

Explicit Binary Tree Codes with Polylogarithmic Size Alphabet

Gil Cohen* Bernhard Haeupler† Leonard J. Schulman‡

Abstract

This paper makes progress on the problem of explicitly constructing a binary tree code with constant distance and constant alphabet size.

For every constant $\delta < 1$ we give an explicit binary tree code with distance δ and alphabet size $\text{poly}(\log n)$, where n is the depth of the tree. This is the first improvement over a two-decade-old construction that has an exponentially larger alphabet of size $\text{poly}(n)$.

As part of the analysis we prove a bound on the number of positive integer roots a real polynomial can have in terms of its *sparsity* with respect to the Newton basis—a result of independent interest.

*Department of Computer Science, Princeton University, Princeton, USA. Part of this work was done while the author was a CMI postdoctoral fellow at The California Institute of Technology, Pasadena, USA. Email: gilc@princeton.edu.

†Department of Computer Science, Carnegie Mellon University, Pittsburgh, USA. Supported in part by NSF grants CCF-1527110, CCF-1618280 and NSF CAREER award CCF-1750808. Email: haeupler@cs.cmu.edu.

‡E&AS Division, California Institute of Technology, Pasadena, USA. Part of this work was done while in residence at the Israel Institute for Advanced Studies; financial support is due to United States NSF grant 1618795 and to a EURIAS Senior Fellowship co-funded by the Marie Skłodowska-Curie Actions under the 7th Framework Programme. Email: schulman@caltech.edu.

1 Introduction

This paper makes progress on the problem of explicitly constructing a binary tree code with constant distance and constant alphabet size.

Tree codes are a powerful but so-far elusive combinatorial structure, defined and proven to exist by [Sch93, Sch96] as a key ingredient for achieving a constant rate interactive coding scheme. Tree codes are the central object for encoding information in the interactive coding theory which developed from the initial papers. They remain a crucial building block in almost all interactive coding schemes [RS94, BR14, BGMO17, GHK⁺16, BK12, BN13, ABE⁺16, BKN14, GH14, GHS14, LV15, AGS16, BK12, GMS14, BN13, BE14, JKL15, SW17]. The absence of an explicit construction that is also efficiently decodable is the only reason why most of these schemes are computationally inefficient, requiring exponential-time computations. Other works have invested significant effort in avoiding the use of (large) tree codes, often at a considerable loss in the fraction or generality of errors that can be tolerated [GMS14, BK12, BN13, Bra12, KR13, Hae14, HV17, BKN14, GH14]. We refer to the excellent survey by Gelles [Gel17] for an in-depth account of the role tree codes hold in the area of interactive coding theory.

In addition, tree codes have important uses as streaming codes for both Hamming errors [FGOS15] and synchronization errors [HS17, HSV17, BGMO17, HS18]. In control theory, although mostly unknown to the computer science community, tree codes are closely connected to anytime reliable codes that are necessary to controlling and stabilizing systems over unreliable channels [Sah01, Sah04, SM06, SH11a, SH11b, Gro16, KHH16]; there, too, the absence of explicit constructions of tree codes was the motivation for an elaborate work-around for certain control applications [ORS09]. Tree codes have also found surprising application in metric embeddings [LdMM12] and complexity theory [Dru11, DL16].

Let us define tree codes and explain why one should think of them as an online version of a regular error correcting block code: A tree code consists of a complete rooted binary tree (either infinite or of finite depth n) in which each edge is labeled by a symbol from an alphabet Σ . There is a natural one-to-one mapping assigning each binary string s to a path starting at the root, where s simply indicates which child is taken in each of the steps. For a tree code, such a path naturally maps to a string over the alphabet Σ , which is formed by concatenating the symbols along the path. This way a tree code T encodes any binary string s into an equally long string $T(s)$ over Σ . This encoding has an online characteristic because the encoding of any prefix does not depend on later symbols. In particular, any two distinct strings that agree in their first k symbols also have encodings that agree in their first k symbols. A tree code is said to achieve distance $\delta \geq 0$ if the encodings of any two strings differ in at least a δ -fraction of the positions after their first disagreement. The rate of a tree code is $\frac{1}{\log_2 |\Sigma|}$. A tree code is said to be asymptotically good if it achieves both constant distance $\delta > 0$ and a constant rate, namely, the alphabet size $|\Sigma| = O(1)$.

Three different proofs were provided in [Sch93, Sch96], showing that for any $\delta < 1$ there exists a binary tree code with a constant-size alphabet achieving distance δ . All of these proofs, as well as a later quantitative improvement by Peczarski [Pec06], rely on the probabilistic method. Interestingly however, in contrast to conventional error correcting block codes, a constant-size-alphabet random tree code is not asymptotically good and has a distance of zero, with high probability.

The problem of giving an explicit construction of asymptotically good tree codes has drawn substantial attention, but has endured as a difficult challenge.¹ Technically, for a depth n tree code to be explicit, we require that there exist a deterministic algorithm running in time $\text{poly}(n)$ which on input $s \in \{0, 1\}^{n'}$ ($n' \leq n$), outputs the label of the last edge on the path s . We remark that one of the existence proofs is based on the Lovász local lemma [EL75] but the algorithmic LLL [MT10, CGH13] does not yield explicit tree codes as it must construct the entire $\exp(n)$ -size tree. In fact, for explicit constructions not much has been known beyond a construction of Evans, Klugerman and Schulman [Sch94] (dating to 1994) that provides a tree code with alphabet size $\text{poly}(n)$. Pudlák [Pud16] studies sufficient and necessary structural results for linear (MDS) tree codes and provides a construction with large arity. Moore and Schulman gave a candidate construction [MS14], but its distance property relies on an open conjecture about certain exponential sums. A tree code construction which reduces the brute-force time of encoding and decoding from exponential, i.e., $\exp(n)$, to sub-exponential, i.e., $\exp(n^\varepsilon)$, at the cost of an alphabet size $\exp(1/\varepsilon)$, was given by Braverman [Bra12].

1.1 Our Results

In this work we obtain the first proven improvement over the two-decade-old construction of [Sch94] by giving an explicit binary tree code with constant distance and an exponentially smaller, i.e., polylogarithmic, alphabet size.

Theorem 1.1. *For every constant $\delta < 1$ and integer $n \geq 1$ there exists an explicit binary tree code $\text{TC}: \{0, 1\}^n \rightarrow \Sigma^n$ with distance δ and $|\Sigma| = (\log n)^{O(1)}$.*

Put differently, Theorem 1.1 gives a binary tree code with rate $\Omega(1/\log \log n)$ and distance δ . We point out that our techniques readily yield a depth n tree code that can be constructed in time $\exp(n^\varepsilon)$ with alphabet size $\text{poly}(1/\varepsilon)$. This should be compared with an alphabet size $\exp(1/\varepsilon)$ that was obtained by Braverman [Bra12] under the same running-time restriction.

We prove Theorem 1.1 in two steps. First, we construct a tree code over the integers as given in Theorem 1.2 below. This tree code has the advantage of being infinite. We then reduce the input alphabet to Boolean.

¹Wigderson in his new book “Mathematics and Computation” calls it “the most elegant open problem of [the] theory [of interactive coding]” [Wig17, page 202, Open Problem 15.34].

Theorem 1.2. *For every constant $\delta < 1$ there exists an explicit tree code $\text{TC}_{\mathbb{Z}}: \mathbb{Z}^{\mathbb{N}} \rightarrow \mathbb{Z}^{\mathbb{N}}$ ² with distance δ . Further, for every $z = (z_t)_{t \in \mathbb{N}} \in \mathbb{Z}^{\mathbb{N}}$ and $t \in \mathbb{N}$,*

$$|\text{TC}_{\mathbb{Z}}(z)_t| \leq 2^{O(t^2)} \cdot (\max(|z_0|, \dots, |z_t|))^{O(1)}.$$

Our construction is at its cleanest form when $\delta = 1/2$. In such case, the dependence on t is also better. Throughout this section we focus on this tree code whose parameters are given in the following theorem.

Theorem 1.3. *There exists an explicit tree code $\text{TC}_{\mathbb{Z}}: \mathbb{Z}^{\mathbb{N}} \rightarrow \mathbb{Z}^{\mathbb{N}}$ with distance $1/2$. Further, for every $z = (z_t)_{t \in \mathbb{N}} \in \mathbb{Z}^{\mathbb{N}}$ and $t \in \mathbb{N}$, $|\text{TC}_{\mathbb{Z}}(z)_t| \leq 2^t \cdot \max(z_0^2, \dots, z_t^2)$.*

We wish to give some remarks regarding the bound on $|\text{TC}_{\mathbb{Z}}(z)_t|$ that is guaranteed by Theorem 1.3. Assume that $|z_t| \leq m$ for all t . Theorem 1.3 gives a bound of $2^t m^2$ on the t 'th output symbol. This should be compared with the trivial bound of m^t and with the bound $m^{\log t}$ that is obtained by adapting the technique of [Sch94] to tree codes over the integers. Although our bound has an exponential dependence on t , the two parameters m and t are decoupled and so one can take t super-constant while keeping the bound polynomial in m . In the second step of our construction, we show that this property suffices to obtain the improved binary tree code claimed by Theorem 1.1 with distance $1/2$. The same argument shows how Theorem 1.2 implies Theorem 1.1 for any constant distance $\delta < 1$.

The proof of Theorem 1.3 is obtained by adapting the Reed Solomon polynomial interpolation framework to the online setting. We give an overview of the proof in Section 2 and the formal proof is the content of Section 5. To analyze the distance, we prove a bound on the number of distinct integral roots a real polynomial can have in terms of its *sparsity* in a certain basis—a result of independent interest on which we elaborate on in Section 1.1.1 below. The alphabet reduction technique we use to deduce Theorem 1.1 is covered in Section 6. Finally, in Section 7, we prove Theorem 1.2.

1.1.1 A Bound on the Number of Integral Roots via Sparsity

The fundamental theorem of algebra asserts that a degree $d > 0$ polynomial with complex coefficients has exactly d complex roots when counted with multiplicities. More generally, over any field \mathbb{F} , a degree $d > 0$ polynomial $f \in \mathbb{F}[x]$ has at most d roots in \mathbb{F} (and exactly d roots in the algebraic closure of \mathbb{F}).

The sparsity of a polynomial, however, cannot be used to bound the number of its distinct roots. There are natural examples of sparse polynomials with many roots even in the base field, e.g., $x^p - x$ in $\mathbb{F}_p[x]$. Nevertheless, for the analysis of our tree code construction, we provide a meaningful bound on the number of positive *integer* roots (that is, roots in \mathbb{N}) a real polynomial can have in terms of its sparsity.

²Throughout the paper, it will be convenient to use the notation $\mathbb{N} = \{0, 1, 2, \dots\}$.

Unlike the notion of degree, sparsity is, of course, basis dependent. The basis for which our bound holds is not the standard basis $\{1, x, x^2, \dots\}$ but rather the Newton basis which consists of polynomials of the form $\binom{x}{k} \in \mathbb{R}[x]$ for $k \in \mathbb{N}$, where

$$\binom{x}{k} = \frac{x(x-1)\cdots(x-(k-1))}{k!}.$$

It is easy to verify that for every $d \in \mathbb{N}$, the set $\{\binom{x}{k} \mid k = 0, 1, \dots, d\}$ forms a basis for the space of univariate real polynomials of degree at most d .

Of course, with respect to this basis, the sparsity cannot be taken as a bound on the number of distinct roots a polynomial can have. Indeed, for any $d \in \mathbb{N}$, consider the degree d polynomial $\binom{x}{d}$ which has sparsity $s = 1$ in the Newton basis. Evidently, $\binom{x}{d}$ has d distinct roots at $x = 0, 1, \dots, d-1$. Thus, one cannot hope to prove a general bound on the number of roots in terms of sparsity even when restricting to integral roots and not accounting for multiplicities.

Consider a polynomial with sparsity $s = 2$ in the Newton basis. Such a polynomial has the form

$$f(x) = \gamma \binom{x}{c} + \delta \binom{x}{d},$$

where $0 \leq c < d$ are integers and γ, δ are nonzero real numbers. Clearly, $0, 1, \dots, c-1$ are all roots of f , and c can be taken much larger than 2 – the sparsity of f . More generally, if f is a polynomial with sparsity s and $c = c(f)$ is the least integer such that $\binom{x}{c}$ appears in the expansion of f in the Newton basis then f will surely have $0, 1, \dots, c-1$ as its roots. Again, it may be the case that $c \gg s$.

We prove that but for these c “trivial” integral roots, f has at most $s-1$ roots in \mathbb{N} . This holds regardless of the degree of f . More precisely, we prove the following lemma which can be interpreted as an uncertainty principle for the Newton basis.

Lemma 1.4. *Let $f \in \mathbb{R}[x]$ be a nonzero polynomial of sparsity $s \geq 1$ in the Newton basis. Let $c \geq 0$ be the least integer such that $f(c) \neq 0$. Then, f has at most $s-1$ distinct roots in $[c, \infty) \cap \mathbb{Z}$.*

Observe that the restriction to integral roots is necessary, that is, one cannot strengthen the result by arguing about non-integral roots in $[c, \infty)$. To see this, take any integer $d > 1$ and consider the polynomial with integral coefficients $f_C(x) = \binom{x}{1} + C \binom{x}{d}$, where $C \in \mathbb{N}$ is chosen sufficiently large. The polynomial f_C has sparsity 2, degree d and, with the notation above, $c(f_C) = 1$. However, f_C has the same roots as $f_\varepsilon(x) = \varepsilon \binom{x}{1} + \binom{x}{d}$ where $\varepsilon = 1/C \approx 0$. However, $f_\varepsilon(x) \approx \binom{x}{d}$ and so for ε small enough in absolute value, f_ε has $d-1$ distinct roots in $[1, \infty)$.

We prove Lemma 1.4 in Section 4. The proof makes use of the beautiful Lindström-Gessel-Viennot Lemma (see Lemma 3.9). Given the usefulness of the degree bound on the number of roots, we believe that Lemma 1.4 should find further applications.

2 Overview of the Construction

The polynomial interpolation framework is at the heart of several important constructions of error correcting block codes such as the Reed Solomon code. Our construction is based on identifying a suitable adjustment of the polynomial interpolation framework to the online setting. To motivate our construction, we start by highlighting the difficulties in pursuing such an approach. To this end, we recall the definition of the Reed Solomon code. Let n be an integer. Assume, for simplicity, that n is prime, and let \mathbb{F} be the field of n elements. For an integer $k \leq n$, the Reed Solomon code $\text{RS}: \mathbb{F}^k \rightarrow \mathbb{F}^n$ is defined as follows. Define the polynomial $f_m(x) = \sum_{i=0}^{k-1} m_i x^i \in \mathbb{F}[x]$. The encoding of m is defined by $\text{RS}(m) = (f_m(0), f_m(1), \dots, f_m(n-1))$.

As f_m is linear in m , the analysis of the distance of RS proceeds by proving an upper bound on the number of zero entries of a codeword that corresponds to a nonzero message. By construction, these entries correspond to the number of distinct roots of f_m in \mathbb{F} . Here is where the degree bound on the number of roots of f_m is invoked.

An obvious difficulty in adapting the above idea to the construction of tree codes arises from the latter's online nature. As we do not have the entire message available to us up until the very end, there is no clear sense as to which polynomial we should work with. However, the challenge is more significant. Even given the entire message, one still needs to gain nonzero output symbols starting from the index of the first nonzero entry of the message. That is, not only that one has to work with partial information, a tree code must also gain distance as soon as a disagreement, or “split”, occurs and to keep the distance above a certain threshold from that point on. Restricting to the Reed Solomon construction, this means that even given the message m , the polynomial f_m defined above should somehow be evaluated on a carefully chosen sequence of points in the field. We do not know how to implement such an approach or even if it is possible in principle. Anyhow, one does not have the message in its entirety.

Having these difficulties in mind motivates our construction which we present next. Although our construction is based on polynomial interpolation, and so it is inherently algebraic, it can be motivated both using a combinatorial reasoning and from an algebraic perspective. We start by presenting the combinatorial point of view in Section 2.1. We then discuss the algebraic perspective in Section 2.2. The formal proof, given in Section 5, is presented and analyzed only via the algebraic perspective, nevertheless, we believe that the combinatorial point of view gives a natural motivation for our construction.

2.1 The combinatorial point of view

In this section we motivate and describe the tree code construction $\text{TC}_{\mathbb{Z}}: \mathbb{Z}^{\mathbb{N}} \rightarrow \mathbb{Z}^{\mathbb{N}}$ from Theorem 1.3. Let $z = (z_t)_{t \in \mathbb{N}} \in \mathbb{Z}^{\mathbb{N}}$ be a message that we want to encode. For $t \in \mathbb{N}$, define $f_t \in \mathbb{R}[x]$ to be the polynomial of least degree such that $f_t(i) = z_i$ for all $i \in \{0, 1, \dots, t\}$. Note

that f_t is fully determined by z_0, \dots, z_t . Further, observe that f_t is linear in z . Therefore, so long as we define $\text{TC}_{\mathbb{Z}}(z)_t$ as a linear combination of the evaluation of the polynomials f_0, \dots, f_t on fixed points, it will follow that $\text{TC}_{\mathbb{Z}}(z)_t = \text{TC}_{\mathbb{Z}}(z')_t$ if and only if $\text{TC}_{\mathbb{Z}}(z - z')_t = \text{TC}_{\mathbb{Z}}(\bar{0})_t = 0$; i.e., for purposes of distance analysis, it suffices to compare every nonzero message z against the all-zeros message.

To recap, while the Reed Solomon code interprets the message as a polynomial, in our construction (which has not yet been presented) every prefix of the message is interpreted as a polynomial, and so z , chosen by the “adversary”, induces an infinite sequence of polynomials f_0, f_1, f_2, \dots .

Consider a scenario in which the adversary makes it so that $f_t = f_{t+1} = \dots = f_{t+\ell}$. Intuitively, such a scenario is favorable for us. Indeed, one can imagine how useful it would be if the adversary is committed to a single polynomial f for a long interval of time while (just as in Reed Solomon) outputting evaluations of f in the interval. In some sense, it is as if the tree code is only required to work against an off-line input on that interval; the fundamental theorem of algebra can come into play and prevent having many 0s in the output during this interval. Thus, a natural idea is to penalize the adversary when switching to a new polynomial from time $t - 1$ to time t . This can be done by outputting, at time t , the value

$$\delta_t = f_t(t) - f_{t-1}(t).$$

In control theory, one would call δ_t the “innovation”. A complexity-theoretic point of view would interpret δ_t as a consistency checking procedure. Indeed, unless the adversary sticks to his polynomial f_{t-1} at time t , he pays in distance as then $\delta_t \neq 0$. This consistency checking symbol at time t is concatenated with $f_t(t)$, i.e., z_t . To summarize, we define $\text{TC}_{\text{comb}}: \mathbb{Z}^{\mathbb{N}} \rightarrow (\mathbb{Z}^2)^{\mathbb{N}}$ by ³

$$\text{TC}_{\text{comb}}(z)_t = (z_t, \delta_t).$$

It is not immediately clear from this representation why δ_t should be an integer but as it turns out, that is the case and $|\delta_t|$ can be bounded as we discuss at the end of Section 2.2.

How large is the distance of TC_{comb} ? It seems that an adversary that sticks to a polynomial for not-too-long intervals may pay very little as we do not gather sufficient amount of information during a short interval. On the other hand, it is intuitive that something is gained by this approach. How would an adversary work against TC_{comb} ?

A potential attack. One attack might work as follows. Fix some $d \geq 1$. The adversary will choose z_0, \dots, z_d such that the degree d polynomial f_d has roots at $d + 1, \dots, 2d$. Now by choosing just one new nonzero value z_{2d+1} , the adversary obtains a new polynomial f_{2d+1} of degree $2d + 1$ which shares the already-recorded d roots but also has $d + 1$ new roots,

³Note that the output symbols are pairs of integers rather than integers as stated in Theorem 1.3. This, of course, is a non-issue and is only meant for a cleaner presentation.

potentially at $2d+2, \dots, 3d+2$; the adversary then uses $z_{2d+2} = \dots = z_{3d+2} = 0$ and in this case we have $\delta_{2d+2} = \dots = \delta_{3d+2} = 0$. The adversary can again use a nonzero z_{3d+3} , obtaining a new polynomial f_{3d+3} with $d+2$ new roots, potentially at $3d+4, \dots, 4d+5$, and the adversary then makes the choice $z_{3d+4} = \dots = z_{4d+5} = 0$. If this process can be repeated in this manner (i.e., if the described polynomials exist), the result will be a branch of the tree code which at depth ℓ has weight only $O(\sqrt{\ell})$.

Even if there are no polynomials that perfectly interpolate as required by the above attack, it is not obvious that one can rule out a quantitatively-relaxed version of such an attack. What one can see however, is that the adversary cannot beat $\Omega(\sqrt{\ell})$: after the s 'th nonzero value, say it is z_t , there cannot be a run of zeros $0 = f_t(t+1) = f_t(t+2) = \dots = f_t(t+s)$ as this would contradict Lemma 3.9, which establishes that the relevant minor in the linear transformation is nonsingular.

Though certainly a nontrivial bound, an $\Omega(1/\sqrt{\ell})$ distance is far from the constant distance that we are shooting for. Interestingly, TC_{comb} has, in fact, distance $1/2$! In particular, the above attack is far from feasible. To prove that, we consider an algebraic point of view on the construction of TC_{comb} . Taking the algebraic perspective, we can prove that the adversary has a budget of roots that is bounded by the sparsity (with respect to a certain basis) of the polynomials rather than by their degree.

2.2 The algebraic perspective

So far, when working with the polynomials $(f_t)_t$, we did not pay attention to the basis in which the polynomials are represented. Generally, the polynomial interpolation framework works over any basis as the degree, which is typically used in the analysis, is basis invariant. However, for our purpose, the standard basis $\{1, x, x^2, \dots\}$ has the following drawback. Let $y_0, \dots, y_n \in \mathbb{R}$. Let $f(x) = \sum_{i=0}^n a_i x^i$ be the least degree polynomial that interpolates on the points $(0, y_0), (1, y_1), \dots, (n, y_n)$. Then generally, given a new point $(n+1, y_{n+1})$, the least degree polynomial, $g(x) = \sum_{i=0}^{n+1} b_i x^i$, that interpolates on $(0, y_0), \dots, (n+1, y_{n+1})$ will have a completely different sequence of coefficients (i.e., $a_i \neq b_i$).

By contrast, using the Newton basis, the coefficients that were already “recorded” stay intact given the new point $(n+1, y_{n+1})$. More precisely, if $f(x) = \sum_{i=0}^n \gamma_i \binom{x}{i}$ then $g(x) = f(x) + \gamma_{n+1} \binom{x}{n+1}$ for some $\gamma_{n+1} \in \mathbb{R}$. Classically, this fact makes the Newton basis attractive for numerical stability and was used for obtaining structural results for polynomials [vZGR97, ST11, CST11]. For constructing tree codes, this property is attractive as it means that for every t , the coefficient γ_t is determined by y_0, y_1, \dots, y_t .

The above discussion suggests a second construction of tree codes over \mathbb{Z} . Let γ_t be the coefficient of $\binom{x}{t}$ in the expansion of f_t , as defined in Section 2.1. We define the tree code $\text{TC}_{\text{alg}}: \mathbb{Z}^{\mathbb{N}} \rightarrow (\mathbb{Z}^2)^{\mathbb{N}}$ by

$$\text{TC}_{\text{alg}}(z) = (z_t, \gamma_t).$$

Interestingly, one can show that $\gamma_t = \delta_t$ for every t and so TC_{comb} and TC_{alg} are one and the same! The algebraic point of view on the tree code will allow us to prove a bound of $1/2$ on the distance as we now explain.

Let c be the least integer such that $z_c \neq 0$. Let $\ell \geq 1$ and set $t = c + \ell - 1$. Observe that the number of non-zeros in the sequence $\gamma_0, \gamma_1, \dots, \gamma_t$ is precisely the sparsity of f_t in the Newton basis. This, together with the fact that for every $i \leq t$, $z_i = f_t(i)$, implies that to “break” the construction TC_{alg} , the adversary must come up with a sparse polynomial f_t that has many roots in $I = \{c, c + 1, \dots, t\}$. Indeed, if f_t is not sparse, then many of the γ -entries of $(\text{TC}_{\text{alg}}(z))_{i \in I}$ will be nonzero. On the other hand, if f_t has only few roots in I then many of the z -entries are nonzero.

To give the quantitative bound we invoke Lemma 1.4 which implies that if the sparsity of f_t is s then there can be at most $s - 1$ zeros among the z -entries of $\{\text{TC}_{\text{alg}}(z)_i\}_{i \in I}$. So the combined number of nonzero integers among the 2ℓ integers in $(\text{TC}_{\text{alg}}(z))_{i \in I}$ is at least $\ell + 1$. Thus, at least half of the pairs are nonzero pairs, establishing a distance of $1/2$.

Another issue that can be handled via the algebraic perspective is related to the integrality of the output symbols. The symbol z_t is clearly an integer. However, it is not a priori clear that $\gamma_t = \delta_t$ is an integer. The Newton basis has another useful property we use—if z_0, \dots, z_t are all integers, so are the coefficients $\gamma_0, \dots, \gamma_t$. Note that this property does not hold for the standard basis. Moreover, there is a closed formula for γ_t as a function of z_0, \dots, z_t (see Lemma 3.8) which allows us to prove the desired bound on $|\gamma_t|$.

2.3 Tree codes for any distance $\delta < 1$

It is fairly straightforward to adapt the ideas described above to obtain any distance $\delta < 1$. A natural strategy is to use more evaluation points. This idea, however, should be executed with some care. In this section we sketch how to obtain distance $\delta = 2/3$. The idea can be easily generalized to yield any distance $\delta < 1$ (Section 7).

Let us suggestively denote the input message by $z = (z_0, z_2, z_4, \dots) \in \mathbb{Z}^{2\mathbb{N}}$. As before, we define a sequence of real polynomials f_0, f_1, f_2, \dots . However, to obtain the improved distance, we also *define* a sequence of integers z_1, z_3, z_5, \dots inductively on $t \in \mathbb{N}$, as follows. For even t we define f_t , as before, to be the least degree real polynomial such that $\forall i \in \{0, 1, \dots, t\}$, $f_t(i) = z_i$. We then define $f_{t+1} = f_t$ and compute $z_{t+1} = f_{t+1}(t + 1)$.

For $t \in \mathbb{N}$, let γ_t be the coefficient of $\binom{x}{t}$ in the expansion of f_t . We define $\text{TC}_{\text{alg}}^{2/3}: \mathbb{Z}^{2\mathbb{N}} \rightarrow (\mathbb{Z}^3)^{\mathbb{N}}$ by

$$\text{TC}_{\text{alg}}^{2/3}(z_0, z_2, z_4, \dots)_t = (\gamma_{2t}, z_{2t}, z_{2t+1}).$$

One can show that $\text{TC}_{\text{alg}}^{2/3}$ is a linear online function. To argue about the distance, let c be the least integer such that $z_{2c} \neq 0$. Let $\ell \geq 1$. Set $t = c + \ell - 1$ and denote $I = \{c, c + 1, \dots, t\}$. Observe that $\gamma_i = 0$ for every odd i , and so the number of non-zeros in the sequence $\gamma_0, \gamma_2, \dots, \gamma_{2t}$ is the sparsity of f_{2t} in the Newton basis. By Lemma 1.4, among

the evaluation points $\{2c, 2c + 1, \dots, 2t, 2t + 1\}$, at most $s - 1$ are roots of f_{2t} . Thus, the number of nonzero triplets among $(\text{TC}_{\text{alg}}^{2/3}(z)_i)_{i \in I}$ is at least

$$\frac{s + 2\ell - (s - 1)}{3} \geq \frac{2}{3}\ell,$$

proving that the distance is $2/3$.

One concern that must be addressed is the bound on the γ symbols. Unlike the $1/2$ distance construction, now γ_t depends on the *computed* value z_{t-1} which, in turn, depends on γ_{t-2} . Thus, potentially, the γ symbols can grow much faster and, indeed, the bound we give for distance $2/3$ is weaker than the corresponding bound for distance $1/2$. Nevertheless, as it turns out, for deducing Theorem 1.1, the weaker bound suffices. Further, the bound does not degrade substantially when considering any constant distance $2/3 < \delta < 1$.

3 Preliminaries

Let $n \geq 1$ be an integer and Σ some (finite or infinite) set. For a string $x = (x_1, \dots, x_n) \in \Sigma^n$ and integers $1 \leq a \leq b \leq n$, we let $x_{[a,b]}$ denote the substring (x_a, \dots, x_b) . If $\sigma \in \Sigma$ then σ^n denotes the string $(\sigma, \dots, \sigma) \in \Sigma^n$. Given $x, y \in \Sigma^n$, we write $\text{dist}(x, y)$ for their Hamming distance.

For an integer $n \geq 1$ write $[n]$ for $\{1, 2, \dots, n\}$. We use the convention that 0 is a natural number and denote $\mathbb{N} = \{0, 1, 2, \dots\}$. We also follow the standard convention that $\binom{a}{b} = 0$ for integers $0 \leq a < b$.

3.1 Error correcting block codes

Definition 3.1. A function $\text{ECC}: \Sigma_{\text{in}}^k \rightarrow \Sigma_{\text{out}}^n$ is an error correcting block code with distance δ if for every distinct $x, y \in \Sigma_{\text{in}}^k$, $\text{dist}(\text{ECC}(x), \text{ECC}(y)) \geq \delta n$. The rate of ECC is given by $(k \log_2 |\Sigma_{\text{in}}|) / (n \log_2 |\Sigma_{\text{out}}|)$.

For the proof of Theorem 1.1, it is convenient to consider error correcting block codes whose output length is shorter than their input length and with output alphabet consisting of binary strings of a certain length. We make use of the following construction of error correcting block codes. The construction and its proof are given in Appendix A.

Lemma 3.2. For every constant $0 < \delta < 1$ and constant integer $t \geq 1$ there exists an integer $c = c(t, \delta)$ such that for every large enough integer n there exists an explicit error correcting block code $\text{ECC}: \{0, 1\}^n \rightarrow (\{0, 1\}^c)^{n/t}$ with distance δ .

3.2 Tree codes

Tree codes, as their name suggest, are trees with certain distance properties. However, in this paper, we use an equivalent definition of tree codes that more directly specifies their online characteristic compared to the one given in the original papers [Sch93, Sch96] and, in particular, does not involve trees. This will be more convenient for presenting our construction.

Definition 3.3. A function $f: \Sigma_{\text{in}}^n \rightarrow \Sigma_{\text{out}}^n$ is said to be online if for every $i \in [n]$ and $x \in \Sigma_{\text{in}}^n$, $f(x)_i$ is determined by x_1, \dots, x_i .

Definition 3.4. For a pair of distinct $x, y \in \Sigma^n$, we define $\text{split}(x, y)$ as the least integer $s \in [n]$ such that $x_s \neq y_s$.

Definition 3.5 ([Sch93, Sch96]). An online function $\text{TC}: \Sigma_{\text{in}}^n \rightarrow \Sigma_{\text{out}}^n$ is a tree code with distance δ if for every distinct $x, y \in \Sigma_{\text{in}}^n$, with $s = \text{split}(x, y)$, and every $\ell \in \{0, 1, \dots, n - s\}$,

$$\text{dist}(\text{TC}(x)_{[s, s+\ell]}, \text{TC}(y)_{[s, s+\ell]}) \geq \delta(\ell + 1).$$

We refer to n as the depth of TC . We refer to $\Sigma_{\text{in}}, \Sigma_{\text{out}}$ as the input alphabet and output alphabet, respectively.

We remark that the terms depth and split are coming from the original point of view of tree codes as trees with certain distance properties. The depth is simply the depth of the tree and the split is the level at which the pair of paths diverge. We borrow this terminology even though we do not explicitly view tree codes as trees in this work. We are interested in some further properties of tree codes.

Definition 3.6. Let $\text{TC}: \Sigma_{\text{in}}^n \rightarrow \Sigma_{\text{out}}^n$ be a tree code.

- We say that TC is a binary tree code if $\Sigma_{\text{in}} = \{0, 1\}$.
- Assume $\Sigma_{\text{in}}, \Sigma_{\text{out}}$ are rings. TC is said to be linear if for every $t \in [n]$, $\text{TC}(x)_t$ is a linear function of x .
- We say that TC is explicit if it can be evaluated on every input $m \in \Sigma_{\text{in}}^n$ in polynomial time in the bit complexity of m .

We also consider the stronger notion of infinite tree codes, as was done in the original papers [Sch93, Sch96]. For a set Σ , we denote the set of all sequences $(x_i)_{i \in \mathbb{N}}$, where $x_i \in \Sigma$, by $\Sigma^{\mathbb{N}}$. One can extend the notion of a split and of online functions to functions of the form $f: \Sigma_{\text{in}}^{\mathbb{N}} \rightarrow \Sigma_{\text{out}}^{\mathbb{N}}$ in the natural way.

Definition 3.7 ([Sch93, Sch96]). An online function $\text{TC}: \Sigma_{\text{in}}^{\mathbb{N}} \rightarrow \Sigma_{\text{out}}^{\mathbb{N}}$ is a tree code with distance δ if for every distinct $x, y \in \Sigma_{\text{in}}^{\mathbb{N}}$, with $s = \text{split}(x, y)$, and every integer $\ell \geq 0$,

$$\text{dist}(\text{TC}(x)_{[s, s+\ell]}, \text{TC}(y)_{[s, s+\ell]}) \geq \delta(\ell + 1).$$

Note that an infinite tree code with distance δ yields, for every integer $n \geq 1$, a tree code of depth n with distance δ . We extend, in the natural way, the property of linearity for infinite tree codes. We say that an infinite tree code is explicit if for every $t \in \mathbb{N}$, the restriction of TC to its first t coordinates is explicit as a finite tree code. Such a restriction is well-defined as TC is an online function.

3.3 The Newton basis

For $k \in \mathbb{N}$, the Newton polynomial $\binom{x}{k} \in \mathbb{R}[x]$ is defined by

$$\binom{x}{k} = \frac{x(x-1)\cdots(x-(k-1))}{k!}.$$

As mentioned, $\{\binom{x}{k}\}_{k=0}^n$ is a basis for the space of polynomials of degree at most n , over \mathbb{R} . In fact any function $f: \mathbb{N} \rightarrow \mathbb{R}$ can be expanded as a pointwise-converging power series over this basis. The following lemma gives a formula for the coefficients of the expansion.

Lemma 3.8. *Let $f: \mathbb{N} \rightarrow \mathbb{R}$. Then, for $x \in \mathbb{N}$,*

$$f(x) = \sum_{k=0}^x \gamma_k \binom{x}{k} = \sum_{k \geq 0} \gamma_k \binom{x}{k} \quad (3.1)$$

where

$$\gamma_k = \sum_{i=0}^k (-1)^{k-i} \binom{k}{i} f(i) = \sum_{i \geq 0} (-1)^{k-i} \binom{k}{i} f(i).$$

Furthermore, if $f \in \mathbb{R}[x]$ is a polynomial of degree n then it equals the sum of the first $n+1$ terms of expansion (3.1).

This is simply the inversion formula for a triangular matrix; for a proof see, e.g., Appendix A in [CST11].

3.4 The Lindström-Gessel-Viennot Lemma

Our analysis relies on the following corollary of the beautiful lemma by Lindström [Lin73] and by Gessel and Viennot [GV85]. See also [Aig07], Chapter 5.4 or [AZ10], Chapter 25.

Lemma 3.9 ([GV85], Corollary 2). *Let $0 \leq a_1 < a_2 < \cdots < a_n$ and $0 \leq b_1 < b_2 < \cdots < b_n$ be integers. Define the $n \times n$ matrix M by $M_{i,j} = \binom{a_i}{b_j}$. If $a_i \geq b_i$ for each $i \in [n]$ then $\det M \neq 0$.*

4 A Bound on the Number of Integral Roots via Sparsity

Proof of Lemma 1.4. The proof is by contradiction. Let $s \geq 1$ be least integer such that there is a counterexample: a polynomial $f \in \mathbb{R}[x]$ with sparsity s , specified by integers $0 \leq c_1 < \dots < c_s$ and non-zero real numbers $\gamma_1, \dots, \gamma_s$ such that

$$f(x) = \sum_{i=1}^s \gamma_i \binom{x}{c_i}$$

(note that $c = c_1$), and such that there exist integers t_i with $c \leq t_1 < \dots < t_s$ such that all $f(t_i) = 0$. Necessarily $s > 1$ as the case $s = 1$ merely reflects that $\binom{x}{c} \neq 0$ for $x \geq c$.

Now we argue that all $t_j > c_j$; This clearly holds for $j = 1$ as $f(c_1) = \gamma_1 \neq 0$. By way of contradiction, let $j \geq 2$ be the least counterexample. Then, $t_{j-1} \leq c_{j-1} \leq c_j - 1$. By Lemma 3.8, the polynomial $\sum_{i=1}^{j-1} \gamma_i \binom{x}{c_i}$ agrees with f on $\{0, \dots, c_j - 1\}$, so it has roots t_1, \dots, t_{j-1} (all $> c_1$) and sparsity $j - 1$; since $j - 1 \leq s - 1 < s$, this contradicts the minimality of f .

Finally, consider the $s \times s$ matrix A with entries $A_{i,j} = \binom{t_i}{c_j}$ for $i, j \in [s]$. Let $\gamma \in \mathbb{R}^s$ be the vector with entries γ_j . Then $f(t_i) = (A\gamma)_i$ so, all t_i being roots, $A\gamma = \bar{0}$. However, since $c_1 < \dots < c_s$, $t_1 < \dots < t_s$, and all $c_j < t_j$, the hypothesis of Lemma 3.9 is met for A and so $\det(A) \neq 0$, a contradiction to $\gamma \neq \bar{0}$. \square

5 Infinite Tree Codes over the Integers

In this section we prove Theorem 1.3. We start by defining the construction of $\text{TC}_{\mathbb{Z}}$.

The construction of $\text{TC}_{\mathbb{Z}}$. Define the function $\text{TC}_{\mathbb{Z}}: \mathbb{Z}^{\mathbb{N}} \rightarrow (\mathbb{Z}^2)^{\mathbb{N}}$ as follows. For $z = (z_t)_{t \in \mathbb{N}} \in \mathbb{Z}^{\mathbb{N}}$ let $f: \mathbb{N} \rightarrow \mathbb{N}$ be such that $f(t) = z_t$ for all $t \in \mathbb{N}$. By Lemma 3.8, one can expand f in the Newton basis

$$f(x) = \sum_{t \in \mathbb{N}} \gamma_t \binom{x}{t}.$$

For $t \in \mathbb{N}$, define $\text{TC}_{\mathbb{Z}}(z)_t = (z_t, \gamma_t)$.

Analysis. By Lemma 3.8, for all $t \in \mathbb{N}$, γ_t is a \mathbb{Z} -linear combination of z_0, \dots, z_t and so $\text{TC}_{\mathbb{Z}}$ is a linear online function with the asserted range. For $t \in \mathbb{N}$, define $m_t = \max(|z_i|: i \in \{0, 1, \dots, t\})$. By Lemma 3.8,

$$|\gamma_t| \leq \sum_{i=0}^t \binom{t}{i} |z_i| \leq 2^t m_t$$

which proves the asserted bound on the output symbols.

We turn to analyze the distance of $\text{TC}_{\mathbb{Z}}$. As $\text{TC}_{\mathbb{Z}}$ is linear, it suffices to consider a nonzero sequence $z = (z_t)_{t \in \mathbb{N}}$. Assume that $c \in \mathbb{N}$ is the least integer such that $z_c \neq 0$. Let $\ell \geq 0$. Set $t = c + \ell$ and define $I = \{c, c + 1, \dots, t\}$. Let $f_t \in \mathbb{R}[x]$ be the polynomial

$$f_t(x) = \sum_{i=0}^t \gamma_i \binom{x}{i}.$$

Observe that for every $i \in I$, the first entry of the pair $\text{TC}_{\mathbb{Z}}(z)_i$ equals $z_i = f(i) = f_t(i)$. Let s be the sparsity of f_t in the Newton basis. Note that precisely s of the pairs $(\text{TC}_{\mathbb{Z}}(z)_i)_{i \in I}$ have a nonzero second entry. On the other hand, by Lemma 1.4, f_t has at most $s - 1$ roots in I . Hence, the number of indices $i \in I$ for which the first entry of $\text{TC}_{\mathbb{Z}}(z)_i$ is 0 is bounded above by $s - 1$. Thus, the number of indices $i \in I$ for which $\text{TC}_{\mathbb{Z}}(z)_i$ is nonzero (as a pair) is bounded below by

$$\max(s, \ell + 1 - (s - 1)) \geq \frac{\ell + 1}{2}.$$

This completes the proof of Theorem 1.3.

We remark that obtaining a construction with output symbols that are bounded by $\text{poly}(t)$ rather than the exponential dependence that was obtained above, would yield asymptotically good binary tree code.

6 Binary Tree Codes with Polylogarithmic Size Alphabet

In this section we prove Theorem 1.1. We do so for distance $\delta = 1/3$ based on Theorem 1.3. In Section 7, we explain how to achieve any distance $\delta < 1$ based on Theorem 1.2. We start by deducing the following corollary from Theorem 1.3.

Corollary 6.1. *For every integer $\ell \geq 1$ there exists an explicit tree code*

$$\text{TC}_{\ell}: (\{0, 1\}^{\ell})^{\ell} \rightarrow (\{0, 1\}^{3\ell})^{\ell}$$

with distance $1/2$.

We remark that by using Pudlák's construction ([Pud16], Lemma 6.1), one can obtain an explicit tree code $\text{TC}: (\{0, 1\}^{\ell^3})^{\ell} \rightarrow (\{0, 1\}^{O(\ell^3)})^{\ell}$ with constant distance. This construction too would have sufficed in place of Lindström-Gessel-Viennot as a starting point for the proof of Theorem 1.1, albeit with weaker parameters.

Proof of Corollary 6.1. The tree code TC_ℓ is obtained by restricting $\text{TC}_\mathbb{Z}$, defined in Section 5, to its first ℓ coordinates where, for an integer b , we identify $\{0, 1\}^b$ with $\{0, 1, \dots, 2^b - 1\}$.

It will be more convenient to start the index set of TC_ℓ from 1 rather than 0 as was done in $\text{TC}_\mathbb{Z}$. Note that the input symbols, when represented as integers, are bounded by $2^\ell - 1$. Therefore, by Theorem 1.3, for every $t \in [\ell]$, the t 'th output symbol is bounded in absolute value by

$$2^{t-1}(2^\ell - 1)^2 \leq 2^{\ell-1}(2^\ell - 1)^2 \leq 2^{3\ell-1} - 1,$$

and so 3ℓ bits suffice to represent the output symbols of $\text{TC}_\mathbb{Z}$, including the sign of γ_t . Thus TC_ℓ inherits the distance $1/2$ bound from $\text{TC}_\mathbb{Z}$. \square

For the proof of Theorem 1.1, it will be convenient to introduce a relaxed notion of tree codes which we call *lagged tree codes*. Roughly, these are tree codes that are only required to gain distance after some lag from the split.

Definition 6.2. An online function $\text{TCLag}: \Sigma_{\text{in}}^n \rightarrow \Sigma_{\text{out}}^n$ is a lagged tree code with distance δ and lag L if for every distinct $x, y \in \Sigma_{\text{in}}^n$, with $s = \text{split}(x, y)$, and every integer $L \leq \ell \leq n - s$,

$$\text{dist}(\text{TCLag}(x)_{[s, s+\ell]}, \text{TCLag}(y)_{[s, s+\ell]}) \geq \delta(\ell + 1).$$

We borrow the terminology used for tree codes for lagged tree codes. That is, we refer to n as the depth of TCLag and to $\Sigma_{\text{in}}, \Sigma_{\text{out}}$ as the input alphabet and output alphabet, respectively. We also extend Definition 3.6 to lagged tree codes in the natural way. Note that a tree code is a lagged tree code with lag $L = 0$. At the other extreme, using error correcting block codes, it is not hard to obtain explicit binary lagged tree codes with a trivial lag of ℓ (that is, the guarantee on the distance holds only after reading the entire codeword), constant distance, and $|\Sigma_{\text{out}}| = O(1)$. In the following claim, we obtain explicit binary lagged tree codes with a constant distance, $|\Sigma_{\text{out}}| = O(1)$, and depth that is quadratic in the lag.

Claim 6.3. There exists a constant $c_{\text{lag}} \geq 1$ such that for every integer $\ell \geq 1$ there exists an explicit lagged tree code $\text{TCLag}_\ell: \{0, 1\}^\ell \rightarrow (\{0, 1\}^{c_{\text{lag}}})^\ell$ with distance $1/3$ and lag $L = 16\sqrt{\ell}$.⁴

Proof. For the construction of TCLag_ℓ we make use of the following building blocks:

- Let $\text{TC}_{\sqrt{\ell}}: (\{0, 1\}^{\sqrt{\ell}})^{\sqrt{\ell}} \rightarrow (\{0, 1\}^{3\sqrt{\ell}})^{\sqrt{\ell}}$ be the tree code from Corollary 6.1. Recall that $\text{TC}_{\sqrt{\ell}}$ has distance $1/2$.
- Let $\text{ECC}: \{0, 1\}^{3\sqrt{\ell}} \rightarrow (\{0, 1\}^{c_{\text{lag}}})^{\sqrt{\ell}}$ be the error correcting block code from Lemma 3.2 set with distance $5/6$. By Lemma 3.2, c_{lag} is a constant.

⁴One can achieve distance $1/2 - \varepsilon$ for any constant $\varepsilon > 0$. This will effect the value of the constant c_{lag} and the constant multiplying $\sqrt{\ell}$ in the lag L .

Let $m \in \{0, 1\}^\ell$. Partition m to $\sqrt{\ell}$ consecutive blocks each consisting of $\sqrt{\ell}$ bits, namely, $m = (m_1, \dots, m_{\sqrt{\ell}})$ where $m_i \in \{0, 1\}^{\sqrt{\ell}}$. Similarly, for $t \in [\sqrt{\ell}]$, we write $\text{TCLag}_\ell(m)_t$ for $\text{TCLag}_\ell(m)$ projected to the t 'th block, where each block consists of $\sqrt{\ell}$ elements of $\{0, 1\}^{\text{clag}}$. Formally, $\text{TCLag}_\ell(m)_t = \text{TCLag}_\ell(m)_{[(t-1)\sqrt{\ell}+1, t\sqrt{\ell}]} \in (\{0, 1\}^{\text{clag}})^{\sqrt{\ell}}$. We define $\text{TCLag}_\ell(m)_1 = (0^{\text{clag}})^{\sqrt{\ell}}$. For $t \in \{2, \dots, \sqrt{\ell}\}$, define

$$\text{TCLag}_\ell(m)_t = \text{ECC}(\text{TC}_{\sqrt{\ell}}(m)_{t-1}),$$

where we interpret m as an element of $(\{0, 1\}^{\sqrt{\ell}})^{\sqrt{\ell}}$ when passing it to $\text{TC}_{\sqrt{\ell}}$.

Observe that TCLag_ℓ is online. We turn to show that TCLag_ℓ has distance $1/3$ and lag $16\sqrt{\ell}$. Let $x, y \in \{0, 1\}^\ell$ be distinct strings with $s = \text{split}(x, y)$. Let $d \in [16\sqrt{\ell}, \ell - s]$. Let $i_1 \in [\sqrt{\ell}]$ be the index for which the blocks $x_{i_1}, y_{i_1} \in \{0, 1\}^{\sqrt{\ell}}$ contain the split s . That is, i_1 is the split of x, y when interpreted as elements of $(\{0, 1\}^{\sqrt{\ell}})^{\sqrt{\ell}}$. Note that $i_1 = \lfloor (s-1)/\sqrt{\ell} \rfloor + 1$. Set $b = \lfloor d/\sqrt{\ell} \rfloor - 1$ and observe that block numbers $i_1 + 1, \dots, i_1 + b - 1$ are all fully contained in $[s, s + d]$. Hence, by construction, the codeword $\text{TCLag}_\ell(x)$ projected to $[s, s + d]$ contains

$$(\text{TCLag}_\ell(x)_{i_1+1}, \dots, \text{TCLag}_\ell(x)_{i_1+b-1}) = (\text{ECC}(\text{TC}_{\sqrt{\ell}}(x)_{i_1}), \dots, \text{ECC}(\text{TC}_{\sqrt{\ell}}(x)_{i_1+b-2}))$$

as a substring. Similarly, the codeword $\text{TCLag}_\ell(y)$ projected to $[s, s + d]$ contains

$$(\text{TCLag}_\ell(y)_{i_1+1}, \dots, \text{TCLag}_\ell(y)_{i_1+b-1}) = (\text{ECC}(\text{TC}_{\sqrt{\ell}}(y)_{i_1}), \dots, \text{ECC}(\text{TC}_{\sqrt{\ell}}(y)_{i_1+b-2}))$$

as a substring in the corresponding indices.

As $\text{TC}_{\sqrt{\ell}}$ is a tree code with distance $1/2$ and i_1 is the split of x, y when considered as elements of $(\{0, 1\}^{\sqrt{\ell}})^{\sqrt{\ell}}$, at least $(b-1)/2$ of the indices $i_1, \dots, i_1 + b - 2$ are such that $\text{TC}_{\sqrt{\ell}}(x)_i \neq \text{TC}_{\sqrt{\ell}}(y)_i$. As ECC has distance $5/6$, each such index i contributes $5\sqrt{\ell}/6$ to the total distance. Thus, the number of disagreements between $\text{TCLag}_\ell(x), \text{TCLag}_\ell(y)$ projected to $[s, s + d]$ is bounded below by

$$\frac{b-1}{2} \cdot \frac{5}{6} \sqrt{\ell} \geq \frac{d}{3},$$

where the last inequality follows as $d \geq 16\sqrt{\ell}$. \square

Let $\ell \leq n$ be integers. Let c_{lag} be the constant from Claim 6.3. Define the function $\text{TCLag}_\ell^n: \{0, 1\}^n \rightarrow (\{0, 1\}^{2c_{\text{lag}}})^n$ as follows. Let $m \in \{0, 1\}^n$. Write $m = (m_1, \dots, m_{2n/\ell})$ ⁵ where each $m_i \in \{0, 1\}^{\ell/2}$. For $i = 1, \dots, n/\ell$ define

$$\begin{aligned} o_i &= \text{TCLag}_\ell(m_{2i-1}, m_{2i}), \\ e_i &= \text{TCLag}_\ell(m_{2i}, m_{\min(2i+1, 2n/\ell)}). \end{aligned}$$
⁶

Note that each of o_i, e_i is an element of $(\{0, 1\}^{c_{\text{lag}}})^\ell$. Let $c_i = (o_i, e_i)$ of which we think of as an element of $(\{0, 1\}^{2c_{\text{lag}}})^\ell$. Define

$$\text{TCLag}_\ell^n(m) = (c_1, c_2, \dots, c_{n/\ell}).$$

⁵For the sake of clarity, we ignore issues of divisibility that can be easily handled.

⁶The minimum with $2n/\ell$ is taken only to make sure that we do not go out of the index set of the message.

Claim 6.4. *Let x, y be distinct n -bit strings and let $s = \text{split}(x, y)$. Assume that $s \leq n - \ell/2$. Then, for every integer $16\sqrt{\ell} \leq d \leq \ell/2$,*

$$\text{dist}(\text{TCLag}_{\ell}^n(x)_{[s, s+d]}, \text{TCLag}_{\ell}^n(y)_{[s, s+d]}) \geq d/3.$$

Proof. Let $i \in [2n/\ell]$ be the index of blocks x_i, y_i that contain the split s . Assume i is odd and let $t = s - (i-1)\ell/2$ be the index within block i of the split. Then $\text{TCLag}_{\ell}^n(x)_{[s, s+d]}$ contains, as a substring, $\text{TCLag}_{\ell}(x_i, x_{i+1})_{[t, t+d]}$, where we have used the fact that $t + d \leq \ell$ as implied by the hypothesis $d \leq \ell/2$ and since, by construction, $t \leq \ell/2$. Similarly, $\text{TCLag}_{\ell}^n(y)_{[s, s+d]}$ contains $\text{TCLag}_{\ell}(y_i, y_{i+1})_{[t, t+d]}$ in the corresponding indices. The proof follows as TCLag_{ℓ} is a lagged tree code with distance $1/3$ and lag $16\sqrt{\ell}$. A similar argument holds for even i 's. \square

We are now ready to prove Theorem 1.1 (for distance $\delta = 1/3$).

Proof of Theorem 1.1. First, observe that it suffices to construct a lagged tree code with distance $1/3$ and lag $L = O(1)$. Such a lagged tree code can be efficiently converted to a tree code with distance $1/3$ and only a constant overhead in alphabet size by including the last L inputs in the encoding at any point of time. Set $j = \log_2 \log_2 n$ and define the sequence of integers ℓ_1, \dots, ℓ_j recursively as follows: $\ell_1 = 2^{20}$ and for $i \geq 1$, $\ell_{i+1} = \ell_i^2/2^{10}$. One can verify that $\ell_i = 2^{10(2^{i-1}+1)}$ and so for every integer $2^{14} \leq d \leq n$ there exists $i \in [j]$ such that $16\sqrt{\ell_i} \leq d \leq \ell_i/2$.

We define $\text{TC}: \{0, 1\}^n \rightarrow \Sigma^n$ as follows. Let $m \in \{0, 1\}^n$. For $i \in [j]$ define $t_i = \text{TCLag}_{\ell_i}^n(m) \in (\{0, 1\}^{2^{\text{clag}_i}})^n$. Define

$$\text{TC}(m) = (t_1, t_2, \dots, t_j) \in (\{0, 1\}^{2^{\text{clag}_j}})^n.$$

Note that $\Sigma = \{0, 1\}^{2^{\text{clag}_j}}$ is an alphabet of size $(\log n)^{O(1)}$. TC is an online function as each of the functions $\text{TCLag}_{\ell_1}^n, \dots, \text{TCLag}_{\ell_j}^n$ is online. We turn to show that TC is a lagged tree code with lag 2^{14} and distance $1/3$. Let $x, y \in \{0, 1\}^n$ be distinct with $s = \text{split}(x, y)$. Let $2^{14} \leq d \leq n - s$. By the above, there exists $i \in [j]$ such that $16\sqrt{\ell_i} \leq d \leq \ell_i/2$. Hence, $\text{TCLag}_{\ell_i}^n$ guarantees that the fraction of disagreements between $\text{TC}(x)$ and $\text{TC}(y)$ projected to $[s, s+d]$ is at least $1/3$. \square

7 Tree Codes for any Distance $\delta < 1$

In this section we prove the following theorem which readily implies Theorem 1.2.

Theorem 7.1. *For every integer $r \geq 1$ there exists an explicit tree code $\text{TC}_{\mathbb{Z}}^r: \mathbb{Z}^{\mathbb{N}} \rightarrow (\mathbb{Z}^{r+1})^{\mathbb{N}}$ with distance $1 - 1/(r+1)$. Further, for every $z = (z_t)_{t \in \mathbb{N}} \in \mathbb{Z}^{\mathbb{N}}$ and $t \in \mathbb{N}$, each of the $r+1$ integers in $\text{TC}_{\mathbb{Z}}^r(z)_t$ is bounded, in absolute value, by $2^{O(t^{2r})} \cdot \max(|z_0|, \dots, |z_t|)$.*

Proof. Let $r\mathbb{N} = \{0, r, 2r, \dots\}$ be the set of all natural numbers that are divisible by r , and let $\bar{r}\mathbb{N} = \mathbb{N} \setminus r\mathbb{N}$. Formally, the tree code $\text{TC}_{\mathbb{Z}}^r$ defined next has domain $\mathbb{Z}^{r\mathbb{N}}$ but, for ease of readability, we identify $\mathbb{Z}^{r\mathbb{N}}$ with $\mathbb{Z}^{\mathbb{N}}$ in the natural way.

The construction of $\text{TC}_{\mathbb{Z}}^r$. Define the function $\text{TC}_{\mathbb{Z}}^r: \mathbb{Z}^{\mathbb{N}} \rightarrow (\mathbb{Z}^{r+1})^{\mathbb{N}}$ as follows. Given $z = (z_t)_{t \in r\mathbb{N}}$, we define a sequence of real polynomials $(f_t)_{t \in \mathbb{N}}$ and a sequence of integers $(z_t)_{t \in r\mathbb{N}}$ inductively with respect to $t \in \mathbb{N}$ as follows:

1. Define f_{tr} to be the least degree real polynomial such that $\forall i \in \{0, 1, \dots, tr\}$, $f_{tr}(i) = z_i$.
2. Define the real polynomials $f_{tr+1}, \dots, f_{tr+r-1}$ to equal f_{tr} .
3. For $i = 1, \dots, r-1$, set $z_{tr+i} = f_{tr+i}(tr+i)$.

By Lemma 3.8, there exists a sequence of integers $(\gamma_i)_{i \in \mathbb{N}}$ such that for every $t \in \mathbb{N}$, $f_t(x) = \sum_{i=0}^t \gamma_i \binom{x}{i}$. For $t \in \mathbb{N}$, define

$$\text{TC}_{\mathbb{Z}}^r(z)_t = (\gamma_{tr}, z_{tr}, z_{tr+1}, z_{tr+r-1}).$$

Analysis. Observe that the γ_t 's as well as z_t for $t \in r\mathbb{N}$ are all \mathbb{Z} -linear combination of the input sequence $(z_t)_{t \in r\mathbb{N}}$, and so $\text{TC}_{\mathbb{Z}}^r$ is a linear function with range \mathbb{Z}^{r+1} . Further, $\text{TC}_{\mathbb{Z}}^r$ is online.

We turn to analyze the distance of $\text{TC}_{\mathbb{Z}}^r$. As $\text{TC}_{\mathbb{Z}}^r$ is linear, it suffices to consider a nonzero sequence z . Assume that $c \in \mathbb{N}$ is the least integer such that $z_{cr} \neq 0$. Let $\ell \geq 0$. Set $t = c + \ell$ and define $I = \{c, c+1, \dots, t\}$. Recall that $f_{tr}(x) = \sum_{i=0}^{tr} \gamma_i \binom{x}{i}$. By (2), $\gamma_i = 0$ for every $i \in r\mathbb{N}$ and so

$$f_{tr}(x) = \sum_{i=0}^t \gamma_{ir} \binom{x}{ir}.$$

Denote the sparsity of f_{tr} by s . Note that s of the γ -entries in $(\text{TC}_{\mathbb{Z}}^r(z)_i)_{i \in I}$ are nonzero. By Lemma 1.4, f_{tr} has at most $s-1$ roots in $[cr, \infty) \cap \mathbb{Z}$ and, in particular, at most $s-1$ roots among $\{cr, cr+1, \dots, tr+r-1\}$. As the evaluation of f_{tr} on these $(\ell+1)r$ points appear as entries in $(\text{TC}_{\mathbb{Z}}^r(z)_i)_{i \in I}$, we have that at least

$$s + (\ell+1)r - (s-1) = (\ell+1)r + 1$$

of the $(\ell+1)(r+1)$ integers in $((\text{TC}_{\mathbb{Z}}^r)_i)_{i \in I}$ are nonzero. Thus, at least

$$\frac{(\ell+1)r + 1}{r+1} > \left(1 - \frac{1}{r+1}\right) (\ell+1)$$

of the indices $i \in I$ are such that $\text{TC}_{\mathbb{Z}}^r(z)_i$ is nonzero as an $(r+1)$ -tuple, establishing the desired bound on the distance.

We turn to bound the output symbols. We start by bounding the γ symbols. Recall that $\gamma_i = 0$ for every i not divisible by r . For $t \in \mathbb{N}$, let $\Gamma_t = \max(|\gamma_i|: i \in \{0, 1, \dots, t\})$ and

define $m_{tr} = \max(|z_i|: i \in r\mathbb{N} \cap [0, tr])$. By Lemma 3.8, for every $t \geq 1$ we have that

$$\begin{aligned}\gamma_{tr} &= \sum_{i=0}^{tr} (-1)^{tr-i} \binom{tr}{i} z_i \\ &= z_{tr} + \sum_{i=0}^{tr-1} (-1)^{tr-i} \binom{tr}{i} \sum_{j=0}^i \gamma_j \binom{i}{j},\end{aligned}$$

and so

$$\begin{aligned}|\gamma_{tr}| &\leq |z_{tr}| + \sum_{i=0}^{tr-1} \binom{tr}{i} \sum_{j=0}^i |\gamma_j| \binom{i}{j} \\ &\leq |z_{tr}| + \sum_{i=0}^{tr-1} \binom{tr}{i} \Gamma_i 2^i \\ &\leq |z_{tr}| + \Gamma_{tr-1} \sum_{i=0}^{tr} \binom{tr}{i} 2^i \\ &= |z_{tr}| + 3^{tr} \Gamma_{tr-1}.\end{aligned}$$

As $\Gamma_{tr-1} = \Gamma_{(t-1)r}$ and since, being an input symbol, $|z_{tr}| \leq m_{tr}$, we have that $\Gamma_{tr} \leq 3^{tr} \Gamma_{(t-1)r} + m_{tr}$, which implies $\Gamma_{tr} \leq m_{tr} \cdot 3^{t^2 r}$.

As for the computed z values, by (2) and (3), for $t \in \mathbb{N}$ and $k \in \{1, 2, \dots, r-1\}$,

$$\begin{aligned}z_{tr+k} &= f_{tr+k}(tr+k) \\ &= f_{tr}(tr+k) \\ &= \sum_{i=0}^t \gamma_{ir} \binom{tr+k}{ir} \\ &\leq m_{tr} \cdot 3^{t^2 r} \cdot 2^{tr+k},\end{aligned}$$

which completes the proof of Theorem 7.1. □

Corollary 7.2. *There exists a universal constant $c \geq 1$ such that for every $\zeta > 0$ and every integer $\ell \geq 1$ there exists an explicit tree code*

$$\text{TC}_\ell: \left(\{0, 1\}^{\ell^2}\right)^\ell \rightarrow \left(\{0, 1\}^{\frac{c}{\zeta^2} \cdot \ell^2}\right)^\ell$$

with distance $\delta = 1 - \zeta$.

Proof. The tree code TC_ℓ is obtained by restricting $\text{TC}_\mathbb{Z}^r$, set with $r = \lceil 1/\zeta \rceil$, to its first ℓ coordinates where, as in the proof of Corollary 6.1, for an integer b , we identify $\{0, 1\}^b$ with $\{0, 1, \dots, 2^b - 1\}$. Note that the input symbols, when represented as integers, are bounded

by 2^{ℓ^2} . Therefore, by Theorem 7.1, there exists a universal constant $c \geq 1$ such that every output symbol is bounded in absolute value by $2^{c(1/\zeta)^2 \ell^2}$ and so $c(1/\zeta)^2 \ell^2$ bits suffice to represent the output symbols of $\text{TC}_{\mathbb{Z}}^r$. Clearly, TC_ℓ inherits the $1 - 1/(r+1) \geq 1 - \zeta$ distance bound from $\text{TC}_{\mathbb{Z}}^r$. \square

The reduction in Section 6, instantiated with Corollary 7.2 instead of Corollary 6.1, and when executed with a suitable choice of parameters, yields a binary tree code with distance $\delta = 1 - \zeta$ and alphabet size $(\log n)^{O(1/\zeta^2)}$.

Acknowledgments

The second author thanks Noga Ron-Zewi and Shachar Lovett for many helpful discussions.

References

- [ABE⁺16] Noga Alon, Mark Braverman, Klim Efremenko, Ran Gelles, and Bernhard Haeupler. Reliable communication over highly connected noisy networks. *ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing (PODC)*, 61, 2016.
- [AGS16] Shweta Agrawal, Ran Gelles, and Amit Sahai. Adaptive protocols for interactive communication. *IEEE International Symposium on Information Theory (ISIT)*, pages 595–599, 2016.
- [Aig07] M. Aigner. *A course in enumeration*. Springer, 2007.
- [AZ10] M. Aigner and G. Ziegler. *Proofs from the Book*. Springer, 4th edition, 2010.
- [BE14] Mark Braverman and Klim Efremenko. List and unique coding for interactive communication in the presence of adversarial noise. In *Proceedings of the IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 236–245, 2014.
- [BGMO17] Mark Braverman, Ran Gelles, Jieming Mao, and Rafail Ostrovsky. Coding for interactive communication correcting insertions and deletions. *IEEE Transactions on Information Theory*, 63(10):6256–6270, 2017.
- [BK12] Zvika Brakerski and Yael Tauman Kalai. Efficient interactive coding against adversarial noise. *Proceedings of the IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 160–166, 2012.
- [BKN14] Zvika Brakerski, Yael Tauman Kalai, and Moni Naor. Fast interactive coding against adversarial noise. *Journal of the ACM (JACM)*, 61(6):35:1–35:30, 2014.

- [BN13] Zvika Brakerski and Moni Naor. Fast algorithms for interactive coding. In *Proceedings of the ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 443–456, 2013.
- [BR14] M. Braverman and A. Rao. Toward coding for maximum errors in interactive communication. *IEEE Transactions on Information Theory*, 60(11):7248–7255, 2014.
- [Bra12] Mark Braverman. Towards deterministic tree code constructions. In *Proceedings of the ACM-SIGACT Innovations in Theoretical Computer Science Conference (ITCS)*, pages 161–167, 2012.
- [CGH13] Karthekeyan Chandrasekaran, Navin Goyal, and Bernhard Haeupler. Deterministic algorithms for the lovász local lemma. *SIAM Journal on Computing*, 42(6):2132–2155, 2013.
- [CST11] G. Cohen, A. Shpilka, and A. Tal. On the degree of univariate polynomials over the integers. In *Electronic Colloquium on Computational Complexity (ECCC)*, number 002, 2011.
- [DL16] Yevgeniy Dodis and Allison Bishop Lewko. Interactive coding for interactive proofs. *International Theory of Cryptography Conference (TCC)*, pages 352–366, 2016.
- [Dru11] Andrew Drucker. Efficient probabilistically checkable debates. In *Proceedings of the International Workshop on Randomization and Computation (RANDOM)*, pages 519–529, 2011.
- [EL75] P. Erdős and L. Lovász. Problems and results on 3-chromatic hypergraphs and some related questions. In *Infinite and Finite Sets*. North Holland, 1975.
- [FGOS15] Matthew Franklin, Ran Gelles, Rafail Ostrovsky, and Leonard J. Schulman. Optimal coding for streaming authentication and interactive communication. *IEEE Transactions on Information Theory*, 61(1):133–145, 2015.
- [Gel17] Ran Gelles. Coding for interactive communication: A survey. *Foundations and Trends in Theoretical Computer Science*, 13(12):1–157, 2017.
- [GH14] Mohsen Ghaffari and Bernhard Haeupler. Optimal Error Rates for Interactive Coding II: Efficiency and List Decoding. In *Proceedings of the IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 394–403, 2014.

- [GHK⁺16] Ran Gelles, Bernhard Haeupler, Gillat Kol, Noga Ron-Zewi, and Avi Wigderson. Towards optimal deterministic coding for interactive communication. *Proceedings of the ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2016.
- [GHS14] Mohsen Ghaffari, Bernhard Haeupler, and Madhu Sudan. Optimal error rates for interactive coding I: Adaptivity and other settings. In *Proceedings of the ACM Symposium on Theory of Computing (STOC)*, pages 794–803, 2014.
- [GMS14] Ran Gelles, Ankur Moitra, and Amit Sahai. Efficient coding for interactive communication. *IEEE Transactions on Information Theory*, 60(3):1899–1913, 2014.
- [Gro16] Leefke Grosjean. *Practical anytime codes*. PhD thesis, KTH Royal Institute of Technology, 2016.
- [GS95] A. Garcia and H. Stichtenoth. A tower of Artin-Schreier extensions of function fields attaining the Drinfeld-Vladut bound. *Inventiones Mathematicae*, 121(1):211–222, 1995.
- [GV85] I. Gessel and G. Viennot. Binomial determinants, paths, and hook length formulae. *Advances in Mathematics*, 58(3):300–321, 1985.
- [Hae14] Bernhard Haeupler. Interactive Channel Capacity Revisited. In *Proceedings of the IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 226–235, 2014.
- [HS17] B. Haeupler and A. Shahrasbi. Synchronization strings: Codes for insertions and deletions approaching the singleton bound. *ACM Symposium on Theory of Computing (STOC)*, 2017.
- [HS18] B. Haeupler and A. Shahrasbi. Synchronization strings: Explicit constructions, local decoding, and applications. *ACM Symposium on Theory of Computing (STOC)*, 2018.
- [HSV17] B. Haeupler, A. Shahrasbi, and E. Vitercik. Synchronization strings: Channel simulations and interactive coding for insertions and deletions. *arXiv:1707.04233*, 2017.
- [HV17] Bernhard Haeupler and Ameya Velingker. Bridging the capacity gap between interactive and one-way communication. In *Proceedings of the ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 2123–2142, 2017.

- [JKL15] Abhishek Jain, Yael Tauman Kalai, and Allison Bishop Lewko. Interactive coding for multiparty protocols. In *Proceedings of the ACM Innovations in Theoretical Computer Science Conference (ITCS)*, pages 1–10, 2015.
- [KHH16] Anatoly Khina, Wael Halbawi, and Babak Hassibi. (Almost) practical tree codes. In *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, pages 2404–2408, 2016.
- [KR13] Gillat Kol and Ran Raz. Interactive channel capacity. In *Proceedings of the ACM Symposium on Theory of Computing (STOC)*, pages 715–724, 2013.
- [LdMM12] James R. Lee, Arnaud de Mesmay, and Mohammad Moharrami. Dimension reduction for finite trees in L_1 . In *Proceedings of the ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 43–50, 2012.
- [Lin73] B. Lindström. On the vector representations of induced matroids. *Bulletin of the London Mathematical Society*, 5(1):85–90, 1973.
- [LV15] Allison Lewko and Ellen Vitercik. Balancing communication for multi-party interactive coding. *CoRR*, abs/1503.06381, 2015.
- [MS14] Cristopher Moore and Leonard J. Schulman. Tree codes and a conjecture on exponential sums. In *Proceedings of the ACM-SIGACT Innovations in Theoretical Computer Science Conference (ITCS)*, pages 145–154, 2014.
- [MT10] Robin A Moser and Gábor Tardos. A constructive proof of the general lovász local lemma. *Journal of the ACM (JACM)*, 57(2):11, 2010.
- [ORS09] R. Ostrovsky, Y. Rabani, and L. J. Schulman. Error-correcting codes for automatic control. *IEEE Trans. Information Theory*, 55(7):2931–2941, 2009.
- [Pec06] Marcin Peczarski. An improvement of the tree code construction. *Information Processing Letters (IPL)*, 99(3):92–95, 2006.
- [Pud16] Pavel Pudlák. Linear tree codes and the problem of explicit constructions. *Linear Algebra and its Applications*, 490:124–144, 2016.
- [RS94] Sridhar Rajagopalan and Leonard J. Schulman. A coding theorem for distributed computation. In *Proceedings of the ACM Symposium on Theory of Computing (STOC)*, pages 790–799, 1994.
- [Sah01] Anant Sahai. *Anytime information theory*. PhD thesis, Massachusetts Institute of Technology, 2001.

- [Sah04] Anant Sahai. The necessity and sufficiency of anytime capacity for control over a noisy communication link. In *Proceedings of the IEEE Conference on Decision and Control*, volume 2, pages 1896–1901, 2004.
- [Sch93] Leonard J. Schulman. Deterministic coding for interactive communication. In *Proceedings of the ACM Symposium on Theory of Computing (STOC)*, pages 747–756, 1993.
- [Sch94] Leonard J. Schulman. Postscript of 21 September 2003 to Coding for Interactive Communication. <http://users.cms.caltech.edu/~schulman/Papers/intercodingpostscript.txt>, 1994.
- [Sch96] Leonard J. Schulman. Coding for interactive communication. *IEEE Transactions on Information Theory*, 42(6):1745–1756, 1996.
- [SH11a] Ravi Teja Sukhavasi and Babak Hassibi. Anytime reliable codes for stabilizing plants over erasure channels. In *Proceedings of the IEEE Conference on Decision and Control and European Control Conference (CDC-ECC)*, pages 5254–5259, 2011.
- [SH11b] Ravi Teja Sukhavasi and Babak Hassibi. Linear error correcting codes with anytime reliability. In *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, pages 1748–1752, 2011.
- [SM06] Anant Sahai and Sanjoy Mitter. The necessity and sufficiency of anytime capacity for stabilization of a linear system over a noisy communication link. *IEEE Transactions on Information Theory (TransInf)*, 52(8):3369–3395, 2006.
- [ST11] Amir Shpilka and Avishay Tal. On the minimal fourier degree of symmetric boolean functions. In *2011 IEEE 26th Annual Conference on Computational Complexity (CCC)*, pages 200–209. IEEE, 2011.
- [Sti09] H. Stichtenoth. *Algebraic function fields and codes*, volume 254. Springer Science & Business Media, 2009.
- [SW17] Alexander A. Sherstov and Pei Wu. Optimal interactive coding for insertions, deletions, and substitutions. In *IEEE Symposium on Foundations of Computer Science (FOCS)*, 2017.
- [vZGR97] Joachim von Zur Gathen and James R Roche. Polynomials with two values. *Combinatorica*, 17(3):345–362, 1997.
- [Wig17] A. Wigderson. *Mathematics and Computation*. 2017. Book draft.

A Proof of Lemma 3.2

For the proof of Lemma 3.2, we make use of algebraic-geometric codes.

Theorem A.1 ([GS95]. See also [Sti09]). *Let p be a prime number and $m \in \mathbb{N}$ even. Set $q = p^m$. For every $0 < \rho < 1$ and a large enough integer n , there exists an explicit rate ρ linear error correcting block code $\text{ECC}: \mathbb{F}_q^n \rightarrow \mathbb{F}_q^{n/\rho}$ with distance*

$$\delta \geq 1 - \rho - \frac{1}{\sqrt{q} - 1}.$$

Proof of Lemma 3.2. Let ℓ be the least integer such that $2^\ell \geq 2/\varepsilon + 1$ and let $q = 2^{2^\ell}$. Let m be the least integer such that $1/m \leq \varepsilon/2$. By Theorem A.1, there exists an explicit error correcting block code $\text{ECC}': \mathbb{F}_q^n \rightarrow \mathbb{F}_q^{mn}$ with distance $1 - \varepsilon$. Identify $\mathbb{F}_q = \mathbb{F}_{2^{2^\ell}}$ with $\{0, 1\}^{2^\ell}$ by fixing a representation for \mathbb{F}_q . Note that $q = O(1/\varepsilon^2) = O(1)$ and so such a representation can be computed in constant time. Set $c = 2\ell tm$. Define, $\text{ECC}: \{0, 1\}^n \rightarrow (\{0, 1\}^c)^{n/t}$ as follows. Given $x \in \{0, 1\}^n$, identify x with an element of $\mathbb{F}_2^n \subseteq \mathbb{F}_q^n$. For $i \in [n/t]$, define

$$\text{ECC}(x)_i = (\text{ECC}'(x)_{(i-1)mt+1}, \text{ECC}'(x)_{(i-1)mt+2}, \dots, \text{ECC}'(x)_{imt}).$$

Note that $c = O((t/\varepsilon) \log(1/\varepsilon))$.

Consider distinct $x, y \in \{0, 1\}^n$ and, as before, identify $\{0, 1\}^n$ with $\mathbb{F}_2^n \subseteq \mathbb{F}_q^n$. Then, the codewords $\text{ECC}'(x), \text{ECC}'(y)$ agree on at most ε -fraction of the coordinates. This readily implies that $\text{ECC}(x), \text{ECC}(y)$ agree on at most ε -fraction of the coordinates. \square