
UNIQUE RECONSTRUCTION FROM SUBSTRINGS

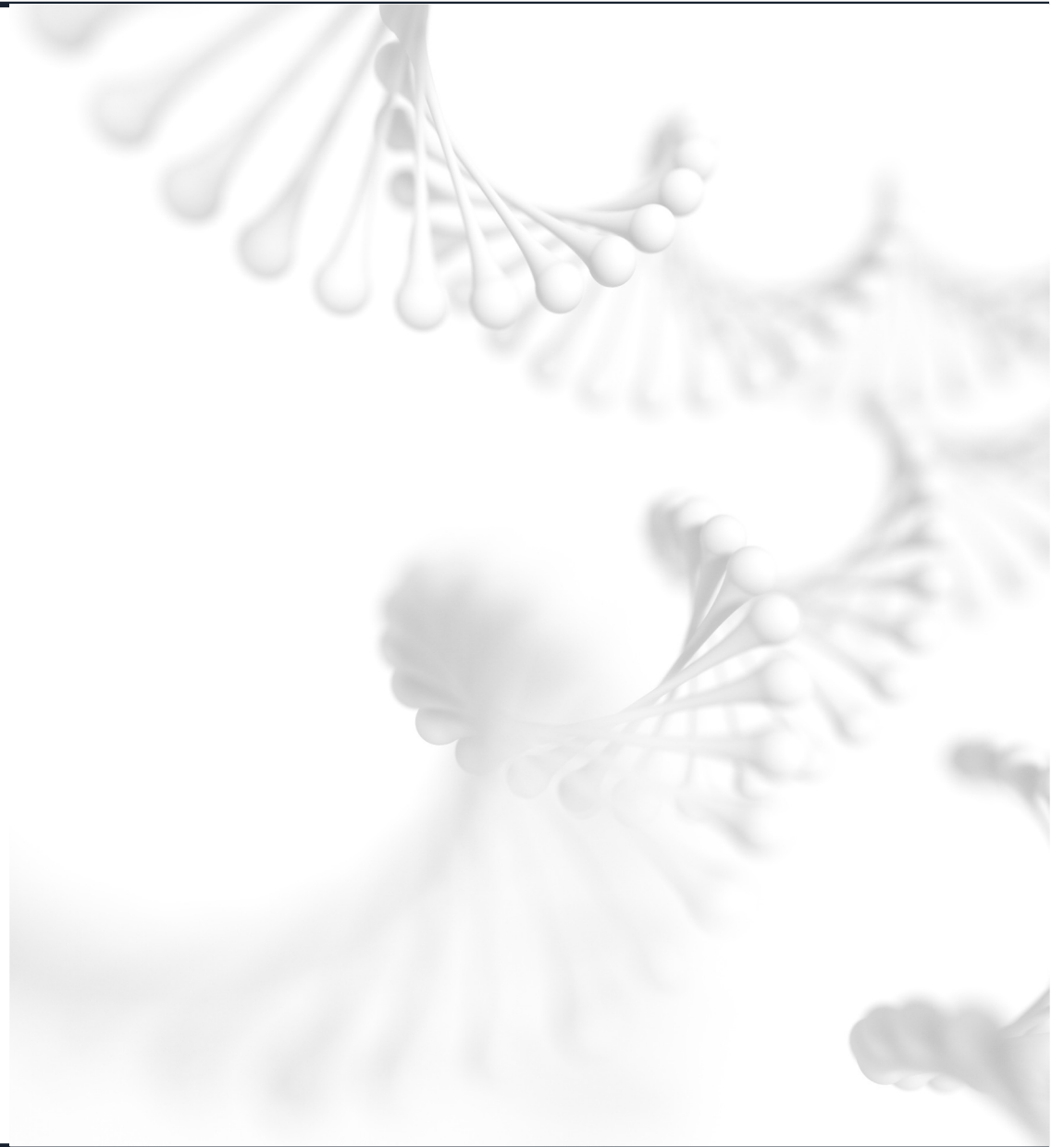
Daniella Bar-Lev | Technion-Israel Institute of Technology

A joint work with:

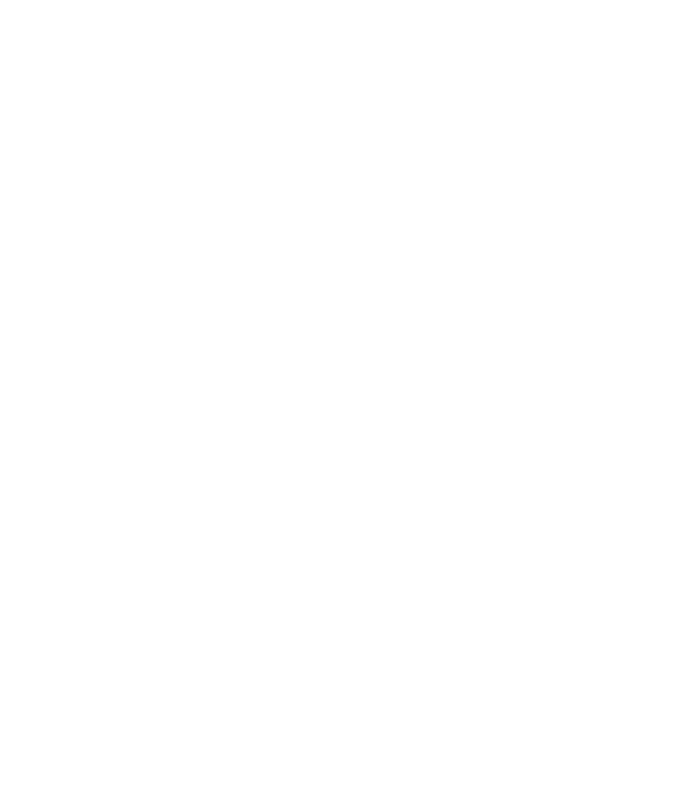
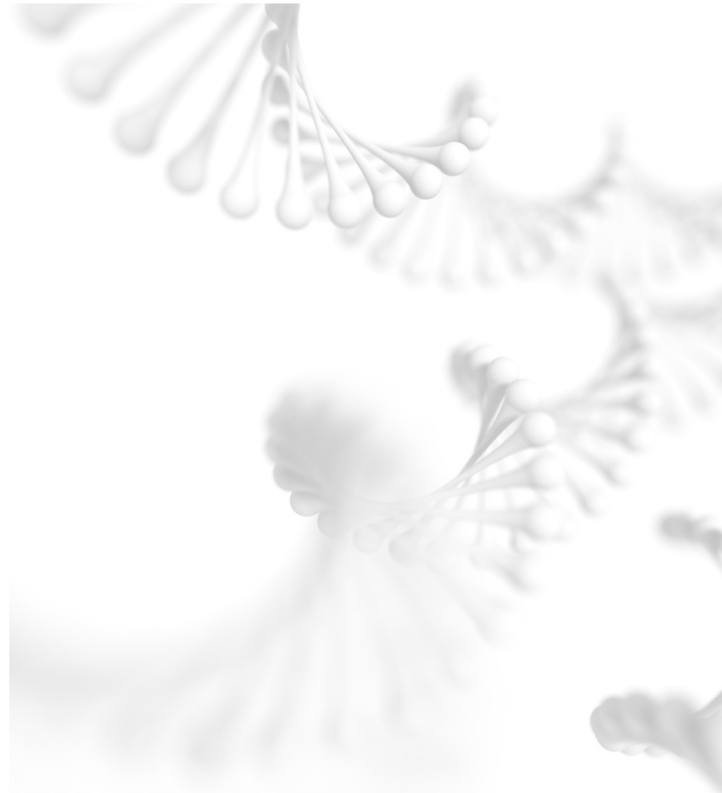
Yonatan Yehezkeally | Technical University of Munich

Sagi Marcovich | Technion-Israel Institute of Technology

Eitan Yaakobi | Technion-Israel Institute of Technology



DNA Storage



String Reconstruction

*Reconstruct string from **multiple, incomplete** and/or **noisy** observations*

Examples:

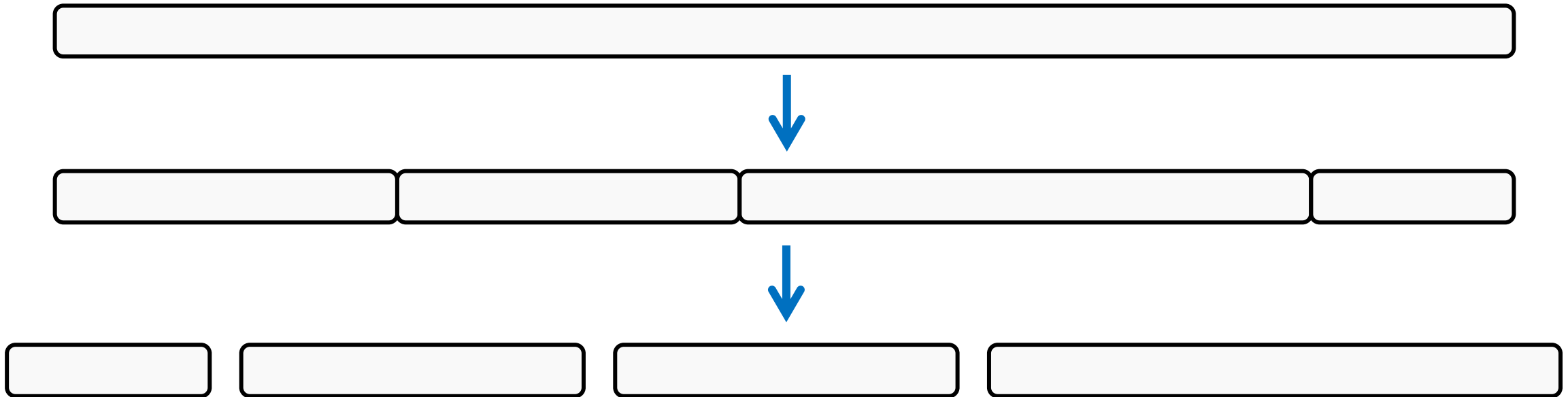
- Levenshtein's *reconstruction problem* [1]
- *Trace reconstruction problem* [2]
- *k-deck problem* [3]

[1] V. I. Levenshtein, "Efficient reconstruction of sequences from subsequences or supersequences," J. Combin. Theory, Feb. 2001

[2] T. Batu, S. Kannan, S. Khanna, and A. McGregor, "Reconstructing strings from random traces," in SODA, Society for Industrial and Applied Mathematics, Jan. 2004.

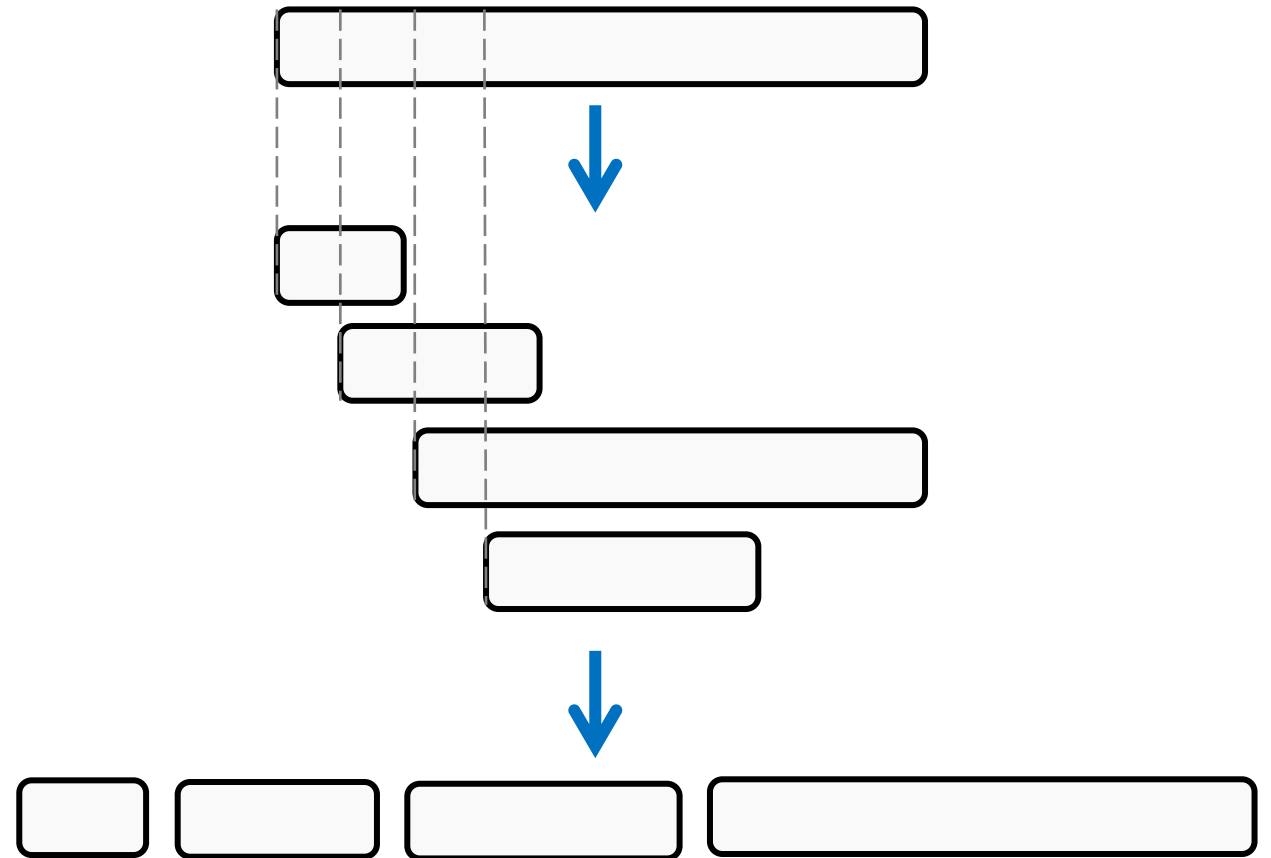
[3] B. Manvel, A. Meyerowitz, A. Schwenk, K. Smith, and P. Stockmeyer, "Reconstruction of sequences," Discrete Mathematics, 1991.

Torn-Paper Reconstruction



D. Bar-Lev, S. Marcovich, E. Yaakobi, and Y. Yehezkeally, "Adversarial torn-paper codes," in *IEEE ISIT*, 2022.
S. Nassirpour, I. Shomorony, and A. Vahid, "Reassembly codes for the chop-and-shuffle channel," in arXiv, 2022.
A. N. Ravi, A. Vahid, and I. Shomorony, "Capacity of the torn paper channel with lost pieces," in *IEEE ISIT*, 2021.
I. Shomorony and A. Vahid, "Torn-paper coding," *IEEE TIT*, Dec. 2021.

Reconstruction from Substring-Composition



R. Gabrys and O. Milenkovic, "Unique reconstruction of coded sequences from multiset substring spectra," *IEEE TIT*, Jun. 2019.

G. Bresler, M. Bresler, and D. Tse, "Optimal assembly for high throughput shotgun sequencing," *BMC Bioinformatics*, Jul. 2013.

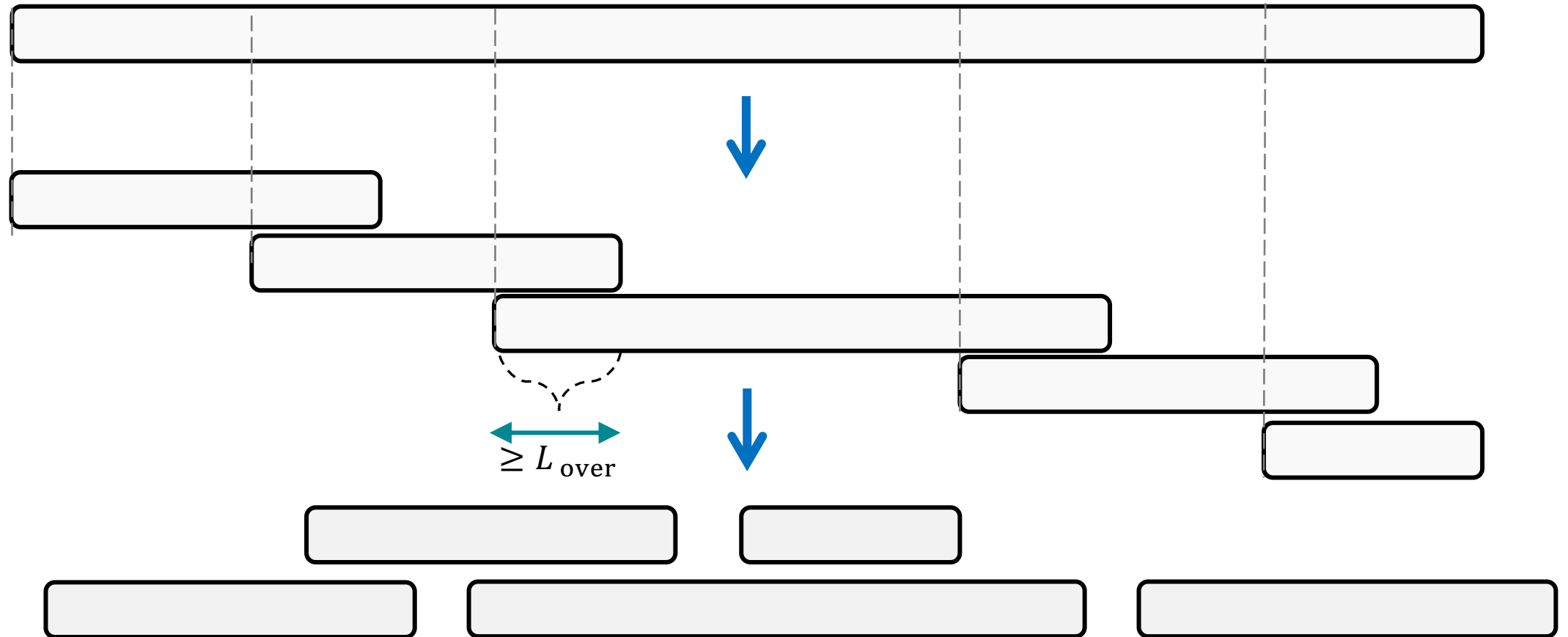
H. M. Kiah, G. J. Puleo, and O. Milenkovic, "Codes for DNA sequence profiles," *IEEE TIT*, Jun. 2016.

S. Marcovich and E. Yaakobi, "Reconstruction of strings from their substrings spectrum," *IEEE TIT*, Jul. 2021.

A. S. Motahari, G. Bresler, and D. N. C. Tse, "Information theory of DNA shotgun sequencing," *IEEE TIT*, Oct. 2013.

Y. Yehezkeally, S. Marcovich, and E. Yaakobi, "Multi-strand reconstruction from substrings," in *IEEE ITW*, 2021

Partial-Overlap Channel



Probabilistic version of this channel:

Ravi, A. N., Vahid, A., and Shomorony, I. (2022). Coded Shotgun Sequencing. *ACM JSAIT*, 3(1), 147-159.

Agenda

01

**Torn-Paper
Channel**

02

**Substring-
Composition
Channel**

03

**Partial Overlap
Channel**

04

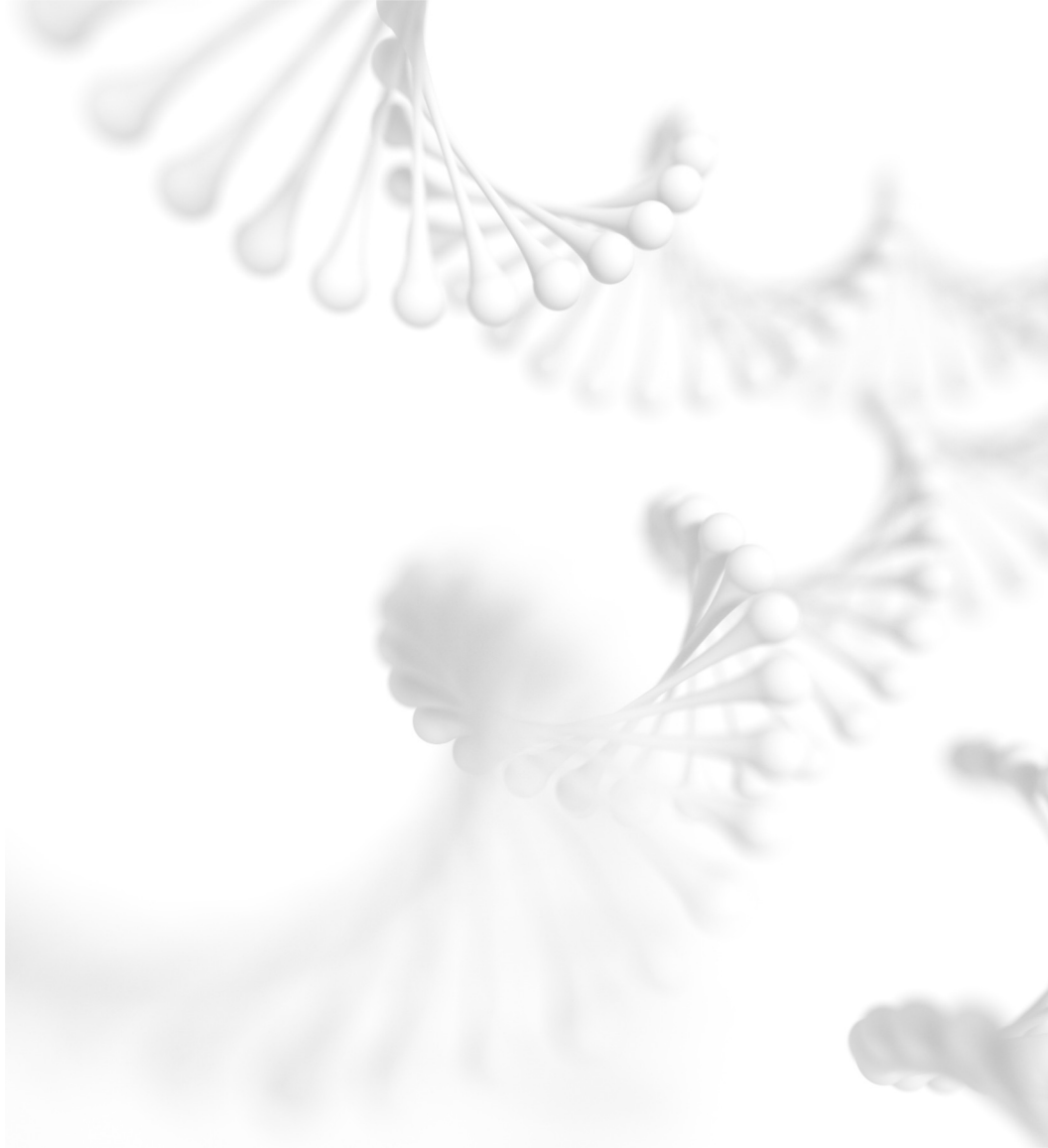
Future Directions

Notations

$[n] = \{0, 1, \dots, n - 1\}$.

$x \circ y$: concatenation of $x, y \in \Sigma^*$.

ℓ -substring of x : substring of length ℓ .

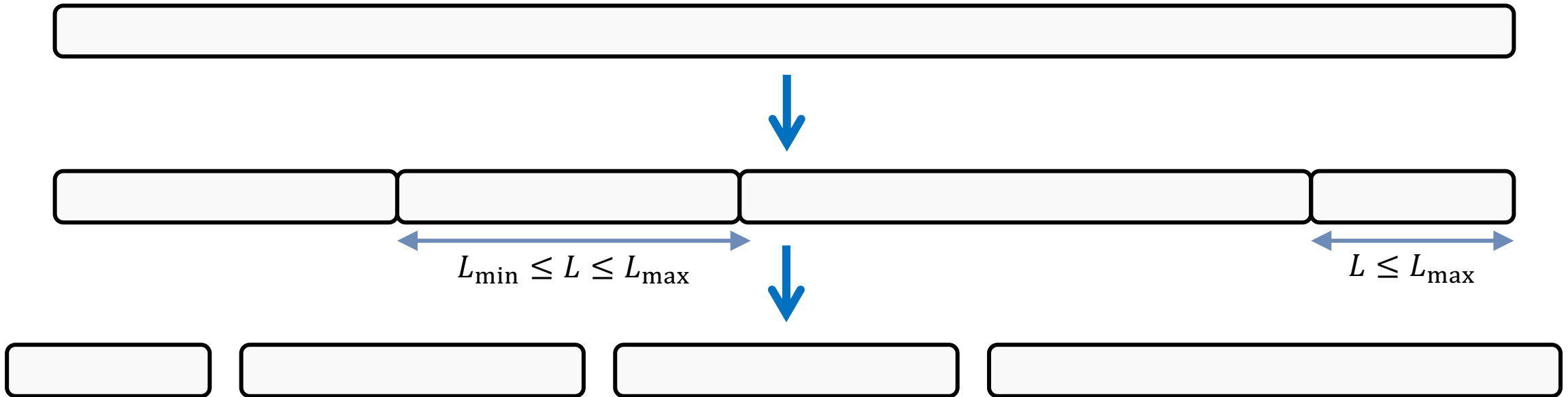


Torn-Paper Channel

Single-Strand



Single-Strand Torn-Paper Reconstruction



D. Bar-Lev, S. Marcovich, E. Yaakobi, and Y. Yehezkeally, "Adversarial torn-paper codes," in *IEEE ISIT*, 2022.
S. Nassirpour, I. Shomorony, and A. Vahid, "Reassembly codes for the chop-and-shuffle channel," in arXiv, 2022.
A. N. Ravi, A. Vahid, and I. Shomorony, "Capacity of the torn paper channel with lost pieces," in *IEEE ISIT*, 2021.
I. Shomorony and A. Vahid, "Torn-paper coding," *IEEE TIT*, Dec. 2021.

Definitions

An (L_{\min}, L_{\max}) -segmentation of $x \in \Sigma^n$ is a multiset

$\{\{\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{m-1}\}\}$ of substrings such that:

(1) $x = \mathbf{u}_0 \circ \mathbf{u}_1 \circ \dots \circ \mathbf{u}_{m-1}$.

(2) $L_{\min} \leq |\mathbf{u}_i| \leq L_{\max}$ for $0 \leq i < m - 1$ and $|\mathbf{u}_{m-1}| \leq L_{\max}$

$\mathcal{J}_{L_{\min}}^{L_{\max}}(x)$: the multiset of all (L_{\min}, L_{\max}) -segmentations of x .

Example: $x=0001011001$

$\{\{0001, 01, 10, 01\}\}$ is a $(2,4)$ -segmentation of x .

Definitions

An (L_{\min}, L_{\max}) -segmentation of $\mathbf{x} \in \Sigma^n$ is a multiset

$\{\{\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{m-1}\}\}$ of substrings such that:

(1) $\mathbf{x} = \mathbf{u}_0 \circ \mathbf{u}_1 \circ \dots \circ \mathbf{u}_{m-1}$.

(2) $L_{\min} \leq |\mathbf{u}_i| \leq L_{\max}$ for $0 \leq i < m - 1$ and $|\mathbf{u}_{m-1}| \leq L_{\max}$

$\mathcal{J}_{L_{\min}}^{L_{\max}}(\mathbf{x})$: the multiset of all (L_{\min}, L_{\max}) -segmentations of \mathbf{x} .

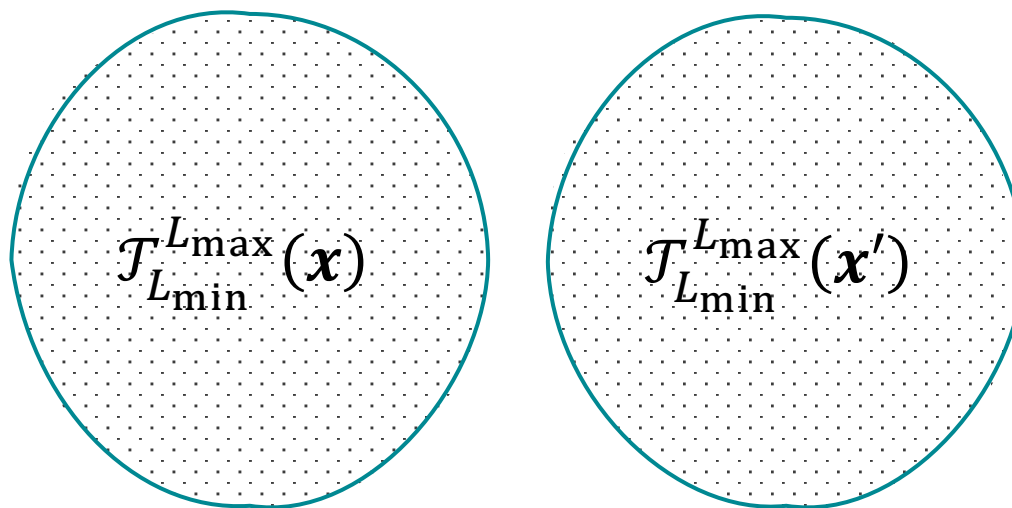
Channel input: $\mathbf{x} \in \Sigma^n$

Channel output: an (L_{\min}, L_{\max}) -segmentation of \mathbf{x} .

Definitions

$(L_{\min}, L_{\text{over}})$ - *single strand torn-paper code*:

code $\mathcal{C} \subseteq \Sigma^n$ such that for any two distinct strings $\mathbf{x}, \mathbf{x}' \in \mathcal{C}$, it holds that $\mathcal{J}_{L_{\min}}^{L_{\max}}(\mathbf{x}) \cap \mathcal{J}_{L_{\min}}^{L_{\max}}(\mathbf{x}') = \emptyset$.



Single-Strand Torn-Paper Codes: Rate

Theorem: If $L_{\min} = a \log n + O_n(1)$, for some value $a > 1$, then for any (L_{\min}, L_{\max}) -single-strand torn-paper code $\mathcal{C} \subseteq \Sigma^n$,

$$R(\mathcal{C}) \leq 1 - \frac{1}{a} + o(1).$$

Preliminaries: Gray Code

A *Gray Code* is an ordering of Σ^n such that any two adjacent strings differ in only one symbol position.

Example:

$\Sigma = \{0,1\}$ and $n = 3$

$|\Sigma^3| = 8$

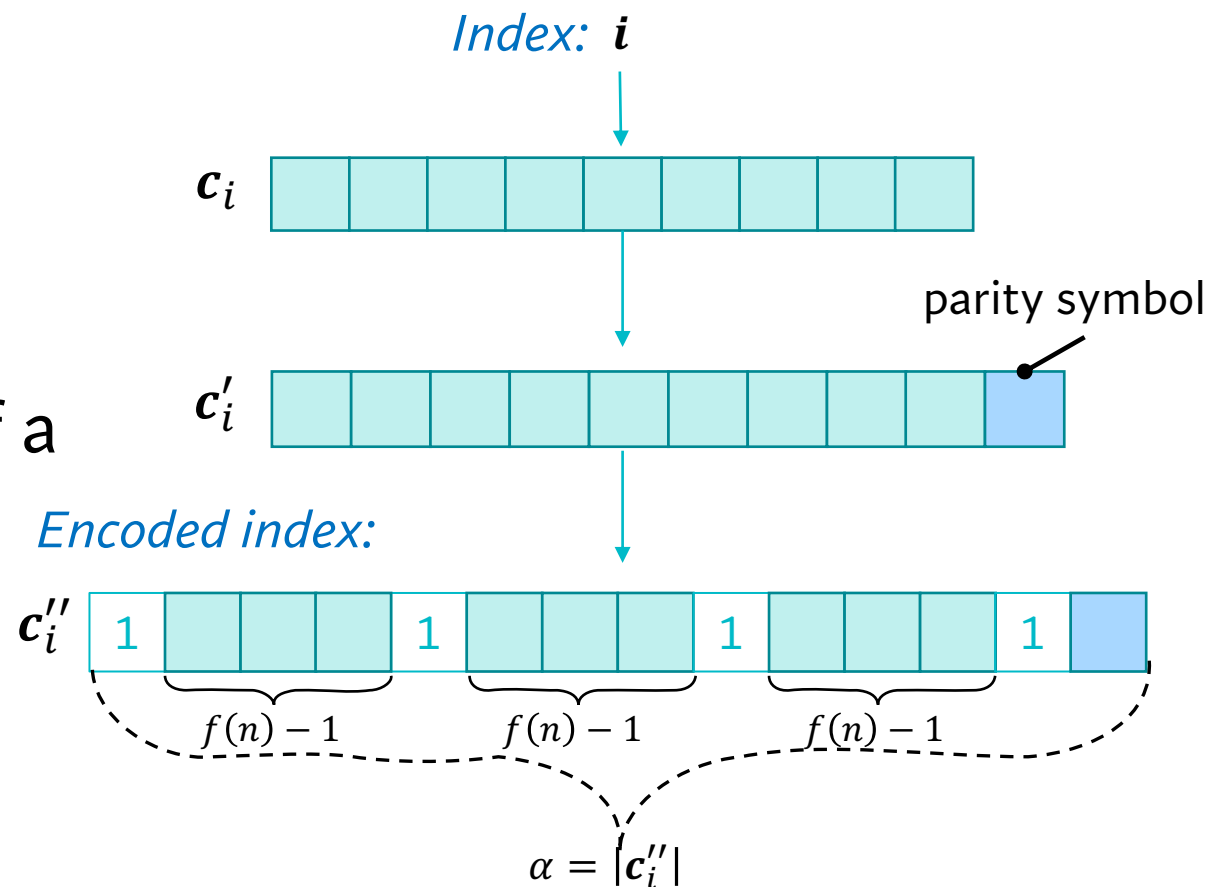
0	000
1	001
2	011
3	010

4	110
5	100
6	101
7	111

Preliminaries: Index Generation

Consider:

- Index length: $I = \lceil \log_q(n/L_{\min}) \rceil$
- Block size: $f(n) = o(\log n)$
- $\mathbf{c}_0, \mathbf{c}_1, \dots, \mathbf{c}_{q^I-1}$: the codewords of a q -ary Gray code, in order.



Preliminaries: Run-Length Limited Encoding

For $t < N$ the set of *run-length limited (RLL)* strings is

$$\mathcal{RLL}_t(N) \triangleq \{x \in \Sigma^N : \text{no } t\text{-length runs of zeros}\}.$$

RLL encoder: input length m , output length $N = N_n(m)$, and $t = f(n)$

$$E_m^{RLL} : \Sigma^m \rightarrow \mathcal{RLL}_{f(n)}(N_n(m)).$$

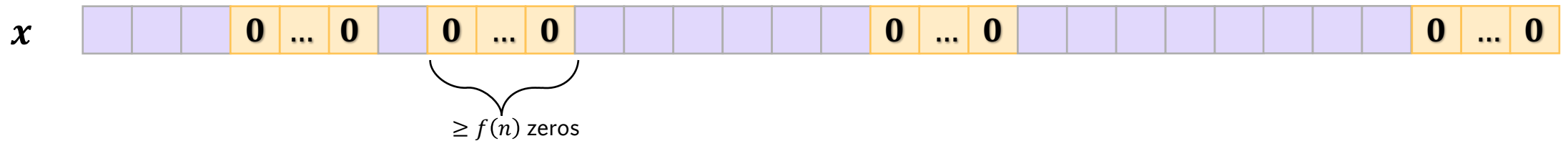
[1] M. Levy and E. Yaakobi, "Mutually uncorrelated codes for DNA storage," IEEE Trans. on Inform. Theory, vol. 65, no. 6, pp. 3671–3691, Jun. 2019.

[2, Lem. 5] Y. Yehezkeally, S. Marcovich, and E. Yaakobi, "Multi-strand reconstruction from substrings," in IEEE ITW, 2021.

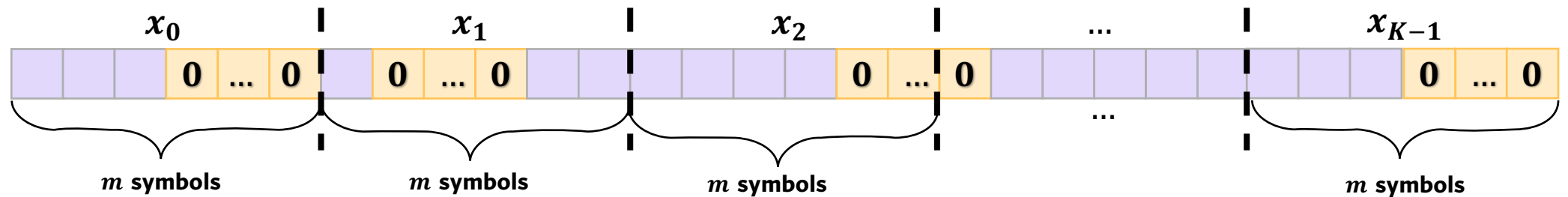
Construction A

Encoder for (L_{\min}, L_{\max}) -single strand Torn-Paper code $\mathcal{C}_A(n)$:

Input: a sequence $x \in \Sigma^{Km}$



Divide x into q^l non-overlapping substrings of length m .

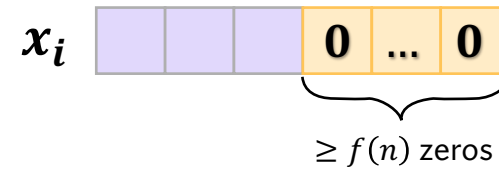


Next, we encode each substring x_i into $z_i \in \Sigma^{L_{\min}}$ independently.

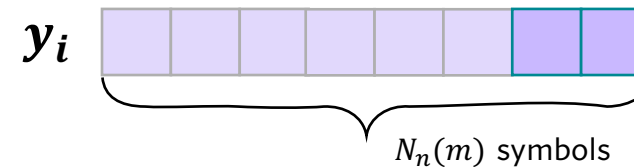
Construction A

Encoder for (L_{\min}, L_{\max}) -single strand Torn-Paper code $\mathcal{C}_A(n)$:

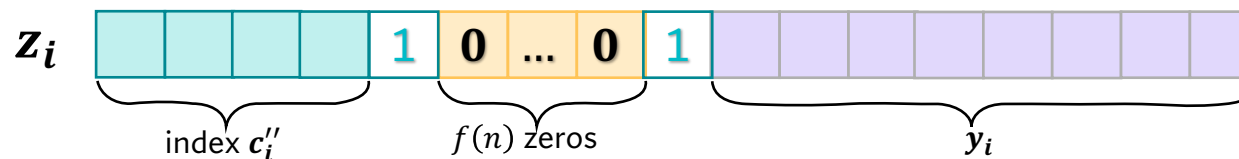
For any $0 \leq i \leq K - 1$:



(1) Encode x_i using the RLL encoder to obtain $y_i = E_m^{RLL}(x_i)$.



(2) Let c_i'' be the i -th encoded index and let $z_i = c_i'' \circ 10^{f(n)} 1 \circ y_i$.

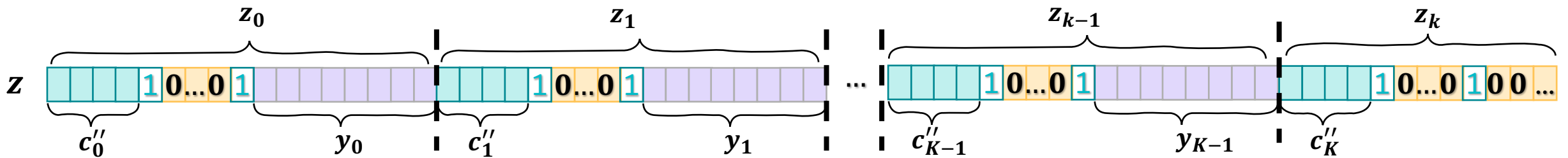


$$|z_i| = L_{\min}.$$

Construction A

Encoder for (L_{\min}, L_{\max}) -single strand Torn-Paper code $\mathcal{C}_A(n)$:

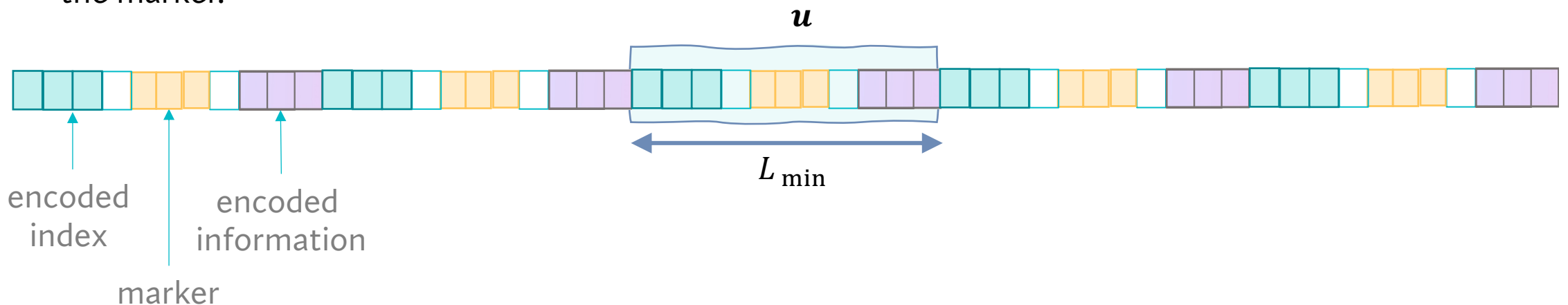
Finally, let $\text{Enc}_A(x) \triangleq z = z_0 \circ z_1 \circ \dots \circ z_{K-1} \circ z_K$



Construction A – Encoder Correctness

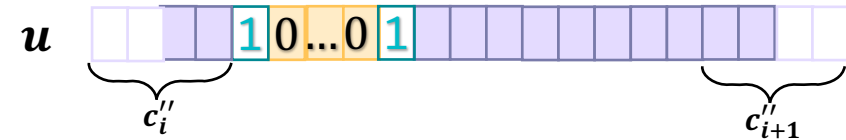
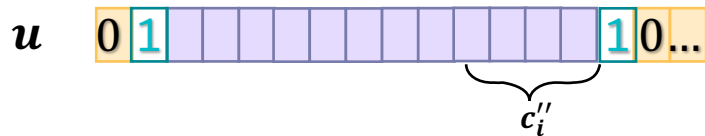
Every L_{\min} -substring u of z :

- does not contain any occurrences of the marker `10001` except those explicitly added after each encoded index.
- either u contains an occurrence of the marker or it has a suffix-prefix pair whose concatenation is the marker.



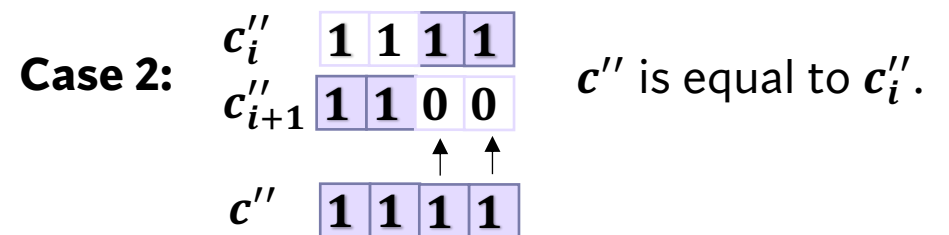
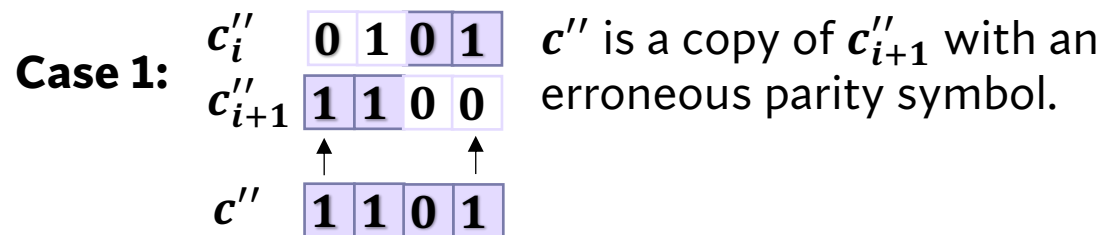
Construction A – Encoder Correctness

Either u contains an occurrence of $\boxed{1}00\boxed{1}$ or it has a suffix-prefix pair whose concatenation is $\boxed{1}00\boxed{1}$.



c''_i and c''_{i+1} differ only at the parity symbol and one more coordinate.

c'' : the concatenation of the $(\alpha - j)$ -suffix of u with the j -prefix of u .



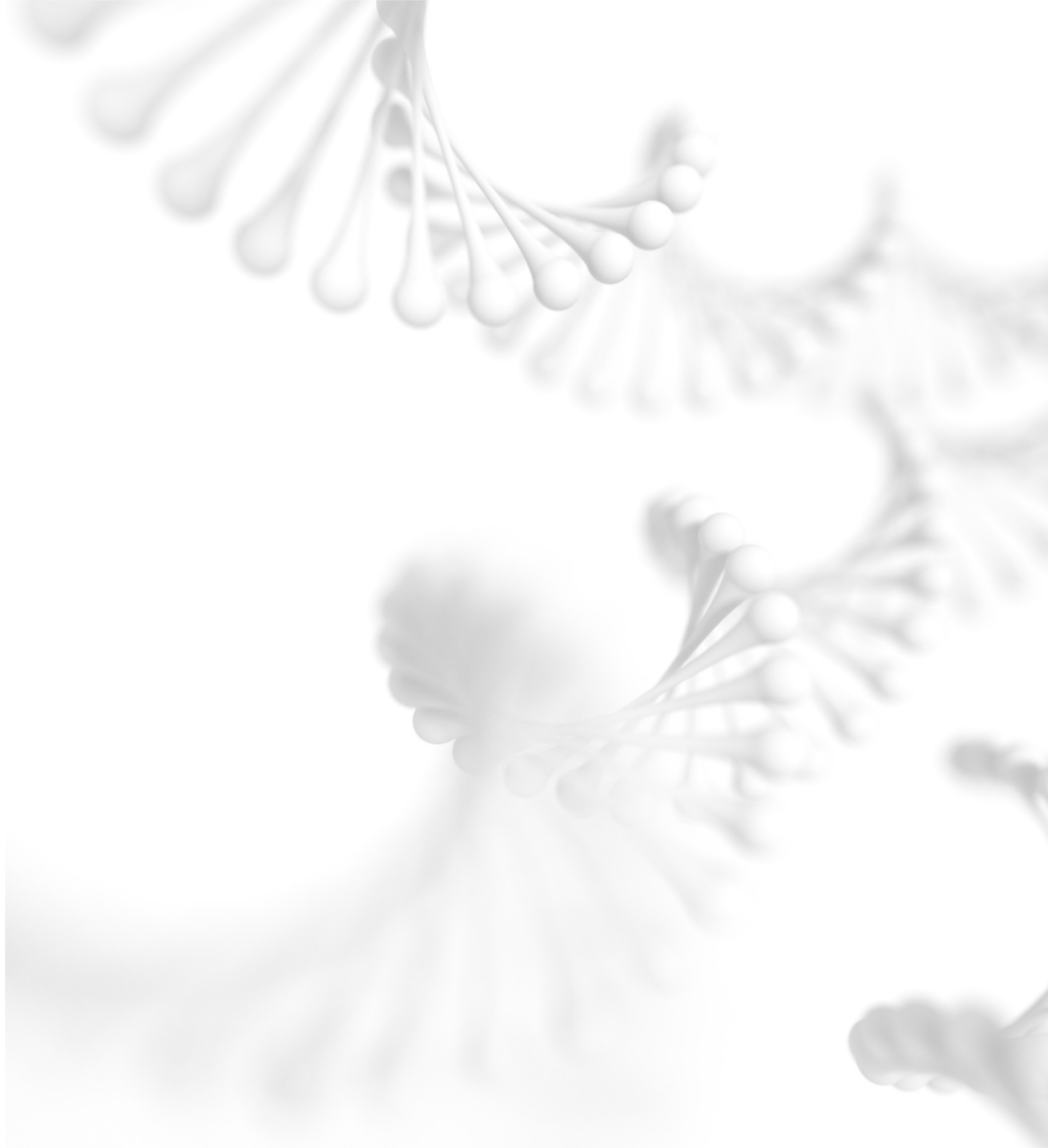
Hence, if the parity symbol is correct then i is the decoding of c'' and otherwise i is the decoding of c'' minus one.

Construction A – Rate

Theorem: Letting $f(n) \triangleq (1 + o(1))\sqrt{\log(n)}$ we have that

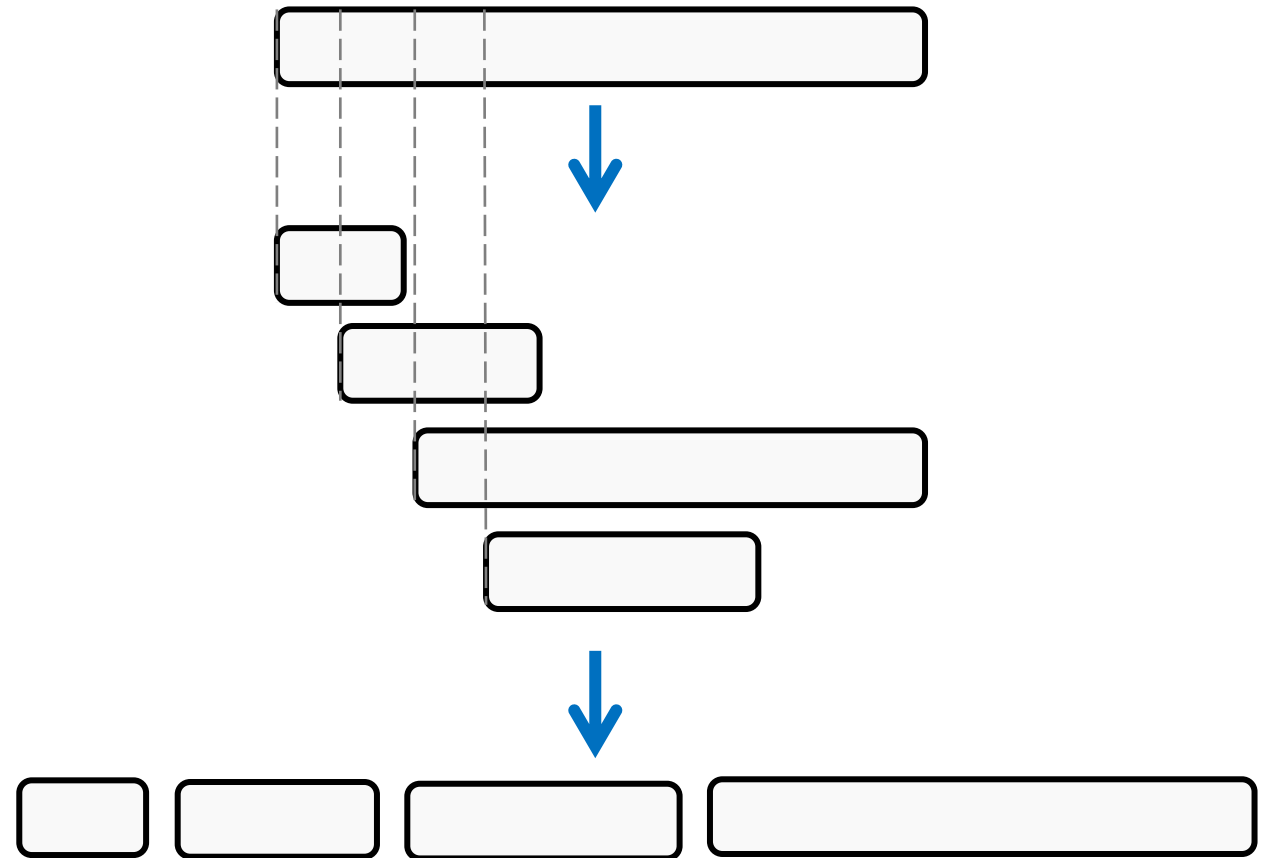
$$\text{red}(\mathcal{C}_A(n)) \leq \frac{n}{a} \left(1 + \frac{2 + o(1)}{\sqrt{\log(n)}} \right).$$

Hence, Construction A asymptotically meets the upper bound.



Multi-Strand Reconstruction from Substring-Composition

Reconstruction from Substring-Composition



R. Gabrys and O. Milenkovic, "Unique reconstruction of coded sequences from multiset substring spectra," *IEEE TIT*, Jun. 2019.

G. Bresler, M. Bresler, and D. Tse, "Optimal assembly for high throughput shotgun sequencing," *BMC Bioinformatics*, Jul. 2013.

H. M. Kiah, G. J. Puleo, and O. Milenkovic, "Codes for DNA sequence profiles," *IEEE TIT*, Jun. 2016.

S. Marcovich and E. Yaakobi, "Reconstruction of strings from their substrings spectrum," *IEEE TIT*, Jul. 2021.

A. S. Motahari, G. Bresler, and D. N. C. Tse, "Information theory of DNA shotgun sequencing," *IEEE TIT*, Oct. 2013.

Y. Yehezkeally, S. Marcovich, and E. Yaakobi, "Multi-strand reconstruction from substrings," in *IEEE ITW*, 2021

Notations

$$\mathcal{X}_{n,k} \triangleq \{S = \{\{x_1, \dots, x_k\}\} : \forall i, x_i \in \Sigma^n\}.$$

$|S|$: number of unique elements in S .

Definitions

An ℓ -trace of $x \in \Sigma^n$ is a multiset of substrings such that:

- (1) All substrings are of length at least ℓ .
- (2) Succeeding substrings overlap is at least $\ell - 1$.
- (3) x is covered by the substrings.

ℓ -trace spectrum of x , denoted by $\mathcal{T}_\ell(x)$: set of all ℓ -traces of x .

Example: $x=1110111$

$\{\{11101, 1101, 10111, 0111\}\}$ is a 4-trace of x .

Definitions

An *ℓ -trace* of $x \in \Sigma^n$ is a multiset of substrings such that:

- (1) All substrings are of length at least ℓ .
- (2) Succeeding substrings overlap is at least $\ell - 1$.
- (3) x is covered by the substrings.

ℓ -trace spectrum of x , denoted by $\mathcal{T}_\ell(x)$: set of all ℓ -traces of x .

Example: $x=1110111$

$\{\{11101, 1101, 10111, 0111\}\}$ is a 4-trace of x .

$$\mathcal{T}_\ell(\mathcal{S}) \triangleq \bigcup_{x \in \mathcal{S}} \mathcal{T}_\ell(x).$$

Definitions

An *ℓ -trace* of $\mathbf{x} \in \Sigma^n$ is a multiset of substrings such that:

- (1) All substrings are of length at least ℓ .
- (2) Succeeding substrings overlap is at least $\ell - 1$.
- (3) \mathbf{x} is covered by the substrings.

ℓ -trace spectrum of \mathbf{x} , denoted by $\mathcal{T}_\ell(\mathbf{x})$: set of all ℓ -traces of \mathbf{x} .

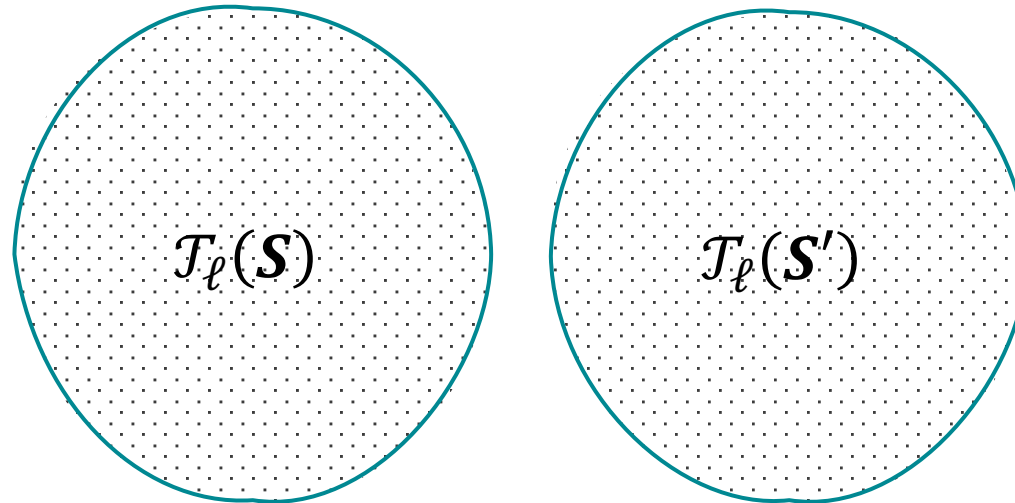
Channel input: $\mathcal{S} \in \mathcal{X}_{n,k}$

Channel output: an ℓ -trace of \mathcal{S} .

$$\mathcal{T}_\ell(\mathcal{S}) \triangleq \bigcup_{\mathbf{x} \in \mathcal{S}} \mathcal{T}_\ell(\mathbf{x}).$$

Definitions

Multi-strand ℓ -trace code: code $\mathcal{C} \subseteq \mathcal{X}_{n,k}$ such that for any two distinct multisets $\mathcal{S}, \mathcal{S}' \in \mathcal{C}$, it holds that $\mathcal{T}_\ell(\mathcal{S}) \cap \mathcal{T}_\ell(\mathcal{S}') = \emptyset$.



Preliminaries: Repeat-Free Encoding

$\mathcal{L}_\ell(\mathbf{x}) \in \mathcal{T}_\ell(\mathbf{x})$ is the ℓ -trace that consists of all the ℓ -substrings of \mathbf{x} .

Example: $\mathbf{x} = 11101110101111$

$$\mathcal{L}_{12}(\mathbf{x}) = \{111011101011, 110111010111, 101110101111\}$$

For $\ell < n$ the set of *ℓ -repeat-free (RF)* strings is

$$\begin{aligned} \mathcal{RF}_\ell(n) &\triangleq \{\mathbf{x} \in \Sigma^n: \text{no } \ell\text{-substring repeats}\} \\ &= \{\mathbf{x} \in \Sigma^n: |\mathcal{L}_\ell(\mathbf{x})| = n - \ell + 1\}. \end{aligned}$$

A *multi-strand ℓ -repeat-free* strings:

$$\mathcal{RF}_\ell(n, k) \triangleq \{\mathbf{S} \in \mathcal{X}_{n,k}: |\mathcal{L}_\ell(\mathbf{S})| = k(n - \ell + 1)\}$$

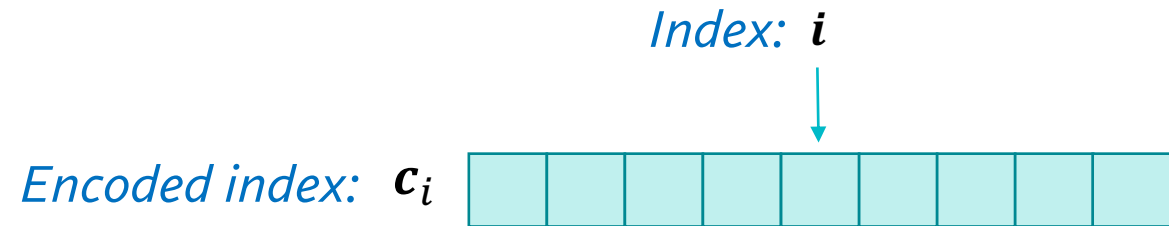
Preliminaries: Repeat-Free Encoding

Lemma: For all $S \in \mathcal{RF}_\ell(n, k)$ there exists an efficient algorithm reconstructing S from any ℓ -trace of S .

Preliminaries: Index Generation

Consider:

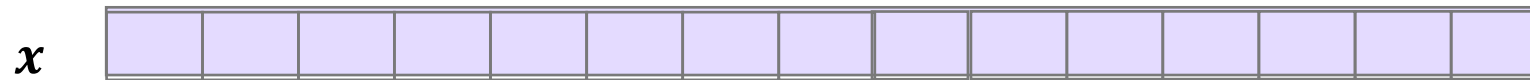
- Index length: $I = \lceil \log_q(k) \rceil$



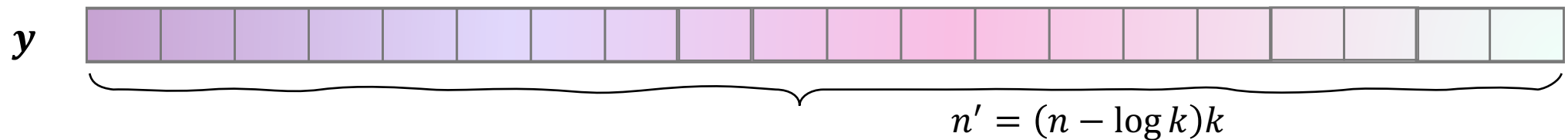
Construction B

Encoder for multi-strand ℓ -trace code \mathcal{C}_B :

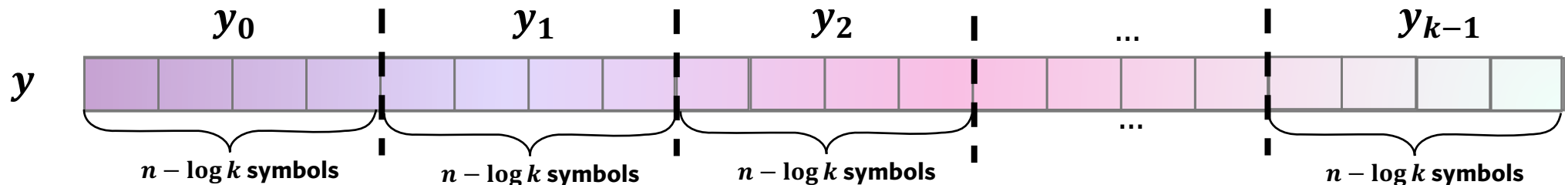
Input: a sequence $x \in \Sigma^m$



Encode x using the RF encoder for $\ell' = \ell - \log k$ to obtain $y = E_{m, \ell'}^{\mathcal{RF}}(x)$.

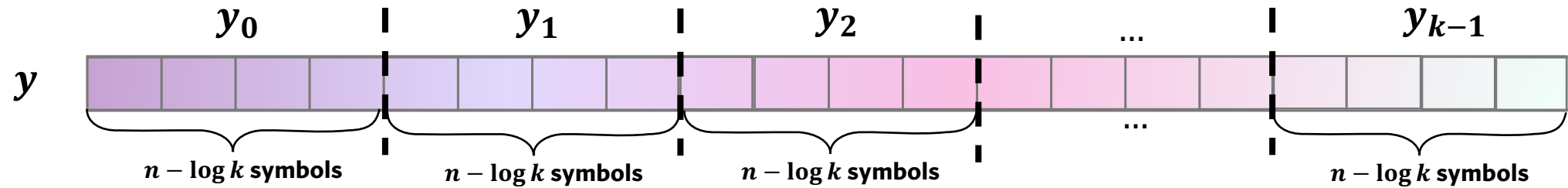


Divide y into k non-overlapping substrings of length $n - \log k$.

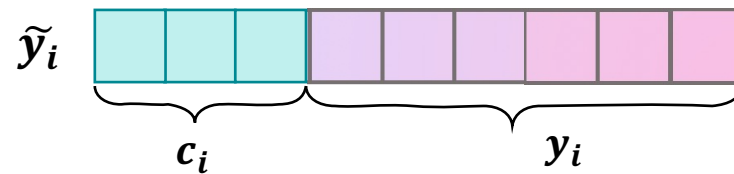


Construction B

Encoder for multi-strand ℓ -trace code \mathcal{C}_B :



For each y_i define $\tilde{y}_i = c_i \circ y_i$

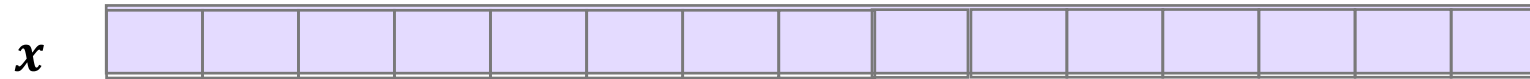


Finally, let $\text{Enc}_B(\mathbf{x}) \triangleq \{\{\tilde{y}_i : i \in [k]\}\} \in \mathcal{X}_{n,k}$.

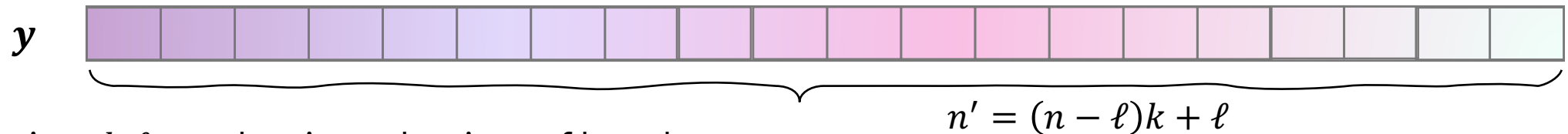
Construction C

Encoder for multi-strand ℓ -trace code \mathcal{C}_C :

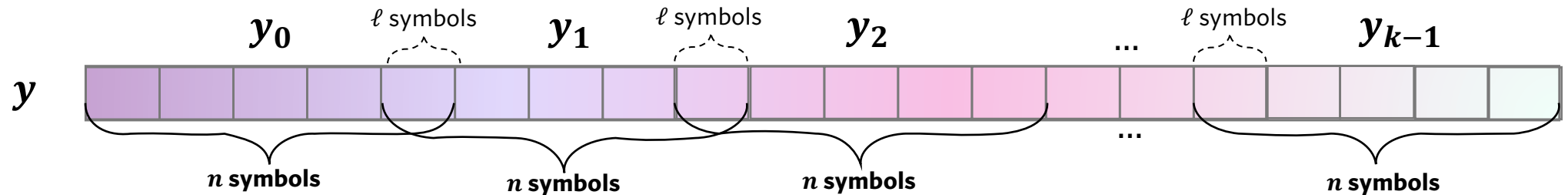
Input: a sequence $\mathbf{x} \in \Sigma^m$



Encode \mathbf{x} using the RF encoder to obtain $\mathbf{y} = E_{m,\ell}^{\mathcal{RF}}(\mathbf{x})$.

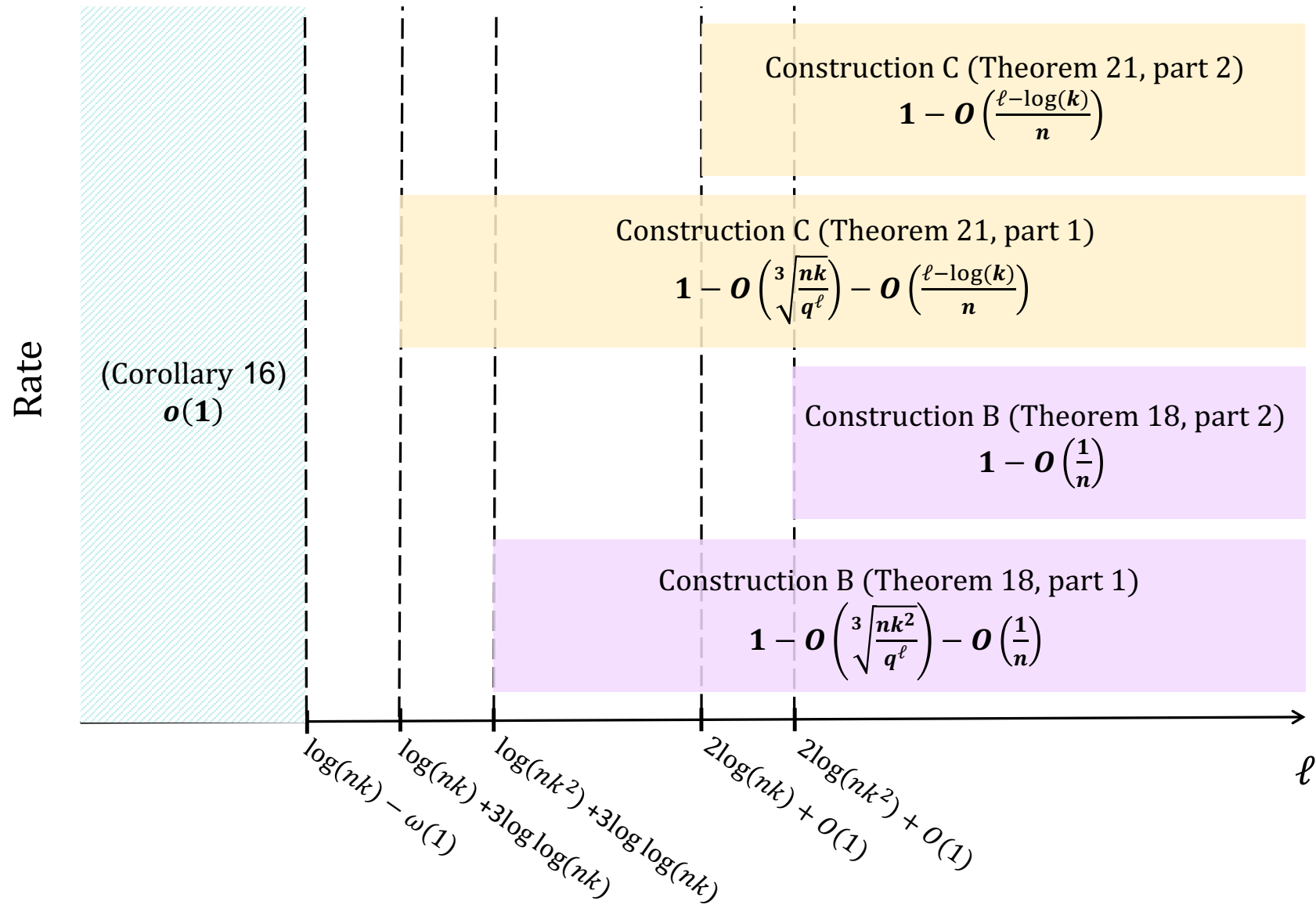


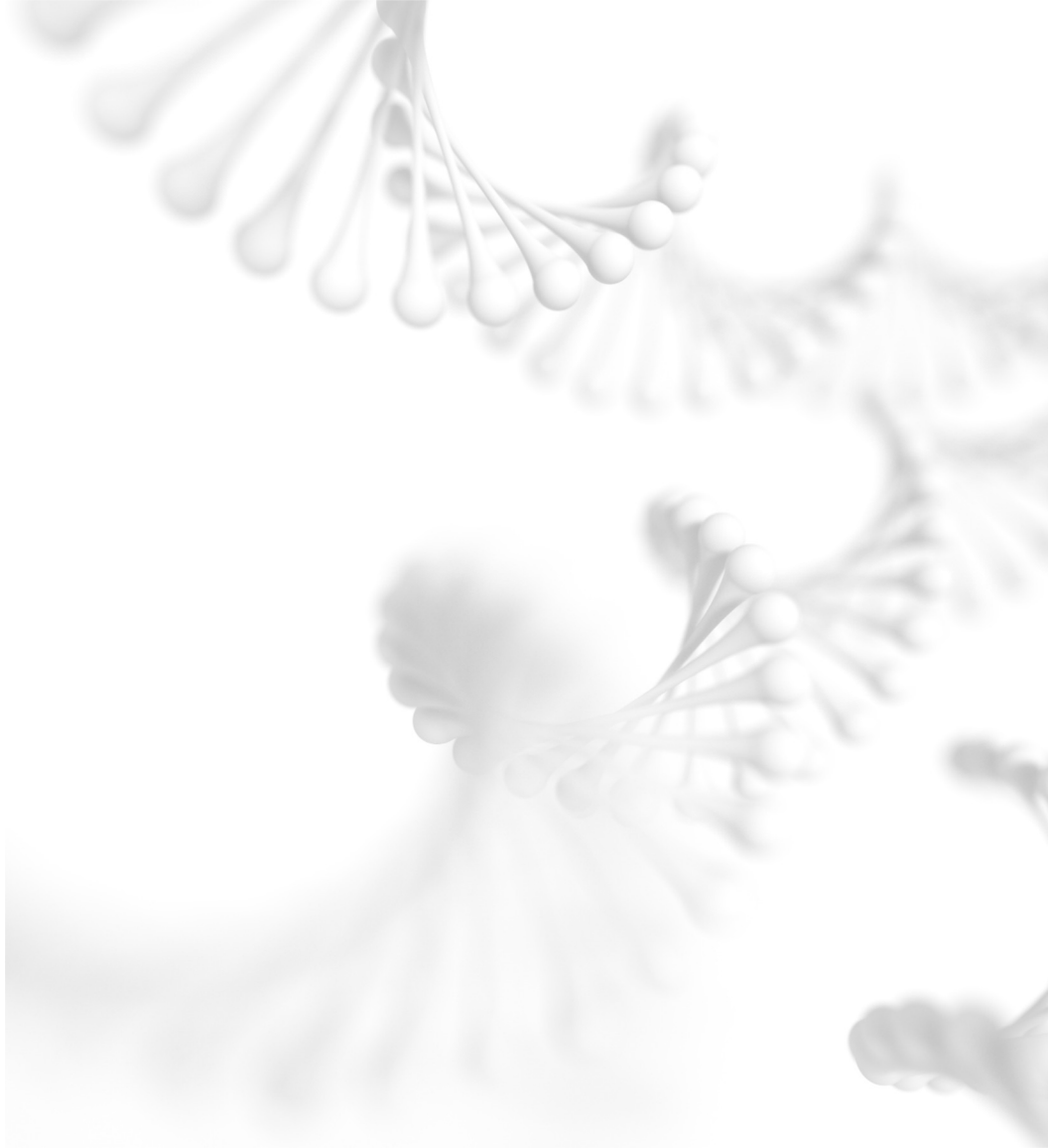
Divide \mathbf{y} into k ℓ -overlapping substrings of length n .



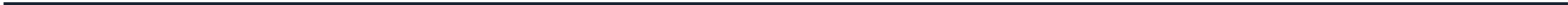
Finally, let $\text{Enc}_C(\mathbf{x}) \triangleq \{\{\mathbf{y}_i : i \in [k]\}\} \in \mathcal{X}_{n,k}$.

Constructions B&C – Rate

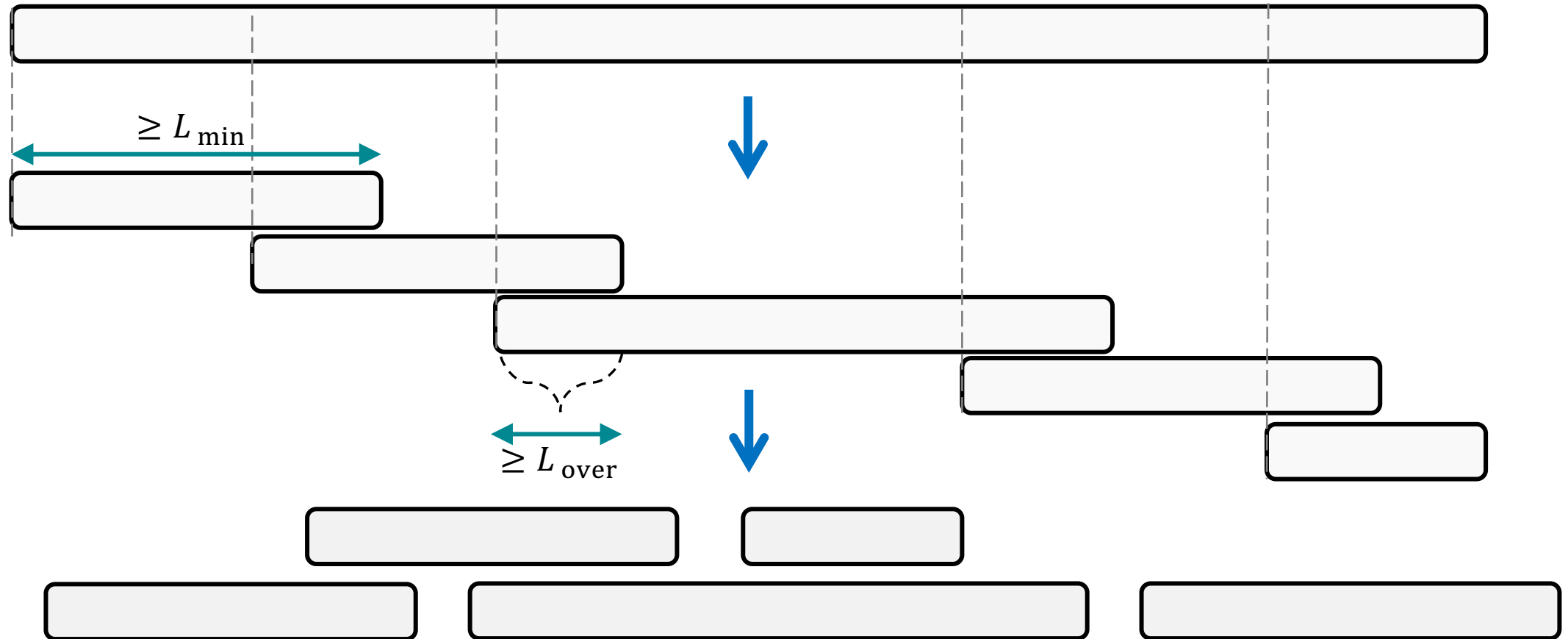




Reconstruction from Substrings with Partial Overlap



Partial-Overlap Channel



Probabilistic version of this channel:

Ravi, A. N., Vahid, A., and Shomorony, I. (2022). Coded Shotgun Sequencing. *ACM JSAIT*, 3(1), 147-159.

Definitions

An $(L_{\min}, L_{\text{over}})$ -trace of $x \in \Sigma^n$ is a multiset of substrings such that:

- (1) All substrings are of length at least L_{\min} .
- (2) Succeeding substrings overlap at least L_{over} .
- (3) x is covered by the substrings.

$(L_{\min}, L_{\text{over}})$ - trace spectrum of x , denoted by $\mathcal{J}_{L_{\min}}^{L_{\text{over}}}(x)$:

set of all $(L_{\min}, L_{\text{over}})$ -traces of x .

Example: $x=111011101011111$

$\{\{1110111, 1110101, 011111\}\}$ is a (6,2)-trace of x .

Definitions

An $(L_{\min}, L_{\text{over}})$ -trace of $x \in \Sigma^n$ is a multiset of substrings such that:

- (1) All substrings are of length at least L_{\min} .
- (2) Succeeding substrings overlap at least L_{over} .
- (3) x is covered by the substrings.

$(L_{\min}, L_{\text{over}})$ - trace spectrum of x , denoted by $\mathcal{J}_{L_{\min}}^{L_{\text{over}}}(x)$:

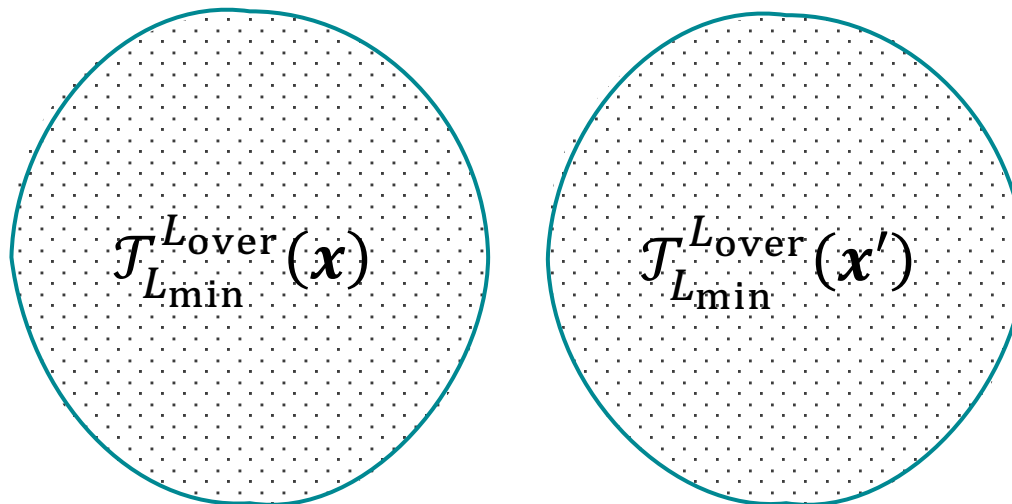
set of all $(L_{\min}, L_{\text{over}})$ -traces of x .

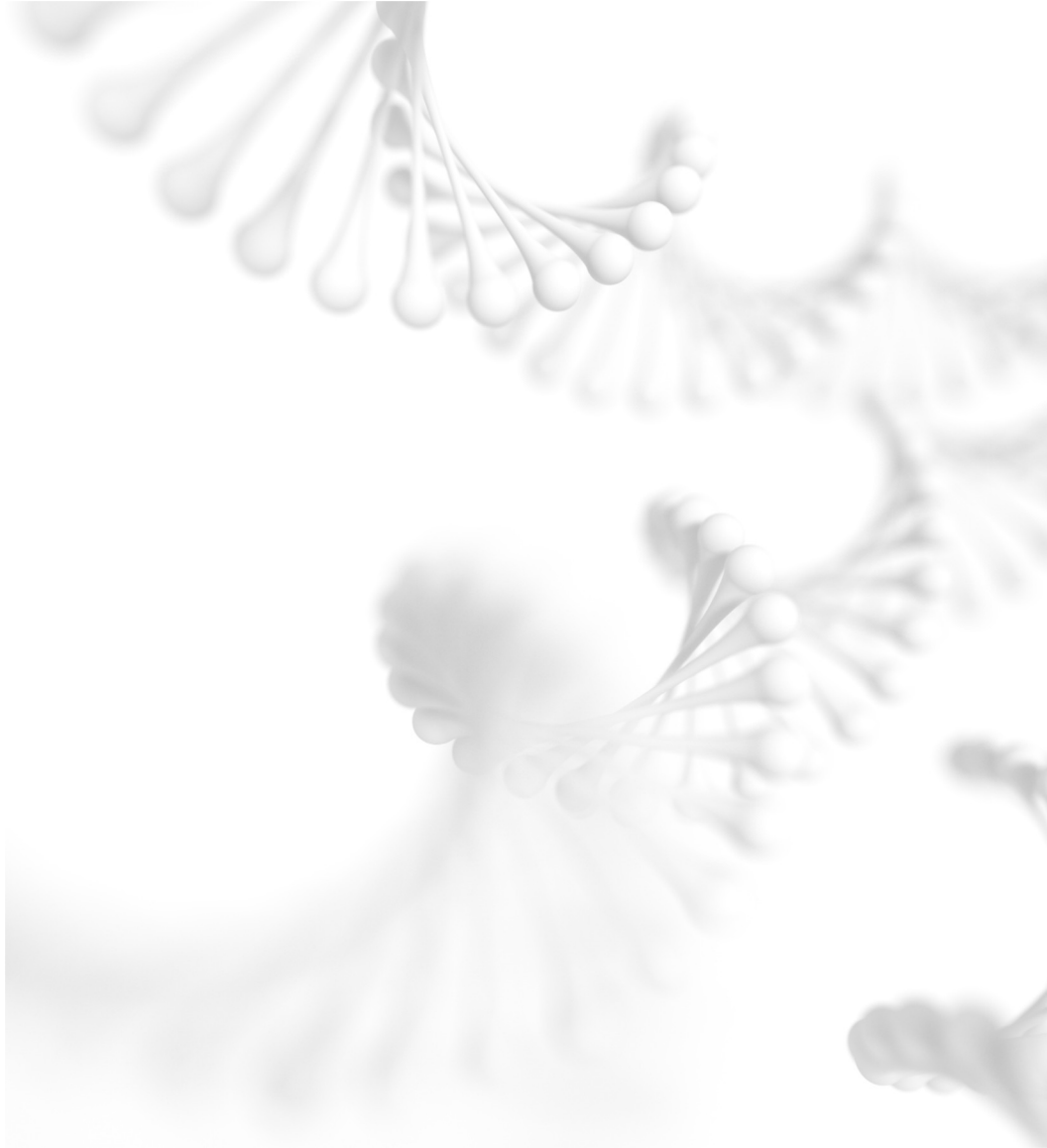
Channel input: $x \in \Sigma^n$

Channel output: an $(L_{\min}, L_{\text{over}})$ -trace of x .

Definitions

$(L_{\min}, L_{\text{over}})$ -trace code: code $\mathcal{C} \subseteq \Sigma^n$ such that for any two distinct strings $\mathbf{x}, \mathbf{x}' \in \mathcal{C}$, it holds that $\mathcal{J}_{L_{\min}}^{L_{\text{over}}}(\mathbf{x}) \cap \mathcal{J}_{L_{\min}}^{L_{\text{over}}}(\mathbf{x}') = \emptyset$.





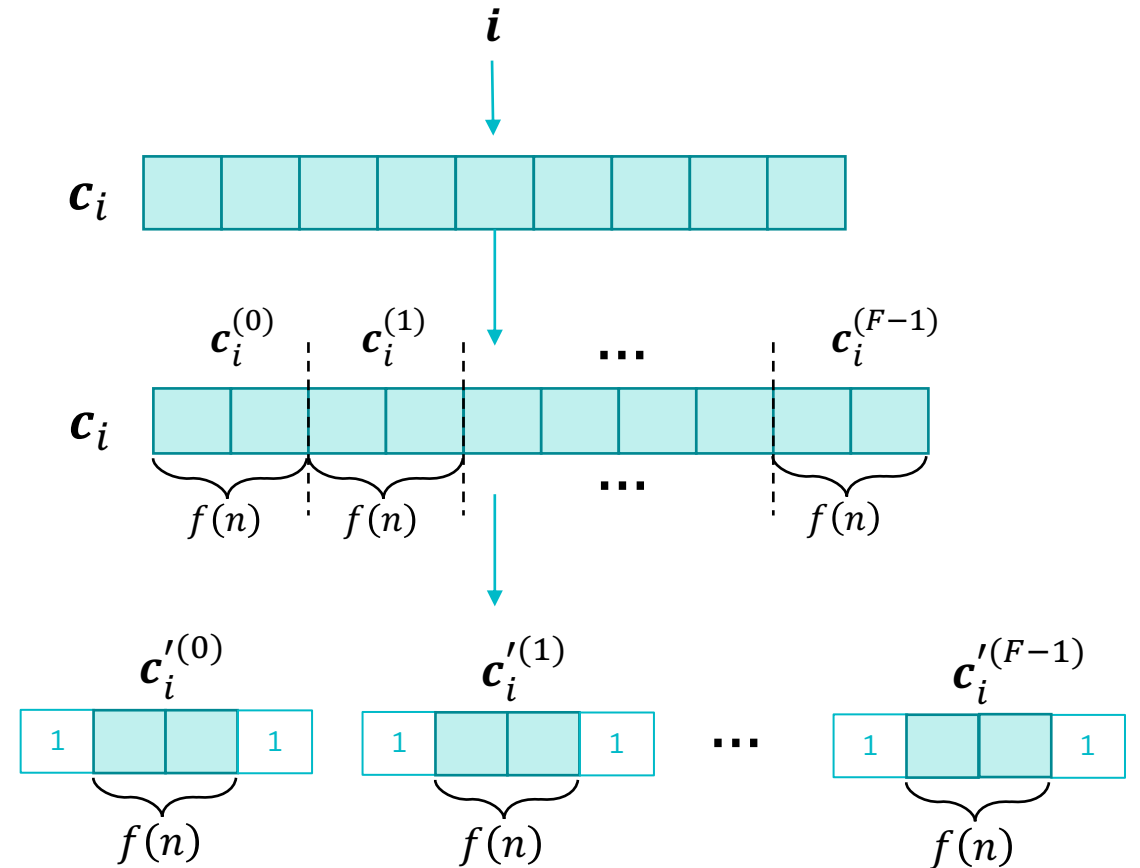
Construction for $(L_{\min}, L_{\text{over}})$ -Trace Codes

$$\begin{aligned}L_{\min} &= a \log(n) + O_n(1), & a > 1 \\L_{\text{over}} &= \gamma L_{\min} + O_n(1), & 0 < \gamma \leq \frac{1}{a}\end{aligned}$$

Preliminaries: Index Generation

Consider:

- Index length: $I \approx \left(1 - \frac{1-\gamma a}{1-\gamma}\right) L_{\min}$
- Block size: $f(n) = o(\log n)$
- # blocks: $F = \frac{I}{f(n)}$



Segments of the encoded index

Preliminaries: Repeat-Free Encoding

For $\ell < N$ the set of *repeat-free (RF)* strings is

$$\mathcal{RF}_\ell(N) \triangleq \{x \in \Sigma^N: \text{no } \ell\text{-substring repeats}\}.$$

Lemma: For $q > 2$ and integers $\ell(N) > \lceil \log(N) \rceil + 3\lceil \log \log(N) \rceil$ and $\lceil \log \log(N) \rceil < t \leq \left\lfloor \frac{\ell(N) - \lceil \log(N) \rceil}{3} \right\rfloor$, there exists an efficient **encoder/decoder** pair of $\mathcal{RF}_{\ell(N)}(N)$ that also does **not contain any t -length runs of zeros, with rate $1 - O_n\left(\frac{t}{n} + q^{-t}\right)$.**

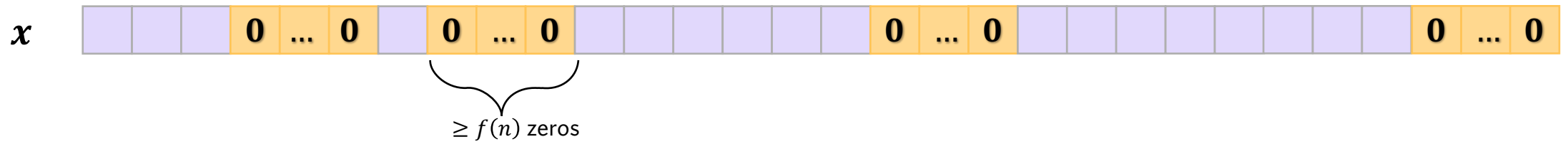
Define encoder $E_{m,\ell}^{\mathcal{RF}}$: input length m , output length $N = N_{n,\ell}(m)$, $t = f(n)$

Similar result was proven in [1] for specific values of $\ell(n)$ and t

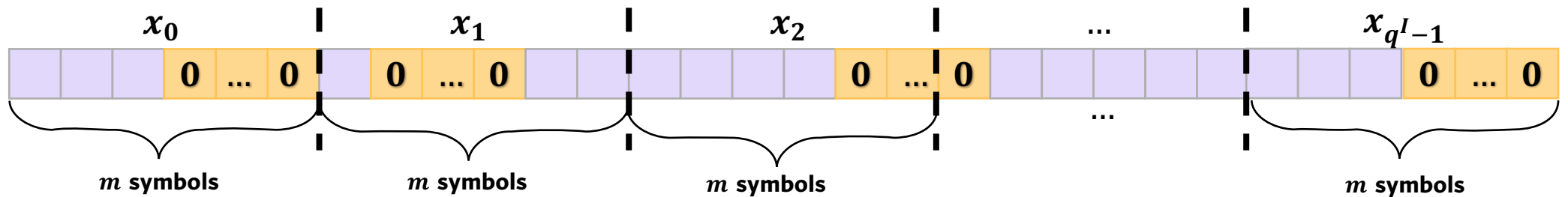
Construction D

Encoder for $(L_{\min}, L_{\text{over}})$ -trace code \mathcal{C}_D :

Input: a sequence $x \in \Sigma^{q^l m}$



Divide x into q^l non-overlapping substrings of length m .



Next, we encode each substring x_i into $z_i \in \Sigma^{n/q^l}$ independently.

Construction D

Encoder for $(L_{\min}, L_{\text{over}})$ -trace code \mathcal{C}_D :

For any $0 \leq i \leq q^l - 1$ we want \mathbf{z}_i to satisfy two properties:

- (1) the index i can be decoded from any L_{\min} -substring of \mathbf{z}_i , and
- (2) the string \mathbf{z}_i can be uniquely reconstructed from any $(L_{\min}, L_{\text{over}})$ -trace of \mathbf{z}_i .

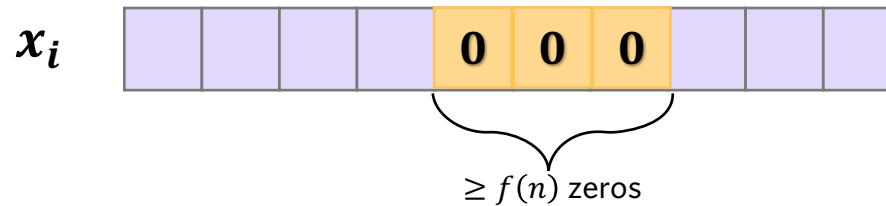
Then, we let

$$\text{Enc}_A(\mathbf{x}) \triangleq \mathbf{z} = \mathbf{z}_0 \circ \cdots \circ \mathbf{z}_{q^l-1}.$$

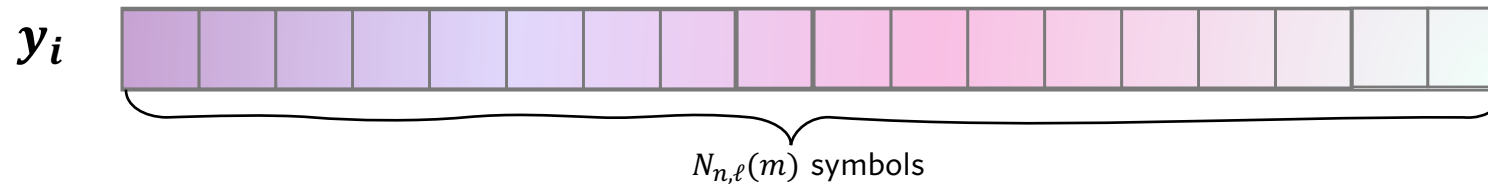
Construction D

Encoder for $(L_{\min}, L_{\text{over}})$ -trace code \mathcal{C}_D :

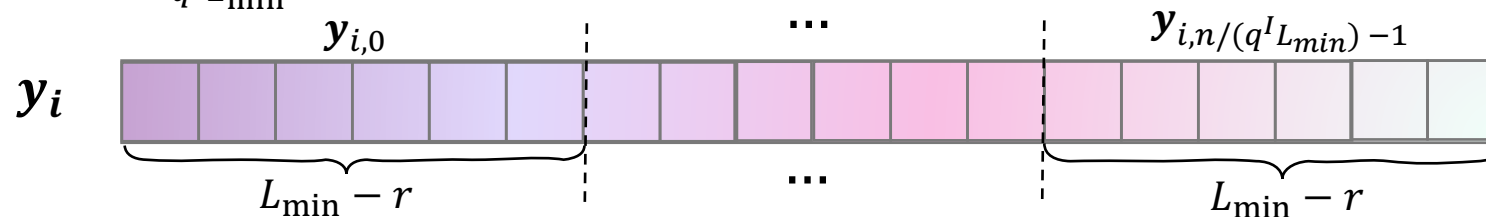
For any $0 \leq i \leq q^l - 1$:



(1) Encode \mathbf{x}_i using the RF encoder to obtain $\mathbf{y}_i = E_{m,\ell}^{\text{RF}}(\mathbf{x}_i)$.



(2) Partition \mathbf{y}_i into $\frac{n}{q^l L_{\min}}$ non-overlapping segments of length $L_{\min} - r$.



$$N_{n,\ell}(m) \approx \frac{1 - \gamma a}{1 - \gamma} n q^{-l}$$

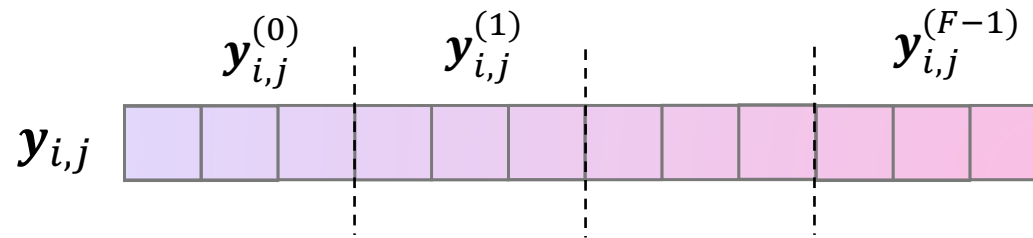
Construction D

Encoder for $(L_{\min}, L_{\text{over}})$ -trace code \mathcal{C}_D :

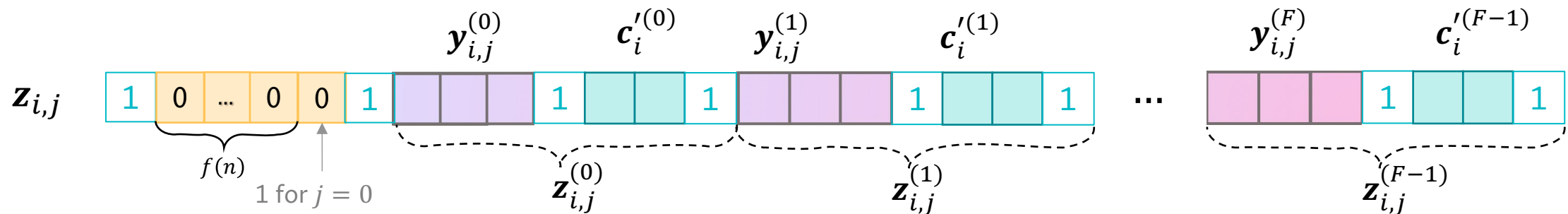
(3) For all $j \in [n/(q^l L_{\min})]$:



(3.1) Partition $y_{i,j}$ into F non-overlapping segments of equal lengths (up to ± 1 , if necessary).



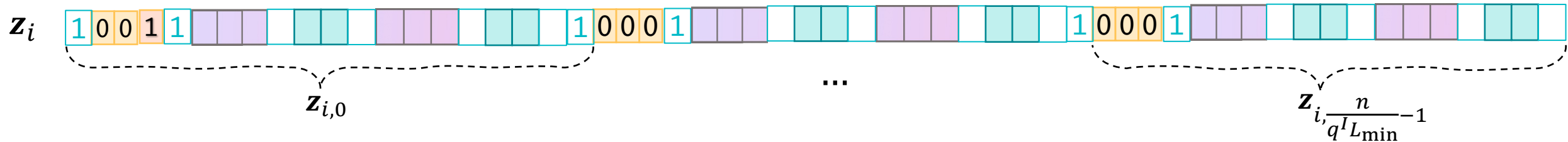
(3.2) For $k \in [F]$, let $z_{i,j}^{(k)} \triangleq y_{i,j}^{(k)} \circ c_i'^{(k)}$



Construction D

Encoder for $(L_{\min}, L_{\text{over}})$ -trace code \mathcal{C}_D :

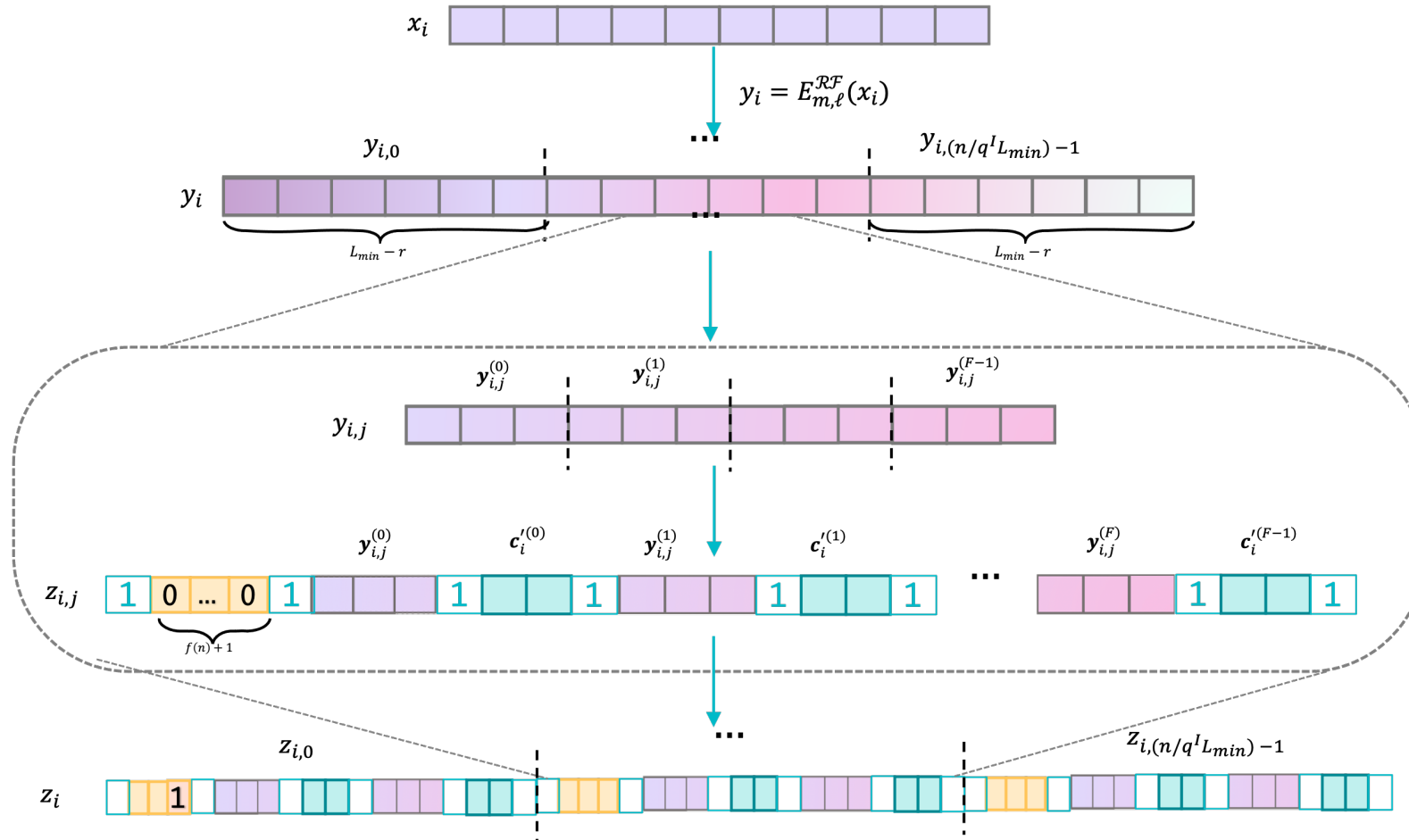
(4) Concatenate $\mathbf{z}_i = \mathbf{z}_{i,0} \circ \mathbf{z}_{i,1} \circ \dots \circ \mathbf{z}_{i, \frac{n}{q^I L_{\min}} - 1}$



Lastly, we concatenate all the \mathbf{z}_i strings, in order, to obtain the code word

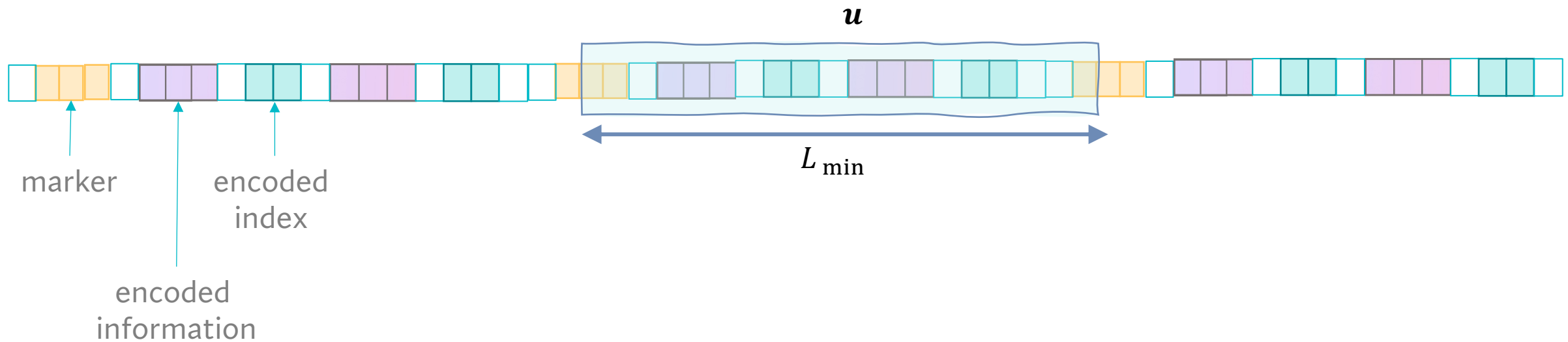
$$\text{Enc}_D(\mathbf{x}) \triangleq \mathbf{z} = \mathbf{z}_0 \circ \dots \circ \mathbf{z}_{q^I - 1}.$$

Construction D – Overview



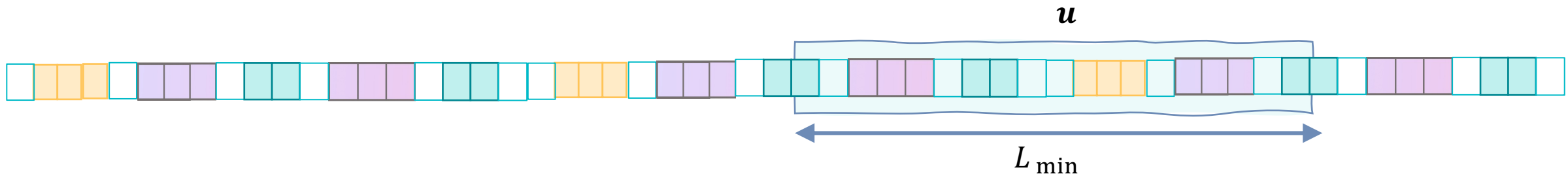
Construction D – Encoder Correctness

Lemma: Every L_{\min} -substring u of z contains as subsequences at least an $(I - \mu)$ -suffix of c_i and μ -prefix of either c_i or c_{i+1} for some $i \in [q^I]$ and $\mu \in [I]$, in identifiable locations.

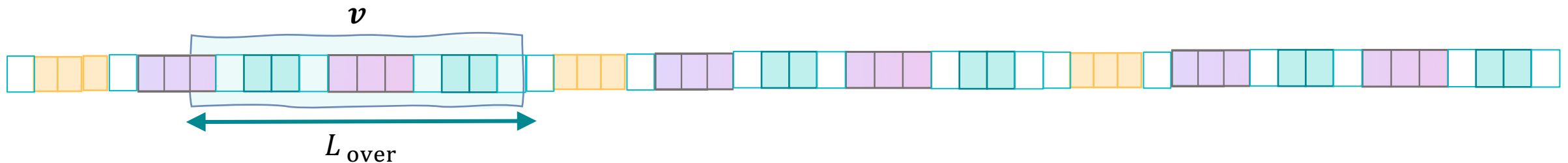


Construction D – Encoder Correctness

Lemma: Every L_{\min} -substring u of z contains as subsequences at least an $(I - \mu)$ -suffix of c_i and μ -prefix of either c_i or c_{i+1} for some $i \in [q^l]$ and $\mu \in [I]$, in identifiable locations.



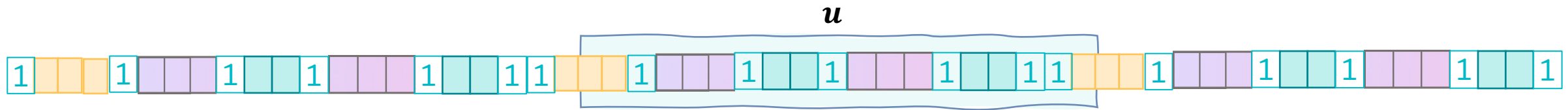
Lemma: Every L_{over} -substring v of z contains at least ℓ consecutive symbols of $y \triangleq y_0 \circ y_1 \circ \dots \circ y_{q^l-1}$.



Construction D – Encoder Correctness

Theorem : for all admissible values of n , the code \mathcal{C}_D is an $(L_{\min}, L_{\text{over}})$ -trace code.

Proof: Take $\mathbf{z} \in \mathcal{C}_D$ and let $T \in \mathcal{T}_{L_{\min}}^{L_{\text{over}}}$ be an $(L_{\min}, L_{\text{over}})$ -trace of \mathbf{z} .



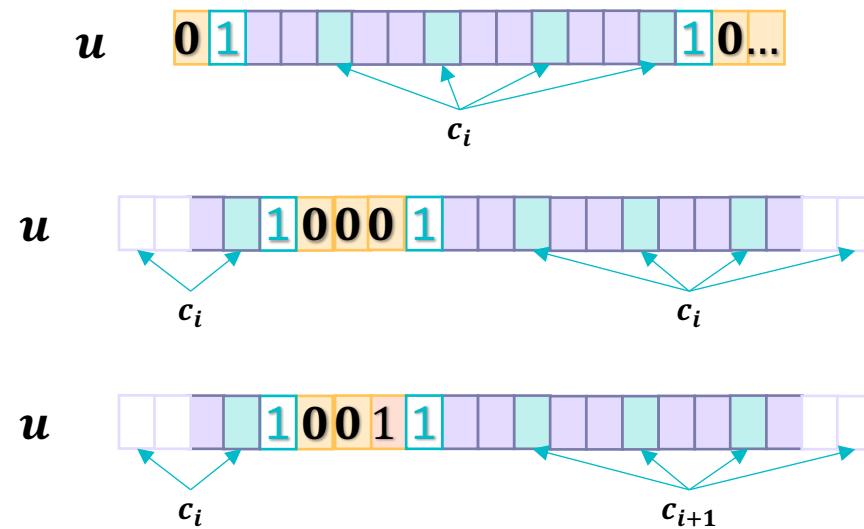
Let $\mathbf{u} \in T$, we can assume w.l.o.g. that $|\mathbf{u}| = L = L_{\min}$.

Note that:

1. The segment \mathbf{u} does not contain any occurrences of the marker $\boxed{1}00\boxed{1}$ except those explicitly added as a prefix of $\mathbf{z}_{i,j}$ (for some i, j).
2. For any i, j , $|\mathbf{z}_{i,j}| \leq L_{\min}$.
3. Either \mathbf{u} contains an occurrence of $\boxed{1}00\boxed{1}$ or it has a suffix-prefix pair whose concatenation is $\boxed{1}00\boxed{1}$.

Construction D – Encoder Correctness

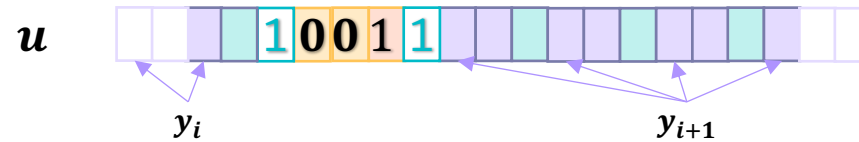
Either u contains an occurrence of $\boxed{1}00\boxed{1}$ or it has a suffix-prefix pair whose concatenation is $\boxed{1}00\boxed{1}$.



We can always correctly deduce i from u .

It is therefore possible to partition T by the index i (corresponding to the substring z_i).

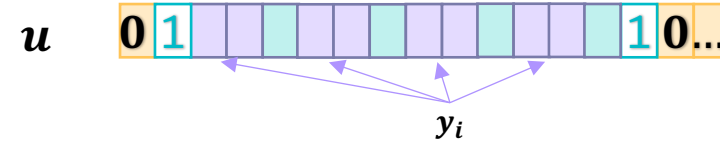
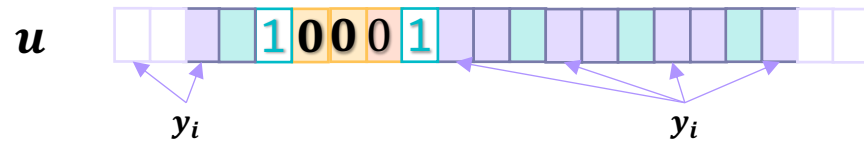
Construction D – Encoder Correctness



If u is intersecting both y_i, y_{i+1} , then u contains the complete synchronization marker 10011 hence its location in u implies the exact location of u in z .

10011

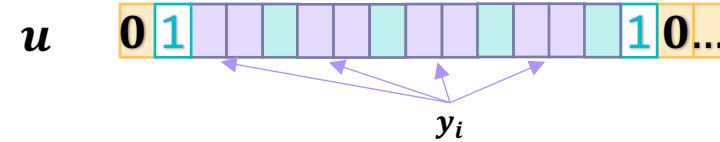
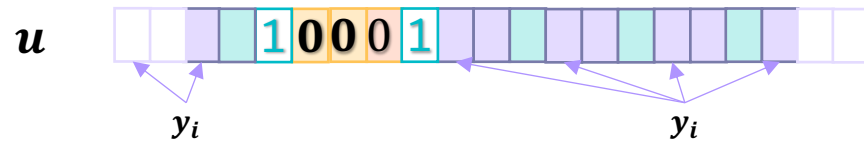
Construction D – Encoder Correctness



Lemma: Every L_{over} -substring v of z contains at least ℓ consecutive symbols of $y \triangleq y_0 \circ y_1 \circ \dots \circ y_{q^l-1}$.

Since each y_i is ℓ -repeat-free, there exist a unique way to match the overlaps of these substrings.

Construction D – Encoder Correctness



Lemma: Every L_{over} -substring v of z contains at least ℓ consecutive symbols of $y \triangleq y_0 \circ y_1 \circ \dots \circ y_{q^I-1}$.

Since each y_i is ℓ -repeat-free, there exist a unique way to match the overlaps of these substrings.

Finally, once z is reconstructed we may extract $\{y_i\}_{i \in [q^I]}$, then decode $\{x_i\}_{i \in [q^I]}$ with the decoder of $E_{m,\ell}^{\mathcal{RF}}$.

Construction D – Rate

Theorem: Letting $f(n) \triangleq \lceil \sqrt{\log(n)} \rceil$ we have that

$$R(\mathcal{C}_D) \geq \frac{1 - \frac{1}{a}}{1 - \gamma} - \frac{\frac{1}{a}}{(\log(n))^{0.5-\epsilon}} - o\left(\frac{1}{\sqrt{\log(n)}}\right).$$

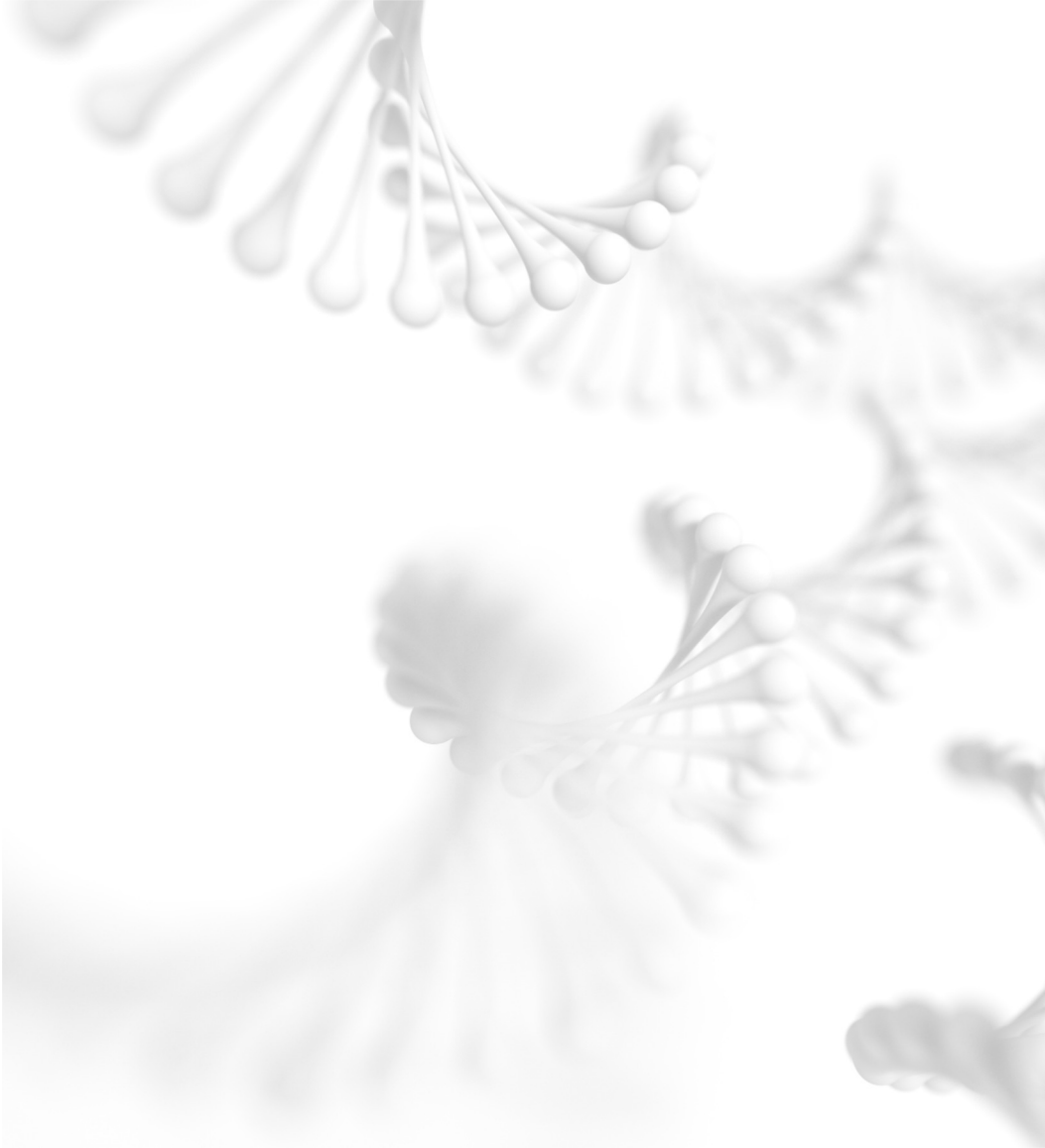
Maximum Asymptotic Rate of $(L_{\min}, L_{\text{over}})$ -Trace Codes

Lemma: If $L_{\min} = a \log(n) + O_n(1)$ and $L_{\text{over}} = \gamma L_{\min} + O_n(1)$, for some $a > 1$ and $0 < \gamma \leq \frac{1}{a}$, then any $(L_{\min}, L_{\text{over}})$ -trace code $\mathcal{C} \subseteq \Sigma^n$ satisfies

$$R(\mathcal{C}) \leq \frac{1 - \frac{1}{a}}{1 - \gamma} + o\left(\frac{\log \log(n)}{\log(n)}\right).$$

Hence, Construction D asymptotically meets the upper bound.

Corollary: If $\limsup_{n \rightarrow \infty} \frac{L_{\min}}{\log(n)} \leq 1$, then $R(\mathcal{C}) = o_n(1)$ for any $(L_{\min}, L_{\text{over}})$ -trace code $\mathcal{C} \subseteq \Sigma^n$.



Future Work

- Erroneous versions of the channels.
 - Multistrand setup for the partial overlap channel.
 - Analyze the worst-case for a more realistic setup.
-

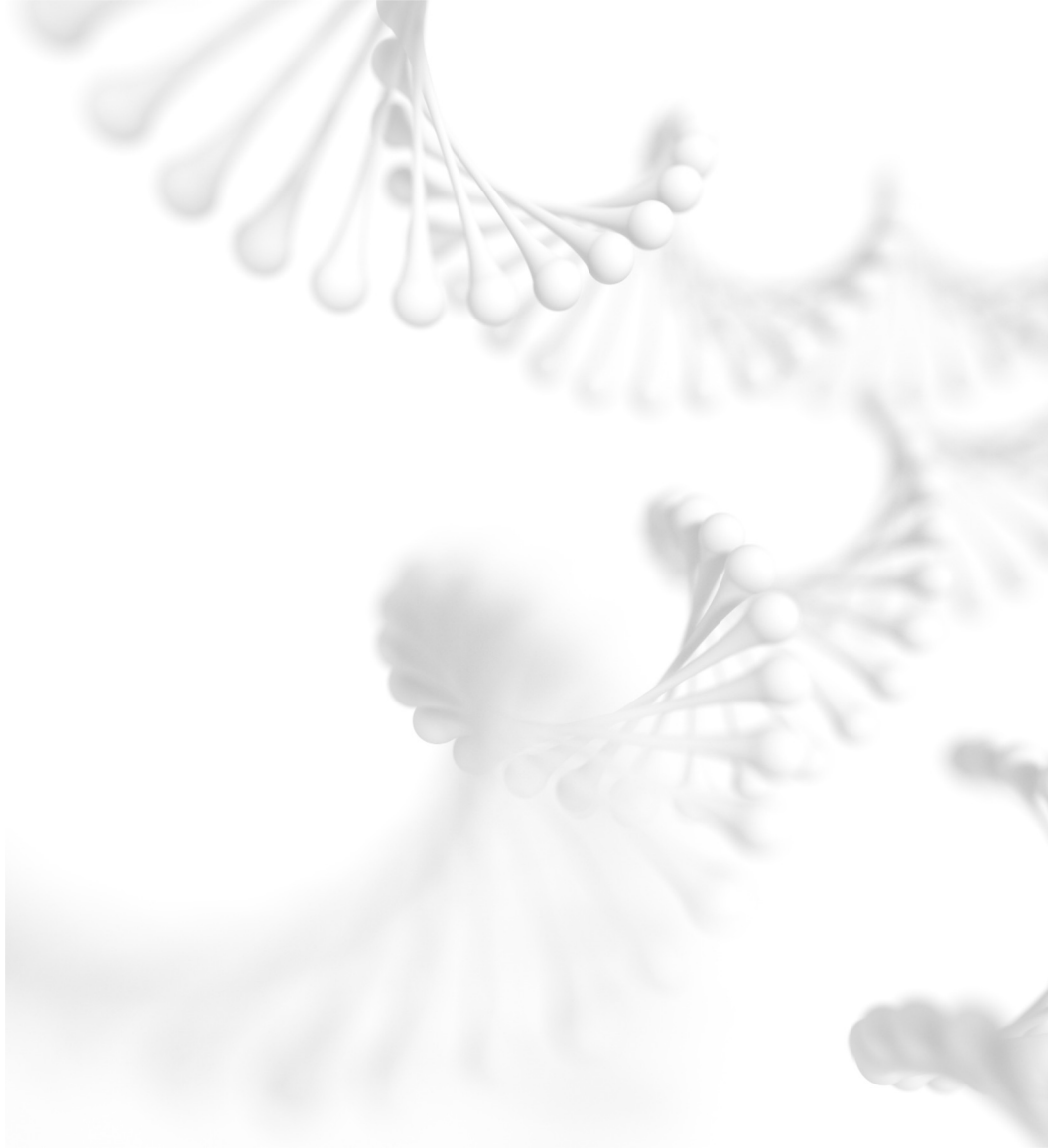


References

- [1] **D. Bar-Lev**, S. Marcovich, E. Yaakobi, and Y. Yehezkeally, “Adversarial Torn-Paper Codes”, submitted to *IEEE Transactions on Information Theory*.
- [2] Y. Yehezkeally, **D. Bar-Lev**, S. Marcovich, and E. Yaakobi, “Generalized Unique Reconstruction from Substrings”, submitted to *IEEE Transactions on Information Theory*.

and also

- [3] Y. Yehezkeally, S. Marcovich, and E. Yaakobi, “Multi-strand reconstruction from substrings,” in *IEEE ITW*, 2021
 - [4] **D. Bar-Lev**, S. Marcovich, E. Yaakobi, Y. Yehezkeally, “Adversarial Torn-Paper Codes,” in *IEEE ISIT*, 2022.
 - [5] Y. Yehezkeally, **D. Bar-Lev**, S. Marcovich, and E. Yaakobi, “Reconstruction from Substrings with Partial Overlap”, in *IEEE ISITA.*, 2022.
-



Thank You!