



On Codes for the Noisy Substring Channel

Yonatan Yehezkeally , *Member, IEEE* and Nikita Polyanski , *Member, IEEE*

Abstract—We consider the problem of coding for the substring channel, in which information strings are observed only through their (multisets of) substrings. Due to existing DNA sequencing techniques and applications in DNA-based storage systems, interest in this channel has renewed in recent years. In contrast to existing literature, we consider a noisy channel model where information is subject to noise *before* its substrings are sampled, motivated by in-vivo storage.

We study two separate noise models, substitutions or deletions. In both cases, we examine families of codes which may be utilized for error-correction and present combinatorial bounds on their sizes. Through a generalization of the concept of repeat-free strings, we show that the added required redundancy due to this imperfect observation assumption is sublinear, either when the fraction of errors in the observed substring length is sufficiently small, or when that length is sufficiently long. This suggests that no asymptotic cost in rate is incurred by this channel model in these cases. Moreover, we develop an efficient encoder for such constrained strings in some cases.

Finally, we show how a similar encoder can be used to avoid formation of secondary-structures in coded DNA strands, even when accounting for *imperfect* structures.

Index Terms—DNA storage, Sequence reconstruction, Error-correcting codes, Insertion/deletion-correcting codes, Constrained codes

I. INTRODUCTION

DNA as a medium for data storage offers high density and longevity, far greater than those of electronic media [1]. Among its applications, data storage in DNA may offer a protected medium for long-period data storage [2], [3]. In particular, it has recently been demonstrated that storage in the DNA of living organisms (henceforth, *in-vivo* DNA storage) is now feasible [4]; the envelope of a living cell affords some level of protection to the data, and even offers propagation, through cell replication. Among its varied usages, in-vivo DNA storage allows watermarking genetically modified organisms (GMOs) [5]–[7] to protect intellectual property, or labeling research material [3], [8]. It may even conceal sensitive information, as it may appear indistinguishable from the organism’s own genetic information [9].

Manuscript received 25 September 2023. This work has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (Grant agreement No. 801434). The work of Yonatan Yehezkeally was supported the Alexander von Humboldt Foundation under a Carl Friedrich von Siemens Post-Doctoral Research Fellowship. The work of Nikita Polyanski was supported by the German Research Foundation (Deutsche Forschungsgemeinschaft, DFG) under Grant No. WA3907/1-1. An earlier version of this paper was presented in part at the 2021 IEEE International Symposium on Information Theory (ISIT) [DOI: 10.1109/ISIT45174.2021.9517943]. (*Corresponding author: Yonatan Yehezkeally.*)

Yonatan Yehezkeally is with the Institute for Communications Engineering, School of Computation, Information and Technology, Technical University of Munich, 80333 Munich, Germany (e-mail: yonatan.yehezkeally@tum.de). Nikita Polyanski is with IOTA Foundation, Berlin, Germany.

Similarly to other media, information stored over this medium is subjected to noise due to mutations, creating errors in data, which accumulate over time and replication cycles. Examples of such noise include symbol insertions or deletion, in addition to substitutions (point-mutations) [10], [11]; the latter is the focus of the vast majority of classical error-correction research, and the former have also been studied. Interestingly, however, the very methods we currently use to store and later retrieve data from DNA inherently introduce new constraints on information reconstruction. While desired sequences may be synthesized (albeit, while suffering from errors, e.g., substitution noise), the process of DNA sequencing, i.e., retrieving the DNA sequence of an organism, only observes that sequence as the (likely incomplete) multiset of its substrings (practically, up to a certain substring length) [12]. Thus, information contained in the order of these substrings might be irrevocably lost. As a result of these constraints, conventional and well-developed error-correction approaches cannot simply be applied.

To overcome these effects, one approach in existing literature is to add redundancy in the form of indexing, in order to recover the order of substrings (see, e.g., [13]–[15]). A different approach, potentially more applicable to in-vivo DNA storage, is to add redundancy in the form of constraints on the long information string, such that it can be uniquely reconstructed by knowledge of its substrings of a given length (or range of lengths). The combinatorial problem of recovering a sequence from its substrings has attracted attention in recent years [16]–[23], and coding schemes involving only these substrings (including the incidence frequency of each substring) were studied [12], [15], [24]–[26].

However, works dedicated to overcoming this obstacle, inherent to the technology we use, have predominantly focused on storage outside of living cells (i.e., *in-vitro* DNA storage). Likewise, works focused on error-correction for in-vivo DNA data storage (e.g., [27]–[29]) have disregarded the technical process by which data is to be read. However, in real applications varied distinct noise mechanisms act on stored data concurrently. Hence, in practice, both sets of challenges have to be collectively overcome in order to robustly store information using in-vivo DNA.

The aim of this work is to protect against errors in the information string (caused by mutations over the replication process of cells), when channel outputs are given by the multisets of their substrings, of a predetermined length, rather than entire strings. This models the process of DNA sequencing, once information needs to be read from the medium. We shall study the required redundancy of this model, and devise coding strategies, under the assumption of two different error types: substitution and deletion noise.

Another application for this line of research is secondary-

structure avoidance. Secondary structures are complex spatial structures that can form in a chemically active single-stranded DNA, as a result of the strand folding upon itself to allow two sub-segments to bond via complementary-base-pair hybridization [30]. Their formation renders the DNA strand chemically inactive and is therefore detrimental for sequencing and DNA-based computation, hence a number of recent works have looked to avoid them through coding [31]–[34]. Herein we focus on relatively long structures, but unlike recent works, we do not consider only perfect structures, but also attempt to avoid ones which contain impairments, i.e., imperfect structures. We show that this problem is closely connected to the above-described channel; thus, we are able to also present an efficient encoder for this setting.

The paper is organized as follows. In Section II, we discuss the main contribution of this paper, in context of related works. In Section III we then present necessary notation. Then, in Section IV we study the suggested model with substitution errors, and in Section V with deletion errors. Finally, in Section VI we develop an encoder for avoiding the formation of even imperfect secondary structures.

II. RELATED WORKS AND MAIN CONTRIBUTION

Given a string of length n , the problem of reconstructing \mathbf{x} from the multiset of (all-, or, in some works, most-) its substrings of a fixed length $\ell \leq n$, has been studied in literature. Assuming no errors occur in \mathbf{x} prior to sampling of its substrings, the problem of interest is identifying a set of constraints on the information string, equivalent of sufficient, for such reconstruction to be achievable.

It was observed in [16] that under certain circumstances, distinct information strings in which repetitions of ℓ -substrings appear in different positions, exhibit the same multisets of $(\ell + 1)$ -substrings. These observations indicate that care must be taken when including code-words which contain repeating ℓ -substrings (indeed, where observations are made via the multiset of ℓ' -substrings, for some $\ell' \leq \ell + 1$). On the other hand, if every ℓ -substring of \mathbf{x} is unique, then \mathbf{x} is uniquely reconstructible from the multiset of its $(\ell + 1)$ -substrings (and in fact, ℓ' -substrings, for all $\ell' > \ell$), as evident from a greedy reconstruction algorithm (which at each stage searches for the next/previous character in the information string). This observation motivates the study of *repeat-free strings*; \mathbf{x} is said to be ℓ -repeat-free if every ℓ -substring of \mathbf{x} is unique (put differently, if \mathbf{x} is of length n , then it contains $n - \ell + 1$ distinct ℓ -substrings).

Focus on repeat-free strings is further justified by the following results. It was observed in [19], via introduction of *profile vectors*, that over an alphabet of size q , where the length of strings n grows, if $\ell < \frac{\log_q(n)}{1+\epsilon}$ then the rate of all existing ℓ -substring multisets vanishes. Conversely, it was demonstrated in [22] using probabilistic arguments that the asymptotic redundancy of the code-book consisting of all ℓ -repeat-free strings of length n (which, as noted above, is an upper bound for the redundancy of a code assuring reconstruction from $(\ell + 1)$ -substrings), is $O(n^{2-\ell/\log_q(n)})$; thus, when $\ell > (1 + \epsilon)\log_q(n)$, the rate of repeat-free strings alone is 1.

In this paper, we extend the setting of previous works by allowing information strings to suffer a bounded number of errors, prior to the sampling of their substrings. We study this model under two separate error models: substitution (Hamming) errors, and deletion errors. In both cases we show (see Theorems 10 and 18) that when $\ell > (1 + \epsilon)\log_q(n)$ and the fraction of errors in the substring length ℓ is sufficiently small, the rate of generalized repeat-free strings dubbed *resilient-repeat-free* suffers no penalty from the process of sampling, or from the presence of noise (when compared to the results of [22]); i.e., the required added redundancy is sub-linear. In the case of Hamming noise, we also show that when the fraction of errors is too large, resilient-repeat-free strings do not exist. However, it is left for future works to determine the precise transition between the two regimes. Further, we develop an efficient encoder for resilient-repeat-free sequences (see Theorem 14), although our encoder does not output sequences of a fixed length n , but rather only guarantees that the output is of length *at most* n .

It should be noted that [20] presented almost explicit en-/decoding algorithms for codes with a similar noise model. However, in that paper's setting, substitution noise affects individual substrings *after* sampling; the codes it constructs are capable of correcting a constant number of errors in each substring, but requires the assumption that errors do not affect the same information symbol in a majority of the substrings that reflect it. Therefore, its setting is incompatible with the one considered herein, whereby each error occurring *before* sampling affects ℓ consecutive substrings. [23] also developed codes with full rate, capable of correcting a fixed number of errors, occurring in substrings independently after sampling. It replaced the aforementioned restriction by a constraint on the number of total erroneous substrings, which is at most logarithmic in the information string's length. Hence, the total number of errors in its setting remains asymptotically smaller than the one incurred in the setting considered here.

Finally, as mentioned above we exploit the similarity between the aforementioned setting and channel model and the problem of avoiding secondary structures. We focus on hairpin-loop structures with long stems (scaling logarithmically in the length of the sequence), and unlike recent works [31]–[34] the encoder we develop prevents the formation of such structures even when the underlying complementary-base-pair hybridization (in a region called the *stem* of the structure) is imperfect, that is, it contains at most a δ -fraction of mismatched nucleobases (which cannot stably hybridize), while asymptotically achieving full rate.

III. PRELIMINARIES

Let Σ^* be the set of finite strings over an alphabet Σ , which for convenience we assume to be a finite unital ring of size q (e.g., \mathbb{Z}_q). For $\mathbf{x} = x(0)x(1) \cdots x(n-1) \in \Sigma^*$, we let $|\mathbf{x}| = n$ denote the *length* of \mathbf{x} . We note that indices in the sequel are numbered $0, 1, \dots$. For $\mathbf{x}, \mathbf{y} \in \Sigma^*$, we let \mathbf{xy} be their concatenation. For $I \subseteq \mathbb{N}$ and $\mathbf{x} \in \Sigma^*$, we denote by \mathbf{x}_I the restriction of \mathbf{x} to indices in I (excluding any indices $|x| \leq i \in I$), ordered according to the naturally inherited order on I .

We let $|A|$ denote the size of a finite set A . For a code $C \subseteq \Sigma^n$, we define its *redundancy* $\text{red}(C) \triangleq n - \log_q |C|$, and *rate* $R(C) \triangleq \frac{1}{n} \log_q |C| = 1 - \frac{\text{red}(C)}{n}$.

For $n \in \mathbb{N}$, denote $[n] \triangleq \{0, 1, \dots, n-1\}$. Although perhaps confusable, for $m \leq n \in \mathbb{N}$ we use the common notation $[m, n] \triangleq \{m, m+1, \dots, n\}$. We shall interpret x_I as enumerated by $[[I]]$, i.e., $x_I(0) = x(\min I)$, etc. Where it is convenient, we will also assume $I \subseteq \mathbb{N}$ to be enumerated by $[[I]]$, such that the order of elements is preserved; i.e., $I = \{I(i) : i \in [[I]]\}$, and for all $i \in [[I] - 1]$ one has $I(i) < I(i+1)$. Under this convention, e.g., $x_I(0) = x(I(0))$. We follow the standard group notation in denoting for $j \in \mathbb{N}$ and $I \subseteq \mathbb{N}$ the *coset* $j + I \triangleq \{j + i : i \in I\}$.

For $\mathbf{x} \in \Sigma^*$ and $i, \ell \in \mathbb{N}$, where $i + \ell \leq |\mathbf{x}|$, we say that $\mathbf{x}_{i+[\ell]}$ is the length ℓ *substring* of \mathbf{x} at index i , or ℓ -*mer* (at index i) for short. Using notation from [16], for $\mathbf{x} \in \Sigma^*$ and $\ell \in \mathbb{N}$ we denote the multiset of ℓ -mers of \mathbf{x} by

$$Z_\ell(\mathbf{x}) \triangleq \{\{\mathbf{x}_{i+[\ell]} : 0 \leq i \leq |\mathbf{x}| - \ell\}\}.$$

We follow [22] in denoting the set of ℓ -*repeat-free* strings

$$\mathcal{RF}_\ell(n) \triangleq \{\mathbf{x} \in \Sigma^n : i < j \implies \mathbf{x}_{i+[\ell]} \neq \mathbf{x}_{j+[\ell]}\}.$$

Assuming an underlying error model, known in context but yet to be determined, we let $B_t(\mathbf{x})$, for some $\mathbf{x} \in \Sigma^*$, be the set of strings $\mathbf{y} \in \Sigma^*$ which may be the product of at most t errors occurring to \mathbf{x} . Using this notation, our aim shall be to study and design codes $C \subseteq \Sigma^n$, such that given $\mathbf{x} \in C$ and $\mathbf{y} \in B_t(\mathbf{x})$, for some fixed (or bounded) t , \mathbf{x} can be uniquely reconstructed given only $Z_\ell(\mathbf{y})$. We shall study constraints which allow unique reconstruction of \mathbf{y} , and state in Corollary 16 specific cases where this in turn allows reconstruction of \mathbf{x} .

IV. SUBSTITUTION NOISE

In this section we consider substitution noise, with error balls $B_t^s(\mathbf{x}) \triangleq \{\mathbf{y} : d_H(\mathbf{y}, \mathbf{x}) \leq t\}$, where $d_H(\mathbf{x}, \mathbf{y})$ denotes the Hamming distance between \mathbf{x} and \mathbf{y} .

We present and study a family of repeat-free strings which are resilient to substitution errors:

Definition 1 We say that $\mathbf{x} \in \Sigma^*$ is (t, ℓ) -resilient repeat free if the result of any t substitution errors to \mathbf{x} is ℓ -repeat-free. More precisely, we define

$$\mathcal{RRF}_{t,\ell}^s(n) \triangleq \{\mathbf{x} \in \Sigma^n : B_t^s(\mathbf{x}) \subseteq \mathcal{RF}_\ell(n)\}.$$

A. Rate of resilient-repeat-free strings

In the following section we dedicate ourselves to study $\text{red}(\mathcal{RRF}_{t,\ell}^s(n))$, where t, ℓ are taken to be functions of n . In particular, we will be interested in developing sufficient (and to a lesser degree, necessary) conditions on t, ℓ that assure $R(\mathcal{RRF}_{t,\ell}^s(n)) = 1 - o_n(1)$.

Recall that [22] showed that if $\ell = a \log(n) + o(\log(n))$, then

$$R(\mathcal{RF}_\ell(n)) = \begin{cases} o_n(1), & a < 1; \\ 1 - o_n(1), & a > 1. \end{cases}$$

Since $\mathcal{RRF}_{t,\ell}^s(n) \subseteq \mathcal{RRF}_{0,\ell}^s(n) = \mathcal{RF}_\ell(n)$, then with the above scaling of ℓ , $a < 1$ implies that $R(\mathcal{RRF}_{t,\ell}^s(n)) = o_n(1)$ as well; we shall see that when $a > 1$, then for sufficiently small t we still have $R(\mathcal{RRF}_{t,\ell}^s(n)) = 1 - o_n(1)$.

A particular notion that will aid in our analysis is the following: for $0 < k \leq \ell$, denote

$$\mathcal{A}_t^\ell(k) \triangleq \left\{ \mathbf{x} \in \Sigma^{\ell+k} : \exists \mathbf{y} \in B_t^s(\mathbf{x}) : \mathbf{y}_{[\ell]} = \mathbf{y}_{k+[\ell]} \right\}.$$

We let $\pi_t^\ell(k) \triangleq q^{-(\ell+k)} |\mathcal{A}_t^\ell(k)|$ (i.e., $\pi_t^\ell(k) = \Pr(\mathbf{x} \in \mathcal{A}_t^\ell(k))$ where $\mathbf{x} \in \Sigma^{\ell+k}$ is chosen uniformly at random). For convenience, when ℓ, t are known from context, we also abbreviate:

$$\pi \triangleq \pi_t^\ell(\ell); \quad \pi' \triangleq \max_{0 < k < \ell} \pi_t^\ell(k). \quad (1)$$

The usefulness of the notation in (1) is substantiated in the following theorem.

Theorem 2 Let $\ell = \ell(n), t = t(n)$ be integer functions, and assume $t \leq \ell \leq n$. If for all sufficiently large n it holds that $\ell^2 \pi' + \ell n \pi \leq 1/2e$, then

$$\text{red}(\mathcal{RRF}_{t,\ell}^s(n)) = O(n \log(n) \pi' + n^2 \pi).$$

Before proving Theorem 2, we present the following result.

Definition 3 For positive $\ell \leq n$, denote $\binom{[n]}{\ell} \subseteq 2^{[n]}$ the collection of ℓ -subsets of $[n]$. A pair of subsets $(I, J) \in \binom{[n]}{\ell}^2$ is said to be *observable* if $I(k) < J(k)$ for all $k \in [\ell]$.

Given a string $\mathbf{x} \in \Sigma^n$, known from context, we will denote for an observable pair $(I, J) \in \binom{[n]}{\ell}^2$

$$\mathbf{u}_{I,J} \triangleq x_I - x_J \in \Sigma^\ell.$$

We also denote $L_I \triangleq \{(P, Q) : (P, Q) \text{ is observable, } (P \cup Q) \cap I = \emptyset\}$. To simplify notation, where some $\ell \leq n$ is also given, we shall abbreviate $\mathbf{u}_{i,j} \triangleq \mathbf{u}_{i+[\ell], j+[\ell]}$ and $L_i \triangleq L_{i+[\ell]}$, for any $0 \leq i < j \leq n - \ell$.

Lemma 4 Take $\ell \leq n$ and an observable pair $(I, J) \in \binom{[n]}{\ell}^2$. Further, let $\mathbf{x} \in \Sigma^n$ be chosen uniformly at random. Then $\mathbf{u}_{I,J}$ is distributed uniformly and mutually independent of $\{\mathbf{u}_{P,Q} : (P, Q) \in L_I\}$.

Proof: First, since $\mathbf{u}_{I,J}$ is the image of \mathbf{x} under a linear map (more precisely, a module homomorphism), the pre-image of any point is a coset of the map's kernel and, thus, of equal size; as a result, $\mathbf{u}_{I,J}$ is distributed uniformly on the map's range. Since (I, J) is observable, the map is surjective onto Σ^ℓ , hence the first part is completed.

Second, observe that \mathbf{x}_I is independent of $\mathbf{x}_{[n] \setminus I}$, hence mutually independent of $\{\mathbf{u}_{P,Q} : (P, Q) \in L_I\}$. Since given $\mathbf{x}_{[n] \setminus I}$, there exist a bijection between \mathbf{x}_I and $\mathbf{u}_{I,J}$, the proof is concluded. ■

With Lemma 4, we can now prove Theorem 2. Our proof strategy is based on Lovász's local lemma (LLL), which we slightly rephrase below.

Theorem 5 [35, Th. 1.1] Let $\{A_{i,j}\}_{i,j}$ be events in a probability space Ω . If for all i, j there exist constants $0 < f_{i,j} < 1$ such that

$$\Pr(A_{i,j}) \leq f_{i,j} \prod_{(p,q) \notin \Gamma} (1 - f_{p,q}),$$

where Γ is such that the event $A_{i,j}$ is mutually independent of events $\{A_{p,q} : (p,q) \notin \Gamma\}$, then

$$\Pr\left(\Omega \setminus \bigcup_{i,j} A_{i,j}\right) \geq \prod_{i,j} (1 - f_{i,j}).$$

To the best of authors' knowledge, this application of the lemma is novel to the conference version of this work; it then also appeared in similar form in a concurrent journal version of [22].

Proof of Theorem 2: We define for all $0 \leq i < j \leq n - \ell$ the sets

$$A_{i,j} \triangleq \left\{ \mathbf{x} \in \Sigma^n : \exists \mathbf{y} \in B_t^s(\mathbf{x}) : \mathbf{y}_{i+[\ell]} = \mathbf{y}_{j+[\ell]} \right\}.$$

Note that $\Sigma^n \setminus \mathcal{RRF}_{t,\ell}^s(n) = \bigcup_{i,j} A_{i,j}$.

We let $\mathbf{x} \in \Sigma^n$ be chosen uniformly at random. Then $\Pr(\mathbf{x} \in A_{i,j}) = \pi_t^\ell(\min\{\ell, j - i\})$. Further,

$$|\mathcal{RRF}_{t,\ell}^s(n)| = q^n \cdot \Pr(\mathbf{x} \in \mathcal{RRF}_{t,\ell}^s(n)),$$

and hence

$$\text{red}(\mathcal{RRF}_{t,\ell}^s(n)) = -\log_q \Pr(\mathbf{x} \in \mathcal{RRF}_{t,\ell}^s(n)).$$

As mentioned above, we shall rely on Theorem 5. Note that, in our notation, $\Pr(\mathbf{x} \notin \bigcup_{i,j} A_{i,j}) = \Pr(\mathbf{x} \in \mathcal{RRF}_{t,\ell}^s(n))$.

To determine Γ , we claim for $0 \leq i < j \leq n - \ell$ that the event $\{\mathbf{x} \in A_{i,j}\}$ is mutually independent of the events $\{\{\mathbf{x} \in A_{p,q}\} : |i - p|, |i - q| \geq \ell\}$. Indeed, Lemma 4 then implies that $\mathbf{u}_{i,j}$ is mutually independent of $\{\mathbf{u}_{p,q} : (p,q) \in L_i\}$. It is left to the reader to verify that there exists a set $B_{i,j} \subseteq \Sigma^\ell$, which depends on i, j but not \mathbf{x} , such that the event $\{\mathbf{x} \in A_{i,j}\}$ can be restated as $\{\mathbf{u}_{i,j} \in B_{i,j}\}$ (and similarly for p, q); therefore, the claim holds. Hence, for any $0 \leq i < j \leq n - \ell$ we let $\Gamma \triangleq \{(p,q) : (p,q) \notin L_i\}$. Observe that the number of such pairs (p,q) satisfying $|p - q| < \ell$ is at most $(2\ell - 1)^2 - \binom{2\ell - 1}{2} = (2\ell - 1)\ell < 2\ell^2$, and the number of other pairs is at most $(2\ell - 1)(n - \ell + 1) < 2\ell n$.

Recalling for $0 < f < 1$ that $1 - f \geq e^{-f/(1-f)}$, and denoting $y = f/(1 - f)$ we observe

$$\begin{aligned} f_{i,j} \prod_{p,q} (1 - f_{p,q}) &= \frac{f_{i,j}}{1 - f_{i,j}} \prod_{p,q \text{ inc. } i,j} (1 - f_{p,q}) \\ &\geq y_{i,j} \exp\left(-\sum_{p,q \text{ inc. } i,j} y_{p,q}\right) \end{aligned}$$

We shall apply an almost symmetric version of LLL, where

$$f_{i,j} = \begin{cases} f, & j - i \geq \ell, \\ f', & j - i < \ell. \end{cases}$$

Then, to satisfy the LLL conditions it suffices that $\pi \leq ye^{-2\ell^2 y' - 2\ell n y}$ and $\pi' \leq y'e^{-2\ell^2 y' - 2\ell n y}$, where y, y' are

defined as above for f, f' respectively. Let $y' \triangleq e\pi'$, $y \triangleq e\pi$. We have,

$$\begin{aligned} ye^{-2\ell^2 y' - 2\ell n y} &= \pi e^{1 - 2e(\ell^2 \pi' + \ell n \pi)} \geq \pi; \\ y'e^{-2\ell^2 y' - 2\ell n y} &= \pi' e^{1 - 2e(\ell^2 \pi' + \ell n \pi)} \geq \pi', \end{aligned}$$

as required.

Finally,

$$\begin{aligned} \text{red}(\mathcal{RRF}_{t,\ell}^s(n)) &= -\log_q \Pr(\mathbf{x} \in \mathcal{RRF}_{t,\ell}^s(n)) \\ &= -\log_q \Pr\left(\mathbf{x} \notin \bigcup_{i,j} A_{i,j}\right) \\ &\leq -\log_q \prod_{i,j} (1 - f_{i,j}) \\ &= \sum_{i,j} \log_q (1 + y_{i,j}) \\ &\leq \frac{1}{\ln(q)} \sum_{i,j} y_{i,j} \\ &\leq \frac{n\ell}{\ln(q)} y' + \frac{n^2}{\ln(q)} y, \end{aligned}$$

which concludes the proof. \blacksquare

Based on the last theorem, it is of interest to bound π, π' from above. To that end, we note the following result.

Lemma 6 Take $t \leq \ell \leq n \in \mathbb{N}$, $\mathbf{x} \in \Sigma^n$. If for all $0 \leq i < j \leq n - \ell$ it holds that

$$d_H(\mathbf{x}_{i+[\ell]}, \mathbf{x}_{j+[\ell]}) > t + \max\{0, \min\{t, \ell - j + i\}\},$$

then $\mathbf{x} \in \mathcal{RRF}_{t,\ell}^s(n)$.

Proof: The proof follows from applying the triangle inequality by cases on $(i + [\ell]) \cap (j + [\ell])$. Assume to the contrary that there exist $\mathbf{y} \in B_t^s(\mathbf{x})$ and $0 \leq i < j \leq n - \ell$ such that $\mathbf{y}_{i+[\ell]} = \mathbf{y}_{j+[\ell]}$. Note that

$$\begin{aligned} d_H(\mathbf{x}_{i+[\ell]}, \mathbf{x}_{j+[\ell]}) &\leq d_H(\mathbf{x}_{i+[\ell]}, \mathbf{y}_{i+[\ell]}) + \\ &\quad d_H(\mathbf{y}_{i+[\ell]}, \mathbf{y}_{j+[\ell]}) + \\ &\quad d_H(\mathbf{y}_{j+[\ell]}, \mathbf{x}_{j+[\ell]}) \\ &= d_H(\mathbf{x}_{i+[\ell]}, \mathbf{y}_{i+[\ell]}) + \\ &\quad d_H(\mathbf{y}_{j+[\ell]}, \mathbf{x}_{j+[\ell]}). \end{aligned}$$

We continue by cases.

If $j - i \geq \ell$ then, since $(i + [\ell]) \cap (j + [\ell]) = \emptyset$, then $d_H(\mathbf{x}_{i+[\ell]}, \mathbf{y}_{i+[\ell]}) + d_H(\mathbf{y}_{j+[\ell]}, \mathbf{x}_{j+[\ell]}) \leq t$, which contradicts the theorem's assumption.

On the other hand, if $j - i \leq \ell - t$, then we may simply bound $d_H(\mathbf{x}_{i+[\ell]}, \mathbf{y}_{i+[\ell]}) + d_H(\mathbf{y}_{j+[\ell]}, \mathbf{x}_{j+[\ell]}) \leq 2t$, again in contradiction.

Finally, suppose $\ell - t < j - i < \ell$. Note that

$$\begin{aligned} &d_H(\mathbf{x}_{i+[\ell]}, \mathbf{y}_{i+[\ell]}) + d_H(\mathbf{y}_{j+[\ell]}, \mathbf{x}_{j+[\ell]}) \\ &= d_H(\mathbf{x}_{[i,j-1]}, \mathbf{y}_{[i,j-1]}) + 2d_H(\mathbf{x}_{[j,i+\ell-1]}, \mathbf{y}_{[j,i+\ell-1]}) + \\ &\quad d_H(\mathbf{y}_{[i+\ell,j+\ell-1]}, \mathbf{x}_{[i+\ell,j+\ell-1]}). \end{aligned}$$

Since $[i, j-1], [j, i+\ell-1], [i+\ell, j+\ell-1]$ are pairwise disjoint,

$$d_H(\mathbf{x}_{[i,j-1]}, \mathbf{y}_{[i,j-1]}) + d_H(\mathbf{x}_{[j,i+\ell-1]}, \mathbf{y}_{[j,i+\ell-1]}) + d_H(\mathbf{y}_{[i+\ell,j+\ell-1]}, \mathbf{x}_{[i+\ell,j+\ell-1]}) \leq t.$$

Hence, denoting $\Delta \triangleq d_H(\mathbf{x}_{[j,i+\ell-1]}, \mathbf{y}_{[j,i+\ell-1]})$, we have

$$d_H(\mathbf{x}_{i+[\ell]}, \mathbf{x}_{j+[\ell]}) \leq (t - \Delta) + 2\Delta = t + \Delta \leq t + (\ell - j + i),$$

once more in contradiction. This concludes the proof. \blacksquare

Corollary 7 $\pi_t^\ell(k) \leq q^{-\ell} \sum_{i=0}^{t+\min\{t,\ell-k\}} \binom{\ell}{i} (q-1)^i$.

Proof: By Lemma 6 we have $\pi_t^\ell(k) \leq \Pr(\text{wt}(\mathbf{u}_{0,k}) \leq t + \min\{t, \ell - k\})$. Since by Lemma 4 $\mathbf{u}_{0,k} \in \Sigma^\ell$ is distributed uniformly, the proof is concluded. \blacksquare

The bound of Corollary 7 can be improved upon in some cases, depending on k (thus improving the upper bound on π'):

Lemma 8

$$\pi_t^\ell(k) \geq q^{-\ell} \sum_{i=0}^t \binom{\ell}{i} (q-1)^i,$$

and

$$\pi_t^\ell(k) \leq \begin{cases} q^{-\ell} \sum_{i=0}^t \binom{\ell+k}{i} (q-1)^i, & k \leq \frac{\ell}{2}; \\ q^{-\ell} \sum_{i=0}^t \binom{2\ell-k}{i} (q-1)^i & k > \frac{\ell}{2}. \end{cases}$$

Proof: Take integers $p \geq 2$ and $0 \leq r < k$ such that $\ell + k = pk + r$. Observe that $\mathbf{y}_{[\ell]} = \mathbf{y}_{k+[\ell]}$ implies that \mathbf{y} can be determined by its first k coordinates, i.e.

$$\mathbf{y} = \underbrace{(\mathbf{y}_{[k]}, \dots, \mathbf{y}_{[k]})}_{p \text{ times}}, \mathbf{y}_{[r]}.$$

Observe for each $\mathbf{y} \in \Sigma^{\ell+k}$, satisfying $\mathbf{y}_{[\ell]} = \mathbf{y}_{k+[\ell]}$, that one may form a unique $\mathbf{x} \in \mathcal{A}_t^\ell(k)$ by changing at most t of the symbols $\mathbf{y}_{k+[\ell]}$. It follows that

$$|\mathcal{A}_t^\ell(k)| \geq q^k \sum_{i=0}^t \binom{\ell}{i} (q-1)^i.$$

On the other hand, it is also straightforward that

$$|\mathcal{A}_t^\ell(k)| \leq q^k \sum_{i=0}^t \binom{\ell+k}{i} (q-1)^i.$$

When $p = 2$ or, equivalently, $k > \frac{\ell}{2}$, we shall improve the above bound. Take $\mathbf{x} \in \mathcal{A}_k$ and $\mathbf{y} \in B_t^s(\mathbf{x})$ satisfying $\mathbf{y}_{[\ell]} = \mathbf{y}_{k+[\ell]}$. Define the intervals $I_1 \triangleq [\ell - k]$, $I_2 \triangleq [\ell - k, k - 1]$, $I_3 \triangleq [k, \ell - 1]$, $I_4 \triangleq [\ell, 2k - 1]$, $I_5 \triangleq [2k, k + \ell - 1]$. Using this notation, we have $\mathbf{y}_{I_2} = \mathbf{y}_{I_4}$ and $\mathbf{y}_{I_1} = \mathbf{y}_{I_3} = \mathbf{y}_{I_5}$, i.e.,

$$\mathbf{y} = \mathbf{y}_{I_1} \mathbf{y}_{I_2} \mathbf{y}_{I_1} \mathbf{y}_{I_2} \mathbf{y}_{I_1}.$$

Consider the string $\mathbf{y}' \triangleq \mathbf{y}_{I_1} \mathbf{x}_{I_2} \mathbf{y}_{I_1} \mathbf{x}_{I_2} \mathbf{y}_{I_1}$ and note that

$$\begin{aligned} d_H(\mathbf{x}, \mathbf{y}') &= d_H(\mathbf{x}_{I_1 \cup I_3 \cup I_5}, \mathbf{y}'_{I_1 \cup I_3 \cup I_5}) + \\ &\quad d_H(\mathbf{x}_{I_2}, \mathbf{y}'_{I_2}) + d_H(\mathbf{x}_{I_4}, \mathbf{y}'_{I_2}) \\ &= d_H(\mathbf{x}_{I_1 \cup I_3 \cup I_5}, \mathbf{y}_{I_1 \cup I_3 \cup I_5}) + \\ &\quad d_H(\mathbf{x}_{I_2}, \mathbf{x}_{I_2}) + d_H(\mathbf{x}_{I_4}, \mathbf{x}_{I_2}) \\ &= d_H(\mathbf{x}_{I_1 \cup I_3 \cup I_5}, \mathbf{y}_{I_1 \cup I_3 \cup I_5}) + \\ &\quad 0 + d_H(\mathbf{x}_{I_4}, \mathbf{x}_{I_2}). \end{aligned}$$

Applying the triangle inequality on the last addend,

$$\begin{aligned} d_H(\mathbf{x}, \mathbf{y}') &\leq d_H(\mathbf{x}_{I_1 \cup I_3 \cup I_5}, \mathbf{y}_{I_1 \cup I_3 \cup I_5}) + \\ &\quad d_H(\mathbf{x}_{I_4}, \mathbf{y}_{I_4}) + d_H(\mathbf{y}_{I_4}, \mathbf{x}_{I_2}) \\ &= d_H(\mathbf{x}_{I_1 \cup I_3 \cup I_5}, \mathbf{y}_{I_1 \cup I_3 \cup I_5}) + \\ &\quad d_H(\mathbf{x}_{I_4}, \mathbf{y}_{I_4}) + d_H(\mathbf{y}_{I_2}, \mathbf{x}_{I_2}) \\ &= d_H(\mathbf{x}, \mathbf{y}) \leq t. \end{aligned}$$

Therefore, for any $\mathbf{x} \in \mathcal{A}_t^\ell(k)$ there exists $\mathbf{y}' \in B_t^s(\mathbf{x})$ of the form

$$\mathbf{y}' = \mathbf{y}_{I_1} \mathbf{x}_{I_2} \mathbf{y}_{I_1} \mathbf{x}_{I_2} \mathbf{y}_{I_1},$$

and in particular $\mathbf{x}_{I_2} = \mathbf{y}'_{I_2}$. This implies an improved upper bound, by freely choosing $\mathbf{y}_{I_1} \mathbf{x}_{I_2} \in \Sigma^k$ and subsequently at most t coordinates from $[\ell + k] \setminus I_2$ to change, as follows

$$|\mathcal{A}_t^\ell(k)| \leq q^k \sum_{i=0}^t \binom{2\ell - k}{i} (q-1)^i. \quad \blacksquare$$

Corollary 9

$$\pi = q^{-\ell} \sum_{i=0}^t \binom{\ell}{i} (q-1)^i \leq q^{-\ell(1 - H_q(\min\{\frac{q-1}{q}, \frac{t}{\ell}\}))},$$

and

$$\begin{aligned} \pi' &\leq q^{-\ell} \sum_{i=0}^t \binom{\lceil 3\ell/2 \rceil}{i} (q-1)^i \\ &\leq q^{-\ell(1 - \frac{3}{2}H_q(\min\{\frac{q-1}{q}, \frac{2t}{3\ell}\}))}, \end{aligned}$$

where H_q is the q -ary entropy function.

Proof: First observe the equality on the first line follows from the lower bound of Lemma 8, together with the upper bound of either Corollary 7 or Lemma 8. Similarly, the first inequality on the second line follows from the upper bound of Lemma 8.

The second inequality on both lines follows from the standard bound on the size of the q -ary Hamming ball (see, e.g., [36, Lem. 4.7]); in particular observe that

$$\sum_{i=0}^t \binom{\lceil 3\ell/2 \rceil}{i} (q-1)^i \leq q^{\lceil 3\ell/2 \rceil H_q(\min\{\frac{q-1}{q}, \frac{t}{\lceil 3\ell/2 \rceil}\})},$$

and since $x \mapsto xH_q(1/x)$ is increasing for $x \geq 1$, the claim follows. \blacksquare

We note before continuing that applying the upper bound of Corollary 7 instead of Lemma 8 would result in an inferior upper bound on π' .

As a matter of convenience, we denote moving forward $\bar{H}_q(\delta) = H_q(\min\{\frac{q-1}{q}, \delta\})$.

Motivated by the discussion at the beginning of this section, we make a few additional notations. First, fix $a > 1$, and denote $\ell_a \triangleq \lfloor a \log_q(n) \rfloor$ as n grows. Further, for a fixed real number $\delta > 0$, we denote $t_\delta \triangleq \lfloor \delta \ell_a \rfloor = \lfloor \delta \lfloor a \log_q(n) \rfloor \rfloor$. By slight abuse of notation, we let $\mathcal{RRF}_{\delta,a}^s(n) \triangleq \mathcal{RRF}_{t_\delta, \ell_a}^s(n)$, and $\mathcal{RRF}_{\delta,a}^s \triangleq \bigcup_{n \in \mathbb{N}} \mathcal{RRF}_{\delta,a}^s(n)$.

Inspired by Corollary 9, we also denote by $\tilde{\delta}_q$ the (unique) real number $0 < \tilde{\delta}_q < \frac{q-1}{q}$ satisfying

$$H_q\left(\frac{2}{3}\tilde{\delta}_q\right) = \frac{2}{3}.$$

We observe by substitution that $\tilde{\delta}_q > \frac{q-1}{2q}$, and provide $\tilde{\delta}_q$ for some small values of q :

q	$\frac{q-1}{2q}$	$\tilde{\delta}_q$	$\frac{q-1}{q}$
2	0.25	0.2609	0.5
3	0.3333	0.3723	0.6667
4	0.375	0.4375	0.75
5	0.4	0.4817	0.8
6	0.4167	0.5141	0.8333

Applying the result of Corollary 9 to Theorem 2, we can now obtain the following result.

Theorem 10 Fix $a > 1$, $0 < \delta < \tilde{\delta}_q$. Then, as $n \rightarrow \infty$,

$$\text{red}(\mathcal{RRF}_{\delta,a}^s(n)) = O(n^{2-\alpha(1-H_q(\delta))}).$$

Proof: If $a \leq (1 - H_q(\delta))^{-1}$ the proposition vacuously holds.

Otherwise, let $x \in \Sigma^{\ell_a+k}$ be chosen uniformly at random. Based on Corollary 9 (recalling again that $x \mapsto xH_q(1/x)$ is increasing for $x \geq 1$), we observe for $\delta < \tilde{\delta}_q$ that

$$\begin{aligned} \pi &\leq q \cdot n^{-\alpha(1-H_q(\delta))}, \\ \pi' &\leq q \cdot n^{-\alpha(1-\frac{3}{2}H_q(\frac{2}{3}\delta))}. \end{aligned}$$

Hence, for sufficiently large n it holds that $\ell_a^2 \pi' + \ell_a n \pi < 1/2e$, satisfying the conditions of Theorem 2. Since we also have $n \log(n) \pi' = o(n^{2-\alpha(1-H_q(\delta))})$, the claim follows from Theorem 2. ■

Corollary 11 Take $0 < \delta < \tilde{\delta}_q$. If $a > (1 - H_q(\delta))^{-1}$ then $R(\mathcal{RRF}_{\delta,a}^s(n)) = 1 - o(1)$, and if $a \geq 2(1 - H_q(\delta))^{-1}$, then $\mathcal{RRF}_{\delta,a}^s(n)$ incurs a constant number of redundant symbols.

The last corollary can be viewed in the context of related works; as mentioned above, [22] demonstrated that if $a > 1$ then $R(\mathcal{RF}_{\ell_a}(n)) = 1 - o_n(1)$, and if $a \geq 2$ then $\text{red}(\mathcal{RF}_{\ell_a}(n)) = O_n(1)$. Corollary 11 demonstrates that if $a > 1$ (respectively $a \geq 2$), then for all sufficiently small $\delta > 0$ it holds that $R(\mathcal{RRF}_{\delta,a}^s(n)) = 1 - o_n(1)$ (respectively, $\text{red}(\mathcal{RRF}_{\delta,a}^s(n)) = O_n(1)$). I.e., resilient-repeat-free sequences for a number of substitutions errors logarithmic in the string length (linear in the substring length) incur no additional asymptotic cost.

Up until here, we have focused on demonstrating conditions sufficient for the rate of resilient-repeat-free strings to be asymptotically optimal. In the sequel, we pursue the converse, i.e., necessary conditions for such strings to obtain non-vanishing rate.

Definition 12 For a real δ , $0 \leq \delta < 1$, and an integer $\ell > 0$, let $M_q(\ell, \delta)$ be the maximum number of code-words in a code $C \subseteq \Sigma^\ell$ such that $d_H(\mathbf{x}, \mathbf{y}) \geq \delta \ell$ for any distinct $\mathbf{x}, \mathbf{y} \in C$. For a given $\delta > 0$, define the maximum achievable rate by

$$R_q(\delta) \triangleq \limsup_{\ell \rightarrow \infty} \frac{1}{\ell} \log_q M_q(\ell, \delta).$$

For completeness, we state the well-known Gilbert-Varshamov and Elias-Bassalygo bounds (see, e.g., [36, Thm.4.9-12]) for $\delta \leq \frac{q-1}{q}$,

$$1 - H_q(\delta) \leq R_q(\delta) \leq 1 - H_q\left(\frac{q-1}{q}\left(1 - \sqrt{1 - \frac{q}{q-1}\delta}\right)\right).$$

The following lemma states a converse bound to Corollary 11.

Lemma 13 If $a < R_q(\delta)^{-1}$, then for sufficiently large $n \in \mathbb{N}$

$$\mathcal{RRF}_{\delta,a}^s(n) = \emptyset.$$

In particular, the statement holds if $t \geq \frac{q-1}{q}\ell_a$, for all a .

Proof: Take, on the contrary, some $\mathbf{x} \in \mathcal{RRF}_{\delta,a}^s(n)$. By Definition 1, the ℓ_a -mers

$$\{\mathbf{x}_{i\ell_a+[a]} : 0 \leq i \leq \lfloor n/\ell_a \rfloor - 1\} \subseteq \Sigma^{\ell_a}$$

form a code of size $\lfloor n/\ell_a \rfloor$ and minimum distance $d > t_\delta$. Hence, $d > \delta \ell_a$. By Definition 12 we obtain

$$\frac{\log \lfloor n/\ell_a \rfloor}{\ell_a} \leq R_q(\delta) + o(1).$$

Recalling $\ell_a = \lfloor a \log n \rfloor$ yields that

$$\frac{1}{a} \leq R_q(\delta) + o(1),$$

in contradiction to the assumption. ■

It should be noted that Lemma 13 specifically pertains to resilient-repeat-free strings, which the reader will observe are not necessarily required for successful reconstruction of information. Nevertheless, it might be conjectured, based on the noiseless case, that resilient-repeat-free sequences may achieve optimum asymptotic rate.

Before concluding, we note that a twofold gap remains between Theorem 10 and the converse of Lemma 13. First, $\text{red}(\mathcal{RRF}_{\delta,a}^s(n))$ is not characterized when $R_q(\delta)^{-1} \leq a \leq (1 - H_q(\delta))^{-1}$; and second, it is not found when $\delta \geq \tilde{\delta}_q$.

B. Encoding resilient-repeat-free codes

In this section, we present an explicit encoder of resilient-repeat-free strings, in the hope that it may then be utilized in constructing error-correcting codes for the noisy substring channel.

We first discuss how elements of the ball $B_t^s(\mathbf{0})$ (which throughout this discussion we assume to contain length- ℓ sequences) may be enumerated. Observe that given any $\mathbf{x}_{[k]} \in \Sigma^k$ with $\text{wt}(\mathbf{x}_{[k]}) \leq t$,

$$\begin{aligned} n(\mathbf{x}_{[k]}) &\triangleq \left| \left\{ \mathbf{y} \in B_t^s(\mathbf{0}) : \mathbf{y}_{[k]} = \mathbf{x}_{[k]} \right\} \right| \\ &= \sum_{j=0}^{t-\text{wt}(\mathbf{x}_{[k]})} \binom{\ell-k}{j} (q-1)^j. \end{aligned}$$

Assuming a total order $<$ on Σ , denote $\|x\| \triangleq |\{y \in \Sigma : y < x\}|$ for all $x \in \Sigma$. It was shown in [37] that the lexicographic index of $\mathbf{x} \in B_t^s(\mathbf{0})$ equals

$$\begin{aligned} i(\mathbf{x}) &= \sum_{k \in [\ell]} \sum_{\alpha < x(k)} n(\mathbf{x}_{[k-1]}\alpha) \\ &= \sum_{k \in [\ell]} \|\mathbf{x}(k)\| \cdot n(\mathbf{x}_{[k]}) + \\ &\quad \text{wt}(\mathbf{x}(k)) \cdot \binom{\ell-k}{t-\text{wt}(\mathbf{x}_{[k-1]})} (q-1)^{t-\text{wt}(\mathbf{x}_{[k-1]})}, \end{aligned}$$

where we set $\text{wt}(\mathbf{x}_{[k-1]}) \triangleq 0$ if $k = 0$, and $\text{wt}(x) = \mathbb{1}_{x \neq 0}$ for $x \in \Sigma$. Computationally, the most taxing expression to calculate in this sum is $n(\mathbf{x}_{[k]})$; however, one might employ a recursive approach to obtaining the sum. Indeed, from the Pascal identity we observe

$$\begin{aligned} n(\mathbf{x}_{[k-1]}) &= \sum_{j=0}^{t-\text{wt}(\mathbf{x}_{[k-1]})} \binom{\ell-(k-1)}{j} (q-1)^j \\ &= \sum_{j=1}^{t-\text{wt}(\mathbf{x}_{[k-1]})} \binom{\ell-k}{j-1} (q-1)^j + \\ &\quad \sum_{j=0}^{t-\text{wt}(\mathbf{x}_{[k-1]})} \binom{\ell-k}{j} (q-1)^j \\ &= q \sum_{j=0}^{t-\text{wt}(\mathbf{x}_{[k-1]})-1} \binom{\ell-k}{j} (q-1)^j + \\ &\quad \binom{\ell-k}{t-\text{wt}(\mathbf{x}_{[k-1]})} (q-1)^{t-\text{wt}(\mathbf{x}_{[k-1]})} \\ &= q \sum_{j=0}^{t-\text{wt}(\mathbf{x}_{[k-1]})} \binom{\ell-k}{j} (q-1)^j - \\ &\quad (q-1) \binom{\ell-k}{t-\text{wt}(\mathbf{x}_{[k-1]})} (q-1)^{t-\text{wt}(\mathbf{x}_{[k-1]})} \end{aligned}$$

Partitioning into cases by $\text{wt}(x(k-1))$, we find

$$\begin{aligned} n(\mathbf{x}_{[k-1]}) &= q \cdot n(\mathbf{x}_{[k]}) - \\ &\quad (-1)^{\text{wt}(x(k-1))} \binom{\ell-k}{t-\text{wt}(\mathbf{x}_{[k-1]})} (q-1)^{t-\text{wt}(\mathbf{x}_{[k]})+1}, \end{aligned}$$

where trivially $n(\mathbf{x}_{[\ell]}) = n(\mathbf{x}) = 1$.

It follows that computing the sum $i(\mathbf{x})$ can be done for $k \in [\ell]$ in descending order, where at each addend it is required to:

- 1) Compute binomial coefficients of the form $\binom{\ell-k}{j}$ for $j \leq t$ (all of order at most $\binom{\ell}{t} \leq \frac{\ell^t}{t!} \leq (\frac{e\ell}{t})^t$). Observe each binomial coefficient requires $\log(\frac{e\ell}{t})^t < t \log(\ell)$ symbols, and at most $t+1$ need to be stored at a time, so that $\{\binom{\ell-k}{j} : j \leq t\}$ could be computed via the Pascal identity from $\{\binom{\ell-k}{j} : j \leq t\}$. This stage hence requires $O(t^2 \log(\ell))$ operations, and $O(t^2 \log(\ell))$ space. Further, obtaining $\{\binom{\ell}{j} : j \leq t\}$ for initialization requires at most $O(t^2 \log(\ell)\ell)$ operations, if it is performed similarly.
- 2) Multiplying a binomial coefficient by at most q^t , which may practically be performed in $O(t \log(\ell) \log \log(\ell) \log \log \log(\ell))$ operations.
- 3) Computing $n(\mathbf{x}_{[k]})$, which has seen above requires $O(t \log(\ell))$ space and $O(t \log(\ell) \log \log(\ell) \log \log \log(\ell))$ operations.
- 4) Summing the results requires $O(t \log(\ell))$ operations and $O(t \log(\ell))$ space.

The entire algorithm therefore requires at most $O(t^2 \log(\ell)\ell)$ operations and $O(t^2 \log(\ell))$ space. That is, if $t, \ell = O(\log(n))$, at most $O((\log(n))^3 \log \log(n))$ operations and $O((\log(n))^2 \log \log(n))$ space.

The inverse operation, obtaining $\mathbf{x} \in B_t^s(\mathbf{0})$ such that $i(\mathbf{x}) = i$, for some given i , is also due to [37]: starting with the empty sequence for $k = 0$, assume $\mathbf{x}_{[k-1]}$ has already been constructed for some $0 < k \leq \ell$. Going over $\alpha \in \Sigma$ in increasing order (assuming the same total order as before), if $i \leq n(\mathbf{x}_{[k-1]}\alpha)$ then set $x(k-1) \triangleq \alpha$ and update $i \leftarrow i - n(\mathbf{x}_{[k-1]}\alpha)$; otherwise increase α and repeat; the maximum element of Σ can be filled in without comparison, if the algorithm arrives at it. Again, the limiting step of the algorithm is obtaining the representation of the binomial coefficients, and while the algorithm might require $t\ell$ steps in the worst case, these do not need to be recalculated unless k is increased. Thus, calculating the inverse also requires at most $O((\log(n))^3 \log \log(n))$ operations and $O((\log(n))^2 \log \log(n))$ space, for $t, \ell = O(\log(n))$.

In summary, we have obtained an explicit and invertible enumerator of $B_t^s(\mathbf{0})$, with the aforementioned complexity, which we denote $\text{en}(\mathbf{x})$. Recall that $|B_t^s(\mathbf{0})| \leq q^{\ell H_q(t/\ell)}$, i.e., $\text{en} : B_t^s(\mathbf{0}) \rightarrow \Sigma^{\lceil \ell H_q(t/\ell) \rceil}$.

Equipped with an (efficient) enumeration algorithm for $B_t^s(\mathbf{0})$, we may now propose an explicit encoder of resilient-repeat-free sequences. Our encoder has the drawback that it produces sequences in $\bigcup_{m \leq n} \mathcal{RRF}_{t,\ell,a}^s(m)$ rather than solely in $\mathcal{RRF}_{\delta,a}^s(n)$; however, in practice this does not seem too onerous for applications, were shorter sequences may be stored just as easily, as long as data is recoverable.

Our construction is summarized in Algorithm 1; its main idea is a generalization of [22, Alg. 3], as follows. Assume $a(1 - H_q(\delta)) > 1$, and choose $\epsilon > 0$ such that $\zeta \triangleq a(1 - H_q(\delta) - \epsilon) - 1 > 0$. Let $z = \lfloor \zeta \log(n) \rfloor$. An information string is first encoded into a length- n string \mathbf{x} containing no 0-run of length z , which may be done in linear time using $\lceil \frac{a}{q-2} n^{1-\zeta} \rceil = O(n^{2-a(1-H_q(\delta)-\epsilon)})$ redundant symbols [38, Lem. 4]. Interestingly, this allows us to achieve redundancy which is arbitrarily close, in orders of magnitude, to the result

Algorithm 1: Resilient-repeat-free Encoder**Input:** $\mathbf{x} \in \Sigma^n$ containing no 0-run of length z **Output:** $\text{Enc}_1(\mathbf{x}) \in \bigcup_{m \leq n} \mathcal{RRF}_{t_\delta, \ell_a}^s(m)$ $j \leftarrow 1$ **while** $j \leq |\mathbf{x}| - \ell_a$ **do** **for** $i = j - 1, \dots, 0$ **do** **if** $\exists \mathbf{y} \in B_{t_\delta}^s(\mathbf{x}) : \mathbf{y}_{i+[\ell_a]} = \mathbf{y}_{j+[\ell_a]}$ **then** Replace $\mathbf{x}_{j+[\ell_a]}$ with \mathbf{s} from (3) $j \leftarrow \max\{0, j - \ell_a + 1\}$ **break** **end** **end** $j \leftarrow j + 1$ **end****return** \mathbf{x}

of Theorem 10. Next, using

$$\ell \triangleq 11 + \lceil 2 \log(\ell_a) + a(1 - \epsilon) \log(n) \rceil \leq \ell_a, \quad (2)$$

where the last inequality holds for all sufficiently large n , it is then iteratively checked whether $\mathbf{x}_{[\ell+j]} \in \mathcal{RRF}_{\delta, a}^s(\ell + j)$, for $j \in [n - \ell + 1]$ in increasing order.

If in some iteration it is determined that $\mathbf{x}_{[\ell+j]} \notin \mathcal{RRF}_{\delta, a}^s(\ell + j)$, then the algorithm deletes $\mathbf{x}_{j+[\ell]}$ from \mathbf{x} and replaces it with a sequence with the following form:

$$\mathbf{s} \triangleq 0^{z'} \circ E(j - i) \circ 10^{z'} \circ E(\text{en}(\mathbf{e})) \circ 1, \quad (3)$$

where

- j, i are the loop-indices at any specific iteration of Algorithm 1, and by abuse of notation we take $(j - i)$ to represent the q -ary expansion of the difference, using only as many symbols as required;
- $\mathbf{e} \triangleq \mathbf{x}_{j+[\ell]} - \mathbf{x}_{i+[\ell]} \in \Sigma^\ell$ when $j - i \geq \ell$, or $\mathbf{e} \triangleq \mathbf{x}_{i+[\ell+k]} - \mathbf{y}_{i+[\ell+k]} \in \Sigma^{\ell+k}$ when $j - i < \ell$; recall, however, that from the proof of Lemma 8 it follows that we may always assume $|\text{supp}(\mathbf{e})| \leq \lfloor 3\ell/2 \rfloor$ is known, even for $k > \lfloor \ell/2 \rfloor$. In both cases $\text{wt}(\mathbf{e}) \leq t_\delta$, and $\mathbf{x}_{j+[\ell]}$ is recoverable from $\mathbf{e}, \mathbf{x}_{i+[\ell-j]}$; and
- $z' \triangleq \lceil \log(\ell_a) \rceil + 2$ and $E(\cdot)$ is the explicit and efficient encoder described in [39, Alg. 1]; it can accept sequences of lengths at most $q\ell_a$, and returns an encoded version containing no 0-runs of length z' , utilizing only a single redundant symbol. Observe that both $\log(j - i), |\text{en}(\mathbf{e})| \leq \frac{3}{2}\ell_a < q\ell_a$.

We note that since $E(\cdot)$ may accept sequences of varying (sufficiently small) lengths, so too is $|\mathbf{s}|$ not constant. Our next aim is to bound $|\mathbf{s}|$ from above, as this affects the correct operation of Algorithm 1, namely, its termination condition.

Theorem 14 *If $\frac{3}{2}\bar{H}_q(\frac{2}{3}\delta) - H_q(\delta) \leq \frac{1}{a}$, then Algorithm 1 terminates, $\text{Enc}_1(\mathbf{x}) \in \bigcup_{m \leq n} \mathcal{RRF}_{t_\delta, \ell_a}^s(m)$, and \mathbf{x} can be decoded from it.*

Proof: First observe that if the last iteration of Algorithm 1 terminates, then its output is resilient-repeat-free.

Next, we will show that the inserted substring in (3) is strictly less than ℓ , hence each replacement that the algorithm

performs *shortens* \mathbf{x} . As a consequence, the algorithm must terminate. Indeed, observe that

$$|\mathbf{s}| < 10 + \log(\ell_a) + \zeta \log(n) + \log(j - i) + |\text{en}(\mathbf{e})|.$$

If $j - i \geq \ell$, then we bound $\log(j - i) \leq \log(n)$, and we have seen that

$$\begin{aligned} |\text{en}(\mathbf{e})| &\leq \lceil \ell \bar{H}_q(t_\delta/\ell) \rceil \leq \lceil \ell_a \bar{H}_q(t_\delta/\ell_a) \rceil \\ &\leq \lceil \ell_a H_q(\delta) \rceil < a H_q(\delta) \log(n) + 1, \end{aligned}$$

hence in this case

$$\begin{aligned} |\mathbf{s}| &< 11 + \log(\ell_a) + (\zeta + 1 + a H_q(\delta)) \log(n) \\ &= 11 + \log(\ell_a) + a(1 - \epsilon) \log(n) < \ell, \end{aligned}$$

as required. Otherwise we bound $\log(j - i) \leq \log(\ell) \leq \log(\ell_a)$ and

$$\begin{aligned} |\text{en}(\mathbf{e})| &\leq \lceil (\ell + k) \bar{H}_q(t_\delta/(\ell + k)) \rceil \\ &< \frac{3}{2} \ell_a \bar{H}_q(\frac{2}{3}\delta) + 1, \end{aligned}$$

where $k \triangleq \min\{j - i, \lfloor \ell/2 \rfloor\}$. Hence,

$$\begin{aligned} |\mathbf{s}| &< 11 + 2 \log(\ell_a) + (\zeta + \frac{3}{2} a \bar{H}_q(\frac{2}{3}\delta)) \log(n) \\ &\leq \ell + (a(\frac{3}{2} \bar{H}_q(\frac{2}{3}\delta) - H_q(\delta)) - 1) \log(n). \end{aligned}$$

Under the assumption of the theorem, it also holds in this case that $|\mathbf{s}| < \ell$.

Lastly, observe that, iterating over $j \in [|\text{Enc}_1(\mathbf{x})| - \ell]$ in decreasing order, the first observed instance of $0^{z'}$ is always the last to have been inserted by Algorithm 1; this holds because after each replacement, j is decreased only by $\ell - 1$, hence any later replacements, say at index j' , either satisfy $j' > j$ or they overwrite the first 0 of $0^{z'}$ (observe that \mathbf{s} ends with a 1). Further, by observing the first instance of $0^{z'}$ following that instance of $0^{z'}$, it is possible to uniquely deduce the coordinates of $E(j - i)$, and therefore to deduce i . Now, given $E(\text{en}(\mathbf{e}))$ one obtains \mathbf{e} , and with $\mathbf{e}, \mathbf{x}_{i+[\ell-j]}$ is uniquely possible to reconstruct the removed segment $\mathbf{x}_{j+[\ell]}$. Since every replacement of Algorithm 1 is reversible, and the process can be tracked in reverse, \mathbf{x} can be reconstructed. ■

Lemma 15 *The run-time of Algorithm 1 is $O(n^2 \log(n)^2)$.*

Proof: In any iteration, if there exists $\mathbf{y} \in B_{t_\delta}^s(\mathbf{x}_{[\ell+j]}) \setminus \mathcal{RRF}_{\ell_a}(\ell + j)$, and j is minimal such that this occurs, then there necessarily exists $i < j$ such that $\mathbf{y}_{i+[\ell]} = \mathbf{y}_{j+[\ell]}$. By Lemma 6, if $i \leq j - \ell$, the existence of such \mathbf{y} is equivalent to $\text{wt}(\mathbf{x}_{j+[\ell]} - \mathbf{x}_{i+[\ell]}) \leq t_\delta$, which may be verified in at most $\ell(j - \ell)$ operations. On the other hand, for each $0 < k < \ell$ we check whether there exists $\mathbf{y} \in B_{t_\delta}^s(\mathbf{x}_{[\ell+j]})$ with $i = j - k$; as seen in the proof of Lemma 8, this implies that $\mathbf{y}_{i+[\ell+k]}$ is k -periodic. The following procedure verifies whether such \mathbf{y} exists: denote for convenience $\mathbf{u} \triangleq \mathbf{x}_{i+[\ell+k]}$; for each $p \in [k]$, we define the multiset $U_p \triangleq \{u(q) : q \equiv p \pmod{k}\}$. If the most frequent element in U_p is some $x \in \Sigma$, denote by t_p the number of occurrences of all other elements in U_p ; clearly, there exists \mathbf{y} with the given i if and only if $\sum_{p \in [k]} t_p \leq t_\delta$. This algorithm requires at most $O(\ell \log(\ell))$ operations for each k (due to the summation). For $\ell = O(\log(n))$, any

iteration requires at most $O(n \log(n))$ operations in total, for both cases.

Finally, observe that there could be at most $n\ell \leq n\ell_a$ iterations of Algorithm 1, completing the proof. \blacksquare

C. Error-correcting codes for the noisy substring channel

Based on Corollary 11, we can demonstrate the existence of error-correcting codes for the noisy substring channel, which achieve at most a constant redundancy over that of classical error-correcting codes for Hamming noise.

Corollary 16 *Let $C \subseteq \Sigma^n$ be an error-correcting code, capable of correcting t_δ substitution errors, and denote, for some $\mathbf{z} \in \Sigma^n$, $\bar{C}_z \triangleq (\mathbf{z} + C) \cap \mathcal{RRF}_{\delta,a}^s(n)$. Then for any $\mathbf{x} \in \bar{C}_z$ and $\mathbf{y} \in B_{t_\delta}^s(\mathbf{x})$, it is possible to uniquely decode \mathbf{x} observing only $Z_{\ell_a+1}(\mathbf{y})$. Further, decoding is possible through a greedy algorithm for reconstruction of \mathbf{y} , followed by application of any decoding scheme for C .*

Finally, in the cases indicated in Corollary 11, where $\text{red}(\mathcal{RRF}_{\delta,a}^s(n)) = O(1)$, there exists \mathbf{z} satisfying $\text{red}(\bar{C}_z) = \text{red}(C) + O(1)$.

Note that Corollary 16 is unfortunately nonconstructive. It is our hope that the encoder of Algorithm 1 may be combined with error-correction techniques to yield explicit code constructions for this channel. However, achieving this goal seems to require new ideas, and we leave it for future study.

V. DELETION NOISE

This section is dedicated to the study of resilient-repeat-free sequences under deletion, rather than Hamming, errors. We demonstrate that the same probabilistic tools can be used to bound from below the redundancy of such sequences. We remark that the same method can be used to study insertion errors, even though the equivalence of insertion/deletion-correction does not extend in a straightforward manner to our setting.

For $\mathbf{x} \in \Sigma^n$, let $S_t^d(\mathbf{x}) \subseteq \Sigma^{n-t}$ denote the set of strings generated from \mathbf{x} by t deletions.

Definition 17 *For integers $t, \ell \leq n$, define a family of repeat-free strings which is resistant to deletion noise:*

$$\mathcal{RRF}_{t,\ell}^d(n) \triangleq \{\mathbf{x} \in \Sigma^n : S_t^d(\mathbf{x}) \subseteq \mathcal{RF}_\ell(n-t)\}.$$

Recall the definitions $\ell_a \triangleq \lceil a \log_q(n) \rceil$ and $t_\delta \triangleq \lfloor \delta \ell_a \rfloor$, for some fixed real numbers $a > 1$ and $\delta > 0$. As in the previous sections, we abbreviate $\mathcal{RRF}_{\delta,a}^d(n) \triangleq \mathcal{RRF}_{t_\delta, \ell_a}^d(n)$. Then we have the following:

Theorem 18 *For all $a > 1$ and $\delta > 0$ it holds that*

$$\text{red}(\mathcal{RRF}_{\delta,a}^d(n)) = O\left(n^{2-a+\frac{2a(1+\delta)}{\log_2(q)}} H_2(\delta/(1+\delta)) / \log(n)\right).$$

Proof: We follow a similar strategy as in Theorem 2, but apply a symmetric bound in Theorem 5, i.e., utilizing a fixed constant $f_{i,j} \equiv f$. Note that a sufficient condition for $\mathbf{x} \in$

$\mathcal{RRF}_{\delta,a}^d(n)$ is that for every observable pair $(I, J) \in \binom{[n]}{\ell_a}^2$, such that

$$I(\ell_a - 1) - I(0) < \ell_a + t_\delta \quad (4)$$

(and similarly for J), it holds that $\mathbf{x}_I \neq \mathbf{x}_J$. For such a pair, denote

$$A_{I,J} \triangleq \{\mathbf{x} \in \Sigma^n : \mathbf{x}_I = \mathbf{x}_J\} = \{\mathbf{x} \in \Sigma^n : \mathbf{u}_{I,J} = 0\}.$$

Again, we let $\mathbf{x} \in \Sigma^n$ be chosen uniformly at random, implying $\text{red}(\mathcal{RRF}_{\delta,a}^d(n)) = -\log_q \Pr(\mathbf{x} \in \mathcal{RRF}_{\delta,a}^d(n))$, and use a symmetric version of Lovász's local lemma (i.e., using a fixed constant $f_{i,j} \equiv f$) to bound $\Pr(\mathbf{x} \in \mathcal{RRF}_{\delta,a}^d(n)) \geq \Pr(\mathbf{x} \notin \bigcup_{I,J} A_{I,J})$ from below.

For any observable pair (I, J) , note that $\Pr(\mathbf{x} \in A_{I,J}) = q^{-\ell_a} \leq q \cdot n^{-a}$, and for convenience denote $\pi_d \triangleq q \cdot n^{-a}$.

Next, we estimate $|\Gamma|$, where Γ is a set of observable pairs (P, Q) also satisfying (4), such that the event $\{\mathbf{x} \in A_{I,J}\}$ is mutually independent of $\{\{\mathbf{x} \in A_{P,Q}\} : (P, Q) \notin \Gamma\}$. Observe by Lemma 4 that it suffices that Γ consists of all $(P, Q) \notin L_I$ satisfying (4). Thus, to determine P , it suffices to choose

- 1) a single element of I (which shall be a member of $P \cap I$);
- 2) an interval of length $\ell_a + t_\delta$ containing the chosen element; and
- 3) any $\ell_a - 1 < \ell_a$ additional elements of the chosen interval.

Then Q can be chosen from any interval of length $\ell_a + t_\delta$. The same holds for a suitable choice of $Q \cap I \neq \emptyset$. Thus, $|\Gamma| \leq \ell_a(\ell_a + t_\delta)n^{\binom{\ell_a + t_\delta}{\ell_a}^2}$.

Now, to apply LLL, we find $0 < f < 1$ such that $\pi_d \leq f(1-f)^{|\Gamma|}$. From $\binom{s}{r} \leq \sqrt{\frac{s}{r(s-r)}} 2^{sH_2(r/s)}$ (a relaxation of, e.g., [41, Ch.10, Sec.11, Lem.7]) we observe that

$$\binom{\ell_a + t_\delta}{t_\delta}^2 \leq \frac{\ell_a + t_\delta}{\ell_a t_\delta} n^{\frac{2a(1+\delta)}{\log_2(q)} H_2(\delta/(1+\delta))}.$$

If $\frac{2(1+\delta)}{\log_2(q)} H_2(\frac{\delta}{1+\delta}) \geq 1$ or $a < (1 - \frac{2(1+\delta)}{\log_2(q)} H_2(\frac{\delta}{1+\delta}))^{-1}$, the theorem vacuously holds. Otherwise, we note that

$$\begin{aligned} (|\Gamma| + 1)\pi_d &\leq qn^{-a} + q \frac{(\ell_a + t_\delta)^2}{t_\delta} n^{1-a+\frac{2a(1+\delta)}{\log_2(q)} H_2(\frac{\delta}{1+\delta})} \\ &= o_n(1). \end{aligned}$$

Denote $y \triangleq e\pi_d$, and for sufficiently large n observe that $ye^{-(|\Gamma|+1)y} = \pi_d e^{1-e^{-(|\Gamma|+1)\pi_d}} > \pi_d$. By further letting $f \triangleq \frac{y}{1+y}$, and recalling for $0 < f < 1$ that $1-f \geq e^{-f/(1-f)}$, we observe

$$f(1-f)^{|\Gamma|} \geq \frac{f}{1-f} e^{-\frac{f}{1-f}(|\Gamma|+1)} = ye^{-(|\Gamma|+1)y} > \pi_d.$$

Finally, one needs also note that the number of observable pairs (I, J) satisfying the given requirements is no more than $\binom{n-\ell_a-t_\delta}{2} \cdot \binom{\ell_a+t_\delta}{\ell_a}^2 < n^2 \binom{\ell_a+t_\delta}{\ell_a}^2$. From LLL it follows that

$$\begin{aligned} \Pr\left(\mathbf{x} \notin \bigcup_{I,J} A_{I,J}\right) &\geq (1-f)^{n^2 \binom{\ell_a+t_\delta}{\ell_a}^2} \\ &= (1+y)^{-n^2 \binom{\ell_a+t_\delta}{\ell_a}^2}, \end{aligned}$$

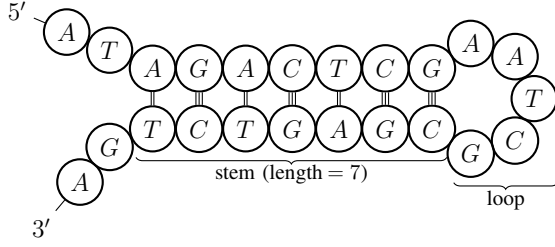


Figure 1. Formation of a hairpin-loop secondary structure in an oligonucleotide.

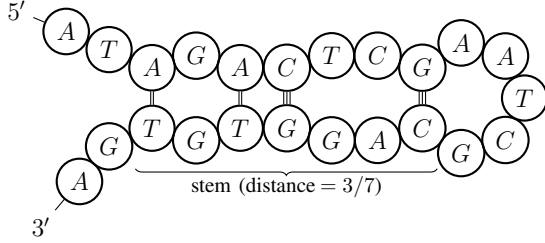


Figure 2. Imperfect stem in a hairpin-loop structure.

and hence

$$\begin{aligned}
\text{red}(\mathcal{RRF}_{\delta,a}^d(n)) &= -\log_q \Pr(x \in \mathcal{RRF}_{\delta,a}^d(n)) \\
&\leq n^2 \binom{\ell_a + t_\delta}{\ell_a}^2 \log_q(1+y) \\
&\leq \frac{1}{\ln(q)} \binom{\ell_a + t_\delta}{\ell_a}^2 n^2 y \\
&\leq \frac{e \cdot q}{\ln(q)} \binom{\ell_a + t_\delta}{\ell_a}^2 n^{2-a},
\end{aligned}$$

which completes the proof. \blacksquare

Corollary 19 *If $a > (1 - \frac{2(1+\delta)}{\log_2(q)} H_2(\frac{\delta}{1+\delta}))^{-1}$, for any $\delta > 0$, then $R(\mathcal{RRF}_{\delta,a}^d(n)) = 1 - o_n(1)$, and if $a \geq 2(1 - \frac{2(1+\delta)}{\log_2(q)} H_2(\frac{\delta}{1+\delta}))^{-1}$ then $\text{red}(\mathcal{RRF}_{\delta,a}^d(n)) = O_n(1)$.*

Note again that if $a > 1$ (respectively $a > 2$), then for all sufficiently small $\delta > 0$ it holds that $R(\mathcal{RRF}_{\delta,a}^d(n)) = 1 - o_n(1)$ (respectively, $\text{red}(\mathcal{RRF}_{\delta,a}^d(n)) = O_n(1)$). Before concluding, we also note that a parallel statement to Corollary 16 holds in this setting, as well.

VI. SECONDARY STRUCTURE AVOIDANCE

In this section, we leverage Algorithm 1 to protect against the formation of secondary structures in coded DNA strands. We focus on a special type of secondary structure, called *hairpin-loop* (see Figure 1). Unlike recent works, our analysis does not require a perfect binding in the stem region of the hairpin structures (see Figure 2), and we show that Algorithm 1 can be utilized to avoid the formation of such structures. However, we rely in this section on the Hamming metric rather than the Levenshtein metric as was suggested in [30], thus we do not consider the formation of so-called *bulge-loops* due

to the elasticity of the DNA sugar-phosphate backbone. We remark that given an efficient enumeration of the Levenshtein ball about any point, these methods can be extended to utilize that metric, too.

In order to define hairpin-loop-avoiding sequences, we first present some notation. An involution on Σ is a mapping $x \mapsto \bar{x}$ such that for all $x \in \Sigma$ it holds that $\bar{\bar{x}} = x$; we now assume Σ to be equipped with such an involution (we allow fixed points, in order to account for odd q , which shall not affect our analysis). For example, DNA is composed of four *nucleotide bases*: the *purines*, adenine (A) and guanine (G), are respectively the complements of the *pyrimidines* thymine (T) and cytosine (C); when forming a *double helix* (or *duplex*) structure, each base can only stably bond (*hybridize*) with its complement. For $x \in \Sigma^n$, denote $\bar{x} \triangleq \bar{x}(0)\bar{x}(1) \cdots \bar{x}(n-1)$.

DNA strands are also oriented: each nucleotide is composed of one of four nitrogenous bases, together with a pentose sugar and a phosphate group; the phosphate groups connect the sugar rings of adjacent nucleotides 5'-end to 3'-end (referring to the five-carbon sites of the sugar rings) to form a long chain (*oligonucleotide*), and thus the orientation can be observed from any segment of the chain. Stable duplexes only form between oligonucleotides of reverse orientations, and therefore coiled-loop secondary structures cannot appear. To capture this notion, we denote for $x \in \Sigma^n$ the *reverse* sequence $x^r \triangleq x(n-1) \cdots x(1)x(0)$.

For integers $t \leq \ell$, we define the set of length- n (t, ℓ)-*hairpin avoiding* strings to contain those strings that do not have the potential for the formation of a loop with stem-length ℓ , of which at least $\ell - t$ symbols are hybridized. More precisely,

$$\mathcal{HA}_{t,\ell}(n) \triangleq \left\{ x \in \Sigma^n : \begin{array}{l} \forall 0 \leq i < j < n \\ \ell - t \leq \ell' \leq \min\{\ell, j - i, n - j\} \\ d_H(x_{i+[\ell']}, (\bar{x}_{j+[\ell']})^r) > t - (\ell - \ell') \end{array} \right\}.$$

Observe in the above definition, that for $0 < i < j < n - 1$,

$$\begin{aligned}
d_H(x_{i-1+[\ell'+1]}, (\bar{x}_{j+[\ell'+1]})^r) &> t - (\ell - \ell') + 1 \\
\implies d_H(x_{i+[\ell']}, (\bar{x}_{j+[\ell']})^r) &> t - (\ell - \ell'),
\end{aligned}$$

hence some conditions in the above definition are redundant.

As before, for fixed real numbers $a > 1$ and $0 < \delta < 1$, we also make the notation $\mathcal{HA}_{\delta,a}(n)$. We will show that when $a > (1 - \bar{H}_q(\delta))^{-1}$ (for $\delta < 1 - \frac{1}{q}$) then Algorithm 1 can, with slight necessary adjustments, encode into $\bigcup_{m \leq n} \mathcal{HA}_{t_\delta, \ell_a}(m)$ with redundancy $O(n^{2-a(1-H_q(\delta)-\epsilon)})$ for arbitrarily small $\epsilon > 0$. We leave the interesting problem of stating an analogue of Lemma 13 in this case for future study.

Indeed, the encoder presented in Algorithm 2 differs from Algorithm 1 only in the type of condition in the inner loop (ranges for j, i are adjusted accordingly); if a replacement is required, instead of necessarily replacing an entire ℓ_a -substring, in the case $i > j - \ell_a$ (i.e., $\ell' < \ell_a$) only the ℓ' -suffix is replaced, with the substring

$$s \triangleq 0^z 1 \circ E(j - i) \circ 10^{z'} 1 \circ E(\text{en}(e)) \circ 1. \quad (5)$$

For convenience we repeat the previous definitions for the following expressions:

- $\zeta = a(1 - H_q(\delta) - \epsilon) - 1 > 0$;

Algorithm 2: Hairpin-avoiding Encoder

Input: $\mathbf{x} \in \Sigma^n$ containing no 0-run of length z
Output: $\text{Enc}_1(\mathbf{x}) \in \bigcup_{m \leq n} \mathcal{RRF}_{t_\delta, \ell_a}^s(m)$
 $j' \leftarrow 2(\ell_a - t_\delta)$
while $j' \leq |\mathbf{x}|$ **do**
 for $i = j' - 2(\ell_a - t_\delta), \dots, 0$ **do**
 $\ell' \leftarrow \min\{\ell_a, \lfloor \frac{j'-i}{2} \rfloor\}$
 $j \leftarrow j' - \ell'$
 if $d((\bar{\mathbf{x}}_{j+[\ell']})^r, \mathbf{x}_{i+[\ell']}) \leq t_\delta - (\ell_a - \ell')$ **then**
 Replace $\mathbf{x}_{j+[\ell']}$ with \mathbf{s} from (5)
 $j' \leftarrow \max\{2(\ell_a - t_\delta), j' - \ell' + 1\}$
 break
 end
 end
 $j' \leftarrow j' + 1$
end
return \mathbf{x}

- $z = \lfloor \zeta \log(n) \rfloor$;
- $j - i$ represents the q -ary expansion of the difference (using only as many symbols as required),

and we adjust the following definitions:

- $\mathbf{e} \triangleq (\bar{\mathbf{x}}_{j+[\ell']})^r - \mathbf{x}_{i+[\ell']}$;
- $z' \triangleq \lceil \log(\ell_a) \rceil + 1$; and, finally,
- $E(\cdot)$ is an explicit and efficient encoder into strings containing no 0-runs of length z' , accepting inputs of lengths at most ℓ_a and requiring a single redundant symbol [39, Alg. 1]. We shall see below that both $\log(j - i), |\text{en}(\mathbf{e})| \leq \ell_a$.

The analysis of Algorithm 2 is much similar to that of Algorithm 1 in Section IV-B. We summarize the result in the following theorem.

Theorem 20 *If $(1 - \bar{H}_q(\delta))^{-1} < a < \delta^{-1}(1 - \bar{H}_q(\delta))^{-1}$, then for sufficiently large n Algorithm 2 terminates, $\text{Enc}_2(\mathbf{x}) \in \bigcup_{m \leq n} \mathcal{HA}_{t_\delta, \ell_a}(m)$, and \mathbf{x} can be decoded from it.*

Proof: We prove only the first part; the latter two follow exactly as in the proof of Theorem 14. As before,

$$|\mathbf{s}| < 9 + \log(\ell_a) + \zeta \log(n) + \log(j - i) + |\text{en}(\mathbf{e})|.$$

Repeating the analysis of Theorem 14, if $j - i \geq \ell_a$ (i.e., $\ell' = \ell_a$), then (bounding $\log(j - i) \leq \log(n)$) we have

$$\begin{aligned} |\text{en}(\mathbf{e})| &\leq \lceil \ell_a \bar{H}_q(t_\delta / \ell_a) \rceil \leq \lceil \ell_a H_q(\delta) \rceil \\ &< a H_q(\delta) \log(n) + 1, \end{aligned}$$

as before, and hence again (for sufficiently large n)

$$\begin{aligned} |\mathbf{s}| &< 9 + \log(\ell_a) + (\zeta + 1 + a H_q(\delta)) \log(n) \\ &= 9 + \log(\ell_a) + a(1 - \epsilon) \log(n) < \ell_a. \end{aligned}$$

It follows that such an iteration of Algorithm 2 also shortens \mathbf{x} .

Otherwise, when $\ell' < \ell_a$ we again bound $\log(j - i) \leq \log(\ell') < \log(\ell_a)$ and

$$\begin{aligned} |\text{en}(\mathbf{e})| &\leq \lceil \ell' \bar{H}_q\left(\frac{t_\delta - (\ell_a - \ell')}{\ell'}\right) \rceil \\ &\leq \ell' \bar{H}_q\left(\frac{\ell' - (1 - \delta)\ell_a}{\ell'}\right) + 1 \\ &< \ell' H_q(\delta) + 1. \end{aligned}$$

Hence,

$$\begin{aligned} |\mathbf{s}| &< 10 + 2 \log(\ell_a) + \zeta \log(n) + H_q(\delta) \ell' \\ &= 10 + 2 \log(\ell_a) - (\epsilon a + 1) \log(n) \\ &\quad + (1 - H_q(\delta)) a \log(n) + H_q(\delta) \ell' \\ &\leq 11 + 2 \log(\ell_a) - (\epsilon a + 1) \log(n) \\ &\quad + (1 - H_q(\delta)) \ell_a + H_q(\delta) \ell' \\ &\leq 11 + 2 \log(\ell_a) - (\epsilon a + 1) \log(n) \\ &\quad + (1 - H_q(\delta)) (\ell_a - \ell') + \ell' \\ &\leq (11 + 2 \log(\ell_a) - \epsilon a \log(n)) + \ell' \\ &\quad + ((1 - H_q(\delta)) \delta a - 1) \log(n) < \ell', \end{aligned}$$

where the last inequality again holds for sufficiently large n , and relies on the theorem's assumption. ■

Corollary 21 *For all $a > 1$ and sufficiently small $\delta > 0$, there exists an efficient (explicit) encoder from $\Sigma^{n-o(n)}$ into $\bigcup_{m \leq n} \mathcal{HA}_{t_\delta, \ell_a}(m)$, for sufficiently large n .*

The problems of encoding directly into $\mathcal{HA}_{\delta, a}(n)$, more precisely bounding its redundancy, as well as generalization to the Levenshtein metric (hence, considering also the formation of bulge-loop secondary structures), are left for future study.

REFERENCES

- [1] G. M. Church, Y. Gao, and S. Kosuri, "Next-generation digital information storage in DNA," *Science*, vol. 337, no. 6102, pp. 1628–1628, 2012.
- [2] F. Balado, "Capacity of DNA data embedding under substitution mutations," *IEEE Trans. Inf. Theory*, vol. 59, no. 2, pp. 928–941, Feb. 2013.
- [3] P. C. Wong, K.-k. Wong, and H. Foote, "Organic data memory using the DNA approach," *Commun. ACM*, vol. 46, no. 1, pp. 95–98, Jan. 2003.
- [4] S. L. Shipman, J. Nivala, J. D. Macklis, and G. M. Church, "CRISPR-Cas encoding of a digital movie into the genomes of a population of living bacteria," *Nature*, vol. 547, p. 345, Jul. 2017.
- [5] M. Arita and Y. Ohashi, "Secret signatures inside genomic DNA," *Biotechnology Progress*, vol. 20, no. 5, pp. 1605–1607, 2004.
- [6] D. Heider and A. Barnekow, "DNA-based watermarks using the DNA-Crypt algorithm," *BMC Bioinf.*, vol. 8, no. 1, pp. 176–185, May 2007.
- [7] M. Liss, D. Daubert, K. Brunner, K. Kliche, U. Hammes, A. Leiberer, and R. Wagner, "Embedding permanent watermarks in synthetic genes," *PLoS ONE*, vol. 7, no. 8, p. e42465, 2012.
- [8] D. C. Jupiter, T. A. Ficht, J. Samuel, Q.-M. Qin, and P. de Figueiredo, "DNA watermarking of infectious agents: Progress and prospects," *PLoS pathogens*, vol. 6, no. 6, p. e1000950, 2010.
- [9] C. T. Clelland, V. Risca, and C. Bancroft, "Hiding messages in DNA microdots," *Nature*, vol. 399, no. 6736, pp. 533–534, 1999.
- [10] R. Heckel, G. Mikutis, and R. N. Grass, "A characterization of the DNA data storage channel," *Scientific reports*, vol. 9, no. 1, p. 9663, Jul. 2019.
- [11] O. Sabary, Y. Orlev, R. Shafir, L. Anavy, E. Yaakobi, and Z. Yakhini, "SOLQC: Synthetic oligo library quality control tool," *Bioinformatics*, vol. 37, no. 5, pp. 720–722, Mar. 2021.
- [12] H. M. Kiah, G. J. Puleo, and O. Milenkovic, "Codes for DNA sequence profiles," *IEEE Trans. Inf. Theory*, vol. 62, no. 6, pp. 3125–3146, Jun. 2016.
- [13] A. Lenz, P. H. Siegel, A. Wachter-Zeh, and E. Yaakobi, "Coding over sets for DNA storage," *IEEE Trans. Inf. Theory*, vol. 66, no. 4, pp. 7682–7696, Apr. 2020.
- [14] J. Sima, N. Raviv, and J. Bruck, "Robust indexing - optimal codes for DNA storage," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Los Angeles, CA, USA, Jun. 2020, pp. 717–722.
- [15] —, "On coding over sliced information," *IEEE Trans. Inf. Theory*, vol. 67, no. 5, pp. 2793–2807, May 2021.
- [16] E. Ukkonen, "Approximate string-matching with q-grams and maximal matches," *Theoretical Computer Science*, vol. 92, no. 1, pp. 191–211, Jan. 1992.

- [17] J. Acharya, H. Das, O. Milenkovic, A. Orlitsky, and S. Pan, "String reconstruction from substring compositions," *SIAM J. Discrete Math.*, vol. 29, no. 3, pp. 1340–1371, 2015.
- [18] I. Shomorony, T. A. Courtade, and D. Tse, "Fundamental limits of genome assembly under an adversarial erasure model," *IEEE Trans. Mol., Bio. and Multi-Scale Commun.*, vol. 2, no. 2, pp. 199–208, Dec. 2016.
- [19] Z. Chang, J. Chrisnata, M. F. Ezerman, and H. M. Kiah, "Rates of DNA sequence profiles for practical values of read lengths," *IEEE Trans. Inf. Theory*, vol. 63, no. 11, pp. 7166–7177, Nov. 2017.
- [20] R. Gabrys and O. Milenkovic, "Unique reconstruction of coded strings from multiset substring spectra," *IEEE Trans. Inf. Theory*, vol. 65, no. 12, pp. 7682–7696, Dec. 2019.
- [21] J. Chrisnata, H. M. Kiah, S. Rao, A. Vardy, E. Yaakobi, and H. Yao, "On the number of distinct k -decks: Enumeration and bounds," in *Proceedings of the 2019 19th International Symposium on Communications and Information Technologies (ISCIT)*, Ho Chi Minh City, Vietnam, Vietnam, Sep. 2019, pp. 519–524.
- [22] O. Elishco, R. Gabrys, M. Médard, and E. Yaakobi, "Repeat-free codes," *IEEE Trans. Inf. Theory*, vol. 67, no. 9, pp. 5749–5764, Sep. 2021.
- [23] S. Marcovich and E. Yaakobi, "Reconstruction of strings from their substrings spectrum," *IEEE Trans. Inf. Theory*, vol. 67, no. 7, pp. 4369–4384, Jul. 2021.
- [24] A. Lenz, P. H. Siegel, A. Wachter-Zeh, and E. Yaakobi, "Anchor-based correction of substitutions in indexed sets," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Paris, France, Jul. 2019, pp. 757–761.
- [25] N. Raviv, M. Schwartz, and E. Yaakobi, "Rank-modulation codes for DNA storage with shotgun sequencing," *IEEE Trans. Inf. Theory*, vol. 65, no. 1, pp. 50–64, Jan. 2019.
- [26] N. Beeri and M. Schwartz, "Improved rank-modulation codes for DNA storage with shotgun sequencing," *IEEE Trans. Inf. Theory*, vol. 68, no. 6, pp. 3719–3730, Jun. 2022.
- [27] S. Jain, F. Farnoud, M. Schwartz, and J. Bruck, "Noise and uncertainty in string-duplication systems," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Aachen, Germany, Jun. 2017, pp. 3120–3124.
- [28] N. Alon, J. Bruck, F. Farnoud, and S. Jain, "Duplication distance to the root for binary sequences," *IEEE Trans. Inf. Theory*, vol. 63, no. 12, pp. 7793–7803, Dec. 2017.
- [29] F. Farnoud, M. Schwartz, and J. Bruck, "Estimation of duplication history under a stochastic model for tandem repeats," *BMC Bioinf.*, vol. 20, no. 1, pp. 64–74, Feb. 2019.
- [30] O. Milenkovic and N. Kashyap, "On the design of codes for DNA computing," in *Proc. Int. Workshop on Coding and Cryptography (WCC)*, 2005, Bergen, Norway, ser. Lecture Notes in Computer Science, Ø. Ytrehus, Ed., vol. 3969. Berlin, Heidelberg: Springer-Verlag Berlin Heidelberg, 2006, pp. 100–119.
- [31] —, "DNA codes that avoid secondary structures," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Adelaide, SA, Australia, Sep. 2005, pp. 288–292.
- [32] K. G. Benerjee and A. Banerjee, "On homopolymers and secondary structures avoiding, reversible, reversible-complement and GC-balanced DNA codes," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Espoo, Finland, Jun. 2022, pp. 204–209.
- [33] T. T. Nguyen, K. Cai, H. M. Kiah, D. T. Dao, and K. A. Schouhamer Immink, "On the design of codes for DNA computing: Secondary structure avoidance codes," *arXiv preprint arXiv:2302.13714v1*, Feb. 2023. [Online]. Available: <https://arxiv.org/abs/2302.13714v1>
- [34] D. Bar-Lev, A. Kobovich, O. Leitersdorf, and E. Yaakobi, "Universal framework for parametric constrained coding," *arXiv preprint arXiv:2212.09314v1*, Apr. 2023. [Online]. Available: <https://arxiv.org/abs/2304.01317v1>
- [35] J. Spencer, "Asymptotic lower bounds for Ramsey functions," *Discrete Math.*, vol. 20, pp. 69–76, 1977.
- [36] R. M. Roth, *Introduction to Coding Theory*. Cambridge Univ. Press, 2006.
- [37] T. M. Cover, "Enumerative source encoding," *IEEE Trans. Inf. Theory*, vol. 19, no. 1, pp. 73–77, Jan. 1973.
- [38] Y. Yehezkeally, D. Bar-Lev, S. Marcovich, and E. Yaakobi, "Generalized unique reconstruction from substrings," *IEEE Trans. Inf. Theory*, vol. 69, no. 9, pp. 5648–5659, Sep. 2023.
- [39] M. Levy and E. Yaakobi, "Mutually uncorrelated codes for DNA storage," *IEEE Trans. Inf. Theory*, vol. 65, no. 6, pp. 3671–3691, Jun. 2019.
- [40] R. Gabrys and O. Milenkovic, "Unique reconstruction of coded sequences from multiset substring spectra," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Vail, CO, US, Jun. 2018, pp. 2540–2544.
- [41] F. J. MacWilliams and N. J. A. Sloane, *The Theory of Error-Correcting Codes*. North-Holland, 1978.