

Efficient Distributed Source Coding of Fragmented Genomic Sequencing Data

Yotam Gershon Yuval Cassuto

The Viterby Faculty of Electrical Engineering, Technion - Israel Institute of Technology

Email: {yotamgr@campus, ycassuto@ee}.technion.ac.il

Abstract—In this paper we present a new compression scheme for genomic read data produced by modern sequencing technologies. In this setting, a reference genome similar to the one being sequenced is available only at the decoder, while the starting index of each read in this reference is unknown. The proposed scheme significantly reduces the encoding complexity relative to known reference-based compression schemes. The results include a code construction based on generalized concatenation coset codes, analysis of the decoding failure probability, and optimization of the scheme parameters for minimal compression rate.

I. INTRODUCTION

Genomic sequencing is a process in which the order of nucleic acids within a DNA molecule (genomic sequence) is analyzed. In most modern sequencing technologies, a large set of short sequence fragments, called *reads*, is produced and represented by a string of characters (usually *A,C,G,T*). In this method, called *shotgun sequencing* [1], each read's location within the sequence is generally unknown. Furthermore, the sequencing machine introduces mutation errors into the reads. Therefore, the sequence assembly from the reads requires a large number of reads and high computational effort. Effective compression of pre-assembly reads data is therefore an essential problem. A large number of read-compression methods are available [2]–[4], partitioned into two main categories: reference-free and reference-based tools, differing by whether a closely similar reference genome is shared by the encoder and decoder. In some applications, genome reads will be produced at an edge node, e.g. a physician's office, and then sent for processing to a central node, e.g. a cloud database. In such scenarios, the cost of known reference-based methods may be too high for the edge node's limited resources, due to the need to store long references and perform the reads' alignment on them. Alternatively, in this paper we propose a compression scheme in which the encoder needs to neither store nor process any reference, while benefiting from a reference available at the central node decoding the reads.

Compression with reference available only at the decoder is an instance of the well known Slepian-Wolf (SW) coding problem [5]. Several works have addressed the problem with explicit constructive codes [6]–[12], some specifically for genomic data [13], [14], but focused on either full-block or streamlined data, and not fragmented reads. When compressing fragmented reads, the known capability of recovering a read from a similar reference is not sufficient, and additional information is needed for aligning the read within the full reference. The construction proposed in this paper offers both

capabilities, using the framework of generalized error-locating (GEL) [15] *coset codes*. GEL coset codes were used in [12] for full-block compression, and the novelty of this present work is in using the GEL's hierarchy of inner codes to combine all the information for read recovery into the same codeword: one layer for alignment, one for similarity reconstruction and one for alignment validation. In the fourth layer, a batch of reads is encoded with an outer code to provide extremely low failure probability. In addition to proposing the code construction in Section III, in Section IV we analyze the success probability of the code, taking into account the effect of read misalignments within the reference. In Section V, we seek the optimization of the various code parameters to minimize the compression rate for a specified success probability. The similarity measure considered in this paper is the rate of *substitutions*, which are the most dominant mutations observed in common sequencing technologies. Extending the scheme to deal with insertions and deletions is left as an interesting future work.

II. PROBLEM FORMULATION

A. Problem Setting

A genome data being sequenced in an edge node is represented by a string $\mathbf{X} \in \mathcal{A}^L$. A closely similar reference of this data, represented by \mathbf{Y} , is stored in a central node, and is unavailable at the edge node. The differences between \mathbf{X} and \mathbf{Y} are assumed to be caused by genetic diversity [16]. A sequencing machine at the edge node is generating a set of M reads from \mathbf{X} , denoted by $\{\mathbf{x}^{(i)}\}_{i=1}^M$. Each read $\mathbf{x}^{(i)} \in \mathcal{A}^n$ is an *approximate substring* of \mathbf{X} , taken from an unknown random location within \mathbf{X} , and introduced with *sequencing errors*. Our goal is to encode the reads data $\{\mathbf{x}^{(i)}\}_{i=1}^M$ to a minimal size, such that the reads can be recovered from the similar reference \mathbf{Y} without loss. The encoder output is communicated to the central node without errors.

B. Error Channel Model

In this work, we assume that both genomic diversity and sequencing errors are modeled by substitutions only, without insertions and deletions (indel errors), but supporting indels is possible with modification of the alignment procedure at the central node. A q -ary, length- L *substitution channel* with substitution probability p_1 , between input and output strings $\mathbf{X} = X_1, \dots, X_L$ and $\mathbf{Y} = Y_1, \dots, Y_L$, is defined as a memoryless q -ary symmetric channel, that is:

$$P(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^L P(Y_j|X_j), \quad P(Y_j|X_j) = \begin{cases} 1 - p_1 & Y_j = X_j \\ p_1/(q-1) & Y_j \neq X_j \end{cases}$$

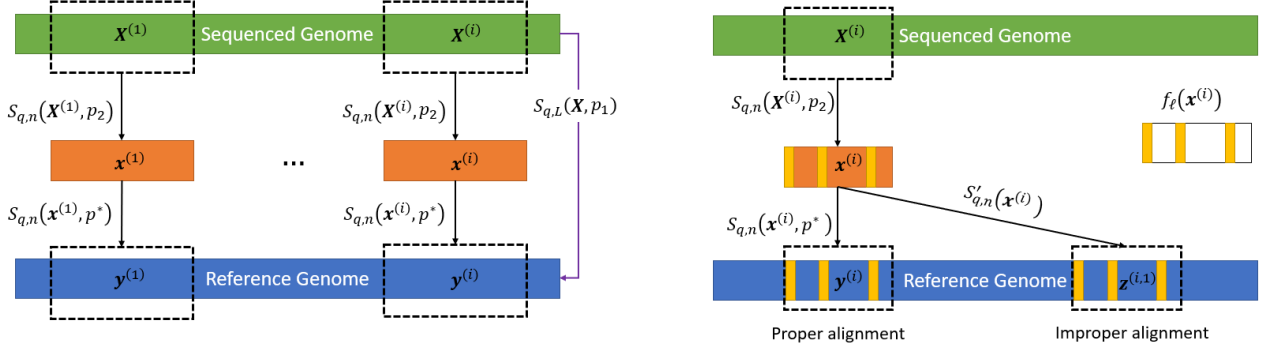


Fig. 1: **Left:** DNA reads generated as fragments of a genome. **Right:** read-identifier-based alignment of a single read.

The channel output string will be denoted by $\mathbf{Y} = S_{q,L}(\mathbf{X}, p_1)$. Let $\mathbf{X}^{(i)} = X_{k_i}, X_{k_i+1}, \dots, X_{k_i+n-1}$ be a substring observed by the sequencer, where k_i is a random, unknown *starting index*. The corresponding read $\mathbf{x}^{(i)}$ is modeled by $\mathbf{x}^{(i)} = S_{q,n}(\mathbf{X}^{(i)}, p_2)$. Let $\mathbf{y}^{(i)} = Y_{k_i}, Y_{k_i+1}, \dots, Y_{k_i+n-1}$ be the substring of \mathbf{Y} with *proper alignment* to the read $\mathbf{x}^{(i)}$. Based on the models above, we get that $\mathbf{y}^{(i)} = S_{q,n}(\mathbf{x}^{(i)}, p^*)$, where $p^* \triangleq p_1 \cdot [1 - p_2/(q-1)] + (1-p_1) \cdot p_2$. That is, the differences between reads and their references are modeled as a single substitution channel with parameter p^* . This model can be extended to richer models (e.g. [17]), potentially exploiting the dependence between two overlapping reads, whose references pass through the same instantiation of the channel $S_{q,L}(\mathbf{X}, p_1)$. The relations between the genome, reference and reads are illustrated in Fig 1.

C. Pre-Decoding Alignment

We wish to encode $\mathbf{x}^{(i)}$ such that only the information required to reconstruct it from $\mathbf{y}^{(i)}$ is transmitted. Nevertheless, for this reconstruction to work, the decoder first needs to know the starting index k_i of this read. Therefore, the encoder must transmit some additional information enabling the decoder to align the read within the reference, while accounting for the substitution errors. This alignment information is an ℓ -bit *read identifier*, output from a function denoted by $f_\ell(\mathbf{x}^{(i)})$. Such identifier and the tradeoffs in setting the value ℓ are discussed in Sections III and IV, respectively.

Since only partial information is provided for alignment, additional *improper alignments*, i.e., erroneous starting indices, are likely to be found. This process provides the decoder with a set $\{\mathbf{z}^{(i,j)} | j = 1, 2, \dots\}$ of length- n substrings of \mathbf{Y} as candidates for read alignment. Every $\mathbf{z}^{(i,j)} \neq \mathbf{y}^{(i)}$ can be regarded as having been obtained from $\mathbf{x}^{(i)}$ passing through a useless channel with zero mutual information. The alignment process is illustrated in Fig 1 (right). Clearly, only the proper alignment is desired for decoding, thus a method for rejecting false candidates is required. This method, described in Section III, is referred to as *validation*.

III. CODE CONSTRUCTION

For simplicity, throughout this section we assume $q = 2$, but the construction can be generalized to any q that is a prime power. Furthermore, cyclic indices will be used, i.e., every

index j in \mathbf{X}, \mathbf{Y} will be taken as $([(j-1) \bmod L] + 1) \in \{1, \dots, L\}$, to avoid edge effects.

A. Read Identifier

A simple *bit sampling* approach is found to be very suitable for read identifiers. Let $1 \leq i_1 < \dots < i_\ell \leq n$ be a predefined set of indices, known to both the encoder and decoder. Now, let $f_\ell(\mathbf{x}^{(i)}) = x_{i_1}, \dots, x_{i_\ell}$ be the read identifier. In this case, an *aligner* at the central node simply correlates this identifier along the reference by evaluating the Hamming distance with respect to each starting index, and produces the set

$$\mathcal{Z}^{(i)} = \left\{ \mathbf{z}^{(i,j)} \mid d_H \left(f_\ell(\mathbf{x}^{(i)}), f_\ell(\mathbf{z}^{(i,j)}) \right) \leq \mathsf{T} \right\}_{j=1}^{K_i}, \quad (1)$$

where $d_H(\cdot, \cdot)$ is the Hamming distance, T is a predefined threshold, and $\left\{ \mathbf{z}^{(i,j)} = [Y_{k_i^{(j)}}, \dots, Y_{k_i^{(j)}+n-1}] \right\}$ is the set of possible alignments of $\mathbf{x}^{(i)}$ within \mathbf{Y} . The remainder of the read, i.e., the indices outside of the identifier, is denoted by $\mathbf{x}_{\mathcal{I}}^{(i)}$, where $\mathcal{I} = \{1, \dots, n\} \setminus \{i_1, \dots, i_\ell\}$, $|\mathcal{I}| = n - \ell$.

B. General Code Construction

Our goal is to design a coding scheme for transmitting reads from $\mathbf{X} \in \{0, 1\}^L$ such that a decoder with access to $\mathbf{Y} = S_{2,L}(\mathbf{X}, p^*)$ will be able to perfectly reconstruct them with high probability.

Definition 1. A $(M, n, \mathcal{R}, p^*, P_s)$ -code is a pair $(\mathcal{E}, \mathcal{D})$ of encoder-decoder for a set $\{\mathbf{x}^{(i)}\}_{i=1}^M$ of length- n reads such that:

- 1) \mathcal{E}, \mathcal{D} have access only to $\{\mathbf{x}^{(i)}\}_{i=1}^M, \mathbf{Y}$, respectively,
- 2) the encoded size satisfies

$$|\mathcal{E}(\{\mathbf{x}^{(i)}\}_{i=1}^M)| = (nM) \cdot \mathcal{R},$$

- 3) the correct decoding probability satisfies

$$\Pr \left\{ \mathcal{D} \left[\mathcal{E}(\{\mathbf{x}^{(i)}\}_{i=1}^M), \mathbf{Y} \right] = \{\mathbf{x}^{(i)}\}_{i=1}^M \right\} \geq P_s.$$

Our general code construction is based on generalized error locating (GEL) codes [15], adapted to use as a source code with alignment-validation capabilities.

Construction 1. Let $\mathcal{C}_1, \mathcal{C}_2$ be a pair of binary linear codes with parameters $[n - \ell, k_i - \ell, d_i]$, $i = \{1, 2\}$, where $k_1 \geq k_2$. Let $\mathbf{H}_1, \mathbf{H}_2$ be parity-check matrices of these codes, respectively, such that they form a nested pair, i.e., all rows of \mathbf{H}_1

appear in H_2 in concatenation with additional $\tau \triangleq k_1 - k_2$ rows denoted by \bar{H}_2 , the *validation matrix*. Let H_c be a matrix such that the concatenation of its rows with H_2 forms a square *full-rank* matrix H . This structure is illustrated in Fig. 2.

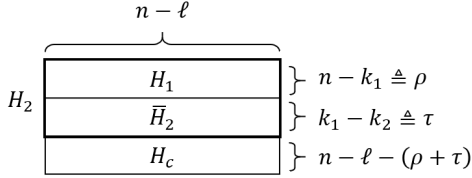


Fig. 2: Structure and sizes of the construction's inner-code parity-check matrices.

Finally, let \mathcal{C}_o be a $[M, k_o, d_o]$ linear code over $\text{GF}(2^\nu)$, with parity-check matrix \mathbf{H}_o , and where $\nu = n - \ell - (\rho + \tau)$. This code will be referred to as the *outer code*.

The encoding and decoding processes are introduced in Algorithms 1 and 2, respectively. We note that the syndrome $\mathbf{s}^{(i)}$ with respect to H_2 , calculated in Algorithm 1, is of the form $\mathbf{s}^{(i)} = [\mathbf{s}_1^{(i)}, \mathbf{s}_2^{(i)}]$, where $\mathbf{s}_1^{(i)}, \mathbf{s}_2^{(i)}$ correspond to H_1, \bar{H}_2 , respectively. In Algorithm 2, we denote by $\mathcal{D}_1(\mathbf{z}, \mathbf{s})$ the result of decoding the word \mathbf{z} with respect to H_1 within the coset of syndrome \mathbf{s} . The same is done for $\mathcal{D}_o(\mathbf{a}, \mathbf{S})$, with a word \mathbf{a} decoded with respect to \mathbf{H}_o to a syndrome \mathbf{S} . We denote by $\mathcal{F}_H(\mathbf{u})$ the linear mapping of \mathbf{u} to the single codeword of syndrome \mathbf{u} in the code defined by H . Finally, we denote by \otimes an erasure occurring if either more than one alignment was validated or all of them failed to be validated.

Algorithm 1: Construction 1 Encoding

Input: $\{\mathbf{x}^{(i)}\}_{i=1}^M, H_2, H_c, \mathbf{H}_o$
for $1 \leq i \leq M$ **do** // Inner Encoding
 Extract $\mathbf{w}^{(i)} = f_\ell(\mathbf{x}^{(i)})$
 Calculate $\mathbf{s}^{(i)} = H_2 [\mathbf{x}_T^{(i)}]^T, \mathbf{a}^{(i)} = H_c [\mathbf{x}_T^{(i)}]^T$
end
Form $\mathbf{a} = [\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(M)}] \in [\text{GF}(2^\nu)]^M$
Calculate $\mathbf{S} = \mathbf{H}_o \mathbf{a}^T$ // Outer Encoding
Output: $\mathcal{E}(\{\mathbf{x}^{(i)}\}_{i=1}^M) = \{\{\mathbf{w}^{(i)}\}_{i=1}^M, \{\mathbf{s}^{(i)}\}_{i=1}^M, \mathbf{S}\}$

Note that the term *validation* in Algorithm 2 and in the sequel includes cases of *misvalidation*, that is, validation of a vector $\mathbf{v} \neq \mathbf{x}_T^{(i)}$, whereas a failed validation, i.e. $\hat{\mathbf{s}}_2^{(i)} \neq \mathbf{s}_2^{(i)}$, forms a *rejection* of the candidate.

Proposition 2. *The rate of Construction 1 is given by*

$$\mathcal{R} = 1 - \frac{k_o}{M} \cdot \frac{n - \ell - (\rho + \tau)}{n} = 1 - \frac{k_o}{M} \cdot \frac{k_2 - \ell}{n}.$$

Proposition 3. *Construction 1 yields a $(M, n, \mathcal{R}, p^*, P_s)$ -code if and only if at the outer decoder output $\Pr\{\hat{\mathbf{a}} = \mathbf{a}\} \geq P_s$.*

By Proposition 3, the success of the scheme depends on the success in decoding the outer code over a channel induced by the inner-decoding outcomes. In the next section we turn to analyze this induced channel, which requires accounting for both the alignment and the inner-decoder performance.

Algorithm 2: Construction 1 Decoding

Input: $\mathcal{E}(\{\mathbf{x}^{(i)}\}_{i=1}^M), \mathbf{Y}, H_1, \bar{H}_2, H_c, \mathbf{H}_o$
for $1 \leq i \leq M$ **do**
 Align $\mathbf{w}^{(i)}$ over \mathbf{Y} , and form $\mathcal{Z}^{(i)}$ (Eq. 1)
 for $1 \leq j \leq |\mathcal{Z}^{(i)}|$ **do** // Inner Decoding
 Set 'found' $\leftarrow 0$
 Decode $\mathbf{v} = \mathcal{D}_1(\mathbf{z}_T^{(i,j)}, \mathbf{s}_1^{(i)})$
 Calculate $\hat{\mathbf{s}}_2^{(i)} = \bar{H}_2 \mathbf{v}^T$
 if $\hat{\mathbf{s}}_2^{(i)} = \mathbf{s}_2^{(i)}$ **then** // Validation
 if 'found' = 0 **then**
 Calculate $\mathbf{b}^{(i)} = H_c \mathbf{v}^T$, Set 'found' $\leftarrow 1$
 else // More Than One Candidate
 Set $\mathbf{b}^{(i)} = \otimes$, break
 end
 end
 end
 if 'found' = 0 **then** // No Candidates
 Set $\mathbf{b}^{(i)} = \otimes$
 end
end
// Outer Decoding
Decode $\hat{\mathbf{a}} = \mathcal{D}_o(\mathbf{b}, \mathbf{S})$, where $\mathbf{b} = [\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(M)}]$
for $1 \leq i \leq M$ **do** // Inverse Mapping
 Map $\hat{\mathbf{x}}_T^{(i)} = \mathcal{F}_H([\mathbf{s}^{(i)}, \hat{\mathbf{a}}^{(i)}])$
 Reconstruct $\hat{\mathbf{x}}^{(i)}$ from $\hat{\mathbf{x}}_T^{(i)}, \mathbf{w}^{(i)}$
end
Output: $\{\hat{\mathbf{x}}^{(i)}\}_{i=1}^M$

IV. SCHEME ANALYSIS

Throughout this section, the sequenced genome is assumed to be a random sequence in which each symbol is drawn i.i.d. from the Bernoulli(1/2) distribution. This assumption is only needed for analyzing the misalignments of the read $\mathbf{x}^{(i)}$, and the scheme's handling of the proper alignment does not rely on it. Extension of the analysis to other genome statistical models can be done based on related studies such as [1], [18].

A. Inner-Code Analysis

The inner decoder is invoked in Algorithm 2 on both the proper alignment of $\mathbf{x}^{(i)}$ in \mathbf{Y} (if found) and, possibly, on improperly aligned $\mathbf{z}^{(i,j)}$ vectors that are not related to $\mathbf{x}^{(i)}$. The scheme's performance depends on the inner-decoding outcomes for both types of inputs, which we now analyze.

Definition 4. *Let $\mathcal{Z}^{(i)}$ be the set of possible alignments of some read $\mathbf{x}_T^{(i)}$. The following terms are defined:*

- $P_{a,s}$ - the probability of the proper alignment being found, i.e., $\mathbf{y}^{(i)} \in \mathcal{Z}^{(i)}$.
- $K_f^{(i)} = |\mathcal{Z}^{(i)} \setminus \{\mathbf{y}^{(i)}\}|$ - the number of improper alignments.

Let $F_b(n, p, t)$ denote the CDF of a binomial random variable with parameters (n, p) , evaluated at the value of t .

Lemma 5. *The following expressions are immediate:*

$$P_{a,s} = F_b(\ell, p^*, T), \quad \mathbb{E}[K_f^{(i)}] \triangleq \bar{K}_f = (L - 1) \cdot F_b(\ell, 1/2, T).$$

Definition 6. For each alignment candidate $z^{(i,j)} \in \mathcal{Z}^{(i)}$ of some read $x_{\mathcal{I}}^{(i)}$, the following probabilities are defined:

- P_{suc} - for successful inner decoding, that is, $v = x_{\mathcal{I}}^{(i)}$.
- P_{miv} - for misvalidated inner decoding, that is, \mathcal{D}_1 returns $v \neq x_{\mathcal{I}}^{(i)}$, and $\hat{s}_2^{(i)} = s_2^{(i)}$.
- $P_{\text{rej}} = 1 - P_{\text{suc}} - P_{\text{miv}}$ - for detected inner-decoding error.

We further add a superscript $(x) \in \{(p), (i.p)\}$ to the above probabilities, corresponding to whether $z^{(i,j)}$ is the proper or improper alignment, respectively. Let $t_1 \triangleq \lfloor \frac{d_1-1}{2} \rfloor$, $n_\ell \triangleq n - \ell$, and $V_n(t)$ denote the volume of a Hamming ball with radius t of length- n words. Then we have the following.

Lemma 7. The following equalities and approximations hold:

$$P_{\text{suc}}^{(p)} = F_b(n_\ell, p^*, t_1), \quad P_{\text{suc}}^{(i.p)} \simeq 0, \quad (2)$$

$$P_{\text{miv}}^{(i.p)} = \frac{V_{n_\ell}(t_1)}{2^{\rho+\tau}}, \quad P_{\text{miv}}^{(p)} \lesssim P_{\text{miv}}^{(i.p)}, \quad (3)$$

$$P_{\text{rej}}^{(i.p)} \simeq 1 - P_{\text{miv}}^{(i.p)}. \quad (4)$$

Proof: In (2) the equality follows from \mathcal{C}_1 having minimum distance d_1 , and the $\simeq 0$ represents the negligible probability that a $z_{\mathcal{I}}^{(i,j)}$ would somehow successfully decode to an unrelated $x_{\mathcal{I}}^{(i)}$. In the left equality of (3), since the improper alignment is not related to the encoder input $x_{\mathcal{I}}^{(i)}$, $P_{\text{miv}}^{(i.p)}$ equals the probability that a random length- n_ℓ vector has a codeword of \mathcal{C}_2 at distance at most t_1 . This can be shown to equal the right-hand side. The \lesssim inequality follows from the fact that the probable proximity of $x_{\mathcal{I}}^{(i)}$ to the proper alignment in general reduces the probability of misvalidation. For this inequality to hold formally, the inner codes need to satisfy a natural property called *properness* [19]. Finally, (4) follows from the second claim in (2) and the definition of P_{rej} . ■

B. Outer-Code Analysis

In the outer code, the syndromes with respect to H_c are treated as symbols in $\text{GF}(2^\nu)$. As seen in Algorithm 2, these symbols may be erased depending on the outcome of the inner decoding. Furthermore, a misvalidated inner-decoder output v may introduce an erroneous symbol. The outer decoded word \underline{b} can thus be modeled as transmitting the correct codeword \underline{a} through an *outer channel*, which produces erasures and errors. The probabilities of these events are directly induced by the coding scheme parameters $\{\ell, k_1, \tau, T\}$. We can now examine the erasure and error probabilities of the outer channel, denoted by $P_{\text{ers}}, P_{\text{err}}$, respectively. For the next lemma, we assume that equation (4) holds by equality.

Lemma 8. Let $P_{p.a} \triangleq P_{a.s} \cdot (1 - P_{\text{rej}}^{(p)})$, $P_{p.m} \triangleq P_{a.s} \cdot P_{\text{miv}}^{(p)}$, then the outer channel probabilities satisfy:

$$\begin{aligned} P_{\text{ers}} &\leq P_{p.a} \left[\bar{K}_f P_{\text{miv}}^{(i.p)} (1 + B_1) \right] + (1 - P_{p.a}) \left[1 - \bar{K}_f P_{\text{miv}}^{(i.p)} (1 + B_1) \right], \\ P_{\text{err}} &\leq P_{p.m} \left[1 - \bar{K}_f P_{\text{miv}}^{(i.p)} (1 + B_1) \right] + (1 - P_{p.a}) \left[\bar{K}_f P_{\text{miv}}^{(i.p)} (1 + B_1) \right], \end{aligned} \quad (5)$$

where B_1 are terms bounded from above by $\bar{K}_f P_{\text{miv}}^{(i.p)}$ when $\bar{K}_f P_{\text{miv}}^{(i.p)} < 1$, and are thus negligible when $\bar{K}_f P_{\text{miv}}^{(i.p)} \ll 1$.

Proof: (sketch): For P_{ers} the events of interest are: 1) having 1 proper validation and 1 or more improper validations (first term), and 2) having 0 proper validations and not exactly 1 improper validation (second term). For P_{err} the events of interest are: 1) having 1 proper *mis*-validation and no improper validations (first term), and 2) having 0 proper validations and 1 improper validation (second term). Each probability is expanded with the Taylor series, and B_1 encapsulates only terms of order at least 1 in $\bar{K}_f P_{\text{miv}}^{(i.p)}$. Since $K_f^{(i)}$ are themselves binomially distributed, marginalizing the probability over K_f gives the raw binomial moments, all of which are order at least 1 in $\bar{K}_f P_{\text{miv}}^{(i.p)}$ [20]. ■

We note that relieving the assumption regarding equation (4), the lemma still holds by simply replacing $P_{\text{miv}}^{(i.p)}$ terms with $1 - P_{\text{rej}}^{(i.p)}$.

It is well known that a code with minimum distance d can correct up to m_1 erasures and m_2 errors as long as $m_1 + 2m_2 \leq d - 1$ [21]. That means the redundancy of the outer code should be set by the following.

Proposition 9. Let \mathcal{C}_o be a maximum distance separable (MDS) code, and define the random variable $W = m_1 + 2m_2$. Then the minimal required redundancy $M - k_o$ of \mathcal{C}_o equals

$$\rho_o^* = \min\{w\} \text{ such that } P(W \leq w) \geq P_s,$$

where

$$P(W = u) = \sum_{\underline{m} \in S(M, u)} \frac{M!}{m_0! m_1! m_2!} (1 - P_{\text{ers}} - P_{\text{err}})^{m_0} \cdot P_{\text{ers}}^{m_1} \cdot P_{\text{err}}^{m_2},$$

for $\underline{m} = (m_0, m_1, m_2)$ and

$$S(M, u) \triangleq \{\underline{m} \mid m_1 + 2m_2 = u, m_0 + m_1 + m_2 = M\}.$$

V. PARAMETER OPTIMIZATION

Given the specified system parameters M, n, L, p^*, P_s , to obtain minimal-rate compression by Construction 1 we need to optimize the free parameters of the construction. More concretely, the optimal parameters are defined in the following.

Definition 10. The optimal parameters of a $(M, n, \mathcal{R}, p^*, P_s)$ -code from Construction 1 are

$$\{k_1, \tau, \ell, T\}^* = \arg \max_{\{k_1, \tau, \ell, T\}} \mathcal{L}(k_1, \tau, \ell, T), \quad (6)$$

where

$$\mathcal{L}(k_1, \tau, \ell, T) \triangleq (k_1 - \tau - \ell) \cdot (M - \rho_o^*). \quad (7)$$

The optimality of $\{k_1, \tau, \ell, T\}^*$ in Definition 10 follows immediately from Propositions 2 and 9. The challenge in finding the optimal values lies in that maximizing (7) needs to account for the complex dependence of ρ_o^* on the other parameters affecting the rate.

A. Choosing the Alignment Parameters

Finding the four parameters of (6) simultaneously is difficult, so we first deal with the alignment parameters ℓ, T separately. We observe that these parameters have importance beyond the above optimization: they control the expected number of improper alignments, and thus the decoding complexity (the expected number of inner-decoder invocations).

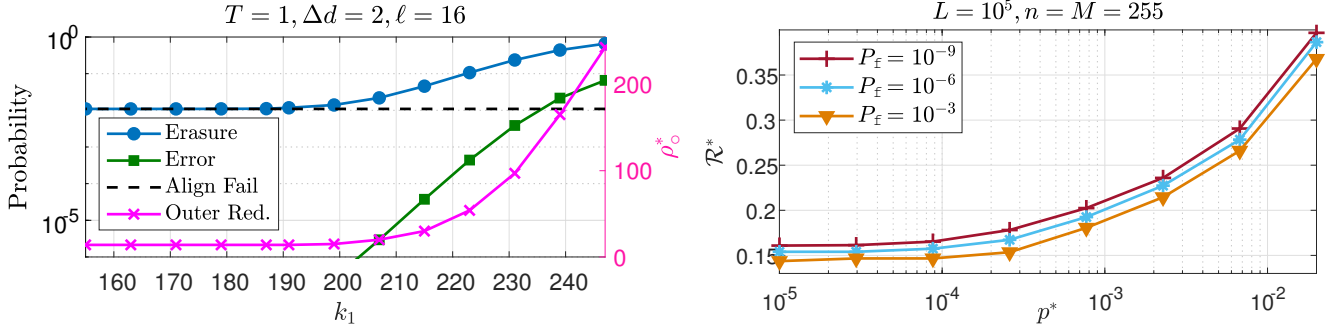


Fig. 3: Example of Construction 1 using BCH inner codes for $L = 10^5, n = M = 255$. **Left:** outer channel erasure and error probabilities and required outer redundancy symbols. **Right:** optimal rate \mathcal{R}^* as a function of p^* for three failure probabilities.

Therefore, we limit ourselves to values that maintain a limit on this expected number, as follows.

Proposition 11. *Given an upper bound K_m for the expected number of improper alignments \bar{K}_f , a valid region for threshold T values, for every read identifier length ℓ , is obtained by*

$$T_{\text{v}}(K_m, \ell) = \{T : V_{\ell}(T) \cdot 2^{-\ell} \leq K_m/L\}.$$

The valid region for ℓ is then $\ell_{\text{v}}(K_m) = \{\ell : T_{\text{v}}(K_m, \ell) \neq \emptyset\}$.

The pairs ℓ, T allowed by Proposition 11 are independent of all other (system+construction) parameters except L , and so can be calculated once and reused for every parameter setup.

B. Approximating the Outer Redundancy

For efficient evaluation of parameter sets during optimization, it is helpful to simplify the dependence of ρ_o^* on $P_{\text{ers}}, P_{\text{err}}$ relative to the exact trinomial CDF in Proposition 9. For that, we approximate the trinomial distribution of W with parameters $M, P_{\text{ers}}, P_{\text{err}}$ by a Normal distribution with the same mean and same variance, as follows.

Lemma 12. *Let $P_{\text{red}} \triangleq P_{\text{ers}} + 2P_{\text{err}}$. The outer redundancy can be approximated by*

$$\rho_o^* \approx \mu_W + Q^{-1}(P_f) \cdot \sigma_W + \beta(P_f), \quad (8)$$

where $\mu_W = MP_{\text{red}}$, $\sigma_W^2 = MP_{\text{red}}(1 - P_{\text{red}}) + 2MP_{\text{err}}$, and $Q(\cdot)$ is the Normal distribution's tail function, $P_f = 1 - P_s$ is the allowed failure probability, and $\beta(\cdot)$ is some empirical correction function of P_f .

The advantage of (8) is that ρ_o^* can be approximated in closed form, using a single fitting parameter β for each P_f , and can be used for any pair $P_{\text{ers}}, P_{\text{err}}$ induced by inner-code parameters.

C. Optimization Procedure

Finding a solution to the problem in (6) can now be done in the following stages: (i.) Given K_m , find the valid regions $\ell_{\text{v}}(K_m), T_{\text{v}}(K_m, \ell)$, (ii.) set k_1 and iterate over valid pairs (ℓ, T) , and numerically find the optimal τ^* by evaluating equations (5), (7) and (8) (neglecting B_1 terms), (iii.) find the optimal set $\{\tau, \ell, T\}^*$ for k_1 , by comparing the results of each pair (ℓ, T) , and (iv.) find the optimal set $\{k_1, \tau, \ell, T\}^*$ by comparing the results of each k_1 . Since the set of available dimensions k_1 of an error-correcting code family (e.g. BCH) is relatively small, the search space in (iv) is manageable.

Similarly, the search space of τ^* in (ii) is relatively small, since the misvalidation probability decays exponentially with the number of validation bits.

VI. NUMERICAL RESULTS

We demonstrate the optimization procedure by showing numerical results for Construction 1 employing binary BCH codes as $\mathcal{C}_1, \mathcal{C}_2$, and a Reed-Solomon code as \mathcal{C}_0 . In this setting, the parameter τ is obtained by choosing the difference in minimum distance between \mathcal{C}_1 and \mathcal{C}_2 , denoted $\Delta d \triangleq d_2 - d_1$. In Fig. 3 (left), we plot the induced outer channel probabilities (left axis) and corresponding outer redundancy (right axis) as a function of k_1 , for $L = 10^5, n = M = 255, p^* = 0.01, P_s = 1 - 10^{-6}$. The tradeoff of redundancy allocation between codes is clear by examining the (non-linear) increase in required outer redundancy with the weakening of the inner code (increase in k_1). The plot reveals a floor value of 10^{-2} for the erasure probability, attributed to failed proper alignment, which cannot be solved by strengthening the inner code. A solution of the rate optimization problem for different values of substitution rate p^* is shown for $P_f = \{10^{-9}, 10^{-6}, 10^{-3}\}$ in Fig. 3 (right). Finally, we compare in Fig. 4 the rates of the proposed scheme with a powerful fixed-rate benchmark that assumes the encoder knows both \mathbf{Y} and the correct alignment of every $\mathbf{x}^{(i)}$ in \mathbf{Y} (and can thus communicate the alignment index plus a fixed-rate encoding of the difference $\mathbf{y}^{(i)} - \mathbf{x}^{(i)}$). It is shown that for $P_f = 10^{-9}$, the proposed scheme outperforms the benchmark for every p^* starting at 10^{-4} , while saving the space and computation effort at the encoder required to store \mathbf{Y} and align the reads to it.

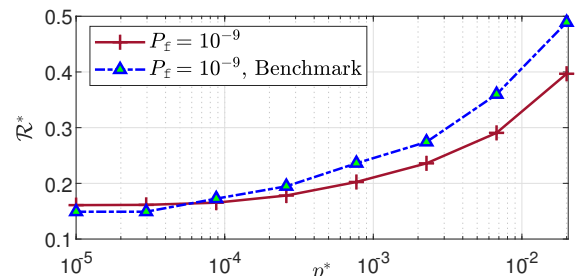


Fig. 4: Optimal (+) vs. benchmark (Δ) rates for $P_f = 10^{-9}$.

VII. ACKNOWLEDGEMENT

This work was supported in part by the US-Israel Binational Science Foundation and by the Israel Science Foundation.

REFERENCES

- [1] A. S. Motahari, G. Bresler, and D. N. C. Tse, "Information Theory of DNA Shotgun Sequencing," *IEEE Trans. on Information Theory*, Vol. 59, No. 10, 2013.
- [2] Z. Zhu, Y. Zhang, Z. Ji, S. He, and X. Yang, "High-Throughput DNA Sequence Data Compression," *Briefings In Bioinformatics*, Vol. 16, No. 1, 1-15, 2013.
- [3] N. S. Bakr and A. A. Sharawi, "DNA Lossless Compression Algorithms: Review," *American Journal of Bioinformatics Research*, 3(3), 2013.
- [4] M. Hosseini, D. Pratas, and A. J. Pinho, "A Survey on Data Compression Methods for Biological Sequences," *Information* 7, 56, 2016.
- [5] D. Slepian and J. K. Wolf, "Noiseless Coding of Correlated Information Sources," *IEEE Trans. on Information Theory*, Vol. IT-19, No. 4, 1973.
- [6] S. S. Pradhan and K. Ramchandran, "Distributed Source Coding Using Syndromes (DISCUS): Design and Construction," *IEEE Trans. on Information Theory*, Vol. 49, No. 3, 2003.
- [7] T. Uyematsu, "An Algebraic Construction of Codes for Slepian-Wolf Source Networks," *IEEE Trans. on Information Theory*, Vol. 47, No. 7, 2001.
- [8] A. Orlicsky and K. Viswanathan, "One-Way Communication and Error-Correcting Codes," *IEEE Trans. on Information Theory*, Vol. 49, No. 7, 2003.
- [9] Y. Minsky, A. Trachtenberg, and R. Zippel, "Set Reconciliation with Nearly Optimal Communication Complexity," *IEEE Trans. on Information Theory*, Vol. 49, No. 9, 2003.
- [10] A. Aaron and B. Girod, "Compression with Side Information Using Turbo Codes," *IEEE Proceedings DCC 2002. Data Compression Conference*, 2002.
- [11] M. Sartipi and F. Fekri, "Distributed Source Coding in Wireless Sensor Networks using LDPC coding: The entire Slepian-Wolf Rate Region," *IEEE Wireless Communications and Networking Conference*, 2005.
- [12] Y. Cassuto and J. Ziv, "Efficient Compression of Long Arbitrary Sequences with No Reference at the Encoder," *IEEE Trans. on Information Theory*, Vol. 67, No. 1, 2021.
- [13] S. Wang, X. Jiang, F. Chen, L. Cui, and S. Cheng, "Streamlined Genome Sequence Compression using Distributed Source Coding," *Cancer Informatics. Supplementary Issue: Computational Advances in Cancer Informatics*, 2014.
- [14] J. J. Shao, "Genome Sequence Compression Algorithm Based on the Distributed Source Coding," *6th International Conference on Machinery, Materials, Environment, Biotechnology and Computer (MMEBC)*, 2016.
- [15] M. Bossert, "Channel Coding for Telecommunications," Wiley, 1999.
- [16] A. F. Wright, "Genetic Variation: Polymorphisms and Mutations," *Encyclopedia of Life Sciences*, John Wiley and Sons, 2005.
- [17] C. Lottaz, C. Iseli, C. Jongeneel, and P. Bucher, "Modeling Sequencing Errors by Combining Hidden Markov Models," *Bioinformatics*, Vol. 19, Suppl. 2, 2003.
- [18] C. Kusters and T. Ignatenko, "DNA Sequence Modeling Based On Context Trees," *Proc. of the 36th WIC Symposium on Information Theory*, 2015.
- [19] S. K. Leung-Yan-Cheong, E. R. Barnes, and D. U. Friedman, "On Some Properties of the Undetected Error Probability of Linear Codes," *IEEE Trans. on Information Theory*, Vol. IT-25, No. 1, 1979.
- [20] A. Knoblauch, "Closed-Form Expression for the Moments of the Binomial Probability Distribution," *SIAM Journal on Applied Mathematics*, Vol. 69, No. 1, 2008.
- [21] R. M. Roth, "Introduction To Coding Theory," Cambridge University Press, 2006.