# Coded trace reconstruction in a constant number of traces

Joshua Brakensiek[*], Ray Li[†], Bruce Spang[‡]

September 15, 2020

## Abstract

The *coded trace reconstruction* problem asks to construct a code $C \subset \{0,1\}^n$ such that any $x \in C$ is recoverable from independent outputs ("traces") of $x$ from a binary deletion channel (BDC). We present binary codes of rate $1 - \varepsilon$ that are efficiently recoverable from $\exp(O_q(\log^{1/3}(\frac{1}{\varepsilon})))$ (a constant independent of $n$) traces of a $\mathrm{BDC}_q$ for any constant deletion probability $q \in (0,1)$. We also show that, for rate $1 - \varepsilon$ binary codes, $\tilde{\Omega}(\log^{5/2}(1/\varepsilon))$ traces are required. The results follow from a pair of black-box reductions that show that average-case trace reconstruction is essentially equivalent to coded trace reconstruction. We also show that there exist codes of rate $1 - \varepsilon$ over an $O_\varepsilon(1)$-sized alphabet that are recoverable from $O(\log(1/\varepsilon))$ traces, and that this is tight.

[*]Department of Computer Science, Stanford University, Stanford, CA. Email: `jbrakens@cs.stanford.edu`. Research supported by an NSF Graduate Research Fellowship.

[†]Department of Computer Science, Stanford University. Research supported by an NSF Graduate Research Fellowship grant DGE-1656518 and NSF grant CCF-1814629. Email: `rayyli@cs.stanford.edu`

[‡]Department of Computer Science, Stanford University, Stanford, CA. Email: `bspang@cs.stanford.edu`.

# 1 Introduction

The *trace reconstruction* problem was first proposed in [Lev01a, Lev01b] and further developed in [BKKM04]. In trace reconstruction, we wish to recover an unknown binary string $x \in \{0,1\}^n$ given a few random subsequences of $x$. Each subsequence, or *trace*, is generated by sending $x$ through the *binary deletion channel with deletion probability* $q$ ($\text{BDC}_q$), which independently deletes each symbol of $x$ with probability $q \in (0,1)$. In particular, the positions of the deleted bits are not known. For example, deleting either the first or second bit of "110" gives the trace "10".

Trace reconstruction has been primarily studied in two settings: *worst-case*, in which the input string $x$ is chosen adversarially, and *average-case*, when the input string $x$ is chosen uniformly at random over all possible $n$-bit strings. The fundamental question in both settings is to determine the minimum number of traces $T = T(n)$ needed in order to recover a length $n$ string $x$ correctly with high probability. In both settings, there is currently an exponential gap (as a function of $n$) for bounding $T(n)$ – see Section 1.1 for the best known bounds.

In this work, we consider an emerging [HM14, CGMR20, AVDiF19] variant of the trace reconstruction known as *coded trace reconstruction*. In this model, we want the smallest $T$ such that there exists a high rate code $C \subset \{0,1\}^n$ such that, for an adversarially chosen $x \in C$, we can recover $x$ with high probability from $T$ traces. This model is directly motivated by DNA storage [YGM17, CGMR20], in which data is stored as multiple encoded strands of DNA. Besides directly generalizing the trace reconstruction problem, coded trace reconstruction also generalizes the well-studied problem of determining the capacity of the binary deletion channel.

In this coded setting, we wish to design codes for trace reconstruction with high *rate*, which is defined[1] to be $\log |C|/n$. We consider the regime in which the rate is $1 - \varepsilon$ (i.e., $|C| \approx 2^{(1-\varepsilon)n}$), where $\varepsilon \in (0,1)$ is a small constant or shrinking as a function of $n$. In particular, the key question we study is as follows.

**Question 1.1.** For a given $\varepsilon \in (0,1)$ and positive integer $n$, what is the smallest $T$ such that we can construct a binary code of rate $1 - \varepsilon$ and length $n$ recoverable from $T$ traces?

**Contributions.** We summarize the main contributions of our work below. See Section 1.2 for formal theorem statements. In all these results, we consider any constant $q \in (0,1)$.

1. **Binary codes with constant number of traces.** For $\varepsilon \in (0,1)$, we construct an infinite family of binary codes of rate $1 - \varepsilon$ efficiently recoverable from a constant number of traces over the $\text{BDC}_q$ (independent of $n$). This follows as an immediate corollary (Corollary 1.5) of the following more general result we prove.

2. **Black-box upper bounds from average-case trace reconstruction.** We show that, if average-case trace reconstruction on length $n$ strings succeeds with sufficiently high probability in $T(n)$ traces, then there exist rate $1 - \varepsilon$ codes that are decodable from $T(\tilde{O}_q(1/\varepsilon))$ traces over the $\text{BDC}_q$ (Theorem 1.4). In particular, by a result in [HPP18], $\exp(O_q(\log^{1/3}(\frac{1}{\varepsilon}))) < \frac{1}{\varepsilon^{o(1)}}$ traces suffice (Corollary 1.5).

3. **Black-box lower bounds from average-case trace reconstruction.** Conversely, we show that if average-case reconstruction on length $n$ strings requires $T(n)$ traces, then reconstruction of any binary code of rate $1 - \varepsilon$ requires $T(\tilde{\Omega}_q(1/\sqrt{\varepsilon}))$ traces over the $\text{BDC}_q$ (Theorem 1.8). In particular, by a recent result [Cha19], $\tilde{\Omega}_q(\log^{5/2}(1/\varepsilon))$ traces are required (Corollary 1.9).

---

[1] All logs and exps are base 2 unless otherwise specified.

4. **Near-equivalence of average-case and coded trace reconstruction.** The two black-box reductions together imply that estimating the optimal number of traces for a code of rate $1 - \varepsilon$ is equivalent to closing the lower and upper bounds within a polynomial for average-case trace reconstruction on strings of length $\text{poly}(1/\varepsilon)$ (Remark 1.11).

5. **Optimal number of traces for constant-sized alphabet.** We also consider the coded trace reconstruction problem over larger alphabets than binary. In particular, we give rate[2] $1 - \varepsilon$ codes over an alphabet of size $O_\varepsilon(1)$ that are efficiently encodable and decodable from $O(\log_{1/q}(1/\varepsilon))$ traces (Theorem 1.12). We show this is optimal up to a constant factor (Theorem 1.13). This shows that coded trace reconstruction is strictly easier for larger alphabets than for binary alphabets. To the best of our knowledge, this is the first non-trivial tight result in *any* model of trace reconstruction for the deletion channel.

## 1.1 Related work

We now discuss how our results are situated at the intersection of the trace reconstruction and coding theory literature.

**Classical trace reconstruction.** One of the main motivations for trace reconstruction is the application to DNA sequencing in computational biology [BKKM04]. When DNA is sequenced, the results may have insertion, deletion, and substitution errors. The original goal of trace reconstruction was to understand a simplified model of how an unknown piece of DNA can be recovered from its sequences. Recently, sequencing has been used for DNA storage [YGM17, CGMR20], in which data is encoded so that it can be stored in DNA. This code needs to be decodable using a trace reconstruction-like process, while being high rate and using as few traces as possible.

The theoretical worst-case setting of trace reconstruction, recovering an arbitrary binary string, was originally studied in [Lev01a, Lev01b, BKKM04, HMPW08]. The current state of the art was derived independently in [DOS17] and [NP17], who show that $\exp(O(n^{1/3}))$ traces suffice for any constant deletion probability $q \in (0, 1)$. A very recent result [Cha20] shows that $\exp(O(n^{1/5}))$ traces suffice for any $q \in (0, 1/2]$. Several works have also considered lower bounds for worst-case trace reconstruction [BKKM04, HMPW08, MPV14a, HL+20, Cha19]. The best known lower bound is $\Omega\left(\frac{n^{3/2}}{\log^{16} n}\right)$ traces [Cha19], which has an exponential gap compared to the best known upper bound. Our work does not use or address worst-case trace reconstruction.

In the average-case setting studied by [HMPW08, MPV14a, PZ17, HPP18], the best upper bound is given by [HPP18], who showed that, for all deletion probabilities $q \in (0, 1)$, a subpolynomial $\exp(O(\log^{1/3} n))$ traces suffice to recover a random string with high probability. Several works have also considered lower bounds for average-case trace reconstruction [MPV14a, HL+20, Cha19]. The current best bound of $\Omega\left(\frac{\log n^{5/2}}{(\log\log n)^{16}}\right)$ traces [Cha19] again has an exponential gap. Our work shows that resolving the optimal number of traces up to a constant factor for coded trace reconstruction is essentially equivalent to average-case reconstruction.

Trace reconstruction over a larger alphabet is less well studied. [MPV14b, DOS17] show that it is possible to turn any trace reconstruction algorithm over a non-binary alphabet into a trace over a binary alphabet and use binary trace reconstruction to solve the problem, at a small cost to the failure probability. For coded trace reconstruction, we show that there is a substantial benefit to using a non-binary alphabet. For constant-sized alphabets, we show a matching upper and lower bound, determining the optimal number of traces up to a constant factor.

---

[2]The rate of a code $|C|$ of length $n$ over an alphabet $\Sigma$ is $\frac{\log_{|\Sigma|} |C|}{n}$

**Coded trace reconstruction.** Coded trace reconstruction generalizes the classical questions above about trace reconstruction. The worst-case trace reconstruction question over a binary alphabet asks how many traces $T(n)$ are needed to achieve error probability $o(1)$ for the code $C = \{0,1\}^n$. As we show in Section 2.2, average-case trace reconstruction is equivalent to asking how many traces $T(n)$ are needed to achieve error probability $o(1)$ for a code $C$ of size $2^n(1-o(1))$. We use this connection to average-case trace reconstruction to construct much longer codes which are recoverable from few traces.

Cheraghchi, Gabrys, Milenkovic, and Ribeiro [CGMR20] formulated the coded trace reconstruction problem considered here. Among other constructions, they give explicit constructions of binary codes of rate $1 - O(\frac{1}{\log\log n})$ recoverable in $\exp(O(\log\log n)^{2/3})$ traces, and rate $1 - O(\frac{1}{\log n})$ code recoverable in $\operatorname{poly}\log n$ traces. Our work improves the number of traces and allows a wider range of rates. For any $\varepsilon \geq n^{-o(1)}$, we show that there exist binary codes of rate $1 - \varepsilon$ recoverable in $\exp(O_q(\log^{1/3}(\frac{1}{\varepsilon})))$ traces. Taking $\varepsilon = \Theta(\frac{1}{\log\log n})$ and $\varepsilon = \Theta(\frac{1}{\log n})$ gives the respective improvements to [CGMR20] in the number of traces. We emphasize that all the constructions of [CGMR20] have polynomial time encoding and decoding, whereas our constructions have polynomial time decoding in all considered parameter settings, but only polynomial time encoding when $\varepsilon \geq \Omega(\frac{\log\log n}{\log n})$.

Although our work deals with a constant fraction of deletions, several prior works considered coding for trace reconstruction for small numbers of deletions. Haeupler and Mitzenmacher [HM14] showed that, for any fixed integer $T$, as the deletion probability $q$ approaches 0, there exists a binary code recoverable from $T$ traces across the $\mathrm{BDC}_q$ with rate $1 - O(H(q^T))$, where $H$ is the binary entropy function. By contrast, our codes handle deletion probabilities arbitrarily close to 1. We show, for example, that there exist binary codes of rate 0.99 recoverable from $T = O(1)$ traces of the $\mathrm{BDC}_{0.99}$. Abroshan, Venkataramanan, Dolecek, and Guillén [AVDiF19] consider coding for channels applying a constant number of deletions. They concatenate $\ell$ Varshamov-Tenengolts [VT65] codes of length $m$ to construct a code of length $m\ell$ and rate $1 - O(\frac{\log m}{m})$ for any $m, \ell \geq 1$. They bound the error probability for recovering for a channel that applies exactly $\ell'$ deletions, when $\ell' < \ell$.

**Other trace reconstruction variants.** There has recently been a variety of work on other problems related to trace reconstruction, which our work does not address. [GM19] considers the problem of recovering a string from the multiset of all its length $L$ substrings. [BCF+19] studies population recovery under the deletion channel, an extension to trace reconstruction where we recover an unknown distribution over input strings, rather than a single input string. In [KMMP19], the authors consider the problems of reconstructing matrices and sparse strings from traces.

**Codes for the deletion channel.** The optimal rate for coded trace reconstruction with one trace is also known as the *capacity* of the binary deletion channel, a well-studied and difficult problem. The capacity of the binary deletion channel with deletion probability $q$ is clearly at most $1 - q$, the capacity of the simpler binary erasure channel. When $q \to 0$, the capacity is known to approach $1 - H(q)$, where $H(q)$ is the binary entropy function (see [DG01] for the lower bound and [KM13, KMS10] for the upper bound). When $q \to 1$, the capacity is known to be $\Theta(1-q)$, but the exact capacity is known only to be roughly between $0.11(1-q)$ [DM06, DM07], and $0.41(1-q)$ [RD15]. A polynomial time encodable/decodable code meeting this up to a constant factor was given in [GL19, CS20]. The current best capacity upper bounds for intermediate $q$ (e.g., $q = 0.5$) are given by [FD10, RD15, Che18]. We incorporate techniques used in constructing codes for the binary deletion channel in our construction of Theorem 1.4. Our work shows that, at $q = 1 - \delta$,

if one is allowed to reconstruct from $O_\delta(1)$ traces of the $\text{BDC}_q$ rather than only one trace, the capacity of the resulting channel improves from $\Theta(\delta)$ to $0.99$.

## 1.2 Main results

We now define the coded trace reconstruction problem formally and state our main theorems. For $q \in (0,1)$ and $x \in \{0,1\}^n$, we let $\text{BDC}_q(x)$ denote the probability distribution of output of $x$ across the $\text{BDC}_q$. We let $\{0,1\}^*$ denote the set of binary strings of any length.

**Definition 1.2.** For $q, \delta \in (0,1)$ and positive integers $n$ and $T$, we say a code $C \subset \{0,1\}^n$ is $(T, q, \delta)$ *trace reconstructible* if there exists a *decoding function* $\text{Dec} : (\{0,1\}^*)^T \to C$ such that, for all $c \in C$,

$$\Pr_{z_1,\ldots,z_T \sim \text{BDC}_q(c)}[\text{Dec}(z_1,\ldots,z_T) \neq c] < \delta.$$

Typically, we desire $\delta \to 0$ as $n \to \infty$. We say $C$ is *decodable* in time $t$ if Dec can be computed in time $t$. We say $C$ is *encodable* in time $t$ if there exists a bijection $\text{Enc} : \{1,\ldots,|C|\} \to C$ that can be evaluated in time $t$. The following notation, denoting the optimal number of traces for average-case trace reconstruction, is used throughout the paper.

**Definition 1.3.** For $m \geq 1, q \in (0,1)$, and $\beta \geq 0$, let $T_{q,\beta}^{(\text{avg})}(m)$ denote the smallest integer $T$ such that there exists a trace reconstruction algorithm for the $\text{BDC}_q$ using $T$ traces that, on a uniformly random string $x$ of length $m$, succeeds with probability (over the randomness of the string and channel) at least $1 - \frac{1}{3m^\beta}$. When $\beta$ is omitted, we take $\beta = 0$.

By repetition of the reconstruction algorithm and subsequently taking a majority vote, we have $T_q^{(\text{avg})}(m) \leq T_{q,\beta}^{(\text{avg})}(m) \leq O(\beta \log m) \cdot T_q^{(\text{avg})}(m)$, so $T_{q,\beta}^{(\text{avg})}(m)$ and $T_q^{(\text{avg})}(m)$ are roughly the same size for constant $\beta$.

**Binary upper bound.** We prove the following upper bound for coded trace reconstruction, which allows bounds for average-case trace reconstruction to be turned into bounds for coded trace reconstruction.

**Theorem 1.4.** *For all $q, \varepsilon \in (0,1)$, there exists constants $n_0 = 1/\varepsilon^{O_q(1)}$, $\beta = \Theta_q(1), n_R = \Theta_q(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon})$, and $\delta = 2^{-\varepsilon^{O_q(1)} n}$ such that, for all $n \geq n_0$, there exists a code $C \subset \{0,1\}^n$ of rate $1 - \varepsilon$ that is $(T_{q,\beta}^{(\text{avg})}(n_R), q, \delta)$ trace reconstructible. Furthermore, the encoding can be done in time $\text{poly}_{\varepsilon,q}(n)$ and trace reconstruction can be done in time $\text{poly}(n)$.*

We can instantiate Theorem 1.4 using the state-of-the-art construction for average-case trace reconstruction of Holden, Pemantle, and Peres [HPP18], which states that $T_q^{(\text{avg})}(\frac{1}{\varepsilon}) \leq \exp(O_q(\log^{1/3} \frac{1}{\varepsilon}))$. Doing so gives the following.

**Corollary 1.5.** *For all $q, \varepsilon \in (0,1)$, there exists constants $n_0 = 1/\varepsilon^{O_q(1)}, T = \exp(O_q(\log^{1/3}(\frac{1}{\varepsilon})))$, and $\delta = 2^{-\varepsilon^{O_q(1)} n}$ such that, for all $n \geq n_0$, there exist codes of length $n$ and rate at least $1 - \varepsilon$ that are $(T, q, \delta)$ trace reconstructible.*

**Remark 1.6.** In coding theory, we are sometimes interested in codes with rate quickly approaching 1, and our bounds on the number of traces hold in this setting as well. For every $q \in (0,1)$, Theorem 1.4 and Corollary 1.5 holds for all integers $n \geq \frac{1}{\varepsilon^{\Omega_q(1)}}$. Thus, we obtain obtain similar results for $\varepsilon$ going to 0 with $n$ so long as $\varepsilon \geq \frac{1}{n^{O_q(1)}}$. Setting $\varepsilon = O(\frac{1}{\log n})$, we have codes of rate $1 - O(\frac{1}{\log n})$

recoverable from $\exp(O_q(\log\log n)^{1/3})$ traces with failure probability $2^{-\tilde{O}_q(n)}$, improving upon the poly $\log n$ number of traces in [CGMR20] needed for the same $\varepsilon$. Our construction also gives a better bound on the number of traces when $\varepsilon = O(\frac{1}{\log\log n})$, improving from $\exp(O_q(\log\log n)^{2/3})$ traces to $\exp(O_q(\log\log\log n)^{1/3})$ traces.

**Remark 1.7.** While we improve on the number of traces in [CGMR20] and also give polynomial time decoding like in [CGMR20], their codes are all polynomial time encodable, whereas ours are only so when $\varepsilon \geq \Omega(\frac{\log\log n}{\log n})$: a careful look at our runtimes shows our code is encodable in time $t_{enc}(\Theta(\frac{1}{\varepsilon}\log\frac{1}{\varepsilon})) \cdot \text{poly}\, n$, where $t_{enc}(n')$ is the amount of time needed to encode a string of length $n'$ used for average-case trace reconstruction, as in Lemma 2.7. Naively we upper bound $t_{enc}(n') \leq 2^{O(n')}$. Thus, when $\varepsilon = O(\frac{1}{\log n})$, while we improve on the number of traces from [CGMR20] and also give polynomial time decoding, only [CGMR20] has codes with both encoding and reconstruction in polynomial time. Furthermore, the constants in our code are quite large, making them currently impractical. Still, we hope the ideas in our construction could be used for future efficient constructions.

**Binary lower bound.** We also prove the following converse, showing that the number of traces needed for rate $1 - \varepsilon$ trace reconstruction is at least the number of traces needed for average-case trace reconstruction on length $\frac{1}{\varepsilon^{1/2-o(1)}}$ strings with failure probability $1/3$.

**Theorem 1.8.** *For all $q, \delta \in (0, 1)$, for sufficiently small $\varepsilon > 0$, there exists $m = \tilde{\Omega}_q(\frac{1}{\varepsilon^{1/2}})$ such that, if $T = T_q^{(\text{avg})}(m)$, all rate $1 - \varepsilon$ codes of sufficiently large length are not $(T - 1, q, \delta)$-trace reconstructible.*

Using Theorem 1.8, we can adapt the state-of-the-art lower bound for average case trace reconstruction into a lower bound for coded trace reconstruction. Recently Chase [Cha19], building off work of Holden and Lyons [HL$^+$20], showed that $T_q^{(\text{avg})}(m) \geq \tilde{\Omega}_q((\log m)^{5/2})$.[3] Applying Theorem 1.8 to this result gives us the following lower bound.

**Corollary 1.9.** *For all $q, \delta \in (0, 1)$ and $\varepsilon > 0$ sufficiently small, there exists $T = \tilde{\Omega}_q((\log\frac{1}{\varepsilon})^{5/2})$ such that all rate $1 - \varepsilon$ codes of sufficiently large length are not $(T, q, \delta)$-trace reconstructible.*

**Remark 1.10.** Theorem 1.8 holds when $n \geq \tilde{\Omega}_q(\frac{1}{\varepsilon^2})$. Hence, similar to Remark 1.6, the lower bound of Theorem 1.8 holds for $\varepsilon$ approaching 0 with $n$, so long as $\varepsilon \geq \Omega_q(\frac{1}{n^{1/2}})$.

**Remark 1.11.** Theorem 1.4 and Theorem 1.8 together show that the optimal number of traces for a code of rate $1 - \varepsilon$ is bounded above and below by the number of traces for average-case trace reconstruction of a string of length $\text{poly}(1/\varepsilon)$. More precisely, there exist $m_1 = \tilde{\Omega}_q(\frac{1}{\sqrt{\varepsilon}})$ and $m_2 = \tilde{O}_q(\frac{1}{\varepsilon})$ such that the optimal number of traces for rate $1 - \varepsilon$ coded trace reconstruction with failure probability $\frac{1}{3}$ is between $T_q^{(\text{avg})}(m_1)$ and $O_q(\log\frac{1}{\varepsilon}) \cdot T_q^{(\text{avg})}(m_2)$. Hence any qualitative improvement to the upper or lower bounds for coded trace reconstruction implies an analogous improvement for average-case trace reconstruction and vice versa.

**Large alphabet upper and lower bounds.** So far, we have focused on codes for binary alphabets. By defining the deletion channel for strings over larger alphabets in the same way as the binary deletion channel, one can ask questions for coded trace reconstruction over larger alphabets.

---

[3]Here, $\tilde{\Omega}(\cdot)$ suppresses log log factors. In fact, they show something stronger: even achieving success probability $\exp(m^{-0.15})$ requires that many traces.

In this setting, our results are stronger in two ways. Firstly, we are able to show matching upper and lower bounds for large alphabet trace reconstruction. Secondly, these constructions are simpler and do not rely on average-case trace reconstruction results.

**Theorem 1.12.** *For all $q, \varepsilon \in (0,1)$ and infinitely many $n$, there exists a rate $1 - \varepsilon$ code over an alphabet of size $2^{O(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon})}$ that is $(T, q, \delta)$ trace reconstructible for $T = O(\log_{1/q} \frac{1}{\varepsilon})$ and $\delta = 2^{-\Omega(n)}$ and which is encodable in time $O(n)$ and decodable in time $O(nT)$.*

And as the following lower bound shows, this is tight in terms of the number of traces.

**Theorem 1.13.** *Any code (over any alphabet) of rate $1 - \varepsilon$ is not $(\lfloor \log_{1/q} \frac{1}{\varepsilon} \rfloor, q, o(1))$ trace reconstructible.*

We do not know if the dependence on $\varepsilon$ for the alphabet size in Theorem 1.12 is optimal. We leave understanding the trade-off between alphabet size and number of traces as an open question for future work.

## 1.3 Techniques

In this section we describe our constructions. We first combine synchronization strings [HS17] and erasure codes [GI05] to give our large alphabet construction (Theorem 1.12), and match this construction with a simple lower bound (Theorem 1.13).

Extending these ideas to our binary code construction (Theorem 1.4) requires more work, and we introduce a novel technique for binary code concatenation, turning our large alphabet code from Theorem 1.12 into a binary code. This concatenation also leverages codes for the binary deletion channel (e.g. [GL19]), and bounds for average-case trace reconstruction [HPP18].

We finish this section by describing our lower bound for coded trace reconstruction for the binary alphabet (Theorem 1.8). Trace reconstruction lower bounds usually find a hard pair of strings and prove that it takes many traces to distinguish these strings. Coded trace reconstruction can simply avoid these hard pairs of strings, which makes applying prior results difficult. Using techniques from information theory, we are able to transfer average-case trace reconstruction lower bounds to the coded setting.

**Large alphabet construction and lower bound.** As a warm-up, first observe that any binary code $C \subset \{0,1\}^n$ can be turned into a code $C'$ over an alphabet of size $2n$ by mapping each codeword $(r_1, \ldots, r_n)$ to a codeword $((r_1, 1), (r_2, 2), \ldots, (r_n, n)) \in (\{0,1\} \times [n])^n$. This code has very low rate, but has the useful property that the deletion channel is essentially turned into an erasure channel: from a received string, we can always recover the indices of the received symbols, and thus the corresponding $r_i$. If $C$ is a code of rate $1 - \varepsilon$ tolerating a $\delta = \text{poly}(\varepsilon)$ fraction of erasures, $C'$ is recoverable from $O(\log_{1/q} \frac{1}{\varepsilon})$ traces: with high probability at most $q^T < \delta$ fraction of symbols are never received, producing less than $\delta n$ erasures, which can be corrected.

Our construction for large alphabets (Theorem 1.12) uses the above intuition, but relies on synchronization strings to avoid ruining the rate of the resulting code. Instead of specifying the exact position of each symbol, we include a symbol of a synchronization string [HS17] from a much smaller alphabet of size $\text{poly}\left(\frac{1}{\varepsilon}\right)$. We take our starting code $C$ to be over a large alphabet of size $2^{O(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon})}$ and tolerate a $\delta = \text{poly}(\varepsilon)$ fraction of erasures [GI05]. Increasing the size of the alphabet beyond that of [GI05] helps ensure the correct rate when combining with the synchronization string. At the cost of a few more erasures, we can convert the outputs on the deletion channel into outputs with erasures and correct the erasures.

For the lower bound (Theorem 1.13), any code of rate $1 - \varepsilon$ recovering from $T$ traces must also be able to recover from the erasure channel with erasure probability $q^T$, which has capacity at most $1 - q^T$. Therefore, $1 - \varepsilon < 1 - q^T$ so $\log_{1/q} \frac{1}{\varepsilon}$ traces are necessary for the erasure channel, and thus the deletion channel.

**Binary alphabet construction.** Our construction for binary alphabets (Theorem 1.4) uses additional ideas beyond those in the large alphabet construction. Again, we use a high rate error correcting code with codewords $(r_1, \ldots, r_{n_{out}}) \in C$ and a synchronization string $(s_1, \ldots, s_{n_{out}})$. Naively, one might "concatenate" the large alphabet construction with a high rate code of length $n_{in} = O(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon})$ recoverable from a $O_\varepsilon(1)$ number of traces (which exists by [HPP18]), so that each pair $(r_i, s_i)$ is encoded in a binary string $a_i$ of length $n_{in}$, and the final codeword is the concatenation $a_1 || \cdots || a_{n_{out}}$. Then, to recover the message, we first use the $T$ traces of the codeword $a_1 || \cdots || a_{n_{out}}$ to recover $T$ traces of each $a_i$. As in [CGMR20], we can make sure we know where the traces of the $a_i$ start and finish by adding buffers of long runs on the ends of each $a_i$. From the traces of each $a_i$, we run the inner trace reconstruction to recover each $a_i$, and thus recover the pair $(r_i, s_i)$. We then run the outer error correction to fix any incorrectly decoded $r_i$'s.

This construction does not work for a subtle reason. Because the length of each $a_i$ is a constant, we expect a (very small) constant fraction of the $a_i$'s buffers to be deleted, and we also expect a (very small) constant fraction of $a_i$'s to have deletions applied so that the interior of the $a_i$ looks like a buffer (we call this a "spurious" buffer). From the $T$ traces of the codeword, we try to recover $T$ traces of each of the $a_i$'s using the buffers, but these $T$ traces, supposedly of $a_i$, might contain some traces of, e.g., $a_{i-5}$ or $a_{i+3}$. Therefore, we need to know the synchronization symbols $s_i$ to determine which substrings of each of the $T$ traces belong to which $a_i$. Thus, recovering the synchronization symbols must happen *before* running trace reconstruction on the $a_i$'s. However, the synchronization symbols $s_i$ are encoded in the $a_i$, so in this construction the synchronization symbols cannot be recovered until *after* the trace reconstruction.

To avoid this issue, our construction crucially encodes the content symbol $r_i$ and the synchronization symbol $s_i$ separately. To our knowledge, this kind of concatenation has not appeared in other constructions of deletion codes. Each content symbol $r_i$ is encoded using a high rate code of length $n_R = \Theta(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon})$ obtained from bounds on average-case trace reconstruction. Each synchronization symbol is encoded in a code of length $n_S = \Theta(\log \frac{1}{\varepsilon})$ decodable in, crucially, 1 trace from the binary deletion channel. We can afford a very low rate code for the synchronization symbols because they are over a much smaller alphabet than the content symbols. Furthermore, we structure the encoded content symbols and encoded synchronization symbols so that they are not easily confused with each other.

For the final decoding algorithm, we first recover the synchronization symbols within each trace. We then use the synchronization strings to determine the parts of each trace that corresponding to traces of a particular $a_i$. We then use these traces of $a_i$ in trace reconstruction to recover each content symbol $r_i$. Finally, we use the error correction of the outer code $C$ to fix any mistakes in this process.

**Binary alphabet lower bound.** Our binary lower bound (Theorem 1.8) reduces coded trace reconstruction to constructing a code over an appropriately chosen memoryless channel, i.e. a channel where each alphabet symbol is corrupted independently or the other symbols. In particular, we partition the input string $x \in \{0,1\}^n$ into $n/m$ substrings of length $m \approx 1/\sqrt{\varepsilon}$. We then upper bound the rate of a code $C \subset (\{0,1\}^m)^{n/m}$ over alphabet $\{0,1\}^m$ recovering a sequence $x$ of length $m$ substrings from $T = T_q^{(avg)}(m)$ independent traces of each of the $n/m$ substrings. This is easier

than recovering $x$ from $T$ independent traces of itself, so any rate upper bound for the code for $n/m$ substrings yields a rate upper bound for the original coded trace reconstruction problem.

Now, we can view the problem as coding over a discrete memoryless channel: we view our binary code as a code of length $n/m$ over the input alphabet $\mathcal{X} = \{0,1\}^m$ and the channel produces outputs in $\mathcal{Y} = (\{0,1\}^*)^T$, corresponding to $T$ independent traces of the elements of $\mathcal{X}$. By Shannon's noisy channel coding theorem [Sha48], the capacity of this channel equals the maximum, over distributions $\lambda$ on $\mathcal{X}$, of the mutual information $I(X_\lambda, Y_\lambda)$, where $X_\lambda \in \mathcal{X}$ is sampled from $\lambda$ and $Y_\lambda \in \mathcal{Y}$ is a tuple of $T$ strings each sampled as an independent trace of $X_\lambda$. Thus, to upper bound the rate of $C$, it suffices to upper bound the mutual information $I(X_\lambda, Y_\lambda)$ for all distributions $\lambda$ on $\mathcal{X}$. If the distribution $\lambda$ is "far" from the uniform distribution, we can upper bound the mutual information by the entropy of $X_\lambda \sim \lambda$. Otherwise, if $\lambda$ is "close" to the uniform distribution, the mutual information is limited by the performance of average-case trace reconstruction. In either case, we get an upper bound on the mutual information which implies an upper bound on the rate of a code correctable from $T$ traces.

## 1.4 Paper organization

In Section 2, we define a few building blocks for our work. These include synchronization strings, codes for the binary deletion channel, and high rate error correcting codes. In Section 3, we present the proofs of our coded trace reconstruction results over large alphabets in Theorems 1.12 and Theorem 1.13. These proofs are simpler and serve as warm-ups for our results over binary alphabets, which require additional ideas. In Section 4.1, we sketch the proof of Theorem 1.4, showing how to convert upper bounds for average-case trace reconstruction into upper bounds for coded trace reconstruction. In the remainder of Section 4, we formally prove Theorem 1.4. In Section 5, we prove Theorem 1.8, giving a black-block reduction from lower bounds for average-case trace reconstruction to lower bounds for coded trace reconstruction. Appendix A fills in various technical details omitted from the main body.

# 2 Preliminaries

## 2.1 Basics

All logs and exps are base 2 unless otherwise specified. For an alphabet $\Sigma$, we let $\Sigma^*$ denote the set of strings over $\Sigma$ of any length. For strings $w, w'$, we let $ww'$ denote the concatenation of strings $w$ and $w'$. We may also denote the concatenation by $w\|w'$ for clarity. For a string $w$ and integer $i$, let $w^i$ denote the string $ww\cdots w$ with $w$ repeated $i$ times. A *substring* is a sequence of consecutive characters in a string. A *run* is a maximal substring of a string all of whose bits are the same. A *partial function* $f : A \nrightarrow B$ is a function from a subset of $A$ to $B$. For $x \in (0,1)$, let $H(x) = -x \log x - (1-x) \log(1-x)$ denote the binary entropy function.

A *code* $C$ of length $n$ over an alphabet $\Sigma$ is a subset of $\Sigma^n$. The elements of $C$ are called *codewords*, and $n$ is called the *length* of the code. If $|\Sigma| = 2$, we say $C$ is a binary code. The *rate* of a code $C$ is defined to be $\frac{\log |C|}{n \log |\Sigma|}$. A code may have an associated *message set* $\mathcal{M}$ and *encoding function* Enc $: \mathcal{M} \to C$, which is an injective map from messages to codewords. By default, $\mathcal{M} = \{1, \ldots, |C|\}$. A code is *decodable under the* $\mathrm{BDC}_q$ *with failure probability* $\delta$ if it is $(1, q, \delta)$ trace reconstructible. To *construct* a code means to produce a description of its encoding and decoding functions. Given two codes $C_1 \subset \Sigma_1^{n_1}$ and $C_2 \subset \Sigma_2^{n_2}$ with $|\Sigma_1| \leq |C_2|$, a *concatenation* of $C_1$ and $C_2$ is a code $C \subset \Sigma_2^{n_1 n_2}$ whose codewords are $\mathrm{Enc}_2(c_1)\|\ldots\|\mathrm{Enc}_2(c_{n_1})$ where $c_1 \cdots c_{n_1} \in C_1$, and where $\mathrm{Enc}_2 : \Sigma_1 \to C_2$ is a fixed injective map.

We use the following forms of the Chernoff bound (e.g., [DP09])

**Lemma 2.1** (Chernoff bound – discrete). *Let $X_1, \ldots, X_n$ be independent and identically distributed random variables with mean $\mu$ supported on $\{0, 1\}$ Then, for $\delta \geq 0$,*

$$\mathbf{Pr}[X_1 + \cdots + X_n \leq (1 - \delta) \cdot n\mu] \leq e^{-\frac{\delta^2}{2} \cdot n\mu} \tag{1}$$

$$\mathbf{Pr}[X_1 + \cdots + X_n \geq (1 + \delta) \cdot n\mu] \leq e^{-\frac{\delta^2}{2+\delta} \cdot n\mu}. \tag{2}$$

**Lemma 2.2** (Chernoff bound – continuous). *Let $X_1, \ldots, X_n$ be independent and identically distributed random variables with mean $\mu$ supported on $[0, 1]$ Then, for $\delta \geq 0$,*

$$\mathbf{Pr}[X_1 + \cdots + X_n \geq (1 + \delta) \cdot n\mu] \leq e^{-2\delta^2 \cdot \mu^2 n}. \tag{3}$$

## 2.2 Short codes from average-case trace reconstruction

In this section, we show a connection between short codes for trace reconstruction and average-case trace reconstruction. We use this connection to construct short, high-rate, trace reconstructible codes, which are building blocks in our main result.

The current state of the art for the optimal number of traces for average-case trace reconstruction is due to Holden, Pemantle, and Peres [HPP18], who show the following bound on $T_{q,\beta}^{(\mathrm{avg})}(n)$.

**Theorem 2.3** ([HPP18]). *For all $q \in (0, 1)$ and $\beta \geq 1$, we have $T_{q,\beta}^{(\mathrm{avg})}(n) \leq \exp(O_{q,\beta}(\log^{1/3} n))$.*

Note that the paper [HPP18] only states Theorem 2.3 for failure probability $1/n$, but their proof works in the same way for any polynomial failure probability $1/n^\beta$. There is also a slick way to amplify the failure probability in average-case trace reconstruction: with polynomially more traces, we can turn failure probability $1/n$ into $1/n^\beta$, by appending random bits to each trace and running trace reconstruction for $n' = n^\beta$ (see e.g., Theorem 3.2 of [BCSS19]).

We now have the following two simple observations that results for average-case trace reconstruction show the existence of codes for coded trace reconstruction and vice versa.

**Claim 2.4.** *If there exists a code of size $2^n(1 - o(1))$ that is $(T, q, o(1))$ trace reconstructible, then average case trace reconstruction can be done in $T$ traces with failure probability $o(1)$.*

*Proof.* The probability that a random string is both in the code and is decoded correctly from $T$ traces is at least $1 - o(1)$. $\square$

And conversely,

**Lemma 2.5.** *Let $\beta > 1$, $q \in (0, 1)$, and $T = T_{q,2\beta}^{(\mathrm{avg})}(n)$. For all positive integers $n$, there exists a code $C$ with $|C| \geq (1 - n^{-\beta})2^n$ that is $(T, q, n^{-\beta})$ trace reconstructible.*

*Proof.* For any string $x \in \{0, 1\}^n$, let $\delta_x$ denote the probability that $x$ is recovered incorrectly using the algorithm solving trace reconstruction for random traces on the $\mathrm{BDC}_q$ in $T$ traces with failure probability $n^{-2\beta}$. By definition of $T = T_{q,2\beta}^{(\mathrm{avg})}(n)$, we have $\mathbf{E}_x[\delta_x] \leq \frac{1}{3}n^{-2\beta}$, so, by Markov's inequality, $\mathbf{Pr}_x[\delta_x \geq n^{-\beta}] < n^{-\beta}$. Setting $C$ to be the set of all $x$ with $\delta_x \leq n^{-\beta}$ and using the same trace reconstruction algorithm gives that $C$ is $(T, q, n^{-\beta})$-trace reconstructible, and has at least $(1 - n^{-\beta})2^n$ codewords. $\square$

We need to combine these short trace reconstruction codes into a longer one in Theorem 1.4. The following notion helps prevent these short codes from being confused with the other components of our construction.

**Definition 2.6.** A string $w$ is *m-protected* if it can be written as $w = 0^m w^\circ 1^m$, where $w^\circ$ starts with a 1, ends with a 0, and every substring of $w^\circ$ of length $m' \geq m/4$ has between $\frac{m'}{4}$ and $\frac{3m'}{4}$ 1s (inclusive). In any $m$-protected string $w$, we let $w^\circ$ denote the string $w$ with the leading $m$ 0s and the trailing $m$ 1s deleted. We refer to $w^\circ$ as the *interior* of $w$. A code is $m$-protected if all of its codewords are $m$-protected.

We use short codes which are both $m$-protected and trace reconstructible in our construction. The following Lemma (see Appendix A.1 for details) shows that these codes exist.

**Lemma 2.7.** *For all $q \in (0,1)$ and $\beta \geq 150$, there exists an absolute constant $\varepsilon_0 = \varepsilon_0(\beta, q) > 0$ such that the following holds. For all $\varepsilon \in (0, \varepsilon_0)$ and $n \geq 8\beta \frac{1}{\varepsilon} \log \frac{1}{\varepsilon}$, if $m = \lfloor \beta \log n \rfloor$ and $T = T_{q,6\beta}^{(avg)}(n)$, there exist codes of length $n$ and rate at least $1 - \frac{\varepsilon}{2}$ that are $m$-protected and $(T, q, n^{-3\beta})$ trace reconstructible.*

## 2.3 Synchronization strings

Synchronization strings [HS17] are useful tools for turning synchronization errors (insertions and deletions) into erasures (replacing symbol with the symbol '?') and substitution errors (replacing symbol with another symbol). Here, we state the construction of synchronization strings that we use and a few useful properties.

**Definition 2.8** (Insertion-deletion distance)**.** Given two strings $S \in \Sigma^n$ and $T \in \Sigma^m$, the *insertion-deletion distance* between $S$ and $T$, denoted $\mathrm{ID}(S,T)$ is the minimum number of characters that needed to be inserted into $S$ and deleted from $S$ to produce $T$.

Insertion-deletion distance is similar to *edit distance* which allows for substitutions at a cost of 1. Observe that if $S$ and $T$ have disjoint character sets, then $\mathrm{ID}(S,T)$ is the sum of their lengths.

**Definition 2.9** ($\eta$-synchronization string)**.** String $S \in \Sigma^n$ is an $\eta$-synchronization string if for every $1 \leq i < j < k \leq n+1$, we have that $\mathrm{ID}(S[i,j], S[j,k]) > (1-\eta)(k-i)$.

**Theorem 2.10** (Theorems 4.5 and 4.7 of [HS18])**.** *For any $\eta \in (0,1)$ and all $n$, one can construct an $\eta$-synchronization string of length $n$ in time $\mathrm{poly}(n)$ over an alphabet of size $6000\eta^{-4}$.*

We now describe some useful properties of synchronization strings. Informally, a *string matching* between two strings describes how to transform one string into the other via insertions and deletions. We use a definition of string matching equivalent to the one introduced in [HS17].

**Definition 2.11** (String matching)**.** For strings $c$ and $c'$ of length $n$ and $n'$, respectively, a *string matching* is a strictly increasing partial function $i^* : [n'] \nrightarrow [n]$ such that, for all $j$ in the domain of $i^*$, we have $c_{i^*(j)} = c'_j$. Given a string matching, an index $j \in [n']$ is called *successfully transmitted* if it is in the domain of $i^*$, and is called an *insertion* otherwise. An element $i \in [n]$ is called a *deletion* if it is not in the codomain of $i^*$.

A $(n, \delta)$-*indexing algorithm for a string $S$* takes as input a string $S'$ of length $n'$ with an unknown string matching $i^* : [n'] \nrightarrow [n]$ having at most $n\delta$ insertions and deletions and outputs an index in $[n] \cup \{\bot\}$ for every index in $[n']$. We say the algorithm *decodes index $j \in [n']$ correctly* under a string matching $i^*$ if it outputs $i^*(j)$ for index $j$ when $i^*(j)$ exists and outputs $\bot$ if it does not exist.

A *misdecoding* of an algorithm is a successfully transmitted, incorrectly decoded index $j \in [n']$. An indexing algorithm is *error free* if every $j \in [n']$ is correctly decoded or is assigned $\perp$.

Haeupler and Shahrasbi proved many results showing that synchronization strings yield indexing algorithms with few misdecodings. In this work, we use the following two results.

**Theorem 2.12** (Theorem 5.10 of [HS17]). *Let $S$ be an $\eta$-synchronization string of length $n$. Then there exists an $(n, \delta)$-indexing algorithm for $S$ guaranteeing at most $\frac{2n\delta}{1-\eta}$ misdecodings. Furthermore, this algorithm runs in time $O(n^4)$*

**Theorem 2.13** (Theorem 6.18 of [HS17]). *Let $S$ be an $\eta$-synchronization string of length $n$. There exists a linear time error-free deletion-only $(n, \delta)$-indexing algorithm for $S$ guaranteeing at most $\frac{\eta}{1-\eta} \cdot n\delta$ misdecodings.*

## 2.4 Binary deletion channel codes

The following lemma gives codes for the $\mathrm{BDC}_q$ with failure probability at most $\delta$ and length $O(\log \delta^{-1})$. In our application, we take $\delta = \mathrm{poly} \frac{1}{\varepsilon}$, where $1 - \varepsilon$ is the rate of our code. A similar construction appears in [GL19] (Proof of Theorem 1). We provide a proof in Appendix A.2 for completeness.

**Lemma 2.14.** *For all $q \in (0, 1)$ and positive integers $K$ and $m$, there exists a binary code $C : [2^K] \to \{0, 1\}^{3Km}$ where every codeword has exactly $2K$ runs, all of which have length either $m$ or $2m$ and decodable in linear time under the $\mathrm{BDC}_q$ with failure probability at most $6K \cdot 2^{-(1-q)m/20}$.*

**Remark 2.15.** The code above has rate $\frac{1}{3m}$, which approaches 0 as $m$ grows. Using a construction similar to [GL19], if we drop the requirement of runs having length exactly $m$ or $2m$, it is possible to achieve a failure probability $2^{-\Omega(m)}$ with a code of rate $c(1 - q)$ for some absolute $c > 0$. We use the result in Lemma 2.14 as the proof is simpler and the result is sufficient.

## 2.5 High rate error correcting codes

Our constructions leverage high rate (rate $1 - \varepsilon$) error correcting codes that are polynomial time encodable and decodable from a $\mathrm{poly}(\varepsilon)$ fraction of worst-case substitution errors. For the details of these constructions and their parameters, see Appendix A.3.

For our binary upper bound, it suffices to use the following variant of a construction by Justesen [Jus72]. Conveniently, it gives codes for *all* sufficiently large $n$, rather than only infinitely many $n$. This property is necessary for Remark 1.6, where we wish to take $\varepsilon \to 0$ as $n \to \infty$ in our binary upper bound construction.

**Proposition 2.16.** *For every $\varepsilon \in (0, \frac{1}{2})$ and $\Sigma$ whose size is a power of 2, there exists an $n_0 = \tilde{\Theta}(\frac{1}{\varepsilon^2})$ such that, for all $n \geq n_0$, there exists a code of length $n$ over alphabet $\Sigma$ of rate $1 - \varepsilon$ that is encodable and decodable in time $O_\varepsilon(n^2)$ from up to a fraction $\frac{\varepsilon^2}{500 \log \frac{1}{\varepsilon}}$ of worst-case substitution errors.*

It would suffice to use Proposition 2.16 for our large alphabet construction (Theorem 1.12) as well. However, using the following error correcting code of Guruswami and Indyk [GI05] allows linear time encoding/decoding of our large alphabet construction.

**Proposition 2.17.** *For every $\varepsilon \in (0, \frac{1}{2})$ and $\Sigma$ whose size is a power of 2, there exist an infinite family of codes over $\Sigma$ of rate $1 - \varepsilon$ encodable in linear time and decodable in linear time from up to a fraction $\frac{1}{40}\varepsilon^3$ of worst-case substitution errors.*

# 3 Optimal number of traces for large alphabet codes

We begin by describing the upper and lower bounds for coded trace reconstruction over a large alphabet. Many of the tools used in this section are important building blocks for the analysis of coded trace reconstruction over a binary alphabet.

## 3.1 Upper bound

*Proof of Theorem 1.12.* We start by defining a few parameters for our construction.

**Parameters.** Let $T = \lceil \log_{1/q} \frac{160}{\varepsilon^3} \rceil$. Let $q' = \frac{1+q}{2}$ and $\eta = \frac{\varepsilon^3}{160T}$. Let $\Sigma_S$ be an alphabet such that there exist $\eta$-synchronization strings over $\Sigma_S$, and assume $|\Sigma_S|$ is a power of 2. We may take $|\Sigma_S| = O_q(\text{poly} \frac{1}{\varepsilon})$ by Theorem 2.10.

**Code.** Let $C_1$ be a length $n$ erasure code over an alphabet $\Sigma_C$ of size $|\Sigma_S|^{\lceil 2/\varepsilon \rceil}$, rate at least $1 - \frac{\varepsilon}{2}$, and decodable from a $\frac{\varepsilon^3}{40}$ fraction of worst-case substitution errors, given by Proposition 2.17. Let $s_1, s_2, \ldots, s_n$ be an $\eta$-synchronization string over alphabet $\Sigma_S$. Let $\Sigma = \Sigma_C \times \Sigma_S$. Let $C$ be a code with encoding $\mathcal{M} \to \Sigma^n$ whose codewords are $(c_1, s_1), \ldots, (c_n, s_n)$ for codewords $(c_1, \ldots, c_n) \in C$.

**Decoding algorithm.** For $t \in [T]$, let $z^{(t)} = (x_1^{(t)}, y_1^{(t)}), \ldots, (x_{n^{(t)}}^{(t)}, y_{n^{(t)}}^{(t)})$ be the $t$th trace, which has length $n^{(t)}$. Call a trace $z^{(t)}$ for $t \in [T]$ *useful* if $n^{(t)} \geq (1 - q') \cdot n$.

1. For every useful trace $z^{(t)}$, run the error-free deletion-only $(n, q')$-indexing algorithm in Theorem 2.13 to obtain indices $i_1^{(t)}, \ldots, i_{n^{(t)}}^{(t)} \in [n] \cup \{\bot\}$.

2. For $i = 1, \ldots, n$, if there exists a useful $t \in [T]$ and index $j \in [n^{(t)}]$ such that $i_j^{(t)} = i$, then let $\hat{c}_i = x_j^{(t)}$. Otherwise, let $\hat{c}_i = \bot$.

3. Run the erasure decoding for $C_1$ on the string $(\hat{c}_1, \ldots, \hat{c}_n)$ to obtain a message in $\mathcal{M}$.

**Efficiency.** The code $C_1$ and synchronization string can each be constructed in polynomial time. Since $C_1$ has linear time encoding, so does our code. Decoding takes time $O(n \log \frac{1}{\varepsilon})$: the indexing algorithm for synchronization strings takes linear time by Theorem 2.10 and we run it $T$ times, and decoding the code $C_1$ from the resulting erasures takes linear time by Proposition 2.17.

**Rate.** The rate of the code $C_1$ is at least $1 - \frac{\varepsilon}{2}$, so there are $|\Sigma_C|^{n(1 - \frac{\varepsilon}{2})} = |\Sigma|^{n(1 - \frac{\varepsilon}{2}) \cdot \frac{\log |\Sigma_C|}{\log |\Sigma|}} \geq |\Sigma|^{n(1-\varepsilon)}$ codewords. The inequality follows as $\frac{\log |\Sigma_C|}{\log |\Sigma|} > 1 - \frac{\varepsilon}{2}$ Hence, the rate of $C$ is at least $1 - \varepsilon$.

**Analysis.** First, the probability that some trace is not useful is equal to the probability that a binomial $B(n, 1 - q)$ is at most $(1 - q')n = \frac{1-q}{2}n$, which, by the Chernoff bound, is at most $e^{-(1-q)n/8}$. Thus, the probability that there exists a trace that is not useful is, by the union bound, at most $T \cdot e^{-(1-q)n/8} \leq 2^{-\Omega(n)}$.

For all useful $t \in [T]$, $z^{(t)}$ is obtained from applying at most $q'n$ deletions to $c$. Thus, the $(n, q')$ indexing-algorithm in Theorem 2.13 succeeds with at most $\frac{\eta}{1-\eta} \cdot nq' < 2\eta n$ misdecodings. Hence, for all $j \in [n^{(t)}]$, we either have $i_j^{(t)} = \bot$ or $j$ is correctly decoded, in which case $x_j^{(t)} = c_{i_j}$. We conclude that, for all $i = 1, \ldots, n$, we either have $\hat{c}_i = c_i$ or $\hat{c}_i = \bot$. We now simply need to lower bound the number of $\hat{c}_i$ that are not $\bot$. If every trace is useful, for each index $i$ with $\hat{c}_i = \bot$, either $(c_i, s_i)$ is deleted in every trace or some trace has a misdecoding at the image of $(c_i, s_i)$. The expected number of symbols $(c_i, s_i)$ deleted in every trace is $q^T n$, so by the Chernoff bound 2, the probability that there are more than $2q^T n$ symbols deleted in every trace is $2^{-\Omega_q(n)}$. Across all traces, the total number of misdecodings is at most $T \cdot 2\eta n$ by above. Thus, with probability at least $1 - 2^{-\Omega_q(n)}$, there are at most $2q^T n + 2T\eta n < \frac{\varepsilon^3}{40}n$ indices $i$ with $\hat{c}_i = \bot$. Hence, as the code $C_1$ tolerates $\frac{\varepsilon^3}{40} \cdot n$ errors (and thus erasures), we decode our message correctly. □

## 3.2 Lower bound

*Proof of Theorem 1.13.* For brevity, let $DC_q$ denote the deletion channel with deletion probability $q$. Let $EC_q$ denote the erasure channel with erasure probability $q$. That is $EC_q$ takes an input string and independently with probability $q$ replaces each symbol with the symbol '?'.

We show that a $(T, q, o(1))$ trace reconstructible code over the $DC_q$ is a code for $EC_{q^T}$ with block error probability $o(1)$. To do this, we show that we can turn an output of $EC_{q^T}$ into $T$ independent outputs of $DC_q$. From a single symbol sent over $EC_{q^T}$, one can produce $T$ independent copies of the symbol sent across $EC_q$: if the output is an erasure, return $T$ erasures, and if the output is the original symbol, return the output of $T$ independent copies of the symbol over $EC_q$, conditioned on not all outputs being erasures. Using the above, from a single output from $EC_{q^T}$, one symbol at a time, produce $T$ independent outputs over $EC_q$, and replace the erasures with deletions to obtain $T$ independent outputs over $DC_q$, as desired. Since the capacity of $EC_{q^T}$ is $1 - q^T$ (see e.g. [Sha48]), we have that our code cannot be $(T, q, o(1))$ trace reconstructible when $1 - \varepsilon > 1 - q^T$, i.e. $T < \log_{1/q} \frac{1}{\varepsilon}$. $\square$

# 4 Upper bound on traces for binary codes

In this section, we prove Theorem 1.4.

## 4.1 Proof sketch

As the proof of Theorem 1.4 is involved, we start with a sketch of the proof. Throughout this proof sketch, fix $q$ to be some constant between 0 and 1. We prove Theorem 1.4 when $n$ is any sufficiently large multiple of a constant (the constant is $n_R + n_S = \Theta_q(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon})$ in the proof). To extend to all sufficiently large $n$ we simply pad the beginning of codewords in an existing code with 0s.

The proof uses concatenation on top of the construction for Theorem 1.12. Recall that the code in Theorem 1.12 is obtained by "zipping" codewords $r_1, \ldots, r_n \in \Sigma_R^n$ from a high-rate error correcting code with a fixed synchronization string $s_1, \ldots, s_n \in \Sigma_S^n$, where $|\Sigma_R| \geq |\Sigma_S|^{\Omega(1/\varepsilon)}$. We call the elements of $\Sigma_R$ *content symbols* and the elements of $\Sigma_S$ *synchronization symbols*.

**A first attempt.** Naively, we could concatenate the code in Theorem 1.12 of rate $1 - \Theta(\varepsilon)$ over the large alphabet $\Sigma_R \times \Sigma_S$ with binary code $C_R$ with encoding $\text{Enc}_R : \Sigma_R \times \Sigma_S \to \{0, 1\}^{n_R}$ of length $n_R = \tilde{\Theta}(\frac{1}{\varepsilon})$ and rate $1 - \Theta(\varepsilon)$ that is decodable from $T = \exp(O(\log^{1/3} \frac{1}{\varepsilon}))$ traces, giving a concatenated code of rate $1 - \Theta(\varepsilon)$ (such a code exists by [HPP18]). In this way, the binary codewords are of the form $\text{Enc}_R(r_1, s_1) || \cdots || \text{Enc}_R(r_n, s_n)$. Then, perhaps, from $T$ traces, we could run the trace reconstruction algorithm for $C_R$ to recover guesses $(\hat{r}_i, \hat{s}_i)$ for $(r_i, s_i)$, and then run the outer decoding to correct any errors/insertions/deletions.

**The problem and the fix.** The problem with the above approach is that we need to recover the synchronization information of the inner codewords *before* we run the inner trace reconstruction algorithm: we do not know, for instance, where the trace of $\text{Enc}_R(r_1, s_1)$ ends and the trace of $\text{Enc}_R(r_2, s_2)$ starts. To fix this, we need the following key idea: separately encode the content symbol $r_i$ and the synchronization symbol $s_i$. Further, in order to ensure that the encoded content bits and the encoded synchronization bits are not confused, we ensure that (1) the encoded synchronization bits only have long runs and (2) the encoded content bits are relatively dense in both 0s and 1s in every small interval (with the exception of one long run at the beginning and end of

the string). This yields the encoding of $(r_1, s_1), \ldots, (r_n, s_n)$ depicted below.

$$\underbrace{\underbrace{0\ldots0}_{m\text{-bit buffer}} \underbrace{\text{interior of } a_1}_{n_R - 2m \text{ bits}} \underbrace{1\ldots1}_{m\text{-bit buffer}}}_{a_1 = \text{Enc}_R(r_1)} \; \Big\| \; \underbrace{\underbrace{0\ldots0}_{k_1 \in \{m, 2m\}} \underbrace{1\ldots1}_{\ell_1} \cdots \underbrace{0\ldots0}_{k_K} \underbrace{1\ldots1}_{\ell_K}}_{b_1 = \text{Enc}_S(s_1): \; 2K \text{ long runs}} \; \Big\| \; \cdots \; \Big\| \; \text{Enc}_R(r_n) \; \Big\| \; \text{Enc}_S(s_n)$$

**Code construction sketch.** We take our outer error correcting code $C_{out}$ to have length $n_{out}$, rate $1 - \Theta(\varepsilon)$, tolerate a $\Theta(\varepsilon^3)$ fraction of worst-case substitution errors, and with alphabet $\Sigma_R$ equal in size to the number of codewords in $C_R$, which we may take to be a power of two by arbitrarily throwing out at most half of the codewords. Such a code exists by Proposition 2.17. We think of $n_{out}$ as growing and all other parameters as fixed. We take a synchronization string $s_1, \ldots, s_{n_{out}}$ with constant synchronization parameter $\eta = \Theta(1)$. We take the length of the encoding $\text{Enc}_R(r_i)$ of $r_i$ to be $\Theta(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon})$, and the length of the encoding $\text{Enc}_S(s_i)$ to be $\Theta(\log \frac{1}{\varepsilon})$, so the rate is at least $1 - \varepsilon$. We ensure that the encoding of $r_i$ is recoverable from $T$ traces with failure probability at most $O(\varepsilon^{100})$. We also ensure that each encoded word of $r_i$ is $m$-protected in the sense of Definition 2.6. The average-case trace reconstruction results of Holden, Pemantle, and Peres [HPP18] implies that such a code exists (see Lemma 2.7). We also ensure that the encoding of $s_i$ is recoverable from *one* trace of the $\text{BDC}_q$ with failure probability $\varepsilon^{100}$. Note that, since the synchronization parameter $\eta$ is a constant, we have that $|\Sigma_S|$ is a constant, so such a code $C_S$ with encoding $\text{Enc}_S : \Sigma_S \to \{0, 1\}^{\Theta(\log(1/\varepsilon))}$ for the $\text{BDC}_q$ exists (see Lemma 2.14).

**Decoding algorithm sketch.** Our decoding algorithm is depicted in Figure 4.1 and divides into three steps.

1. (Trace alignment) For each $t \in [T]$, for all $i \in [n_{out}]$, determine an estimate $\widehat{\tau^{(t)}(a_i)}$ for the bits from the $i$th content symbol

2. (Inner trace reconstruction) For $i \in [n_{out}]$, run the trace reconstruction for the code $C_R$ on $\widehat{\tau^{(1)}(a_i)}, \ldots, \widehat{\tau^{(T)}(a_i)}$ to recover an estimate for $\hat{r}_i$.

3. (Outer error correction) Run the error correction for $C_{out}$ on the estimates $\hat{r}_1, \ldots, \hat{r}_{n_{out}}$.

**Decoding analysis sketch.** Let $a_i = \text{Enc}_R(r_i)$ and $b_i = \text{Enc}_S(s_i)$ be the binary encodings of the $i$th content symbol and $i$th synchronization symbol, respectively. We call $a_i$ a *content block* and $b_i$ a *synchronization block*. For $t \in [T]$ and $i \in [n_{out}]$, let $\tau^{(t)}(a_i)$ and $\tau^{(t)}(b_i)$ denote the images of the $i$th content symbol and $i$th synchronization symbol, respectively, in the $t$th trace.

The key to the analysis is that, by using the indexing algorithm for synchronization strings, with high probability for every trace $t$, the estimates of all but a $O(\varepsilon^{100})$ fraction of the images $\widehat{\tau^{(t)}(a_i)}$ are *exactly* correct, i.e. satisfy $\tau^{(t)}(a_i) = \widehat{\tau^{(t)}(a_i)}$. This is because, with high probability, in each trace, we can find $(1 - O(\varepsilon^{100}))n_{out}$ pairs of strings $(x_j^{(t)}, y_j^{(t)})$ equal to some $(\tau^{(t)}(a_i), \tau^{(t)}(b_i))$ by simply scanning the trace. Then, from the substrings $\tau^{(t)}(b_i)$, we can recover a $1 - O(\varepsilon^{100})$ fraction of the synchronization symbols $s_i$. Using the synchronization symbols, we run the indexing algorithm for the synchronization string to match the pairs $(x_j^{(t)}, y_j^{(t)})$ to the correct index $i \in [n_{out}]$, so that, in each trace, $1 - O(\varepsilon^{100})$ fraction of the pairs are indexed correctly. This produces $(1 - O(\varepsilon^{100}))n_{out}$ accurate estimates $\widehat{\tau^{(t)}(a_i)}$ in every trace.

If the above holds, by the union bound, for all but a $O(T \cdot \varepsilon^{100}) \leq O(\varepsilon^{99})$ (recall $T = \exp(O_q(\log^{1/3}(\frac{1}{\varepsilon}))) = \varepsilon^{-o(1)}$, and assume $\varepsilon$ is sufficiently small) fraction of indices $i \in [n_{out}]$, the image of the inner codeword $a_i$ is correctly determined in *every* trace. Among these indices $i \in [n_{out}]$,
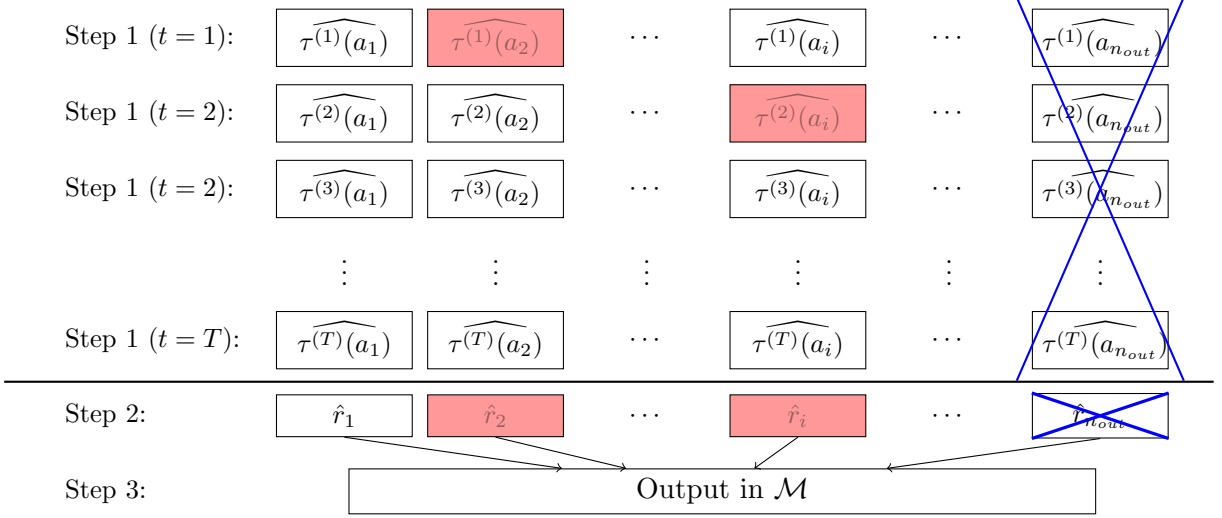
14

Figure 1: Decoding: Inner trace reconstruction and outer error correction. Index $i$ is incorrect only if (i) some trace $t$ incorrectly guesses the image of $a_i$ (shaded red), or (ii) the inner trace reconstruction procedure $\text{Dec}_R$ fails (X-ed out blue).

we expect the inner trace reconstruction algorithm to fail on a $O(\delta_R)$ fraction, and for the rest of these indices,

$$\hat{r}_i = \text{Dec}_{in}(\widehat{\tau^{(1)}(a_i)}, \ldots, \widehat{\tau^{(T)}(a_i)}) = \text{Dec}_{in}(\tau^{(1)}(a_i), \ldots, \tau^{(T)}(a_i)) = r_i.$$

Thus, the fraction of indices $i \in [n_{out}]$ for which $\hat{r}_i \neq r_i$ is $O(\gamma T + \delta_R) = O(\varepsilon^{99})$ with high probability, which, by our choice of parameters, is less than the fraction of substitution errors tolerated by our outer code. Hence, the outer error correction succeeds with high probability, as desired. We point out that all of the big-Os in the outer error fractions do not have a dependence on $q$, so the parameters of the outer code do not need to depend on $q$.

## 4.2 Construction

First we define the code.

**Parameters.**[4] Let $\beta = \frac{10^4}{(1-q)^3}$, let $\tilde{\varepsilon}_0 = \tilde{\varepsilon}_0(\beta, q)$ be given by Lemma 2.7. Let $n_R \stackrel{\text{def}}{=} \lfloor 10^4 \beta \frac{1}{\varepsilon} \log(\frac{1}{\varepsilon}) \rfloor$. This is the length of our inner codeword ("R" for "reconstruction"). Let $T \stackrel{\text{def}}{=} T_{q,6\beta}^{(\text{avg})}(n_R)$. This is the number of traces we use. Let $\delta_R \stackrel{\text{def}}{=} n_R^{-3\beta}$. This is an upper bound on the failure probability of the inner code's reconstruction algorithm. By Theorem 2.3, $T \leq \exp(O(\log^{1/3} n_R)) < \varepsilon^{o(1)}$. Thus, for $\varepsilon$ sufficiently small, we have (i) $T < \frac{1}{\varepsilon}$, (ii) $\varepsilon < \beta^{-1}$, and (iii) $\varepsilon < \tilde{\varepsilon}_0$. For the rest of the proof, assume $\varepsilon$ is such that all three items hold.

Let $m \stackrel{\text{def}}{=} \lfloor \beta \log n_R \rfloor$. This is the size of a "buffer". Let $m' \stackrel{\text{def}}{=} \frac{1}{2}(1-q)m$. This is the threshold for deciding whether a run in a trace is interpreted as a buffer or not. Let $\eta \stackrel{\text{def}}{=} \frac{1}{3}$ be the synchronization parameter. Let $K \stackrel{\text{def}}{=} 20$. This is the number of bits in a synchronization symbol. Let $n_S \stackrel{\text{def}}{=} 60m$. This is the number of bits in an encoded synchronization symbol. Let $\delta_S \stackrel{\text{def}}{=} 6K \cdot 2^{-(1-q)m/40}$. This is an upper bound on the probability a synchronization symbol is decoded correctly. Let

---
[4]We make no attempt to optimize the constants in the proof.

| Parameter | Value | $\lim_{\varepsilon \to 0}$ | Description |
|---|---|---|---|
| $q$ | | | Deletion probability |
| $\varepsilon$ | | | Constructed code has rate $1 - \varepsilon$ |
| $\beta$ | $\frac{10^4}{(1-q)^3}$ | $\Theta_q(1)$ | Large constant |
| $n_R$ | $\lfloor 10^4 \beta \frac{1}{\varepsilon} \log(\frac{1}{\varepsilon}) \rfloor$ | $\tilde{\Theta}_q(\frac{1}{\varepsilon})$ | Content block length |
| $T$ | $T_{q,6\beta}^{(\text{avg})}(n_R)$ | $\varepsilon^{-o_q(1)}$ | Number of traces used |
| $\delta_R$ | $n_R^{-3\beta}$ | $O(\varepsilon^{100})$ | Upper bound on inner code $C_R$'s trace reconstruction failure probability |
| $m$ | $\lfloor \beta \log n_R \rfloor$ | $\Theta_q(\log \frac{1}{\varepsilon})$ | Size of a buffer |
| $m'$ | $\frac{1}{2}(1-q)m$ | $\Theta_q(\log \frac{1}{\varepsilon})$ | Threshold for interpreting an output run as a buffer |
| $\eta$ | $\frac{1}{3}$ | $\Theta(1)$ | Synchronization parameter |
| $K$ | $20$ | $\Theta(1)$ | Number of bits in synchronization symbol: $|\Sigma_S| = 2^K$ |
| $n_S$ | $60m$ | $\Theta_q(\log \frac{1}{\varepsilon})$ | Number of bits in encoded synchronization symbol |
| $\delta_S$ | $6K \cdot 2^{-(1-q)m/40}$ | $O(\varepsilon^{100})$ | Upper bound on inner code $C_S$'s decoding failure probability |
| $\gamma$ | $2^{-(1-q)m/80}$ | $O(\varepsilon^{100})$ | Upper bound on probability content block is "incorrectly parsed" |
| $n_{out}$ | $\to \infty$ | $\to \infty$ | Outer code length |
| $\delta_{out}$ | $\frac{1}{50000}\varepsilon^3$ | $\Omega(\varepsilon^3)$ | Lower bound on the outer code's error tolerance |

$\gamma \overset{\text{def}}{=} 2^{-(1-q)m/80}$. This is an upper bound on the probability an inner codeword is "incorrectly parsed" (defined below). In this way, $\gamma T < \varepsilon^{100}$. Let $\delta_{out} = \frac{1}{50000}\varepsilon^3$ be a bound on the outer code's error tolerance. Throughout this analysis we think of $q, \varepsilon, \beta, m, m', K, n_R, n_S, \delta_R, \delta_S, \delta_{out}, T, \gamma, \eta$ as fixed, and $n_{out}$, the length of the outer code defined below, as growing.

**Inner codes.** By our choice of parameters, $\beta \geq 150$, $\varepsilon < \tilde{\varepsilon}_0$, $n_R \geq 8\beta\frac{1}{\varepsilon}\log\frac{1}{\varepsilon}$, and $T = T_{q,6\beta}^{(\text{avg})}(n_R)$. By Lemma 2.7 there exists a code $C_R$ of length $n_R$ and rate $1 - \frac{\varepsilon}{2}$ with message set $\Sigma_R$ and encoding function $\text{Enc}_R : \Sigma_R \to \{0,1\}^{n_R}$ all of whose codewords are $m$-protected (as $m = \lfloor \beta \log n_R \rfloor$), and that is $(T, q, \delta_R)$ trace reconstructible. By removing at most half of the codewords (arbitrarily), we may assume the alphabet size $|\Sigma_R|$ is a power of 2, and the rate is at least $1 - \frac{\varepsilon}{2} - \frac{1}{n_R}$. Let the corresponding decoding function be $\text{Dec}_R : (\{0,1\}^*)^T \to \Sigma_R$.

Let $\eta = \frac{1}{3}$, and let $s_1, \ldots, s_n$ be an $\eta$-synchronization string of length $n$ over alphabet $\Sigma_S$ of size $2^K$: such strings exist by Theorem 2.10 and are constructible in polynomial time. By Lemma 2.14, there exists a code $C_S$ with encoding function $\text{Enc}_S : \Sigma_S \to \{0,1\}^{n_S}$ and decoding function $\text{Dec}_S : \{0,1\}^* \to \Sigma_S$ that is decodable under the $\text{BDC}_q$ with failure probability at most $\delta_S$, all of whose codewords start with a 0, end with a 1, and have runs of length exactly $m$ or $2m$.

**Outer code.** Let $C_{out} : \mathcal{M} \to \Sigma_R^{n_{out}}$ be a code of length $n_{out}$ and rate $1 - \frac{\varepsilon}{10}$ over the alphabet $\Sigma_R$ correcting a $\frac{(\frac{\varepsilon}{10})^2}{500 \log \frac{10}{\varepsilon}} > \delta_{out}$ fraction of worst-case errors, given by Proposition 2.16.

**Encoding.** Our encoding is as follows.

1. Take a message and encode it with $C_{out}$ to obtain symbols $r_1, \ldots, r_n \in \Sigma_R$.

2. Let $a_i = \text{Enc}_R(r_i)$ and $b_i = \text{Enc}_S(s_i)$. We call $a_i$ a *content block* and $b_i$ a *synchronization*

16

*block.*

3. Concatenate $c = a_1||b_1||a_2||b_2||\cdots||a_n||b_n$.

$$\underbrace{\underbrace{0\ldots0}_{m\text{-bit buffer}}\;\underbrace{\text{interior of } a_1}_{n_R-2m \text{ bits}}\;\underbrace{1\ldots1}_{m\text{-bit buffer}}}_{a_1=\text{Enc}_R(r_1)} \;\bigg|\bigg|\; \underbrace{\underbrace{0\ldots0}_{k_1\in\{m,2m\}}\;\underbrace{1\ldots1}_{\ell_1}\cdots\underbrace{0\ldots0}_{k_K}\;\underbrace{1\ldots1}_{\ell_K}}_{b_1 \,=\, \text{Enc}_S(s_1):\; 2K \text{ long runs}}\;\bigg|\bigg|\;\cdots\;\bigg|\bigg|\;\text{Enc}_R(r_n)\;\bigg|\bigg|\;\text{Enc}_S(s_n)$$

**Length and Rate.** The length is $n_{out} \cdot (n_R + n_S)$. The outer code rate is $(1 - \frac{\varepsilon}{10})$, the inner code $C_R$ rate is $1 - \frac{\varepsilon}{2} - \frac{1}{n_R}$, and the synchronization symbols multiply the rate by $1 - \frac{n_S}{n_R+n_S} > 1 - \frac{60\beta\log n_R}{10^4\beta\frac{1}{\varepsilon}\log\frac{1}{\varepsilon}} = 1 - \frac{\varepsilon}{10}$. The total rate is thus at least $(1 - \frac{\varepsilon}{10})(1 - \frac{\varepsilon}{2} - \frac{1}{n_R})(1 - \frac{\varepsilon}{10}) > 1 - \varepsilon$.

**Decoding.** Let $z^{(1)}, \ldots, z^{(T)}$ be the traces. We use the following notation for the "Trace Alignment" step of the decoding below. The crucial elements of the Trace Alignment step's analysis are given in Definition 4.2, Lemma 4.3, and Lemma 4.6. For every trace, call a (maximal) run of length greater than $m'$ a *decoded buffer*. Call every bit in a decoded buffer a *decoded buffer bit*, and call all other bits *decoded content bits*. For every trace, define a *decoded content block* to be a substring of the form $0^{t_0}w1^{t_1}$, where the first $t_0$ 0s and the last $t_1$ 1s each form a decoded buffer, and $w$ is a nonempty string of decoded content bits. Note in particular that $w$ must start with a 1 and end with a 0. If two decoded content blocks overlapped, say $0^{t_0}w1^{t_1}$ and $0^{t'_0}w'1^{t'_1}$, then $0^{t_0}$ is the same run of bits as $0^{t'_0}$, because $w$ does not consist of any decoded buffers. Likewise, $1^{t_1}$ is the same run of bits as $1^{t'_1}$ so $w = w'$. Therefore, any two decoded content blocks are disjoint.

We can thus enumerate the decoded content blocks of a trace $t$ in order $x_1^{(t)}, \ldots, x_{n^{(t)}}^{(t)}$, where $n^{(t)}$ is the number of decoded content blocks in trace $t$. For each decoded content block $x_j^{(t)}$, we define the associated *decoded synchronization block* $y_j^{(t)}$ as the substring between $x_j^{(t)}$ and $x_{j+1}^{(t)}$ (or the end of the string, if $k = n^{(t)}$).[5] Our decoding algorithm is as follows.

1. (Trace alignment) For each trace $t \in [T]$, compute $\widehat{\tau^{(t)}(a_1)}, \ldots, \widehat{\tau^{(t)}(a_{n_{out}})}$ as follows:

   (a) Compute the decoded content blocks $x_1^{(t)}, \ldots, x_{n^{(t)}}^{(t)}$ of $z^{(t)}$ along with their associated decoded synchronization blocks $y_1^{(t)}, \ldots, y_{n^{(t)}}^{(t)}$.

   (b) For all $j \in [n^{(t)}]$, decode a synchronization symbol $\hat{s}_j^{(t)} \stackrel{\text{def}}{=} \text{Dec}_S(y_j^{(t)}) \in \Sigma_S$ from the decoded synchronization block.

   (c) From the string $\hat{s}_1 \ldots \hat{s}_{n^{(t)}}$, obtain indices $\hat{i}_1^{(t)}, \ldots, \hat{i}_{n^{(t)}}^{(t)}$ using the $(n, 13\gamma)$ indexing algorithm in Theorem 2.12.

   (d) For $j = 1, \ldots, n^{(t)}$, let $\widehat{\tau^{(t)}(a_{\hat{i}_j})} \stackrel{\text{def}}{=} x_j^{(t)}$, and let $\widehat{\tau^{(t)}(a_i)} = \perp$ for $i \notin \{\hat{i}_1, \ldots, \hat{i}_{n^{(t)}}\}$. Here, $\widehat{\tau^{(t)}(a_i)}$ is a string denoting our guess for the image of $a_i$ in the $t$th trace.

2. (Inner trace reconstruction) For $i = 1, \ldots, n$, let $\hat{r}_i = \text{Dec}_R(\widehat{\tau^{(1)}(a_i)}, \ldots, \widehat{\tau^{(T)}(a_i)}) \in \Sigma_R$.

3. (Outer error correction) Run $\text{Dec}_{out}(\hat{r}_1, \ldots, \hat{r}_n)$ to obtain a message in $\mathcal{M}$.

---

[5]Note that there may be some bits at the beginning of the string that are neither in decoded content blocks nor in decoded synchronization blocks.
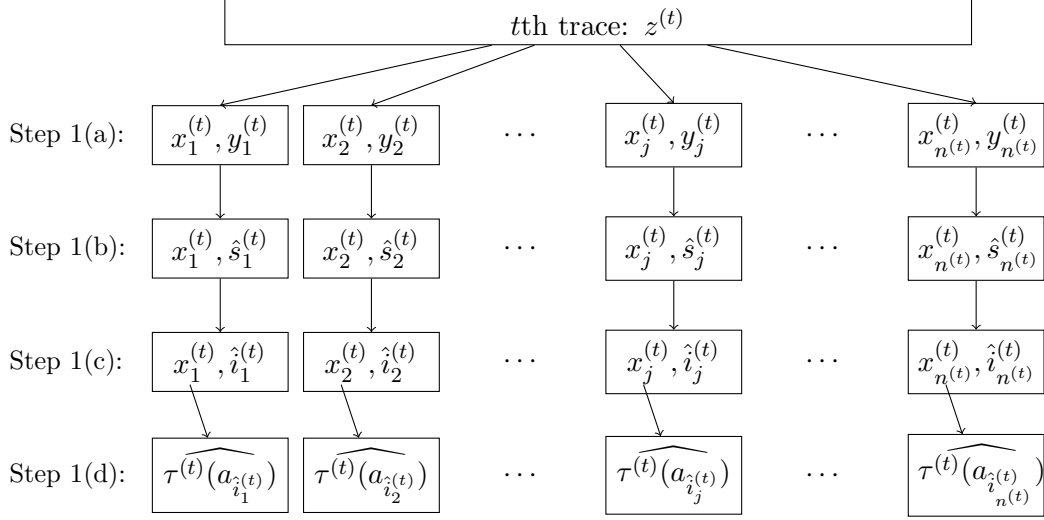
Figure 2: Decoding: Trace alignment

**Run time.** The encoding runs in time $O(n^2)$: the outer encoding runs in time $O(n^2)$, and each of the $O(n)$ inner encodings runs in time $O_\varepsilon(1)$.

The decoding runs in $O_\varepsilon(n^4)$ time. Determining the decoded content blocks and decoded synchronization blocks can be done in linear time $O(T \cdot n_{out}) = O(n_R \cdot n_{out}) \leq O(n)$. Here, we used that $T < n_R$ as $\varepsilon$ is sufficiently small, so in particular, this decoding time has no dependence on $q$. The code $C_S$ is decodable from deletions in linear time by Lemma 2.14, so computing all synchronization symbols $\tilde{s}_j^{(t)}$ takes $O(n)$ time. The indexing step for the synchronization string takes time $O(n_{out}^4)$ in each trace, so all indexing steps take time $O(T \cdot n_{out}^4) \leq O(n^4)$. Thus, the entire trace alignment step takes time $O(n^4)$. Each inner trace reconstruction step takes $n_R^{1+o(1)}$ time [HPP18], so the entire inner trace reconstruction step (Step 2) takes $n_R^{1+o(1)} \cdot n_{out} \leq n^{1+o(1)}$ time by running the decoder for [HPP18]. The outer error correction (Step 3) runs in time $O(n^2)$. The total decoding time is thus $O(n^4)$.

## 4.3 Analysis

### 4.3.1 Notation

Throughout this analysis, we think of all the bits of $c, z^{(1)}, \ldots, z^{(t)}$ as distinct. We may informally refer to $z^{(t)}$ as "trace $t$". Let $\tau^{(t)}$ denote the $t$th *deletion pattern*, i.e. a map from bits in the codeword $c$ to bits in $t$th trace $z^{(t)}$.[6] In this way, $\tau^{(t)}(c) = z^{(t)}$, and if $c'$ is a substring of $c$, then $\tau^{(t)}(c')$ is a substring of $z^{(t)}$. Throughout the analysis, if $w$ and $w'$ are substrings of a trace $z^{(t)}$, we write $w =^{(t)} w'$ to indicate that they are the same substring of trace $z^{(t)}$.

### 4.3.2 Correctly parsed indices

**Definition 4.1** (Spurious buffer). We say that a *spurious buffer* of a trace $z^{(t)}$ is a decoded buffer that is a substring of the image of a content block's interior, i.e. a substring of $\tau^{(t)}(a_i^\circ)$ for some $i$.

---

[6]Formally, $\tau^{(t)}$ is an injective and surjective partial function from the bits of $c$ to the bits of $z^{(t)}$, such that the undeleted bits of $c$ form the domain of $\tau^{(t)}$, and these bits are mapped in order to the bits of $z^{(t)}$.

**Definition 4.2.** For $t \in [T], i \in \{0, \ldots, n_{out} + 1\}$, we say index $i$ is *intact in trace $t$* if $i = 0$, $i = n + 1$, or all of the following hold:

1. At least $m'$ of the $m$ leading 0s of $a_i$ are not deleted in $\tau^{(t)}$

2. At least $m'$ of the $m$ trailing 1s of $a_i$ are not deleted in $\tau^{(t)}$

3. $\tau^{(t)}(a_i)$ has no spurious buffers.

4. The image of all runs of $b_i$ under $\tau^{(t)}$ have length at least $m'$.

5. $\mathrm{Dec}_S(\tau^{(t)}(b_i)) \neq s_i$

For $t \in [T]$, call an index $i \in \{1, \ldots, n_{out}\}$ *correctly parsed in trace $t$* if indices $i - 1, i$, and $i + 1$ are all intact in trace $t$, and *incorrectly parsed in trace $t$* otherwise.

Note that the event "$i$ is intact in trace $t$" depends only on the images of $a_i$ and $b_i$ under $\tau^{(t)}$, and hence all such events are independent. The following lemma justifies the terminology in Definition 4.2.

**Lemma 4.3.** *If index $i \in [n_{out}]$ is correctly parsed in trace $t \in [T]$, there exists an index $j \in [n^{(t)}]$ such that $\tau^{(t)}(a_i) =^{(t)} x_j^{(t)}$, and $\tau^{(t)}(b_i) =^{(t)} y_j^{(t)}$.*

*Proof.* First, we show that, since indices $i - 1$ and $i$ are intact in trace $t$, the image $\tau^{(t)}(a_i)$ of content block $a_i$ forms a decoded content block in $z^{(t)}$. The image of the substring $b_{i-1}||a_i||b_i$ of our codeword $c$ under the $t$th deletion pattern $\tau^{(t)}$ is $\tau^{(t)}(b_{i-1})||\tau^{(t)}(a_i)||\tau^{(t)}(b_i)$. Further, the substring $\tau^{(t)}(b_{i-1})$ ends in a 1 (property 2), the substring $\tau^{(t)}(a_i)$ begins with a 0 and ends with a 1 (properties 1 and 2), and the substring $\tau^{(t)}(b_i)$ starts with a 0 (property 1). Hence $\tau^{(t)}(a_i)$ starts with a decoded buffer of 0s (property 1), ends with a decoded buffer of 1s (property 2), and has no other decoded buffers (property 3), so $\tau^{(t)}(a_i)$ is a decoded content block. Thus, there exists some $j \in [n^{(t)}]$ such that $\tau^{(t)}(a_i) =^{(t)} x_j^{(t)}$.

Similarly, since indices $i$ and $i + 1$ are intact in trace $t$, the substring $\tau^{(t)}(a_{i+1})$ is a decoded content block by the same reasoning. Since index $i$ is intact in trace $t$, all the bits in the image of $b_i$ under $\tau^{(t)}$ are decoded buffer bits (property 4), so there are no decoded content bits between $\tau^{(t)}(a_i)$ and $\tau^{(t)}(a_{i+1})$. Thus, $\tau^{(t)}(a_i)$ and $\tau^{(t)}(a_{i+1})$ are consecutive decoded content blocks, so $\tau^{(t)}(b_i) =^{(t)} y_j^{(t)}$, as desired. $\square$

As an immediate corollary, we have the following lemma.

**Lemma 4.4.** *Let $t \in [T]$. If there are at least $k$ indices $i \in [n_{out}]$ that are correctly parsed in trace $t$, then the sequences of pairs $(\tau^{(t)}(a_1), s_1), \ldots, (\tau^{(t)}(a_n), s_n)$ and $(x_1^{(t)}, \hat{s}_1^{(t)}), \ldots, (x_{n^{(t)}}^{(t)}, \hat{s}_{n^{(t)}}^{(t)})$ have a common subsequence of length $k$.*

*Proof.* By Lemma 4.3, the sequences of pairs of strings $(\tau^{(t)}(a_1), \tau^{(t)}(b_1)), \ldots, (\tau^{(t)}(a_n), \tau^{(t)}(b_n))$ and $(x_1^{(t)}, y_1^{(t)}), \ldots, (x_{n^{(t)}}^{(t)}, y_{n^{(t)}}^{(t)})$ have a common subsequence of length $k$, namely the subsequence corresponding to the $k$ correctly parsed pairs $(\tau^{(t)}(a_i), \tau^{(t)}(b_i))$. The result follows as $\mathrm{Dec}_{in}(\tau^{(t)}(b_i)) = s_i$ and $\mathrm{Dec}_{in}(y_j^{(t)}) = \hat{s}_j^{(t)}$, so we can apply the $\mathrm{Dec}_{in}$ operator to the second element in each pair of each sequence to obtain the desired result. $\square$

### 4.3.3 Bounding incorrectly parsed indices

The following lemma guarantees that for all $t \in [T]$ and $i \in [n_{out}]$, the string $\tau^{(t)}(a_i)$ has no spurious buffers with high probability.

**Lemma 4.5.** *For any $t \in [T], i \in [n_{out}]$, the expected number of spurious buffers in $\tau^{(t)}(a_i)$ is at most $e^{-(1-q)m/40}$.*

*Proof.* Call a substring of a content block's interior $a_i^\circ$ *spurious buffer indicator* if it has length exactly $\lfloor \frac{m}{4} \rfloor$ and at least one of the following occur:

1. All of the 0s are deleted.

2. All of the 1s are deleted.

3. At least $m'$ of the 0s are not deleted.

4. At least $m'$ of the 1s are not deleted.

Consider a spurious buffer of 0s. Let $c'$ denote the minimal substring of $c$ containing the spurious buffer's preimage. In $c'$, all the 1s are deleted and at least $m'$ of the 0s are not deleted. Hence, any substring or superstring of this preimage $c'$ of length $\lfloor \frac{m}{4} \rfloor$ is a spurious buffer indicator. The same holds for the preimage of a spurious buffer of 1s. Hence, we may identify each spurious buffer with a corresponding spurious buffer indicator of $a_i^\circ$ of length $\frac{m}{4}$ (breaking ties arbitrarily). Because distinct spurious buffers are disjoint, the spurious buffer indicator are distinct. Thus, the number of runs of decoded buffer bits is bounded above by the number of spurious buffer indicator. Since $a_i$ is $m$-protected, every substring of length $\lfloor \frac{m}{4} \rfloor$ has at least $\lfloor \frac{m}{16} \rfloor$ 0s and at least $\lfloor \frac{m}{16} \rfloor$ 1s. Thus, the probability all the 0s are deleted is at most $q^{-\lfloor m/16 \rfloor} < e^{-(1-q)m/20}$, and the probability all the 1s are deleted is also at most $e^{-(1-q)m/20}$. Any run of $\lfloor \frac{m}{4} \rfloor$ bits has at most $\lfloor \frac{m}{4} \rfloor$ 0s, and the probability that at least $m'$ 0s are not deleted is bounded above by the probability that the binomial random variable $B(\lfloor \frac{m}{4} \rfloor, 1-q)$ is at least $m' = \frac{m(1-q)}{2}$, which, by the Chernoff bound (2) is at most $e^{-(1-q)\lfloor m/12 \rfloor}$. The same probability holds for the 1s. Hence, the probability that any run of $\lfloor \frac{m}{4} \rfloor$ bits is a spurious buffer indicator is at most $4e^{-(1-q)m/20}$, so the expected number of spurious buffer indicator, and thus the number of spurious buffers, is at most $4n_R \cdot e^{-(1-q)m/20}$ by linearity of expectation, which, by definition of $m$ and $n_R$, is at most $e^{-(1-q)m/40}$. $\qquad\square$

**Lemma 4.6.** *For $t \in [T], i \in [n_{out}]$, the probability index $i$ is intact is at least $1 - \gamma$.*

*Proof.* We bound the probability each property in Definition 4.2 fails.

1. Since $a_i$ begins with $m = \frac{2m'}{1-q}$ leading 0s, the expected number of undeleted 0s among the leading 0s is distributed as the binomial $B(m, 1-q)$, which has mean $2m'$. Hence, the probability that property 1 fails is at most probability that this binomial is less than $m'$, which, by the Chernoff bound (1), is at most $e^{-m'/4}$.

2. By the same reasoning, the probability that property 2 fails is at most $e^{-m'/4}$.

3. If property 3 fails, $\tau^{(t)}(a_i)$ has a spurious buffer. Since $\mathbf{Pr}[X > 0] \le \mathbf{E}[X]$ for all nonnegative integer random variables $X$, the probability $\tau^{(t)}(a_i)$ has a spurious buffer is at most $e^{-(1-q)/40}$ by Lemma 4.5.

4. By construction, $b_j$ has $2K$ runs of length at least $m$. Property 4 fails if some run has less than $m'$ non-deleted bits. The number of non-deleted bits in a run is distributed as one of the binomials $B(m, 1-q)$ or $B(2m, 1-q)$. Hence, by the Chernoff bound (1) and union bound, property 4 fails with probability at most $2K \cdot e^{-m'/4}$.

5. Property 5 fails with probability at most $\delta_S \leq 6K \cdot 2^{-(1-q)m/40}$ by the definition of code $C_S$.

By the union bound, the probability that index $i$ is not intact in trace $t$ is at most

$$\underbrace{e^{-m'/4}}_{\text{property 1}} + \underbrace{e^{-m'/4}}_{\text{property 2}} + \underbrace{4n_R \cdot e^{-(1-q)m/12}}_{\text{property 3}} + \underbrace{2K \cdot e^{-m'/4}}_{\text{property 4}} + \underbrace{6K \cdot 2^{-(1-q)m/40}}_{\text{property 5}}$$

$$\leq (2 + 4n_R + 8K)2^{-(1-q)m/40}.$$

For our choice of $n_R$, we have $2 + 4n_R + 8K < 2^{(1-q)m/80}$, so the probability index $i$ is not intact in trace $t$ is at most $2^{-(1-q)m/80} = \gamma$. □

A simple Chernoff bound gives the following corollary.

**Corollary 4.7.** *The probability that there exists a $t \in [T]$ with more than $6\gamma n_{out}$ incorrectly parsed indices is at most $2^{-\Omega(n_{out})}$.*

*Proof.* For $t \in [T]$ and $i \in [n_{out}]$, let $\mathcal{E}_{t,i}$ denote the event that index $i$ is not intact in trace $t$. The events $\mathcal{E}_{t,i}$ are all independent, and any such event happens with probability at most $\gamma$ by Lemma 4.6. Hence, by the Chernoff bound (2) and the union bound, the probability there exists a $t \in [T]$ with more than $2\gamma n_{out}$ events $\mathcal{E}_{t,i}$ occurring is at most $T \cdot 2^{-\gamma n_{out}/3}$. Hence, as the number of incorrectly parsed indices $i$ in a trace $t$ is at most 3 times the number of non-intact pairs, the probability that there exists a $t$ with more than $6\gamma n_{out}$ incorrectly parsed pairs is also at most $T \cdot 2^{-\gamma n_{out}/3} = 2^{-\Omega(n_{out})}$. □

### 4.3.4 Most traces of content bits are recovered

Note that a trace may have more than $n_{out}$ decoded content blocks if there are spurious buffers. The next lemma shows that, with high probability, this does not happen too much.

**Lemma 4.8.** *For any $t$, with probability $1 - \exp(-\Omega(n_{out}))$, we have $n^{(t)} \leq (1 + \gamma)n_{out}$.*

*Proof.* The number of decoded blocks is bounded above by $n_{out}$ plus the number of spurious buffers. Let $X_i$ denote the number of spurious buffers in block $\tau^{(t)}(a_i)$. By the definition of spurious buffer, $X_1, \ldots, X_{n_{out}}$ are all independent. Additionally, each block has at most $n_R$ bits, so it certainly has at most $n_R$ spurious buffers. Hence, $\frac{X_i}{n_R} \in [0,1]$ for all $i$. By the Lemma 4.5, $\mathbf{E}[\frac{X_i}{n_R}] \leq \frac{e^{-(1-q)m/40}}{n_R} = \frac{\gamma^2}{n_R}$ for all $i = 1, \ldots, n$. Thus, by the Chernoff bound (3) on the variables $\frac{X_1}{n_R}, \ldots, \frac{X_{n_{out}}}{n_R}$, the probability that $X_1 + \cdots + X_{n_{out}} \geq \gamma n_{out}$ is at most $2^{-\gamma^2 n_{out}/3} \leq 2^{-\Omega(n_{out})}$. □

**Lemma 4.9.** *Let $t \in [T]$. If there are at least $(1 - 6\gamma)n_{out}$ correctly parsed indices in trace $t$ and if $n^{(t)} \leq (1 + \gamma)n_{out}$, then there are at least $(1 - 46\gamma)n_{out}$ indices $i$ such that $\tau^{(t)}(a_i) = \widehat{\tau^{(t)}(a_i)}$.*

*Proof.* Suppose there are at most $\gamma n_{out}$ incorrectly parsed indices in trace $t$ and also that $n^{(t)} \leq (1 + \gamma)n$. By Lemma 4.4, there is a common subsequence between $(\tau^{(t)}(a_1), s_1), \ldots, (\tau^{(t)}(a_n), s_n)$ and $(x_1^{(t)}, \hat{s}_1^{(t)}), \ldots, (x_{n^{(t)}}^{(t)}, \hat{s}_{n^{(t)}}^{(t)})$ of length at least $(1 - 6\gamma)n$. Since $n^{(t)} \leq (1 + \gamma)n$, there exists a string matching between $(\tau^{(t)}(a_1), s_1), \ldots, (\tau^{(t)}(a_n), s_n)$ and $(x_1^{(t)}, \hat{s}_1^{(t)}), \ldots, (x_{n^{(t)}}^{(t)}, \hat{s}_{n^{(t)}}^{(t)})$ with at most

21

$6\gamma n$ deletions and at most $7\gamma n$ insertions, for a total of at most $13\gamma n$ insertions or deletions. In particular, there are at least $(1-6\gamma)n_{out}$ correctly transmitted indices. This string matching gives a corresponding string matching between $s_1, \ldots, s_n$ and $\hat{s}_1^{(t)}, \ldots, \hat{s}_{n^{(t)}}^{(t)}$. In this string matching, for the correctly transmitted symbols $j \in [n^{(t)}]$, let $i_j^{(t)}$ be such that $\hat{s}_j = s_{i_j^{(t)}}$, so that $x_j^{(t)} = \tau^{(t)}(a_{i_j^{(t)}})$. By Theorem 2.12, the $(n, 13\gamma)$-indexing algorithm for the $\eta$-synchronization string $s_1, \ldots, s_n$ guarantees that there are at most $\frac{2}{1-\eta} \cdot 13\gamma n < 40\gamma n$ misdecodings. That is, there are at most $40\gamma n$ correctly transmitted indices $j = 1, \ldots, n^{(t)}$ such that $\hat{i}_j^{(t)} \neq i_j^{(t)}$. For the other correctly transmitted indices $j \in n^{(t)}$, we have $\tau^{(t)}(a_{i_j^{(t)}}) = x_j^{(t)} = \widehat{\tau^{(t)}(a_{\hat{i}_j^{(t)}})}$. Hence, there are at least $(1-6\gamma)n - 40\gamma n = (1-46\gamma)n$ indices $i$ such that $\tau^{(t)}(a_i) = \widehat{\tau^{(t)}(a_i)}$. $\qquad\square$

### 4.3.5 Finishing the proof

We now complete the proof. The next lemma shows that, with high probability, most of the inner trace reconstructions succeed "in theory". That is, they succeed assuming the trace alignment steps recovered the images of the $a_i$'s successfully in all traces.

**Lemma 4.10.** , *With probability $1 - \exp(-\Omega(n_{out}))$, for all but at most $2\delta_R$ fraction of indices $i \in [n_{out}]$, we have $\mathrm{Dec}_R(\tau^{(1)}(a_i), \ldots, \tau^{(T)}(a_i)) = r_i$.*

*Proof.* Call an index $i$ *incorrect* if $\mathrm{Dec}_R(\tau^{(1)}(a_i), \ldots, \tau^{(T)}(a_i)) \neq r_i$. The probability an index $i$ is incorrect is at most $\delta_R$ as $C_R$ is $(T, q, \delta_R)$ trace reconstructible and $\tau^{(1)}(a_i), \ldots, \tau^{(T)}(a_i)$ are independent traces of $a_i$. Furthermore, for all $i$, the events that $i$ is incorrect are independent of each other. The expected number of incorrect $i$ is at most $\delta_R n$, so by the Chernoff bound (2), the probability the number of incorrect $i$ is larger than $2\delta_R n$ is at most $2^{-\delta_R n/3}$, as desired. $\qquad\square$

Now we can prove Theorem 1.4.

*Proof of Theorem 1.4.* By Corollary 4.7, Lemma 4.8, and Lemma 4.10, with probability $1 - 2^{-\Omega(n)}$, all the following occur:

1. Every trace has at most $6\gamma n_{out}$ incorrectly parsed indices.

2. For all $t \in [T]$, we have $n^{(t)} \leq (1 + \gamma)n_{out}$.

3. All but at most $2\delta_R n_{out}$ indices $i \in [n_{out}]$ satisfy $\mathrm{Dec}_R(\tau^{(1)}(a_i), \ldots, \tau^{(T)}(a_i)) = r_i$.

We show that, when all of the above occur, the decoding succeeds. By Lemma 4.9 and properties 1 and 2 above, there are at least $(1 - 46\gamma)n_{out}$ indices $i$ such that $\tau^{(t)}(a_i) = \widehat{\tau^{(t)}(a_i)}$. Thus, there are at least $(1 - 46\gamma T)n_{out}$ indices $i$ such that $\tau^{(t)}(a_i) = \widehat{\tau^{(t)}(a_i)}$ for all $t \in [T]$. Hence, by property 3 above, there are at least $(1 - 46\gamma T - 2\delta_R)n_{out} > (1 - \delta_{out})n_{out}$ indices $i$ such that $\tau^{(t)}(a_i) = \widehat{\tau^{(t)}(a_i)}$ for all $t \in [T]$ and $\mathrm{Dec}_R(\tau^{(1)}(a_i), \ldots, \tau^{(T)}(a_i)) = r_i$. For all such indices $i$, we have

$$\hat{r}_i = \mathrm{Dec}_R\left(\widehat{\tau^{(1)}(a_i)}, \ldots, \widehat{\tau^{(T)}(a_i)}\right) = \mathrm{Dec}_R\left(\tau^{(1)}(a_i), \ldots, \tau^{(T)}(a_i)\right) = r_i.$$

As $C_{out}$ tolerates $\delta_{out} n_{out}$ errors, the outer decoding finds the correct message in $\mathcal{M}$, as desired. $\qquad\square$

## 4.4 Extending to all sufficiently large $n$

For some $n_0 = \text{poly} \frac{1}{\varepsilon}$, the above construction gives error probability $\delta < \frac{1}{3}$ for all $n \geq n_0$ that are multiples of $n_R + n_S = \Theta(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon})$: synchronization strings exist for all lengths $n_{out}$, the codes in Proposition 2.16 exist for all lengths $n_{out}$ at least $\Omega(\frac{1}{\varepsilon^3})$, and the overall error probability is bounded by $2^{-\Omega(\gamma^2 n)}$, so it suffices to take $n \gg \Omega(\frac{1}{\gamma^2}) = \text{poly} \frac{1}{\varepsilon}$.

To extend to larger values of $n$, we take our constructed code of length $n - (n \mod (n_R + n_S))$ and pad the beginning of all codewords with $(n \mod (n_R + n_S))$ 0s. This multiplies the rate by at least $1 - \frac{n_R + n_S}{n} < 1 - o(\varepsilon)$, so the rate is still at least $1 - \varepsilon$. Lemma 4.3 still holds for all indices $i \in [n_{out}]$ except possibly the first ($x_1^{(t)}$ may have some extra 0s padded to $\tau^{(t)}(a_1)$), so the number of incorrect parsed indices is now bounded by $6\gamma n_{out} + 1$. For $n$ (and thus $n_{out}$) a sufficiently large polynomial in $\frac{1}{\varepsilon}$, the total fraction of incorrect (outer) content symbols $\hat{r}_i \neq r_i$ is still less than $\delta_{out}$ with high probability, so the decoding still succeeds with high probability.

# 5 Lower bound on traces for binary codes

In this section, we prove the following theorem, which implies Theorem 1.8.

**Theorem 5.1.** *Let $q \in (0, 1)$ and $\varepsilon < \frac{1}{4}$. Let $m = \lfloor \sqrt{\frac{1/\varepsilon}{128 \log(1/\varepsilon)}} \rfloor$ and $T = T_{q,0}^{(\text{avg})}(m) - 1$. Then, for all $\delta \in (0, 1)$, there exists $n_0 = O_\delta(1/\varepsilon^2)$ such that all rate $1 - \varepsilon$ codes of length at least $n_0$ are not $(T, q, \delta)$-trace reconstructible.*

## 5.1 Mutual information and Shannon's theorem

Recall that the entropy of a random variable $X$ is $H(X) \overset{\text{def}}{=} -\sum_x \mathbf{Pr}[X = x] \log \mathbf{Pr}[X = x]$. For two random variables $X$ and $Y$ their conditional entropy of $Y$ given $X$ is defined to be $H(X|Y) \overset{\text{def}}{=} \sum_y \mathbf{Pr}[Y = y] \cdot H(X|Y = y)$, where $H(X|Y = y)$ is the entropy of the random variable $X$ given that $Y = y$. From this, we can define their mutual information $I(X, Y)$ to be $I(X, Y) \overset{\text{def}}{=} H(X) - H(X|Y)$. A *discrete memoryless channel* has finite input alphabet $\mathcal{X}$ and finite output alphabet $\mathcal{Y}$, and is given by a matrix $w(y|x)$, denoting, for each $x \in \mathcal{X}$, a distribution over received symbols $y \in \mathcal{Y}$. With $w$, any probability distribution over $\mathcal{X}$ gives a joint distribution on $\mathcal{X}, \mathcal{Y}$.

Given a discrete memoryless channel $w$, we say a code $C \subset \mathcal{X}^n$ is *decodable with failure probability at most $\delta$* if there exists a map $f : \mathcal{Y}^n \to \mathcal{X}^n$ such that, for all $x_1 \cdots x_n \in C$, we have

$$\mathbf{Pr}_{y_i \sim w(\cdot|x_i)} [f(y_1, \ldots, y_n) \neq x_1 \cdots x_n] \leq \delta.$$

We need the following result, which provides a strong converse to Shannon's noisy channel coding theorem [Sha48].

**Theorem 5.2** (e.g. Theorem 3.3.1 of [Wol78]). *Let $w(\cdot|\cdot)$ define a discrete memoryless channel with inputs $\mathcal{X}$ and outputs $\mathcal{Y}$. Let*

$$R_{cap} \overset{\text{def}}{=} \max_{p(x)} I(X, Y), \tag{4}$$

*where the maximum is taken over probability distributions on $\mathcal{X}$, and let $\gamma > 0$. Then, for all $\delta \in (0, 1)$, there exists $n_0 = O_\delta(\frac{1}{\gamma^2})$ such that, for all $n \geq n_0$ there do not exist codes of rate $\frac{R_{cap} + \gamma}{\log |\mathcal{X}|}$*

*decodable with failure probability at most $\delta$ under the channel $w(\cdot|\cdot)$.*[78]

A classic result known as Fano's inequality can be used to lower bound the mutual information $I(X,Y)$ in (4) with a quantity involving the probability of error. The following result by Tebbe and Dwyer [TI68] helps bound the mutual information $I(X,Y)$ in the other direction, and is useful in our proof.

**Lemma 5.3** ([TI68]). *Let $\delta \in (0,1)$. Suppose we are given a probability distribution $\mathcal{D}$ over $\mathcal{X} \times \mathcal{Y}$ such that, for all maps $f : \mathcal{Y} \to \mathcal{X}$, we have $\mathbf{Pr}_{X,Y}[f(Y) \neq X] \geq \delta$. Then $H(X|Y) \geq \frac{\delta}{2}$.*

## 5.2   Random to coded lower bound

Let $\mathcal{X}_m \stackrel{\text{def}}{=} \{0,1\}^m$ and $\mathcal{Y}_{m,T} \stackrel{\text{def}}{=} (\{0,1\}^{\leq m})^T$. For all $q$, $m$ and $T$, there is a natural channel with inputs $\mathcal{X}_m$ and outputs $\mathcal{Y}_{m,T}$. We induce a joint probability distribution on $\mathcal{X}_m, \mathcal{Y}_{m,T}$ as follows. Let $\lambda$ be a probability density function on $\mathcal{X}_m$. Let $X_\lambda \sim \mathcal{X}_m$ be the distribution where $x$ is sampled with probability $\lambda(x)$. We let $Y_\lambda$ be the output of $T$ independent traces of the sampled $x \sim X_\lambda$ across the $\text{BDC}_q$.

Note that since $H(X_\lambda) \leq m$, for any distribution $X_\lambda \sim \mathcal{X}_m$, we have that $I(X_\lambda, Y_\lambda) \leq m$. We show in Lemma 5.4 that if $T \approx T_{q,0}^{(\text{avg})}(m)$, then this upper bound can be improved by a significant amount. This upper bound is subsequently used in Theorem 5.1 to show a limitation of the capacity of coded trace reconstruction.

**Lemma 5.4.** *Let $\beta \geq 1$. Suppose $T = T_{q,0}^{(\text{avg})}(m) - 1$ for $m \geq 32$. For all probability distributions $X_\lambda$ on $\mathcal{X}_m$, if $Y_\lambda \in \mathcal{Y}_{m,T}$ is distributed as $T$ independent traces of $X_\lambda$, then*

$$I(X_\lambda, Y_\lambda) \leq m - \frac{1}{32m \log m}.$$

*Proof.* Let $\mathcal{X}'$ be the elements of $\mathcal{X}$ with $\lambda(x) \geq \frac{1}{(m \log m)2^m}$. We consider two cases.

**Case 1: $|\mathcal{X}'| \leq 2^{m-1/3}$.** We have

$$
\begin{aligned}
I(X_\lambda, Y_\lambda) \;&\leq\; H(X_\lambda) \\
&=\; \sum_{x \in \mathcal{X}'} \lambda(x) \log \frac{1}{\lambda(x)} + \sum_{x \notin \mathcal{X}'} \lambda(x) \log \frac{1}{\lambda(x)} \\
&\leq\; \log |\mathcal{X}'| + \sum_{x \notin \mathcal{X}'} \lambda(x) \log \frac{1}{\lambda(x)} \hspace{3cm} (5) \\
&\leq\; \log |\mathcal{X}'| + \sum_{x \notin \mathcal{X}'} \frac{1}{(m \log m)2^m} \cdot \log((m \log m)2^m) \hspace{1.2cm} (6) \\
&\leq\; \log |\mathcal{X}'| + 2^m \cdot \frac{1}{m(\log m)2^m} \log(m(\log m)2^m) \\
&=\; m - \frac{1}{3} + \frac{m + \log m + \log\log m}{m \log m} \;<\; m - \frac{1}{3} + \frac{1}{4} \;<\; m - \frac{1}{32m \log m}. \hspace{0.5cm} (7)
\end{aligned}
$$

In (5) we used that $\sum_{x \in \mathcal{X}} \lambda(x) \leq 1$ and that $z \log \frac{1}{z}$ is concave. In (6) we used that $z \log \frac{1}{z}$ is increasing for $z < 1/3$. In (7) we used that $m$ is sufficiently large.

---

[7]The quantity $R$ is often referred to as the *capacity* of the channel

[8]Typically the normalizing term $\frac{1}{\log |\mathcal{X}|}$ is not present when stating Shannon's capacity theorem. This is because the "rate" used in Shannon capacity is often defined as $\frac{\log |C|}{n}$, whereas the rate for us is defined as $\frac{\log |C|}{n \log |\mathcal{X}|}$.

**Case 2:** $|\mathcal{X}'| \geq 2^{m-1/3}$.

For this case, a similar argument appears in [HL$^+$20] (Proposition 4.1). Let $\sigma(x)$ be the uniform distribution on the elements of $\mathcal{X}$. Let $\mu(x)$ be the uniform distribution on the elements of $\mathcal{X}'$. Consider any trace reconstruction algorithm $f : \mathcal{Y}_{m,T} \to \mathcal{X}_m$. Note that

$$\mathbf{Pr}[f(Y_\sigma) \neq X_\sigma] \leq \frac{|\mathcal{X} \setminus \mathcal{X}'|}{|\mathcal{X}|} + \frac{|\mathcal{X}'|}{|\mathcal{X}|} \mathbf{Pr}[f(Y_\mu) \neq X_\mu].$$

By definition, $T \stackrel{\text{def}}{=} T_{q,0}^{(\text{avg})}(m) - 1$ and $|\mathcal{X}'| \geq 2^{-1/3}|\mathcal{X}|$, so

$$\mathbf{Pr}[f(Y_\mu) \neq X_\mu] \geq 2^{-1/3} \mathbf{Pr}[f(Y_\sigma) \neq X_\sigma] - (1 - 2^{-1/3}) \geq \frac{2^{-1/3}}{3} - 1 + 2^{-1/3} > \frac{1}{8}.$$

Let $\nu(x)$ be the probability distribution on $\mathcal{X}$ given by

$$\nu(x) \stackrel{\text{def}}{=} \frac{\lambda(x) - \frac{1}{2m \log m}\mu(x)}{1 - \frac{1}{2m \log m}}.$$

We have $|\mathcal{X}'| \geq \frac{1}{2}|\mathcal{X}|$, so $\mu$ assigns probability at most $\frac{2}{2^m}$ to each element of $\mathcal{X}'$. Since $\lambda$ assigns probability at least $\frac{1}{(m \log m)2^m}$ to each element of $\mathcal{X}'$, $\nu(x) \geq 0$ for all $x$. Furthermore, it is easy to check that $\sum_{x \in \mathcal{X}} \nu(x) = 1$, so $\nu(x)$ is a legitimate probability distribution. We can sample from $\lambda$ as follows: with probability $\frac{1}{2m \log m}$ sample from $\mu$, otherwise, sample from $\nu$. Thus, for any recovery algorithm $f : \mathcal{Y}_{m,T} \to \mathcal{X}_m$.

$$\begin{aligned}
\mathbf{Pr}[f(Y_\lambda) \neq X_\lambda] &= \frac{1}{2m \log m} \mathbf{Pr}[f(Y_\mu) \neq X_\mu] + \left(1 - \frac{1}{2m \log m}\right) \mathbf{Pr}[f(Y_\nu) \neq X_\nu] \\
&\geq \frac{1}{2m \log m} \cdot \mathbf{Pr}[f(Y_\mu) \neq X_\mu] \\
&\geq \frac{1}{2m \log m} \cdot \frac{1}{8} = \frac{1}{16m \log m}.
\end{aligned}$$

The last inequality is uses (5.2). Thus, $H(X_\lambda|Y_\lambda) \geq \frac{1}{32m \log m}$ by Lemma 5.3. We thus may bound

$$I(X_\lambda, Y_\lambda) = H(X_\lambda) - H(X_\lambda|Y_\lambda) \leq \log|\mathcal{X}| - H(X_\lambda|Y_\lambda) \leq m - \frac{1}{32m \log m}.$$

This covers all cases, completing the proof. $\qquad\square$

*Proof of Theorem 5.1.* Recall $m = \lfloor \sqrt{\frac{1/\varepsilon}{128 \log(1/\varepsilon)}} \rfloor$ and $T = T_{q,0}^{(\text{avg})}(m) - 1$. Let $n_0'$ be the constant given by Theorem 5.2 with the parameter $\gamma \stackrel{\text{def}}{=} \varepsilon m$. Let $n_0 \stackrel{\text{def}}{=} m \cdot n_0' \leq O(\frac{1}{\varepsilon^2})$.

We first prove that codes of rate $1 - 2\varepsilon$ are not $(T, q, \delta)$ trace reconstructible when $n$ is any sufficiently large multiple of $m$. Let $C$ be a code that is $(T, q, \delta)$ trace reconstructible when $n \geq n_0$ is a multiple of $m$. We show $C$ must have rate less than $1 - 2\varepsilon$. Let $n_{out} \stackrel{\text{def}}{=} \frac{n}{m}$. For each $i \in [n_{out}]$, given a codeword $c = (c_1, \ldots, c_n) \in C$, let $X_i$ denote the string

$$X_i \stackrel{\text{def}}{=} c_{(i-1)m+1}, c_{(i-1)m+2}, \ldots, c_{im}.$$

Let $Y_i \in \mathcal{Y}_{m,T}$ be a tuple of $T$ of strings distributed as independent traces of $X_i$ under the BDC$_q$. By assumption of our code, it is possible to recover $c$ from $Y_1, \ldots, Y_{n_{out}}$ with failure probability at most

$\delta$: take the trace-wise concatenation of $Y_1, \ldots, Y_{n_{out}}$ and use the trace reconstruction algorithm that is assumed. Hence, the code $C$, when interpreted as a code in $\mathcal{X}^{n_{out}}$, achieves failure probability $\delta$ on the memoryless channel $w(\cdot|\cdot)$ with inputs $\mathcal{X}_m$ and outputs $\mathcal{Y}_{m,T}$ where $Y$ is distributed as $T$ independent traces of $X$. By Lemma 5.4, we have

$$\max_{\lambda \text{ on } \mathcal{X}_m} I(X_\lambda, Y_\lambda) \leq m \left(1 - \frac{1}{32m^2 \log m}\right) \leq m(1 - 4\varepsilon),$$

since $\varepsilon$ sufficiently small. By Theorem 5.2, since $\gamma \overset{\text{def}}{=} \varepsilon m$ and $n_{out} \geq n_0'$, our code $C$, when interpreted as a code in $\mathcal{X}^{n_{out}}$, must have rate less than

$$\frac{1}{\log |\mathcal{X}|} \left(\max_{\lambda \text{ on } \mathcal{X}_m} I(X_\lambda, Y_\lambda) + \gamma\right) < 1 - 2\varepsilon,$$

as desired.

Now suppose $n$ is not a multiple of $m$. Then, suppose for contradiction that $C \subset \{0,1\}^n$ is a code of length $n$ and rate $1-\varepsilon$ that is $(T, q, \delta)$ trace reconstructible. By a simple counting argument, there exists a code $C' \subset \{0,1\}^{n'}$ of rate $1 - \varepsilon - \frac{\varepsilon n'}{n-n'} > 1 - 2\varepsilon$ and a string $w$ such that $c'||w \in C$ for all $c' \in C'$. Furthermore, recovering all codewords of $C$ requires recovering all codewords of the form $c'||w$ for $c' \in C'$. The failure probability of recovering $c'$ from $T$ traces of $c'||w$ is at least the failure probability of recovering $c'$ from $T$ traces of $c'$, which, as we showed, is more than $\delta$, a contradiction. $\qquad\square$

# 6 Conclusion and Open Problems

In this paper, we considered the coded trace reconstruction problem. We obtain lower and upper bounds on the problem which show that the average-case trace reconstruction problem is essentially equivalent to the coded trace reconstruction problem. Even with this contribution, there are still many questions left unanswered.

1. The most fundamental open question in this space is closing the exponential gaps for the worst-case trace reconstruction and average-case trace-reconstruction. For worst-case trace reconstruction, the optimal number of traces is between $\tilde{\Omega}(n^{3/2})$ and $\exp(O(n^{1/3}))$ (or $\exp(O(n^{1/5}))$ for $q \leq 1/2$), and for average-case trace reconstruction, the optimal number of traces is between $\tilde{\Omega}(\log^{5/2} n)$ and $\exp(O(\log^{1/3} n))$.

2. One way to generalize the coded trace reconstruction model considered in this paper to consider a more general synchronization channel, such as with insertions and deletions. For example, such a model could insert $k$ random bits between $x_i$ and $x_{i+1}$ with probability $(1-q)q^k$ and then apply i.i.d. deletions with probability $q$. See the recent survey by Cheraghchi and Ribeiro [CR20] for an overview of various models for random insertions, deletions, substitutions and replications. The authors suspect that similar primitives to those used in this paper could be useful in these more general settings.

3. Another combinatorial variant of this question is *necklace reconstruction*. This question is similar to ordinary trace reconstruction, except a random cyclic shift is also applied to each trace, and the original string needs to be recovered up to an arbitrary cyclic shift. Many protocols for the traditional trace reconstruction problem exploit that the initial prefix of the trace can be easily determined by looking at the prefixes of the traces. For necklace reconstruction, this strategy would no longer work (due to the random shift), so new techniques

26

need to be developed. Even beating $O((1-q)^{-n})$ traces, the probability of receiving the whole necklace as a trace, seems nontrivial. A recent paper [NR20] studies this problem.

4. A challenging question in the context of coded trace reconstruction is formulating other interesting models beyond i.i.d. deletions. Adversarial deletions is not an interesting model because the adversary could delete the same bits on each trace, reducing the problem to the deletion code problem. One possibility of such a model would be adversarial deletions subject to some global constraints–such as the distribution of deletions being approximately $k$-wise independent.

5. Another challenge is coming up with deletion models and codes that more accurately correspond to practical use cases and string lengths. Trace reconstruction as used in DNA computing often considers string of approximately length 100 (e.g., [OAC+18]). Constructing such codes may require different techniques than those used in this paper.

6. We do not know if Theorem 1.12 achieves the smallest alphabet size for $O(\log_{1/q} \frac{1}{\varepsilon})$ traces. It would be interesting to determine the trade-off between alphabet size and number of traces.

# 7  Acknowledgements

# References

[AVDiF19]  Mahed Abroshan, Ramji Venkataramanan, Lara Dolecek, and Albert Guillén i Fàbregas. Coding for deletion channels with multiple traces. In *2019 IEEE International Symposium on Information Theory (ISIT)*, pages 1372–1376. IEEE, 2019.

[BCF+19]  Frank Ban, Xi Chen, Adam Freilich, Rocco A Servedio, and Sandip Sinha. Beyond trace reconstruction: Population recovery from the deletion channel. In *2019 IEEE 60th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 745–768. IEEE, 2019.

[BCSS19]  Frank Ban, Xi Chen, Rocco A Servedio, and Sandip Sinha. Efficient average-case population recovery in the presence of insertions and deletions. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM 2019)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.

[BKKM04]  Tugkan Batu, Sampath Kannan, Sanjeev Khanna, and Andrew McGregor. Reconstructing strings from random traces. *SODA*, 2004.

[CGMR20]  Mahdi Cheraghchi, Ryan Gabrys, Olgica Milenkovic, and Joao Ribeiro. Coded trace reconstruction. *IEEE Transactions on Information Theory*, 2020.

[Cha19]     Zachary Chase. New Lower Bounds for Trace Reconstruction. *arXiv.org*, May 2019.

[Cha20]     Zachary Chase. New upper bounds for trace reconstruction. *arXiv preprint arXiv:2009.03296*, 2020.

[Che18]     Mahdi Cheraghchi. Capacity upper bounds for deletion-type channels. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018*, pages 493–506, 2018.

[CR20]      Mahdi Cheraghchi and João Ribeiro. An overview of capacity results for synchronization channels. *IEEE Transactions on Information Theory*, 2020. to appear.

[CS20]      Roni Con and Amir Shpilka. Explicit and efficient constructions of coding schemes for the binary deletion channel. In *2020 IEEE International Symposium on Information Theory (ISIT)*, pages 84–89. IEEE, 2020.

[DG01]      Suhas Diggavi and Matthias Grossglauser. On transmission over deletion channels. In *Proceedings of the 39th Annual Allerton Conference on Communication, Control, and Computing*, pages 573–582, 2001.

[DM06]      Eleni Drinea and Michael Mitzenmacher. On lower bounds for the capacity of deletion channels. *IEEE Trans. Information Theory*, 52(10):4648–4657, 2006.

[DM07]      Eleni Drinea and Michael Mitzenmacher. Improved lower bounds for the capacity of i.i.d. deletion and duplication channels. *IEEE Trans. Information Theory*, 53(8):2693–2714, 2007.

[DOS17]     Anindya De, Ryan O'Donnell, and Rocco A Servedio. Optimal mean-based algorithms for trace reconstruction. *STOC*, pages 1047–1056, 2017.

[DP09]      Devdatt P. Dubhashi and Alessandro Panconesi. *Concentration of Measure for the Analysis of Randomized Algorithms*. Cambridge University Press, 2009.

[FD10]      Dario Fertonani and Tolga M. Duman. Novel bounds on the capacity of the binary deletion channel. *IEEE Trans. Information Theory*, 56(6):2753–2765, 2010.

[GI05]      Venkatesan Guruswami and Piotr Indyk. Linear-time encodable/decodable codes with near-optimal rate. *IEEE Trans. Information Theory*, 51(10):3393–3400, 2005.

[GL19]      Venkatesan Guruswami and Ray Li. Polynomial time decodable codes for the binary deletion channel. *IEEE Trans. Information Theory*, 65(4):2171–2178, 2019.

[GM19]      Ryan Gabrys and Olgica Milenkovic. Unique reconstruction of coded strings from multiset substring spectra. *IEEE Transactions on Information Theory*, 65(12):7682–7696, 2019.

[HL⁺20]     Nina Holden, Russell Lyons, et al. Lower bounds for trace reconstruction. *Annals of Applied Probability*, 30(2):503–525, 2020.

[HM14]      Bernhard Haeupler and Michael Mitzenmacher. Repeated deletion channels. *2014 IEEE Information Theory Workshop (ITW 2014)*, pages 152–156, August 2014.

[HMPW08]    Thomas Holenstein, Michael Mitzenmacher, Rina Panigrahy, and Udi Wieder. Trace reconstruction with constant deletion probability and related results. *SODA*, 2008.

[HPP18]     Nina Holden, Robin Pemantle, and Yuval Peres. Subpolynomial trace reconstruction for random strings and arbitrary deletion probability. In *Conference On Learning Theory*, pages 1799–1840, 2018.

[HS17]       Bernhard Haeupler and Amirbehshad Shahrasbi. Synchronization Strings: Codes for Insertions and Deletions Approaching the Singleton Bound. *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing - STOC 2017*, pages 33–46, 2017.

[HS18]       Bernhard Haeupler and Amirbehshad Shahrasbi. Synchronization strings: Explicit constructions, local decoding, and applications. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 841–854, 2018.

[Jus72]      Jørn Justesen. Class of constructive asymptotically good algebraic codes. *IEEE Trans. Information Theory*, 18(5):652–656, 1972.

[KM13]      Yashodhan Kanoria and Andrea Montanari. Optimal coding for the binary deletion channel with small deletion probability. *IEEE Trans. Information Theory*, 59(10):6192–6219, 2013.

[KMMP19] Akshay Krishnamurthy, Arya Mazumdar, Andrew McGregor, and Soumyabrata Pal. Trace reconstruction: Generalized and parameterized. In *27th Annual European Symposium on Algorithms (ESA 2019)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.

[KMS10]     Adam Kalai, Michael Mitzenmacher, and Madhu Sudan. Tight asymptotic bounds for the deletion channel with small deletion probabilities. In *2010 IEEE International Symposium on Information Theory*, pages 997–1001. IEEE, 2010.

[Lev01a]     Vladimir I Levenshtein. Efficient reconstruction of sequences. *IEEE Trans. Information Theory*, 47(1):2–22, 2001.

[Lev01b]     Vladimir I Levenshtein. Efficient Reconstruction of Sequences from Their Subsequences or Supersequences. *Journal of Combinatorial Theory*, 2001.

[MPV14a]   Andrew McGregor, Eric Price, and Sofya Vorotnikova. Trace Reconstruction Revisited. In Andreas S. Schulz and Dorothea Wagner, editors, *Algorithms - ESA 2014*, volume 8737, pages 689–700. Springer Berlin Heidelberg, Berlin, Heidelberg, 2014.

[MPV14b]   Andrew McGregor, Eric Price, and Sofya Vorotnikova. Trace Reconstruction Revisited. *ESA*, 8737(2):689–700, 2014.

[NP17]       Fedor Nazarov and Yuval Peres. Trace reconstruction with exp(O(n1/3)) samples. *STOC*, 2017.

[NR20]       Shyam Narayanan and Michael Ren. Circular trace reconstruction. *arXiv preprint arXiv:2009.01346*, 2020.

[OAC+18]   Lee Organick, Siena Dumas Ang, Yuan-Jyue Chen, Randolph Lopez, Sergey Yekhanin, Konstantin Makarychev, Miklos Z Racz, Govinda Kamath, Parikshit Gopalan, Bichlien Nguyen, et al. Random access in large-scale dna data storage. *Nature biotechnology*, 36(3):242, 2018.

[PZ17]   Yuval Peres and Alex Zhai. Average-case reconstruction for the deletion channel: Subpolynomially many traces suffice. In *58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, Berkeley, CA, USA, October 15-17, 2017*, pages 228–239, 2017.

[RD15]   Mojtaba Rahmati and Tolga M. Duman. Upper bounds on the capacity of deletion channels using channel fragmentation. *IEEE Trans. Information Theory*, 61(1):146–156, 2015.

[Sha48]  Claude Shannon. A mathematical theory of communication. *Bell Systems Technical Journal*, 27(3):379423, 1948.

[TI68]   D. Tebbe and Samuel J. Dwyer III. Uncertainty and the probability of error (corresp.). *IEEE Trans. Information Theory*, 14(3):516–518, 1968.

[VT65]   RR Varshamov and GM Tenengolts. Codes which correct single asymmetric errors (in russian). *Automatika i Telemkhanika*, 161(3):288–292, 1965.

[Wol78]  Jacob Wolfowitz. *Coding theorems of information theory*. Springer-Verlag, 1978.

[YGM17]  S. M. Hossein Tabatabaei Yazdi, Ryan Gabrys, and Olgica Milenkovic. Portable and Error-Free DNA-Based Data Storage. *Scientific Reports*, 7:5011, 2017.

# A   Omitted Details

## A.1   Proof of Lemma 2.7

The following lemma shows that a significant fraction of strings of length $n$ are $m$-protected for $m \geq \Omega(\log n)$.

**Lemma A.1.** *Let $m \geq 10^3$ be an integer and $n \in [3m, 2^{m/150}]$, the number of $m$-protected codewords is at least $2^{n-2m-3}$.*

*Proof.* There are $2^{n-2m-2}$ strings of the form $s = 0^m s^\circ 1^m$. Choose one such string at random, so that $s^\circ$ is a uniformly random string in $1||\{0,1\}^{n-2m-2}||0$. For a substring of length $m'$, the probability it has at least $3/4$ or at most $1/4$ fraction of 1s is, by the Chernoff bound (1), at most $2 \cdot 2^{-m'/16}$. Since $m' \geq m/4$, this is at most $2 \cdot 2^{-m/64}$. Since there are at most $n^2$ substrings of $s^\circ$ of length at least $m/4$, by the union bound, the resulting string is *not* $m$-protected with probability at most

$$2n^2 \cdot 2^{-m/64} \leq 2 \cdot 2^{2m/150 - m/64} = 2 \cdot 2^{-m/300} < \frac{1}{2}.$$

Hence, at least half of all strings of the form $0^m s^\circ 1^m$ are $m$-protected, as desired.   □

With the protected strings from Lemma A.1 and the codes for trace reconstruction from Lemma 2.5, we can prove Lemma 2.7. Intuitively, both the codes with protected codewords and codes which are efficiently trace reconstructible are both very large, so we can find the desired code in their intersection.

*Proof of Lemma 2.7.* By Lemma 2.5 with parameter $\beta' = 3\beta$, there exists a code $C_1$ with $|C_1| \geq (1 - n^{-3\beta})2^n = 2^n - 2^{n-3m}$ that is $(T, q, n^{-3\beta})$ trace reconstructible.

Let $C_2$ be the set of length $n$ strings that are $m$-protected. Assume $\varepsilon$ is sufficiently small so that $n \geq 6\frac{1}{\varepsilon}m > 3m$. Note also that by our choice of $n$,

$$2^{m/150} = 2^{\lfloor \beta \log n \rfloor/150} \geq 2^{\log n} = n.$$

Then, by Lemma A.1, we have $|C_2| \geq 2^{n-2m-3}$.

Let $C = C_1 \cap C_2$. We have $|C| = |C_1 \cap C_2| \geq 2^{n-2m-3} - 2^{n-3m} > 2^{n-3m}$, so $C$ has rate at least $1 - \frac{3m}{n} > 1 - \frac{\varepsilon}{2}$. Furthermore, since $C_1$ is $(T, q, n^{-3\beta})$ trace reconstructible, $C$ is as well. $\qquad \square$

## A.2 Proof of Lemma 2.14

In this section, we show how to construct codes for the binary deletion channel of length $O(\log \frac{1}{\delta})$ and failure probability at most $\delta$.

*Proof of Lemma 2.14.* **Encoding.** Map every element $\sigma \in [2^K]$ to a string $\tilde{c}_\sigma \in \{0,1\}^{3K}$ that starts with a 0, ends with a 1, has $K$ runs are length 1, and has $K$ runs are length 2. There are $\binom{2K}{K} \geq 2^K$ such strings as each string is uniquely determined by its sequence of run lengths, so each $\sigma$ can be assigned to a distinct string. Let $c_\sigma$ be $\tilde{c}_\sigma$ with every symbol duplicated $m$ times.

**Decoding.** To decode a received word $s$ under the $\text{BDC}_q$, we first recover $\tilde{c}_\sigma$, and then recover $\sigma$. To recover $\tilde{c}_\sigma$, suppose $s$ is of the form $0^{k_1'} 1^{\ell_1'} \cdots 0^{k_K'} 1^{\ell_K'}$ where $k_i', \ell_i' \geq 1$ for all $i$. If $s$ is not of this form, return an arbitrary symbol in $[2^K]$ (give up). For each $i = 1, \ldots, K$, if $k_i' \geq 1.4(1-q)m$, let $x_i' = 2$, and otherwise let $x_i' = 1$. Similarly, if $\ell_i' \geq 1.4(1-q)m$, let $y_i' = 2$, and otherwise let $y_i' = 1$. The decoding returns the symbol $\sigma'$ such that

$$\tilde{c}_{\sigma'} = 0^{x_1'} 1^{y_1'} \cdots 0^{x_K'} 1^{y_K'}.$$

**Analysis.** The decoding is clearly linear time. To prove correctness, suppose our input symbol $\sigma$ satisfies $c_\sigma = 0^{x_1} 1^{y_1} \cdots 0^{x_K} 1^{y_K}$, where $x_i, y_i \in \{1, 2\}$ for all $i$. Let $k_1, \ell_1, \ldots, k_K, \ell_K$ denote the number of bits not deleted in the corresponding runs $0^{x_1}, 1^{y_1}, \ldots, 1^{y_K}$. We bound the probability each of the following happen.

1. There exists some $i$ such that $k_i = 0$ ($\ell_i = 0$)

2. There exists some $i$ with $x_i = 1$ ($y_i = 1$) such that $k_i \geq 1.4(1-q)m$ ($\ell_i \geq 1.4(1-q)m$).

3. There exists some $i$ with $x_i = 2$ ($y_i = 2$) such that $k_i < 1.4(1-q)m$ ($\ell_i < 1.4(1-q)m$).

If $x_i = 1$, then $k_i$ is distributed as the binomial distribution $B(m, 1-q)$. If $x_i = 2$, then $k_i$ is distributed as the binomial distribution $B(2m, 1-q)$. In either case, we have

$$\mathbf{Pr}[k_i = 0] = \mathbf{Pr}[\ell_i = 0] \leq q^m < e^{-(1-q)m} < 2^{-(1-q)m/20}$$

By the Chernoff bound (2), for $i$ such that $x_i = 1$, we have

$$\mathbf{Pr}[x_i \neq x_i'] = \mathbf{Pr}[k_i \geq 1.4(1-q)m] \leq e^{-\frac{0.4^2}{2+0.4}(1-q)m} < 2^{-(1-q)m/20}$$

On the other hand, for $i$ such that $x_i = 2$, we have, by the Chernoff bound (1),

$$\mathbf{Pr}[x_i \neq x_i'] = \mathbf{Pr}[k_i < 1.4(1-q)m] \leq e^{-\frac{0.3^2}{2}(1-q)m} < 2^{-(1-q)m/20}$$

The same probabilities hold for $y_i$'s. Hence the probability any of events 1, 2, or 3 happen is at most $6K \cdot 2^{-(1-q)m/20}$, as desired. However, if event 1 does not happen then the decoding guarantees that $k_i' = k_i$ and $\ell_i' = \ell_i$ for all $i$. If additionally, events 2 and 3 do not happen, the decoding guarantees that $x_i' = x_i$ and $y_i' = y_i$, and hence $\sigma' = \sigma$. Thus, the decoding fails with probability at most $6K \cdot 2^{-(1-q)m/40}$, as desired. $\qquad \square$

## A.3 High rate error correcting codes

In this section, we show how Proposition 2.17 follows from the construction of Guruswami and Indyk [GI05]. Guruswami and Indyk prove the following.

**Theorem A.2** (Theorem 5 of [GI05]). *For every $\varepsilon > 0$ and any $R \in (0,1)$, there exists a family of binary codes of rate $R$ encodable in linear time and decodable in linear time from up to a fraction $\delta$ of substitution errors, where*

$$\delta \geq \max_{R < r < 1} \frac{(1 - r - \varepsilon)H^{-1}(1 - R/r))}{2}.$$

By setting $R = 1 - \varepsilon'$ and $\varepsilon = \frac{\varepsilon'}{10}$, and taking $r = 1 - \frac{\varepsilon'}{2}$, we have

$$\delta \geq \frac{2\varepsilon'}{5} \cdot H^{-1}\left(\frac{\varepsilon'}{2(1 - \varepsilon'/2)}\right) \geq \frac{2\varepsilon'}{20} \cdot \left(\frac{\varepsilon'}{2(1 - \varepsilon'/2)}\right)^2 \geq \frac{2\varepsilon'}{20} \cdot \left(\frac{\varepsilon'}{2}\right)^2 \geq \frac{(\varepsilon')^3}{40}.$$

Here we used that $H^{-1}(x) \geq \frac{x^2}{4}$ for all $x \in (0,1)$.

Further for every $\Sigma$ whose size is $2^\ell$ a power of 2, every family binary codes of rate $R$ and tolerating a $\delta$ fraction of worst-case substitution errors can be made into a family of codes over $\Sigma$ with the same asymptotic rate and error tolerance: pad each codeword so that its length is a multiple of $\ell$ (this has a negligible effect on the asymptotic rate and error tolerance), then map each length $\ell$ string $b_1, \ldots, b_\ell$ to a unique element of $\Sigma$. For a codeword $c = (c_1, \ldots, c_n) \in \{0,1\}^n$, create a codeword over $\Sigma^{n/\ell}$ whose $i$th symbol is the image of $c_{(i-1)\ell+1}, c_{(i-1)\ell+2}, \ldots, c_{i\ell}$ under this mapping. Then to correct a string in $\Sigma^{n/\ell}$, interpret it as a binary string of length $n$: $\delta$ fraction of substitution errors in a codeword in $\Sigma^{n/\ell}$ yields at most a $\delta$ fraction of worst-case substitution errors over the underlying binary codeword, which can be corrected by assumption.

We now prove Proposition 2.16. [Jus72]

**Lemma A.3.** *For all positive integers $s \leq m$, there exists a linear code $C : \mathbb{F}_{2^m} \to \mathbb{F}_2^{m+s}$ of dimension $m$ and length $m + s$ tolerating $\frac{1}{2}\lfloor (m + s) \cdot H^{-1}(\frac{s}{m+s})\rfloor$ errors. Furthermore, such a code can be found in time $\tilde{O}(2^{2m})$.*

*Proof.* Since $\mathbb{F}_{2^m}$ is a $\mathbb{F}_2$ vector space, there exists a linear bijection $\sigma : \mathbb{F}_{2^m} \to \mathbb{F}_2^m$. Let $\sigma' : \mathbb{F}_{2^m} \to \mathbb{F}_2^s$ be given by taking the first $s$ bits of $\sigma(x)$. Let $e \overset{\text{def}}{=} \lfloor (m + s) \cdot H^{-1}(\frac{s}{m+s})\rfloor$.

For $\alpha \in \mathbb{F}_{2^m}$, let $C_\alpha$ be the code given by the encoding $\text{Enc}_\alpha : \mathbb{F}_2^m \to \mathbb{F}_2^{m+s}$ with

$$x \mapsto (x, \sigma'(\alpha \cdot \sigma^{-1}(x))).$$

Since multiplication by $\alpha$ is bijective and $\mathbb{F}_2$ linear, and $\sigma$ and $\sigma'$ are linear, all such codes $C_\alpha$ are linear. For any $x \in \mathbb{F}_2^m$, for a random $\alpha \in \mathbb{F}_{2^m}$, we have $\alpha\sigma^{-1}(x)$ is uniform on $\mathbb{F}_2^s$, so $\sigma'(\alpha\sigma^{-1}(x))$ is uniform on $\mathbb{F}_2^s$. Thus, each element of $\mathbb{F}_2^{m+s}$ appears exactly $2^{m-s}$ times in $\{C_\alpha : \alpha \in \mathbb{F}_q\}$ Let $X_{bad}$ denote the set of nonzero element of $\mathbb{F}_2^{m+s}$ with Hamming weight at most $e$. This set has size at most $\sum_{i=1}^{e} \binom{m+s}{i} < 2^{(m+s)H(\frac{e}{m+s})}$. Thus, for a uniformly random $\alpha \in \mathbb{F}_2^m$, the probability that there exists a nonzero element of $X_{bad}$ in $C_\alpha$

$$\frac{|X_{bad}| \cdot 2^{m-s}}{2^m} < \frac{2^{(m+s)H(\frac{e}{m+s})}}{2^s} \leq 1$$

Hence, there exists some $\alpha$ such that $C_\alpha$ has no elements in $X_{bad}$. We can find such an $\alpha$ by brute force in time $\tilde{O}(2^{2m})$: each $\alpha$ takes time $\tilde{O}(2^m)$ to compute all codewords and check their hamming weight, and there are $2^m$ such $\alpha$. In $C_\alpha$, any two codewords have Hamming distance at least $e$, so it tolerates up to $\frac{e}{2}$ errors. $\square$

*Proof or Proposition 2.16.* Let $m$ be the smallest integer larger than $\frac{12}{\varepsilon}$ such that $m \cdot 2^m \geq n$. Let $s$ be the largest integer such that $\frac{m}{m+s} \geq 1 - \frac{\varepsilon}{3}$, so that $\frac{s}{m+s} \geq \frac{\varepsilon}{4}$.

By Lemma A.3, there exists a code $C_{in} : \mathbb{F}_{2^m} \to \mathbb{F}_2^{m+s}$ of dimension $m$ and length $m + s$ with minimum distance $\lfloor (m+s)H^{-1}(\frac{s}{m+s}) \rfloor$. Let $C_{out}$ : be a Reed Solomon code over $\mathbb{F}_{2^m}$ of length $n' \stackrel{\text{def}}{=} \lfloor n/(m+s) \rfloor$ and dimension $k' = \lceil n'(1 - \frac{\varepsilon}{3}) \rceil$. Let $C \subset \{0,1\}^n$ be the concatenation of $C_{in}$ and $C_{out}$ with $n - n'm$ 0s padded on the end. The code $C_{in}$ has rate $\frac{m}{m+s} > 1 - \frac{\varepsilon}{3}$. The code $C_{out}$ has rate at least $1 - \frac{\varepsilon}{3}$. The padding of 0s multiplies the rate by $\frac{n'(m+s)}{n} \geq 1 - \frac{m+s}{n} > 1 - \frac{\varepsilon}{3}$. Thus, the total rate is at least $(1 - \frac{\varepsilon}{3})^3 > 1 - \varepsilon$.

To decode a received word $c_1, \ldots, c_n$, we first run the inner decoding to obtain symbols $\alpha_i \in \mathbb{F}_{2^m}$ for $i = 1, \ldots, n'$, where $\alpha_i$ is the decoding of $c_{(i-1)(m+s)+1}, \ldots, c_{i(m+s)}$ under $C_{in}$. Then, we run the outer decoding on $\alpha_1, \ldots, \alpha_{n'}$ to obtain the message. The inner decoding can be computed by brute force in time $O(m2^m) < O_\varepsilon(n)$. The outer decoding can be computed in time $O(n^2)$ using the Berlekamp-Massey algorithm. Thus, the total decoding run time is $O_\varepsilon(n^2)$. The encoding takes time $O_\varepsilon(n^2)$, and construction takes time $\tilde{O}(2^{2m}) = O_\varepsilon(n^2)$ because we need to construct the inner code.

The outer code tolerates $n' - k' > \frac{\varepsilon n'}{4}$ errors. The inner code tolerates up to $\frac{1}{2}\lfloor (m+s)H^{-1}(\frac{s}{m+s}) \rfloor$ errors, and thus every incorrect $\alpha_i$ accounts for at least $\frac{1}{2}(m+s)H^{-1}(\frac{s}{m+s}) > \frac{1}{2}(m+s)H^{-1}(\frac{\varepsilon}{4})$ errors. Thus, for the outer decoding to fail, we need at least $\frac{\varepsilon n'}{4} \cdot \frac{1}{2}(m+s) \cdot H^{-1}(\frac{\varepsilon}{4}) > \frac{\varepsilon^2 n}{500 \log \frac{1}{\varepsilon}}$ errors. Here, we used that $(m+s)n' > 0.9n$ and that $H^{-1}(x) > \frac{x}{2 \log(6/x)}$, so $H^{-1}(\frac{\varepsilon}{4}) > \frac{\varepsilon}{48 \log(1/\varepsilon)}$. $\qquad\square$