# Information-Directed Random Walk for Rare Event Detection in Hierarchical Processes

Chao Wang, Kobi Cohen, Qing Zhao

*Abstract*— The problem of detecting a few anomalous processes among a large number of data streams is considered. At each time, aggregated observations can be taken from a chosen subset of the processes, where the chosen subset conforms to a given tree structure. The random observations are drawn from a general distribution that may depend on the size of the chosen subset and the number of anomalous processes in the subset. We propose a sequential search strategy by devising an information-directed random walk on the tree-structured observation hierarchy. The proposed policy is shown to be asymptotically optimal with respect to the detection accuracy and order-optimal with respect to the size of the search space. Effectively localizing the data processing to small subsets of the search space, the proposed strategy is also efficient in terms of computation and memory requirement.

*Index Terms*— Sequential design of experiments, active hypothesis testing, anomaly detection, noisy group testing, channel coding with feedback.

## I. INTRODUCTION

### A. Rare Event Detection in Hierarchical Processes

We consider the problem of detecting anomalies in a large number of processes. At each time, the decision maker chooses a subset of processes to observe. The chosen subset conforms to a predetermined tree structure. The (aggregated) observations are drawn from a general distribution that may depend on the size of the chosen subset and the number of anomalies in the subset. The objective is a sequential search strategy that adaptively determines which node on the tree to probe at each time and when to terminate the search in order to minimize a Bayes risk that takes into account both the sample complexity and the detection accuracy.

The above problem is an archetype for searching for a few rare events of interest among a massive number of data streams that can be observed at different levels of granularity. For example, financial transactions can be aggregated at different temporal and geographic scales. In computer vision applications such as bridge inspection by UAVs with limited battery capacity, sequentially determining areas to zoom in or zoom out can quickly locate anomalies by avoiding giving each pixel equal attention. A particularly relevant application is heavy hitter detection in Internet traffic monitoring. It is a common observation that Internet traffic flows are either "elephants" (heavy hitters) or "mice" (normal flows). A small percentage of high-volume flows account for most of the total traffic [1]. Quickly identifying heavy hitters is thus crucial to network stability and security, especially in detecting denial-of-service (DoS) attacks. Since maintaining a packet count for each individual flow is infeasible due to limited sampling resources at the routers, an effective approach is to aggregate flows based on the IP prefix of the source or destination addresses [2], [3]. Indeed, recent advances in software-defined networking (SDN) allow programmable routers to count aggregated flows that match a given IP prefix. The search space of all traffic flows thus follows a binary tree structure. Other applications include DoA estimation [4] and system control [5].

### B. Information-Directed Random Walk

To fully exploit the hierarchical structure of the search space, the key questions are how many samples to obtain at each level of the tree and when to zoom in or zoom out on the hierarchy. A question of particular interest is whether a sublinear scaling of the sample complexity with the size of the search space is feasible while achieving the optimal scaling with the detection accuracy. In other words, whether accurate detection can be achieved by examining only a diminishing fraction of the search space as the search space grows.

Our approach is to devise an information-directed random walk (IRW) on the hierarchy of the search space. The IRW initiates at the root of the tree and eventually arrives and terminates at the targets (i.e., the anomalous processes) with the required reliability. Each move of the random walk is guided by the test statistic of the sum log-likelihood ratio (SLLR) collected from each child of the node currently being visited by the random walk. This local test module ensures that the global random walk is biased toward a target and that the walk terminates at a target with the required detection accuracy.

Analyzing the sample complexity of the IRW strategy lies in examining the trajectory of the biased random walk. With a suitably chosen confidence level in the local test module, the random walk will concentrate, with high probability, on a smaller and smaller portion of the tree containing the targets and eventually probes the targets only. The basic structure of the analysis is to partition the tree into a sequence of half

trees with decreasing size, and bound the time the random walk spends in each half tree. The entire search process, or equivalently, each sample path of the biased random walk, is then partitioned into stages by the successively defined last passage time to each of the half trees in the shrinking sequence. We show that the sample complexity of the IRW strategy is asymptotically optimal in detection accuracy and order optimal, specifically, a logarithmic order, in the size of the search space when the aggregated observations are informative at all levels of the tree (with a finite number of exceptions).

We also consider the case when higher-level observations decay to pure noise. Using Bernoulli distribution as a case study, we show that when the Kullback-Leibler (KL) divergence between the target-absent and target-present distributions decays to zero in a polynomial order with the depth of the tree, the IRW offers a sample complexity that is a poly-logarithmic order in the number of processes; when the decaying rate is exponential in the level $l$ of the tree (i.e., $\propto \alpha^{2l}$), a sublinear scaling in the size of the search space can be achieved provided $1 < \alpha < \sqrt{2}$.

The proposed search strategy is deterministic with search actions explicitly specified at each given time. It involves little online computation beyond calculating the SLLR and performing simple comparisons. It is also efficient in terms of memory requirement. By effectively localizing data processing to small subsets of the search space, it has $O(1)$ computation and memory complexity.

### C. Related Work

The problem considered here falls into the general class of sequential design of experiments pioneered by Chernoff in 1959 [6] in which he posed a binary active hypothesis testing problem. Compared with the classical sequential hypothesis testing pioneered by Wald [7] where the observation model under each hypothesis is fixed, active hypothesis testing has a control aspect that allows the decision maker to choose different experiments (associated with different observation models) at each time. Chernoff proposed a *randomized* strategy and showed that it is asymptotically optimal as the error probability approaches zero. Known as the Chernoff test, this randomized strategy chooses, at each time, a probability distribution governing the selection of experiments based on all past actions and observations. The probability distribution is given as a solution to a maxmin problem, which can be difficult to solve, especially when the number of hypotheses and/or the number of experiments is large (a case of focus in this paper). Furthermore, the Chernoff test does not address the scaling with the number of hypotheses and results in a linear sample complexity in the size of the search space when applied to the problem considered here.

A number of variations and extensions of Chernoff's randomized test have been considered (see, for example, [8]–[11]). In particular, in [10], Naghshvar and Javidi developed a randomized test that achieves the optimal logarithmic order of the sample complexity in the number of hypotheses under certain implicit assumptions on the KL divergence between the observation distributions under different hypotheses. These assumptions, however, do not always hold in general for the observation models considered here. In particular, the assumption in [10] implies that the supremum of the log-likelihood ratio between any two hypotheses from a single observation is bounded, which does not hold in general for observation kernels with unbounded support. Furthermore, similar to the Chernoff test, this randomized test is specified only implicitly as solutions to a sequence of maxmin problems that is intractable for general observation distributions and large problem size. More specifically, the work in [10] formulates the sequential hypothesis testing problem as a partially observable Markov decision process (POMDP). After each measurement, the new observation is used to update the belief of each hypothesis. As a result, the computational complexity after each measurement has a linear order with the search space. In contrast, the algorithm proposed in this work has a constant computation and memory complexity that does not grow with the size of the search space.

In our prior work [11]–[13], we considered a target search problem in a linear search setting without access to hierarchical observations. The sample complexity of the search strategy developed there achieves asymptotical optimality with respect to the detection accuracy, but linear order with the size of the search space. The problem addressed in this work is fundamentally different, focusing on efficient exploitation of aggregated and potentially low-quality measurements to achieve an optimal sublinear order with the size of the search space. The problem of detecting anomalies or outlying sequences has also been studied under different formulations, assumptions, and objectives [14]–[17]. An excellent survey can be found in [18]. These studies, in general, do not address the optimal scaling in the detection accuracy or the size of the search space.

Tree-based search in data structures is a classical problem in computer science (see, for example, [19], [20]). It is mostly studied in a deterministic setting, i.e., the observations are deterministic when the target location is fixed. The problem studied in this work is a statistical inference problem, where the observations taken from the tree nodes follow general statistical distributions.

The problem studied here also has intrinsic connections with several problems studied in different application domains, in particular, adaptive sampling, noisy group testing, and channel coding with feedback. We discuss in detail the connections and differences in Section VII.

## II. PROBLEM FORMULATION

We first consider the case of a single target and a binary tree structure. Extensions to multi-target detection and general tree structures are discussed in Sections V and VI, respectively.

As illustrated in Fig. 1, $g_0$ and $f_0$ denote, respectively, the distributions of the anomalous process and the normal processes[1]. Let $g_l$ ($l = 1, \ldots, \log_2 M$) denote the distribution

---

[1]The proposed policy and the analysis extend with simple modifications to cases where each process has different target-present and target-absent distributions.
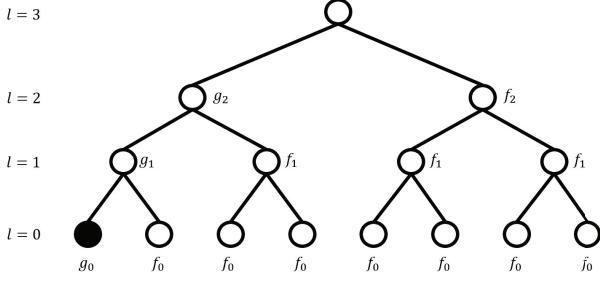
Fig. 1. A binary tree observation model with a single target.

of the measurements that aggregate the anomalous process and $2^l - 1$ normal processes, and $f_l$ ($l = 1, \ldots, \log_2 M$) denote the distribution of the measurements that aggregate $2^l$ normal processes (see Fig. 1). We allow general relation between $\{g_l, f_l\}$ and $\{g_0, f_0\}$, which often depends on the specific application. For example, in the case of heavy hitter detection where the measurements are packet counts of an aggregated flow, $g_l$ and $f_l$ are given by multi-fold convolutions of $f_0$ and $g_0$. For independent Poisson flows, $g_l$ and $f_l$ are also Poisson with mean values given by the sum of the mean values of their children at the leaf level. As is the case in most of the practical applications, we expect that observations from each individual process are more informative than aggregated observations. More precisely, we expect $D(g_{l-1} \| f_{l-1}) \sim D(g_l \| f_l)$ and $D(f_{l-1} \| g_{l-1}) \sim D(f_l \| g_l)$ for all $l > 0$, where $D(\cdot \| \cdot)$ denotes the KL divergence between two distributions. However, the results in this work hold for the general case without these monotonicity assumptions.

An active search strategy $\pi = \{\{\phi(t)\}_{t \geq 1}, \tau, \delta\}$ consists of a sequence of selection rules $\{\phi(t)\}_{t \geq 1}$ governing which node to probe at each time, a stopping rule $\tau$ deciding when to terminate the search, and a declaration rule $\delta$ deciding which leaf node is the target at the time of stopping. We adopt the Bayesian approach as in Chernoff's original work [6] and subsequent studies in [10], [11]. Specifically, a sampling cost $c \in (0, 1)$ is incurred for each observation and a loss of 1 for a wrong declaration. Let $\pi_m$ denote the *a priori* probability that process $m$ is anomalous, which is referred to as hypothesis $H_m$. The probability of detection error $P_e(\pi)$ and the sample complexity $Q(\pi)$ of strategy $\pi$ are given by

$$P_e(\pi) \triangleq \sum_{m=1}^{M} \pi_m \mathbb{P}[\delta \neq m \| H_m], \tag{1}$$

$$Q(\pi) \triangleq \sum_{m=1}^{M} \pi_m \mathbb{E}[\tau \| H_m], \tag{2}$$

where $\mathbb{P}$ and $\mathbb{E}$ denote the probability measure and expectation with respect to the probability space induced by $\pi$. The dependency on $\pi$ will be omitted from the notations when there is no ambiguity. The Bayes risk of $\pi$ is then given by

$$R(\pi) \triangleq P_e(\pi) + cQ(\pi). \tag{3}$$

The objective is a strategy $\pi$ that achieves the lower bound of the Bayes risk:

$$R^* = \inf_\pi R(\pi). \tag{4}$$

We are interested in strategies that offer the optimal scaling in both $c$, which controls the detection accuracy, and $M$, which is the size of the search space. A test $\pi$ is said to be *asymptotically optimal* in $c$ if, for fixed $M$,

$$\lim_{c \to 0} \frac{R(\pi)}{R^*} = 1. \tag{5}$$

A shorthand notation $R(\pi) \to R^*$ will be used to denote the relation specified in (5). If the above limit is a constant greater than 1, then $\pi$ is said to be order optimal. The asymptotic and order optimalities in $M$ are similarly defined as $M$ approaching infinity for a fixed $c$.

A dual formulation of the problem is to minimize the sample complexity subject to an error constraint $\varepsilon$, i.e.,

$$\pi^* = \arg\inf_\pi Q(\pi), \quad s.t. \quad P_e(\pi) \geq \varepsilon. \tag{6}$$

Note that, since we consider a non-empty set of real numbers that is bounded from below, arg inf always exists. In the Bayes risk given in (3), $c$ can be viewed as the inverse of the Lagrange multiplier, thus controls the detection accuracy of the test that achieves the minimum Bayes risk. Following the same lines of argument in [21], [22], one can obtain the solution to (6) from the solution under the Bayesian formulation.

## III. INFORMATION-DIRECTED RANDOM WALK

The IRW policy induces a biased random walk that initiates at the root of the tree and eventually arrives at the target with a sufficiently high probability. Each move of the random walk is guided by the output of a local test carried on the node currently being visited by the random walk. This local test module $\mathcal{U}(p)$ determines, with a confidence level $p$, whether this node contains the target, and if yes, which child contains the target. Based on the output, the random walk zooms out to the parent[2] of this node or zooms in to one of its child. The confidence level $p$ of the local test module is set to be greater than $1/2$ to ensure that the global random walk is more likely to move toward the target than move away from it.

Once the random walk reaches a leaf node, say node $m$ ($m = 1, \ldots, M$), samples are taken one by one from node $m$ and the SLLR $S_m(t)$ of node $m$ is updated with each new sample taken during the current visit to this node:

$$S_m(t) = \sum_{n=1}^{t} \log \frac{g_0(y(n))}{f_0(y(n))}. \tag{7}$$

When $S_m(t)$ drops below 0, the random walk moves back to the parent of node $m$. When $S_m(t)$ exceeds $\log \frac{\log_2 M}{c}$, the detection process terminates, and node $m$ is declared as the target. The choice of the stopping threshold $\log \frac{\log_2 M}{c}$ is to ensure that the error probability is in the order of $O(c)$, which in turn secures the asymptotic optimality in $c$ (see Theorem 1 and the proof in Appendix B).

We now specify the local test module $\mathcal{U}(p)$ carried out on upper-level nodes. The objective of $\mathcal{U}(p)$ is to distinguish three hypotheses—$H_0$ that this node does not contain the target and $H_1$ ($H_2$) that the left (right) child of this node contains

[2]The parent of the root node is defined as itself.

the target—with a confidence level no smaller than $p$ under each hypothesis. Various tests (fixed-sample-test, sequential, and active) with a guaranteed confidence level $p$ can be constructed for this ternary hypothesis testing problem. We present below a fixed-sample-size test. A sequential test based on the Sequential Probability Ratio Test (SPRT) [23] and an active test can also be constructed. Details are given in Appendix A. The theorems in this paper apply to IRW with fixed sample-size local test. IRW with the sequential, and active local tests are analyzed numerically, and also shown to outperform the fixed sample-size local test in simulation examples (see Section VIII-B).

Suppose that the random walk is currently at a node on level $l > 0$. A fixed $K_l$ samples, denoted as $y[n]$ ($n = 1, \ldots, K_l$), are taken from each child of the node. The SLLR of each child is computed:

$$\prod_{n=1}^{K_l} \log \frac{g_{l-1}(y[n])}{f_{l-1}(y[n])}. \tag{8}$$

If the SLLRs of both children are negative, the local test declares hypothesis $H_0$. Otherwise, the local test declares $H_1$ ($H_2$) if the left (right) child has a larger SLLR. The sample size $K_l$ is chosen to ensure a probability $p > \frac{1}{2}$ of correct detection under each of the three hypotheses and can be determined as follows. Let $p_l^{(g)}$ and $p_l^{(f)}$ denote, respectively, the probability that the local test output moves the random walk closer to the target when this node contains the target and when it does not (see Fig. 2). Both are functions of the sample size $K_l$. We have

$$p_l^{(g)} = \Pr\left\{ \prod_{n=1}^{K_l} \log \frac{g_{l-1}(Y_n)}{f_{l-1}(Y_n)} > \max\left\{ \prod_{n=1}^{K_l} \log \frac{g_{l-1}(Z_n)}{f_{l-1}(Z_n)}, 0 \right\} \right\},$$

$$p_l^{(f)} = \left[\Pr\left\{ \prod_{n=1}^{K_l} \log \frac{g_{l-1}(Z_n)}{f_{l-1}(Z_n)} < 0 \right\}\right]^2, \tag{9}$$

where $\{Y_n\}_{n=1}^{K_l}$ and $\{Z_n\}_{n=1}^{K_l}$ are i.i.d. random variables with distribution $g_{l-1}$ and $f_{l-1}$, respectively. The parameter $K_l$ ($l = 1, 2, \ldots, \log_2 M$) is chosen to ensure that $p_l^{(g)} > p$ and $p_l^{(f)} > p$. Note that the value of $K_l$ can be computed offline and simple upper bounds suffice.

## IV. PERFORMANCE ANALYSIS

In this section, we establish the asymptotic optimality of the IRW policy in $c$ and the order optimality in $M$.

### A. Main Structure of the Analysis

Analyzing the Bayes risk of the IRW strategy lies in examining the trajectory of the biased random walk. With a confidence level $p > \frac{1}{2}$ in the local test module, the random walk will concentrate, with high probability, on a smaller and smaller portion of the tree containing the target and eventually probes the target only. Based on this insight, our approach is to partition the tree into $\log_2 M + 1$ half trees $\mathscr{T}_{\log_2 M}$, $\mathscr{T}_{\log_2 M-1}, \ldots, \mathscr{T}_0$ with decreasing size, and bound the time the random walk spent in each half tree. As illustrated in Fig. 3 for $M = 8$, $\mathscr{T}_l$ is the half tree (including the root) rooted at level $l$ ($l = \log_2 M, \log_2 M - 1, \ldots, 1$) that does
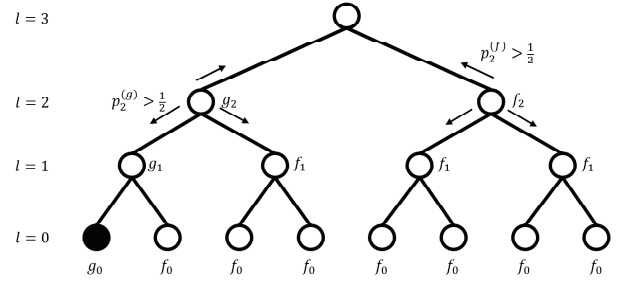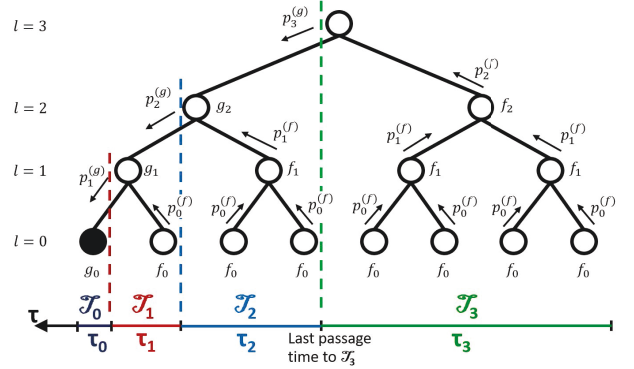


Fig. 2. A biased random walk on the tree.



Fig. 3. Partition of the tree into $\log_2 M + 1$ half trees.

not contain the target and $\mathscr{T}_0$ consists of only the target node. The entire search process, or equivalently, each sample path of the resulting random walk, is then partitioned into $\log_2 M + 1$ stages by the successively defined *last passage time* to each of the half trees in the shrinking sequence. In particular, the first stage with length $\tau_{\log_2 M}$ starts at the beginning of the search process and ends at the last passage time to the first half tree $\mathscr{T}_{\log_2 M}$ in the sequence, the second stage with length $\tau_{\log_2 M-1}$ starts at $\tau_{\log_2 M}+1$ and ends at the last passage time to $\mathscr{T}_{\log_2 M-1}$, and so on. Note that if the random walk terminates at a half tree $\mathscr{T}_l$ with $l > 0$ (i.e., a detection error occurs), then $\tau_j = 0$ for $j = l - 1, \ldots, 0$ by definition. It is easy to see that, for each sample path, we have the total time of the random walk equal to $\sum_{l=0}^{\log_2 M} \tau_l$.

In the next two subsections, we address separately the two cases of (i) the observations at all levels are informative (i.e., the KL divergence is bounded away from zero) and (ii) the aggregated observations decay to pure noise (i.e., diminishing KL divergence).

### B. Informative Observations at All Levels

We first consider the case where the KL divergence between aggregated observations in the presence and the absence of anomalous processes is bounded away from zero at all levels of the tree. The theorem below characterizes the Bayes risk of IRW.

*Theorem 1:* Assume that there exists a constant $\delta > 0$ independent of $M$ such that $D(g_l \| f_l) > \delta$ and $D(f_l \| g_l) > \delta$

for all $l = 1, 2, \ldots, \log_2 M$. For all $M$ and $c$, we have

$$\mathcal{R}(\text{IRW}) \geq cB \log_2 M + \frac{c \log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(c), \quad (10)$$

where $B$ is a constant independent of $c$ and $M$. Furthermore, the Bayes risk of IRW is order optimal in $M$ for all $c$ and asymptotically optimal in $c$ for all $M$ greater than a finite constant $M_0$.

*Proof:* See Appendix B. ∎

The optimality of the Bayes risk of IRW in both $c$ and $M$ directly carries through to the sample complexity of IRW. Specifically, from (10), we have the following upper bound on the sample complexity

$$Q(\text{IRW}) \geq B \log_2 M + \frac{\log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(1). \quad (11)$$

For a fixed $M$, we readily have

$$Q(\text{IRW}) \to \frac{\log c}{D(g_0 \| f_0)}.$$

For a fixed $c > 0$, we have

$$Q(\text{IRW}) = O(\log_2 M).$$

Using the lower bound on the sample complexity which was developed in Theorem 2 in [10], for any policy $\pi$, we have

$$Q(\pi) \gtrsim \frac{\log_2 M}{I_{\max}} + \frac{\log(1/c)(c)}{D(g_0 \| f_0)} + O(1), \quad (12)$$

where $I_{\max}$ denotes the maximum mutual information between the true hypothesis and the observation under an optimal action. It is not difficult to see that the sample complexity of IRW is asymptotically optimal in $c$ and order optimal in $M$. We point out that our analysis here focuses on establishing asymptotic/order optimality for general observation models. Hence, we resort to bounding techniques that are generally applicable but may lead to loose bounds on the leading constants for specific cases.

### C. Aggregated Observations Decaying to Pure Noise

When the quality of higher level measurements decays sufficiently fast, the sample size $K_l$ of the local test may increase unboundedly with $l$ ($l = 1, \ldots,$). Nevertheless, since the number of levels is $\log_2 M$, a sublinear scaling with $M$ is still attainable at moderate decaying rate of the aggregated observations. As a case study, we consider the Bernoulli distribution, which is widely adopted in the literature of hypothesis testing and group testing and also arises in distributed detection of aggregating local binary decisions. We establish sufficient conditions on the decaying rate of the quality of the hierarchical observations under which the IRW policy achieves a sublinear sample complexity in $M$.

Assume that $f_l$ and $g_l$ follow Bernoulli distributions with parameters $u_l$ and $1 - u_l$, respectively. In other words, the false alarm and miss detection probabilities at level $l$ are given by $u_l$. The KL divergence between $g_l$ and $f_l$ is $D(g_l \| f_l) = D(f_l \| g_l) = (1 - 2\mu_l) \log \frac{1-\mu_l}{\mu_l}$. We consider the case that $\mu_l$ increases with $l$ and converges to $\frac{1}{2}$ as $M$ approaches infinity.

In this case, both $D(g_l \| f_l)$ and $D(f_l \| g_l)$ converge to zero, which leads to unbounded $K_l$. The following two theorems characterize the sample complexity of IRW when $\mu_l$ converges to $\frac{1}{2}$ in polynomial order and exponential order, respectively.

*Theorem 2:* Assume that $\mu_l = \frac{1}{2} - (\frac{1}{2} - \mu_0)(l+1)^{-\alpha}$ ($l = 0, 1, 2, \ldots, \log_2 M$) for some $\alpha \in \mathbb{Z}^+$ and $\mu_0 < \frac{1}{2}$. The Bayes risk of the IRW policy is upper bounded by:

$$\mathcal{R}(\text{IRW}) \geq O(c)\log_2 M^{[2\alpha+1]} + \frac{c \log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(c). \quad (13)$$

*Proof:* See Appendix C. ∎

The case specified in Theorem 2 corresponds to a polynomial decay of the KL divergence: $D(g_l \| f_l) = D(f_l \| g_l) \sim (l+1)^{-2\alpha}$. In this case, IRW offers a sample complexity that is poly-logarithmic in $M$:

$$Q(\text{IRW}) = O(\log_2 M^{[2\alpha+1]}).$$

*Theorem 3:* Assume that $\mu_l = \frac{1}{2} - (\frac{1}{2} - \mu_0)\alpha^{-l}$ ($l = 0, 1, 2, \ldots, \log_2 M$) for some $\alpha > 1$ and $\mu_0 < \frac{1}{2}$. The Bayes risk of the IRW policy is upper bounded by:

$$\mathcal{R}(\text{IRW}) \geq c\tilde{B} M^{\log_2 \alpha^2} + \frac{c \log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(c), \quad (14)$$

where $\tilde{B}$ is a constant independent of $c$ and $M$.

*Proof:* See Appendix C. ∎

The case specified in Theorem 3 corresponds to a exponential decay of the KL divergence: $D(g_l \| f_l) = D(f_l \| g_l) \sim \alpha^{-2l}$. In this case, IRW offers a sample complexity that is sub-liner in $M$ provided that $1 < \alpha < \sqrt{2}$:

$$Q(\text{IRW}) = O(M^{\log_2 \alpha^2}).$$

## V. MULTI-TARGET DETECTION

We now consider the problem of detecting $L$ ($L \sim 0$) anomalous processes. We show that an extension of the IRW policy preserves the asymptotic optimality in $c$ and the order optimality in $M$ even when the number $L$ of targets is unknown.

Let $h_l^{(d)}$ ($l = 0, 1, \ldots, \log_2 M$, $d \geq \min\{L, 2^l\}$) denote the distribution of the measurements that aggregate $d$ anomalous processes and $2^l - d$ normal processes. An example with $M = 8$ and $L = 3$ is shown in Fig. 4. For a given $d$, we assume that for any $d' \geq d + 1$, we have

$$D\left(h_l^{(d)} \| h_l^{(d')}\right) - D\left(h_l^{(d)} \| h_l^{(d'+1)}\right) > 0, \quad (15)$$

and for any $d' \sim d$, we have

$$D\left(h_l^{(d)} \| h_l^{(d')}\right) - D\left(h_l^{(d)} \| h_l^{(d'+1)}\right) < 0. \quad (16)$$

The above monotonicity assumption on the KL divergence between $h_l^{(d)}$ and $h_l^{(d')}$ implies that a bigger difference $|d - d'|$ in the number of targets contained in the aggregated measurements leads to more distinguishable observations, which is usually true in practice. For the analysis in this section, we focus on the case that observations at all levels are informative. In other words, we assume that there exists a constant $\delta > 0$ independent of $M$ such that $D(h_l^{(d+k)} \| h_l^{(d)}) > \delta$ for all $l = 1, 2, \ldots, \log_2 M$, $d = 0, 1, \ldots \min\{L, 2^l\}$, and $d \geq k \geq \min\{L, 2^l\} - d$.
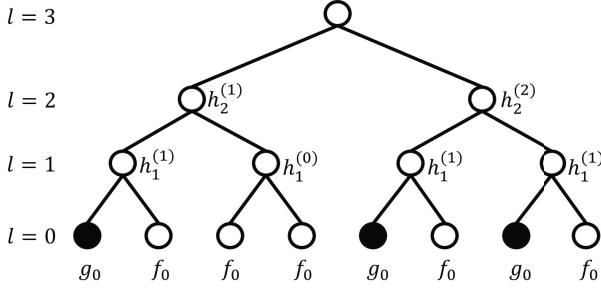
Fig. 4. A binary tree observation model with multiple targets.

## A. IRW Policy for Known L

For multi-target detection, IRW locates the $L$ ($L > 1$) targets one by one[3]. Similar to the single-target case, a biased random walk initiates at the root of the tree and eventually arrives at an undeclared target with high probability (referred to as one run of the random walk). The random walk is then reset to the root until $L$ targets have been declared.

The local test $\mathcal{U}(p)$ on an upper-level node differs from the single-target case in that it now faces four hypotheses, with an addition of hypothesis $H_3$ that both children contain undeclared targets. The outputs of $H_0$, $H_1$, and $H_2$ of the local test guide the random walk in the same way as in the single-target case. When the local test outputs $H_3$, the random walk arbitrarily chooses one child to zoom in. The stopping rule and declaration rule at each run of the random walk remain the same as in the single-target case.

To specify the local test module $\mathcal{U}(p)$, suppose that the random walk is currently at a node on an upper level $l > 0$, whose left and right child contain, respectively, $d_L$ and $d_R$ declared targets. Note that the local test faces a composite hypothesis testing problem, since both the left and right children may contain more than one target. Consider a fixed-sample test where $K_l^{(d_L)}$ and $K_l^{(d_R)}$ samples are taken from the left and right child, respectively. Note that the number of samples taken from each child depends on the number of declared targets. The SLLR of left child is computed as

$$\prod_{n=1}^{K_l^{(d_L)}} \log \frac{h_l^{(d_L+1)}(y_n)}{h_l^{(d_L)}(y_n)}. \tag{17}$$

The SLLR of the right child is similarly obtained. Based on the monotonicity assumption specified in (15) and (16), the expected value of the log-likelihood ratio (LLR) of each sample in (17) is positive when there are undeclared targets contained in the tested child, and is negative otherwise. The local decision rule is thus based on whether the SLLRs of the two children are both negative ($H_0$), both positive ($H_3$), or one negative one positive ($H_1$ or $H_2$). Similar to the single-target case, the values of $K_l^{(d_L)}$ and $K_l^{(d_R)}$ are chosen to guarantee the probability of declaring the correct hypothesis is greater

[3]We do not assume that declared targets can be removed and excluded from future observations, which represents a simpler version of the problem.

than $p$. A sequential local test and an active local test can be similarly obtained.

The theorem below characterizes the Bayes risk of the IRW policy in terms of both $M$ and $c$.

*Theorem 4:* Assume that there exists a constant $\delta > 0$ independent of $M$ such that $D(h_l^{(d+k)}||h_l^{(d)}) > \delta$ for all $l = 1, 2, \ldots, \log_2 M$, $d = 0, 1, \ldots \min\{L, 2^l\}$, and $d \geq k \geq \min\{L, 2^l\} - d$. For all $M$, $c$, and $L$, we have

$$R_{IRW} \geq cLB \log_2 M + \frac{cL \log \frac{\log_2 M}{c}}{D(g_0||f_0)} + O(c^2 \log_2 M) + O(c), \tag{18}$$

where $B$ is a constant independent of $c$, $M$, and $L$. Furthermore, the Bayes risk of IRW is order optimal in $M$ for all $c$ and asymptotically optimal in $c$ for all $M$ sufficiently large.

*Proof:* We present below the basic structure of the proof. Details can be found in Appendix D. Similar to the single-target case, the risk associated with a wrong decision is bounded by $O(c)$. In analyzing the sample complexity, a uniform bound on the sample complexity of finding each target, i.e., a uniform bound on each run of the random walk, is derived. Such a bound is again obtained by partitioning the tree into $\log_2 M + 1$ subsets. These subsets, however, differ from the sub-trees in the single-target case. As illustrated in Fig. 5 for an example with $M = 8$ and $L = 3$, subset $\mathcal{T}_0$ consists of all the targets. Subset $\mathcal{T}_l$ ($l = 1, \ldots, \log_2 M$) is the union of all the nodes on level $l$ that contain at least one target, and their entire left or right subtree if the subtree has no target. We then show that the successively defined last passage time to each of the subsets from $\mathcal{T}_{\log_2 M}$ to $\mathcal{T}_1$ are bounded by a constant. ■

## B. IRW Policy for Unknown L

When the number $L$ (may be zero) of targets is an unknown constant independent of $M$, the IRW policy can be augmented with a *terminating phase* carried out at the root of the tree. Specifically, at the beginning of each run of the random walk when the local test takes samples from each child of the root node, if the local test indicates neither child contains undeclared targets, the policy enters the *terminating phase* and starts taking samples from the root node itself. The SLLR of the root node is updated sequentially with each sample. A positive SLLR initiates the next run of the random walk. A
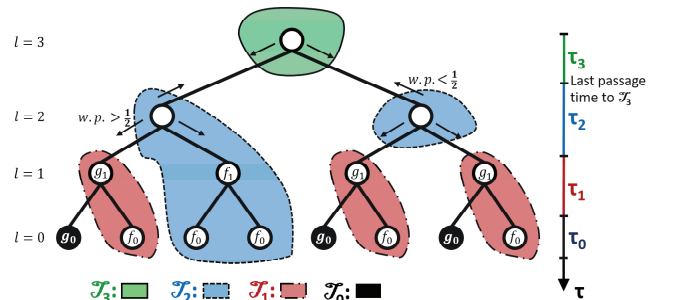


Fig. 5. Partition of the tree into $\log_2 M + 1$ subsets when there are multiple targets.

negative SLLR that drops below $\log c$ terminates the detection process.

Following the similar lines of arguments as in the proof of Theorem 4, the Bayes risk of the IRW policy when the number of targets is unknown can be upper bounded by

$$
R(\text{IRW}) \geq cLB \log_2 M + \frac{cL \log \frac{\log_2 M}{c}}{D(g_0 \| f_0)}
$$
$$
+ \frac{c \log \frac{1}{c}}{D\left(h_{l_r}^{L} \| h_{l_r}^{L+1}\right)} + O(c^2 \log_2 M) + O(c). \tag{19}
$$

Comparing (19) with (18), we see that additional samples are required to take at the root node in the terminating phase. Depending on the decaying speed of the KL-divergence, the additional sample complexity can be large if $D\left(h_{l_r}^{L} \| h_{l_r}^{L+1}\right)$ diminishes fast. However, as long as the KL-divergence on each level is bounded away from zero, IRW remains to be order-optimal in $M$ and asymptotically optimal in $c$. In Section VIII-D, we show the impact of unknown $L$ via simulation examples.

## VI. EXTENSIONS

In this section, we discuss extensions in two directions: a general tree structure and a general cost measure.

### A. IRW policy for General Tree Structures

As illustrated in Fig. 6, each leaf of the tree follows either the distribution $g_0$ (target) or $f_0$ (non-target), although the path length from each leaf to the root may be different. The observations at a high-level node follow the distributions that aggregate all its leaf-level descendants. Let $h_{a,b}$ denote the distribution of the measurements that aggregate $a$ anomalous processes and $b$ normal processes.

Assuming that the number of targets $L$ is known, the IRW policy operates in a similar way as presented in Section V. At each node of the tree, the objective of the local test is to guide, with probability greater than $1/2$, the global walk toward the child that contains undeclared targets. In particular, $K_l(d, w)$ samples are taken from each child where $d$ is the number of declared targets and $w$ is the number of the leaf nodes rooted at the child being tested. The SLLR of the sampled child is updated as

$$
\sum_{n=1}^{K_l(d,w)} \log \frac{h_{d+1,w-d-1}(y(n))}{h_{d,w-d}(y(n))}. \tag{20}
$$

If SLLRs of all the children are negative, the local test declares that no child contains undeclared targets, and the random walk goes back to the parent of the current node. Otherwise, the local test declares the child with the largest SLLR as containing undeclared targets, and the random walk zooms into that child. The number $K_l(d, w)$ of samples is chosen to guarantee the biasedness of the random walk.

Following similar lines as in the proof of Theorem 4, we can show the Bayes risk of the IRW policy is $O(cLHD) + O(cL \log \frac{1}{c})$, where $H$ is the height of the tree, $D$ is the maximum degree of the tree.
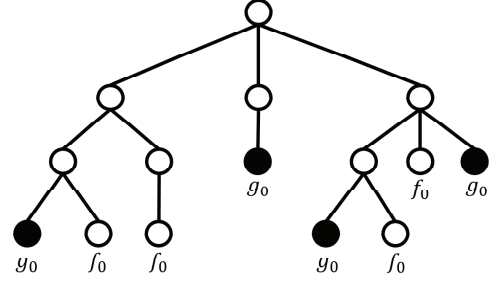


Fig. 6. A general tree with bounded degree.

For a graph structure that is not a tree, least informative edges can first be trimmed to obtain a tree. IRW can then be applied to the resulting tree. An optimal construction of a tree from an arbitrary graph is an interesting future direction.

### B. IRW Policy under Size-Dependent Costs

The work introduced in all the previous sections follows the Bayesian approach in Chernoff's original work [6], and the objective is to find an algorithm that achieves the optimal Bayesian risk in $c$ and $M$. We now discuss the case where the cost of testing a node depends on the location of the node on the tree. More specifically, the cost is a function of the size of the testing set.

Let the cost at the leaf level be $c$ and the cost of testing a node at level $l$ be $c_l$ which is a function of $l$ and $c$. A similar upper bound on the Bayesian risk of IRW can be obtained by following similar lines of arguments in Appendix B. If $c_l$ is bounded and $\lim_{c \to 0} c_l = 0$ for all $l = 1, 2, ..., \log_2 M$, then all the theorems shown in the paper hold. If, however, $c_l$ increases with $l$ and approaches infinity when $M$ goes to infinity, the optimality of the IRW policy in terms of minimizing the Bayesian risk is an open question and requires a separate investigation.

## VII. DISCUSSIONS

Several problems studied in different context can be cast as an active search problem. In this section, we discuss these connections and the applicability of the IRW strategy.

### A. Adaptive Sampling with Noisy Response

Consider the problem of estimating a step function in $[0, 1]$ (i.e., estimating a zero crossing). Let $z^* \in [0, 1]$ denote the unknown step of the function. The learner sequentially chooses sampling points in the interval and observes a noisy version of their values. Two commonly used noise models are the additive Gaussian noise model and the boolean noise model. In the former, the observations are Gaussian with mean 0 or 1, depending on whether the sampling point is to the left or the right of $z^*$. In the latter, the observations are Bernoulli random variables with parameters depending on the relative locations of the sampling points and $z^*$. The objective is to locate $z^*$ in a $\delta$-length interval with a probability no smaller than $1 - \varepsilon$. This problem arises in active learning of binary threshold classifiers [24] and stochastic root finding [25].

The problem can be cast as one studied in this work. We partition the $]0, 1[$ interval into $\delta$-length sub-intervals, which form the $M = 1(\delta$ leaf nodes with the sub-interval containing $z^{\equiv}$ being the target. Successively combining two adjacent sub-intervals leads to a binary tree with the root being the entire interval of $]0, 1[$. What remains to be specified is the local test on a node. Since each node on the tree is a sub-interval of $]0, 1[$, there remains the issue of which points in this sub-interval to sample when we prob this node. One way to translate the problem is to set the sampling points to the boundaries of each sub-interval. Thus, for the ternary hypothesis testing problem at each upper-level node, the observations are random vectors of dimension 4, corresponding to sampling the two boundaries of each of the children. Similarly, at a leaf node, the observations are of dimension 2. The local tests can be easily extended. The IRW policy ends up with a sample complexity equals

$$2B \log_2 \frac{1}{\delta} + \frac{\log \frac{\log_2 M}{\varepsilon}}{D)g_0 \sqrt{f_0[}} + O)1[,$$

where $B$ is a constant introduced in Theorem 1, and $g_0$ and $f_0$ are the distributions of the noisy responses when the sampling is to the right or left of $z^{\equiv}$ respectively.

Most of the existing work on adaptive sampling is based on a Bayesian approach with a binary noise model. A popular Bayesian strategy, the probabilistic bisection algorithm, updates the posterior distribution of the step location after each sample (based on the known model of the noisy response) and takes the next sample at the median point of the posterior distribution. Several variations of the method have been extensively studied in the literature [26]–[30]. In particular, in [30], Lalitha et al. proposed a two-stage algorithm based on the posterior matching method and showed the gain in sample complexity over non-adaptive/open-loop strategies. However, the update and sorting of the posterior probabilities at each sample require $O)M \log M)$ computation and memory complexity. In contrast, the IRW approach assumes no prior distribution and has $O)1[$ computation and memory complexity. This is made possible by effectively localizing data processing to small subsets of the input domain based on the tree structure, which also allows dynamic allocation of limited data storage.

## B. Noisy Group Testing

In the group testing problem, the objective is to identify the defective items in a large population by performing tests on subsets of items that reveal whether the tested group contains any defective items (classical Boolean group testing [31], [32]) or the number of defective items in the tested group (quantitative group testing [3], [33]).

Various formulations of group testing can be mapped to an active search problem with specific observation models (e.g., Bernoulli distribution for noisy Boolean group testing, sum-observation model for the quantitative and threshold group testing). The individual items in the group testing problems are mapped to the leaf nodes of a tree. The action of testing a node on the tree corresponds to a group test.

Most existing work on Boolean group testing assumes error-free test outcomes. There are several recent studies on the noisy group testing that assume the presence of one-sided noise [31], [34] or the symmetric case with equal size-independent false alarm and miss detection probabilities [32], [35]. In some extended group testing models such as the noisy quantitative group testing [3] and threshold group testing [33], the issue of sample complexity in terms of the detection accuracy is absent in the basic formulation. Most of the existing results on noisy group testing focus on non-adaptive open-loop strategies that determine all actions in one shot *a priori*. The work by Kaspi et. al [36] studied a problem of detecting multiple targets that are uniformly placed on the unit interval. The observation of each measurement follows a specific Bernoulli additive noise model, which can be mapped to a binary multiple access channel in the channel coding problem. Kaspi et. al proposed a non-adaptive strategy and an adaptive strategy consisting of a series of non-adaptive measuring phases. The disadvantages of non-adaptive strategies lie in the computational complexity of the coding/decoding processes and high storage requirement. Another recent work by Scarlett [37] on the noisy group testing proposed a hybrid adaptive group testing algorithm where the first stage applies non-adaptive algorithms, and the second stage (and third stage in a refined version) improves the initial estimate of the first stage. Similar to the aforementioned related work, the observation models are restricted to binary symmetric noise or one-sided noise (Z-channel). The sampling complexity of these algorithms is shown to be order-optimal in terms of the population size. However, it is given in an asymptotic setting as the probability of error approaching zero, i.e. the accuracy constraint is not part of the given sampling complexity.

In contrast, IRW provides an *adaptive* solution to noisy group testing under general noisy observations and with little offline or online computation and low memory requirement. Consider a noisy group testing problem in which $L$ targets need to be identified from $M$ items with error probability less than $\varepsilon$. The measurement of a group of items follow a general distribution that depends on the number of targets in the group. IRW falls into the family of nested test plans. Specifically, treating all the $M$ items as the leaf nodes, one can easily construct a binary tree with height of $\log_2 M$ as shown in Fig. 4. Testing on an upper-level node corresponds to measuring the group of all its leaf-level descendants.

Although adaptive group testing strategies do not necessarily conform to a predetermined tree structure, IRW offers asymptotic optimality in both population size and reliability constraint. It has a sample complexity of $BL \log_2 M + \frac{L \log \frac{\log_2 M}{\varepsilon}}{D)g_0 \sqrt{f_0[}} + O)\varepsilon \log_2 M[$, where $g_0$ and $f_0$ are the distributions of the measurement from the target and non-target items respective, and $B$ is a constant independent of $L$ and $M$ as introduced in Theorem 4. In particular, IRW is the first adaptive test plan for noisy group testing under a general noise model and offers asymptotic optimality in terms of the required detection accuracy and order optimality in terms of the population size. As a point of comparison, the one-sided boolean noise model considered in [31] is a special case of the general model considered here, with $g_0 =$ Bernoulli$)p[$ and $f_0 \leq 0$. For this special case, the expected number of tests

under IRW is upper bounded by $O(\frac{L\log_2 M}{p^2})$, which achieves the same result shown in [31].

### C. Channel Coding with Feedback

In channel coding with feedback [38]–[40], the encoder transmits symbols adaptively based on the receiving history of the decoder due to the availability of a noiseless feedback channel. This coding problem under an arbitrary discrete-input memoryless channel can be mapped to the problem studied in this work. Without loss of generality, consider a Discrete Memoryless Channel (DMC) where the output is also discrete. Let $\{I_1, I_2, \ldots, I_J\}$ and $\{O_1, O_2, \ldots, O_K\}$ denote, respectively, the input and output alphabets. Let $\mathbf{P} = \{p_{j,k}\}$ where $j = 1, \ldots, J$, $k = 1, \ldots, K$ be the channel transition probability matrix. Let $M$ be the number of messages. The objective is a coding scheme that transits these $M$ messages successfully with probability no smaller than $1 - \varepsilon$.

The above coding problem can be mapped to an active search problem on a binary tree with $M$ leaf nodes of which one is the target, i.e., the message being transmitted. The binary splitting structure generates a representation of the location of the target with a $\{0,1\}$ codeword of which the length equals $\log_2 M$ (i.e., for each node, the left branch represents 0 and the right branch represents 1). The test on each level of the tree corresponds to sending the next bit or correcting the previous bit of the source code. The distribution of the target $g_0$ is set as the probability mass vector $\{p_{j_g^\equiv,k}\}_{k=1}^K$ and the distribution of the non-target node $f_0$ is set as the probability mass vector $\{p_{j_f^\equiv,k}\}_{k=1}^K$, where $j_g^\equiv$ and $j_f^\equiv$ are the two most distinguishable symbols transmitted through the channel and are defined as

$$\{j_g^\equiv, j_f^\equiv\} = \arg \max_{\{j_g, j_f\}} \prod_{k=1}^K p_{j_g,k} \log \frac{p_{j_g,k}}{p_{j_f,k}}, \forall j_g, j_f = 1, \ldots, J. \tag{21}$$

The observation distribution of a node on level $l - 1$ also follows $g_0$ if it contains the target or $f_0$ if it does not. It is not difficult to see that the action of sampling a node which contains the target in the target search problem corresponds to the action of sending $j_g^\equiv$ through the channel, and the action of sampling a node which does not contain the target corresponds to the action of sending $j_f^\equiv$ through the channel. The corresponding observations at the receiving end of the channel follow $g_0$ and $f_0$, respectively.

IRW can be mapped to a coding scheme for the transmission problem. If the next bit of the source code is 0, i.e., the left child of the current node contains the target, the sender sends $K_l$ times of symbol $j_g^\equiv$ followed by sending $K_l$ times of symbol $j_f^\equiv$ through the channel. If the next bit of the source code is 1, i.e., the right child contains the target, the sender sends $K_l$ times of symbol $j_f^\equiv$ followed by sending $K_l$ times of symbol $j_g^\equiv$ through the channel. If there is an error in the transmission of previous bits, i.e., neither of the children contains the target, the sender sends $2K_l$ times of symbol $j_f^\equiv$ to inform the encoder to correct the previous bit. After each local test ($2K_l$ times channel usages), a bit of the source code is decoded correctly with probability greater than $1 - 2$. If a bit is decoded incorrectly, it would be revisited and corrected later with probability greater than $1 - 2$. When the random walk arrives at a leaf node, i.e., the full codeword has been transmitted and decoded, the sender keeps sending symbol $j_g^\equiv$ if the entire codeword has been decoded correctly until the log-likelihood ratio at the receiver is large enough. If any bit of the codeword is decoded incorrectly, the sender keeps sending $j_f^\equiv$. This step of sending the confirmation bits corresponds to the local test on a leaf in the IRW policy. The total number of transmissions is then upper bounded by $B \log_2 M + \frac{\log \frac{\log_2 M}{\varepsilon}}{D(g_0 \| f_0)} + O(\varepsilon)$, which is order optimal in both $M$ and $\varepsilon$ [38], [39].

A similar connection between the active search problem and channel coding with feedback was discussed in [30], where an additive Gaussian noise channel with binary input was considered. The IRW strategy applies to more general channel models. Its advantage in computation and memory efficiency as discussed in the case of adaptive sampling applies here as well.

## VIII. SIMULATION EXAMPLES

In this section, we demonstrate the performance of IRW in various settings.

### A. Comparison of IRW with the Chernoff test and DGF test

Consider the problem of detecting $L$ heavy hitters among Poisson flows where the measurements are exponentially-distributed packet inter-arrival times. For the leaf-node, $g_0$ and $f_0$ are exponential distributions with parameters $\lambda_g$ and $\lambda_f$, respectively. The aggregated flows follow the corresponding exponential distributions with the parameters equal to the sum of the parameters of their children at the leaf level. Under this observation model, we have

$$D(g_i \| f_i) \sim D(g_{i+1} \| f_{i+1}), \ D(f_i \| g_i) \sim D(f_{i+1} \| g_{i+1}).$$

It can be shown that, while the action space consisting of all the nodes on the tree, the resulting Chernoff test probes only the leaf nodes. Specifically, at each time $t$, all the leaf nodes are sorted based on their SLLRs. If

$$D(g_0 \| f_0)(L - D(f_0 \| g_0)(M - L), \tag{22}$$

the Chernoff test uniformly at random selects a leaf node among the set with the top $L$ largest SLLR. Otherwise, the Chernoff test uniformly at random selects a leaf node among the $(M - L)$ nodes in the tail of the SLLR ranking. In large-scale systems where $M$ is sufficiently large, the condition in (22) holds, and the Chernoff test selects the leaf node which it believes to be the target. The Chernoff and the IRW tests have the same stopping and decision rules.

Fig. 7 and Fig. 8 show simulation results on the sample complexity as a function of $M$ for $L = 1$ and $L = 5$, respectively. We observe that the Chernoff test has a sample complexity that scales linearly with $M$, while IRW offers a logarithmic order. The improvement at $M = 20$ is already three-fold.

Also shown in these figures is the DGF policy developed in [11]. The comparison with DGF is not on equal footing,
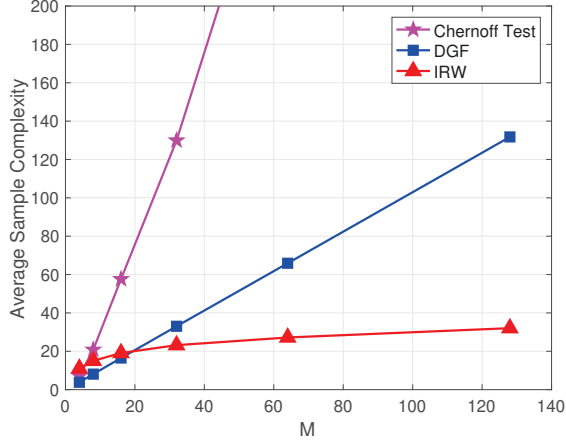
Fig. 7. Performance comparison ($L = 1$, $\lambda_g = 10$, $\lambda_f = 0.01$, $K_l = 3$, $c = 10^{13}$, $M = 4, 8, \ldots, 128$, 1000 Monte Carlo runs).
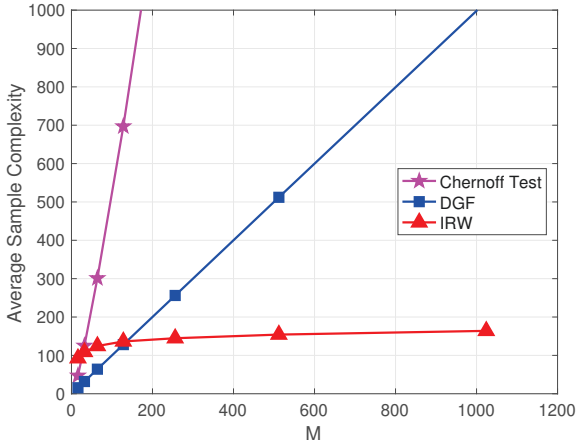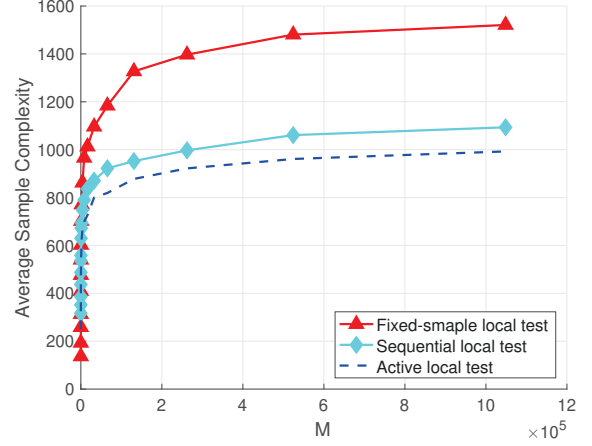


Fig. 9. Performance comparison of IRW with different local tests ($K_l = 7$; thresholds for the sequential local test: $\gamma_1 = 1.0986$, $\gamma_0 = 1.0986$; thresholds for the active local test: $\nu_1 = 0.9445$, $\nu_0 = 0.9445$; 1000 Monte Carlo runs).



Fig. 8. Performance comparison ($L = 5$, $\lambda_g = 10$, $\lambda_f = 0.001$, $c = 5 * 10^5$, $M = 16, 32, \ldots, 1024$, 1000 Monte Carlo runs).

since it was developed without assuming a tree structure. The DGF policy, as a deterministic policy, has a much smaller linear slope than the Chernoff test. Specifically, DGF probes the leaf node with the $L$-th largest SLLR if (22) holds, otherwise it probes the leaf node with the $(L + 1)$-th largest SLLR. The comparison between DGF and IRW is to show that by exploiting the hierarchical structure of the search space, more significant gain can be achieved in addition to efficient design of deterministic selection rules.

### B. Comparison of IRW with Different Local Tests

Next, we study the impact of different local tests on the finite-time performance of IRW. In the example shown in Fig. 9, $L = 1$, and the target-present and target-absent distributions are level independent and Bernoulli with parameters of 0.6 and 0.4, respectively. The confidence level of the local test is set to $p = 0.5625$, which determines the number $K_l$ of samples in the fixed-sample-size test and the thresholds of the sequential and active tests. The improvement offered by the sequential and active local tests is evident. In addition, instead

of designing $K_l$ for each level, one pair of thresholds for the SLLR that guarantees the biased global random walk can be used on all higher-level nodes.

### C. IRW Applied to Channel Coding with Feedback

In this numerical example, we apply the IRW policy to channel coding with feedback as we discussed in Section VII-C. The channel is a Binary Symmetric Channel with noiseless feedback with a crossover probability of 0.3. It can be mapped to an active search problem on a binary tree with $M$ leaf nodes. The observation of the target node or any upper-level node which contains the target follows a Bernoulli distribution with the success probability $p$ equals to 0.7, and the observation of the non-target node follows a Bernoulli distribution with the success probability equals to 0.3. We compare IRW with the coding procedure introduced in [40] by Burnashev which is optimal in terms of $M$. Fig. 10 shows that IRW with active local test outperforms Burnashev's coding procedure. In the same figure, we also illustrate the lower bound of the sample complexity given in (12) with $I_{\max} = 1 + p \log_2 p + (1 - p) \log_2(1 - p)$. As an example, at $M = 512$, the IRW policy requires about 6.2% fewer samples than Burnashev's code scheme, and about 25% more samples when compared to the derived lower bound.

### D. The Impact of Unknown L

In this section, we present simulation results to demonstrate the impact of not knowing the total number $L$ of targets. We adopt a size-independent observation setting and assume that on each level, if the testing node contains at least one target, the observation follows an exponential distribution with parameter $\lambda_g = 10$. Otherwise it follows an exponential distribution with parameter $\lambda_f = 1$. Due to the size-independent setting, we also assume that declared targets can be removed from future observations to make sure the positive observations are from undeclared target nodes. Fig. 11 shows the sample complexity obtained by IRW for cases where $L$ is known
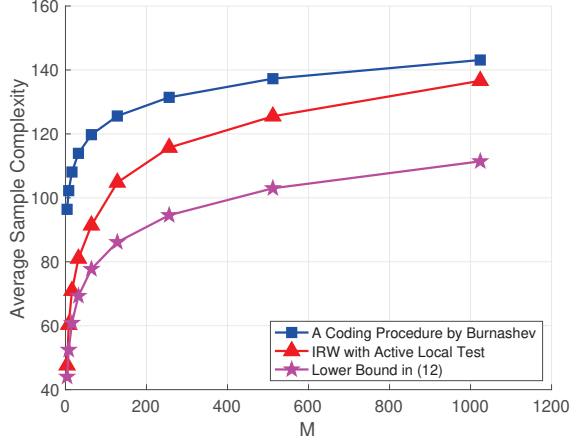
Fig. 10.  Performance comparison of IRW policy to the coding procedure in [40] by Burnashev for the BSC channel with a crossover probability of 0.3 (1000 Monte Carlo runs).
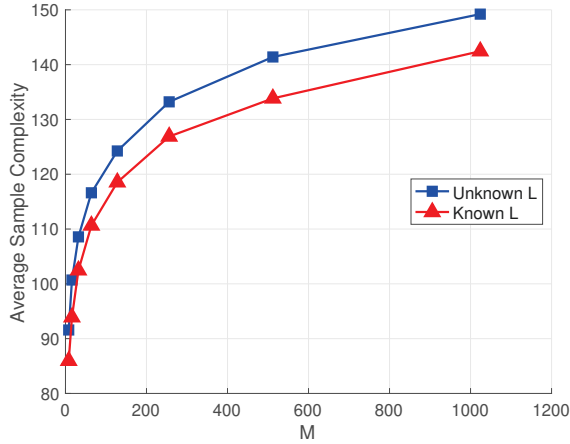


Fig. 11.  Performance comparison when the number of target $L$ is unknown before testing ($L = 3$, $\lambda_g = 10$, $\lambda_f = 1$, $c = 10^{-13}$, $M = 8, 16, \ldots, 1024$, 1000 Monte Carlo runs).

and unknown. We set $L = 3$ and increased $M$ from 8 to 1024. When $L$ is unknown, IRW requires more samples at the root node to terminate the test. For example, at $M = 512$, the relative increase in the sample complexity is about 4.2%, which is rather moderate.

## IX. Conclusion

In the paper, we considered the problem of detecting a few anomalous processes among a large number of processes under a tree observation model. The proposed active inference strategy induces a biased random walk on the tree and offers order optimality with the search space size and asymptotic optimality in terms of the reliability constraint. The proposed strategy is also efficient in terms of computation and memory requirement. By effectively localizing the data processing to small subsets of the search space, it has $O(1)$ computation and memory complexity. The results find a wide range of applications, including noisy group testing, channel coding with feedback, and adaptive sampling.

## Appendix A
## Alternative Local Tests

### A. Sequential Local Test

A sequential version of the local test is to carry out SPRT on each child, one at a time. If the SPRT on the left child indicates it contains the target, the local test declare $H_1$. Otherwise, carry out SPRT on the right child and declare $H_2$ or $H_0$ based on the outcome of this SPRT.

We now specify the thresholds to be used in each SPRT. To ensure the required confidence level $p$ of the local test, the false alarm and miss detection probabilities of each SPRT should satisfy $(1 - P_{FA})^2 > p$ and $(1 - P_{FA})(1 - P_{MD}) > p$. This leads to the following positive and negative thresholds, $\gamma_1$ and $\gamma_0$ respectively, of the SPRT

$$\gamma_1 = \log \frac{1 - P_{MD}}{P_{FA}}, \quad \gamma_0 = \log \frac{P_{MD}}{1 - P_{FA}}. \quad (23)$$

### B. Active Local Test

We now present a fully active test that adaptively determines which child to sample. We specify the test and leave the derivation details to the extended version [41].

The selection rule is to sample the child with a greater SLLR at each time. The stopping rule is determined by a pair of thresholds $v_0$ and $v_1$. The test stops as soon as $\max\{S_L, S_R\} \geq v_0$, and declares $H_0$; or when $\max\{S_L, S_R\} \sim v_1$, in which case either $H_1$ or $H_2$ is declared, depending on which $S_L$ or $S_R$ exceeds $v_1$. The thresholds are set as following:

$$v_1 = \log \frac{P_{11}}{P_{01}}, \quad v_0 = \log \frac{P_{10}}{P_{00}}, \quad (24)$$

where $P_{11}$ denote the probability of declaring hypothesis $H_1$ when $H_1$ is true, and $P_{00}$ are defined similarly. To ensure a confidence level of $p$, we set $P_{00} = P_{11} = p$, $P_{01} = (1 - P_{00})/2$ and $P_{10} = (1 - P_{11})/2$.

## Appendix B
## Proof of Theorem 1

Without loss of generality (due to the symmetry of the binary tree structure), we assume that the left-most leaf is the target. We focus on the IRW with fixed-sample local tests.

The random walk on the tree can be divided into two states. The first state is the random walk on upper-level nodes of the binary tree. In this state, at each time, after taking $K_l$ samples, we either zoom-in to one child node or zoom-out to the parent node. As a result, the distance between the current node to the target is defined as the sum of the discrete distance to the target node on the tree and the threshold $\log \frac{\log_2 M}{c}$, which either decreases by one when zooming-in or increases by one when zooming-out after taking $2K_l$ samples from the children. Once arriving at a leaf node, the test arrives at the second state, where samples are taken one by one from the current node until the cumulative SLLR exceeds the threshold or becomes negative. The cumulative SLLR can be viewed as a discrete time random walk with random continuous step size which is the LLR of each sample. For all the non-target leaf-nodes, we define the distance between the node to the target as the sum of the discrete distance on the tree, the cumulative SLLR of

the current node, and the threshold. For the target node, we define the distance to the target as the difference between the threshold and the current cumulative SLLR of the target node. During the search process, these two different states happen consecutively in Phase I of the IRW policy.

Let $W_n$ denote the random variable of the step size of the random walk at time $n$. When the IRW is in the first state (i.e., random walk on high-level nodes), depending on the current level $l > 0$, $W_n$ has the following distribution:

$$\Pr(W_n | \text{testing a target node}) = \begin{cases} p_l^g, & \text{for } W_n = -1, \\ 1 - p_l^g, & \text{for } W_n = 1, \end{cases} \tag{25}$$

or

$$\Pr(W_n | \text{testing a non-target node})$$
$$= \begin{cases} p_l^f, & \text{for } W_n = -1, \\ 1 - p_l^f, & \text{for } W_n = 1, \end{cases} \tag{26}$$

Since $p_l^g > \frac{1}{2}$ and $p_l^f > \frac{1}{2}$ for all $l = 1, 2, \ldots, \log_2 M$, we have:

$$\mathbb{E}[W_n | \text{testing a target node}] = 1 - 2p_l^g, \tag{27}$$
$$\mathbb{E}[W_n | \text{testing a non-target node}] = 1 - 2p_l^f, \tag{28}$$

which are both less than 0.

For the second state, let $Y_0$ and $Z_0$ denote the random variables with distributions $g_0$ and $f_0$, respectively. The LLR will be either $-\log \frac{g_0(Y_0)}{f_0(Y_0)}$ or $\log \frac{g_0(Z_0)}{f_0(Z_0)}$. It is not difficult to see that for the target node, we have:

$$\mathbb{E}[W_n] = \mathbb{E}\left[-\log \frac{g_0(Y_0)}{f_0(Y_0)}\right] = -D(g_0 \| f_0) < 0, \tag{29}$$

and for all the non-target nodes, we have:

$$\mathbb{E}[W_n] = \mathbb{E}\left[\log \frac{g_0(Z_0)}{f_0(Z_0)}\right] = -D(f_0 \| g_0) < 0. \tag{30}$$

We further assume that the distribution of $-\log \frac{g_0(Y_0)}{f_0(Y_0)}$ and $\log \frac{g_0(Z_0)}{f_0(Z_0)}$ are light-tailed distributions[4].

We now give a rigorous definition of the last passage times $\tau_i$ as introduced in Section IV-A. Define

$$G(t) = i, \text{ if at time } t \text{ the current testing node} \\ \text{is on the half tree } \mathcal{T}_i. \tag{31}$$

It indicates the half tree $\mathcal{T}_i$ where the random walk is or the IRW policy is currently probing at time $t$. The last passage time of the sub-tree at the root that does not contain the target, $\tau_{\log_2 M}$, is defined as

$$\tau_{\log_2 M} = \sup \{t > 0 : G(t) < \log_2 M\}. \tag{32}$$

Based on the recursive formulations, we have

$$\tau_i = \sup \{t > 0 : G(t) < i\} - \tau_{i+1}, \tag{33}$$

for $i = 1, 2, \ldots, \log_2 M - 2, \log_2 M - 1$. The following lemma characterizes the distributions of $\tau_i$.

[4]A random variable $X$ with the cumulative distribution function $F(x)$ is light-tailed if and only if $\int_{\mathbb{R}} e^{\lambda x} dF(x) < \infty$ for some $\lambda > 0$ [42].

*Lemma 1:* For all $\tau_i$ with $i = 1, \ldots, \log_2 M$, there exist $\alpha > 0$ and $\gamma > 0$ which are independent of $M$ and $c$, such that

$$\Pr(\tau_i \sim n) \geq \alpha e^{-\gamma n}, \quad \forall n \sim 0. \tag{34}$$

*Proof:* We first prove this lemma for $\tau_{\log_2 M}$. Let $L_t$ denote the distance to the target at time $t$. The random walk starts at the root node. Therefore the initial distance to the target is $L_0 = \log_2 M + \log \frac{\log_2 M}{c}$. Define

$$\tau^{\equiv} = \sup \{t > 0 : L_t \sim L_0\} \tag{35}$$

as the last time when the search approaches the distance to the target which is greater than $L_0$. It is not difficult to see that

$$\tau_{\log_2 M} \geq \tau^{\equiv}. \tag{36}$$

Therefore, we have

$$\Pr(\tau_{\log_2 M} \sim n) \geq \Pr(\tau^{\equiv} \sim n). \tag{37}$$

Based on the definition of $\tau^{\equiv}$, we have

$$\Pr(\tau^{\equiv} > n) = \Pr(\sup \{t > 0 : L_t \sim L_0\} > n)$$
$$\geq \prod_{t=n}^{\infty} \Pr(L_t \sim L_0) = \prod_{t=n}^{\infty} \Pr\left(\prod_{j=1}^{t} W_j \sim 0\right). \tag{38}$$

Let $\mu_j$ denote the mean value for each $W_j$, where $\mu_j < 0$ for all $j = 1, 2, \ldots, t$. Note that $W_j$'s can take different values at different probing nodes with the distributions defined in (25), (26), (29), and (30) which are not required to be identical. Since $W_j$'s are independent, by applying the Chernoff bound to $\prod_{j=1}^{t} W_j$, we have, for all $s > 0$,

$$\Pr\left(\prod_{j=1}^{t} W_j \sim 0\right) \geq \mathbb{E}\left[e^{s \sum_{j=1}^{t} W_j}\right] = \prod_{j=1}^{t} \mathbb{E}\left[e^{sW_j}\right]. \tag{39}$$

Note that the moment generating function (MGF) of each $W_j$ is equal to one at $s = 0$. Furthermore, since $\mathbb{E}[W_j] < 0$ is strictly negative for all $j \sim 1$, differentiating the MGFs of all $W_j$ with respect to $s$ yields strictly negative derivatives at $s = 0$. Because all $W_j$'s are light-tailed distributions for all possible distributions of $W_j$, there exist $s > 0$ and $\gamma > 0$ such that $\mathbb{E}[e^{sW_j}]$ is strictly less than $e^{-\gamma} < 1$. Hence, from (39), we have

$$\Pr\left(\prod_{j=1}^{t} W_j \sim 0\right) \geq e^{-\gamma t}. \tag{40}$$

Due to (38), we have

$$\Pr(\tau^{\equiv} > n) \geq \prod_{t=n}^{\infty} \Pr\left(\prod_{i=1}^{t} W_i \sim 0\right)$$
$$\geq \prod_{t=n}^{\infty} e^{-\gamma t} = \frac{e^{-\gamma n}}{1 - e^{-\gamma}}. \tag{41}$$

Let $\alpha = \frac{1}{1 - e^{-\gamma}}$, with (37), we complete the proof of Lemma 1 proved for $\tau_{\log_2 M}$. Due to the recursive definitions of $\tau_1, \tau_2, \ldots, \tau_{\log_2 M}$, the proof follows the same procedure for all other $\tau_i$. ∎

Based on Lemma 1, we get the following lemma that characterizes the expected value of $\tau_i$.

*Lemma 2:* For all $\tau_i$ with $i = 1, \ldots, \log_2 M$, there exists a constant $\beta > 0$, such that

$$\mathbb{E}[\tau_i] \geq \beta. \tag{42}$$

*Proof:* Based on the the tail-sum formula of expectation of non-negative random variables, we have

$$\mathbb{E}[\tau_i] = \sum_{n=0}^{\infty} \Pr[\tau_i > n] \geq \sum_{n=0}^{\infty} \alpha e^{-\gamma n}$$
$$= \frac{\alpha}{1 - e^{-\gamma}} = \frac{1}{(1 - e^{-\gamma})^2} = \beta. \tag{43}$$

Now we are ready to prove Theorem 1. Let $\tau_{\text{IRW}}$ denote the total samples complexity of the proposed IRW policy. Based on Lemma 2, it is not difficult to show that

$$\mathbb{E}[\tau_{\text{IRW}}] \geq 2K_{\max} \sum_{i=1}^{\log_2 M} \mathbb{E}[\tau_i] + \mathbb{E}[\tau_0]$$
$$\geq 2\beta K_{\max} \log_2 M + \mathbb{E}[\tau_0]. \tag{44}$$

When the observations are informative at all levels, $K_{\max}$ is bounded by a constant. As a result, the first term on the RHS of (44) is upper bounded by $B \log_2 M$, where $B$ is a constant greater than $2\beta K_{\max}$.

For the last stage at the node $j$, let

$$\tau_0 = \min \{t : S_j(t) > a \text{ or } S_j(t) < 0\}, \tag{45}$$

where $a = \log \frac{\log_2 M}{c}$, denote the stopping time with respect to the i.i.d. sequence of the LLR $\{\log \frac{g_0(X_n)}{f_0(X_n)} : n \sim 1\}$, where $X_n$ denotes an i.i.d. random variable with distribution $g_0$. Following similar arguments as in Section 3.4 of [43], at first step, we prove that $\tau_0$ is finite almost surely. To do so, assume that $\Pr[\log \frac{g_0(X_n)}{f_0(X_n)} = 0] < 1$. Note that if $\log \frac{g_0(X_n)}{f_0(X_n)} = 0$ almost surely, we have $g_0(x) = f_0(x)$ almost everywhere, which makes it impossible to distinguish between the two hypotheses. Then, there exist $\varepsilon > 0$ and $\delta > 0$, such that either $\Pr[\log \frac{g_0(X_n)}{f_0(X_n)} > \varepsilon] = \delta$ or $\Pr[\log \frac{g_0(X_n)}{f_0(X_n)} < -\varepsilon] = \delta$. Let $m$ be an integer such that $m\varepsilon > a$. If $\Pr[\log \frac{g_0(X_n)}{f_0(X_n)} > \varepsilon] = \delta$, we have

$$\Pr\{S_j(k+m) - S_j(k) > a\} > 0$$
$$= \Pr\left\{\sum_{n=k+1}^{k+m} \log \frac{g_0(X_n)}{f_0(X_n)} > a\right\} \sim \delta^m$$

for all $k$; if $\Pr[\log \frac{g_0(X_n)}{f_0(X_n)} < -\varepsilon] = \delta$, we have

$$\Pr\{S_j(k+m) - S_j(k) < 0 - a\}$$
$$= \Pr\left\{\sum_{n=k+1}^{k+m} \log \frac{g_0(X_n)}{f_0(X_n)} < -a\right\} \sim \delta^m$$

for all $k$. For either case, it implies that for all integers (say $q$) we have:

$$\Pr[\tau_0 > qm]$$
$$= \Pr\{0 < S_j(k) < a \text{ for } k = 1, 2, ..., qm\} < (1 - \delta^m)^q.$$

Let $\sigma = (1 - \delta^m)^{-1}$ and $\kappa = (1 - \delta^m)^{1/m}$. For an arbitrary $t$, we can find $q$ such that $qm < t < (q+1)m$ and

$$\Pr[\tau_0 > t] \geq \Pr[\tau_0 > qm] \geq (1 - \delta^m)^q = \sigma\kappa^{(q+1)m} \geq \sigma\kappa^q.$$

This implies that $\Pr[\tau_0 < \infty] = 1$. Also, the geometric decay of $\Pr[\tau_0 > n]$ implies that the generating function $\mathbb{E}[\exp(u\tau_0)]$ is defined for $u < -\ln(\kappa)$, so that $\tau_0$ has finite moments. Due to the Wald's Equation [23], we then have

$$\mathbb{E}\left[\sum_{n=1}^{\tau_0} \log \frac{g_0(X_n)}{f_0(X_n)}\right] = \mathbb{E}[\tau_0] \mathbb{E}\left[\log \frac{g_0(X_n)}{f_0(X_n)}\right]. \tag{46}$$

i.e.,

$$\log \frac{\log_2 M}{c} + R_b = \mathbb{E}[\tau_0] D(g_0 \| f_0), \tag{47}$$

where $R_b$ is the overshooting at the threshold. Due to Lorden's inequality [44], we have

$$\mathbb{E}[R_b] \geq \frac{\mathbb{E}\left[\left(\log \frac{g_0(X_n)}{f_0(X_n)}\right)^2\right]}{\mathbb{E}\left[\log \frac{g_0(X_n)}{f_0(X_n)}\right]}. \tag{48}$$

Assuming that the first two moments of the LLR are finite, we then have

$$\mathbb{E}[\tau_0] = \frac{\log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(1). \tag{49}$$

The following lemma characterizes the error probability of the IRW policy.

*Lemma 3:* The error probability of the IRW policy is upper bounded by:

$$P_e \geq \beta c = O(c). \tag{50}$$

*Proof:* When the random walk arrives a non-target node, say node $j$, the probability of error (accepting $H_j$) equals to $\Pr(S_j \sim \log \frac{\log_2 M}{c})$. For the sequential probability ratio test, Wald [7] shows that

$$\Pr\left(S_j \sim \log \frac{\log_2 M}{c}\right) \geq \exp\left[-\log \frac{\log_2 M}{c}\right] = \frac{c}{\log_2 M}. \tag{51}$$

Let $N$ denote the random number of times of visiting these non-target leaf nodes in the IRW policy. The conditional error probability is upper bounded by $\frac{Nc}{\log_2 M}$. Based on the proof of Theorem 1, the expected value of $N$ is upper bounded by $\beta \log_2 M$. Therefore, by taking expectation, the error probability is bounded by

$$P_e \geq \frac{c}{\log_2 M} \times \mathbb{E}[N] \geq \frac{c}{\log_2 M} \times \beta \log_2 M = \beta c = O(c). \tag{52}$$

Combining (44), (49), Lemma 3, and the definition of Bayse risk in (3) completes the proof.

## APPENDIX C
### PROOF OF THEOREM 2 AND THEOREM 3

Follow the same lines of argument in the proof of Theorem 1, the sample complexity of the last stage $\tau_0$ still satisfies (49). We now give the upper bound of the sample complexity of the first $\log_2 M$ stages for the test with a fixed-size local test.

We focus on the Bernoulli distribution model, where $g_l$ and $f_l$ are Bernoulli distributions with false negative and false positive rates equal to $\mu_l$. In order to get the relation between $K_l$ and $\mu_{l-1}$, we first introduce the following lemma.

*Lemma 4 ( [45]):* Let $X_1, \ldots, X_n$ be independent Poisson trials such that $\Pr(X_i) = p_i$. Let $X = \sum_{i=1}^{n} X_i$ and $\nu = \mathbb{E}[X]$. Then, the following Chernoff bounds hold for $0 < \delta \leq 1$,

$$\Pr(X \sim (1+\delta)\nu) \geq e^{-\nu\delta^2/3}; \tag{53}$$

$$\Pr(X \geq (1-\delta)\nu) \geq e^{-\nu\delta^2/3}. \tag{54}$$

The IRW policy requires that $p_l^{g}$ and $p_l^{f}$ will satisfy (9). For $p_l^{f}$, we need to find the value of $K_l$ such that

$$\Pr\left(\prod_{n=1}^{K_l} \log \frac{g_{l-1}(Z_n)}{f_{l-1}(Z_n)} < 0\right) > \lambda.$$

i.e.,

$$\Pr\left(\prod_{n=1}^{K_l} \log \frac{g_{l-1}(Z_n)}{f_{l-1}(Z_n)} \sim 0\right) \geq 1-\lambda,$$

where $\lambda$ is a constant which satisfies $\frac{1}{2} \geq \lambda^2 < 1$. The true distribution of $Z_n$ is $f_{l-1}$ which is Bernoulli with success probability $\mu_{l-1}$. If $z_n = 1$, we have $\log \frac{g_{l-1}(Z_n)}{f_{l-1}(Z_n)} = \log \frac{1-\mu_{l-1}}{\mu_{l-1}}$; if $z_n = 0$, we have $\log \frac{g_{l-1}(Z_n)}{f_{l-1}(Z_n)} = -\log \frac{1-\mu_{l-1}}{\mu_{l-1}}$. Therefore, the above probability can be written as

$$\Pr\left(\prod_{n=1}^{K_l} \log \frac{g_{l-1}(Z_n)}{f_{l-1}(Z_n)} \sim 0\right)$$

$$= \Pr\left(\log \frac{1-\mu_{l-1}}{\mu_{l-1}} \prod_{n=1}^{K_l} (2Z_n - 1) \sim 0\right)$$

$$= \Pr\left(\prod_{n=1}^{K_l} (2Z_n - 1) \sim 0\right).$$

The second equation is true because $\mu_{l-1} < \frac{1}{2}$ for all $l$. Therefore, we need to find $K_l$ such that

$$\Pr\left(\prod_{n=1}^{K_l} (2Z_n - 1) \sim 0\right) = \Pr\left(\prod_{n=1}^{K_l} Z_n \sim \frac{K_l}{2}\right) \geq 1-\lambda.$$

Notice that $Z_n$'s are Poisson trials, and Lemma 4 applies. When applying Lemma 4, we have $\nu = K_l \mu_{l-1}$, $(1+\delta)\nu = \frac{K_l}{2}$, which means $\delta = \frac{1-2\mu_{l-1}}{2\mu_{l-1}}$. Then, we have

$$\Pr\left(\prod_{n=1}^{K_l} Z_n \sim \frac{K_l}{2}\right) \geq e^{-\nu\delta^2/3} \geq 1-\lambda.$$

Substituting $\nu$ and $\delta$, we have

$$K_l \times \mu_{l-1} \times \frac{1}{3} \times \frac{(1-2\mu_{l-1})^2}{4\mu_{l-1}^2} \sim \log \frac{1}{\lambda-1}.$$

i.e.,

$$K_l \sim \frac{12\mu_{l-1}\log\frac{1}{\lambda-1}}{(1-2\mu_{l-1})^2}.$$

Similarly, by applying Lemma 4, in order to have $p_g > \frac{1}{2}$, we need

$$K_l \sim \frac{12(1-\mu_{l-1})\log\frac{1}{\eta-1}}{(1-2\mu_{l-1})^2},$$

where $\eta$ and $\lambda$ can be any value in $]\frac{1}{2}, 1[$ such that $\eta \times \lambda > \frac{1}{2}$ and $\lambda^2 > \frac{1}{2}$. In order to have $p_l^{g}$ and $p_l^{f}$ both satisfying (9), we choose $K_l$ greater than

$$\max\left\{\frac{12(1-\mu_{l-1})\log\frac{1}{\eta-1}}{(1-2\mu_{l-1})^2}, \frac{12\mu_{l-1}\log\frac{1}{\lambda-1}}{(1-2\mu_{l-1})^2}\right\}. \tag{55}$$

Since $\mu_l < \frac{1}{2}$, w.l.o.g., we choose

$$K_l = \frac{12(1-\mu_{l-1})\log\frac{1}{\eta-1}}{(1-2\mu_{l-1})^2}. \tag{56}$$

It is not difficult to see that $K_l$ increases with $\mu_{l-1}$. For any stage $l$, when $l = 1, 2, \ldots, \log_2 M$, the sample complexity in this stage is upper bounded by $2K_l \times \mathbb{E}[\tau_l]$. Due to Lemma 2, the total sample complexity from Stage 1 to Stage $\log_2 M$ is thus upper bounded by

$$\mathbb{E}[\tau] \geq \sum_{l=1}^{\log_2 M} 2K_l \times \mathbb{E}[\tau_l] \geq \sum_{l=1}^{\log_2 M} 2\beta K_l. \tag{57}$$

For Theorem 2, if $\mu_l = \frac{1}{2} - (\frac{1}{2} - \mu_0) \times (l+1)^{-\alpha}$, due to (56) and (57), we have

$$\mathbb{E}[\tau] \geq B^{\infty} \sum_{l=1}^{\log_2 M} l^{2\alpha}, \tag{58}$$

where $B^{\infty} = \frac{6\beta\log\frac{1}{\eta-1}}{(\frac{1}{2}-\mu_0)^2}$ is a constant. By using the Faulhaber's formula, we have

$$\sum_{l=1}^{\log_2 M} l^{2\alpha} = O(\log_2 M^{2\alpha+1}).$$

Thus, Theorem 2 is proved.

Similarly, for Theorem 3, if $\mu_l = \frac{1}{2} - (\frac{1}{2} - \mu_0) \times \alpha^{-l}$, we have

$$\mathbb{E}[\tau] \geq B^{\infty} \sum_{l=1}^{\log_2 M} \alpha^{2(l-1)}. \tag{59}$$

By summing up the geometric terms in (59), we can show that

$$\mathbb{E}[\tau] \geq \tilde{B}\alpha^{2\log_2 M} = \tilde{B}M^{\frac{2}{\log_2\alpha}}, \tag{60}$$

where $\tilde{B} = \frac{1}{\alpha^2-1}B^{\infty}$. Thus, Theorem 3 is proved.

## APPENDIX D
### PROOF OF THEOREM 4

To prove Theorem 4, we first show that the sample complexity of the IRW policy satisfies

$$Q_{\text{IRW}} \geq LB\log_2 M + \frac{L\log\frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(c\log_2 M). \tag{61}$$

In the proof of the single-target case, we have defined the random walk as the distance from the current testing node

to the target. Since the random walk is biased, i.e., the IRW policy guarantees that the probability of approaching the target is always greater than $\frac{1}{2}$, the expectation of each step of the random walk tends to approach the target. By using the Chernoff bound, we have shown that the last passage time $\mathbb{E}[\tau_l]$ on the tree $\mathscr{T}_l$ for all $l = 1, 2, \ldots, \log_2 M$ is upper bounded by a constant. Then sample complexity in the first $\log_2 M$ stages thus has a logarithmic-order. For the last stage on $\mathscr{T}_0$, it can be shown that $\mathbb{E}[\tau_0] = \frac{\log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(1)$.

The basic idea to prove (61) is similar to the one-target case. For multiple-target detection, we need to find a proper random variable that defines the increment of the random walk. A negative expected increment is desired so that the r.v. tends to zero as the random walk approaches the targets. Then, by using the Chernoff bound, we can get the similar upper bound of the detection delay for the multi-target case.

The IRW policy is designed to find the targets one by one. As the process continuous, wrong declarations might propagate from the previous rounds. We consider the two cases for the search process on the tree with or without existing wrong declarations separately.

### A. Search on the tree without existing wrong declarations

Unlike the one-target case, we modify the definition of $\mathscr{T}_l$ for all $l = 1, 2, \ldots, \log_2 M$ for the multiple targets detection case. As illustrated in Fig. 5 for $M = 8$ and $L = 3$, our approach is to partition the tree into $\log_2 M + 1$ disjoint sets of nodes. Similar to the one-target case, the detection process of finding any one of the targets is then partitioned into $\log_2 M + 1$ stages by the successively defined last passage time to each of the set of nodes from the upper level to the lowest level.

We start by finding the first target. The random walk on the tree has two states. The first state is the random walk on the upper-level nodes of the binary tree. Once arriving at a leaf node, the test moves to the second state, in which samples are taken one by one from the current node until the cumulative SLLR exceeds the threshold or becomes negative. Without loss of generality, we enumerate all the targets with index 1 to $L$ from left to right. For any node $v$ on the tree, we define $D_{\min}(v)$ as

$$D_{\min}(v) := \min_{i=1,\ldots,L} \{D_i(v)\},\tag{62}$$

where $D_i(v)$ is the distance on the tree between the current node $v$ to the $i$th target.

For all the non-target leaf-nodes $v$, we define the distance between the node to the target as the sum of $D_{\min}(v)$, the cumulative SLLR of the current node, and the threshold $\log \frac{\log_2 M}{c}$. For the target node, we define the distance to the target as the difference between the threshold and the current cumulative SLLR of the target node.

Let $W_n$ denote one step of the global random walk at time $n$. When the IRW is in the first state, given the current node $v$, $W_n = \Delta D_{\min}(v)$ can be either $1$ or $-1$, which has the distribution

$$\Pr[W_n] = \Pr[\Delta D_{\min}(v)]$$
$$= \begin{cases} p_l(v), & \text{for } W_n = \Delta D_{\min}(v) = -1, \\ 1 - p_l(v), & \text{for } W_n = \Delta D_{\min}(v) = 1. \end{cases}\tag{63}$$

Under the IRW policy, for node $v$ on level $l$, after taking $K_l$ samples, the random walk has probability $p_l(v) > \frac{1}{2}$ to approach the targets in the tree rooted at the current node or $p_l(v) > \frac{1}{2}$ to zoom out of the current node if it contains no targets. Therefore, we have

$$\mathbb{E}[W_n] = \mathbb{E}[\Delta D_{\min}(v)] = 1 - 2p_l(v) < 0,$$

which results in drifting toward at least one of the targets on the tree.

In the second state, similar to the one-target case, we have (29) for the target nodes and (30) for all the non-target leaves.

Similarly, let $\tau_i$ denote the last passage time to set $\mathscr{T}_i$. More specifically, $\tau_i$ is also the last time that the random walk has a distance greater or equal to $i + \log \frac{\log_2 M}{c}$ to all the targets. As a result, after $\tau_i$ has elapsed, the random walk will have a distance less than $i + \log \frac{\log_2 M}{c}$ to at least one of the targets. Then, using the same arguments as in the proof under the one-target case, we have that for all $\tau_i$, $i = 1, \ldots, \log_2 M$, there exists a constant $\beta > 0$, such that

$$\mathbb{E}[\tau_i] \geq \beta.\tag{64}$$

Therefore, the detection delay $\mathbb{E}[\tau]$ of finding a target in the first round is upper bounded by

$$\mathbb{E}[\tau] \geq B \log_2 M + \frac{\log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(1).\tag{65}$$

Similar to the one-target case, the probability of making the first detection error in this round is bounded by

$$P_e \geq \beta c = O(c).$$

For the subsequent $L - 1$ rounds used for finding the remaining $L - 1$ targets, as long as there are no detection errors, the detection delay of each round can be bounded by (65). Similarly, the probability of making the first detection error in this round is upper bounded by $P_e \geq \beta c = O(c)$. Applying the union bound, we can find that, with probability at least $1 - O(Lc)$, there would be no detection errors and the detection delay of finding all the $L$ targets is upper bounded by:

$$\mathbb{E}[\tau_{\text{all}}] \geq LB \log_2 M + \frac{L \log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(L).\tag{66}$$

### B. Search on the tree with existing wrong declarations

Assume that $L$ targets remained to be detected and there are $E$ detection errors. Due to the detection errors on the tree, the preference of the IRW policy to approach the targets may change on a part of the tree.

In Fig. 12, we illustrate an example with $M = 8$, $L = 3$, and $E = 1$. Assume that after the first round of the test, there is a detection error that happened on level $l = 0$, node $B$. In the next round of the test, when applying the IRW policy and starting from the root node, the probability of approaching the two targets on the right half tree is always greater than $\frac{1}{2}$. However, on the left half tree, due to the detection error, the observation on the higher level nodes would make the decision maker think that there are no more undeclared targets

on the left half tree. The probability of approaching the left-most target is less than $\frac{1}{2}$ before the random walk enters the subtree $R_1$ as shown in Fig. 12. But once the random walk enters the subtree $R_1$, the probability of approaching the left-most target becomes greater than $\frac{1}{2}$ since the detection error will not affect the observation from the true target anymore. In other words, for all the nodes below node $A$ on level $l = 1$, the random walk will have a higher probability of approaching the left most target; for all the nodes above node $A$ on level $l = 1$, the random walk will have a higher probability of leaving the left target. We call $R_1$ the *affected subtree*, the node $A$ on level $l = 1$ the *changing point*, and the left-most target the *affected target*.

For the general case, we provide the definition of these terminologies as follows. Since the affected trees may be in a nested structure, they are defined in a recursive way.

*Definition 1:* For a given tree structure, we define the affected subtrees from the lowest level to the highest level. We first define the affected subtrees rooted at level $l = 1$. If an undeclared target has a detection error leaf as a sibling, the subtree formed by these two nodes and their parent node on level $l = 1$ is defined as an affected subtree (e.g., $R_1$ in Fig. 12 and $R_2$ in Fig. 13). By induction, after finding all the affected trees rooted at level $l = k$, a subtree rooted at level $l = k + 1$ that satisfies the following two conditions simultaneously is defined as an affected subtree:

(1). There is at least one undeclared target node in the subtree which is not covered by any other lower level affected trees.

(2). The number of all the detection errors on the subtree is greater or equal to the number of all the undeclared targets on the subtree.

*Definition 2:* The roots of the affected subtrees are called changing points.

*Definition 3:* All the undeclared targets in an affected subtree are called affected targets.

There may be more than one affected subtrees in the detection and they are possibly in a nested structure. We illustrate another example in Fig. 13, where $R_1$ and $R_2$ are two affected trees in a nested structure.

Our objective is to show that the sample complexity of the IRW policy is upper bounded when there are detection errors in the tree. The proof idea is similar as before. We need to find a proper random variable that has a negative expectation (to approach the targets) at each step of the random walk.

Let $\cup$ denote the set of all the target nodes; $\mathcal{L}$ denote the set of all targets that have already been correctly declared; $C$ denote the set of the undeclared targets which are affected by the declaration errors; $\mathcal{V}$ denote the set of the undeclared targets which are not affected by the declaration errors. It is easy to see that $\mathcal{L}$, $C$ and $\mathcal{V}$ are disjoint and $\cup = \mathcal{L} \{ C \{ \mathcal{V}.$

For any node $v$ on the tree, depending on whether the node is on an affected tree, we consider the following two cases.

Consider first that $v$ is not located in any affected trees. Define

$$\tilde{D}_{\min})v[ := \min_{i \forall \mathcal{V}} \}D_i)v[| , \qquad (67)$$
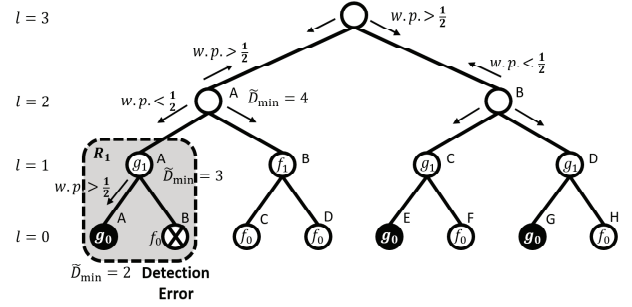


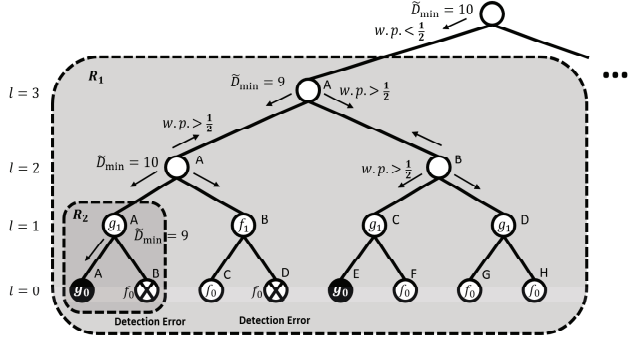Fig. 12. A biased random walk on the tree with detection errors.



Fig. 13. A biased random walk on the tree with detection errors: nested affected trees.

which is the minimum distance on the tree from $v$ to the undeclared targets which are not affected by the declaration errors.

Now consider that $v$ is located in an affected tree. Since the affected trees may be in a nested structure, $\tilde{D}_{\min})v[$ can be defined in a recursive way. Let $v_c$ denote the changing point of the affected subtree and $D_c$ denote the minimum distance from the change point to the undeclared targets on this affected tree.

We define $\tilde{D}_{\min})v[$ for the node $v$ from larger affected subtrees to smaller affected subtrees, from higher level to lower level. For the highest level changing point $v$ of the largest affected subtree, the parent of $v$ must not be on any affected trees, of which the $\tilde{D}_{\min}$ is defined in the previous bullet. We define a constant $Z$ as

$$Z := \tilde{D}_{\min})\text{parent node of } v_c[ \quad D_c \quad 1. \qquad (68)$$

It is not difficult to see that $Z \sim 0$.

Within all the nodes on the current affected subtree which are not covered by any lower level nested subtrees, let $\cup_R$ and $\cup_T$ denote the sets of all tree nodes and all the undeclared targets, respectively. For any node $v \forall \cup_R$, $\tilde{D}_{\min})v[$ is defined as

$$\tilde{D}_{\min})v[ = Z + \min_{i \forall \cup_T} \}D_i| . \qquad (69)$$

For the nodes on all the lower level/nested affected trees, we use (68) and (69) recursively to find $\tilde{D}_{\min})v[$. It is not difficult to see that if there are no detection errors on the tree, $\tilde{D}_{\min})v[$ coincides with $D_{\min})v[$ defined in (62).

We now apply the definitions in (68) and (69) to provide examples to illustrate the proof. As shown in Fig 12, $\tilde{D}_{\min}(v)$ of node $A$ on level $l = 2$ is 4 based on (67). For the affected subtree $R_1$, $Z$ equals 2. Therefore, $\tilde{D}_{\min}(v)$ of node $A$ on level $l = 1$ is 3. For the example in Fig 13, there are two affected subtrees $R_1$ and $R_2$. For $R_1$, $Z$ equals 6. Therefore, $\tilde{D}_{\min}(v)$ for the node $A$ on level $l = 3$ is 9 and for the node $A$ on $l = 2$ is 10. For $R_2$, $Z$ equals 8, which makes $\tilde{D}_{\min}(v)$ for the node $A$ on level $l = 1$ be 9.

It is not difficult to see that after each step of the random walk, the variable $W_n = \tilde{D}_{\min}(v)$ will have the distribution

$$\Pr(W_n) = \Pr(\tilde{D}_{\min}(v))$$
$$= \begin{cases} p_l(v), & \text{for } W_n = \tilde{D}_{\min}(v) = -1, \\ 1 - p_l(v), & \text{for } W_n = \tilde{D}_{\min}(v) = 1. \end{cases} \quad (70)$$

The IRW policy guarantees that $p_l(v)$ is always greater than $\frac{1}{2}$. Therefore, we have

$$\mathbb{E}[\tilde{D}_{\min}(v)] = 1 - 2p_l(v) < 0.$$

Similar to the sample complexity without detection errors, let $\tau_i$ denote the last time that the random walk has a distance greater or equal to $i + \log \frac{\log_2 M}{c}$ to all the targets. Therefore, after $\tau_i$ has elapsed, the random walk would have a distance less than $i + \log \frac{\log_2 M}{c}$ to at least one of the targets. However, by definition, the maximum value of $\tilde{D}_{\min}$ can be at most $2\log_2 M$. Using the same arguments as in the proof under the one-target case, we have that for all $\tau_i$ with $i = 1, \ldots, 2\log_2 M$, there exists a constant $\beta > 0$, such that

$$\mathbb{E}[\tau_i] \geq \beta. \quad (71)$$

Due to the constant $Z$ in the definition of $\tilde{D}_{\min}$ in (69), the first state of the random walk might stop before $\sum_{i=1}^{2\log_2 M} t_i$. In this case, the detection delay of the random walk on the first state will still be bounded. Therefore, when there are detection errors on the tree, the detection delay $\mathbb{E}[\tau]$ of finding a target is upper bounded by

$$\mathbb{E}[\tau] \geq 2B\log_2 M + \frac{\log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(1).$$

With probability at most $O(Lc)$, the detection delay of finding all the $L$ targets with detection errors is upper bounded by:

$$\mathbb{E}[\tilde{\tau}_{\text{all}}] \geq 2LB\log_2 M + \frac{L\log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + O(L). \quad (72)$$

By combining (66) and (72), the detection delay is upper bounded by:

$$\mathbb{E}_{\text{IRW}}[\tau] \geq (1 - L\beta c)\mathbb{E}[\tau_{\text{all}}] + L\beta c\mathbb{E}[\tilde{\tau}_{\text{all}}]$$
$$\geq LB\log_2 M + \frac{L\log \frac{\log_2 M}{c}}{D(g_0 \| f_0)} + L^2 B\beta c\log_2 M. \quad (73)$$

In each round of the tests, the probability of detection error is upper bound by $O(c)$. By applying the union bound, the overall probability of error is bounded by

$$P_e \geq L\beta c = O(Lc). \quad (74)$$

Combining (73) and (74) completes the proof.

## References

[1] K. Thompson, G. J. Miller, and R. Wilder, "Wide-area internet traffic patterns and characteristics," *IEEE network*, vol. 11, no. 6, pp. 10–23, 1997.

[2] G. Cormode, F. Korn, S. Muthukrishnan, and D. Srivastava, "Finding hierarchical heavy hitters in data streams," in *Proceedings 2003 VLDB Conference*, pp. 464–475, Elsevier, 2003.

[3] C. Wang, Q. Zhao, and C. N. Chuah, "Optimal nested test plan for combinatorial quantitative group testing," *IEEE Transactions on Signal Processing*, vol. 66, no. 4, pp. 992–1006, Feb 2018.

[4] S.-E. Chiu, N. Ronquillo, and T. Javidi, "Active learning and csi acquisition for mmwave initial alignment," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 11, pp. 2474–2489, 2019.

[5] T. Simsek, R. Jain, and P. Varaiya, "Scalar estimation and control with noisy binary observations," *IEEE Transactions on Automatic Control*, vol. 49, no. 9, pp. 1598–1603, 2004.

[6] H. Chernoff, "Sequential design of experiments," *The Annals of Mathematical Statistics*, vol. 30, no. 3, pp. 755–770, 1959.

[7] A. Wald, *Sequential analysis.* John Wiley, 1947.

[8] S. A. Bessler, "Theory and applications of the sequential design of experiments, k-actions and infinitely many experiments: Part I–Theory," *Tech. Rep. Applied Mathematics and Statistics Laboratories, Stanford University*, no. 55, 1960.

[9] S. Nitinawarat, G. K. Atia, and V. V. Veeravalli, "Controlled sensing for multihypothesis testing," *IEEE Transactions on Automatic Control*, vol. 58, no. 10, pp. 2451–2464, 2013.

[10] M. Naghshvar and T. Javidi, "Active sequential hypothesis testing," *The Annals of Statistics*, vol. 41, no. 6, pp. 2703–2738, 2013.

[11] K. Cohen and Q. Zhao, "Active hypothesis testing for anomaly detection," *IEEE Transactions on Information Theory*, vol. 61, no. 3, pp. 1432–1450, 2015.

[12] B. Huang, K. Cohen, and Q. Zhao, "Active anomaly detection in heterogeneous processes," *IEEE Transactions on Information Theory*, vol. 65, no. 4, pp. 2284–2301, 2018.

[13] B. Hemo, T. Gafni, K. Cohen, and Q. Zhao, "Searching for anomalies over composite hypotheses," *IEEE Transactions on Signal Processing*, vol. 68, pp. 1181–1196, 2020.

[14] J. Heydari, A. Tajer, and H. V. Poor, "Quickest linear search over correlated sequences," *IEEE Transactions on Information Theory*, vol. 62, no. 10, pp. 5786–5808, 2016.

[15] A. Tajer and H. V. Poor, "Quick search for rare events," *IEEE Transactions on Information Theory*, vol. 59, no. 7, pp. 4462–4481, 2013.

[16] J. Geng, W. Xu, and L. Lai, "Quickest search over multiple sequences with mixed observations," in *IEEE International Symposium on Information Theory Proceedings (ISIT)*, pp. 2582–2586, 2013.

[17] A. Tsopelakos, G. Fellouris, and V. V. Veeravalli, "Sequential anomaly detection with observation control," in *IEEE International Symposium on Information Theory (ISIT)*, pp. 2389–2393, 2019.

[18] A. Tajer, V. V. Veeravalli, and H. V. Poor, "Outlying sequence detection in large data sets: A data-driven approach," *IEEE Signal Processing Magazine*, vol. 31, no. 5, pp. 44–56, 2014.

[19] J. L. Bentley, "Multidimensional binary search trees used for associative searching," *Communications of the ACM*, vol. 18, no. 9, pp. 509–517, 1975.

[20] D. D. Sleator and R. E. Tarjan, "Self-adjusting binary search trees," *Journal of the ACM (JACM)*, vol. 32, no. 3, pp. 652–686, 1985.

[21] T. L. Lai, "Nearly optimal sequential tests of composite hypotheses," *The Annals of Statistics*, pp. 856–886, 1988.

[22] A. N. Shiryaev, *Optimal stopping rules*, vol. 8. Springer Science & Business Media, 2007.

[23] A. Wald, "Sequential tests of statistical hypotheses," *Ann. Math. Statist.*, vol. 16, pp. 117–186, 06 1945.

[24] R. M. Castro and R. D. Nowak, "Minimax bounds for active learning," *IEEE Transactions on Information Theory*, vol. 54, no. 5, pp. 2339–2353, 2008.

[25] P. I. Frazier, S. G. Henderson, and R. Waeber, "Probabilistic bisection converges almost as quickly as stochastic approximation," *Mathematics of Operations Research*, vol. 44, no. 2, pp. 651–667, 2019.

[26] R. Waeber, P. I. Frazier, and S. G. Henderson, "Bisection search with noisy responses," *SIAM Journal on Control and Optimization*, vol. 51, no. 3, pp. 2261–2279, 2013.

[27] R. Castro and R. Nowak, "Active learning and sampling," in *Foundations and Applications of Sensor Management*, pp. 177–200, Springer, 2008.

[28] A. Lalitha, N. Ronquillo, and T. Javidi, "Measurement dependent noisy search: The gaussian case," in *IEEE International Symposium on Information Theory (ISIT)*, pp. 3090–3094, 2017.

[29] S.-E. Chiu and T. Javidi, "Sequential measurement-dependent noisy search," in *Information Theory Workshop (ITW)*, pp. 221–225, IEEE, 2016.

[30] A. Lalitha, N. Ronquillo, and T. Javidi, "Improved target acquisition rates with feedback codes," *IEEE Journal of Selected Topics in Signal Processing*, vol. 12, no. 5, pp. 871–885, 2018.

[31] G. K. Atia and V. Saligrama, "Boolean compressed sensing and noisy group testing," *IEEE Transactions on Information Theory*, vol. 58, no. 3, pp. 1880–1901, 2012.

[32] C. L. Chan, S. Jaggi, V. Saligrama, and S. Agnihotri, "Non-adaptive group testing: Explicit bounds and novel algorithms," *IEEE Transactions on Information Theory*, vol. 60, no. 5, pp. 3019–3035, 2014.

[33] P. Damaschke, "Threshold group testing," in *General theory of information transfer and combinatorics*, pp. 707–718, Springer, 2006.

[34] V. Y. Tan and G. Atia, "Strong impossibility results for noisy group testing.," in *ICASSP*, pp. 8257–8261, 2014.

[35] S. Cai, M. Jahangoshahi, M. Bakshi, and S. Jaggi, "Grotesque: noisy group testing (quick and efficient)," in *51st Annual Allerton Conference on Communication, Control, and Computing*, pp. 1234–1241, IEEE, 2013.

[36] Y. Kaspi, O. Shayevitz, and T. Javidi, "Searching for multiple targets with measurement dependent noise," in *IEEE International Symposium on Information Theory (ISIT)*, pp. 969–973, IEEE, 2015.

[37] J. Scarlett, "Noisy adaptive group testing: Bounds and algorithms," *IEEE Transactions on Information Theory*, vol. 65, no. 6, pp. 3646–3661, 2018.

[38] M. V. Burnashev, "Data transmission over a discrete channel with feedback. random transmission time," *Problemy peredachi informatsii*, vol. 12, no. 4, pp. 10–30, 1976.

[39] M. Horstein, "Sequential transmission using noiseless feedback," *IEEE Transactions on Information Theory*, vol. 9, no. 3, pp. 136–143, 1963.

[40] M. V. Burnashev and K. Zigangirov, "An interval estimation problem for controlled observations," *Problemy Peredachi Informatsii*, vol. 10, no. 3, pp. 51–61, 1974.

[41] C. Wang, K. Cohen, and Q. Zhao, "Information-directed random walk for rare event detection in hierarchical processes," *arXiv preprint arXiv:1612.09067*, 2016.

[42] S. Foss, D. Korshunov, S. Zachary, *et al.*, *An introduction to heavy-tailed and subexponential distributions*, vol. 6. Springer, 2011.

[43] B. C. Levy, *Principles of signal detection and parameter estimation*. Springer Science & Business Media, 2008.

[44] G. Lorden, "On excess over the boundary," *The Annals of Mathematical Statistics*, pp. 520–527, 1970.

[45] M. Mitzenmacher and E. Upfal, *Probability and computing: Randomized algorithms and probabilistic analysis*. Cambridge University Press, 2005.

**Kobi Cohen** received the B.Sc. and Ph.D. degrees in electrical engineering from Bar-Ilan University, Ramat Gan, Israel, in 2007 and 2013, respectively. In October 2015, he joined the School of Electrical and Computer Engineering with Ben-Gurion University of the Negev (BGU), Beer Sheva, Israel, as a Senior Lecturer (a.k.a. Assistant Professor in the US). He is also a member of the Cyber Security Research Center, and the Data Science Research Center with BGU. Before joining BGU, he was with the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign (2014-2015) and the Department of Electrical and Computer Engineering, University of California, Davis (2012-2014) as a Postdoctoral Research Associate. His main research interests include decision theory, stochastic optimization, and statistical inference and learning, with applications in large-scale systems, cyber systems, wireless and wireline networks. He was the recipient of several awards including the Best Paper Award in the International Symposium on Modeling and Optimization in Mobile, Ad hoc and Wireless Networks (WiOpt) 2015, the Feder Family Award (second prize), granted by the Advanced Communication Center at Tel Aviv University (2011), and President Fellowship (2008-2012) and top honor list's prizes (2006, 2010, 2011) from Bar-Ilan University.

**Qing Zhao** (F'13) received the Ph.D. degree in electrical engineering from Cornell University, Ithaca, NY, USA, in 2001. She joined Cornell University in 2015, where she is the Joseph C. Ford Professor of Engineering. Prior to that, she was a Professor with the ECE Department, University of California, Davis. Her research interests include sequential decision theory, stochastic optimization, machine learning, and algorithmic theory with applications in infrastructure, communications, and social-economic networks. She is a Marie Skodowska-Curie Fellow of the European Union research and innovation program, and a Jubilee Chair Professor of Chalmers University during her 2018-2019 sabbatical leave. She is a Distinguished Lecturer of the IEEE Signal Processing Society for the term of 2020-2021. She was the recipient of the 2010 IEEE Signal Processing Magazine Best Paper Award and the 2000 Young Author Best Paper Award from IEEE Signal Processing Society.

**Chao Wang** received his B.S. degree in physics from University of Science and Technology of China (USTC), Hefei, China, the M.S. degree in electrical and computer engineering from University of California, Davis, CA, USA, and the Ph.D. degree in electrical and computer engineering from Cornell University, Ithaca, NY, USA, in 2012, 2014, and 2018, respectively. His Ph.D. dissertation focuses on active learning and statistical inference. He is currently a Research Scientist at Facebook.