Sequential Outlier Hypothesis Testing under Universality Constraints

Jun Diao and Lin Zhou

School of Cyber Science and Technology (CST), Beihang University, China Emails: {jundiao,lzhou}@buaa.edu.cn

Abstract—We revisit sequential outlier hypothesis testing and derive bounds on the achievable exponents. Specifically, the task of outlier hypothesis testing is to identify the set of outliers that are generated from an anomalous distribution among all observed sequences where most are generated from a nominal distribution. In the sequential setting, one obtains a sample from each sequence per unit time until a reliable decision could be made. We assume that the number of outliers is known while both the nominal and anomalous distributions are unknown. For the case of exactly one outlier, our bounds on the achievable exponents are tight, providing exact large deviations characterization of sequential tests and strengthening a previous result of Li, Nitinawarat and Veeravalli (2017). In particular, we propose a sequential test that has bounded average sample size and better theoretical performance than the fixed-length test, which could not be guaranteed by the corresponding sequential test of Li, Nitinawarat and Veeravalli (2017). Our results are also generalized to the case of multiple outliers.

Index Terms—Error Exponent, Large Deviations, Hypothesis testing, Anomaly Detection

I. INTRODUCTION

Outlier hypothesis testing is a popular statistical inference problem [1]–[4], where one is asked to identify a set of outliers among a given number M of observed sequences. The majority of sequences are generated i.i.d. from a nominal distribution and the rest are generated i.i.d. from an anomalous distribution different from the nominal distribution. Both the nominal distribution and anomalous distributions are *unknown*.

The number of outliers can be assumed either known or unknown. When the number of outliers is known, the task is relatively simpler and corresponds to a generalization of classification [5], [6]. When the number of outliers is unknown, one could estimate the number of outliers and subsequently identify the set of outliers using a test for the known number case. As a compromise, one could also consider the case of at most T outliers, where an upper bound T on the number of outliers is known. When T = 1, the case is termed at most one [1], [3]. In this paper, for simplicity, we consider the case of known number of outliers and our results could be generalized to the case of unknown number of outliers by having an additional step to estimate the number of outliers.

Depending on the test design, a test could be fixed-length or sequential. When the sample size of each observed sequence is fixed, the corresponding test is a fixed-length test. When the sample size is a random variable depending on particular observations of sequences, the corresponding test is a sequential test. In a sequential test, one obtains a new sample from each sequence per unit time until one is confident to make a decision. The expected value of the sample size is also known as the expected stopping time. Since the generating distributions of sequences are unknown, for sequential tests, naturally, one can put a universal constraint either on the error probability or the average stopping time [7, Def. 2 and 3]. Specifically, for any pair of nominal and anomalous distributions, the error probability universality constraint requires the test to have the error probability bounded by a tolerable value $\beta \in (0, 1)$ under each hypothesis while the expected stopping time universality constraint requires that the expected stopping time under each hypothesis is bounded. Correspondingly, for fixed-length tests, only error probability universality constraint is valid since the sample size is fixed a-priori.

For both cases of at most one outlier and of at most T outliers, Li, Nitinawarat and Veeravalli [1] proposed generalized likelihood (GL) tests and proved the optimality of the test by having largest exponential decay rates of error probabilities when the number M of observed sequences tends to infinity. Subsequently, Li, Nitinawarat and Veeravalli [2] generalized the above results to the sequential setting under the error probability universality constraint. However, there are several limitations for the results in [2]. Firstly, only achievability results under the error probability universality constraint were derived. Without a matching converse result, the optimality of error exponents could not be guaranteed. Furthermore, the expected stopping time constraint was not considered, which leads to the undesired fact that the sequential tests might stop at very large sample sizes. Finally, it was only numerically shown that the sequential test outperforms the fixed-length test only when the average stopping time is relatively large [2, Figs. 1 and 2]. Without a theoretical guarantee, the benefit of the sequential design is not fully uncovered.

In this paper, for a slightly easier setting of exactly one outlier, we address all above limitations. Furthermore, we generalize our results to the case of multiple outliers when the number of outliers is known. Our main contribution is summarized in the following subsection.

A. Main Contributions

For the case of exactly one outlier, we refine the result in [2, Theorem 3.2] by deriving a matching converse result and re-

proving a simpler achievability part under the error probability universality constraint. Furthermore, we derive the exact error exponents under the expected stopping time constraint. In particular, in the achievability part, we propose a sequential test that has bounded average sample size under any pair of nominal and anomalous distributions and analytically show that the test could have strictly better performance than the fixed-length test in [1]. To compare the performance of the sequential tests under both universality constraints, we provide numerical examples to illustrate the achievable exponents and the expected stopping time of both tests, which imply that our proposed test under the expected stopping time universality constraint has much smaller average sample size and even better performance in certain scenarios. We also generalize our results to the case of multiple outliers and derive bounds on achievable exponents.

B. Other Related Works

We briefly recall other (non-exhausting) related works on outlier hypothesis testing. Bu, Zou and Veeravalli [4] proposed a low-complexity test for outlier hypothesis testing and showed that the test ensures exponential decay of error probabilities. Zhou, Wei and Hero [3] proposed an optimal threshold-based test under the generalized Neyman-Pearson criterion [5] and derived a second-order asymptotic approximation to the finite sample size performance. Zou *et al.* [8] used the maximum mean discrepancy metric to design a test for outlier hypothesis testing of continuous sequences and showed that the test is exponentially consistent

II. PROBLEM FORMULATION AND EXISTING RESULTS

Notation

We use \mathbb{R} , \mathbb{R}_+ , \mathbb{N} to denote the set of real numbers, nonnegative real numbers, and natural numbers respectively. Given any two integers $(a, b) \in \mathbb{N}^2$, we use [a : b] to denote the set of integers $\{a, a + 1, \ldots, b\}$ and use [a] to denote [1 : a]. Random variables and their realizations are denoted by upper case variables (e.g., X) and lower case variables (e.g., x), respectively. All sets are denoted in calligraphic font (e.g., \mathcal{X}). Given any integer $N \in \mathbb{N}$, let $X^N := (X_1, \ldots, X_N)$ be a random vector of length N and let $x^N = (x_1, \ldots, x_N)$ be a particular realization of X^N . The set of all probability distributions on a finite set \mathcal{X} is denoted as $\mathcal{P}(\mathcal{X})$.

A. Problem Formulation

Consider a set of M observed sequences $\mathbf{X}^{\tau} := \{X_1^{\tau}, \ldots, X_M^{\tau}\}$, where τ is a random stopping time with respect to the filtration $\{\mathcal{F}_n\}_{n\in\mathbb{N}}$ and \mathcal{F}_n is generated by σ -algebra $\sigma\{X_1, X_2, \ldots, X_n\}$. Most sequences are generated i.i.d. from an unknown nominal distribution P_N while the rest of few sequences known as outliers are generated i.i.d. from an unknown anomalous distribution P_A . We first consider the case of exactly one outlier and then generalize the results to multiple outliers with the number of outliers known.

When there is exactly one outlier, the task is to design a test $\Phi = \{\tau, \phi_{\tau}\} : \mathcal{X}^{M\tau} \to \{H_1, H_2, \dots, H_M\}$ that consists

of a random stopping time τ and a decision rule ϕ_{τ} to classify among the following M hypotheses:

• $H_i, i \in [M]$: the *i*-th sequence is the outlier.

To evaluate the performance of a test, we consider the misclassification error probability and the expected stopping time of a sequential test. Specifically, for each $i \in [M]$, the misclassification error probability is defined as follows:

$$\beta_i(\Phi|P_{\mathcal{A}}, P_{\mathcal{N}}) := \mathbb{P}_i\{\Phi(\mathbf{X}^{\tau}) \neq \mathcal{H}_i\}, i \in [M], \qquad (1)$$

where we define $\mathbb{P}_i(\cdot) := \Pr\{\cdot | H_i\}$ to denote the joint distribution of observed sequences \mathbf{X}^{τ} , where X_i^{τ} is generated i.i.d. from the anomalous distribution P_A and for each $j \in \mathcal{M}_i := \{j \in [M] : j \neq i\}, X_j^{\tau}$ is generated i.i.d. from the nominal distribution P_N . Furthermore, the expected stopping time under hypothesis H_i satisfies

$$\mathbb{E}_i[\tau] = \sum_{k=1}^{\infty} \mathbb{P}_i\{\tau > k\}.$$
(2)

Since there are two performance criteria, one could put a universal constraint on either one. Motivated by the analyses for sequential binary classification [7, Def. 2 and 3], we define the following two universality constraints on sequential tests.

Definition 1. (Universality Constraint on the Error Probability): Given $\beta \in (0, 1)$ and a sequential test Φ , we say that Φ satisfies the universality constraint on the error probability with β if for any pair of distributions $(P_N, P_A) \in \mathcal{P}(\mathcal{X})^2$,

$$\max_{i \in [M]} \beta_i(\Phi | P_{\mathcal{A}}, P_{\mathcal{N}}) \le \beta.$$
(3)

For a sequential test satisfying the error probability universality constraint, we are interested in the following error exponent for each $i \in [M]$:

$$E_i(\Phi|P_{\rm A}, P_{\rm N}) := \liminf_{n \to \infty} \frac{-\log \beta}{\mathbb{E}_i[\tau]}.$$
 (4)

Definition 2. (Universality Constraint on the Expected Stopping Time): Given $n \in \mathbb{N}$ and a sequential test Φ , we say that Φ satisfies the universality constraint on the expected stopping time with n if for any pair of distributions $(P_N, P_A) \in \mathcal{P}(\mathcal{X})^2$,

$$\max_{i \in [M]} \mathbb{E}_i[\tau] \le n.$$
(5)

For a sequential test satisfying the expected stopping time universality constraint, we are interested in the following error exponent for each $i \in [M]$:

$$E_i(\Phi|P_{\mathcal{A}}, P_{\mathcal{N}}) := \liminf_{n \to \infty} \frac{-\log \beta_i(\Phi|P_{\mathcal{A}}, P_{\mathcal{N}})}{n}.$$
 (6)

B. Existing Results

To compare the performance of sequential tests and fixedlength test, we first recall the results of the fixed-length test Φ_{Li} by Li, Nitinawarat and Veeravalli [1]. To present the test, we need the following definition [3, Eq. (4)]. Given a tuple of distributions $\mathbf{Q} = (Q_1, \ldots, Q_M) \in \mathcal{P}(\mathcal{X})^M$, for each $i \in [M]$, define the following linear combination of KL divergence terms between each single distribution and a mixture distribution:

$$G_i(\mathbf{Q}) := \sum_{j \in \mathcal{M}_i} D\left(Q_j \left\| \frac{\sum_{l \in \mathcal{M}_i} Q_l}{M - 1} \right),\tag{7}$$

where $\mathcal{M}_i = \{j \in [M] : j \neq i\}$. Note that $G_i(\mathbf{Q})$ is used to measure the similarity of distributions \mathbf{Q} except Q_i . The measure $G_i(\mathbf{Q}) = 0$ if and only if $Q_j = Q$ for all $j \in \mathcal{M}_i$ for an arbitrary $Q \in \mathcal{P}(\mathcal{X})$.

For the case of exactly one outlier, the test in [1, Eq. (15)] applies the following minimal scoring function decision rule:

$$\Phi_{\mathrm{Li}}(\mathbf{x}^n) = \mathrm{H}_j, \text{ if } j = \operatorname*{arg\,min}_{i \in [M]} \mathrm{S}_i(\mathbf{x}^n), \tag{8}$$

where $S_i(\mathbf{x}^n) = G_i(\hat{T}_{x_1^n}, \dots, \hat{T}_{x_M^n}).$

Li, Nitinawarat and Veeravalli derived the following result [1, Theorem 2].

Theorem 1. Given any pair of distributions $(P_N, P_A) \in \mathcal{P}(\mathcal{X})^2$, the achievable error exponent of the fixed-length test satisfies that for each $i \in [M]$,

$$E_{i}(\Phi_{\mathrm{Li}}|P_{\mathrm{A}}, P_{\mathrm{N}}) = \min_{\substack{\mathbf{Q}\in\mathcal{P}(\mathcal{X})^{M}:\\G_{1}(\mathbf{Q})\geq G_{2}(\mathbf{Q})}} D(Q_{1}||P_{\mathrm{A}}) + \sum_{j\in[2,M]} D(Q_{j}||P_{\mathrm{N}}).$$
(9)

III. MAIN RESULTS FOR EXACTLY ONE OUTLIER

We characterize the optimal error exponent of sequential outlier hypothesis testing under both the expected stopping time universality constraint and the error probability universality constraint. Specifically, for each case, we propose a corresponding sequential test using the empirical distributions of observed sequences and derive exact large deviations for error probabilities.

A. Error Probability Universality

1) Test Design and Intuition: Given any $\beta \in (0,1)$ and $k \in \mathbb{N}$, define the set

$$\Psi_k(\mathbf{x}^k) := \left\{ l \in [M] : \ \mathcal{S}_l(\mathbf{x}^k) > g(\beta, k) \right\},$$
(10)

where the scoring function $S_l(\mathbf{x}^k) = G_l(\hat{T}_{x_1^k}, \dots, \hat{T}_{x_M^k})$ and the threshold satisfies

$$g(\beta, k) := \frac{-\log(\beta(|\mathcal{X}| - 1))}{k} + \frac{(M+1)|\mathcal{X}|\log(k+1)}{k}.$$
(11)

Under the error probability universality constraint, our sequential test $\Phi_{\rm Ep} = (\tau, \phi_{\tau})$ consists of a random stopping time and the decision rule. The stopping time τ satisfies

$$\tau := \inf \left\{ k \in \mathbb{N} : |\Psi_k(\mathbf{x}^k)| \ge M - 1 \right\}, \tag{12}$$

Note that for each $l \in [M]$, $S_l(\mathbf{x}^k)$ measures the closeness of types of all sequences except the *l*-th sequence. Thus, sequential test Φ_{Ep} stops if the types of nominal samples and the outlier are far away for all M - 1 possibilities of such mixtures. The threshold $g(\beta, k)$ determines how "far away" is measured, which increases with β and decreases with k.

At stopping time τ , our test uses the following decision rule:

$$\phi_{\tau}(\mathbf{x}^{\tau}) = \mathbf{H}_{i}, \text{ if } i = [M] \backslash \Psi_{\tau}(\mathbf{x}^{\tau}).$$
(13)

The above test generalizes the test for sequential classification in [9, Eq. (24)].

We now explain the intuitive reason why the above test works using the weak law of large numbers. Under hypothesis H_i , for each $j \in \mathcal{M}_i$, as the sample size k increases, the empirical distribution $\hat{T}_{x_j^k}$ of a nominal sequence x_j^k tends to the nominal distribution P_N while the empirical distribution $\hat{T}_{x_i^n}$ of the outlier tends to P_A . Thus, as k increases, the scoring function $S_i(\mathbf{x}^k)$ tends to zero and scoring functions $S_j(\mathbf{x}^k)$ for each $j \in \mathcal{M}_i$ tend to a positive real number. When k is sufficiently large, it follows from the weak law of large numbers that there exists M-1 scoring functions with positive values greater than the vanishing value of $g(\beta, k)$. Therefore, a correct decision could always be made asymptotically.

2) Main Results and Discussion: We need the following definition to present our results. Given any two distributions $(P,Q) \in \mathcal{P}(\mathcal{X})^2$ and any positive real number $\alpha \in \mathbb{R}_+$, the generalized Jensen-Shannon divergence [6, Eq. (2.3)] is defined as

$$GJS(P,Q,\alpha) = \alpha D\left(P \left\|\frac{\alpha P + Q}{1 + \alpha}\right) + D\left(Q \left\|\frac{\alpha P + Q}{1 + \alpha}\right)\right).$$
(14)

Theorem 2. For any pair of distributions $(P_N, P_A) \in \mathcal{P}(\mathcal{X})^2$, our sequential test satisfies the error probability universality constraint with $\beta \in (0, 1)$ and the error exponent of our test satisfies that for each $i \in [M]$,

$$E_i(\Phi_{\rm Ep}|P_{\rm A}, P_{\rm N}) \ge {\rm GJS}(P_{\rm N}, P_{\rm A}, M-2).$$
(15)

Conversely, for any sequential test Φ_{β} satisfying the error probability universality constraint with $\beta \in (0, 1)$, under any pair of distributions $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^2$, the error exponent satisfies that for each $i \in [M]$,

$$E_i(\Phi_\beta | P_A, P_N) \le \text{GJS}(P_N, P_A, M - 2).$$
(16)

The proof of Theorem 3 is provided in Appendix A and B, which is inspired by the proof of [7, Theorem 2] for sequential binary classification under the error probability universality constraint. In the achievability proof, we show that for our sequential test, the random variable $\frac{\tau}{\log\beta}$ is uniformly integrable and subsequently we obtain the desired exponent by analyzing the convergence properties of $\frac{-\log\beta}{\mathbb{E}[\tau]}$. In the converse part, we use the binary KL divergence and apply the data processing inequality to upper bound error exponents.

Theorem 2 strengthens [2, Theorem 3.2] by deriving a matching converse result. We manage to do so for a slightly easier setting of exactly one outlier by excluding the null hypothesis where there might be no outliers. Furthermore, we propose another sequential test in addition to [2, Eq. (3.11)] and provide a relatively simpler achievability proof using our sequential test.

Since $GJS(P, Q, \alpha)$ increases in α , it follows that as the number M of observed sequences increases, the achievable error exponent increases. This is consistent with our intuition because with more samples, our estimation of the nominal distribution is more accurate. Thus, it is easier to identify the outlier. In the extreme case of $M \to \infty$, the exponent equals to $D(P_A || P_N)$, which is exactly the performance of knowing the nominal distribution [2, Prop. 3.1].

B. Expected Stopping Time Universality

1) Test Design and Intuition: Under the expected stopping time universality constraint, our sequential test $\Phi_{\text{Est}} = (\tau, \phi_{\tau})$ consists of a random stopping time and a decision rule. The stopping time τ satisfies

$$\tau := \inf\{k \ge n - 1 : \exists i \in [M] \text{ s.t. } \mathcal{S}_i(\mathbf{x}^k) \le f(k)\}, \quad (17)$$

where the scoring function $S_i(\mathbf{x}^k) = G_i(\hat{T}_{x_1^k}, \dots, \hat{T}_{x_M^k})$ and the threshold is given by $f(k) = \frac{(M+1)|\mathcal{X}|\log(k+1)}{k}$.

Note that the sequential test $\Phi_{\rm Est}$ stops if the types of all sequences except for outlying sequences are "close enough" to each other, where the threshold f(k) is used to characterize the closeness level.

At the stopping time τ , using M observed sequences \mathbf{x}_{M}^{τ} , our test applies the following minimal scoring function decision rule:

$$\phi_{\tau}(\mathbf{x}^{\tau}) = \mathbf{H}_i, \text{ if } i = i^*(\mathbf{x}^{\tau}), \tag{18}$$

where $i^*(\mathbf{x}^{\tau})$ is the index of the scoring function with smallest value, i.e.,

$$i^*(\mathbf{x}^{\tau}) := \operatorname*{arg\,min}_{i \in [M]} \mathcal{S}_i(\mathbf{x}^{\tau}). \tag{19}$$

Our test generalizes the sequential classification test under expected stopping time universality in [7, Def. 7].

We now explain the intuitive reason why the above test works using the weak law of large numbers. As discussed below (13), under hypothesis H_i , the scoring function $S_i(\mathbf{x}^k)$ tends to zero and scoring functions $S_j(\mathbf{x}^k)$ for each $j \in \mathcal{M}_i$ tend to a positive real number as the sample size k increases. Therefore, when k is sufficiently large, for each $i \in [M]$, if *i*th sequence is the outlier, our test stops and makes the correct decision H_i .

We remark that if the sample size is not large enough, the empirical distributions could be rather different from generating distributions, which might lead to decision error. To avoid such errors, similarly to [7, Def. 7], we set the minimal stopping time as n - 1 for some integer $n \in \mathbb{N}$.

2) Main Results and Discussions: Given any pair of distributions $(P,Q) \in \mathcal{P}(\mathcal{X})^2$ and any $\alpha \in \mathbb{R}_+$, the Rényi Divergence of order α [10, Eq. (1)] is defined as

$$D_{\alpha}(P||Q) := \frac{1}{\alpha - 1} \log \sum_{x \in \mathcal{X}} P(x)^{\alpha} Q(x)^{1 - \alpha}.$$
 (20)

The Rényi Divergence has the following variational form [7, Eq. (7)]:

$$D_{\frac{\alpha}{1+\alpha}}(P||Q) := \min_{V \in \mathcal{P}(\mathcal{X})} \alpha D(V||P) + D(V||Q).$$
(21)

Theorem 3. Under any pair of distributions $(P_N, P_A) \in \mathcal{P}(\mathcal{X})^2$, our sequential test satisfies the expected stopping time universality constraint and the error exponent of our test satisfies: for each $i \in [M]$,

$$E_i(\Phi_{\rm Est}|P_{\rm A}, P_{\rm N}) \ge D_{\frac{M-2}{M-1}}(P_{\rm N}||P_{\rm A}).$$
 (22)

Conversely, for any sequential test Φ_n satisfying the expected stopping time universality constraint, under any pair of distributions $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^2$, the error exponent satisfies that for each $i \in [M]$,

$$E_i(\Phi_n | P_A, P_N) \le D_{\frac{M-2}{M-1}}(P_N || P_A).$$
 (23)

The proof of Theorem 3 is provided in Appendix C and D. In the achievability part, we derive the lower bound of error exponents using method of types. In the converse part, we use the binary KL divergence and apply the data processing inequality to upper bound error exponents.

We make several remarks.

Firstly, Theorem 3 strengthens [2, Theorem 3.2] by considering sequential tests satisfying the expected stopping time universality constraint and deriving exact exponential decay rates of error probabilities. Compared with the error probability universality constraint in [2, Sec. 3.1.1] and Sec. III-A1, our proposed sequential test in this subsection has the property of having a bounded expected stopping time under any pair of nominal and anomalous distributions. This property is highly desired in practice since one desires the tests to stop early while the tests in [2, Sec. 3.1.1] and Sec. III-A1 could stop at very large sample sizes (cf. Section III-C for a numerical example).

Secondly, in contrast to lack of theoretical evidence that the sequential test of [2, Sec. 3.1.1] has better performance than the fixed-length test in (8), our sequential test has better theoretic performance than the fixed-length test in (8). This property is desired since the motivation of using sequential tests is to yield better performance. To clarify, we compare the exponents in Theorems 1 and 3. Given any distribution $Q \in \mathcal{P}(\mathcal{X})$, the distributions $\mathbf{Q} = (Q, Q, \dots, Q, P_N)$ satisfy the constraints that $G_1(\mathbf{Q}) \geq G_2(\mathbf{Q})$, it follows from (9) that $E_i(\Phi_{\text{Li}}|P_A, P_N) \leq \min_{\substack{Q \in \mathcal{P}(\mathcal{X})}} D(Q||P_A) + (M 2)D(Q||P_N) = D_{\frac{M-2}{M-1}}(P_N||P_A) = E_i(\Phi_{\text{Est}}|P_A, P_N)$. We numerically verify that $E_i(\Phi_{\text{Est}}|P_A, P_N)$ could be strictly greater than $E_i(\Phi_{\text{Li}}|P_A, P_N)$. Specifically, Under the nominal distribution $P_N = [0.3, 0.7]$ and the anomalous distribution $P_A = [0.1, 0.9]$, we have $E_i(\Phi_{\text{Est}}|P_A, P_N) = 0.0934 >$ $E_i(\Phi_{\text{Li}}|P_A, P_N) = 0.0471$.

Finally, we compare the optimal error exponents in Theorems 2 and 3, i.e., $\operatorname{GJS}(P_{\mathrm{N}}, P_{\mathrm{A}}, M-2)$ and $D_{\frac{M-2}{M-1}}(P_{\mathrm{N}}||P_{\mathrm{A}})$ respectively. It follows from [7, Remark 3] that for any $\alpha \in \mathbb{R}_+$ and any $(P_0, P_1) \in \mathcal{P}(\mathcal{X})^2$, $D_{\frac{\alpha}{1+\alpha}}(P_1||P_0) \geq$ $\operatorname{GJS}(P_0, P_1, \alpha)$. Thus, setting $\alpha = M - 2$ leads to $D_{\frac{M-2}{M-1}}(P_{\mathrm{N}}||P_{\mathrm{A}}) \geq \operatorname{GJS}(P_{\mathrm{A}}, P_{\mathrm{N}}, M - 2)$, which implies that when M = 3, the sequential test under the expected stopping time constraint achieves better error exponent because $D_{\frac{1}{\alpha}}(P_{\mathrm{N}}||P_{\mathrm{A}}) \geq \operatorname{GJS}(P_{\mathrm{A}}, P_{\mathrm{N}}, 1) = \operatorname{GJS}(P_{\mathrm{N}}, P_{\mathrm{A}}, 1)$.

TABLE I Comparison of Achievable Error Exponents of Two Sequential Tests Under Different Universality Constraints

Parameters	$D_{\frac{M-2}{M-1}}(P_{\mathrm{N}} P_{\mathrm{A}})$	$\mathrm{GJS}(P_{\mathrm{N}}, P_{\mathrm{A}}, M-2)$
M = 3 $P_{\rm N} = [0.2, 0.8]$ $P_{\rm A} = [0.4, 0.6]$	0.0493	0.0483
$M = 4 P_{\rm N} = [0.2, 0.8] P_{\rm A} = [0.4, 0.6]$	0.0642	0.0659
$M = 4 P_{\rm N} = [0.3, 0.7] P_{\rm A} = [0.1, 0.9]$	0.0939	0.0830

However, since $GJS(P_A, P_N, M-2) \neq GJS(P_N, P_A, M-2)$ when M > 3, the performance comparison of sequential tests under two universality constraints depends on the nominal and anomalous distributions. To illustrate, in Table I, we calculate the exponents in Theorems 2 and 3 for various pairs of distributions.

C. Numerical Example

Consider the binary alphabet $\mathcal{X} = \{0, 1\}$ and set M = 4. We set the nominal distribution $P_{\rm N} = [0.25, 0.75]$ and the anomalous distribution $P_{\rm A} = [0.3, 0.7]$. We simulate the expected stopping times of our tests under two universality constraints using 5×10^4 independent experiments. The simulation shows that the test $\Phi_{\rm Est}$ under the expected stopping time universality has expected stopping time of 6000 while the test $\Phi_{\rm Ep}$ under the error probability universality test has expected stopping time of 20015. Thus, the numerical results verifies the advantage of sequential tests under the expected stopping time universality in terms of sample complexity.

IV. GENERALIZATION TO MULTIPLE OUTLIERS

Let S denote the set of all subsets of [M] whose size is T, i.e., $S := \{\mathcal{B} \subset [M] : |\mathcal{B}| = T\}$. Our task now is to design a test $\Phi = \{\tau, \phi_{\tau}\} : \mathcal{X}^{M_{\tau}} \to \{\{H_{\mathcal{B}}\}_{\mathcal{B} \in S}\}$ with a stopping time τ and a corresponding decision rule ϕ_{τ} to classify among the following |S| hypotheses:

H_B where B ∈ S: the set of outlying sequences are sequences X^τ_j with j ∈ B.

To evaluate the performance of a test, we use the following misclassification error exponent under each hypothesis $H_{\mathcal{B}}$ with $\mathcal{B} \in \mathcal{S}$:

$$E_{\mathcal{B}}(\Phi|P_{\mathcal{A}}, P_{\mathcal{N}}) := \liminf_{n \to \infty} \frac{-\log \mathbb{P}_{\mathcal{B}}\{\Phi(\mathbf{X}^{\tau}) \neq \mathcal{H}_{\mathcal{B}}\}}{n}.$$
 (24)

To present our tests, we need the following definition [3, Eq. (42)]. Given a tuple of distributions $\mathbf{Q} = (Q_1, \ldots, Q_M) \in \mathcal{P}(\mathcal{X})^M$, for each $\mathcal{B} \in \mathcal{S}$, define

$$G_{\mathcal{B}}(\mathbf{Q}) = \sum_{j \in \mathcal{M}_{\mathcal{B}}} D\left(Q_{j} \left\| \frac{\sum_{l \in \mathcal{M}_{\mathcal{B}}} Q_{l}}{M - |\mathcal{B}|} \right),$$
(25)

where $\mathcal{M}_{\mathcal{B}} := [M] \setminus \mathcal{B} = \{i \in [M] : i \notin \mathcal{B}\}$. Note that $G_{\mathcal{B}}(\mathbf{Q})$ measures the similarity of distributions \mathbf{Q} except $\{Q_i\}_{i \in \mathcal{B}}$.

The measure $G_{\mathcal{B}}(\mathbf{Q}) = 0$ if and only if $Q_j = Q$ for all $j \in \mathcal{M}_{\mathcal{B}}$ for an arbitrary distribution $Q \in \mathcal{P}(\mathcal{X})$.

Under the expected stopping time universality constraint, our sequential test $\Phi_{\rm Est} = (\tau, \phi_{\tau})$ consists of the random stopping time and the decision rule. The stopping time τ is defined as follows:

$$\tau := \inf\{k \ge n - 1 : \exists \ \mathcal{B} \in \mathcal{S} \text{ s.t. } S_{\mathcal{B}}(\mathbf{x}^k) \le f(k)\}.$$
(26)

where the scoring function $S_{\mathcal{B}}(\mathbf{x}^k) = G_{\mathcal{B}}(\hat{T}_{x_1^k}, \dots, \hat{T}_{x_M^k})$ and the threshold $f(k) = \frac{(M+1)|\mathcal{X}|\log(k+1)|}{k}$. Our test applies the following minimal scoring function decision rule:

$$\phi_{\tau}(\mathbf{x}^{\tau}) = \operatorname*{arg\,min}_{\mathcal{B}\in\mathcal{S}} S_{\mathcal{B}}(\mathbf{x}^{\tau}).$$
(27)

To present our results, we need the following two error exponent functions. Given any $\mathcal{B} \in \mathcal{S}$, for any nominal distribution $P_{\rm N}$ and anomalous distribution $P_{\rm A}$, define

$$\operatorname{LD}_{\mathcal{B}}(P_{\mathrm{N}}, P_{\mathrm{A}}, M) := \min_{Q \in \mathcal{P}(\mathcal{X})} \min_{\mathcal{C} \in S_{\mathcal{B}}} |\mathcal{B} \cap \mathcal{M}_{\mathcal{C}}|D(Q||P_{\mathrm{A}}) \\
 + (M - |\mathcal{B} \cup \mathcal{C}|)D(Q||P_{\mathrm{N}}), \quad (28)$$

$$\tilde{\operatorname{LD}}_{\mathcal{B}}(P_{\mathrm{N}}, P_{\mathrm{A}}, M) := \min_{(Q_{1}, Q_{2}) \in \mathcal{P}(\mathcal{X})^{2}} \min_{\mathcal{C} \in S_{\mathcal{B}}}$$

$$|\mathcal{B} \cap \mathcal{M}_{\mathcal{C}}|D(Q_{1}||P_{A}) + (M - |\mathcal{B} \cup \mathcal{C}|)D(Q_{1}||P_{N}) + |\mathcal{C} \cap \mathcal{M}_{\mathcal{B}}|D(Q_{2}||P_{N}) + |\mathcal{B} \cap \mathcal{C}|D(Q_{2}||P_{A}).$$
(29)

Theorem 4. Under any nominal distribution P_N and anomalous distribution P_A , our sequential test satisfies the expected stopping time universality constraint and the error exponent of our test satisfies that for each $\mathcal{B} \in \mathcal{S}$,

$$E_{\mathcal{B}}(\Phi_{\text{Est}}|P_{\text{A}}, P_{\text{N}}) \ge \text{LD}_{\mathcal{B}}(P_{\text{N}}, P_{\text{A}}, M).$$
(30)

Conversely, for any sequential test Φ_n satisfying the expected stopping time universality constraint, under any pair of nominal distribution P_N and anomalous distribution P_A , the error exponent satisfies that for each $\mathcal{B} \in \mathcal{S}$,

$$E_{\mathcal{B}}(\Phi_n | P_{\mathcal{A}}, P_{\mathcal{N}}) \le \tilde{\text{LD}}_{\mathcal{B}}(P_{\mathcal{N}}, P_{\mathcal{A}}, M).$$
(31)

The proof of Theorem 4 is similar to that of Theorem 3. The upper and lower bounds are not tight in general. However, when T = 1, the result in Theorem 4 specializes to Theorem 3.

V. CONCLUSION

We revisited sequential outlier hypothesis testing by proposing tests under two universality constraints and deriving bounds on the achievable error exponents. For the case of exactly one outlier, our results strengthen a previous result of [2] by having a matching converse result and analytically demonstrating the advantage of our sequential test. We also generalized our results to the case of multiple outliers when the number of outliers is known. In future, one can generalize our results to the case with an unknown number of outliers [3], [11], consider continuous observed sequences [8], [12] or propose low complexity tests [4], [13] that achieve performance close to the theoretical benchmarks derived in this paper.

REFERENCES

- [1] Y. Li, S. Nitinawarat, and V. V. Veeravalli, "Universal outlier hypothesis testing," *IEEE Trans. Inf. Theory*, vol. 60, no. 7, pp. 4066–4082, 2014.
- [2] vol. 36, no. 3, pp. 309-344, 2017.
- [3] L. Zhou, Y. Wei, and A. O. Hero, "Second-order asymptotically optimal outlier hypothesis testing," IEEE Trans. Inf. Theory, vol. 68, no. 6, pp. 3585-3607, 2022.
- [4] Y. Bu, S. Zou, and V. V. Veeravalli, "Linear-complexity exponentiallyconsistent tests for universal outlying sequence detection," IEEE Trans. Signal Process., vol. 67, no. 8, pp. 2115-2128, 2019.
- [5] M. Gutman, "Asymptotically optimal classification for multiple tests with empirically observed statistics," IEEE Trans. Inf. Theory, vol. 35, no. 2, pp. 401-408, 1989.
- [6] L. Zhou, V. Y. F. Tan, and M. Motani, "Second-order asymptotically optimal statistical classification," Information and Inference: A Journal of the IMA, vol. 9, no. 1, pp. 81-111, 2020.
- [7] C. Y. Hsu, C. F. Li, and I. H. Wang, "On universal sequential classification from sequentially observed empirical statistics," in IEEE ITW, 2022, pp. 642-647.
- [8] S. Zou, Y. Liang, H. V. Poor, and X. Shi, "Nonparametric detection of anomalous data streams," IEEE Trans. Signal Process., vol. 65, no. 21, pp. 5785-5797, 2017.
- [9] M. Haghifam, V. Y. F. Tan, and A. Khisti, "Sequential classification with empirically observed statistics," IEEE Trans. Inf. Theory, vol. 67, no. 5, pp. 3095-3113, 2021.
- [10] T. Van Erven and P. Harremos, "Rényi divergence and kullback-leibler divergence," IEEE Trans. Inf. Theory, vol. 60, no. 7, pp. 3797-3820, 2014.
- [11] L. Zhou, Y. Wei, and A. Hero, "Asymptotics for outlier hypothesis testing," in IEEE ISIT, 2022, pp. 3303-3308.
- [12] Y. Bu, S. Zou, Y. Liang, and V. V. Veeravalli, "Universal outlying sequence detection for continuous observations," in IEEE ICASSP, 2016, pp. 4254-4258.
- [13] L. Xiong, B. Póczos, and J. Schneider, "Group anomaly detection using flexible genre models," Advances in neural information processing systems, vol. 24, 2011.

A. Achievability Proof of Theorem 2

1) Error probability universality constraint: We first prove our sequential test in Sec. III-A1 satisfies the error probability universality constraint with β . Specifically, given any $\beta \in (0, 1)$, for each $i \in [M]$,

$$\beta_i(\phi_\tau | P_{\rm A}, P_{\rm N}) = \mathbb{P}_i\{\phi_\tau(\mathbf{X}^\tau) \neq \mathbf{H}_i\}$$
(32)

$$\leq \sum_{k=1}^{\infty} \mathbb{P}_i \{ \phi_\tau(\mathbf{X}^k) \neq \mathbf{H}_i \}$$
(33)

$$\leq \sum_{k=1}^{\infty} \mathbb{P}_i \{ \mathbf{S}_i(\mathbf{X}^k) > g(\beta, k) \}$$
(34)

$$\leq \sum_{k=1}^{\infty} \sum_{\substack{\mathbf{Q} \in \mathcal{P}_{k}(\mathcal{X})^{M}:\\ G_{i}(\mathbf{Q}) > g(\beta, k)}} \exp\left\{-k\left(D(Q_{i}||P_{A}) + \sum_{t \in \mathcal{M}_{i}} D(Q_{t}||P_{N})\right)\right\}$$
(35)

$$\leq \sum_{k=1}^{\infty} \sum_{\substack{\mathbf{Q}\in\mathcal{P}_{k}(\mathcal{X})^{M}:\\G_{i}(\mathbf{Q})>g(\beta,k)}} \exp\left\{-k\left(D(Q_{i}||P_{\mathbf{A}})+(M-1)D\left(\frac{\sum_{t\in\mathcal{M}_{i}}Q_{t}}{M-1}\left\|P_{\mathbf{N}}\right)+g(\beta,k)\right)\right\}$$
(36)

$$\leq \sum_{k=1}^{\infty} (k+1)^{M|\mathcal{X}|} \exp\{-kg(\beta,k)\}$$
(37)

$$=\sum_{k=1}^{\infty} (k+1)^{M|\mathcal{X}|} \beta(|\mathcal{X}|-1)(k+1)^{-(M+1)|\mathcal{X}|}$$
(38)

$$\leq \beta(|\mathcal{X}| - 1) \sum_{k=1}^{\infty} (k+1)^{-|\mathcal{X}|}$$
(39)

$$\leq \beta(|\mathcal{X}|-1) \int_0^\infty (u+1)^{-|\mathcal{X}|} \mathrm{d}u \tag{40}$$

$$=\beta(|\mathcal{X}|-1)\frac{1}{-|\mathcal{X}|+1}(u+1)^{-|\mathcal{X}|+1}\Big|_{u=0}^{u=\infty}$$
(41)

$$=\beta,$$
(42)

where (35) follows from the upper bound on the probability of the type class, (36) follows from $G_i(\mathbf{Q}) > g(\beta, k)$ and

$$\sum_{j \in \mathcal{M}_i} D(Q_j || P_{\mathrm{N}}) = (M-1)D\left(\frac{\sum_{t \in \mathcal{M}_i} Q_t}{M-1} \Big\| P_{\mathrm{N}}\right) + \mathcal{G}_i(\mathbf{Q}),\tag{43}$$

(37) follows from the fact that the number of the set of types of length n satisfies $|\mathcal{P}^n(\mathcal{X})| \leq (n+1)^{|\mathcal{X}|}$, (38) follows from the definition of $g(\beta, k)$ in Eq. (11), (40) follows from the similar manner as [7, Appendix C: 1)].

Then we have given any $\beta \in (0, 1)$, for each $i \in [M]$, our sequential test satisfies $\beta_i(\phi_\tau | P_A, P_N) \leq \beta$.

2) Achievable error exponents: Given $i \in [M]$, define the stopping time

$$\tau_i := \inf \left\{ k \in \mathbb{N} : \forall j \in \mathcal{M}_i, \mathcal{S}_j(\mathbf{x}^k) > g(\beta, k) \right\}.$$
(44)

We first prove $-\frac{\tau_i}{\log \beta}$ converge to $\frac{1}{\operatorname{GJS}(P_{\mathrm{N}}, P_{\mathrm{A}}, M-2)}$ in probability.

Proof. Under hypothesis H_i , by the definition of the stopping time τ_i , we have

$$\min_{j \in \mathcal{M}_i} \mathcal{S}_j(\mathbf{x}^{\tau_i}) > g(\beta, \tau_i).$$
(45)

In the following, we show that $\tau_i \to \infty$ as $\beta \to 0$. Given $k \in \mathbb{N}$, we have

$$\mathbb{P}_{i}\{\tau_{i} \leq k\} = \mathbb{P}_{i}\{\forall j \in \mathcal{M}_{i}, \mathcal{S}_{j}(\mathbf{x}^{\tau_{i}}) > g(\beta, \tau_{i}), \tau_{i} \leq k\}$$

$$(46)$$

$$\leq \mathbb{P}_i \left\{ kM \log(M-1) > -\log\left(\beta(|\mathcal{X}|-1)\right) \right\}$$
(47)

where (46) follows from Eq. (45) and (47) follows from $\forall j \in \mathcal{M}_i, \sum_{t \in \mathcal{M}_j} D\left(Q_t \left\| \frac{\sum_{l \in \mathcal{M}_j} Q_l}{M-1} \right) \le M \log(M-1) \text{ and } g(\beta, \tau_i) \ge g(\beta, k) \ge \frac{-\log(\beta(|\mathcal{X}|-1))}{k}$. Thus we can obtain that for each $i \in [M]$,

$$\mathbb{P}_i\{\tau_i \le k\} = 0, \quad \forall k < \frac{-\log\left(\beta(|\mathcal{X}| - 1)\right)}{M\log(M - 1)}.$$
(48)

Now we have $\tau_i \to \infty$ as $\beta \to 0$.

Furthermore, given $i \in [M]$, we obtain $\forall t \in \mathcal{M}_i$, $\hat{T}_{x_i^{\tau_i}} \to P_N$ and $\hat{T}_{x_i^{\tau_i}} \to P_A$ under hypothesis H_i when $\tau_i \to \infty$ as $\beta \to 0$ by the strong law of large numbers. Then with the continuity of KL-divergence, when $\tau_i \to \infty$, we have for each $j \in \mathcal{M}_i$,

$$S_{j}(\mathbf{x}^{\tau_{i}}) \to D\left(P_{A} \left\|\frac{P_{A} + (M-2)P_{N}}{M-1}\right) + (M-2)D\left(P_{N} \left\|\frac{P_{A} + (M-2)P_{N}}{M-1}\right)\right.$$

= GJS(P_N, P_A, M - 2). (49)

Recall the definition of the stopping time τ_i , we have

$$\min_{j \in \mathcal{M}_i} S_j(\mathbf{x}^{\tau_i - 1}) \le g(\beta, \tau_i - 1).$$
(50)

When $\tau_i \to \infty$, both $g(\beta, \tau_i)$ and $g(\beta, \tau_i - 1) \to -\frac{\log \beta}{\tau_i}$. Combining (50) and (49), we have

$$\lim_{\beta \to 0} -\frac{\log \beta}{\tau_i} \ge \text{GJS}(P_{\text{N}}, P_{\text{A}}, M - 2).$$
(51)

Then combining (45) and (49), we have

$$\lim_{\beta \to 0} -\frac{\log \beta}{\tau_i} \le \text{GJS}(P_{\text{N}}, P_{\text{A}}, M - 2).$$
(52)

Consequently, we conclude that

$$\lim_{\beta \to 0} -\frac{\tau_i}{\log \beta} = \frac{1}{\operatorname{GJS}(P_{\mathrm{N}}, P_{\mathrm{A}}, M - 2)}.$$
(53)

To go from the convergence in probability to convergence in mean, it suffices to prove that the sequence of random variables $-\frac{\tau_i}{\log \beta}$ is uniformly integrable as $\beta \to 0$. To prove the uniformly integrable, we need the following lemma.

Lemma 5. Given $k \in \mathbb{N}$, $k \ge 1$, there exists $(n', c) \in \mathbb{R}^2_+$ such that $\mathbb{P}_i\{\tau_i \ge k\} \le \frac{1}{\beta} \exp\{-ck\}$ for any $k \ge n'$.

Proof. Given $\varepsilon \in \mathbb{R}_+$ and $i \in [M]$, we have

$$\mathbb{P}_{i}\{\tau_{i} \geq k\} \leq \mathbb{P}_{i}\left\{\exists j \in \mathcal{M}_{i}, S_{j}(\mathbf{x}^{k-1}) \leq g(\beta, k-1)\right\}$$
(54)

$$\leq \sum_{j \in \mathcal{M}_i} \mathbb{P}_i \Big\{ S_j(\mathbf{x}^{k-1}) \leq g(\beta, k-1) \Big\}$$
(55)

$$\leq \sum_{j \in \mathcal{M}_{i}} \mathbb{P}_{i} \Big\{ S_{j}(\mathbf{x}^{k-1}) \leq g(\beta, k-1) \text{ and } D(\hat{T}_{x_{i}^{k-1}} \| P_{N}) \leq \varepsilon \text{ and } D(\hat{T}_{x_{t}^{k-1}} \| P_{N}) \leq \varepsilon, \forall t \in \mathcal{M}_{i} \Big\} \\ + \mathbb{P}_{i} \Big\{ D(\hat{T}_{x_{i}^{k-1}} \| P_{N}) > \varepsilon \text{ or } D(\hat{T}_{x_{t}^{k-1}} \| P_{N}) > \varepsilon, \exists t \in \mathcal{M}_{i} \Big\}.$$
(56)

The first term of (56) can be upper bounded as follows:

$$\sum_{j \in \mathcal{M}_{i}} \mathbb{P}_{i} \Big\{ S_{j}(\mathbf{x}^{k-1}) \leq g(\beta, k-1) \text{ and } D(\hat{T}_{x_{i}^{k-1}} \| P_{A}) \leq \varepsilon \text{ and } D(\hat{T}_{x_{t}^{k-1}} \| P_{N}) \leq \varepsilon, \forall t \in \mathcal{M}_{i} \Big\}$$

$$\leq \sum_{j \in \mathcal{M}_{i}} \mathbb{P}_{i} \Big\{ S_{j}(\mathbf{x}^{k-1}) \leq g(\beta, k-1) + S_{i}(\mathbf{x}^{k-1}) \text{ and } D(\hat{T}_{x_{i}^{k-1}} \| P_{A}) \leq \varepsilon \text{ and } D(\hat{T}_{x_{t}^{k-1}} \| P_{N}) \leq \varepsilon, \forall t \in \mathcal{M}_{i} \Big\}$$
(57)

$$\leq (M-1)\mathbb{P}_{i}\left\{\alpha(P_{\mathrm{N}}, P_{\mathrm{A}}) \leq g(\beta, k-1) + \mathcal{S}_{i}(\mathbf{x}^{k-1})\right\}$$
(58)

$$\leq \frac{M-1}{\beta(|\mathcal{X}|-1)} k^{(2M+1)|\mathcal{X}|} \exp\{-(k-1)\alpha(P_{\rm N}, P_{\rm A})\},\tag{59}$$

where (57) follows from $S_i(\mathbf{x}^k) \ge 0$, (58) follows from that when ε is sufficiently small, if $D(\hat{T}_{x_i^k}||P_A) \le \varepsilon$ and $D(\hat{T}_{x_i^k}||P_N) \le \varepsilon$ for all $t \in \mathcal{M}_i$, $\exists \alpha(P_N, P_A) > 0$ such that $S_j(\mathbf{x}^k) > \alpha(P_N, P_A)$ [2, Lemma B.1] and (59) follows from the upper bound on the probability of the type class and the number of the set of types. The second term of (56) can be upper bounded as follows using the upper bound of the probability of a type class:

$$\mathbb{P}_{i}\left\{D\left(\hat{T}_{x_{i}^{k-1}} \| P_{\mathcal{A}}\right) > \varepsilon \text{ or } D\left(\hat{T}_{x_{t}^{k-1}} \| P_{\mathcal{N}}\right) > \varepsilon, \exists t \in \mathcal{M}_{i}\right\} \le Mk^{|\mathcal{X}|} \exp\{-(k-1)\varepsilon\}.$$

$$(60)$$

Thus combining (59) and (60), we have that for all $k \ge n'$ and for some c and n' > 0,

$$\mathbb{P}_i\{\tau_i \ge k\} \le \frac{1}{\beta} \exp\{-ck\}.$$
(61)

Using Lemma 5 and [7, Lemma 5], we have that $\{\tau_i/\log\beta\}_{\beta\in(0,0.9]}$ is uniformly integrable. Therefore, we can obtain the convergence in mean of $\{\tau_i/\log\beta\}_{\beta\in(0,0.9]}$ from the convergence in probability in Eq. (53). Furthermore, we have $\tau \leq \tau_i$ by definition in (44). Then we have

$$E_i(\phi_\tau | P_{\rm A}, P_{\rm N}) = \liminf_{\beta \to 0} \frac{-\log \beta}{\mathbb{E}_i[\tau]}$$
(62)

$$\geq \liminf_{\beta \to 0} \frac{-\log \beta}{\mathbb{E}_i[\tau_i]} \tag{63}$$

$$= \mathrm{GJS}(P_{\mathrm{N}}, P_{\mathrm{A}}, M - 2).$$
(64)

B. Converse Proof of Theorem 2

Given $(p,q) \in (0,1)^2$, define the binary KL-divergence as follows:

$$d(p,q) := p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}.$$
(65)

The first-order derivative of d(p,q) on q is:

$$\frac{\partial d(p,q)}{\partial q} = \frac{q-p}{q(1-q)}.$$
(66)

Thus, d(p,q) is increasing in q when q > p and decreasing in q when p > q.

Given $j \in [M]$ and any $i \in \mathcal{M}_j$, define the event $\mathcal{W} := \{\phi_\tau(\mathbf{x}^\tau) = i\}$. For any two pairs of distributions $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^2$ and $(\tilde{P}_A, \tilde{P}_N) \in \mathcal{P}(\mathcal{X})^2$, and any sequential test $\Phi = (\tau, \phi_\tau)$, we have

$$d\left(\mathbb{P}_{i}(\mathcal{W}),\tilde{\mathbb{P}}_{j}(\mathcal{W})\right) \leq D(\mathbb{P}_{i}||\tilde{\mathbb{P}}_{j})|_{\mathcal{F}_{\tau}}$$

$$(67)$$

$$\leq \mathbb{E}_{i} \left[\sum_{t \in [M]: t \neq i, j} \sum_{k=1}^{r} \log \frac{P_{N}(X_{t,k})}{\tilde{P}_{N}(X_{t,k})} + \sum_{k=1}^{r} \log \frac{P_{A}(X_{i,k})}{\tilde{P}_{N}(X_{i,k})} + \sum_{k=1}^{r} \log \frac{P_{N}(X_{j,k})}{\tilde{P}_{A}(X_{j,k})} \right]$$
(68)

$$\leq (M-2)\mathbb{E}_{i}[\tau]D(P_{\mathrm{N}}||\tilde{P}_{\mathrm{N}}) + \mathbb{E}_{i}[\tau]D(P_{\mathrm{A}}||\tilde{P}_{\mathrm{N}}) + \mathbb{E}_{i}[\tau]D(P_{\mathrm{N}}||\tilde{P}_{\mathrm{A}}),$$
(69)

where (67) follows from data processing inequality of divergence and (69) follows from Doob's Optional Stopping Theorem. Since $\beta_i(\phi_\tau | P_A, P_N) \to 0$ for each $i \in [M]$ as $n \to \infty$, we have

$$\mathbb{P}_{i}(\mathcal{W}) = \mathbb{P}_{i}(\phi_{\tau}(\mathbf{x}^{\tau}) = i) = 1 - \beta_{i}(\Phi | P_{\mathrm{A}}, P_{\mathrm{N}}) \to 1,$$
(70)

$$\tilde{\mathbb{P}}_{i}(\mathcal{W}) = \tilde{\mathbb{P}}_{i}(\phi_{\tau}(\mathbf{x}^{\tau}) = i) \le \tilde{\beta}_{i}(\Phi | \tilde{P}_{A}, \tilde{P}_{N}) \to 0.$$
(71)

Thus we obtain $\mathbb{P}_i(\mathcal{W}) > \tilde{\mathbb{P}}_j(\mathcal{W})$ and $1 - \beta_i(\Phi|P_A, P_N) > \tilde{\beta}_j(\Phi|\tilde{P}_A, \tilde{P}_N)$. Furthermore, we have

$$d\left(\mathbb{P}_{i}(\mathcal{W}), \tilde{\mathbb{P}}_{j}(\mathcal{W})\right) \geq d\left(1 - \beta_{i}(\Phi|P_{\mathrm{A}}, P_{\mathrm{N}}), \tilde{\beta}_{j}(\Phi|\tilde{P}_{\mathrm{A}}, \tilde{P}_{\mathrm{N}})\right)$$
(72)

$$= -\log \tilde{\beta}_j(\Phi|\tilde{P}_{\rm A}, \tilde{P}_{\rm N}) \tag{73}$$

$$\geq -\log\beta. \tag{74}$$

where (72) follows from that d(p,q) is decreasing in q when p > q, (73) follows from [7, Lemma 2] when $n \to \infty$ and (74) follows from the definition of error probability universality constraint.

Given $i \in [M]$, combining (69) and (74), we have

$$-\log\beta \leq (M-2)\mathbb{E}_{i}[\tau]D(P_{\mathrm{N}}||\tilde{P}_{\mathrm{N}}) + \mathbb{E}_{i}[\tau]D(P_{\mathrm{A}}||\tilde{P}_{\mathrm{N}}) + \mathbb{E}_{i}[\tau]D(P_{\mathrm{N}}||\tilde{P}_{\mathrm{A}}).$$

$$\tag{75}$$

Therefore, for any sequential test $\Phi = (\tau, \phi_{\tau})$ satisfying the expected stopping time universality constraint and any two pairs of distributions $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^2$ and $(\tilde{P}_A, \tilde{P}_N) \in \mathcal{P}(\mathcal{X})^2$, the type-*i* error exponent satisfies

$$E_{i}(\Phi|P_{\rm A}, P_{\rm N}) \le (M-2)D(P_{\rm N}||\tilde{P}_{\rm N}) + D(P_{\rm A}||\tilde{P}_{\rm N}) + D(P_{\rm N}||\tilde{P}_{\rm A}).$$
(76)

Since (76) is true for all $(\tilde{P}_A, \tilde{P}_N) \in \mathcal{P}(\mathcal{X})^2$, we can minimize the RHS of (76) with respect to $(\tilde{P}_A, \tilde{P}_N)$ and obtain

$$E_i(\Phi|P_{\mathcal{A}}, P_{\mathcal{N}}) \le \min_{Q \in \mathcal{P}(\mathcal{X})} (M-2)D(P_{\mathcal{N}}||Q) + D(P_{\mathcal{A}}||Q)$$
(77)

$$= GJS(P_N, P_A, M - 2).$$
 (78)

C. Achievability Proof of Theorem 3

1) Expected stopping time universality constraint: We first prove our sequential test in Sec. III-B1 satisfies the expected stopping time universality constraint. The average stopping time can be expressed as the following form: for each $i \in [M]$

$$\mathbb{E}_{i}[\tau] = \sum_{k=1}^{\infty} \mathbb{P}_{i}\{\tau > k\} = n - 1 + \sum_{k=n-1}^{\infty} \mathbb{P}_{i}\{\tau > k\}.$$
(79)

Then we upper bound the term $\sum_{k=n-1}^{\infty} \mathbb{P}_i\{\tau > k\}$ as follows,

$$\mathbb{P}_{i}\{\tau > k\} \leq \mathbb{P}_{i}\left\{\bigcap_{t=1}^{k} S_{i}(\mathbf{X}^{t}) \ge f(t)\right\}$$
(80)

$$\leq \mathbb{P}_{i} \{ \mathbf{S}_{i}(\mathbf{X}^{k}) \geq f(k) \}$$

$$(81)$$

$$\sum_{k=1}^{n} \mathbb{P}_{i}(\mathbf{x}^{k}) \times \left(\prod_{k=1}^{n} \mathbb{P}_{i}(\mathbf{x}^{k}) \right)$$

$$(82)$$

$$= \sum_{\mathbf{x}^k \in \mathcal{X}^{Mk}: \mathbf{S}_i(\mathbf{x}^k) \ge f(k)} P_{\mathbf{A}}(x_i^k) \times \Big(\prod_{j \in \mathcal{M}_i} P_{\mathbf{N}}(x_j^k)\Big)$$
(82)

$$= \sum_{\mathbf{Q}\in\mathcal{P}_{k}(\mathcal{X})^{M}:\mathbf{G}_{i}(\mathbf{Q})\geq f(k)} P_{\mathbf{A}}(\mathcal{T}_{Q_{i}}^{k}) \times \left(\prod_{j\in\mathcal{M}_{i}} P_{\mathbf{N}}(\mathcal{T}_{Q_{j}}^{k})\right)$$
(83)

$$\leq \sum_{\mathbf{Q}\in\mathcal{P}_{k}(\mathcal{X})^{M}:\mathbf{G}_{i}(\mathbf{Q})\geq f(k)} \exp\left\{-kD(Q_{i}||P_{\mathbf{A}})-k\sum_{j\in\mathcal{M}_{i}}D(Q_{j}||P_{\mathbf{N}})\right\}$$
(84)

$$\leq \sum_{\mathbf{Q}\in\mathcal{P}_{k}(\mathcal{X})^{M}:\mathbf{G}_{i}(\mathbf{Q})\geq f(k)} \exp\left\{-k\left(D(Q_{i}||P_{\mathbf{A}})+(M-1)D\left(\frac{\sum_{t\in\mathcal{M}_{i}}Q_{t}}{M-1}\Big\|P_{\mathbf{N}}\right)+f(k)\right)\right\}$$
(85)

$$\leq \sum_{\mathbf{Q}\in\mathcal{P}_{k}(\mathcal{X})^{M}:\mathbf{G}_{i}(\mathbf{Q})\geq f(k)} \exp\{-kf(k)\}$$
(86)

$$\leq (k+1)^{M|\mathcal{X}|} (k+1)^{-(M+1)|\mathcal{X}|}$$
(87)

$$=(k+1)^{-|\mathcal{X}|},$$
 (88)

where (84) follows from the upper bound on the probability of the type class, (85) follows from the equation in Eq. (43) and $G_i(\mathbf{Q}) \ge f(k)$ when $\tau > t$, (87) follows from the fact that the number of the set of types of length n satisfies $|\mathcal{P}^n(\mathcal{X})| \le (n+1)^{|\mathcal{X}|}$.

Thus for $n \ge 2$, we have

$$\sum_{k=n-1}^{\infty} \mathbb{P}_{i}\{\tau > k\} \le \sum_{k=n-1}^{\infty} (k+1)^{-|\mathcal{X}|} \le \int_{n-2}^{\infty} (u+1)^{-|\mathcal{X}|} \mathrm{d}u = \frac{1}{-|\mathcal{X}|+1} (u+1)^{-|\mathcal{X}|+1} \Big|_{u=n-2}^{u=\infty}$$
(89)

$$=\frac{(n-1)^{-(|\mathcal{X}|-1)}}{|\mathcal{X}|-1} \le 1.$$
(90)

Combining with Eq. (79), we obtain that for each $i \in [M]$, under hypothesis H_i , $\mathbb{E}_i[\tau] \leq n$.

2) Achievable error exponents: Given any pair of distributions $(P_N, P_A) \in \mathcal{P}(\mathcal{X})^2$ and any $n \in \mathbb{N}, i \in [M]$, define the error exponent function

$$\Delta_i(n, P_{\mathcal{N}}, P_{\mathcal{A}}) := \min_{\substack{j \in \mathcal{M}_i \\ G_j(\mathbf{Q}) \le f(n)}} \min_{\substack{\mathbf{Q} \in \mathcal{P}_n(\mathcal{X})^M: \\ G_j(\mathbf{Q}) \le f(n)}} D(Q_i || P_{\mathcal{A}}) + \sum_{t \in \mathcal{M}_i} D(Q_t || P_{\mathcal{N}}).$$
(91)

We upper bound the misclassification error probability as follows: under hypothesis H_i ,

$$\beta_i(\phi_\tau | P_{\rm A}, P_{\rm N}) = \mathbb{P}_i\{\phi_\tau(\mathbf{X}^\tau) \neq i\}$$
(92)

$$= \mathbb{P}_i \Big\{ \bigcup_{j \in \mathcal{M}_i} \mathcal{S}_j(\mathbf{X}^{\tau}) \le f(\tau) \Big\}$$
(93)

$$\leq \mathbb{P}_{i} \Big\{ \bigcup_{j \in \mathcal{M}_{i}} \bigcup_{k=n-1}^{\infty} \mathrm{S}_{j}(\mathbf{X}^{k}) \leq f(k) \Big\}$$
(94)

$$\leq (M-1) \max_{j \in \mathcal{M}_i} \mathbb{P}_i \Big\{ \bigcup_{k=n-1}^{\infty} S_j(\mathbf{X}^k) \leq f(k) \Big\}$$
(95)

$$\leq (M-1) \max_{j \in \mathcal{M}_i} \sum_{k=n-1}^{\infty} \mathbb{P}_i \{ \mathcal{S}_j(\mathbf{X}^k) \leq f(k) \}$$
(96)

$$\leq (M-1) \max_{j \in \mathcal{M}_i} \sum_{k=n-1}^{\infty} \sum_{\mathbf{Q} \in \mathcal{P}_k(\mathcal{X})^M : G_j(\mathbf{Q}) \leq f(k)} \exp\left\{-k \left(D(Q_i||P_A) + \sum_{t \in \mathcal{M}_i} D(Q_t||P_N)\right)\right\}$$
(97)

$$\leq (M-1)\sum_{k=n-1}^{\infty} \exp\left\{-k\left(\Delta_i(k, P_{\mathrm{N}}, P_{\mathrm{A}}) - \frac{M|\mathcal{X}|\log(k+1)}{k}\right)\right\}$$
(98)

$$\leq (M-1) \sum_{k=n-1}^{\infty} \exp\Big\{-k\Big(\Delta_i(n-1, P_{\rm N}, P_{\rm A}) - \frac{M|\mathcal{X}|\log n}{n-1}\Big)\Big\},\tag{99}$$

where (97) follows from the same argument as Eq. (84) and (99) follows from both the function $\Delta_i(n, P_N, P_A)$ is increasing and $\frac{M|\mathcal{X}|\log n}{n-1}$ is decreasing in n.

As
$$f(n-1) \to 0$$
, $\frac{M|\mathcal{X}|\log n}{n-1}$ tends to 0 and $\Delta_i(n-1, P_{\mathrm{N}}, P_{\mathrm{A}})$ tends to
$$\min_{(Q_1, Q_2) \in \mathcal{P}(\mathcal{X})^2} D(Q_1||P_{\mathrm{A}}) + (M-2)D(Q_1||P_{\mathrm{N}}) + D(Q_2||P_{\mathrm{N}})$$

$$= \min_{Q \in \mathcal{P}(\mathcal{X})} D(Q||P_{\mathrm{A}}) + (M-2)D(Q||P_{\mathrm{N}})$$
(100)

$$= D_{\frac{M-2}{M-1}}(P_{\rm N}||P_{\rm A}), \tag{101}$$

where the last equality follows from the variational form of the Rényi Divergence in (21). Then (99) can be upper bounded as

$$(M-1)\sum_{k=n-1}^{\infty} \exp\left\{-kD_{\frac{M-2}{M-1}}(P_{\rm N}||P_{\rm A})\right\} = (M-1)\frac{\exp\left\{-(n-1)D_{\frac{M-2}{M-1}}(P_{\rm N}||P_{\rm A})\right\}}{1-\exp\left\{-D_{\frac{M-2}{M-1}}(P_{\rm N}||P_{\rm A})\right\}}.$$
(102)

Finally, we obtain the error exponent of misclassification error probability satisfies: for each $i \in [M]$,

$$E_{i}(\phi_{\tau}|P_{\rm A}, P_{\rm N}) \geq \liminf_{n \to \infty} \left\{ \frac{n-1}{n} D_{\frac{M-2}{M-1}}(P_{\rm N}||P_{\rm A}) - \frac{\log(M-1)}{n} + \frac{1}{n} \log\left(1 - \exp\left\{D_{\frac{M-2}{M-1}}(P_{\rm N}||P_{\rm A})\right\}\right) \right\}$$
(103)
$$= D_{\frac{M-2}{M-1}}(P_{\rm N}||P_{\rm A}).$$
(104)

D. Converse Proof of Theorem 3

Recall the definition of the binary KL-divergence in Eq. (65) and its increasing and decreasing in Appendix B.

Given $j \in [M]$ and any $i \in \mathcal{M}_j$, define the event $\mathcal{W} := \{\phi_\tau(\mathbf{x}^\tau) = i\}$. For any two pairs of distributions $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^2$ and $(\tilde{P}_A, \tilde{P}_N) \in \mathcal{P}(\mathcal{X})^2$, and any sequential test $\Phi = (\tau, \phi_\tau)$, we have

$$d\left(\mathbb{P}_{i}(\mathcal{W}),\tilde{\mathbb{P}}_{j}(\mathcal{W})\right) \leq (M-2)\mathbb{E}_{i}[\tau]D(P_{N}||\tilde{P}_{N}) + \mathbb{E}_{i}[\tau]D(P_{A}||\tilde{P}_{N}) + \mathbb{E}_{i}[\tau]D(P_{N}||\tilde{P}_{A})$$
(105)

$$\leq (M-2)nD(P_{\mathrm{N}}||\tilde{P}_{\mathrm{N}}) + nD(P_{\mathrm{A}}||\tilde{P}_{\mathrm{N}}) + nD(P_{\mathrm{N}}||\tilde{P}_{\mathrm{A}}),$$
(106)

where (105) follows from the similar argument as Eq. (69) and (106) follows from the definition of expected stopping time universality constraint.

Analogously to Eq. (73) in Appendix B, we have

$$d(\mathbb{P}_{i}(\mathcal{W}), \tilde{\mathbb{P}}_{j}(\mathcal{W})) \geq -\log \tilde{\beta}_{j}(\Phi | \tilde{P}_{A}, \tilde{P}_{N}).$$
(107)

Given $j \in [M]$, combining (106) and (107), we have

$$-\log\tilde{\beta}_j(\Phi|\tilde{P}_{\mathcal{A}},\tilde{P}_{\mathcal{N}}) \le (M-2)nD(P_{\mathcal{N}}||\tilde{P}_{\mathcal{N}}) + nD(P_{\mathcal{A}}||\tilde{P}_{\mathcal{N}}) + nD(P_{\mathcal{N}}||\tilde{P}_{\mathcal{A}}).$$
(108)

Therefore, for any sequential test $\Phi = (\tau, \phi_{\tau})$ satisfying the expected stopping time universality constraint and any two pairs of distributions $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^2$ and $(\tilde{P}_A, \tilde{P}_N) \in \mathcal{P}(\mathcal{X})^2$, the type-*j* error exponent satisfies

$$\tilde{E}_{j}(\Phi|\tilde{P}_{A},\tilde{P}_{N}) \le (M-2)D(P_{N}||\tilde{P}_{N}) + D(P_{A}||\tilde{P}_{N}) + D(P_{N}||\tilde{P}_{A}).$$
(109)

Since (109) is true for all $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^2$, we can minimize the RHS of (109) with respect to (P_A, P_N) and obtain

$$\tilde{E}_{j}(\Phi|\tilde{P}_{A},\tilde{P}_{N}) \leq \min_{Q \in \mathcal{P}(\mathcal{X})} (M-2)D(Q||\tilde{P}_{N}) + D(Q||\tilde{P}_{A})$$
(110)

$$= D_{\frac{M-2}{M-1}}(P_{\rm N}||P_{\rm A}). \tag{111}$$

E. Achievability Proof of Theorem 4

1) Expected stopping time universality constraint: The average stopping time can be expressed as the following form: for each $\mathcal{B} \in \mathcal{S}$,

$$\mathbb{E}_{\mathcal{B}}[\tau] = \sum_{k=1}^{\infty} \mathbb{P}_{\mathcal{B}}\{\tau > k\} = n - 1 + \sum_{k=n-1}^{\infty} \mathbb{P}_{\mathcal{B}}\{\tau > k\}.$$
(112)

Analogously to Appendix C1, we have that for each $\mathcal{B} \in \mathcal{S}$,

$$\mathbb{P}_{\mathcal{B}}\{\tau > k\} \\ \leq \mathbb{P}_{\mathcal{B}}\{\mathbf{S}_{i}(\mathbf{X}^{k}) \ge f(k)\}$$
(113)

$$= \sum_{\mathbf{Q}\in\mathcal{P}_{k}(\mathcal{X})^{M}:\mathbf{G}_{\mathcal{B}}(\mathbf{Q})\geq f(k)} \left(\prod_{i\in\mathcal{B}} P_{\mathbf{A}}(\mathcal{T}_{Q_{i}}^{k})\right) \times \left(\prod_{j\in\mathcal{M}_{\mathcal{B}}} P_{\mathbf{N}}(\mathcal{T}_{Q_{j}}^{k})\right)$$
(114)

$$\leq \sum_{\mathbf{Q}\in\mathcal{P}_{k}(\mathcal{X})^{M}:\mathbf{G}_{\mathcal{B}}(\mathbf{Q})\geq f(k)} \exp\left\{-k\sum_{i\in\mathcal{B}} D(Q_{i}||P_{\mathbf{A}}) - k\sum_{j\in\mathcal{M}_{\mathcal{B}}} D(Q_{j}||P_{\mathbf{N}})\right\}$$
(115)

$$\leq \sum_{\mathbf{Q}\in\mathcal{P}_{k}(\mathcal{X})^{M}:\mathbf{G}_{\mathcal{B}}(\mathbf{Q})\geq f(k)} \exp\left\{-k\left(\sum_{i\in\mathcal{B}} D(Q_{i}||P_{\mathbf{A}}) + (M-T)D\left(\frac{\sum_{t\in\mathcal{M}_{\mathcal{B}}} Q_{t}}{M-T} \left\|P_{\mathbf{N}}\right) + f(k)\right)\right\}$$
(116)

$$\leq \sum_{\mathbf{Q}\in\mathcal{P}_{k}(\mathcal{X})^{M}:\mathbf{G}_{\mathcal{B}}(\mathbf{Q})\geq f(k)} \exp\{-kf(k)\}$$
(117)

$$= (k+1)^{-|\mathcal{X}|},$$
(118)

where (116) follows from

$$\sum_{j \in \mathcal{M}_{\mathcal{B}}} D(Q_j || P_{\mathrm{N}}) = (M - T) D\left(\frac{\sum_{t \in \mathcal{M}_{\mathcal{B}}} Q_t}{M - T} \Big\| P_{\mathrm{N}}\right) + \mathcal{G}_{\mathcal{B}}(\mathbf{Q}),$$
(119)

and $G_{\mathcal{B}}(\mathbf{Q}) \geq f(k)$ when $\tau > t$. Then we obtain that

$$\sum_{k=n-1}^{\infty} \mathbb{P}_{\mathcal{B}}\{\tau > k\} \le \sum_{k=n-1}^{\infty} (k+1)^{-|\mathcal{X}|} \le 1,$$
(120)

and thus

$$\mathbb{E}_{\mathcal{B}}[\tau] = n - 1 + \sum_{k=n-1}^{\infty} \mathbb{P}_{\mathcal{B}}\{\tau > k\} \le n.$$
(121)

2) Achievable error exponents: We need the following definitions to present our proof. Recall the definition of $\mathcal{M}_{\mathcal{B}} = [M] \setminus \mathcal{B}$. Given any nominal distribution P_N and anomalous distributions $P_{\mathcal{B}}$, define

$$\Delta_{\mathcal{B}}(n, P_{\mathrm{N}}, P_{\mathrm{A}}) := \min_{\mathcal{C} \in S_{\mathcal{B}}} \min_{\substack{\mathbf{Q} \in \mathcal{P}_{n}(\mathcal{X})^{M}:\\\mathbf{G}_{\mathcal{C}}(\mathbf{Q}) \le f(n)}} \sum_{i \in \mathcal{B}} D(Q_{i} || P_{\mathrm{A}}) + \sum_{t \in \mathcal{M}_{\mathcal{B}}} D(Q_{t} || P_{\mathrm{N}}).$$
(122)

Define $S_{\mathcal{B}} := \{ \mathcal{C} \in \mathcal{S} : \mathcal{C} \neq \mathcal{B} \}$ for any $\mathcal{B} \in \mathcal{S}$.

Analogously to Appendix C2, we upper bound the misclassification error probability under hypothesis H_B as follows,

$$\beta_{\mathcal{B}}(\phi_{\tau}|P_{\mathrm{A}}, P_{\mathrm{N}}) = \mathbb{P}_{\mathcal{B}}\Big\{\bigcup_{\mathcal{C}\in\mathcal{S}_{\mathcal{B}}}\mathcal{S}_{\mathcal{C}}(\mathbf{X}^{\tau}) \le f(\tau)\Big\}$$
(123)

$$\leq \mathbb{P}_{\mathcal{B}} \Big\{ \bigcup_{\mathcal{C} \in \mathcal{S}_{\mathcal{B}}} \bigcup_{k=n-1}^{\infty} \mathcal{S}_{\mathcal{C}}(\mathbf{X}^{k}) \leq f(k) \Big\}$$
(124)

$$\leq (|S|-1) \max_{\mathcal{C} \in \mathcal{S}_{\mathcal{B}}} \sum_{k=n-1}^{\infty} \mathbb{P}_{\mathcal{B}} \{ \mathcal{S}_{\mathcal{C}}(\mathbf{X}^k) \leq f(k) \}$$
(125)

$$\leq (|S|-1) \max_{\mathcal{C} \in S_{\mathcal{B}}} \sum_{k=n-1}^{\infty} \sum_{\mathbf{Q} \in \mathcal{P}_{k}(\mathcal{X})^{M}: G_{\mathcal{C}}(\mathbf{Q}) \leq f(k)} \exp\left\{-k\left(\sum_{i \in \mathcal{B}} D(Q_{i}||P_{A}) + \sum_{t \in \mathcal{M}_{\mathcal{B}}} D(Q_{t}||P_{N})\right)\right\}$$
(126)

$$\leq (|S|-1)\sum_{k=n-1}^{\infty} \exp\left\{-k\left(\Delta_{\mathcal{B}}(k, P_{\mathrm{N}}, P_{\mathrm{A}}) - \frac{M|\mathcal{X}|\log(k+1)}{k}\right)\right\}$$
(127)

$$\leq (|S|-1) \sum_{k=n-1}^{\infty} \exp\Big\{-k\Big(\Delta_{\mathcal{B}}(n-1, P_{\rm N}, P_{\rm A}) - \frac{M|\mathcal{X}|\log n}{n-1}\Big)\Big\}.$$
(128)

As $f(n-1) \to 0$, $\frac{M|\mathcal{X}|\log n}{n-1} \to 0$ and $\Delta_{\mathcal{B}}(n-1, P_{\mathrm{N}}, P_{\mathrm{A}}) \to 0$

$$\min_{\substack{(Q,\{Q_j\}_{j\in\mathcal{C}})\in\mathcal{P}(\mathcal{X})^{T+1} \mathcal{C}\in\mathcal{S}_{\mathcal{B}}}} \min_{\mathcal{C}\in\mathcal{S}_{\mathcal{B}}} |\mathcal{B}\cap\mathcal{M}_{\mathcal{C}}|D(Q||P_{A}) + (M-|\mathcal{B}\cup\mathcal{C}|)D(Q||P_{N}) + \sum_{l\in\mathcal{C}\cap\mathcal{M}_{\mathcal{B}}} D(Q_{l}||P_{N}) + \sum_{t\in\mathcal{B}\cap\mathcal{C}} D(Q_{t}||P_{A})$$

$$= \min_{\substack{Q\in\mathcal{P}(\mathcal{X})\mathcal{C}\in\mathcal{S}_{\mathcal{B}}}} \min_{\mathcal{C}\in\mathcal{S}_{\mathcal{B}}} |\mathcal{B}\cap\mathcal{M}_{\mathcal{C}}|D(Q||P_{A}) + (M-|\mathcal{B}\cup\mathcal{C}|)D(Q||P_{N})$$
(129)

$$:= \mathrm{LD}_{\mathcal{B}}(P_{\mathrm{N}}, P_{\mathrm{A}}, M).$$
(130)

Then (128) can be upper bounded as follows:

$$(|S|-1)\sum_{k=n-1}^{\infty} \exp\{-k \cdot \operatorname{LD}(P_{\mathrm{N}}, P_{\mathrm{A}}, M)\} = (|S|-1)\frac{\exp\{-(n-1)\operatorname{LD}(P_{\mathrm{N}}, P_{\mathrm{A}}, M)\}}{1 - \exp\{-\operatorname{LD}(P_{\mathrm{N}}, P_{\mathrm{A}}, M)\}}.$$
(131)

Thus we obtain that the achievable error exponent of misclassification error probability satisfies: for each $\mathcal{B} \in \mathcal{S}$,

$$E_{\mathcal{B}}(\phi_{\tau}|P_{\mathcal{B}}, P_{\mathrm{N}}) \ge \mathrm{LD}_{\mathcal{B}}(P_{\mathrm{N}}, P_{\mathrm{A}}, M).$$
(132)

F. Converse Proof of Theorem 4

We have stated that the function d(p,q) in Eq. (65) is increasing in q when q > p and decreasing in q when p > q in Appendix D. Given $\mathcal{B} \in \mathcal{S}$, define the event $\mathcal{W} := \{\phi_{\tau}(\mathbf{x}^{\tau}) = \mathcal{B}\}$. Given $\mathcal{B} \in \mathcal{S}$ and $\mathcal{C} \in S_{\mathcal{B}}$, for any two pairs of distributions $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^2$ and $(\tilde{P}_A, \tilde{P}_N) \in \mathcal{P}(\mathcal{X})^2$, and any sequential test $\Phi = (\tau, \phi_{\tau})$, we have

$$d\left(\mathbb{P}_{\mathcal{B}}(\mathcal{W}), \tilde{\mathbb{P}}_{\mathcal{C}}(\mathcal{W})\right) \leq D\left(\mathbb{P}_{\mathcal{B}}||\tilde{\mathbb{P}}_{\mathcal{C}})|_{\mathcal{F}_{\tau}}$$

$$\leq \mathbb{E}_{\mathcal{B}}\left[\sum_{t \in \mathcal{M}_{\mathcal{B}\cup\mathcal{C}}} \sum_{k=1}^{\tau} \log \frac{P_{N}(X_{t,k})}{\tilde{P}_{N}(X_{t,k})} + \sum_{i \in \mathcal{B}\cap\mathcal{M}_{\mathcal{C}}} \sum_{k=1}^{\tau} \log \frac{P_{A}(X_{i,k})}{\tilde{P}_{N}(X_{i,k})} + \sum_{j \in \mathcal{C}\cap\mathcal{M}_{\mathcal{B}}} \sum_{k=1}^{\tau} \log \frac{P_{N}(X_{j,k})}{\tilde{P}_{A}(X_{j,k})} + \sum_{l \in \mathcal{B}\cap\mathcal{C}} \sum_{k=1}^{\tau} \log \frac{P_{A}(X_{l,k})}{\tilde{P}_{A}(X_{l,k})}\right]$$

$$(133)$$

$$(134)$$

$$\leq \mathbb{E}_{\mathcal{B}}[\tau] \Big((M - |\mathcal{B} \cup \mathcal{C}|) D(P_{\mathrm{N}} || \tilde{P}_{\mathrm{N}}) + |\mathcal{B} \cap \mathcal{M}_{\mathcal{C}}| D(P_{\mathrm{A}} || \tilde{P}_{\mathrm{N}}) + |\mathcal{C} \cap \mathcal{M}_{\mathcal{B}}| D(P_{\mathrm{N}} || \tilde{P}_{\mathrm{A}}) + |\mathcal{B} \cap \mathcal{C}| D(P_{\mathrm{A}} || \tilde{P}_{\mathrm{A}}) \Big)$$
(135)

$$\leq n\Big((M - |\mathcal{B} \cup \mathcal{C}|)D(P_{\mathrm{N}}||\tilde{P}_{\mathrm{N}}) + |\mathcal{B} \cap \mathcal{M}_{\mathcal{C}}|D(P_{\mathrm{A}}||\tilde{P}_{\mathrm{N}}) + |\mathcal{C} \cap \mathcal{M}_{\mathcal{B}}|D(P_{\mathrm{N}}||\tilde{P}_{\mathrm{A}}) + |\mathcal{B} \cap \mathcal{C}|D(P_{\mathrm{A}}||\tilde{P}_{\mathrm{A}})\Big),$$
(136)

where (133) follows from data processing inequality of divergence, (135) follows from Doob's Optional Stopping Theorem and (136) follows from the definition of expected stopping time universality constraint.

Since $\beta_{\mathcal{B}}(\Phi|P_A, P_N) \to 0$ for each $i \in [M]$ as $n \to \infty$, we have

$$\mathbb{P}_{\mathcal{B}}(\mathcal{W}) = \mathbb{P}_{\mathcal{B}}(\phi_{\tau}(\mathbf{x}^{\tau}) = \mathcal{B}) = 1 - \beta_{\mathcal{B}}(\Phi | P_{\mathrm{A}}, P_{\mathrm{N}}) \to 1,$$
(137)

$$\tilde{\mathbb{P}}_{\mathcal{C}}(\mathcal{W}) = \tilde{\mathbb{P}}_{\mathcal{C}}(\phi_{\tau}(\mathbf{x}^{\tau}) = \mathcal{B}) \le \tilde{\beta}_{\mathcal{C}}(\Phi | \tilde{P}_{A}, \tilde{P}_{N}) \to 0.$$
(138)

Thus we obtain $\mathbb{P}_{\mathcal{B}}(\mathcal{W}) > \tilde{\mathbb{P}}_{\mathcal{C}}(\mathcal{W})$ and $1 - \beta_{\mathcal{B}}(\Phi|P_A, P_N) > \tilde{\beta}_{\mathcal{C}}(\Phi|\tilde{P}_A, \tilde{P}_N)$. Furthermore, we have

$$d\left(\mathbb{P}_{\mathcal{B}}(\mathcal{W}), \tilde{\mathbb{P}}_{\mathcal{C}}(\mathcal{W})\right) \ge d\left(1 - \beta_{\mathcal{B}}(\Phi|P_{\mathcal{A}}, P_{\mathcal{N}}), \tilde{\beta}_{\mathcal{C}}(\Phi|\tilde{P}_{\mathcal{A}}, \tilde{P}_{\mathcal{N}})\right)$$
(139)

$$= -\log \tilde{\beta}_{\mathcal{C}}(\Phi|\tilde{P}_{\mathrm{A}},\tilde{P}_{\mathrm{N}}), \tag{140}$$

where (139) follows from that d(p,q) is decreasing in q when p > q and (140) follows from [7, Lemma 2] when $n \to \infty$. Given $C \in S$, combining (136) and (140), we have

$$-\log \tilde{\beta}_{\mathcal{C}}(\Phi|\tilde{P}_{A},\tilde{P}_{N}) \leq n \min_{\mathcal{B}\in S_{\mathcal{C}}} (M - |\mathcal{B}\cup\mathcal{C}|) D(P_{N}||\tilde{P}_{N}) + |\mathcal{B}\cap\mathcal{M}_{\mathcal{C}}|D(P_{A}||\tilde{P}_{N}) + |\mathcal{C}\cap\mathcal{M}_{\mathcal{B}}|D(P_{N}||\tilde{P}_{A}) + |\mathcal{B}\cap\mathcal{C}|D(P_{A}||\tilde{P}_{A}).$$
(141)

Therefore, for any sequential test $\Phi = (\tau, \phi_{\tau})$ satisfying the expected stopping time universality constraint and any two pairs of distributions $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^2$ and $(\tilde{P}_A, \tilde{P}_N) \in \mathcal{P}(\mathcal{X})^2$, the error exponent satisfies: given $\mathcal{C} \in \mathcal{S}$

$$\tilde{E}_{\mathcal{C}}(\phi_{\tau}|\tilde{P}_{A},\tilde{P}_{N}) \leq \min_{\mathcal{B}\in S_{\mathcal{C}}} (M - |\mathcal{B}\cup\mathcal{C}|)D(P_{N}||\tilde{P}_{N}) + |\mathcal{C}\cap\mathcal{M}_{\mathcal{B}}|D(P_{N}||\tilde{P}_{A}) + |\mathcal{B}\cap\mathcal{M}_{\mathcal{C}}|D(P_{A}||\tilde{P}_{N}) + |\mathcal{B}\cap\mathcal{C}|D(P_{A}||\tilde{P}_{A}).$$
(142)

Since (142) is true for all $(P_A, P_N) \in \mathcal{P}(\mathcal{X})^{T+1}$, we can minimize the RHS of (142) with respect to (P_A, P_N) and obtain

$$\tilde{E}_{\mathcal{C}}(\phi_{\tau}|\tilde{P}_{A},\tilde{P}_{N}) \leq \min_{(Q_{1},Q_{2})\in\mathcal{P}(\mathcal{X})^{2}} \min_{\mathcal{B}\in\mathcal{S}_{\mathcal{C}}} (M - |\mathcal{B}\cup\mathcal{C}|)D(Q_{1}||\tilde{P}_{N}) + |\mathcal{C}\cap\mathcal{M}_{\mathcal{B}}|D(Q_{1}||\tilde{P}_{A})
+ |\mathcal{B}\cap\mathcal{M}_{\mathcal{C}}|D(Q_{2}||\tilde{P}_{N}) + |\mathcal{B}\cap\mathcal{C}|D(Q_{2}||\tilde{P}_{A})$$

$$= \tilde{LD}_{\mathcal{C}}(\tilde{P}_{N},\tilde{P}_{A},M).$$
(143)