

• Supplementary File •

Error Exponent for Concatenated Codes in DNA Data Storage under Substitution Errors

Yuxuan SHI¹, Shuo SHAO¹, Xiaohang ZHANG^{2*}, Yongjian WANG² & Yongpeng WU³

¹ School of Cyber and Engineering, Shanghai Jiao Tong University, Shanghai 200240, China ;

² The National Computer Network Emergency Response Technical Team. Coordination Center of China, Beijing 100029, China;

³Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai 200240, China

Appendix A Symbols Explanation

Symbols	Definition
\mathbf{X}_i	Input sequence
\mathcal{X}^M	Set of input sequences
\mathbf{Y}_j	Output sequence
\mathcal{Y}^N	Set of output sequences
I_j	Cluster label
Z_i	Cluster which contains the sequences from the input X_i
E_j	Error sequence
M	Cardinality of the input sequences/Outer code length
L	Length of the input sequences/Inner code length
N	Cardinality of the output sequences
c	Average times that each sequence has been drawn $c = \frac{N}{M}$
p	Crossover probability of BSC
\hat{X}_i	Output sequence after bit-wised majority decision
D_i	number of times the i -th input sequence has been drawn (realization d)
U_d	number of input sequences has been drawn same d times (realization u)
β	rate loss due to index coding
C_d	Capacity of the sub-channel has been drawn d times
C	Capacity of the overall channel
R	Outer code rate
r	Inner code rate
\mathcal{R}	Overall rate
K	Dimension of outer code
k	Dimension of inner code
d_{out}	Minimum distance of outer code
d_{in}	Minimum distance of inner code
θ	Pre-decoding parameters
$E_G(p, r)$	Gallager error Exponent
η_i	Indicator of the sub-block
δ	The relative minimum distance ($\delta = \frac{d_{\text{out}}}{M}$)
\mathcal{J}	Set of the subscript after pre-decoding

Appendix B Important Definitions and Lemmas

We bound the error probability when both erasure and error probabilities converge to 0, in which the overall rate is shown as a function of inner code rate and coverage depth. We classify the decoding failures into erasures and errors using the specific decoder, where the error exponent is restricted by the Chernoff bound on independent indicators. The necessary definitions and lemmas are presented as follows.

* Corresponding author (email: zhangxiaohang@cert.org.cn)

i	η^M	\mathbf{d}^M	\Rightarrow	d	$\sum \eta_i$	\mathbf{u}^N
1	η_1	1		0	η_9	1
2	η_2	3		1	$\eta_1 + \eta_8$	2
3	η_3	3		2	η_7	1
4	η_4	4		3	$\eta_2 + \eta_3 + \eta_5$	3
5	η_5	3		4	$\eta_4 + \eta_6$	2
6	η_6	4		5	0	0
7	η_7	2		6	0	0
8	η_8	1		\vdots	\vdots	\vdots
9	η_9	0		21	0	0

Table B1 Regroup of the binary indicators in terms of d for $M = 9$ and $N = 21$

Definition 1. Given the input strands $\{\mathbf{X}_i\}_1^M$ and the output $\{\mathbf{Y}_j\}_1^N$, random variable D_i refers to the numbers of the i -th strands has been drawn, with realization d . Moreover the number of input strands which have been drawn exactly d times, is defined as random variable U_d with realization u

$$D_i = |\{j: I_j = i, i \in [1, M]\}| \quad \mathbf{D}^M = \{D_1, D_2, \dots, D_M\}$$

$$U_d = |\{i: D_i = d, d \in [0, N]\}| \quad \mathbf{U}^N = \{U_0, U_1, \dots, U_N\}$$

Apparently $\sum_{i=1}^M D_i = N$ and $\sum_{d=0}^N U_d = M$.

Definition 2. Given $\hat{\mathbf{X}}_i$ and \mathbf{X}_i . Define such a indicator η_i implies whether the i -th sub-block is decoded successfully or not, i.e.

$$\eta_i = \begin{cases} 1 & \mathbb{D}_{\text{in}}(\hat{\mathbf{X}}_i) \neq \mathbf{X}_i \\ 0 & \mathbb{D}_{\text{in}}(\hat{\mathbf{X}}_i) = \mathbf{X}_i \end{cases}$$

where \mathbb{D}_{in} refers to the inner decoder.

Lemma 1. With notation as above, fix $0 < \delta < 1$ and $\mathcal{J} = \{i: |\mathbf{Z}_i| \geq \theta\}$. Then the error probability of outer code could be bounded by

$$P_w = \sum_{\mathbf{u}^N} \mathbb{P} \left\{ \sum_{i \in \mathcal{J}} \eta_i \geq M\delta \mid \mathbf{D}^M = \mathbf{d}_{\mathbf{u}}^M \right\} \cdot \mathbb{P}\{\mathbf{U}^N = \mathbf{u}^N\}, \quad (\text{B1})$$

where $\mathbf{d}_{\mathbf{u}}^M \in \mathcal{D}^M(\mathbf{u}^N) = \{\mathbf{d}^M: \Pi(\mathbf{d}^M) = \mathbf{u}^N\}$ and $\Pi(\cdot)$ is the function which is used to map vector \mathbf{d}^M to \mathbf{u}^N . The function $\Pi(\cdot)$ and the regroup of the indicators η^M is illustrated in Table.B1

Proof. We abbreviate $\mathbb{P}\{\mathbf{D}^M = \mathbf{d}^M\} = \mathbb{P}\{\mathbf{d}^M\}$, the same as \mathbf{u}^N . The error probability can be reformulated as follows:

$$\begin{aligned} P_w &= \sum_{\mathbf{d}^M} \mathbb{P} \left\{ \sum_{i \in \mathcal{J}} \eta_i \geq M\delta \mid \mathbf{d}^M \right\} \mathbb{P}\{\mathbf{d}^M\} \\ &= \sum_{\mathbf{u}^N} \sum_{\mathbf{d}^M \in \mathcal{D}^M} \underbrace{\mathbb{P} \left\{ \sum_{i \in \mathcal{J}} \eta_i \geq M\delta \mid \mathbf{d}_{\mathbf{u}}^M \right\}}_{\mathcal{F}(\mathbf{u}^N)} \mathbb{P}\{\mathbf{d}^M\} \\ &\stackrel{(a)}{=} \sum_{\mathbf{u}^N} \mathcal{F}(\mathbf{u}^N) \sum_{\mathbf{d}^M \in \mathcal{D}^M} \mathbb{P}\{\mathbf{d}^M\} \\ &= \sum_{\mathbf{u}^N} \mathcal{F}(\mathbf{u}^N) \mathbb{P}\{\mathbf{u}^N\} \end{aligned}$$

(a) holds since $\mathbb{P} \left\{ \sum_{i \in \mathcal{J}} \eta_i \geq M\delta \mid \mathbf{d}_{\mathbf{u}}^M \right\} = \mathcal{F}(\mathbf{u}^N)$ is a function of \mathbf{u}^N rather than \mathbf{d}^M .

Remark 1. The specific set $\mathcal{D}^M(\mathbf{u}^N)$, contains all \mathbf{d}^M vectors which could be mapped to \mathbf{u}^N vector. Given the condition that \mathbf{D}^M is in this set, the error probability could be represented by a function depending on \mathbf{u}^N rather than \mathbf{d}^M . This representation is attributed to that the conditional probability relies on the sum of binary indicators but not care about the permutation. It avoids the cumbersome distribution without losing the independence on these indicators.

Lemma 2 (Forney's Exponent [1]). There exists a concatenated code $\mathcal{C}_{\text{cont}}$ whose outer code is a MDS code, with overall length M_0 and overall rate \mathcal{R} . Consider parallel BSC sub-channels with the same flip rate p , its probability of decoding error $P_{\text{err}}(\mathcal{C}_{\text{cont}})$ is given by

$$P_{\text{err}}(\mathcal{C}_{\text{cont}}) \leq \exp\{-M_0 E_C(p, \mathcal{R})\} \quad (0 \leq \mathcal{R} < 1 - H_q(p)) \quad (\text{B2})$$

where

$$E_F(p, \mathcal{R}) = \max_{0 < r < 1 - H_q(p)} \left(1 - \frac{\mathcal{R}}{r}\right) E_G(p, r) \quad (\text{B3})$$

which is called the concatenation exponent or Forney's exponent. $E_G(p, r)$ is the Gallager's random coding error exponent. The overall decoding complexity for the code \mathcal{C} is given by at most

$$O\left(M_0^2 \log^2 M_0\right)$$

due to the fact that the outer code (which is always a Reed-Solomon code) used in [?] has a high decoding complexity.

Lemma 3 (Poissonization [2]). For any fixed $c > 0$ and the poisson distribution $p_c(d) = e^{-c} c^d / d!$ there exists a set of functions $f_d(M), d = 0, \dots, N$ with $f_0(M) + f_1(M) + \dots + f_N(M) = o(M)$ such that for $M \rightarrow \infty$

$$\mathbb{P}(|U_d - Mp_c(d)| \leq f_d(M), d = 0, \dots, N) \rightarrow 1 \quad (\text{B4})$$

Remark 2. This lemma implies the deviation between a random variable U_d and the expectation $Mp_c(d)$ goes to $o(M)$ as $M \rightarrow \infty$. It means in this case $U_d \approx Mp_c(d)$ with large enough M . This effect is known as Poissonization although D_i are statistically dependent and each variable U_d could be viewed as no longer random but almost deterministic. Moreover, the bounding functions $f_d(M)$ are used to describe the "almost".

Appendix C Proof of Theorem 1

Proof. We start the proof by counting the number of erasures and errors during the transmission. Erasures come from the reject to the clusters of small sizes $C_d < r$, since they are not able to transfer reliable information from the perspective of decoder. Consider the unique decoding, erasures consume the same amount of parity check of outer code. Moreover, errors occur when the left redundancy cannot afford the error correction. These two random variables can be given as:

$$v = \sum_{d < \theta} U_d; \quad w = \sum_{i \in \mathcal{J}} \eta_i$$

The error probability of the proposed concatenated codes $P_{\text{err}}(\mathcal{C}_{\text{cont}})$ can be bounded as

$$\begin{aligned} P_{\text{err}}(\mathcal{C}_{\text{cont}}) &= \mathbb{P}\{\mathbf{h} \neq \hat{\mathbf{h}}\} \\ &= \mathbb{P}\{\underbrace{v + 2w > M(1 - R)}_{\mathcal{A}}\} \\ &= \sum_{\mathbf{d}^M} \mathbb{P}\{v + 2w > M(1 - R) | \mathbf{d}^M\} \mathbb{P}\{\mathbf{d}^M\} \\ &\stackrel{(b)}{=} \sum_{\mathbf{u}^N \in \mathcal{T}} \mathbb{P}\{v + 2w > M(1 - R) | \mathbf{d}_{\mathbf{u}}^M\} \mathbb{P}\{\mathbf{u}^N\} + \sum_{\mathbf{u}^N \notin \mathcal{T}} \mathbb{P}\{v + 2w > M(1 - R) | \mathbf{d}_{\mathbf{u}}^M\} \mathbb{P}\{\mathbf{u}^N\} \\ &\stackrel{(c)}{\leq} \mathcal{P}(1 - \epsilon) + \epsilon \end{aligned}$$

where in (b) we use Lemma 1 and split $P_{\text{err}}(\mathcal{C}_{\text{cont}})$ into two parts, namely conditional probabilities given asymptotic and non-asymptotic distribution, where we define the asymptotic set $\mathcal{T} = \{\mathbf{u}^N : |U_d - Mp_c(d)| \leq f_d(M), d = 1, 2, \dots, N\}$ and $f_d(M)$ is the deviation function dependent on M . From Lemma 3, given $M \rightarrow \infty$, we conclude the probability $\mathbb{P}\{\mathbf{u}^N \notin \mathcal{T}\} \rightarrow 0$, i.e. the marginals D_i approach Poisson distributions and the random variables $\frac{U_d}{M} \rightarrow p_c(d)$ can be viewed as asymptotically deterministic, in spite of the statistical dependence among D_i . In (c) we bound the asymptotic error probability by \mathcal{P} the non-asymptotic one by 1.

The conditional probability with $\mathbf{u}^N \in \mathcal{T}$ can be stated as:

$$\begin{aligned} &\mathbb{P}\{v + 2w > M(1 - R) | \mathbf{d}_{\mathbf{u}}^M\} \\ &= \mathbb{P}\left\{2w > M(1 - R - \sum_{d < \theta} p_c(d)) | \mathbf{d}_{\mathbf{u}}^M\right\} \end{aligned}$$

Next given the $\mathbf{d}_{\mathbf{u}}^M$ sampling vector, $\{\eta_i\}_{i \in \mathcal{J}}$ are independent binomial random variables with different P_i respectively, where $P_i = \mathbb{P}\{d_{\text{H}}(\mathbf{X}_i, \hat{\mathbf{X}}_i) \geq \frac{d_{\text{H}}}{2}\}$. Define $\bar{P} = \max_{1 \leq i \leq M} P_i$ we denote the maximum error probability of inner decoder. Then we bound this conditional error probability with Lemma 2

$$\mathbb{P}\left\{2w > M(1 - R - \sum_{d < \theta} p_c(d))\right\} \leq \sum_{i=M\tau}^M \binom{M}{i} \bar{P}^i (1 - \bar{P})^{M-i}$$

$$\begin{aligned}
&\leq \sum_{i=\tau}^M \binom{M}{i} \bar{P}^i \leq \bar{P}^\tau \sum_{i=\tau}^M \binom{M}{i} \\
&\leq 2^M \cdot \bar{P}^M \frac{\tau}{M} \\
&\stackrel{(d)}{\leq} q^{-LM(E_g(p^*, r) \frac{\tau}{M} - \epsilon)}
\end{aligned}$$

where

$$\begin{aligned}
\tau &= \frac{M}{2} (1 - R - \sum_{d < \theta} p_c(d)) \\
p^* &= \sum_{i > \theta/2}^{\theta} \binom{\theta}{i} p^i (1-p)^{(\theta-i)}
\end{aligned}$$

In (d) we bound \bar{P}^M with the Shannon coding theorem for q -ary symmetric channel, and $E_g(p^*)$ is the Gallager error exponent [3]. The claim of Theorem 1 follows with $M \rightarrow \infty$, $\mathcal{P} = q^{-LM(E_g(p^*, r) \frac{\tau}{M} - o(1))}$, and the artificial small parameter ϵ .

References

- 1 S. Hirasawa and M. Kasahara, "Exponential error bounds and decoding complexity for block concatenated codes with tail biting trellis inner codes," *Journal of Discrete Mathematical Sciences and Cryptography*, vol. 9, no. 2, pp. 307–320, 2006.
- 2 A. Lenz, P. Siegel, A. Wachter-Zeh, and E. Yaakobi, "An upper bound on the capacity of the DNA storage channel," in *IEEE Inf. Theory Workshop (ITW)*, Visby, Sweden, Aug. 2019.
- 3 R. Gallager, "A simple derivation of the coding theorem and some applications," *IEEE Trans. Inf. Theory*, vol. 11, no. 1, pp. 3–18, 1965.