

Error Correction by Structural Simplicity: Correcting Samplable Additive Errors

Kenji Yasunaga
Kanazawa University
yasunaga@se.kanazawa-u.ac.jp

August 26, 2018

Abstract

This paper explores the possibilities and limitations of error correction by the structural simplicity of error mechanisms. Specifically, we consider channel models, called *samplable additive channels*, in which (a) errors are efficiently sampled without the knowledge of the coding scheme or the transmitted codeword; (b) the entropy of the error distribution is bounded; and (c) the number of errors introduced by the channel is unbounded. For the channels, several negative and positive results are provided. Assuming the existence of one-way functions, there are samplable additive errors of entropy n^ϵ for $\epsilon \in (0, 1)$ that are pseudorandom, and thus not correctable by efficient coding schemes. It is shown that there is an oracle algorithm that induces a samplable distribution over $\{0, 1\}^n$ of entropy $m = \omega(\log n)$ that is not pseudorandom, but is uncorrectable by efficient schemes of rate less than $1 - m/n - o(1)$. The results indicate that restricting error mechanisms to be efficiently samplable and not pseudorandom is insufficient for error correction. As positive results, some conditions are provided under which efficient error correction is possible.

1 Introduction

In the theory of error-correcting codes, two of the most-studied channel models are probabilistic channels and worst-case channels. In probabilistic channels, errors are introduced through stochastic processes, and the most well-known one is the binary symmetric channel (BSC). In worst-case (or adversarial) channels, errors are introduced adversarially by considering the choice of coding schemes and transmitted strings, called *codewords*, under the restriction of the *error rate*, which is the ratio of the number of errors to the length of the codeword. In his seminal work [28], Shannon showed that reliable communication can be achieved over BSC if the *coding rate*, which represents the efficiency of transmission, is less than $1 - h_2(p)$, where $h_2(\cdot)$ is the binary entropy function, and p is the crossover probability of BSC. In contrast, it is known that reliable communication cannot be achieved over worst-case channels when the error rate is at least $1/4$ unless the coding rate tends to zero [25].

If we view the introduction of errors as *computation* of the channel, probabilistic channels perform low-cost computation with little knowledge about the coding scheme and the transmitted codeword, while worst-case channels perform high-cost computation with the full-knowledge. As intermediate channels between probabilistic channels and worst-case channels, Lipton [21] introduced *computationally-bounded channels*, where errors are introduced by polynomial-time computation. It

has been demonstrated that reliable communication for these intermediate channels can be achieved by more efficient schemes than those for worst-case ones [21, 24, 15].

Although there have been many studies on error correction in various channel models, the most basic principle for error correction has been the same: the number of errors occurred in communications is restricted. Namely, an upper bound on the error rate is a priori provided, and coding schemes are designed with the knowledge of the bound. However, since qualitatively better error correction is possible for computationally-bounded channels, it may be possible to correct errors without resorting to the bound on the error rate. Thus, we ask the following question:

Is it possible to correct errors based on the structural simplicity of error mechanisms?

In this work, we partially answer the above question. Specifically, we introduce a novel channel model, called *samplable additive channel*, and investigate the error-correction capabilities for the channel. In samplable additive channels, (a) errors are efficiently sampled without the knowledge of the coding scheme or the transmitted codeword; (b) the entropy of the error distribution is bounded; and (c) the error rate is unbounded. Condition (a) means that there is a polynomial-time algorithm that generates samples according to the error distribution, where the algorithm is designed without the knowledge of the coding scheme, and does not receive the transmitted codeword as input. The restriction of efficient samplability is also employed in the previous studies of computationally-bounded channels. For a structural simplicity, condition (b) is necessary since high-entropy distributions can generate unpredictable and complex errors. Condition (c) is employed for removing the effect of bounding the error rate for error correction. In addition, for the ease of analysis, we mainly consider *flat* distributions, which are the uniform distributions over the supports.

Samplable additive channels are relatively simple channel models since the error distributions are identical for every coding scheme and transmitted codeword. The binary symmetric channel is an example. The setting is incomparable to previous notions of error correction against computationally-bounded channels. Our model is stronger because we do not restrict the error rate, but is weaker because the channel cannot see the code or the transmitted codeword.

1.1 The Setting

Before presenting our results, we briefly describe our problem setting.

The coding scheme, also referred to as code, consists of two functions $\text{Enc} : \{0, 1\}^{Rn} \rightarrow \{0, 1\}^n$ and $\text{Dec} : \{0, 1\}^n \rightarrow \{0, 1\}^{Rn}$, where $R \in [0, 1]$ is called *coding rate*, or simply *rate*. A message $x \in \{0, 1\}^{Rn}$ is encoded to $\text{Enc}(x)$, and transmitted to the channel. The channel introduces an error $z \in \{0, 1\}^n$, and the decoder receives the string $\text{Enc}(x) + z$, where $+$ is the bit-wise addition modulo 2, and outputs $\tilde{x} = \text{Dec}(\text{Enc}(x) + z)$. Decoding is done successfully if $\tilde{x} = x$. It is desirable to construct a high-rate code that can successfully correct errors. The *error rate* is $w_z/n \in [0, 1]$, where w_z is the number of non-zero elements in z .

For every samplable additive channel, an error distribution Z over $\{0, 1\}^n$ is associated. If Z is flat, each element in the support of Z is sampled with probability $1/|Z|$, and thus the entropy of Z is $\log |Z|$, which takes values in $[0, n]$. In our setting, we assume that sampling from Z can be simulated by a polynomial-time algorithm, the entropy of Z is bounded, and the error rate is not bounded. In addition, the coding scheme can be chosen based on the knowledge of Z , and an error vector z is sampled from Z without knowing $\text{Enc}(x)$. We would like to know which Z is (un)correctable, especially by efficient coding schemes.

1.2 Main Results

We investigate the error correction capabilities of samplable additive errors.

Errors from flat distributions. By simple combinatorial arguments, it can be shown that any flat Z can be corrected by a code of rate $R \leq 1 - m/n - o(1)$, but cannot for rate $R > 1 - m/n - o(1)$, where m is the entropy of Z . Thus, the rate $1 - m/n$ is essentially optimal. In addition, if $m = O(\log n)$ and $R \leq 1 - m/n - o(1)$, both the encoding and the decoding can be done in polynomial time.

Pseudorandom errors. Assuming the existence of one-way functions, there exists pseudorandom generators [17], which generate distributions that look random to every efficient algorithm. Hence, no efficient scheme can correct errors from such distributions. It follows from this fact that, assuming the existence of one-way functions, there exists an error distribution Z with entropy n^ϵ for any constant $\epsilon \in (0, 1)$ that are not efficiently correctable.

Errors with membership test. To avoid the impossibility of correcting pseudorandom errors, we consider samplable distributions for which membership test can be done efficiently. Such distributions are not pseudorandom since the membership test can be used to distinguish them from the uniform distribution.

We show the existence of an uncorrectable distribution with membership test for *low-rate* codes. As sampling algorithms, we employ *oracle algorithms* [2], which can make a black-box use of an external entity called *oracle*. We present an oracle algorithm that induces a samplable distribution Z of entropy $m = \omega(\log n)$ that is not correctable by efficient coding schemes of rate $R < 1 - m/n - o(1)$. The result complements the impossibility of correcting flat distributions for rate $R > 1 - m/n - o(1)$. Also, the entropy of $\omega(\log n)$ is optimal since any flat Z with entropy $O(\log n)$ can be corrected by a polynomial-time coding scheme.

Positive results. As a positive result, we show that if the set of error vectors forms a linear subspace, then every error in the set can be corrected by an efficient coding scheme with optimal rate.

Also, we derive a computational condition under which samplable additive errors can be corrected. Intuitively, the condition is that a variant of the sampling algorithm of Z is efficiently “invertible”. See Section 6.2 for the details. The result implies that the existence of one-way functions is necessary for proving the impossibility results for correcting samplable errors.

We summarize our main results in Table 1, where R denotes the rate of coding schemes, and m the entropy of Z .

1.3 Related Work

The notion of computationally-bounded channel was introduced by Lipton [21]. He showed that if the sender and the receiver can share secret randomness, then the Shannon capacity can be achieved for this channel. Micali et al. [24] considered a similar channel model in a public-key setting, and gave a coding scheme based on list-decodable codes and digital signature. Guruswami and Smith [15] gave constructions of capacity achieving codes for worst-case additive-error channel and time/space-bounded channels. They also gave an impossibility result for bit-

Table 1: Correctability of Samplable Additive Error Z

Entropy m	Correctabilities	Assumptions	References
—	\forall flat Z , (1) \exists code correcting Z for $R \leq 1 - m/n - o(1)$ in time $O(n^2 2^m)$ (2) not correctable for $R > 1 - m/n - o(1)$	None	Proposition 3
n^ϵ ($0 < \epsilon < 1$)	$\exists Z$ not efficiently correctable for any R	OWF	Proposition 6
$\omega(\log n)$	$\exists Z$ with membership test, not efficiently correctable for $R < 1 - m/n - o(1)$	Oracle access	Corollary 1
—	\forall flat Z over a linear subspace of dim. m , \exists code correcting Z for $R \leq 1 - m/n$	None	Proposition 7
—	\forall flat $Z = f(U_r)$ is efficiently correctable for $R \leq 1 - m/n - o(1)$	g is not distOWF	Theorem 2

fixing channels, but their result can be applied to channels that use the information on the code and the transmitted codeword. Shaltiel and Silbak [27] gave explicit list-decodable codes for computationally-bounded channels based on complexity assumptions. Cryptographic treatment of codes against computationally-bounded channels was studied in [35]. Note that all the previous work of computationally-bounded channels assumes that the error rate is bounded above by some constant $p \in [0, 1]$, and codes can be constructed based on the knowledge of p .

Additive-error channels have been studied in the literature [7, 8, 20, 15]. In the previous studies, the error rate is bounded, and the channel cannot see the transmitted codeword, but can depend on the coding scheme. In the present study, stronger obliviousness is considered, in which the channel cannot depend on the code.

Samplable distributions were also studied in the context of data compression [14, 30, 34], randomness extractor [29, 33, 9], and randomness condenser [10]. Samplable distributions with membership test appeared in the study of efficient compressibility of samplable sources [14, 30, 34].

1.4 Organization

In Section 2, we give the formal model of the problem, and introduce several notions of error-correcting codes. Several results on the correctabilities of flat distributions are presented in Section 3. The negative result of correcting pseudorandom errors appears in Section 4. In Section 5, we show the existence of an error distribution with membership test that is not efficiently correctable. The positive results are provided in Section 6.

2 Preliminaries

For $n \in \mathbb{N}$, we write $[n]$ as the set $\{1, 2, \dots, n\}$. For a distribution X , we write $x \sim X$ to indicate that x is chosen according to X . We may use X also as a random variable distributed according to X . The *support* of X is $\text{Supp}(X) = \{x : \Pr_X(x) \neq 0\}$, where $\Pr_X(x)$ is the probability that X assigns to x . The *Shannon entropy* of X is given by $-\sum_{x \in \text{Supp}(X)} \Pr_X(x) \log \Pr_X(x)$. For flat distributions, the Shannon entropy is equal to the *min-entropy*, which is defined to be

$\min_{x \in \text{Supp}(X)} \{-\log \Pr_X(x)\}$. Thus, we simply say that a flat distribution Z has entropy m if its Shannon entropy is m . For $n \in \mathbb{N}$, we write U_n as the uniform distribution over $\{0, 1\}^n$.

We define the notion of additive-error correcting codes.

Definition 1 (Additive-error correcting codes). *For two functions $\text{Enc} : \mathbb{F}^k \rightarrow \mathbb{F}^n$ and $\text{Dec} : \mathbb{F}^n \rightarrow \mathbb{F}^k$, and a distribution Z over \mathbb{F}^n , where $k \leq n$ and \mathbb{F} is a finite field, we say (Enc, Dec) corrects (additive error) Z with error probability ϵ if for any $x \in \mathbb{F}^k$, we have that*

$$\Pr_{z \sim Z} [\text{Dec}(\text{Enc}(x) + z) \neq x] \leq \epsilon.$$

When $\epsilon = 0$, we simply say (Enc, Dec) corrects Z . The rate of (Enc, Dec) is k/n .

Definition 2. *A distribution Z is said to be correctable with rate R and error probability ϵ if there is a pair of functions (Enc, Dec) of rate R that corrects Z with error probability ϵ .*

We call a pair (Enc, Dec) a *coding scheme* or simply *code*. The coding scheme is called *efficient* if Enc and Dec can be computed in polynomial-time in n . The code is called *linear* if Enc is a linear mapping, that is, for any $x, y \in \mathbb{F}^k$ and $a, b \in \mathbb{F}$, $\text{Enc}(ax + by) = a \text{Enc}(x) + b \text{Enc}(y)$. If $|\mathbb{F}| = 2$, we may use $\{0, 1\}$ instead of \mathbb{F} . For any linear code (Enc, Dec) , there is a matrix $G \in \mathbb{F}^{k \times n}$, called a *generator matrix*, such that $\text{Enc}(x) = x \cdot G$ for all $x \in \mathbb{F}^k$. We usually assume that G has full rank, namely, the rank of G is k . A matrix $H \in \mathbb{F}^{(n-k) \times n}$ is called a *parity-check matrix* if for any $c \in \mathbb{F}^n$, $c = \text{Enc}(x)$ for some $x \in \mathbb{F}^k$ if and only if $Hc^T = 0$. Then, it holds that $GH^T = 0$. It is well-known that given a parity-check matrix $H \in \mathbb{F}^{(n-k) \times n}$ of a code, one can compute a generator matrix G of the code by finding a basis of the kernel of H . See [23, 26] for the basic properties of linear codes.

Next, we define syndrome decoding for linear codes. Suppose $v = \text{Enc}(x) + z \in \mathbb{F}^n$ is received, where $\text{Enc}(x)$ is the transmitted codeword and z is the error vector. Syndrome decoder associates e with $v \cdot H^T$, called the *syndrome* of v , because it holds that $v \cdot H^T = (\text{Enc}(x) + z) \cdot H^T = x \cdot G \cdot H^T + z \cdot H^T = z \cdot H^T$. If there is a way for recovering z from $z \cdot H^T$, we can recover x from v .

Definition 3. *For a linear code (Enc, Dec) , Dec is said to be a syndrome decoder if there is a function $\text{Rec} : \mathbb{F}^{n-k} \rightarrow \mathbb{F}^n$ such that $\text{Dec}(v) = (v - \text{Rec}(v \cdot H^T)) \cdot G^{-1}$, where G and H are a generator matrix and a parity-check matrix of the code, respectively, and $G^{-1} \in \mathbb{F}^{n \times k}$ is a right inverse matrix of G (i.e., $GG^{-1} = I$).*

Note that for any full-rank matrix G , there always exist a right inverse matrix of G . In the above definition, Rec recovers the error vector e from the syndrome $v \cdot H^T$. Since $v - e = \text{Enc}(x)$, x is obtained by multiplying G^{-1} by $\text{Enc}(x) = x \cdot G$.

We consider a computationally-bounded analogue of additive-error channels. We introduce the notion of samplable distributions.

Definition 4. *A distribution family $Z = \{Z_n\}_{n \in \mathbb{N}}$ is said to be samplable if there is a probabilistic polynomial-time algorithm S such that $S(1^n)$ is distributed according to Z_n for every $n \in \mathbb{N}$, where 1^n is the string consisting of n ones.*

We consider the setting in which coding schemes can depend on the sampling algorithm of Z , but not on its random coins, and Z does not use any information on the coding scheme or transmitted codewords. In this setting, the randomization of coding schemes does not help much.

Proposition 1. *Let (Enc, Dec) be a randomized coding scheme that corrects a distribution Z with error probability ϵ . Then, there is a deterministic coding scheme that corrects Z with error probability ϵ .*

Proof. Assume that Enc uses at most ℓ -bit randomness. Since (Enc, Dec) corrects Z with error probability ϵ , we have that for every $x \in \mathbb{F}^k$, $\Pr_{z \sim Z, r \sim U_\ell}[\text{Dec}(\text{Enc}(x; r) + z) \neq x] \leq \epsilon$. By the averaging argument, for every $x \in \mathbb{F}^k$, there exists $r_x \in \{0, 1\}^\ell$ such that $\Pr_{z \sim Z}[\text{Dec}(\text{Enc}(x; r_x) + z) \neq x] \leq \epsilon$. Thus, by defining $\text{Enc}'(x) = \text{Enc}(x; r_x)$, the deterministic coding scheme $(\text{Enc}', \text{Dec})$ corrects Z with error probability ϵ . \square

The above result reveals the crucial difference between our setting and that of Guruswami and Smith [15], where the channels can use the information on the coding scheme, but not transmitted codewords. They present a randomized coding scheme with optimal rate $1 - h_2(p)$ for worst-case additive-error channels, for which deterministic coding schemes are only known to achieve rate $1 - h_2(2p)$ for $p < 1/2$, where p is the error rate of the channels.

3 Errors from Flat Distributions

We investigate the correctability of general flat distributions over $\{0, 1\}^n$.

First, we show that, for any flat distribution Z of entropy m , a random linear code of length n and rate R can correct Z with error probability $1 - 2^{-\Omega(n)}$ for $R < 1 - m/n - o(1)$. Consider a random linear code of length n and rate R such that the parity check matrix H is chosen uniformly at random from $\mathcal{H}_{n,R} = \{0, 1\}^{(n-Rn) \times n}$, and a generator matrix $G \in \{0, 1\}^{Rn \times n}$ is obtained by finding a basis of the kernel of H . We use the syndrome decoding. Specifically, for a received word v , the function Rec of the decoder is defined such that it maps $v \cdot H^T$ to unique $z \in \text{Supp}(Z)$ satisfying $v \cdot H^T = z \cdot H^T$ by searching for all possible z . If there are multiple candidates for z , it outputs the decoding failure. Let $\mathcal{C}_{n,R}$ be the set of codes in which each code is defined by using each element in $\mathcal{H}_{R,n}$ as a parity-check matrix.

Proposition 2. *For any flat distribution Z over $\{0, 1\}^n$ of entropy m , a $(1 - \sqrt{\epsilon})$ -fraction of codes in $\mathcal{C}_{n,R}$ corrects Z with error probability $\sqrt{\epsilon}$ for $R < 1 - m/n$, where $\epsilon = 2^{-n(1-R-m/n)}$.*

Proof. The probability that a random code from $\mathcal{C}_{n,R}$ fails the decoding is

$$\begin{aligned} & \Pr_{H \in \mathcal{H}_{n,R}, z \sim Z} [\exists z' \in \text{Supp}(Z) \setminus \{z\} : z \cdot H^T = z' \cdot H^T] \\ &= \Pr_{H,z} [\exists z' \in \text{Supp}(Z) \setminus \{z\} : \forall i \in [n - Rn], h_i \cdot (z - z') = 0] \end{aligned} \quad (1)$$

$$\begin{aligned} &= \sum_{z \in \text{Supp}(Z)} \Pr[z \sim Z] \Pr_H \left[\bigcup_{z' \in \text{Supp}(Z) \setminus \{z\}} [\forall i \in [n - Rn], h_i \cdot (z - z') = 0] \right] \\ &\leq \sum_{z \in \text{Supp}(Z)} \Pr[z \sim Z] \sum_{z' \in \text{Supp}(Z) \setminus \{z\}} \Pr_H [\forall i \in [n - Rn], h_i \cdot (z - z') = 0] \\ &= \sum_{z \in \text{Supp}(Z)} 2^{-m} \cdot (2^m - 1) \cdot \left(\frac{1}{2}\right)^{n-Rn} \end{aligned} \quad (2)$$

$$\leq \epsilon, \quad (3)$$

where $H^T = (h_1^T, \dots, h_{n-Rn}^T)$ in (1), the first inequality follows from a union bound, and (2) follows from the facts that for any non-zero $a \in \{0, 1\}^n$ and random $h \in \{0, 1\}^n$, $\Pr_h[h \cdot a = 0] = 1/2$, and that $|\text{Supp}(Z)| = 2^m$ for flat Z . Then, a $(1 - \sqrt{\epsilon})$ -fraction of codes in $\mathcal{C}_{n,R}$ can correct Z with error probability $\sqrt{\epsilon}$, since otherwise (3) does not hold. Therefore, the statement follows. \square

Thus, we have the following proposition.

Proposition 3. *Let Z be any flat distribution over $\{0, 1\}^n$ of entropy m . There is a linear code of rate R that corrects Z with error probability ϵ for $R = 1 - m/n - 2\log(\epsilon^{-1})/n$. The decoding complexity is at most $O(n^2 2^m)$.*

Proof. The existence of such a code immediately follows from Proposition 2. Given a received word v , the brute-force decoder checks if $(v - z) \cdot H^T = 0$ for all $z \in \text{Supp}(Z)$, where H is the parity check matrix. If so, output x satisfying $x \cdot G = v - z$. Thus, the decoding is done in time $O(n^2) \cdot |\text{Supp}(Z)|$. \square

Proposition 3 implies that for any flat Z of entropy $O(\log n)$, there is a code that corrects Z with error probability ϵ in polynomial time. Although the construction is not fully explicit, we can obtain such a code with probability $1 - \epsilon$.

Conversely, we can show that the rate achieved in Proposition 3 is almost optimal.

Proposition 4. *Let Z be any flat distribution over $\{0, 1\}^n$ of entropy m . If a code of rate R corrects Z with error probability ϵ , then $R \leq 1 - m/n - \log(1 - \epsilon)/n$.*

Proof. Let (Enc, Dec) be a code that corrects Z with error probability ϵ . For $x \in \{0, 1\}^{Rn}$, define $D_x = \{v \in \{0, 1\}^n : \text{Dec}(v) = x\}$. That is, D_x is the set of inputs that are decoded to x by Dec . Since the code corrects the flat distribution Z with error probability ϵ , $|D_x| \geq (1 - \epsilon)2^m$ for every $x \in \{0, 1\}^{Rn}$. Since each D_x is disjoint, $\sum_{x \in \{0, 1\}^{Rn}} |D_x| \leq 2^n$. Therefore, we have that $(1 - \epsilon)2^m \cdot 2^{Rn} \leq 2^n$, which implies the statement. \square

By Proposition 3, one may hope to construct a *single* code that corrects errors from any flat distribution with the same entropy, as constructed in [5] for the case of binary symmetric channels by using Justesen's construction [19]. However, it is impossible to construct such codes. We show that for every deterministic coding scheme of rate k/n , there is a flat distribution Z of entropy m that is not correctable by the scheme for any $1 \leq m \leq k$. By combining this result with Proposition 1, we can conclude that there is no coding scheme of rate k/n that corrects every flat distribution of entropy m with $1 \leq m \leq k$. For deterministic coding schemes, we assume that the encoding function Enc is injective. Namely, for any distinct $x, x' \in \mathbb{F}^k$, $\text{Enc}(x) \neq \text{Enc}(x')$. The assumption is necessary because otherwise the code always fails to decode either x or x' such that $\text{Enc}(x) = \text{Enc}(x')$.

Proposition 5. *For any deterministic code of rate k/n and any m with $1 \leq m \leq k$, there is a flat distribution of entropy m that is not correctable by the code with error probability $\epsilon < 1/2$.*

Proof. By assumption, the code contains 2^k codewords that are all distinct. Define a flat distribution to be a uniform distribution over 2^m distinct codewords c_1, \dots, c_{2^m} . If the input to the decoder is $c_i + c_j$ for $i, j \in [2^m]$, the decoder cannot distinguish the two cases where the transmitted codewords are c_i and c_j . Thus, the decoder outputs the wrong answer with probability at least $1/2$ for at least one of the two cases. \square

4 Pseudorandom Errors

We show that no efficient coding scheme can correct pseudorandom errors, which can be sampled by pseudorandom generators.

Proposition 6. *Assume that a one-way function exists. Then, for any constant $\epsilon \in (0, 1)$, there is a samplable distribution Z over $\{0, 1\}^n$ of entropy n^ϵ such that no polynomial-time algorithms (Enc, Dec) can correct Z .*

Proof. If a one-way function exists, there is a pseudorandom generator $G : \{0, 1\}^{n^\epsilon} \rightarrow \{0, 1\}^n$ secure for any polynomial-time algorithm [17]. Namely, the distribution $G(U_{n^\epsilon})$ is indistinguishable from the uniform distribution U_n for any polynomial-time algorithm. Then, a distribution $Z = G(U_{n^\epsilon})$ is not correctable by polynomial-time algorithms (Enc, Dec). If so, we can construct a polynomial-time distinguisher for pseudorandom generator by employing (Enc, Dec), and thus a contradiction follows. \square

5 Errors with Membership Test

In this section, we present samplable distributions that are not pseudorandom, but cannot be corrected by efficient coding schemes. For such distributions, we consider distributions for which the membership test can be done efficiently. A distribution Z is called a *distribution with membership test* if there is a polynomial-time algorithm D such that $D(z) = 1 \Leftrightarrow z \in \text{Supp}(Z)$. Since the algorithm D can distinguish Z from the uniform distribution, Z is not pseudorandom.

We consider a sampling algorithm/circuit that can access an *oracle*, which, on querying input x , responds with the value $f(x)$ for an a priori specified function f . It is assumed that the algorithm can obtain the value of any input in a single step, and that other algorithms can also access the same oracle. Such oracle algorithms/circuits are used in the field of computational complexity [2]. It is said to be relative to an oracle if algorithms are allowed to access it.

We show that there is an oracle relative to which there exists a samplable distribution with membership test that is not correctable by efficient coding schemes with low rate. In the proof, we use the technique called *reconstruction paradigm* [13, 34]. Before starting the proof, we briefly describe the technique of [13, 34] and our proof idea.

In [13], the paradigm is used to prove that a random function is one-way with high probability. Roughly speaking, it is shown that if f is not one-way, then f has a “short” description. Here, we say f has a short description if, given access to some oracle A , the truth table of f can be reconstructed by using A and some short information. It is shown in [13] that there are not so many functions with short description, and thus a random function is one-way with high probability. The technique was used in [34] to show the existence of incompressible samplable source with low pseudoentropy.

Here, we use the technique to prove the existence of uncorrectable samplable distributions with membership test. By following the approach of [34], we define a class of functions `correctf`, which contains functions f that can be efficiently corrected by some coding scheme. Then, we show that `correctf` has a short description. Since there are not so many functions with short description, we can show that there is a function f that cannot be corrected with efficient coding schemes.

We now begin the formal proof. Let $N = 2^n, K = 2^k, M = 2^m$. Let \mathcal{F} be the set of injective

functions $f : \{0, 1\}^m \rightarrow \{0, 1\}^n$. For each $f \in \mathcal{F}$, define an oracle \mathcal{O}_f such that

$$\mathcal{O}_f(b, y) = \begin{cases} \mathcal{O}_f^S(y) & \text{if } b = 0, y \in \{0, 1\}^m \\ \mathcal{O}_f^M(y) & \text{if } b = 1, y \in \{0, 1\}^n \\ \perp & \text{otherwise} \end{cases},$$

$$\mathcal{O}_f^M(y) = \begin{cases} 1 & \text{if } y \in f(\{0, 1\}^m) \\ 0 & \text{if } y \notin f(\{0, 1\}^m) \end{cases},$$

$$\mathcal{O}_f^S(y) = f(y).$$

Namely, on input y , the oracle \mathcal{O}_f^M outputs the membership of y under f , and \mathcal{O}_f^S samples $f(y)$. Thus, sampling $f(y)$ can be performed by querying $(0, y)$ to \mathcal{O}_f , and membership test of y can be done by $(1, y)$.

Let correctf be the set of functions $f \in \mathcal{F}$ for which there exist oracle circuits (Enc, Dec) that make q queries to oracle \mathcal{O}_f in total and correct $f(U_m)$ with rate k/n . For each $f \in \text{correctf}$, fix a pair (Enc, Dec) that make q queries to \mathcal{O}_f and correct $f(U_m)$ with rate k/n . We define

$\text{invert}_f = \{y \in \{0, 1\}^m : \text{for any } x \in \{0, 1\}^k, \text{ on input } \text{Enc}(x) + f(y), \text{ Dec queries } \mathcal{O}_f^S \text{ on } y\}$,

$\text{forge}_f = \{y \in \{0, 1\}^m : \text{for some } x \in \{0, 1\}^k, \text{ on input } \text{Enc}(x) + f(y), \text{ Dec does not query } \mathcal{O}_f^S \text{ on } y\}$.

Note that invert_f and forge_f is a partition of $\{0, 1\}^m$. We also define

$$\text{invertible} = \{f \in \text{correctf} : |\text{invert}_f| > \epsilon \cdot 2^m\},$$

$$\text{forgeable} = \{f \in \text{correctf} : |\text{forge}_f| \geq \delta \cdot 2^m\},$$

where ϵ and δ are any positive constants satisfying $\epsilon + \delta = 1$. Note that $\text{correctf} = \text{invertible} \cup \text{forgeable}$. Roughly speaking, if $f \in \text{invertible}$, an ϵ -fraction of $f(y)$ can be corrected by querying \mathcal{O}_f^S on y , which implies that $f(y)$ is invertible by Dec. For $f \in \text{forgeable}$, a δ -fraction of $f(y)$ can be corrected without querying \mathcal{O}_f^S on y , which implies that $f(y)$ is generated by Dec with no help from the oracle.

In the following, we show that every $f \in \text{correctf} = \text{invertible} \cup \text{forgeable}$ has a short description. For each $f \in \text{correctf}$, we present a way for constructing the truth table of f by employing (Enc, Dec) . If $f \in \text{invertible}$, it is done by computing $\text{Enc}(x) + f(y)$ and monitoring oracle queries that Dec($\text{Enc}(x) + f(y)$) makes to \mathcal{O}_f^S . For y on which Dec does not query \mathcal{O}_f^S , the pair $(y, f(y))$ is stored in a look-up table. Similarly, if $f \in \text{forgeable}$, then Dec corrects $f(y)$ without querying \mathcal{O}_f^S on y . This means that $f(y)$ can be described using Dec and $\text{Enc}(x) + f(y)$, and thus if $\text{Enc}(x) + f(y)$ has a short description, the size of forgeable is small.

First, we show that $f \in \text{invertible}$ has a short description.

Lemma 1. *Take any $f \in \text{invertible}$ and a pair of oracle circuits (Enc, Dec) that makes at most q queries to \mathcal{O}_f in total and corrects $f(U_m)$ with rate k/n . Then f can be described using at most*

$$\log \binom{N}{c} + \log \binom{M}{c} + \log \left(\binom{N-c}{M-c} (M-c)! \right)$$

bits, given (Enc, Dec) , where $c = \epsilon M/q$.

Proof. First, consider an oracle circuit A such that, on input z , A picks any $x \in \{0, 1\}^k$ and simulates Dec on input $\text{Enc}(x) + z$. Then, for any $y \in \text{invert}_f$, on input $f(y)$, A outputs y by making at most q queries to \mathcal{O}_f .

Next, we show that for any $f \in \text{invertible}$, f has a short description given A . Without loss of generality, we assume that A makes distinct queries to \mathcal{O}_f^S . We also assume that on input $f(y)$, A always queries \mathcal{O}_f^S on y before it outputs y . We will show that there is a subset $T \subseteq f(\text{invert}_f)$ such that f can be described given T , $B(T)$, $f|_{\{0,1\}^m \setminus B(T)}$, where $f|_X$ denotes the set $\{(x, f(x)) : x \in X\}$ and $B(T) = \{y \in \{0, 1\}^m : y \leftarrow A(z), z \in T\}$.

We describe how to construct T below.

CONSTRUCT- T :

1. Initially, T is empty, and all elements in $T^* = f(\text{invert}_f)$ are candidates for inclusion in T .
2. Choose the lexicographically smallest z from T^* , put z in T , and remove z from T^* .
3. Simulate A on input z , and halt the simulation immediately after A queries \mathcal{O}_f^S on y . Let y'_1, \dots, y'_p be the queries that A makes to \mathcal{O}_f^S , where $y'_p = y$ and $p \leq q$.
 - Remove $f(y'_1), \dots, f(y'_{p-1})$ from T^* . (This means that these elements will never belong to T , and in simulating $A(z)$ in the recovering phase, the answers to these queries are made by using the look-up table for f .)
 - Continue to remove the lexicographically smallest z from T^* until we have removed exactly $q - 1$ elements in Step 3.
4. Return to Step 2.

Next, we describe how to reconstruct f from T , $B(T)$, and $f|_{\{0,1\}^m \setminus B(T)}$. We show how to recover the look-up table for f on values in $B(T)$.

RECOVER- f :

1. Choose the lexicographically smallest element $z \in T$, and remove it from T .
2. Simulate A on input z , and halt the simulation immediately after A queries \mathcal{O}_f^S on y for which the answer does not exist in the look-up table for f . Since the query y satisfies that $y = f^{-1}(z)$, add the entry (y, z) to the look-up table.
3. Return to Step 1.

We explain why we can correctly simulate $A(z)$ in Step 2 of RECOVER- f . Since $B(T)$ and $f|_{\{0,1\}^m \setminus B(T)}$ are given, we can answer all queries to \mathcal{O}_f^M . For any query y' to \mathcal{O}_f^S , it must be either (1) $y' \notin B(T)$, or (2) y' is the output of A on input z' such that $z' \in W$ and z' is lexicographically smaller than z . In either case, the look-up table has the corresponding entry, and thus we can answer the query.

In each iteration in CONSTRUCT- T , we add one element to T and remove exactly q element from T^* . Since initially the size of $T^* = f(\text{invert}_f)$ is ϵM , the size of T in the end is $c = \epsilon M/q$.

Finally, we evaluate the number of bits to describe T , $B(T)$, and $f|_{\{0,1\}^m \setminus B(T)}$. Since T can be specified by choosing c elements from the set $f(\text{invert}_f) \subseteq \{0, 1\}^n$ of size at most N , T can be described using at most $\log \binom{N}{c}$ bits. Similarly, $B(T)$ can be specified by choosing at most $|T|$ elements from the set $\{0, 1\}^m$ of size M . Hence, $B(T)$ can be described using $\log \binom{M}{c}$ bits. The look-up table for $f|_{\{0,1\}^m \setminus B(T)}$ consists of $\{(y, f(y)) : y \in \{0, 1\}^m \setminus B(T)\}$, which can be specified by

first choosing the elements of $\{f(y) : y \in \{0, 1\}^m \setminus B(T)\}$ from $\{0, 1\}^n \setminus f(B(T))$ and then ordering them lexicographically with respect to the input $y \in \{0, 1\}^m \setminus B(T)$. Since f is injective, we have that $|B(T)| = |T| = c$, and thus $|\{f(y) : y \in \{0, 1\}^m \setminus B(T)\}| = M - c$, $|\{0, 1\}^n \setminus f(B(T))| = N - c$, and $|\{0, 1\}^m \setminus B(T)| = M - c$. Thus, the look-up table for $f|_{\{0, 1\}^m \setminus B(T)}$ can be described using $\log\left(\binom{N-c}{M-c}(M-c)!\right)$ bits. Therefore, the statement follows. \square

We show that the fraction of $f \in \mathcal{F}$ for which $f \in \text{invertible}$ and $f(U_m)$ is correctable is small.

Lemma 2. *If $m > 3 \log s + \log n + O(1)$, then the fraction of functions $f \in \mathcal{F}$ such that $f \in \text{invertible}$ and $f(U_m)$ can be corrected by a pair of oracle circuits (Enc, Dec) of total size s is less than $2^{-(sn \log s + 1)}$ for all sufficiently large n .*

Proof. It follows from Lemma 1 that, given (Enc, Dec) , the fraction is

$$\frac{|\text{invertible}|}{\binom{N}{M}M!} \leq \frac{\binom{N}{c}\binom{M}{c}\binom{N-c}{M-c}(M-c)!}{\binom{N}{M}M!} = \frac{\binom{M}{c}}{c!},$$

where $c = \epsilon M / (qK)$. By using the fact that $q \leq s$ and the inequalities $\binom{n}{k} < \left(\frac{en}{k}\right)^k$ and $n! > \left(\frac{n}{e}\right)^n$, the expression is upper bounded by

$$\left(\frac{\epsilon M}{c}\right)^c \left(\frac{e}{c}\right)^c = \left(\frac{e^2 q^2}{\epsilon^2 M}\right)^{\epsilon M/q} < \left(\frac{1}{2}\right)^{ns \log s + 1}$$

for all sufficiently large n . The last inequality follows from the fact that

$$\frac{e^2 q^2}{\epsilon^2 M} < \frac{e^2 q^2}{\epsilon^2 \Omega(s^3 n)} < \frac{1}{2} \quad \text{and} \quad \frac{\epsilon M}{q} > \frac{\epsilon \Omega(s^3 n)}{q} > ns \log s + 1.$$

\square

Next, we show that forgeable has a short description.

Lemma 3. *Take any $f \in \text{forgeable}$ and a pair of oracle circuits (Enc, Dec) that make at most q queries to \mathcal{O}_f in total and corrects $f(U_m)$ with rate k/n . Then f can be described using at most*

$$\log \binom{M}{d} + \log \left(\binom{N-d}{M-d} (M-d)! \right) + d(k + m + \log q)$$

bits, given (Enc, Dec) , where $d = \delta M / q$.

Proof. First, consider an oracle circuit A such that, on input w , A obtains x by simulating Dec on input w , queries \mathcal{O}_f^M on $w - \text{Enc}(x)$, and outputs \perp if $\mathcal{O}_f^M(w - \text{Enc}(x)) = 0$, and x otherwise. Then, A satisfies that, on input w , A outputs \perp if $w \notin \text{Enc}(\{0, 1\}^k) + f(\{0, 1\}^m)$, and $\text{Dec}(w)$ otherwise.

Next, we show that for any $f \in \text{forgeable}$, f has a short description given A . Without loss of generality, we assume that A makes distinct queries to \mathcal{O}_f^S and \mathcal{O}_f^M . We also assume that for $x \in \{0, 1\}^k$ and $y \in \{0, 1\}^m$, $A(\text{Enc}(x) + f(y))$ always queries \mathcal{O}_f^M on $f(y)$ before it outputs x . Note that for $y \in \text{forge}_f$, there is some $x \in \{0, 1\}^k$ such that, on input $\text{Enc}(x) + f(y)$, A does not query \mathcal{O}_f^S on y .

We will show that there is a subset $Y \subseteq \text{forge}_f$ such that f can be described given Y , $f|_{\{0,1\}^m \setminus Y}$, and $\{(x_y, a_y, b_y) \in \{0,1\}^k \times [M] \times [q] : y \in Y\}$ of a set of advice strings. For $x \in \{0,1\}^k$, we define $D(x) = \{\text{Enc}(x) + f(y) : y \in \{0,1\}^m\}$. Note that $|D(x)| = M$ for any $x \in \{0,1\}^k$.

We describe how to construct Y below.

CONSTRUCT- Y :

1. Initially, Y is empty. All elements in $Y^* = \text{forge}_f$ are candidates for inclusion in Y . For every $x \in \{0,1\}^k$, set $D_x = \{\text{Enc}(x) + f(y) : y \in \text{forge}_f\}$. We write $\mathcal{D}_k = \bigcup_{x \in \{0,1\}^k} D_x$.
2. Choose the lexicographically smallest y from Y^* , put y in Y , and remove y from Y^* .
3. Choose the lexicographically smallest w from the set of $\text{Enc}(x) + f(y) \in D_x$ such that A does not query \mathcal{O}_f^S on y . If $w = \text{Enc}(x) + f(y)$, set $x_y = x$. Then, for every $x' \in \{0,1\}^k$, remove $\text{Enc}(x') + f(y)$ from $D_{x'}$. (This removal means that hereafter there are no elements in \mathcal{D}_k for which A outputs some x such that $f(y)$ is the error vector.) When w is the lexicographically t -th smallest element in $D(x)$, set $a_y = t$ (so that we can recognize that the a_y -th element in $D(x)$ is w in the recovering phase).
4. Simulate A on input w , and halt the simulation immediately after A queries \mathcal{O}_f^M on $f(y)$. Let y'_1, \dots, y'_p be the queries that A makes to \mathcal{O}_f^S , and $z'_1, \dots, z'_r = f(y)$ be the queries that A makes to \mathcal{O}_f^M . Set $b_y = r$ (so that we can recognize that the b_y -th query that Dec makes to \mathcal{O}_f^M is $f(y)$ in the recovering phase).
 - (a) For every $x' \in \{0,1\}^k$, remove $\text{Enc}(x') + f(y'_1), \dots, \text{Enc}(x') + f(y'_p)$ from $D_{x'}$.
 - (b) For every $i \in [p]$, if $z'_i \in f(\text{forge}_f)$, then for every $x' \in \{0,1\}^k$, remove $\text{Enc}(x') + z'_i$ from $D_{x'}$, and otherwise, do nothing.
 - (c) Continue to remove the elements $\text{Enc}(x') + f(y)$ from $D_{x'}$ for every $x' \in \{0,1\}^k$ for the lexicographically smallest $w = \text{Enc}(x) + f(y) \in \mathcal{D}_k$ until we have removed exactly $(q-1)K$ elements from \mathcal{D}_k in Step 4.
5. Return to Step 2.

Next, we describe how to construct f from Y , $f|_{\{0,1\}^m \setminus Y}$, and $\{(x_y, a_y, b_y) \in \{0,1\}^k \times [M] \times [q] : y \in Y\}$. We show how to recover the look-up table for f on values in Y .

RECOVER- f :

1. Choose the lexicographically smallest $y \in Y$, and remove it from Y . Then, choose the lexicographically a_y -th smallest element w from $D(x_y)$.
2. Simulate A on input w , and halt the simulation immediately after A makes the b_y -th query to \mathcal{O}_f^M . Since the b_y -th query is $f(y)$, add the entry $(y, f(y))$ to the look-up table.
3. Return to Step 1.

We explain why we can correctly simulate $A(w)$ in Step 2 of RECOVER- f . For any query y' to \mathcal{O}_f^S , it must be either (1) $y' \notin Y$ or (2) y' is lexicographically smaller than y . In case (1), we can answer the query by using $f|_{\{0,1\}^m \setminus Y}$. In case (2), since y was chosen as the lexicographically smallest element such that A does not query \mathcal{O}_f^S on y , the look-up table has the answer to the query. Consider any of the first $b_y - 1$ queries z' to \mathcal{O}_f^M . If $z' \in f(\{0,1\}^m)$, namely $z' = f(y')$ for some y' , then it must be either (1) $y' \notin Y$ or (2) y' is lexicographically smaller than y . In either

case, the look-up table has the entry (y', z') . If $z' \notin f(\{0, 1\}^m)$, there is no entry for z' in the look-up table. Thus, we can answer the query by saying “yes” if z' is in the look-up table, and “no” otherwise.

In each iteration in `CONSTRUCT- Y` , we add one element to Y and remove exactly qK elements from \mathcal{D}_k . Since initially the size of \mathcal{D}_k is at least δKM , the size of Y in the end is at least $d = \delta M/q$.

Finally, we evaluate the number of bits to describe Y , $f|_{\{0,1\}^m \setminus Y}$, and $\{(x_y, a_y, b_y) \in \{0, 1\}^k \times [M] \times [q] : y \in Y\}$. Since Y can be specified by choosing d elements from the set `forge $_f$` $\subseteq \{0, 1\}^m$ of size at most M , Y can be described using $\log \binom{M}{d}$ bits. By the same argument in the proof of Lemma 1 for $f|_{\{0,1\}^m \setminus B(T)}$, we can show that the look-up table for $f|_{\{0,1\}^m \setminus Y}$ can be described using $\log \left(\binom{N-d}{M-d} (M-d)! \right)$ bits. By simply listing the elements, the set $\{(x_y, a_y, b_y) \in \{0, 1\}^k \times [M] \times [q] : y \in Y\}$ can be described using $d(k + m + \log q)$ bits. Therefore, the statement follows. \square

We show that the fraction of $f \in \mathcal{F}$ for which $f \in \text{forgeable}$ and $f(U_m)$ is correctable is small.

Lemma 4. *If $m > 3 \log s + \log n + O(1)$ and $m < n - k - 2 \log s - O(1)$, then the fraction of functions $f \in \mathcal{F}$ such that $f \in \text{forgeable}$ and $f(U_m)$ can be corrected by a pair of oracle circuits (Enc, Dec) of total size s is less than $2^{-(sn \log s + 1)}$ for all sufficiently large n .*

Proof. It follows from Lemma 3 that, given (Enc, Dec) , the fraction is

$$\frac{|\text{forgeable}|}{\binom{N}{M} M!} \leq \frac{\binom{M}{d} \binom{N-d}{M-d} (M-d)!}{\binom{N}{M} M!} 2^{d(k+m+\log q)} = \frac{\binom{M}{d}}{\binom{N}{d} d!} (qKM)^d,$$

where $d = \delta M/q$. By using the fact that $q \leq s$ and the inequalities $\binom{n}{k} < \left(\frac{en}{k}\right)^k$, $\binom{n}{k} > \left(\frac{n}{k}\right)^k$, and $n! > \left(\frac{n}{e}\right)^n$, the expression is upper bounded by

$$\left(\frac{eM}{d}\right)^d \left(\frac{d}{N}\right)^d \left(\frac{e}{d}\right)^d (qKM)^d = \left(\frac{e^2 q^2 KM}{\delta N}\right)^{\delta M/q} < \left(\frac{1}{2}\right)^{ns \log s + 1}$$

for all sufficiently large n . The last inequality follows from the fact that

$$\frac{e^2 q^2 KM}{\delta N} < \frac{e^2 q^2}{\delta \Omega(s^2 n)} < \frac{1}{2} \quad \text{and} \quad \frac{\delta M}{q} > \frac{\delta \Omega(s^3 n)}{q} > ns \log s + 1.$$

\square

We obtain the main result of this section.

Theorem 1. *For any m and k satisfying $3 \log s + \log n + O(1) < m < n - k - 2 \log s - O(1)$, there exist injective functions $f : \{0, 1\}^m \rightarrow \{0, 1\}^n$ such that, given oracle access to \mathcal{O}_f , (1) $f(U_m)$ is a samplable distribution with membership test of entropy m , and (2) $f(U_m)$ cannot be corrected with rate k/n by oracle circuits of size s .*

Proof. Since `correctf` = `invertible` \cup `forgeable`, it follows from Lemmas 2 and 4 that for a fixed (Enc, Dec) of size s , the fraction of functions $f \in \mathcal{F}$ such that (Enc, Dec) corrects $f(U_m)$ with rate k/n is less than $2^{-(sn \log s)}$. Since there are at most $2^{sn \log s}$ circuits of size s , there are functions $f \in \mathcal{F}$ such that $f(U_m)$ cannot be corrected with rate k/n by oracle circuits of size s . Given oracle access to \mathcal{O}_f , $f(U_m)$ is samplable. Since f is injective, $f(U_m)$ has entropy m . \square

The following corollary immediately follows.

Corollary 1. *For any m and k satisfying $\omega(\log n) < m < n - k - \omega(\log n)$, there exists an oracle relative to which there exists a samplable distribution with membership test of entropy m that cannot be corrected with rate k/n by polynomial size circuits.*

6 Positive Results

In this section, we present positive results for correcting samplable additive errors.

6.1 Errors from Linear Subspaces

We show that if the set of error vectors forms a linear subspace, every error can be corrected by a linear code with optimal rate. Let $Z' = \{z_1, z_2, \dots, z_m\} \subseteq \mathbb{F}^n$ be a set of linearly independent vectors. We construct a linear code that corrects additive errors from the linear span of Z' .

Proposition 7. *Let Z be the uniform distribution over the linear span of Z' , which has entropy m . There is a linear code of rate $1 - m/n$ that corrects Z by syndrome decoding.*

Proof. Consider $n - m$ vectors $w_{m+1}, \dots, w_n \in \mathbb{F}^n$ such that the set $\{z_1, z_2, \dots, z_m, w_{m+1}, \dots, w_n\}$ forms a basis of \mathbb{F}^n . Then, there is a linear transformation $T : \mathbb{F}^n \rightarrow \mathbb{F}^m$ such that $T(z_i) = e_i$ and $T(w_i) = 0$, where e_i is the vector with 1 in the i -th position and 0 elsewhere. Let H be the matrix in $\mathbb{F}^{m \times n}$ such that $xH^T = T(x)$, and consider a code with parity check matrix H . Let $z = \sum_{i=1}^m a_i z_i$ be a vector in the linear span of Z' , where $a_i \in \mathbb{F}$. Since $z \cdot H^T = (\sum_{i=1}^m a_i z_i) \cdot H^T = \sum_{i=1}^m a_i e_i = (a_1, \dots, a_m)$, the code can correct the error z by syndrome decoding. Since $H \in \mathbb{F}^{m \times n}$ is the parity check matrix, the rate of the code is $(n - m)/n$. \square

6.2 A Computational Condition

Let $Z = f(U_r)$ be a flat distribution over $\{0, 1\}^n$ of entropy m associated with a samplable additive channel, where f is an efficiently computable function, and $r \geq m$. We give a computational condition under which Z is efficiently correctable. Roughly speaking, we show that if a variant of the function f is efficiently “invertible”, then Z can be efficiently corrected.

Specifically, for function $f : \{0, 1\}^r \rightarrow \{0, 1\}^n$, we define function $g : \{0, 1\}^r \times \mathcal{H} \rightarrow \mathcal{H} \times \{0, 1\}^{m+2c \log n}$ such that

$$g(y, h) = (h, h(f(y))),$$

where $\mathcal{H} = \{h : \{0, 1\}^n \rightarrow \{0, 1\}^{m+2c \log n}\}$ is a family of *linear universal hash functions* [4, 31], and c is a positive constant. The universality means that for any distinct $x, x' \in \{0, 1\}^n$,

$$\Pr_{h \in \mathcal{H}} [h(x) = h(x')] \leq 2^{-(m+2c \log n)},$$

and the linearity means that for any $x, x' \in \{0, 1\}^n$ and $a, b \in \{0, 1\}$, $h(ax + bx') = ah(x) + bh(x')$.

As efficient “invertibility”, we introduce the notion of *distributionally one-way function* [18]. Intuitively, such a function g guarantees that the distribution $(x, g(x))$ for random x is difficult to be simulated by efficient algorithms given only $g(x)$.

Definition 5. A function g is said to be distributionally one-way if it is computable in polynomial time and there exists a constant $c > 0$ such that for every probabilistic polynomial-time algorithm A , the statistical distance between $(x, g(x))$ and $(A(g(x)), g(x))$ is at least $1/n^c$, where $x \sim U_n$.

We show that if the function g defined above is not distributionally one-way, then the error distribution $Z = f(U_r)$ is efficiently correctable. Before giving the formal proof, we describe the underlying idea.

We employ the technique used in the proof of [34, Theorem 6.3] that shows the necessity of one-way functions for separating pseudoentropy and compressibility. In the proof, it is observed that every samplable distribution Z can be optimally compressed by a hash function h such that a sample z from Z is simply compressed to $h(z)$. In addition, if distributionally one-way functions do not exist, it is shown that z can be recovered from $h(z)$ by a polynomial-time algorithm. We observe that a family of *linear* hash functions is used for giving a recovering algorithm. Since a linear compression function is a dual object of a linear code that corrects additive errors [3], we can construct a linear code correcting additive errors of Z . More specifically, we construct an efficient syndrome decoder that recovers $z = f(y)$ from a syndrome $h(z)$ by assuming that f is not distributionally one-way. A problem for recovering z from $h(z)$ is that there may be exponentially many preimages of $h(z)$, and we need to choose z that is in the support of Z . To solve the problem, we define the function $g(y, h) = (h, h(f(y)))$, and construct a code by assuming that g is not distributionally one-way.

Theorem 2. If $g(y, h) = (h, h(f(y)))$ is not distributionally one-way, then the flat distribution $Z = f(U_r)$ over $\{0, 1\}^n$ of entropy m can be corrected by a polynomial-time coding scheme of rate $1 - m/n - (2c \log n)/n$ and error probability $O(n^{-c})$.

Proof. For each $h \in \mathcal{H}$, define $C_h = \{z \in \text{Supp}(Z) : \exists z' \in \text{Supp}(Z) \text{ s.t. } z' \neq z \wedge h(z) = h(z')\}$. Namely, C_h is the set of inputs with collisions under h . By a union bound, it holds that for any $z \in \text{Supp}(Z)$,

$$\Pr_{h \in \mathcal{H}} [\exists z' \in \text{Supp}(Z) : z' \neq z \wedge h(z') = h(z)] \leq \frac{|\text{Supp}(Z)|}{2^{m+2c \log n}} \leq \frac{1}{n^{2c}}.$$

Thus, $\mathbb{E}_{h \in \mathcal{H}} [|C_h|] \leq 2^m/n^{2c}$. We say $h \in \mathcal{H}$ is good if $|C_h| \leq 2^m/n^c$. By Markov's inequality, we have that

$$\Pr_{h \in \mathcal{H}} \left[|C_h| > \frac{2^m}{n^c} \right] < \frac{1}{n^c}.$$

By the assumption that g is not distributionally one-way, there is a polynomial-time algorithm A such that the statistical distance between $((y, h), g(y, h))$ and $(A(g(y, h)), g(y, h))$ is at most n^{-c} , where $y \sim U_r$ and $h \in \mathcal{H}$. Then, it holds that

$$\Pr_{A, y, h} [g(A(g(y, h))) = g(y, h)] \geq 1 - \frac{1}{n^c},$$

where the probability is taken over the random coins of A , $y \sim U_r$, and $h \in \mathcal{H}$. Thus, we have that

$$\Pr_{A, y, h} [g(A(g(y, h))) = g(y, h) \wedge h \text{ is good}] \geq 1 - \frac{2}{n^c}.$$

By fixing the coins of A and $h \in \mathcal{H}$, it holds that there are deterministic algorithm A' and $h_0 \in \mathcal{H}$ such that h_0 is good and

$$\Pr_{y \sim U_r} [g(A'(g(y, h_0))) = g(y, h_0)] \geq 1 - \frac{2}{n^c}.$$

For $y \in \{0, 1\}^r$ satisfying $g(A'(g(y, h_0))) = g(y, h_0)$, we write $A'(g(y, h_0)) = (A'_1(g(y, h_0)), A'_2(g(y, h_0))) = (y', h')$. Then, it holds that $h' = h_0$ and $h_0(f(y)) = h_0(f(y'))$. Furthermore, since h_0 is good, $\Pr_y[f(y) \notin C_{h_0}] \geq 1 - 1/n^c$. Since \mathcal{H} is a set of linear hash functions, there is a matrix $H_0 \in \{0, 1\}^{(m+2c \log n) \times n}$ such that $xH_0^T = h_0(x)$ for $x \in \{0, 1\}^n$. Consider a linear coding scheme in which H_0 is employed as the parity check matrix, and A'_1 is employed for recovering errors from syndromes. That is, $\text{Enc}(x) = xG$ for a full-rank matrix $G \in \{0, 1\}^{(n-m-2c \log n) \times n}$ satisfying $GH_0^T = 0$, and $\text{Dec}(y) = (y - f(A'_1(h_0, yH_0^T)))G^{-1}$, where $G^{-1} \in \{0, 1\}^{n \times (n-m-2c \log n)}$ is a right inverse matrix of G . Then, for any $x \in \{0, 1\}^m$,

$$\begin{aligned} & \Pr_{y \sim U_r} [\text{Dec}(\text{Enc}(x) + f(y)) = x] \\ &= \Pr_{y \sim U_r} [\text{Enc}(x) + f(y) - f(A'_1(h_0, (\text{Enc}(x) + f(y))H_0^T)) = xG] \\ &= \Pr_{y \sim U_r} [f(A'_1(g(y, h_0))) = f(y)], \end{aligned}$$

where we use the property that $GG^{-1} = I$, $\text{Enc}(x) = xG$, $GH_0^T = 0$, and $xH_0^T = h_0(x)$. Since the probability that $g(A_0(g(y, h_0))) = g(y, h_0)$ is at least $1 - 2/n^c$, and for any $y \in \{0, 1\}^r$ satisfying $g(A_0(g(y, h_0))) = g(y, h_0)$, $\Pr_y[f(y) \notin C_{h_0}] \geq 1 - 1/n^c$, we have that

$$\Pr_{y \sim U_r} [f(A'_1(g(y, h_0))) = f(y)] \geq 1 - \frac{3}{n^c}.$$

Thus, Z can be corrected with error probability $O(n^{-c})$. Since f is efficiently computable, both Enc and Dec can be computed in time polynomial in n . Hence the statement follows. \square

It is known that if one-way functions do not exist, then neither do distributionally one-way functions [18]. Thus, Theorem 2 implies the following corollary.

Corollary 2. *If one-way functions do not exist, every samplable flat distribution over $\{0, 1\}^n$ of entropy m can be corrected by an efficient coding scheme of rate $1 - m/n - (2c \log n)/n$ and error probability $O(n^{-c})$ for any constant $c > 0$.*

The above corollary indicates the necessity of one-way functions for proving the impossibility results of Sections 4 and 5.

7 Conclusions

In this work, we study the correctability of samplable additive errors with unbounded error rate. We have considered a relatively simple setting in which the error distribution is identical for every coding scheme and codeword. The results imply that even when a distribution is not pseudorandom by membership test, it is difficult to correct every such samplable distribution by efficient coding schemes. Nevertheless, a positive result can be obtained if we consider much more structured errors such as errors from linear subspaces. We present some possible future work of this study.

Further study on the correctability. In this work, we have mostly discussed impossibility results. Thus, showing non-trivial possibility results is interesting. A possible direction is to consider more structured errors than samplable errors. One can consider *computationally* structured errors such as errors computed by log-space machines, constant-depth circuits, or monotone circuits. Also, one can consider other types of structures, e.g., errors are introduced in a *split-state* manner. Namely, an error vector is split into several parts, and each part is independently computed. This model has been well-studied in the context of leakage-resilient cryptography [11, 22] and non-malleable codes [12, 6, 1]. BSC can be seen as an extreme of this type of channels in which each error bit is computed by the same biased-sampler.

Characterizing correctability. We have investigated the correctability of samplable additive errors using entropy as a criterion. There may be another better criterion for characterizing the correctability of these errors, which might be related to efficient computability, to which samplability is directly related. Since we have considered general distributions as error distributions, the information-spectrum approach [32, 16] may be more plausible.

Acknowledgment

This research was supported in part by Japan Society for the Promotion of Science (JSPS) and the Ministry of Education, Culture, Sports, Science and Technology (MEXT) Grant-in-Aid for Scientific Research Numbers 25106509, 15H00851, 16H01705, 17H01695. We thank anonymous reviewers for their helpful comments.

References

- [1] D. Aggarwal, Y. Dodis, and S. Lovett. Non-malleable codes from additive combinatorics. In D. B. Shmoys, editor, *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 - June 03, 2014*, pages 774–783. ACM, 2014.
- [2] S. Arora and B. Barak. *Computational Complexity - A Modern Approach*. Cambridge University Press, 2009.
- [3] G. Caire, S. Shamai, and S. Verdú. Noiseless data compression with low density parity check codes. In P. Gupta, G. Kramer, and A. J. van Wijngaarden, editors, *DIMACS Series in Discrete Mathematics and Theoretical Computer Science, Piscataway, NJ, March 17 - 19, 2003*, pages 263–284, 2004.
- [4] L. Carter and M. N. Wegman. Universal classes of hash functions. *J. Comput. Syst. Sci.*, 18(2):143–154, 1979.
- [5] M. Cheraghchi. Capacity achieving codes from randomness conductors. In *IEEE International Symposium on Information Theory, ISIT 2009, June 28 - July 3, 2009, Seoul, Korea, Proceedings*, pages 2639–2643. IEEE, 2009.
- [6] M. Cheraghchi and V. Guruswami. Non-malleable coding against bit-wise and split-state tampering. *J. Cryptology*, 30(1):191–241, 2017.

- [7] I. Csiszár and P. Narayan. Arbitrarily varying channels with constrained inputs and states. *IEEE Trans. Information Theory*, 34(1):27–34, 1988.
- [8] I. Csiszár and P. Narayan. Capacity and decoding rules for classes of arbitrarily varying channels. *IEEE Trans. Information Theory*, 35(4):752–769, 1989.
- [9] A. De and T. Watson. Extractors and lower bounds for locally samplable sources. *TOCT*, 4(1):3, 2012.
- [10] Y. Dodis, T. Ristenpart, and S. P. Vadhan. Randomness condensers for efficiently samplable, seed-dependent sources. In R. Cramer, editor, *Theory of Cryptography - 9th Theory of Cryptography Conference, TCC 2012, Taormina, Sicily, Italy, March 19-21, 2012. Proceedings*, volume 7194 of *Lecture Notes in Computer Science*, pages 618–635. Springer, 2012.
- [11] S. Dziembowski and K. Pietrzak. Leakage-resilient cryptography. In *49th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2008, October 25-28, 2008, Philadelphia, PA, USA*, pages 293–302. IEEE Computer Society, 2008.
- [12] S. Dziembowski, K. Pietrzak, and D. Wichs. Non-malleable codes. In A. C. Yao, editor, *Innovations in Computer Science - ICS 2010, Tsinghua University, Beijing, China, January 5-7, 2010. Proceedings*, pages 434–452. Tsinghua University Press, 2010.
- [13] R. Gennaro, Y. Gertner, J. Katz, and L. Trevisan. Bounds on the efficiency of generic cryptographic constructions. *SIAM J. Comput.*, 35(1):217–246, 2005.
- [14] A. V. Goldberg and M. Sipser. Compression and ranking. *SIAM J. Comput.*, 20(3):524–536, 1991.
- [15] V. Guruswami and A. D. Smith. Optimal rate code constructions for computationally simple channels. *J. ACM*, 63(4):35:1–35:37, 2016.
- [16] T. S. Han. *Information-Spectrum Methods in Information Theory*. Springer, 2003. The original Japanese edition was published from Baifukan in 1998.
- [17] J. Håstad, R. Impagliazzo, L. A. Levin, and M. Luby. A pseudorandom generator from any one-way function. *SIAM J. Comput.*, 28(4):1364–1396, 1999.
- [18] R. Impagliazzo and M. Luby. One-way functions are essential for complexity based cryptography (extended abstract). In *30th Annual Symposium on Foundations of Computer Science, Research Triangle Park, North Carolina, USA, 30 October - 1 November 1989*, pages 230–235. IEEE Computer Society, 1989.
- [19] J. Justesen. Class of constructive asymptotically good algebraic codes. *IEEE Transactions on Information Theory*, 18(5):652–656, 1972.
- [20] M. Langberg. Oblivious communication channels and their capacity. *IEEE Transactions on Information Theory*, 54(1):424–429, 2008.
- [21] R. J. Lipton. A new approach to information theory. In P. Enjalbert, E. W. Mayr, and K. W. Wagner, editors, *STACS 94, 11th Annual Symposium on Theoretical Aspects of Computer Science, Caen, France, February 24-26, 1994, Proceedings*, volume 775 of *Lecture Notes in Computer Science*, pages 699–708. Springer, 1994.

- [22] F. Liu and A. Lysyanskaya. Tamper and leakage resilience in the split-state model. In R. Safavi-Naini and R. Canetti, editors, *Advances in Cryptology - CRYPTO 2012 - 32nd Annual Cryptology Conference, Santa Barbara, CA, USA, August 19-23, 2012. Proceedings*, volume 7417 of *Lecture Notes in Computer Science*, pages 517–532. Springer, 2012.
- [23] F. MacWilliams and N. Sloane. *The Theory of Error-Correcting Codes*. North-Holland, 2nd edition, 1978.
- [24] S. Micali, C. Peikert, M. Sudan, and D. A. Wilson. Optimal error correction for computationally bounded noise. *IEEE Transactions on Information Theory*, 56(11):5673–5680, 2010.
- [25] M. Plotkin. Binary codes with specified minimum distance. *IRE Transactions on Information Theory*, 6(4):445–450, 1960.
- [26] R. M. Roth. *Introduction to coding theory*. Cambridge University Press, 2006.
- [27] R. Shaltiel and J. Silbak. Explicit list-decodable codes with optimal rate for computationally bounded channels. In K. Jansen, C. Mathieu, J. D. P. Rolim, and C. Umans, editors, *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2016, September 7-9, 2016, Paris, France*, volume 60 of *LIPICs*, pages 45:1–45:38. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2016.
- [28] C. E. Shannon. A mathematical theory of communication. *Bell Systems Technical Journal*, 27:379–423,623–656, 1948.
- [29] L. Trevisan and S. P. Vadhan. Extracting randomness from samplable distributions. In *41st Annual Symposium on Foundations of Computer Science, FOCS 2000, 12-14 November 2000, Redondo Beach, California, USA*, pages 32–42. IEEE Computer Society, 2000.
- [30] L. Trevisan, S. P. Vadhan, and D. Zuckerman. Compression of samplable sources. *Computational Complexity*, 14(3):186–227, 2005.
- [31] T. Tsurumaru and M. Hayashi. Dual universality of hash functions and its applications to quantum cryptography. *IEEE Trans. Information Theory*, 59(7):4700–4717, 2013.
- [32] S. Verdú and T. S. Han. A general formula for channel capacity. *IEEE Transactions on Information Theory*, 40(4):1147–1157, 1994.
- [33] E. Viola. Extractors for circuit sources. In R. Ostrovsky, editor, *IEEE 52nd Annual Symposium on Foundations of Computer Science, FOCS 2011, Palm Springs, CA, USA, October 22-25, 2011*, pages 220–229. IEEE Computer Society, 2011.
- [34] H. Wee. On pseudoentropy versus compressibility. In *19th Annual IEEE Conference on Computational Complexity (CCC 2004), 21-24 June 2004, Amherst, MA, USA*, pages 29–41. IEEE Computer Society, 2004.
- [35] K. Yasunaga. Error-correcting codes against chosen-codeword attacks. In A. C. A. Nascimento and P. S. L. M. Barreto, editors, *Information Theoretic Security - 9th International Conference, ICITS 2016, Tacoma, WA, USA, August 9-12, 2016, Revised Selected Papers*, volume 10015 of *Lecture Notes in Computer Science*, pages 177–189, 2016.