



Optimal Coding for Streaming Authentication and Interactive Communication

Matthew Franklin* Ran Gelles† Rafail Ostrovsky‡ Leonard J. Schulman§

June 6, 2013

Abstract

Error correction and message authentication are well studied in the literature, and various efficient solutions have been suggested and analyzed. This is however not the case for *data streams* in which the message is very long, possibly infinite, and not known in advance to the sender. Trivial solutions for error-correcting and authenticating data streams either suffer from a long delay at the receiver's end or cannot perform well when the communication channel is noisy.

In this work we suggest a constant-rate error-correction scheme and an efficient authentication scheme for data streams over a noisy channel (one-way communication, no feedback) in the shared-randomness model. Our first scheme does not assume shared randomness and (non-efficiently) recovers a $(1 - 2c)$ -fraction prefix of the stream sent so far, assuming the noise level is at most $c < 1/2$. The length of the recovered prefix is tight.

To be able to overcome the $c = 1/2$ barrier we relax the model and assume the parties pre-share a secret key. Under this assumption we show that for any given noise rate $c < 1$, there exists a scheme that correctly decodes a $(1 - c)$ -fraction of the stream sent so far with high probability, and moreover, the scheme is efficient. Furthermore, if the noise rate exceeds c , the scheme aborts with high probability. We also show that no constant-rate authentication scheme recovers more than a $(1 - c)$ -fraction of the stream sent so far with non-negligible probability, thus the relation between the noise rate and recoverable fraction of the stream is tight, and our scheme is optimal.

Our techniques also apply to the task of interactive communication (two-way communication) over a noisy channel. In a recent paper, Braverman and Rao [STOC 2011] show that any function of two inputs has a constant-rate interactive protocol for two users that withstands a noise rate up to $1/4$. By assuming that the parties share a secret random string, we extend this result and construct an interactive protocol that succeeds with overwhelming probability against noise rates up to $1/2$. We also show that no constant-rate protocol exists for noise rates above $1/2$ for functions that require two-way communication. This is contrasted with our first result in which computing the “function” requires only one-way communication and the noise rate can go up to 1.

Keywords: data stream, private codes, adversarial noise, authentication, tree codes, interactive communication.

*CS Department, UC Davis, franklin@cs.ucdavis.edu

†CS Department, UCLA, gelles@cs.ucla.edu

‡UCLA Department of Computer Science and Department Mathematics, 3732D Boelter Hall, Los Angeles CA 90095-1596, U.S. Email: rafail@cs.ucla.edu. Supported in part by NSF grants 0830803, 09165174, 1065276, 1118126 and 1136174, US-Israel BSF grant 2008411, OKAWA Foundation Research Award, IBM Faculty Research Award, Xerox Faculty Research Award, B. John Garrick Foundation Award, Teradata Research Award, and Lockheed-Martin Corporation Research Award. This material is based upon work supported by the Defense Advanced Research Projects Agency through the U.S. Office of Naval Research under Contract N00014-11-1-0392. The views expressed are those of the author and do not reflect the official policy or position of the Department of Defense or the U.S. Government.

§E&AS Division, Caltech. Email: schulman@caltech.edu. Supported in part by NSF Award 1038578.

1 Introduction

The tasks of *error-correction* and of *authentication* are well studied in the literature. In both cases, a sender (Alice) wishes to send a message over a one-way, noisy channel to a receiver (Bob). To do so, Alice produces a longer, redundant message and sends it over the channel. The added redundancy helps Bob in recovering the original message if possible, or aborting otherwise. The overhead of this process is the amount of redundancy added to each message; in this work we focus on *constant-rate* schemes, i.e., schemes in which the transmitted message is at most constant-times longer.

Interestingly, in all known authentication schemes (and in many of the error-correction codes) there are two important assumptions: (1) the message to be communicated has a given length n and (2) the message is fully known to the sender in advance. These two assumptions don't hold anymore when the information to be transmitted is in the form of a *data stream*, which is a long, possibly infinite, sequence of symbols x_1, x_2, \dots over some alphabet Σ , where each x_i arrives at the sender's end at time i and is unknown beforehand.

In this paper, we investigate the question of transmitting data streams over an adversarially noisy channel. Within this framework we consider two related questions, namely, error-correction and authentication of data streams. Loosely speaking, in error-correction schemes, the receiver decodes the correct message as long as the noise level is below some threshold (but possibly outputs a wrong message if the noise exceeds that threshold). In authentication schemes, the receiver's task is to indicate whether or not the received (decoded) message is indeed the one sent to him. To see the relation between these two tasks note that if the corruption level of an adversary is guaranteed to be lower than the threshold, any error-correction guarantees that the receiver decodes the original message. However, while no constant-rate error-correction scheme can withstand a noise level higher than $1/2$, this is not the case for authentication schemes that are capable of indicating a change in the message even when the adversary has a full control of the channel. On the other hand for the task of authentication, it is generally assumed that the parties pre-share a secret key.

Standard error-correction and authentication methods do not apply directly to the model of data streams. The straightforward method to perform error-correction (or authentication) of a data stream is to cut the stream into chunks and separately encode each chunk. The problem now is that while the adversary is limited to some *global* noise rate, there is no restriction on the noise level of any local part of the stream. Specifically, the adversary can corrupt a single chunk in its entirety (while not exceeding the global amount of allowed noise), and cause Bob to decode this chunk in a wrong way. Even if this event is noticed by Bob since the chunk fails the authentication, the information carried within this chunk is lost unless Bob requests a retransmission of that chunk, i.e., unless the communication is interactive. The same problem exists (with high probability) when the noise is random rather than adversarial, given that the stream is long enough or infinite.

A possible mitigation to the above is to increase the chunks' size. This, however, has an undesirable side effect—Bob needs to wait until receiving a complete chunk in order to decode and authenticate it. This means that the information received in the very recent bits is inaccessible to Bob until the chunk is completely received. Our goal is thus, to construct a constant-rate scheme that can withstand a constant fraction of errors (globally) and still guarantee the correct decoding and authenticity of the information received so far. To the best of our knowledge, no such solution is known.

1.1 Our Results

In this work we construct optimal encoding schemes for both interactive and non-interactive (streaming) communication, and show a dramatic difference between these two cases in the following

sense. For each case, we show an upper bound on the noise rates that allow a successful constant-rate communication, and construct a protocol that achieves the bound. Interestingly, the bound for one-way communication is different from the interactive one.

Specifically, our result for one-way communication is a constant-rate coding scheme for data streams that withstands noise rates of less than $1/2$. Informally, as long as the global noise rate up to some time n does not exceed some parameter $c < 1/2$, a fraction of $1 - 2c$ of the stream sent up to time n can be recovered (see Section 4). For constant-rate schemes, it is clear that $c < 1/2$ is a hard limit and no scheme can succeed when the noise is higher. In order to achieve schemes that withstand higher noise rates we must relax the model and give the users more resources. Indeed, with the use of shared randomness (i.e., a shared secret key) we can break the $c = 1/2$ barrier. To emphasize the fact that the parties are allowed to share a secret key, we refer schemes in this model as *authentication schemes* rather than error-correction schemes, based on the relation of these two tasks mentioned above (codes that assume a private shared key are also known as *private codes* [Lan04], see Related Work).

This leads to our first main result: we construct a constant-rate authentication scheme for data streams sent over a noisy (possibly adversarial) channel. For any constant fraction of noise c less than 1, our scheme succeeds in decoding at least a $(1 - c)$ -fraction of the stream so far, with high probability. The decoded part is always the *prefix* of the stream. The decoded prefix is authenticated, meaning that there is only a negligible probability that the scheme outputs a different string. Furthermore, our scheme is *efficient*. More formally (see formal theorems in Section 5), we show that for any noise rate $0 \leq c < 1$ and small constant $\varepsilon > 0$:

- There exists an efficient constant-rate scheme that, at time n , decodes a prefix of length at least $(1 - c)n - \varepsilon n$ of the stream sent so far.
- Any constant-rate protocol that decodes a prefix of length $(1 - c)n + \varepsilon n$ succeeds with probability at most $2^{-\Omega(\varepsilon n)}$ in the worst case.

Our scheme is unconditionally secure and does not make any (cryptographic) assumptions, other than pre-sharing a secret random string. The amount of randomness utilized by the scheme grows with the message length, and can be unbounded if the data stream is infinite. However, if we only consider a computationally bounded adversary, the required amount of randomness is relatively small (polynomial in the security parameter). With the aid of a pseudo-random generator, the parties only need to pre-share a small seed, from which they generate randomness at will. Moreover, such a solution scales to the multiparty case by a simple public-key infrastructure construction. Each user generates a pair of a public and a secret key, and any pair of users perform Diffie-Hellman key-exchange [DH76] to obtain a secret shared authentication-key used as the pseudo-random generator's seed.

We apply the same techniques used in our streaming-authentication scheme onto the task of *interactive communication* to get our second main result. In the interactive communication scenario, two parties perform an arbitrary interactive protocol over a noisy channel, while keeping the amount of exchanged data only a constant factor more than an equivalent protocol for a noiseless channel (i.e., the encoding is constant-rate). This question was initially considered for both random and adversarial noise by Schulman [Sch92, Sch93, Sch96] who showed a constant-rate encoding scheme that copes with a noise rate of up to $1/240$, and recently revisited by Braverman and Rao [BR11] who showed how to deal with noise rates less than $1/4$. In addition, Braverman and Rao show that $1/4$ is the highest error rate any protocol can withstand, as long as the protocol defines whose turn it is to speak at every round regardless of the observed noise. The fascinating open question left by the work of Braverman and Rao is whether other methods could extend the $1/4$ bound.

In this work we improve the bound obtained by [BR11] by allowing the parties to pre-share a secret key. Specifically, we show how to convert any interactive protocol (for noiseless channel) into a constant-rate protocol that withstands any adversarial noise level smaller than $1/2$, given pre-shared randomness. We also show that for higher noise rates, no constant-rate interactive protocol exists for tasks that depend on inputs of *both* parties. Similarly to previous results for interactive communication with adversarial noise [Sch96, BR11, GMS11], our decoding scheme is inefficient. Very recently, Brakerski and Kalai [BK12] showed how to augment previous results of interactive communication protocols and achieved *efficient* schemes that withstand adversarial noise (the computation efficiency was further improved by Brakerski and Naor [BN13] to $O(N \log N)$). Note that the bounds (on adversarial noise) obtained by [BK12] are improved by our work as well, since we improve the bounds of the underlying schemes used by [BK12].

1.2 Our Methods

The Blueberry code. The main ingredient of our construction is an error-detection code we name the *Blueberry code*¹. The Blueberry code uses the shared randomness in order to detect corruptions made by the channel, and marks them as *erasures*. One can think about this code as a weak message authentication code (MAC) that authenticates each symbol separately with a constant probability (see [Gol04] for a formal definition of MAC). To this end, each symbol of the input alphabet Σ is randomly and independently mapped to a larger alphabet Γ (the channel alphabet). This means that only a small subset of the channel alphabet is meaningful and the other symbols serve as “booby-traps”. Since each symbol is encoded independently, any corruption is caught with constant probability $\frac{|\Sigma|-1}{|\Gamma|-1}$ and marked with a special sign \perp to denote it was deleted by the channel. Most of the corruptions made by an adversary become erasures and only a small fraction (arbitrarily small, controlled by the size of $|\Gamma|$) turns into errors.

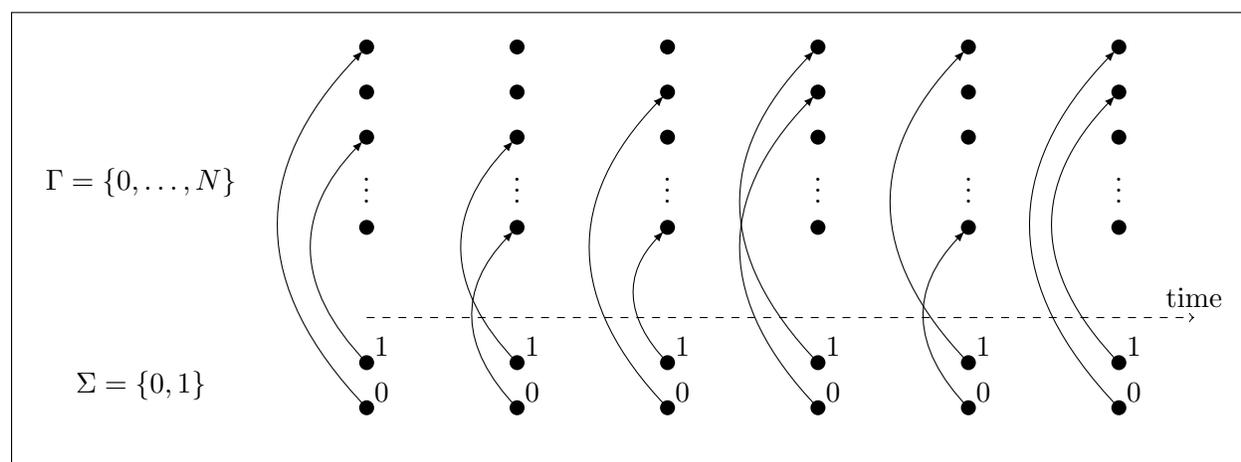


Figure 1: A demonstration of the Blueberry code: at any given time each symbol in Σ is randomly mapped to a symbol of Γ . Symbols of Γ with no incoming arrow are “booby-traps”, which detect corruptions.

The main insight that leads to our results is the different ways error correction codes deal with errors and erasures. We observe that, in terms of Hamming distance, the impact of a single error is twice as harmful as a single erasure. Indeed, assume that the Hamming distance of two

¹The name of the Blueberry code is inspired by the children’s book “The case of the hungry stranger” [Bon63] in which a blueberry pie is gone missing, and the thief (who turns out to be the dog) is identified by his big blue grin.

strings, x and y , is m . Then if x was communicated but y is decoded it means that at least $m/2$ errors have occurred, or alternatively, at least m erasures. More generally, assuming we decode by minimizing the Hamming distance, then our decoding fails if the number of errors e and the number of erasures d satisfy $2e + d \geq m$.

Combining Blueberry codes and tree codes. The second ingredient of our work is encoding via *tree codes* [Sch96], an online encoding that has a “self-healing” property: when decoding a stream at time n , the tree will decode correctly up to a particular time t such that the stream suffix between times t and n is the longest suffix in which the error rate is high. This means, for instance, that even if all the transmissions until some time t' were corrupted (and thus the decoding failed at those times), if the noise rate up to time $n > t'$ is low enough, not only can we decode between t' and n , but we will also be able to decode the *entire* stream up to time n .

Encoding via both a tree code and a Blueberry code immediately gives a streaming authentication method: the Blueberry code prevents the adversary from corrupting too many transmissions without being noticed, and given that the noise level is low enough, the tree code correctly decodes a prefix of the stream whose length is determined by the average noise level up to that time.

Efficient constructions. The only caveat of the above construction is that tree code decoding is not necessarily efficient and may be in the worst case exponential in the length of the received transmission. We obtain an efficient authentication scheme by splitting the stream into small segments and repeatedly sending random segments of the history. That way, even if some part of the transmission was changed by the channel, the same information will keep being retransmitted at random future times, and eventually (with high probability) will be received at the other side intact.

Roughly speaking, we use $n/\log n$ tree codes to encode chunks of the stream (each of length roughly $\log n$). Note that as n grows, so does the number of the trees in use, and the expected depth of each tree. At each time step, we randomly select one of the $n/\log n$ trees and transmit the next label of the path defined by the corresponding chunk of the stream. For most of the trees, the expected number of labels transmitted is $\Theta(\log n)$, and the decoding of the specific chunk succeeds except with polynomially small probability. Since each tree code is used to encode a word of length $O(\log n)$, the decoding can be performed efficiently by an exhaustive search.

1.3 Other Related Works

The works of Even, Goldreich and Micali [EGM90] and Gennaro and Rohatgi [GR97] consider authentication of data streams, however the focus of these schemes is not only to authenticate the message but also to prevent the sender from denying having signed the information. These constructions rely on cryptographic primitives such as one-time signatures. Another related line of research [PCTS00, MS01, GM01] pursues authentication of streams over *lossy* channels, usually in the multicast setting.

Coding schemes that assume the parties pre-share some randomness (also known as *Private Codes* [Lan04]) first appeared in [Sha58], and were greatly analyzed since. The main advantage of such codes is that they can deal with *adversarial noise*, rather than a random noise. Langberg [Lan04] considers private codes for adversarial channels that approach Shannon’s bound and require only $O(\log n)$ randomness for block size n , as well as an $\Omega(\log n)$ lower bound for the needed randomness. The construction of Langberg also implies an efficient code with $O(n \log n)$ randomness. This result was improved to $n + o(n)$ randomness by Smith [Smi07]. Explicit constructions with $o(n)$ randomness are yet unknown (see [Smi07]).

Error correction codes for *computationally bounded* noise models were first addressed by Lipton [Lip94] who constructs error-correction codes given pre-shared randomness and later considered by Micali, Peikert, Sudan and Wilson [MPSW05] who only assume sharing a short public-key, and recently by the surprising result of Guruswami and Smith [GS10] who assume no shared setup between the users. Locally Decodable codes with constant-rate in the public-key model were introduced by Hemenway and Ostrovsky [HO08] and later improved by Hemenway, Ostrovsky, Strauss and Wootters [HOSW11].

2 Preliminaries, Model and Definitions

We denote the set $\{1, 2, \dots, n\}$ by $[n]$, and for a finite set Σ we denote by $\Sigma^{\leq n}$ the set $\cup_{k=1}^n \Sigma^k$. The Hamming distance $\Delta(x, y)$ of two strings $x, y \in \Sigma^n$ is the number of indices i for which $x_i \neq y_i$. Throughout the paper, $\log()$ denotes the binary logarithm (base 2) and $\ln()$ denotes the natural logarithm (base e).

Shared randomness model. We assume the following *shared-randomness model*. The legitimate users (Alice and Bob) have access to a random string R of unbounded length, which is unknown to the adversary (Eve). Protocols in this model are thus *probabilistic*, and are required to succeed with high probability over the choice of R . We assume that all the randomness comes from R and that for a fixed R the protocols are deterministic.

Tree codes. A d -ary tree code [Sch96] over alphabet Σ is a rooted d -regular tree of arbitrary depth N whose edges are labeled with elements of Σ . For any string $x \in [d]^{\leq N}$, a d -ary tree code \mathcal{T} implies an *encoding* of x , $\text{TCenc}_{\mathcal{T}}(x) = w_1 w_2 \dots w_{|x|}$ with $w_i \in \Sigma$, defined by concatenating the labels along the path defined by x , i.e., the path that begins at the root and whose i -th node is the x_i -th child of the $(i-1)$ -th node. We usually omit the subscript \mathcal{T} when the tree is clear from the context. Note that tree code encoding is *online*: to communicate $\text{TCenc}(x\sigma)$ where $\sigma \in [d]$ given that $\text{TCenc}(x)$ was already communicated, we only need to send one symbol of Σ . Hence, if $|\Sigma| = O(1)$ the encoding scheme has a constant rate.

For any two paths (strings) $x, y \in [d]^{\leq N}$ of the same length n , let ℓ be the longest common prefix of both x and y . Denote by $\text{anc}(x, y) = n - |\ell|$ the distance from the n -th level to the least common ancestor of paths x and y . A tree code has distance α if for any $k \in [N]$ and any distinct $x, y \in [d]^k$, the Hamming distance of $\text{TCenc}(x)$ and $\text{TCenc}(y)$ is at least $\alpha \cdot \text{anc}(x, y)$.

For a string $w \in \Sigma^n$, decoding w using the tree code \mathcal{T} means returning the string $x \in [d]^n$ whose encoding minimizes the Hamming distance to the received word, namely,

$$\text{TCdec}_{\mathcal{T}}(w) = \underset{x \in [d]^n}{\text{argmin}} \Delta(\text{TCenc}_{\mathcal{T}}(x), w).$$

A theorem by Schulman [Sch96] proves that for any d and $\alpha < 1$ there exists a d -ary tree code of unbounded depth and distance α over alphabet of size $d^{O(1/(1-\alpha))}$. However, no efficient construction of such a tree is yet known. For a given depth N , Peczarski [Pec06] gives a randomized construction for a tree code with $\alpha = 1/2$ that succeeds with probability at least $1 - \epsilon$, and requires alphabet of size at least $d^{O(\sqrt{\log \epsilon^{-1}})}$. Braverman [Bra12] gives a sub-exponential (in N) construction of a tree code, and Gelles, Moitra and Sahai [GMS11] provide an efficient construction of a randomized relaxation of a tree code of depth N , namely a *potent tree code*, which is powerful enough as a substitute for a tree code in most applications.

Communication model. Our communication model consists of a channel $ch : \Sigma \rightarrow \Sigma$ subject to corruptions made by an adversary (or by the channel itself). The noise model is such that any symbol σ sent through the channel can turn into another symbol $\tilde{\sigma} \in \Sigma$. It is not allowed to insert or delete symbols. For all of our applications we assume that one symbol $\sigma_i \in \Sigma$ is sent at any time slot i .² We say that the adversarial *corruption rate* is c if for n transmissions, at most cn symbols were corrupted.

3 The Blueberry Code

Definition 3.1. For $i \geq 1$ let $B_i : [L + 1] \rightarrow [L + 1]$ be a random and independently chosen permutation. The Blueberry code maps a string x of arbitrary length n to

$$B(x) = B_1(x_1)B_2(x_2) \cdots B_n(x_n).$$

We denote such a code as $B : [L + 1]^* \rightarrow [L + 1]^*$.

We use the Blueberry code in the shared-randomness model where the legitimate parties share the random permutations B_i , unknown to the adversary (these kind of codes, determined by a random string unknown to the channel are referred to as *private codes* by [Lan04]). Although B_i is a permutation on $[L + 1]$, we actually use it to encode strings over a smaller alphabet $[S + 1]$ with $S < L$; that is, we focus on the induced mapping $B : [S + 1]^* \rightarrow [L + 1]^*$. The adversary does not know the specific permutations B_i , and has probability of at most S/L to change a transmission into a symbol whose pre-image is in $[S + 1]$.

Definition 3.2. Assume that at some time i , $y_i = B_i(x_i)$ is transmitted and $\tilde{y}_i \neq y_i$ is received. If $B_i^{-1}(\tilde{y}) \notin [S + 1]$, we mark the transmission as an erasure (specifically, the decoding algorithm outputs \perp); otherwise, this event is called an error.

Corollary 3.3. Let $x \in [S + 1]^n$ and assume $B(x)$ is communicated over a noisy channel. Every symbol altered by the channel will cause either an error with probability S/L , or an erasure with probability $1 - S/L$.

Assuming $S \ll L$, most of the corruptions done by the channel are marked as erasures, and only a small fraction of the corruptions percolate through the Blueberry code and cause an error.

Lemma 3.4. Let $S, L \in \mathbb{N}$ be fixed and assume a Blueberry code $B : [S + 1]^* \rightarrow [L + 1]^*$ is used to transmit a string $x \in [S + 1]^n$ over a noisy channel. For any constant $0 \leq c \leq 1$, if the channel's corruption rate c , then with probability $1 - 2^{-\Omega(n)}$ at least a $(1 - 2\frac{S}{L})$ -fraction of the corruptions are marked as erasures.

Proof. Denote by z_i the random variable which is 1 if the i -th corrupted-transmission is marked as an erasure and 0 otherwise. These are independent Bernoullis with probability $1 - \frac{S}{L}$. Let $Z = \sum_i z_i$ and note that $\mathbb{E}[Z] = cn(1 - \frac{S}{L})$. By Chernoff-Hoeffding inequality,

$$\Pr_R \left[\frac{1}{n} \sum_i z_i < c \left(1 - 2\frac{S}{L}\right) \right] < e^{-2n(cS/L)^2}.$$

□

²The channel time slots need not correspond with the times in which stream symbols are received. I.e, it is possible that between the arrival of stream elements x_i and x_{i+1} , several channel-symbols are transmitted.

Corollary 3.5. *Let $S, L \in \mathbb{N}$ be fixed. If out of n received transmissions, cn were marked as erasures by a Blueberry code $B : [S + 1]^* \rightarrow [L + 1]^*$, then except with probability $2^{-\Omega(n)}$ over the shared randomness, the adversarial corruption rate is at most $c/(1 - 2\frac{S}{L})$.*

We will use the Blueberry code concatenated with another (outer) code that is less sensitive to erasures than to errors. From the outer code’s point of view, this effectively increases the channel’s “error rate resilience” from $1 - 2c$ to $1 - c(1 + S/L)$. The construction of the code B from independent B_i ’s allows us to encode and decode each x_i independently, which is crucial for on-line applications in which the message x to be sent is not fully known in advance.

4 Error Correction of Data Streams

As a simple exposition to our main result, we begin with a simple, non-efficient, constant-rate error-correction scheme for data streams that withstands noise $c < 1/2$ and decodes a prefix of length $1 - 2c$ of the stream sent so far. The scheme is obtained by simply encoding the stream via a tree code \mathcal{T} with large enough distance parameter $\alpha \in (0, 1)$ and a constant-size alphabet, which depends on α .

Theorem 4.1. *For any constants $c < 1/2$ and $\varepsilon > 0$ there exists a constant-rate error-correction scheme for data stream x_1, x_2, \dots such that at any given time n the receiver outputs a string x'_1, x'_2, \dots, x'_n , and if the noise rate until time n is at most c , then*

$$x'_1, x'_2, \dots, x'_{(1-2c)n-\varepsilon n} = x_1, x_2, \dots, x_{(1-2c)n-\varepsilon n}$$

that is, a prefix of the stream of length at least $(1 - 2c)n - \varepsilon n$ is correctly decoded.

Proof. Assume Alice encodes each stream symbol using $\text{TCenc}_{\mathcal{T}}()$ using some tree code \mathcal{T} whose parameters we fix shortly.

For a specific time n , consider a string $\tilde{x} \in \{0, 1\}^n$, such that $\text{anc}(x, \tilde{x}) \geq (2c + \varepsilon)n$. Due to the tree distance property, the Hamming distance between $\text{TCenc}(\tilde{x})$ and $\text{TCenc}(x)$ is at least $\alpha(2c + \varepsilon)n$. Assume Eve causes e errors, a maximal-likelihood decoding will prefer x over \tilde{x} as long as $\lfloor \alpha(2c + \varepsilon)n \rfloor > 2e$. Since Eve’s corruption rate is limited to c , we know that $e \leq cn$. By setting $\alpha > \frac{2c}{2c + \varepsilon}$ we guarantee that $\alpha(2c + \varepsilon)n > 2e$, and Bob decodes a string x' such that $\text{anc}(x, x') < (2c + \varepsilon)n$ with certainty. \square

5 Perpetual Authentication

Sending a data stream over a noisy channel is not a simple task, especially when the noise model is adversarial. Our goal is to design an encoding and decoding scheme such that the encoding has a constant rate and the decoding recovers the encoded transmitted stream, or else aborts. Furthermore, we wish an “authentication” guarantee, that is, if the decoding scheme did not abort, it decodes the *correct* data with high probability (note that the probability that the scheme aborts potentially differs from the probability that the decoding scheme outputs incorrect data). The amount of recoverable data depends on the noise and the goal is to output (and authenticate) the longest possible prefix of the stream, given a constant corruption rate.

Definition 5.1. *A $(c(n), \gamma(n), \kappa(n))$ -Streaming Authentication Scheme with constant rate r is an encoding $e : \{0, 1\}^* \times \{0, 1\}^* \rightarrow \{0, 1\}^r$ that encodes a stream x_1, x_2, \dots into a stream $y_1 = e(x_1, R)$, $y_2 = e(x_1x_2, R)$, \dots , $y_i = e(x_1 \cdots x_i, R)$, and a decoding $d : \{0, 1\}^* \times \{0, 1\}^* \rightarrow \{0, 1\}^* \cup \{\perp\}$ such*

that the following holds. For any n , and for any adversary $\text{Adv}(x_1 \cdots x_n, y_1 \cdots y_n) = y'_1 \cdots y'_n$, either $d(y'_1 \cdots y'_n, R) = x'_1 x'_2 \cdots x'_n$ or $d(y'_1 \cdots y'_n, R) = \perp$, and if at most $c(n)$ transmissions were corrupted, then

1. the scheme aborts with probability at most $\kappa(n)$,

$$\Pr_R[d(y'_1 \cdots y'_n, R) = \perp] < \kappa(n).$$

2. if not aborted, the probability to decode an incorrect $\gamma(n)$ -prefix of the stream is at most $\kappa(n)$,

$$\Pr_R[d(y'_1 \cdots y'_n, R) \neq \perp \wedge x'_1 \cdots x'_{\gamma(n)} \neq x_1 \cdots x_{\gamma(n)}] < \kappa(n).$$

Eve is given both the raw stream and the channel transmissions, however she does not know the shared random string R used as the secret authentication key. It is desired that as long as Eve corrupts only a small fraction of the transmissions, Bob will be able to correctly decode a prefix of the stream, or otherwise be aware of the adversarial intervention and abort.

We show the following dichotomy: If the adversarial corruption rate is some constant c , then there exists a streaming authentication stream that decodes a prefix of at most $(1 - c)$ -fraction of the stream received so far. In addition, there does not exist a streaming authentication scheme that is capable of decoding a longer prefix with non-negligible probability.

Theorem 5.2. *In the shared-randomness model, for every constants c, ε such that $0 \leq c < 1$ and $0 < \varepsilon \leq (1 - c)/2$ there exists a constant-rate $(cn, (1 - c)n - \varepsilon n, 2^{-\Omega(n)})$ -Streaming Authentication Scheme. Moreover, there exists an efficient constant-rate $(cn, (1 - c)n - \varepsilon n, 2^{-\Omega(\log n)})$ -Streaming Authentication Scheme.*

For any constant $c_{th} > c$, if the adversarial corruption rate exceeds c_{th} , the schemes abort with overwhelming probability over the shared randomness.

Theorem 5.3. *Assume that a bitstream x_1, x_2, \dots is communicated using some encoding protocol with a constant rate, and assume that at time n the receiver decodes the bitstring x'_1, \dots, x'_n . If the rate of adversarial corruptions is $0 \leq c \leq 1$, then for any constant $\varepsilon > 0$,*

$$\Pr[x'_1 \cdots x'_{(1-c)n+\varepsilon n} = x_1 \cdots x_{(1-c)n+\varepsilon n}] \leq 2^{-\Omega(\varepsilon n)}$$

where the probability is over the coin tosses of the decoding algorithm, assuming $\{x_i\}$ are uniformly, independently distributed.

We now prove Theorem 5.3 and then construct the protocols guaranteed by Theorem 5.2

Proof. Consider an adversary that, starting at time $(1 - c)n$, corrupts all the transmissions. It is easy to verify that the corruption rate is c . Clearly, from time $(1 - c)n$ and on, the effective capacity of the channel is 0. This means that the decoder has no use of transmissions of times $\geq (1 - c)n$ and he decodes only using transmissions received up to time $(1 - c)n$. However, due to the streaming nature of the model, transmissions at times $< (1 - c)n$ depend only on $x_1, \dots, x_{(1-c)n}$ (the suffix of the stream is yet unknown to the sender). The receiver has no information about any bit x_i with $i > (1 - c)n$ and his best strategy is to guess them. The probability to correctly guess the last εn bits is at most $2^{-\lfloor \varepsilon n \rfloor}$. \square

In order to construct a streaming authentication scheme, we use two concatenated layers of online codes. The inner code is a Blueberry $B : [S + 1]^* \rightarrow [L + 1]^*$ code with constant S, L , and the outer code A is an online code that allows a prefix decoding in the presence of errors and erasures. The entire process can be described by

$$(x_1, \dots) \xrightarrow{A} (y_1, \dots) \xrightarrow{B} (z_1, \dots) \xrightarrow{\text{channel}} (\tilde{z}_1, \dots) \xrightarrow{B^{-1}} (\tilde{y}_1, \dots) \xrightarrow{A^{-1}} (\tilde{x}_1, \dots)$$

We begin with a simple and elegant construction which, although not efficient, demonstrates the power of the Blueberry code.

Proposition 5.4. *Let c, ε be constants $0 \leq c < 1$, $0 < \varepsilon \leq (1 - c)$ and let $A = \text{TCenc}()$ be an encoding using a binary tree code and B a Blueberry code with constant parameters determined by c, ε . The concatenation of A and B is a $(cn, (1 - c)n - \varepsilon n, 2^{-\Omega(n)})$ -streaming authentication scheme.*

Proof. Assume that in order to encode the bitstream x_1, x_2, \dots , we use a binary tree code over alphabet $[S + 1]$ with distance α to be determined later, concatenated with a Blueberry-code $B : [S + 1]^* \rightarrow [L + 1]^*$. We show that if at time n we decode a string $\tilde{x}_1 \cdots \tilde{x}_n$ whose prefix $\tilde{x}_1 \cdots \tilde{x}_{(1-c-\varepsilon)n}$ differs from $x_1 \cdots x_{(1-c-\varepsilon)n}$, then the corruption rate was larger than c .

For a specific time n , consider a string $\tilde{x} \in \{0, 1\}^n$, such that $\text{anc}(x, \tilde{x}) \geq (c + \varepsilon)n$. Due to the tree distance property, the Hamming distance between $\text{TCenc}(\tilde{x})$ and $\text{TCenc}(x)$ is at least $\alpha(c + \varepsilon)n$. Assume Eve causes d erasures and e errors, a maximal-likelihood decoding will prefer x over \tilde{x} as long as $\lfloor \alpha(c + \varepsilon)n \rfloor > 2e + d$.

If Eve's corruption rate is limited to c , Lemma 3.4 implies that with overwhelming probability at most $2cnS/L$ of these corruptions become errors and the rest are marked as erasures. Setting $\alpha > \frac{c}{c+\varepsilon}(1 + \frac{2S}{L})$ we guarantee that $\alpha(c + \varepsilon)n > 2 \cdot 2cnS/L + cn(1 - 2S/L)$,³ thus Bob decodes with overwhelming probability a string \tilde{x} such that $\text{anc}(x, \tilde{x}) < (c + \varepsilon)n$, as claimed.

Note that the actual fraction of adversarial corruptions can be estimated out of the number of erasures marked by the Blueberry code. We abort the decoding if at a specific time n the number of erasures exceeds cn . Lemma 3.4 guarantees that if the adversary corrupts more than a $c/(1 - \frac{2S}{L})$ -fraction of the transmissions, she will cause at least cn erasures, except with negligible probability. Choosing L such that $(1 - \frac{2S}{L}) \geq \frac{c}{c_{th}}$ completes the proof for the non-efficient case of Theorem 5.2. \square

We note that although in the above proof we require ε to be constant, for the case of $c = 0$ (i.e., when the channel is not inherently noisy) we can let ε be smaller. For instance, if we let $\varepsilon = \kappa/n$ for a security parameter κ , the scheme is comparable to a (non-streaming) authentication scheme with the same security parameter: in order to change even a single bit in a prefix of length n , after $n + \kappa$ symbols were transmitted, the adversary must change at least $\alpha\kappa/2$ transmissions, and will be caught except with probability $2^{-\Omega(\kappa)}$. Since the above holds for any time n , we get a perpetual authentication of the stream.

The case where $c > 0$ has a meaning of communicating over a noisy channel (regardless of the adversary). The users do not abort the authentication scheme although they know the message was changed by the channel. Instead, the scheme features both error-correction and authentication abilities and the parties succeed to recover (a prefix of) the original message with high probability.

³It is required to have $\alpha < 1$, thus the choice of (the constant) L should depend on ε and c , specifically, $L > 2S \frac{c}{\varepsilon}$. Also note that S depends on α , however L is independent of both. For a fixed value of α (and $S = d^{O(1/(1-\alpha))}$) there is always a way to choose a constant L that satisfies the conditions.

5.1 Efficient Streaming Authentication

We now complete the proof of Theorem 5.2 by defining an efficient randomized code A_{eff} for prefix-decoding in the presence of errors and erasures. The protocol partitions the stream into words of logarithmic size and encodes each using a tree code. At any time n , one of the $O(n/\log n)$ words is chosen at random and its next encoded symbol is transmitted. The value n increases as the protocol progresses which means that the length of each encoded word increases as well. This however causes no problem: each word is encoded by a tree code (rather than, say, a block code), which is performed in an online manner without assuming knowledge of the word's length. Decoding can be performed efficiently by an exhaustive search since each word is of logarithmic length in the current time n . We note that the parties hold the entire stream in their memory throughout the protocol, which is different from the common practice of streaming algorithms in which there is only a single party (rather than two) which aim to compute some statistics of the stream using poly-logarithmic memory.

Proposition 5.5. *For any constants $0 \leq c < 1$, $0 < \varepsilon \leq (1 - c)/2$ and a constant $c_1 > 0$, there exist efficient constant-rate encoding and decoding scheme such that, for any set of infinite strings $\{\mathbf{x}^1, \mathbf{x}^2, \dots\}$ the following holds for any sufficiently large time n except with polynomially small probability in n .*

If the corruption rate at time n is at most c then the scheme correctly decodes a prefix of length $c_1 \log n$ of each one of the strings \mathbf{x}^k with $k \in \{\lceil \frac{\varepsilon n}{4} \rceil, \dots, \lceil \frac{(1-c-\varepsilon)n}{\log(1-c-\varepsilon)n} \rceil\}$. Moreover, up to time n the encoding scheme assumes knowledge of only strings \mathbf{x}^k with $k \leq n/\log n$.

Let $0 \leq c < 1$ and $0 < \varepsilon < (1 - c)/2$ be fixed parameters of the protocol. Let c_0, c_1 be some constants which depend on c and ε . Let \mathcal{T} be a tree code over alphabet $[S + 1]$ with distance α to be set later.

A_{eff} Encoding: For every $k > 0$ set $\text{count}_k = 0$.

At any time $n > 1$, repeat the following process for $j = 1, 2, \dots, c_0$:

- (a) randomly choose $k \in \{1, \dots, \lfloor n/\log n \rfloor\}$.
- (b) set $\text{count}_k = \text{count}_k + 1$.
- (c) transmit $y_{n,j} \in [S + 1]$, the next symbol of the encoding of \mathbf{x}^k using \mathcal{T} , that is, the last symbol of

$$\text{TCenc}(\mathbf{x}_1^k \cdots \mathbf{x}_{\text{count}_k}^k) = \text{TCenc}(\mathbf{x}_1^k \cdots \mathbf{x}_{\text{count}_k - 1}^k) \circ y_{n,j}.$$

A_{eff} Decoding: For every $(i, j) \in \mathbb{N} \times [c_0]$ we denote by $\text{ID}(i, j)$ the identifier k of the string \mathbf{x}^k used at iteration (i, j) . For each time n , mark all the transmissions $y_{i,j}$ with $i < \varepsilon n/4$ as erasures, and decode \mathbf{x}^k for $\lceil \frac{\varepsilon n}{4} \rceil \leq k \leq \lceil \frac{(1-c-\varepsilon)n}{\log(1-c-\varepsilon)n} \rceil$:

let $Y_k = \{(i, j) \mid \text{ID}(i, j) = k\}$. Decode the received string indexed by Y_k . That is, set

$$\hat{\mathbf{x}}^k = \text{TCdec}(y_{|Y_k}),$$

where $y_{|Y_k}$ is the string given by concatenating all $y_{i,j}$ with $(i, j) \in Y_k$, where $y_{i,j}$ comes before $y_{i',j'}$ if $i < i'$ or $(i = i') \wedge (j < j')$. Consider a prefix of length $c_1 \log n$ of $\hat{\mathbf{x}}^k$ and ignore the rest.

Protocol 1: An efficient protocol for communicating a logarithmic prefix of $\{\mathbf{x}^1, \mathbf{x}^2, \dots\}$.

In Appendix A.1 we show that Protocol 1 concatenated with a Blueberry code $B : [S + 1]^* \rightarrow [L + 1]^*$ satisfies the requirements of Proposition 5.5. We show that with high probability, $\Theta(\log n)$ symbols of $\text{TCenc}(\mathbf{x}^k)$ are transmitted by time n for every k in the range $K_n \triangleq \{\lceil \frac{\varepsilon n}{4} \rceil, \dots, \lceil \frac{(1-c-\varepsilon)n}{\log(1-c-\varepsilon)n} \rceil\}$.

Moreover, at least a constant fraction of these transmissions were not corrupted by the adversary. Therefore, we can use Proposition 5.4 to decode a prefix of length $O(\log n)$ of each of the codewords indexed by K_n , with high probability.

The remaining hurdle is to split the stream x_1, x_2, \dots into words $\{\mathbf{x}^1, \mathbf{x}^2, \dots\}$. In Appendix A.2 we show how to construct the set $\{\mathbf{x}^1, \mathbf{x}^2, \dots\}$ such that for any time n , the entire prefix $x_1, \dots, x_{(1-c-\varepsilon)n}$ appears in a $c_1 \log n$ -prefix of strings $\{\mathbf{x}^k\}$ with $k \in K_n$. This gives an efficient $(cn, (1-c)n - \varepsilon n, 2^{-\Omega(\log n)})$ -authentication scheme and completes the proof of Theorem 5.2.

5.2 Extensions for Streaming Authentication

There are several possible extensions to the above results, which we briefly discuss here.

Efficient streaming authentication scheme with exponentially small error. It is possible to improve the efficient scheme of Theorem 5.2 so that it aborts with polynomially small probability, however, given that it did not abort, the probability that the decoded prefix is incorrect is *exponentially small*. More accurately, the ‘trust’ Bob has in the decoded string *increases* with the amount of received transmissions. Thus, except for the last fraction of the stream, the decoded stream is equal to the one sent by Alice with overwhelming probability.

Theorem 5.6. *For any $0 \leq c < 1$, $0 < \varepsilon \leq \frac{1}{2}(1-c)$ there exists an efficient $(cn, (1-c)n - \varepsilon n, 2^{-\Omega(\log n)})$ -streaming authentication protocol that, for any time n in which the decoding procedure did not abort, for any $1 \leq \ell \leq (1-c-\varepsilon)n$ it holds that*

$$\Pr[x'_\ell \neq x_\ell] < 2^{-\Omega(n)}.$$

See proof in Appendix B.

Decoding a prefix longer than $(1-c)n$. Although our scheme decodes a prefix of length at most $(1-c)n$ in the worst case, the successfully decoded prefix can be in fact longer. The worst case, as demonstrated by Theorem 5.3, happens when the adversary blocks the suffix of the transmitted stream. On the other hand, if the adversary blocks the prefix of the transmissions, then the scheme of Proposition 5.4 correctly decodes the entire stream! In fact, the protocol succeeds to decode the entire prefix for any time n that satisfies the following γ -suffix condition, if the tree distance satisfies $\alpha > \gamma$.

Definition 5.7. *For any constant $0 \leq \gamma < 1$, we say that time n satisfies the γ -suffix condition if any suffix $x_t \dots x_n$ has at most $\gamma(n-t)$ corrupted transmissions.*

Definition 5.8. *Let $c < 1$ and $\gamma \in (c, 1)$ be given. For any time n let $N_\gamma(n)$ be the latest index that satisfies the γ -suffix condition. When n is clear from the context, we denote $N_\gamma(n)$ simply as N_γ .*

The following Lemma guarantees that, for any $\gamma \in (c, 1)$ it holds that $(1-c/\gamma)n \leq N_\gamma(n) \leq n$.

Lemma 5.9. *For every corruption rate c and constant $1 < \xi < 1/c$ there exist a time $t > (1 - \frac{1}{\xi})n$ that satisfies the $c\xi$ -suffix condition.*

See proof in Appendix D.

For a corruption rate c and any $\varepsilon > 0$, and for any time n , if the decoding algorithm did not decode up to time n , then that time n did not satisfy the suffix condition for $\gamma = c/(c + \varepsilon)$ (see formal proof in Appendix D), but then, by Lemma 5.9, there must exist a time $N_\gamma > (1 - c - \varepsilon)n$ that satisfies the γ -suffix condition, and at that time the protocol correctly decoded the entire stream (up to time N_γ). Bob does not know the value of N_γ but he can estimate it by checking the number of erasures marked by the Blueberry code.

Proposition 5.10. *Bob can efficiently compute a (lower-bound) estimation N'_γ for N_γ , such that $N'_\gamma > (1 - c - \epsilon)n$ and*

$$\Pr[N'_\gamma > N_\gamma] < 2^{-\Omega(N'_\gamma - N_\gamma)}.$$

See appendix D for proof and discussion.

Reducing the amount of shared randomness. Our schemes rely on the fact that the parties share a secret random string whose length increases with the size of the information to be communicated. This assumption is sometimes not satisfied in practical applications, especially when considering a multiparty setting in which any two parties run a separate instance of the scheme.

We can mitigate the need for a long shared randomness if the adversary is assumed to be polynomial, assuming standard cryptographic assumptions (specifically, hardness of DDH). To this end, each user generates a pair (sk, pk) of a secret and a public key, broadcasts the public key pk and keeps sk secret. When two users initiate an authentication scheme instance, they first perform a Diffie-Hellman [DH76] key exchange and obtain an authentication key. They both use the authentication key as a seed to a pseudo-random-generator that generates a long random string for the authentication scheme. Under the DDH assumption, a polynomially-bounded adversary has only negligible information about the authentication key nor the generated randomness, and the authentication scheme remains secure. See Appendix C for more details and a proof sketch.

6 Interactive Communication

In this section we extend our discussion to the 2-way communication model of *interactive communication*. We show that for adversarial corruption rate of $1/2$ or higher, no constant-rate protocol can compute functions that require interaction between the parties, while with the usage of the Blueberry code we show how to construct a protocol for any function assuming adversarial corruption rate below $1/2$. We begin by defining the interactive communication model.

The interactive-communication model. Assume that Alice and Bob wish to compute some function $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Z}$, where Alice holds $x \in \mathcal{X}$ and Bob holds $y \in \mathcal{Y}$ in the shared-randomness model. The computation is performed interactively: at each round, both parties communicate a message which depends on their input and previous transmissions. At the end of the computation Alice outputs $z_A \in \mathcal{Z}$ and Bob outputs $z_B \in \mathcal{Z}$, and we say that f was correctly computed if $z_A = z_B = f(x, y)$. Without loss of generality we assume the output is a single bit, $|\mathcal{Z}| = 2$.

We show the following separation theorems,

Theorem 6.1. *For any function f which depends on both x and y , the following holds. If the adversarial corruption rate is $\frac{1}{2}$ or higher then no constant-rate interactive protocol correctly computes f with probability higher than the probability of guessing $f(x, y)$ given only the input x (or only the input y).*

Theorem 6.2. *For any constants $\epsilon > 0$ and for any function f and inputs x, y , there exists an interactive protocol with constant overhead such that if the adversarial corruption rate is at most $c = \frac{1}{2} - \epsilon$, the protocol outputs $f(x, y)$ with overwhelming probability over the shared random string R .*

Proof. (**Theorem 6.1**) Assume that the protocol takes T rounds. Furthermore, recall that in our model it is assumed that at each round both parties send exactly one message.⁴ Hence and without loss of generality, Alice is the sender of at most $T/2$ of the transmissions. Eve corrupts all the transmissions originated by Alice (causing at least an erasure in each one of these transmissions). Effectively, the unidirectional channel from Alice to Bob has a zero capacity, and it cannot be that Bob correctly computes $f(x, y)$ with probability higher than guessing $f(x, y)$ given only y . \square

It is interesting to note that if f only depends on *one* of its inputs, then only 1-way communication is required and $c = \frac{1}{2}$ is no longer a limit, as discussed in Section 5.

We now construct a protocol that correctly computes any $f(x, y)$ with overwhelming probability as long as the adversarial corruption rate is $\frac{1}{2} - \varepsilon$ for $\varepsilon > 0$. To this end, we concatenate an online protocol for computing $f(x, y)$ over an adversarial noisy channel with a Blueberry code.

Let us recall how to construct a constant-rate protocol for computing $f(x, y)$ over a noisy channel out of an interactive protocol π for the same task that assumes a noiseless channel [BR11]. We assume that π consists of T rounds in which Alice and Bob send a single bit according to their input and previous transmissions. Without loss of generality, we assume that Alice sends her bits at odd rounds while Bob transmits at even rounds. We can view the computation of π as a root-leaf walk along a binary tree in which odd levels correspond to Alice's messages and even levels to Bob's, see Figure 2.

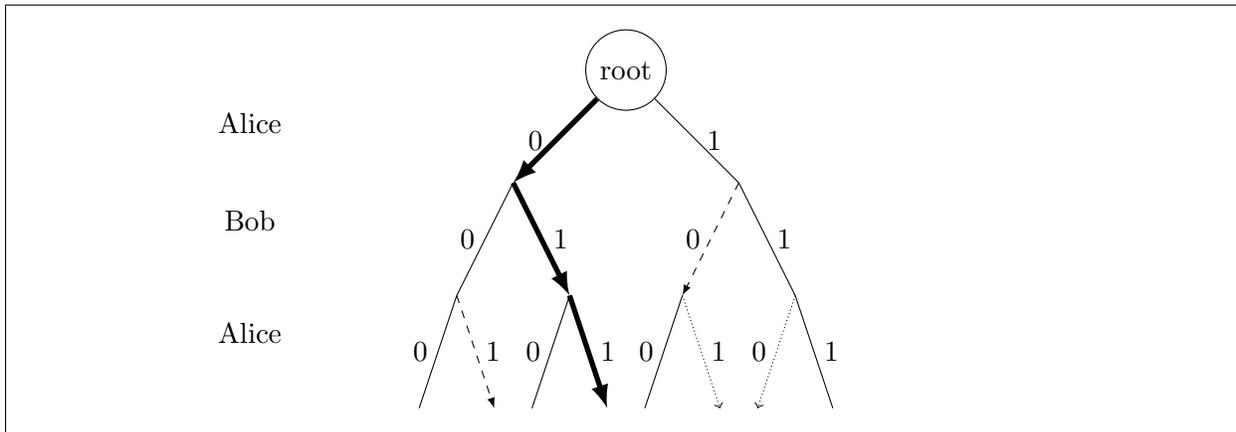


Figure 2: A π -tree showing the path P (bold edges) taken by Alice and Bob for computing $f(x, y)$. Dashed edges represent the hypothetical reply of Alice and Bob given that a different path P' was taken (when such replies are defined).

In order to obtain a protocol that withstands (a low rate of) channel noise, Alice and Bob *simulate* the construction of path P along the π -tree. The users transmit edges of P one by one, where each user transmits the next edge that extends the partial path transmitted so far.⁵ This process is repeated for $N = \lceil T/(1 - \alpha) \rceil$ times, for some constant $\alpha < 1$ to be set later. In [BR11] it

⁴The proof also holds for protocols for which there exists a function $\text{Next}(i)$ which defines, for each round i , which of the parties sends a message, and is independent of the messages sent by now (these kind of protocols are called *oblivious* in [BR11]). In that case there exists one party that communicates at most $T/2$ messages at rounds known to Eve in advance.

On the other hand, the proof does not hold for the most general model, in which the protocol adaptively determines who is next to speak, possibly according to the noise observed so far (so that the party that suffers from higher noise rate gets more transmission slots, etc.)

⁵ The users transmit \perp when they do not know how to extend P based on current information.

is shown that unless the noise rate exceeds $1/4$, after N rounds both parties will decode the entire path P . We refer the reader to [BR11] for a full description of the protocol and correctness proof. We now extend the analysis for the case of channels with errors and erasures.

To simplify the explanation, assume that each transmission is over the alphabet $\Sigma = \{0, \dots, N\} \times \{0, 1\}^{\leq 2}$. Intuitively, the transmission $(e, s) \in \Sigma$ means “extend the path P by taking at most two steps defined by s starting at the child of the edge I have transmitted at transmission number e ”. Although this alphabet is not of constant size, it is easy to obtain a constant size alphabet by encoding each (e, s) into a delimited binary string (see Section 6 in [BR11]). Each symbol (e, s) is communicated to the other side via a $|\Sigma|$ -ary tree code with distance α and alphabet Γ , that is at time n Alice sends $a_n \in \Gamma$, the last symbol of $\text{TCenc}((e, s)_1, \dots, (e, s)_n) = a_1 a_2 \dots a_n$, and Bob receives $\tilde{a}_n \in \Gamma \cup \{\perp\}$, possibly with added noise or an erasure mark (similarly, Bob sends $b_n \in \Gamma$, and Alice receives \tilde{b}_n). Let $\text{TCdec}(\tilde{a}_1, \dots, \tilde{a}_n)$ denote the string Bob decodes at time n (similarly, Alice decodes $\text{TCdec}(\tilde{b}_1, \dots, \tilde{b}_n)$). For every $i > 0$, we denote with $m(i)$ the largest number such that the first $m(i)$ symbols of $\text{TCdec}(\tilde{a}_1, \dots, \tilde{a}_i)$ equal to $a_1, \dots, a_{m(i)}$ and the first $m(i)$ symbols of $\text{TCdec}(\tilde{b}_1, \dots, \tilde{b}_i)$ equal to $b_1, \dots, b_{m(i)}$.

Define $\mathcal{N}(i, j)$ to be the “effective” number of adversarial corruptions in interval $[i, j]$: the number of erasures plus twice the number of errors in the $[i, j]$ interval of the simulation (for both users).

Definition 6.3. Let $\mathcal{N}_a(i, j) = |\{k \mid i \leq k \leq j, \tilde{a}_k = \perp\}| + 2|\{k \mid i \leq k \leq j, \tilde{a}_k \notin \{a_k, \perp\}\}|$, and similarly define $\mathcal{N}_b(i, j)$. The effective number of corruptions in interval $[i, j]$ is $\mathcal{N}(i, j) = \mathcal{N}_a(i, j) + \mathcal{N}_b(i, j)$.

We begin by showing that if $m(i) < i$ then many corruptions must have happened in the interval $[m(i) + 1, i]$.

Lemma 6.4. $\mathcal{N}(m(i) + 1, i) \geq \alpha(i - m(i))$.

Proof. Assume that at time i Bob decodes the string a'_1, \dots, a'_i . By the definition of $m(i)$, $a'_1, \dots, a'_{m(i)} = a_1, \dots, a_{m(i)}$, and assume without loss of generality that $a'_{m(i)+1} \neq a_{m(i)+1}$. Note that the Hamming distance between $\text{TCenc}(a_1, \dots, a_i)$ and $\text{TCenc}(a'_1, \dots, a'_i)$ must be at least $\alpha(i - m(i))$. It is immediate that for Bob to make such a decoding error, $\mathcal{N}_a \geq \alpha(i - m(i))$. \square

Lemma 6.5 ([BR11]). Let $t(i)$ be the earliest time such that both users announced the first i edges of P within their transmissions. For $i \geq 0$, $k \geq 1$, if $t(k) > i + 1$, then $t(k - 1) > m(i)$.

Proof. The proof is taken from [BR11]: Without loss of generality, assume that the k th edge of P describes Alice’s move. Suppose $t(k - 1) \leq m(i)$ and $t(k) > i + 1$. Then it must be the case that the first $k - 1$ edges of P have already been announced within the first $m(i)$ transmissions of both parties, yet the k th edge has not. By the protocol definition, Alice will announce this edge at round $i + 1$, in contradiction to our assumption that $t(k) > i + 1$. \square

Next we show that if at some time i the length of the proposed P is not long enough (less than k), then many transmissions must have been corrupted.

Lemma 6.6. For $i \geq -1$, $k \geq 0$, if $t(k) > i + 1$, then $\mathcal{N}(1, i) \geq \alpha(i - k + 1)$.

Proof. We prove by induction. The claim vacuously holds for $k = 0$ and trivially holds for $i \leq 0$ since $\mathcal{N}(1, i)$ is non-negative. Otherwise, we have

$$\mathcal{N}(1, i) = \mathcal{N}(1, m(i)) + \mathcal{N}(m(i) + 1, i).$$

The second term, by Lemma 6.4 gives $\mathcal{N}(m(i) + 1, i) \geq \alpha(i - m(i))$. For the first term, Lemma 6.5 suggests that $t(k - 1) > m(i)$ and we can use the inductive hypothesis with $i' = m(i) - 1$ and $k' = k - 1$ to get

$$\mathcal{N}(1, m(i)) \geq \mathcal{N}(1, i') \geq \alpha(i' - k' + 1) = \alpha(m(i) - k + 1).$$

□

The above Lemmas allow us to complete the proof of Theorem 6.2 by showing that if the simulation of P failed, there must have been “too many” corruptions.

Proof. (Theorem 6.2) Assume an unsuccessful run of the simulation protocol. That is, the simulation of the path P has failed, $m(N) < t(T)$. The number of adversarial corruptions throughout the protocol is given by $\mathcal{N}(1, N) = \mathcal{N}(1, m(N)) + \mathcal{N}(m(N) + 1, N)$ which by Lemma 6.6 and Lemma 6.4 is lower bounded by

$$\mathcal{N}(1, N) \geq \alpha(m(N) - T + 1) + \alpha(N - m(N)) \geq \alpha(N - T) \geq \alpha(N - (1 - \alpha)N) = \alpha^2 N.$$

Yet, assume the adversary is restricted to corrupt at most $c = 1/2 - \varepsilon$ fraction of the $2N = 2\lceil \frac{T}{1-\alpha} \rceil$ transmissions, then Lemma 3.4 guarantees that with overwhelming probability there will be at least $2cN(1 - 2S/L)$ erasures, assuming a Blueberry code $B : [S + 1]^* \rightarrow [L + 1]^*$ (here, $S + 1 = |\Gamma|$). This implies that with overwhelming probability $\mathcal{N}(1, N) \leq 2cN(1 + 2S/L)$. For any $0 < \varepsilon \leq 1/2$ we can choose constants $\alpha < 1$ and $L > S$ such that $\alpha^2 > (1 - 2\varepsilon)(1 + 2S/L)$ and conclude that the protocol succeeds with overwhelming probability over the shared randomness. □

References

- [BK12] Zvika Brakerski and Yael Tauman Kalai. Efficient interactive coding against adversarial noise. *Foundations of Computer Science, IEEE Annual Symposium on*, pages 160–166, 2012.
- [BN13] Zvika Brakerski and Moni Naor. Fast algorithms for interactive coding. In *Proceedings of the 24th Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '13*, pages 443–456, 2013.
- [Bon63] C. Bonsall. *The case of the hungry stranger*. HarperCollins, 1963.
- [BR11] Mark Braverman and Anup Rao. Towards coding for maximum errors in interactive communication. In *Proceedings of the 43rd annual ACM symposium on Theory of computing, STOC '11*, pages 159–166, New York, NY, USA, 2011. ACM.
- [Bra12] Mark Braverman. Towards deterministic tree code constructions. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pages 161–167. ACM, 2012.
- [DH76] W. Diffie and M. Hellman. New directions in cryptography. *Information Theory, IEEE Transactions on*, 22(6):644 – 654, nov 1976.
- [EGM90] Shimon Even, Oded Goldreich, and Silvio Micali. On-line/off-line digital signatures. In Gilles Brassard, editor, *Advances in Cryptology — CRYPTO '89 Proceedings*, volume 435 of *Lecture Notes in Computer Science*, pages 263–275. Springer Berlin / Heidelberg, 1990.

- [GM01] P. Golle and N. Modadugu. Authenticating streamed data in the presence of random packet loss. In *ISOC Network and Distributed System Security Symposium, NDSS'01*, 2001.
- [GMS11] Ran Gelles, Ankur Moitra, and Amit Sahai. Efficient and explicit coding for interactive communication. In *Foundations of Computer Science (FOCS), 2011 IEEE 52nd Annual Symposium on*, pages 768–777, oct. 2011.
- [Gol04] Oded Goldreich. *Foundations of cryptography. Vol II: Basic applications*. Cambridge University Press, New York, 2004.
- [GR97] Rosario Gennaro and Pankaj Rohatgi. How to sign digital streams. In Burton Kaliski, editor, *Advances in Cryptology — CRYPTO '97*, volume 1294 of *Lecture Notes in Computer Science*, pages 180–197. Springer Berlin / Heidelberg, 1997.
- [GS10] Venkatesan Guruswami and Adam Smith. Codes for computationally simple channels: Explicit constructions with optimal rate. *Foundations of Computer Science, IEEE Annual Symposium on*, pages 723–732, 2010.
- [HO08] Brett Hemenway and Rafail Ostrovsky. Public-key locally-decodable codes. In David Wagner, editor, *Advances in Cryptology – CRYPTO 2008*, volume 5157 of *Lecture Notes in Computer Science*, pages 126–143. Springer Berlin / Heidelberg, 2008.
- [HOSW11] Brett Hemenway, Rafail Ostrovsky, Martin Strauss, and Mary Wootters. Public key locally decodable codes with short keys. In Leslie Goldberg, Klaus Jansen, R. Ravi, and José Rolim, editors, *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, volume 6845 of *Lecture Notes in Computer Science*, pages 605–615. Springer Berlin / Heidelberg, 2011.
- [Lan04] Michael Langberg. Private codes or succinct random codes that are (almost) perfect. In *Proceedings of the 45th Annual IEEE Symposium on Foundations of Computer Science, FOCS '04*, pages 325–334, Washington, DC, USA, 2004. IEEE Computer Society.
- [Lip94] Richard Lipton. A new approach to information theory. In Patrice Enjalbert, Ernst Mayr, and Klaus Wagner, editors, *STACS 94*, volume 775 of *Lecture Notes in Computer Science*, pages 699–708. Springer Berlin / Heidelberg, 1994.
- [MPSW05] Silvio Micali, Chris Peikert, Madhu Sudan, and David Wilson. Optimal error correction against computationally bounded noise. In Joe Kilian, editor, *Theory of Cryptography*, volume 3378 of *Lecture Notes in Computer Science*, pages 1–16. Springer Berlin / Heidelberg, 2005.
- [MS01] S. Miner and J. Staddon. Graph-based authentication of digital streams. In *Security and Privacy, 2001, IEEE Symposium on*, pages 232–246, 2001.
- [PCTS00] A. Perrig, R. Canetti, J.D. Tygar, and Dawn Song. Efficient authentication and signing of multicast streams over lossy channels. In *Security and Privacy, 2000, IEEE Symposium on*, pages 56–73, 2000.
- [Pec06] Marcin Pezarski. An improvement of the tree code construction. *Information Processing Letters*, 99(3):92–95, 2006.

- [Sch92] Leonard J. Schulman. Communication on noisy channels: a coding theorem for computation. *Foundations of Computer Science, Annual IEEE Symposium on*, 0:724–733, 1992.
- [Sch93] Leonard J. Schulman. Deterministic coding for interactive communication. In *STOC '93: Proceedings of the twenty-fifth annual ACM symposium on Theory of computing*, pages 747–756, New York, NY, USA, 1993. ACM.
- [Sch96] Leonard J. Schulman. Coding for interactive communication. *IEEE Transactions on Information Theory*, 42(6):1745–1756, 1996.
- [Sha58] Claude E. Shannon. A note on a partial ordering for communication channels. *Information and Control*, 1(4):390 – 397, 1958.
- [Smi07] Adam Smith. Scrambling adversarial errors using few random bits, optimal information reconciliation, and better private codes. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms, SODA '07*, pages 395–404, Philadelphia, PA, USA, 2007. Society for Industrial and Applied Mathematics.

APPENDIX

A Construction of an Efficient Authentication Scheme

A.1 Proof of Proposition 5.5

In this appendix we show that the concatenation of Protocol 1 with a Blueberry code $[S+1]^* \rightarrow [L+1]^*$ satisfies the conditions of Proposition 5.5. We begin by showing that for any time n , if we look at k 's which are not too close to the start or to the end, that is, k 's in $K_n = \{\lceil \frac{\varepsilon n/4}{\log \varepsilon n/4} \rceil, \dots, \lceil \frac{(1-c-\varepsilon)n}{\log(1-c-\varepsilon)n} \rceil\}$, then every string \mathbf{x}^k is selected by the encoding scheme $\Theta(\log n)$ times in expectation. In addition, a constant fraction of a specific tree's transmissions is received intact, while the expected number of errors is a smaller fraction, controlled by L . Therefore, a logarithmic prefix of the string \mathbf{x}^k can be decoded with high probability.

Lemma A.1. *Let c, ε be given. For a given time n and for every $k \in K_n$,*

1. *the expected number of transmissions (i, j) with $\text{ID}(i, j) = k$ is $\Theta(c_0 \log n)$.*
2. *if the corruption rate is at most c , then for transmissions with $\text{ID}(i, j) = k$, the expected number of transmissions not corrupted by the adversary is $\Theta(c_0 \log n)$ and the expected number of errors is $\Theta(c_0 \log n/L)$.*

Proof. Fix a $k \in K_n$, and recall that $Y_k = \{(i, j) \mid \text{ID}(i, j) = k\}$. It is easy to verify that

$$\Pr((i, j) \in Y_k) = \begin{cases} 0 & i/\log i < k \\ \frac{1}{\lceil \frac{i}{\log i} \rceil} & i/\log i \geq k \end{cases} ,$$

where the probability is over the shared randomness R . Assume that the channel's (Eve's) noise pattern is $P = p_1, \dots, p_{cn}$, with $p_i \in [n] \times [c_0]$. First let us bound the number of erroneous transmissions of symbols from Y_k . For a specific instance of the scheme, let

$$\text{ERR}_k = \left| \{(i, j) \in Y_k \mid \tilde{y}^{i,j} \text{ is an error}\} \right| .$$

We are interested in the (adversarial) corruption pattern that maximizes the expected number of these errors,

$$\max_P \mathbb{E} [\text{ERR}_k].$$

Since the decoding process ignores the first $\varepsilon n/4$ transmissions, and since the expected number of errors is a $\frac{S}{L}$ -fraction of the corrupted transmissions in Y_k , this equals to

$$\max_P \mathbb{E} \left[\frac{S}{L} \left| \{(i, j) \mid \text{ID}(i, j) = k, i > \varepsilon n/4 \text{ and } \tilde{y}^{i,j} \neq y^{i,j}\} \right| \right],$$

where the expectation is over the shared randomness R . Note that $\Pr((i, j) \in Y_k)$ is monotonically decreasing for $i \geq t(k)$, and zero otherwise. The pattern P that maximizes Eve's probability to hit transmissions in Y_k is $P_\circ \triangleq \{t(k), t(k) + 1, \dots, t(k) + cn\}$. However, the decoding algorithm ignores the first $\varepsilon n/4$ indices and Eve has no use in attacking them. Therefore, if $P_\circ \cap [\varepsilon n/4] \neq \emptyset$, Eve's best strategy is to shift her attack to the window $P = \{\lfloor \varepsilon n/4 \rfloor, \lfloor \varepsilon n/4 \rfloor + 1, \dots, \lceil (\varepsilon/4 + c)n \rceil\}$.

$$\begin{aligned} \max_P \mathbb{E} [\text{ERR}_k] &= \max_P \frac{S}{L} \sum_{(i,j) \in P} \Pr((i, j) \in Y_k) \leq \frac{c_0 S}{L} \sum_{i=t(k)}^{cn+t(k)} \frac{1}{\log i - 1} \leq \frac{c_0 S}{L} \sum_{i=t(k)}^{cn+t(k)} \frac{\log i}{i - \log i} \\ &< \frac{c_0 S \log n}{L} \sum_{i=\varepsilon n/4}^n \frac{1}{i - \log i} < \frac{c_0 S \log n}{L} \sum_{i=\varepsilon n/4}^n \frac{1}{c'' i} = \frac{c_0 S \log n}{c'' L} (H_n - H_{\varepsilon n/4}) \end{aligned}$$

where the first inequality on the second line applies to both cases of empty and non-empty $P_\circ \cap [\varepsilon n/4]$, and $c'' < 1$ is some constant such that $c'' i \leq i - \log i$ for $i \geq \varepsilon n/4$, for a sufficiently large n . H_n is the n th Harmonic number, and it holds that $0 < H_n - \ln(n) < 1$. We get $\max_P \mathbb{E} [\text{ERR}_k] = O(c_0 \log n/L)$.

On the other hand, we can lower bound the amount of uncorrupt transmissions in Y_k . In a similar way to the above we define $\text{INTACT}_k = |\{(i, j) \in Y_k \mid \tilde{y}^{i,j} = y^{i,j}\}|$, and wish to lower bound the quantity $\min_P \mathbb{E} [\text{INTACT}_k]$. It is easy to verify that Eve's strategy from above is optimal for this case as well, thus

$$\begin{aligned} \min_P \mathbb{E} [\text{INTACT}_k] &\geq \min_P \sum_{\substack{(i,j) \notin P \\ i > \varepsilon n/4}} \Pr((i, j) \in Y_k) \geq \sum_{i=(c+\varepsilon/4)n+t(k)}^n \frac{c_0 \log i}{i} \\ &\geq c_0 ((1 - c - \varepsilon/4)n - t(k)) \frac{\log((c + \varepsilon/4)n + t(k))}{n}, \end{aligned}$$

which is $\Omega(c_0 \log n)$ for $t(k) \leq (1 - c - \varepsilon)n$, hence the claim holds for $k \leq (1 - c - \varepsilon)n / \log(1 - c - \varepsilon)n$.

Finally, define $\text{TOTAL}_k = |Y_k|$ to be the total amount of transmissions with $\text{ID}(i, j) = k$ (erasures, errors, and intact). The expected amount of this quantity is at least

$$\mathbb{E}[\text{TOTAL}_k] \leq c_0 \sum_{i=\varepsilon n/4}^n \frac{\log i}{i - \log i} \leq c_0 \sum_{\varepsilon n/4}^n \frac{\log i}{c'' i} \leq \frac{c_0 \log n}{c''} (H_n - H_{\varepsilon n/4}) = O(c_0 \log n)$$

with some small constant $c'' < 1$ for a sufficiently large n . The sum begins from $\varepsilon n/4$ since \mathbf{x}^k is declared only at time $t(k) \geq \varepsilon n/4$, if $k \geq \varepsilon n/4 \log(\varepsilon n/4)$. Since the number of intact transmissions is $\Omega(c_0 \log n)$, the total amount of transmissions is lower-bounded by the same quantity, thus $\mathbb{E}[\text{TOTAL}_k] = \Theta(c_0 \log n)$. \square

Lemma A.2. *Let c, ε be given. If, for time n , at most cn of the transmissions were corrupted, then there exist constants c_0, L and a constant distance α such that for every constant $c_1 > 0$, the first $c_1 \log n$ elements of any \mathbf{x}^k with $k \in \{\lceil \frac{\varepsilon n/4}{\log \varepsilon n/4} \rceil, \dots, \lceil \frac{(1-c-\varepsilon)n}{\log(1-c-\varepsilon)n} \rceil\}$ are correctly decoded with polynomial computational effort, except with polynomially small probability over the shared randomness.*

Proof. At time n , assume $\max_{k \in K_n} \mathbb{E}[\text{TOTAL}_k] < C_T \log n$ and $\min_{k \in K_n} \mathbb{E}[\text{INTACT}_k] > C_I \log n$, and define $\beta = C_I/C_T$. Note that β is independent of n and c_1 . Fix $k \in K_n$. Denote by \mathcal{BAD}_1 the event that there were too many erasures and errors for the k th codeword, i.e. $\text{INTACT}_k/\text{TOTAL}_k < \beta/2$, and by \mathcal{BAD}_2 the event that there were not enough transmissions for the k th codeword, $\text{TOTAL}_k < \frac{2c_1}{\beta} \log n$.

By an appropriate choice of $c_0 = O(c, \varepsilon, c_1, 1/\beta)$, we can bound the probability of any bad event to be polynomially small. For large enough c_0 we can assure that $\mathbb{E}[\text{TOTAL}_k] > \frac{4c_1}{\beta} \log n$, and thus by Chernoff, $\Pr[\mathcal{BAD}_2] \leq 2^{-\Omega(\log n)}$. Furthermore, a union bound gives

$$\Pr[\mathcal{BAD}_1] < \Pr[\text{TOTAL}_k > \sqrt{2}C_T \log n] + \Pr[\text{INTACT}_k < \frac{1}{\sqrt{2}}C_I \log n],$$

and by Chernoff inequality,

$$\Pr[\mathcal{BAD}_1] \leq 2^{-\Omega(\mathbb{E}[\text{TOTAL}_k])} + 2^{-\Omega(\mathbb{E}[\text{INTACT}_k])} = 2^{-\Omega(\log n)}.$$

Conditioned on the fact that \mathcal{BAD}_1 and \mathcal{BAD}_2 did not occur, we know that $n^* > \frac{2c_1}{\beta} \log n$ symbols of $\text{TCenc}(\mathbf{x}^k)$ were transmitted, and the adversary has corrupted at most $c^* < (1 - \beta/2)$ fraction of these transmissions. Proposition 5.4 suggests that for an appropriate choice of constant L , we are able to decode a prefix of length at least $\approx (1 - c^*)n^* = c_1 \log n$ except with probability $\exp(-\Omega(n^*)) = \exp(-\Omega(\log n))$. A union bound over all the possible $k \in K_n$ completes the proof.

As for the efficiency, since each codeword is of length $O(\log n)$, decoding via exhaustive search can be performed with polynomial computational effort. Hence, it is easy to verify that both the encoding and decoding can be done efficiently. \square

A.2 Construction of $\{\mathbf{x}^1, \mathbf{x}^2, \dots\}$

For every k , define \mathbf{x}^k to be the string that contains the stream prefix $x_{t(k)}$ down to x_1 concatenated with as many zeros as needed, $\mathbf{x}^k = x_{t(k)}x_{t(k)-1} \dots x_2x_1000 \dots$, where $t(k)$ is defined to be the minimal time such that $t(k)/\log t(k) > k$. We say \mathbf{x}^k is *declared* at time $t(k)$, meaning that only from this time and on the algorithm may choose to send symbols of the encoding of \mathbf{x}^k . It is easy to verify that the string \mathbf{x}^k is well defined at the time it is declared (the corresponding x_i 's are known).

If some string \mathbf{x}^k is declared at time $t(k)$ then \mathbf{x}^{k+1} will be declared at time $t(k+1) \approx t(k) + \log t(k) + O(\log \log t(k))$. By setting $c_1 = 2$ we are guaranteed that, for every $\varepsilon n/4 \leq \ell \leq (1-c-\varepsilon)n$, x_ℓ appears in a correctly decoded $c_1 \log n$ -prefix of some \mathbf{x}^k with $k \in K_n$.

Lemma A.3. *If \mathbf{x}^k is the latest string declared at time $i > 8$, then \mathbf{x}^{k+1} is declared at time sooner than $i + 2 \log i$.*

Proof. Let $f(i) = \frac{i+2 \log i}{\log(i+2 \log i)} - \frac{i}{\log i}$. f is monotonically increasing, and $f(8) > 1$. \square

Corollary A.4. *For any time $n > 8$, and any ℓ , the bit x_ℓ is within the first $2 \log n$ symbols of $\mathbf{x}^{\lceil \ell/\log \ell \rceil}$. Hence, every x_ℓ with $\varepsilon n/4 \leq \ell \leq (1-c-\varepsilon)n$, appears in a $2 \log n$ -prefix of (at least) one of the strings $\{\mathbf{x}^k\}_{k \in K_n}$.*

Unfortunately, with the above choice of \mathbf{x}^k s, only part of the stream, namely $x_{\varepsilon n/4}, \dots, x_{(1-c-\varepsilon)n}$, is decoded by the protocol. In order to communicate the prefix $x_1, \dots, x_{\varepsilon n/4}$ we run another instance of the scheme guaranteed by Proposition 5.5 for the following set of infinite strings $\{\mathbf{v}^1, \mathbf{v}^2, \dots\}$. (We explain how to combine these two instances below). Define \mathbf{v}^k in the following way

$$\mathbf{v}_i^k = \begin{cases} x_1 & k = 1, \forall i \\ x_{1+(\ell \bmod \lceil t(k)/2 \rceil + 1)} & k > 1, i = 1 \text{ and } \mathbf{v}_{2^{\log t(k-1)}}^{k-1} = x_\ell \\ x_{1+(\ell \bmod \lceil t(k)/2 \rceil + 1)} & k > 1, i > 1 \text{ and } \mathbf{v}_{i-1}^k = x_\ell \end{cases}$$

It is easy to verify that at time n , the string $\mathbf{v}^{\lfloor n/\log n \rfloor}$ is well defined and known to the encoder.

Lemma A.5. *For every time $n > 256/(1-c-\varepsilon)$, any bit x_ℓ with $1 \leq \ell \leq \varepsilon n/4$ appears in a $2 \log n$ -prefix of (at least) one of the strings $\{\mathbf{v}^k\}_{k \in K_n}$.*

Proof. Note that the concatenation of $O(\log n)$ -prefix of the \mathbf{v}^k s gives a string of the form $V \triangleq x_1 x_2 \dots x_{\lceil t(k_1)/2 \rceil} x_1 x_2 \dots x_{\lceil t(k_2)/2 \rceil} x_1 x_2 \dots$, and V is decoded by Protocol 1 with high probability.⁶ By taking $c_1 = 2$ and recalling that $\varepsilon < (1-c)/2$, (and thus, $(1-c-\varepsilon)n/4 > \varepsilon n/4$) the length of V is lower bounded by the amount of indices in prefixes of size $2 \log \frac{1}{4}(1-c-\varepsilon)n$ of $\{\mathbf{v}^{(1-c-\varepsilon)n/4}, \dots, \mathbf{v}^{(1-c-\varepsilon)n}\}$,

$$\begin{aligned} 2 \log \frac{1}{4}(1-c-\varepsilon)n \left(\frac{(1-c-\varepsilon)n}{\log(1-c-\varepsilon)n} - \frac{\frac{1}{4}(1-c-\varepsilon)n}{\log \frac{1}{4}(1-c-\varepsilon)n} \right) &\geq \frac{3}{2}(1-c-\varepsilon)n - 4 \frac{(1-c-\varepsilon)n}{\log(1-c-\varepsilon)n} \\ &\geq (1-c-\varepsilon)n \end{aligned}$$

where the last inequality holds for $n > \frac{256}{1-c-\varepsilon}$. Consider the latest place in V where x_1 appears. If that place is at least $(1-c-\varepsilon)/4$ indices from the end of V , it is clear that $x_1 \dots x_{(1-c-\varepsilon)/4}$ appears in the $(1-c-\varepsilon)/4$ -suffix of the decoded V . For the other case, let the bit that precedes this x_1 be x_ℓ . By the way we defined \mathbf{v}^k it follows that $\frac{3}{8}(1-c-\varepsilon) \leq \ell \leq \frac{1}{2}(1-c-\varepsilon)$ which means that $x_1 \dots x_{(1-c-\varepsilon)/4}$ must appear in a prefix of size $3/4 \cdot (1-c-\varepsilon)n$ of V . Since $(1-c-\varepsilon)n/4 > \varepsilon n/4$, the claim holds. \square

One cannot run Protocol 1 twice, once for $\{\mathbf{x}\}$ and once for $\{\mathbf{v}\}$. Indeed, Eve can block all the transmissions of one of the instances, thus prevent the correct decoding of the stream with probability one, while her corruption rate does not exceed $c = 1/2$. One possible solution is to set $c_1 = 4$ and interleave the transmitted data, that is, define the set $\{\mathbf{z}^1, \mathbf{z}^2, \dots\}$ where $\mathbf{z}^k = \mathbf{x}_1^k \mathbf{v}_1^k \mathbf{x}_2^k \mathbf{v}_2^k \dots$, etc.

Corollary A.6. *Let c, ε be constants $0 \leq c < 1$, $0 < \varepsilon \leq (1-c)/2$, and let B be a Blueberry code with constant parameters determined by c, ε . For the strings $\{\mathbf{z}^1, \mathbf{z}^2, \dots\}$ defined above, the concatenation of A_{eff} with B is an efficient $(cn, (1-c)n - \varepsilon n, 2^{-\Omega(\log n)})$ -streaming authentication scheme.*

B Efficient Authentication Scheme with Exponentially Small Error

In this section we show how to change the scheme given by Proposition 5.5 such that the error probability is exponentially small. That is, for any time n the decoding scheme aborts with at most

⁶To be more accurate, V is a substring of the string decoded by the scheme.

polynomially small probability in n , but on the event that the scheme did not abort, the decoded string matches the stream sent by Alice with overwhelming probability $1 - 2^{-\Omega(n)}$.

To this end, we add a parallel transmission of random hash values of the entire stream (up to time n), where the hash length is logarithmic in n .⁷ More formally, define an additional set of infinite strings $\{\mathbf{h}^1, \mathbf{h}^2, \dots\}$. We identify a string $a = a_1 a_2 \dots a_n$ with the n -dimensional vector (a_1, a_2, \dots, a_n) , and define \mathbf{h}^k in the following way. Randomly pick a matrix $R_k \in \{0, 1\}^{\log t(k) \times t(k)}$ and a vector $V_k \in \{0, 1\}^{\log t(k)}$ and set $\mathbf{h}^k = R_k \cdot (x_1, x_2, \dots, x_{t(k)})^T + V_k$, concatenated with as many zeroes as needed. The strings $\{\mathbf{h}^k\}$ are interleaved with the strings $\{\mathbf{x}^k\}$ and $\{\mathbf{v}^k\}$, and c_1 increases as explained in Appendix A.2.

Proposition 5.5 guarantees that except with polynomially small probability in n all the hash values $\{\mathbf{h}^k\}_{k \in K_n}$ are correctly decoded. The decoding at time n aborts if any of the hash values $\{\mathbf{h}^k\}_{k \in K_n}$ mismatch the corresponding prefix $x_1 \dots x_{t(k)}$.

Proposition B.1. *Given that Protocol 1 with hash testing did not abort, let x' be the decoded stream, then for every $\ell \leq (1 - c - \varepsilon)n$,*

$$\Pr_R[x'_\ell \neq x_\ell] < 2^{-\Omega(\max\{(1-c-\varepsilon)n-\ell, \log n\})}.$$

Proof. Eve is oblivious of R_i and V_i , thus for any two vectors $\tilde{\mathbf{x}} \in \{0, 1\}^{t(i)}$, $\tilde{\mathbf{h}} \in \{0, 1\}^{\log t(i)}$ chosen by Eve, $\Pr_{R_i}[\tilde{\mathbf{h}} = R_i \cdot \tilde{\mathbf{x}}^T + V_i] < 2^{-\log t(i)}$.

Clearly, the smaller ℓ is, the more hash values that are checked to be consistent with the decoded x'_ℓ . For $\ell > \varepsilon n/4$ there are at least $((1 - c - \varepsilon)n - \ell)/2 \log n$ independent hash values of stream prefixes longer than ℓ , where the smallest hash length is $\approx \log(\varepsilon n/4)$. Hence, the probability that $x'_\ell \neq x_\ell$ yet the decoding procedure did not abort is at most

$$2^{-\Omega\left(\log(\varepsilon n/4) \frac{(1-c-\varepsilon)n-\ell}{2 \log n}\right)} = 2^{-\Omega((1-c-\varepsilon)n-\ell)}.$$

Clearly, for $\ell < \varepsilon n/4$ there are as many hash tests as for $\ell = \varepsilon n/4$, thus the probability to incorrectly decode x_ℓ with $\ell < \varepsilon n/4$ is exponentially small in n as well. Finally, for the case where $((1-c-\varepsilon)n-\ell) < \log n$ we note that at least one hash value must be consistent, $\mathbf{h}^{\lceil (1-c-\varepsilon)n/\log(1-c-\varepsilon)n \rceil}$. The probability to incorrectly decode x_ℓ and pass the hash check is at most $2^{-\Omega(\log n)}$, which completes the proof. \square

The proof of Theorem 5.6 is immediate from the above Proposition.

Proof. (Theorem 5.6) Let c, ε be fixed. Perform Protocol 1 with hash testing with parameters $c, \varepsilon' = \varepsilon/2$. By Proposition B.1, every decoded x'_ℓ with $\ell \leq (1 - c - \varepsilon)n$ satisfies

$$\Pr[x'_\ell \neq x_\ell] < 2^{-\Omega((1-c-\varepsilon')n-\ell)} \leq 2^{-\Omega(\varepsilon n/2)}.$$

\square

C Computationally Secure Perpetual Authentication

For this section we assume basic knowledge of cryptographic primitives and assumptions.

Multi-Party Shared-Key Setting: The m -party shared-key setting for $(c(n), \gamma(n), \kappa(n))$ -Streaming Authentication is as follows. There are m parties, and parties i and j have shared random string

⁷This method is similar to the classic authentication method of splitting a stream into chunks of size $\log n$ and adding a MAC of logarithmic size after each chunk.

$R_{i,j}$, for all $1 \leq i, j \leq m$. The adversary may corrupt any number of parties, and learn all of their shared random strings. The adversary may invoke any sequence of streaming authentication sessions involving parties of his choice (for sender and receiver), messages of his choice, and transmission lengths of his choice. If the adversary has corrupted the sender, then he may violate the encoding procedure in an arbitrary manner. The actions of the adversary may be adaptive (e.g., corrupt some additional parties after invoking some streaming authentication sessions, etc.), and interleaved (e.g., corrupt some additional parties before some streaming authentication sessions have finished).

Definition C.1. A $(c(n), \gamma(n), \kappa(n))$ -Streaming Authentication Scheme is secure in the m -party shared-key setting if the adversary cannot induce any decoding error (as described in Definition 5.1) for any streaming authentication session for which he has not corrupted one of the parties, except with probability at most $s \cdot \kappa(n)$, where s is the total number of sessions initiated.

Claim C.2. The efficient scheme from the proof of Theorem 5.2 is secure in the m -party shared-key setting if the R_{ij} are independent (proof by union bound).

Multi-Party PKI Setting: The m -party PKI setting for $(c(n), \gamma(n), \kappa(n))$ -Streaming Authentication is as follows. There are public parameters $params \leftarrow ParamGen(1^\lambda)$. There are m parties, and each party i has generated a public key pair $(e_i, d_i) \leftarrow KeyGen(params)$. The adversary can see all of the public keys, and can learn any private key by corrupting that party. Otherwise, the adversary can perform any attack as described for the multi-party shared-key setting, with the following two differences: (1) If party i is in a streaming authentication session with party j then he uses d_i and e_j to compute the shared random string, and (2) the adversary is limited to computation that can be performed in probabilistic polynomial time in λ, n, m, s , where s is the total number of streaming authentication sessions initiated.

Definition C.3. A $(c(n), \gamma(n), \kappa(n))$ -Streaming Authentication Scheme is computationally secure in the m -party PKI setting if the adversary cannot induce any decoding error (as described in Definition 5.1) for any streaming authentication session for which he has not corrupted one of the parties, except with probability at most $s \cdot \kappa(n)$, where s is the total number of sessions initiated.

Theorem C.4. The efficient scheme from the proof of Theorem 5.2 is computationally secure in the m -party PKI setting under the Decision Diffie Hellman (DDH) assumption if $ParamGen(1^\lambda) \rightarrow g, q$ where g is a generator of prime order $q > 2^\lambda$; $KeyGen(g, q) \rightarrow (e_i, d_i)$ where $d_i \leftarrow [1 \dots q - 1]$ and $e_i = g^{d_i}$; For their l th streaming authentication session, parties i and j use $R_{ij} = G_{ij\ell}(g^{x_i x_j})$, where $\{G_{ij\ell}\}_{i,j,\ell}$ is a family of pseudorandom generators.

Proof. (sketch) Let A be a probabilistic polynomial-time adversary that can break the $(c(n), \gamma(n), \kappa(n))$ -Streaming Authentication Scheme. Without loss of generality, A corrupts all but two of the m parties during its attack. We use A to construct a ppt adversary B for the DDH problem. Let (g, u, v, w) be a DDH challenge tuple. That is, g is a generator of order q for suitably large prime q , $u = g^x$ for $x \leftarrow [1 \dots q - 1]$, $v = g^y$ for $y \leftarrow [1 \dots q - 1]$, and $w = g^z$ (where $z = xy \bmod q$ with probability 1/2, and $z \leftarrow [1 \dots q - 1]$ with probability 1/2). Given a DDH challenge tuple, adversary B proceeds as follows:

1. Choose $i^*, j^* \leftarrow [1 \dots m]$ (distinguished parties that B hopes will be targeted for attack).

2. For every party $k \neq i^*, j^*$, generate its public key and private key in the usual way. That is, the private key of party k is $x_k \leftarrow [1 \dots q - 1]$, and its public key is g^{x_k} .
3. For party i^* , its public key will be u (and its private key will be unknown). For party j^* , its public key will be v (and its private key will be unknown).
4. If party i^* and party j^* start their ℓ th streaming authentication session, then run the efficient scheme from the proof of Theorem 5.2 using the output of $G_{i^*j^*\ell}(w)$ as the shared random string.
5. If party i^* and party $j \neq i^*$ start their ℓ th streaming authentication session, then run the efficient scheme from the proof of Theorem 5.2 using the output of $G_{i^*j\ell}(u^{x_j})$ as the shared random string.
6. If party j^* and party $i \neq i^*$ start their ℓ th streaming authentication session, then run the efficient scheme from the proof of Theorem 5.2 using the output of $G_{ij^*\ell}(v^{x_i})$ as the shared random string.
7. If party $i \neq i^*$ and party $j \neq j^*$ start a streaming authentication session, then run the efficient scheme from the proof of Theorem 5.2 using the output of $G_{ij\ell}(g^{x_i x_j})$ as the shared random string.
8. If adversary A corrupts party $k \neq i^*, j^*$, then give x_k to the adversary.
9. If adversary A ever attempts to corrupt party i^* or party j^* then halt and flip a coin to determine if (g, u, v, w) is a valid DDH tuple.
10. If adversary A succeeds in violating the decoding condition for any streaming authentication session involving uncorrupted parties, then determine that (g, u, v, w) is a valid DDH tuple.
11. If adversary A fails to violate the decoding condition for any streaming authentication involving uncorrupted parties, then determine that (g, u, v, w) is not a valid DDH tuple.

Let δ be the maximum success probability of any attack against the scheme in the multi-party shared-key setting (with independent keys). Let $\epsilon > \delta$ be the success probability of A against the scheme in the multi-party PKI setting. Here are the cases (and probabilities) where adversary B makes a correct determination of whether or not (g, u, v, w) is a valid DDH tuple:

Case 1: The adversary A never corrupts i^* or j^* , and the DDH challenge tuple is invalid. If the attack by A fails, then B makes the correct determination about (g, u, v, w) . The probability of this case occurring is at least $(2/n(n-1))(1/2)(1-\delta) = (1-\delta)/((n(n-1)))$.

Case 2: The adversary A never corrupts i^* or j^* , and the DDH challenge tuple is valid. If the attack by A succeeds, then B makes the correct determination. The probability of this case occurring is at least $(2/n(n-1))(1/2)(\epsilon) = \epsilon/(n(n-1))$.

Case 3: The adversary A attempts to corrupt i^* or j^* , and B halts and flips a coin. If this coin flip matches the validity of the DDH challenge tuple then B makes the correct determination. The probability of this case occurring is $(1 - (2/n(n-1)))(1/2) = 1/2 - (1/(n(n-1)))$.

These cases are independent, so the total probability that B makes the correct determination is the sum of their probabilities, which is at least $1/2 + (\epsilon - \delta)/(n(n-1))$. If ϵ is non-negligibly greater than δ , then B distinguishes valid from invalid DDH challenge tuples with non-negligible advantage, which is a contradiction. This completes the proof sketch. \square

D The Suffix Condition

Lemma D.1. *Let c and ϵ be given and let $\gamma = c/(c+\epsilon)$, then for the protocol given in Proposition 5.4, and for every noise level Bob's guessed string $x'_1, \dots, x'_{N_\gamma}$ is identical to the transmitted stream x_1, \dots, x_{N_γ} .*

Proof. Assume towards contradiction that Bob recovers the stream correctly only until time $t < N_\gamma$. Thus it follows that the errors and erasures in the $[t, N_\gamma]$ suffix (e and d respectively), satisfy $\alpha(N_\gamma - t) < 2e + d$. However, we know that N_γ satisfies the γ -suffix condition so $e + d < \gamma(N - t)$. Except with probability $2^{-\Omega(N-t)}$ it holds that $2e + d < (1 + 2\frac{S}{L}) \cdot \gamma(N - t)$, thus except with the same probability, $\alpha < (1 + 2\frac{S}{L})\frac{c}{c+\epsilon}$ which is a contradiction to the way we choose α . \square

However Bob does not know N_γ . We now show how to estimate this value.

D.1 Proof of Proposition 5.10 and Discussion

Consider the following procedure for estimating N_γ .

Let $c \in [0, 1]$ be given. For an input $\gamma \in (c, 1)$, at time n , Bob tries to find the longest suffix that satisfies the γ -suffix condition. To this end Bob performs the following.

1. Set $i = n$.
2. Check all the suffixes x_t, \dots, x_i , with $t < i$.
 - (a) If, in all such suffixes, the number of erasures is less than $\gamma(i-t)(1 - 2\frac{S}{L})$, output $N'_\gamma = i$.
 - (b) Otherwise set $i \leftarrow i - 1$ and repeat. If $i < (1 - c/\gamma)n$ break and output $N'_\gamma = \lfloor (1 - c/\gamma)n \rfloor$.

Proof. (Proposition 5.10). Assume that Bob outputs $N'_\gamma > N_\gamma$. Since N_γ is the latest index that satisfies the γ -suffix condition, there must exist some time t such that $[t, N'_\gamma]$ has more than $\gamma(N' - t)$ corruptions, yet the number of erasures in that interval is less than $\gamma(N'_\gamma - t)(1 - 2\frac{S}{L})$. This happens with probability exponentially small in $N'_\gamma - t$. If $t < N_\gamma$ then $N'_\gamma - t > N'_\gamma - N_\gamma$ which proves the claim for this case.

For the case where $t > N_\gamma$, we note that time t does not satisfy the γ -suffix condition, therefore there must exist some time $t_1 < t$, for which the number of corruptions in the interval $[t_1, t]$ is more than $\gamma(t - t_1)$. If $t_1 > N_\gamma$, then there must exist time t_2 and interval $[t_2, t_1]$ that doesn't satisfy the γ -suffix condition. We repeat this reasoning until we find the first interval $[t_j, t_{j-1}]$ such that $t_j < N$. By considering the union of all these intervals, it follows that the number of corruptions in $[t_j, N'_\gamma]$ is more than $\gamma(N'_\gamma - t_j) > \gamma(N'_\gamma - N_\gamma)$. However, Bob's estimation process found at most $\gamma(N'_\gamma - N_\gamma)(1 - 2\frac{S}{L})$ erasures when it examined the suffix interval $[t_j, N'_\gamma]$ (otherwise, it would have failed the check in Step 2a). By Lemma 3.4/Corollary 3.5, this happens with probability exponentially small in $(N'_\gamma - N_\gamma)$. \square

Bob learns a lower-bound on N_γ and it is guaranteed that $x'_1, \dots, x'_{N'_\gamma}$ is the same as $x_1, \dots, x_{N'_\gamma}$, maybe except for the last bits of the stream, in case N'_γ exceeds N_γ .

Bob can repeat the same procedure and compute a value N''_γ which usually upper-bounds N_γ , by finding the latest time i whose every suffix $[t, i]$ has less than $\gamma(i-t)(1 - \frac{1}{2}\frac{S}{L})$. As above, the probability of the bad event that $N''_\gamma < N_\gamma$ is exponentially small in $(N_\gamma - N''_\gamma)$.

D.2 Proof of Lemma 5.9

Proof. Look at a suffix y_{t+1}, \dots, y_n for which the number of corruptions is strictly larger than $c\xi(n-t)$. If no such suffix exists then the lemma is true for y_1, \dots, y_n . Otherwise, discard y_{t+1}, \dots, y_n , and repeat the process with y_1, \dots, y_t . At each iteration we remove more than $c\xi(n-t)$ corrupted transmissions and shorten the string by $n-t$ symbols. Assume that the process stops with some prefix of length $L < (1 - \frac{1}{\xi})n$, then we have removed at least $c\xi(n-L) > c\xi(n - n + n/\xi) > cn$ corruptions which is a contradiction. Therefore, the entire process must stop with some prefix of length at least $(1 - \frac{1}{\xi})n$. \square