
&= \sum_{w=1}^{|\mathcal{W}|} \sum_{n=0}^{\infty} \frac{1}{n(n+1/P_w)} \\
&= \sum_{n=1}^{\infty} \Big(\frac{1}{n} - \sum_{w=1}^{|\mathcal{W}|} \frac{P_w}{n+1/P_w}\Big) \\
&\geq \sum_{n=1}^{\infty} \left( \frac{1}{n} - \lfloor \frac{1}{1-\Delta} \rfloor \frac{(1-\Delta)^2}{1+n(1-\Delta)} - \frac{1-(1-\Delta)\lfloor \frac{1}{1-\Delta} \rfloor}{1+n(1-(1-\Delta)\lfloor \frac{1}{1-\Delta} \rfloor)} \right),
\end{aligned}
\tag{D.36}
$$

where the final inequality follows from noting the convexity of $g : x \mapsto \sum_{i=1}^d \frac{x_i}{n+1/x_i}$, $\sum_{i=1}^d x_i = 1$, and noting that the optimum value of $g$ on the set $\mathcal{P}_\Delta$ occurs at the extrema defined by

$$
P_\Delta^\star = \Big( \underbrace{1-\Delta,\dots,1-\Delta}_{\lfloor \frac{1}{1-\Delta} \rfloor \text{ times}},\, 1-(1-\Delta)\cdot\lfloor \tfrac{1}{1-\Delta} \rfloor,\, 0,\dots \Big).
$$

$\square$

*Proof of Proposition 3.* Denote $\Lambda(x) := \sup_{\theta \geq 0}\{\theta\rho x - \Psi(\theta)\}$. Note that, an argument same as (D.10) shows that $\Psi'_+(0) = 0$, where $\Psi'_+(\cdot)$ denote the right derivative. Therefore, in light of $\tilde{d} \geq c$ for some constant $c > 0$, there exist $\theta_0 > 0$ such that $\frac{|\Psi(\theta)|}{\theta} \leq \frac{\rho\tilde{d}}{2}$ for all $\theta \in (0, \theta_0)$. Therefore,

$$
\Lambda(\tilde{d}) \geq 2^{-1}\theta_0\rho\tilde{d} - 4^{-1}\theta_0\rho c \geq 4^{-1}\theta_0\rho c,
$$

which immediately implies (D.34). Moving on, we work with the additional assumption that $\sup_{P\in\mathcal{P}} \mathbb{E}_{1,P}[\exp(r|\varepsilon|)] \leq \exp(r^2/2)$. This immediately implies that for all $\theta \in [0, \frac{\eta}{2}]$,

$$
\max\{\log \sup_P \mathbb{E}_{1,P}[\exp(2\theta\varepsilon)], \log \sup_P \mathbb{E}_{1,P}[\exp(-2\theta\varepsilon)]\} \leq 2\theta^2.
$$

Therefore, for all $\theta \in [0, \frac{\eta}{2}]$ it must hold that

$$
\Psi(\theta) \leq \frac{5}{2}\theta^2.
$$

Consequently, in light of $\rho\tilde{d} < \frac{5}{2}\eta$, one obtains,

$$
\Lambda(x) \geq \sup_{\theta\in[0,\frac{\eta}{2}]} \{\theta\rho x - \frac{5}{2}\theta^2\} = \frac{\rho^2 x^2}{10},
$$

which establishes (D.35). $\square$

*Proof of Proposition 4.* Our proof has two key steps: firstly, we will prove that if there is no watermarking in the entire sequence, then

$$
\max_{1\leq k\leq\lceil n/b \rceil} \frac{S_k}{b} \xrightarrow{\mathbb{P}} \mu_0.
\tag{D.37}
$$

Subsequently, we follow an argument similar to the proof of equation (29) in Li et al. (2025a), with crucial tweaks to accommodate the maximum over the block means. Let us first work towards (D.37). We note that a similar result (for the $p$-th moments) appears in Proposition E.2 in Deb et al. (2020)

but without proof. For the sake of completion, we provide an independent proof of (D.37) without invoking the aforementioned result. Fix $\varepsilon > 0$. Note that

$$\mathbb{P}_0\big(\max_{1 \leq k \leq \lceil n/b \rceil} b^{-1}(S_k - \mu_0) > \varepsilon\big) \leq \frac{n}{b}\mathbb{P}(b^{-1}(S_1 - \mu_0) > \varepsilon), \tag{D.38}$$

where for the last inequality we use that $S_k$'s are i.i.d. under $H_0$, i.e. no watermarking. Moving on, we apply the Fuk-Nagaev inequality (Corollary 4, Fuk & Nagaev (1971)),

$$\mathbb{P}_0(b^{-1}(S_1 - \mu_0) > \varepsilon) \leq c_{1,\delta}\frac{n}{(b\varepsilon)^{3+\delta}}\mathbb{E}_0[|X - \mu|^{3+\delta}] + \exp(-c_{2,\delta}\frac{b\varepsilon^2}{\sigma^2}), \ \sigma^2 := \mathbb{E}_0[X^2], \tag{D.39}$$

where $c_{1,\delta}, c_{2,\delta} > 0$ are constants depending solely on $\delta$. Note that $b^2 \asymp n$, and hence $\frac{n^{3/2}}{(b\varepsilon)^{3+\delta}} \to 0$ as $n \to \infty$. On the other hand, $\sqrt{n}\exp(-c_{2,\delta}\frac{b\varepsilon^2}{\sigma^2}) \to 0$ as $n \to \infty$. Therefore, from (D.38) and (D.39), one obtains (D.37).

Now suppose that $\limsup_{n\to\infty} \mathcal{Q}/b > \mu_0$. Then there exists $\gamma > 0$ and a strictly increasing sequence $\{n_k\} \subseteq \mathbb{N}$ such that $\mathcal{Q}_{n_k}/b_{n_k} > \mu_0 + \gamma$ for all sufficiently large $k \in \mathbb{N}$. Since (D.37) implies that

$$\max_{1 \leq l \leq \lceil n_k/b_{n_k} \rceil} \frac{S_{n_k}}{b_{n_k}} \xrightarrow{\mathbb{P}} \mu_0, \text{ as } k \to \infty,$$

therefore, there exists a strictly increasing sub-sequence $\{n_{k_r}\} \subseteq \{n_k\}$ such that

$$\max_{1 \leq l \leq \lceil n_{k_r}/b_{n_{k_r}} \rceil} \frac{S_{n_{k_r}}}{b_{n_{k_r}}} \xrightarrow{\text{a.s.}} \mu_0, \text{ as } r \to \infty, \text{ and } \mathcal{Q}_{n_{k_r}}/b_{n_{k_r}} > \mu_0 + \gamma \text{ for all sufficiently large } r.$$

Therefore, by the dominated convergence theorem,

$$\alpha = \lim_{r\to\infty} \mathbb{P}\left(\max_{1 \leq l \leq \lceil n_{k_r}/b_{n_{k_r}} \rceil} \frac{S_{n_{k_r}}}{b_{n_{k_r}}} > \frac{\mathcal{Q}_{n_{k_r}}}{b_{n_{k_r}}}\right) \leq \lim_{r\to\infty} \mathbb{P}\left(\max_{1 \leq l \leq \lceil n_{k_r}/b_{n_{k_r}} \rceil} \frac{S_{n_{k_r}}}{b_{n_{k_r}}} > \mu_0 + \gamma\right)$$
$$= \mathbb{P}(\mu_0 > \mu_0 + \gamma) = 0, \tag{D.40}$$

which is a contradiction. Hence, $\limsup_{n\to\infty} \mathcal{Q}/b \leq \mu_0$. Very similarly one can show $\liminf_{n\to\infty} \mathcal{Q}/b \geq \mu_0$, which completes the proof. $\qquad\square$

*Proof of Proposition 5.* Let $p_n = \mathbb{P}_0(S_k > \mathcal{Q})$. Clearly, by definition of $\mathcal{Q}$ it follows that

$$\alpha = \mathbb{P}_0(\max_k S_k > \mathcal{Q}) = 1 - (1 - p_n)^{\lceil n/b \rceil} \geq 1 - \exp(-c\sqrt{n}p_n),$$

which implies that $\sqrt{n}p_n = O(1)$. Note that $\sum_{k=1}^{\lceil n/b \rceil} I\{S_k > \mathcal{Q}\} \sim \text{Bin}(\lceil n/b \rceil, p_n)$. Therefore, using Chernoff bound, one obtains

$$\mathbb{P}_0\big(\sum_{k=1}^{\lceil n/b \rceil} I\{S_k > \mathcal{Q}\} \geq C_0\sqrt{\log n}\big) \leq \exp(-2^{-1}(1 + o(1))(\log\log n)\log n) \to 0,$$

which completes the proof. $\qquad\square$