

Insertion and Deletion Correction in Polymer-based Data Storage

Anisha Banerjee¹

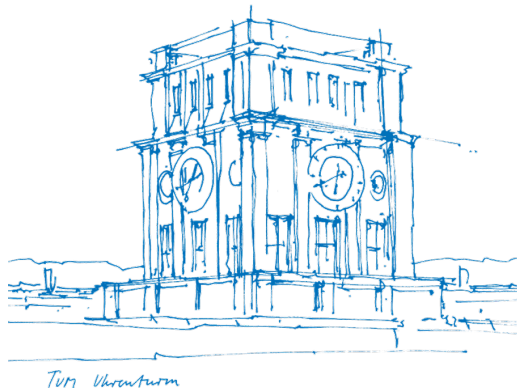
Joint work with

Antonia Wachter-Zeh¹, Eitan Yaakobi²

¹Technical University of Munich
Institute for Communications Engineering

²Technion – Israel Institute of Technology
Department of Computer Science

June 28, 2022



Introduction

Problem Statement & Results

Conclusion

Motivation

- Need for high-density and reliable data storage media.
 - ▶ Molecular storage paradigms, e.g. DNA storage



[OACL17] [A. Al Ouahabi](#), [J.-A. Amalian](#), [L. Charles](#), and [J.-F. Lutz](#), “Mass spectrometry sequencing of long digital polymers facilitated by programmed inter-byte fragmentation”, *Nature Communications*, vol. 8, no. 1, p. 967, Dec. 2017

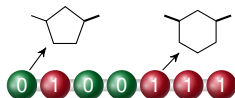
Motivation

- Need for high-density and reliable data storage media.
 - ▶ Molecular storage paradigms, e.g. DNA storage
- One potential candidate is polymer-based data storage[OACL17].
 - ▶ Low readout latency.
 - ▶ Low storage cost.
 - ▶ Common analytical equipment usable.



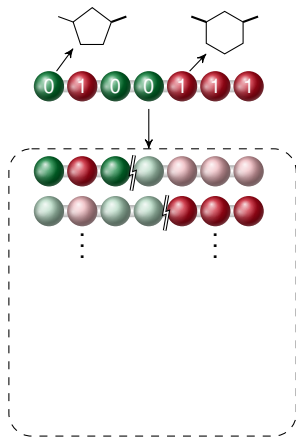
[OACL17] [A. Al Ouahabi](#), [J.-A. Amalian](#), [L. Charles](#), and [J.-F. Lutz](#), “Mass spectrometry sequencing of long digital polymers facilitated by programmed inter-byte fragmentation”, *Nature Communications*, vol. 8, no. 1, p. 967, Dec. 2017

Background



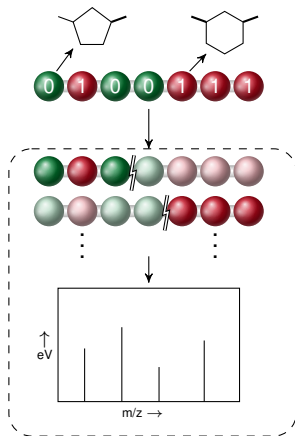
- Encode bits into synthetic macromolecules of distinct masses.

Background



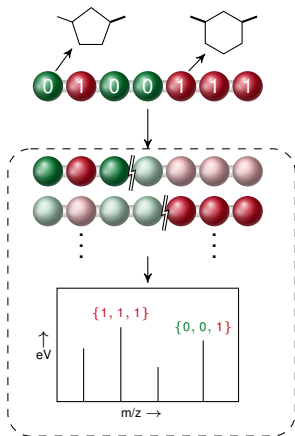
- Encode bits into synthetic macromolecules of distinct masses.
- Read using tandem mass spectrometer.
 - ▶ Breaks links & generates fragments.

Background



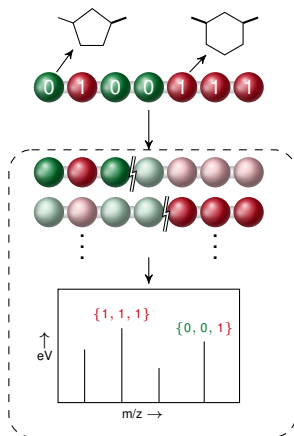
- Encode bits into synthetic macromolecules of distinct masses.
- Read using tandem mass spectrometer.
 - ▶ Breaks links & generates fragments.
 - ▶ Reports masses of fragments via peak series.

Background



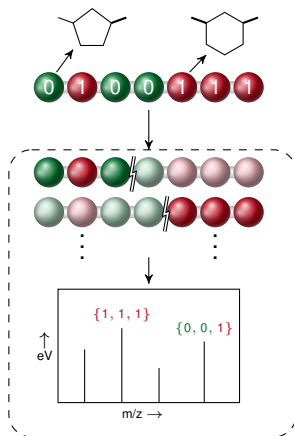
- Encode bits into synthetic macromolecules of distinct masses.
- Read using tandem mass spectrometer.
 - ▶ Breaks links & generates fragments.
 - ▶ Reports masses of fragments via peak series.
 - ▶ Mass suggests number of 0s and 1s in fragment.

Background



- Encode bits into synthetic macromolecules of distinct masses.
- Read using tandem mass spectrometer.
 - ▶ Breaks links & generates fragments.
 - ▶ Reports masses of fragments via peak series.
 - ▶ Mass suggests number of 0s and 1s in fragment.
composition

Background



- Encode bits into synthetic macromolecules of distinct masses.
- Read using tandem mass spectrometer.
 - ▶ Breaks links & generates fragments.
 - ▶ Reports masses of fragments via peak series.
 - ▶ Mass suggests number of 0s and 1s in fragment.
composition
- Assume we get compositions of all contiguous substrings.

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s})$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\},$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\},$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\}, \{0, 1\}\}.$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\}, \{0, 1\}\}.$$

$$C(\mathbf{s})$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\}, \{0, 1\}\}.$$

$$C(\mathbf{s}) = \{\{0\}, \{1\}, \{0\}, \{1\}, \leftarrow C_1(\mathbf{s})\}$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\}, \{0, 1\}\}.$$

$$C(\mathbf{s}) = \{\{0\}, \{1\}, \{0\}, \{1\}, \\ \{0, 1\}, \{1, 0\}, \{0, 1\}, \leftarrow C_2(\mathbf{s})\}$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\}, \{0, 1\}\}.$$

$$C(\mathbf{s}) = \{\{0\}, \{1\}, \{0\}, \{1\}, \\ \{0, 1\}, \{1, 0\}, \{0, 1\}, \\ \{0, 1, 0\}, \{1, 0, 1\}, \leftarrow C_3(\mathbf{s})$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\}, \{0, 1\}\}.$$

$$\begin{aligned} C(\mathbf{s}) = & \{\{0\}, \{1\}, \{0\}, \{1\}, \\ & \{0, 1\}, \{1, 0\}, \{0, 1\}, \\ & \{0, 1, 0\}, \{1, 0, 1\}, \{0, 1, 0, 1\} \leftarrow C_4(\mathbf{s}) \end{aligned}$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\}, \{0, 1\}\}.$$

$$C(\mathbf{s}) = \{\{0\}, \{1\}, \{0\}, \{1\}, \\ \{0, 1\}, \{1, 0\}, \{0, 1\}, \\ \{0, 1, 0\}, \{1, 0, 1\}, \{0, 1, 0, 1\}\}.$$

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\}, \{0, 1\}\}.$$

$$\begin{aligned} C(\mathbf{s}) = & \{\{0\}, \{1\}, \{0\}, \{1\}, \\ & \{0, 1\}, \{1, 0\}, \{0, 1\}, \\ & \{0, 1, 0\}, \{1, 0, 1\}, \{0, 1, 0, 1\}\}. \end{aligned}$$

Note: If $\tilde{\mathbf{s}} = (1, 0, 1, 0)$, $C(\tilde{\mathbf{s}}) = C(\mathbf{s})$.

Model

- For $\mathbf{s} \in \{0, 1\}^n$, define **multiset** $C_k(\mathbf{s})$ as

$$C_k(\mathbf{s}) = \{\{s_i, s_{i+1}, \dots, s_{i+k-1}\} : 1 \leq i \leq n - k + 1\}$$

- Define **composition multiset** $C(\mathbf{s})$ as

$$C(\mathbf{s}) = \{C_k(\mathbf{s}) : 1 \leq k \leq n\}$$

Example

Consider $\mathbf{s} = (0, 1, 0, 1)$.

$$C_2(\mathbf{s}) = \{\{0, 1\}, \{1, 0\}, \{0, 1\}\}.$$

$$\begin{aligned} C(\mathbf{s}) = & \{\{0\}, \{1\}, \{0\}, \{1\}, \\ & \{0, 1\}, \{1, 0\}, \{0, 1\}, \\ & \{0, 1, 0\}, \{1, 0, 1\}, \{0, 1, 0, 1\}\}. \end{aligned}$$

Note: If $\tilde{\mathbf{s}} = (1, 0, 1, 0)$, $C(\tilde{\mathbf{s}}) = C(\mathbf{s})$. → A string and its reversal have identical composition multisets!

Previous Work

- [ADMOP15] showed that strings of length one less than a prime or twice a prime are uniquely reconstructable, up to reversal.

[ADMOP15] J. Acharya, H. Das, O. Milenkovic, A. Orlitsky, and S. Pan, “String Reconstruction from Substring Compositions”, *SIAM Journal on Discrete Mathematics*, vol. 29, no. 3, pp. 1340–1371, Jan. 2015

Previous Work

- [ADMOP15] showed that strings of length one less than a prime or twice a prime are uniquely reconstructable, up to reversal.
- [PGM19] proposed a code of $O(\log n)$ redundancy to allow unique reconstruction for all lengths.

Construction

$$\begin{aligned} \mathcal{S}_R(n) = \{ & \mathbf{s} \in \{0, 1\}^n, s_1 = 0, s_n = 1, \text{ and} \\ & \exists I \subset \{2, \dots, n-1\} \text{ such that} \\ & \text{for all } i \in I, s_i \neq s_{n+1-i}, \text{ for all } i \notin I, s_i = s_{n+1-i}, \\ & \mathbf{s}_{\lfloor n/2 \rfloor \cap I} \text{ has each prefix with strictly more 0s than 1s.} \end{aligned}$$

[PGM19] [S. Pattabiraman, R. Gabrys, and O. Milenkovic](#), “Reconstruction and Error-Correction Codes for Polymer-Based Data Storage”, in *2019 IEEE Information Theory Workshop (ITW)*, Visby, Sweden: IEEE, Aug. 2019, pp. 1–5

Previous Work

- [GPM20] proposed codes to correct composition substitutions.

[GPM20] R. Gabrys, S. Pattabiraman, and O. Milenkovic, “Mass Error-Correction Codes for Polymer-Based Data Storage”, in *2020 IEEE International Symposium on Information Theory (ISIT)*, Los Angeles, CA, USA: IEEE, Jun. 2020, pp. 25–30

[GPM20a] R. Gabrys, S. Pattabiraman, and O. Milenkovic, *Reconstructing Mixtures of Coded Strings from Prefix and Suffix Compositions*, Number: arXiv:2010.11116 arXiv:2010.11116 [cs, math], Oct. 2020

[GM22] U. Gupta and H. Mahdavifar, “A New Algebraic Approach for String Reconstruction from Substring Compositions”, *arXiv:2201.09955 [cs, math]*, Jan. 2022, arXiv: 2201.09955

Previous Work

- [GPM20] proposed codes to correct composition substitutions.
- [GPM20a] introduced codes to reconstruct multiple strings from a mixture of prefix and suffix compositions, simultaneously.

[GPM20] R. Gabrys, S. Pattabiraman, and O. Milenkovic, “Mass Error-Correction Codes for Polymer-Based Data Storage”, in *2020 IEEE International Symposium on Information Theory (ISIT)*, Los Angeles, CA, USA: IEEE, Jun. 2020, pp. 25–30

[GPM20a] R. Gabrys, S. Pattabiraman, and O. Milenkovic, *Reconstructing Mixtures of Coded Strings from Prefix and Suffix Compositions*, Number: arXiv:2010.11116 arXiv:2010.11116 [cs, math], Oct. 2020

[GM22] U. Gupta and H. MahdaviFar, “A New Algebraic Approach for String Reconstruction from Substring Compositions”, *arXiv:2201.09955 [cs, math]*, Jan. 2022, arXiv: 2201.09955

Previous Work

- [GPM20] proposed codes to correct composition substitutions.
- [GPM20a] introduced codes to reconstruct multiple strings from a mixture of prefix and suffix compositions, simultaneously.
- [GM22] devised a new algebraic algorithm to reconstruct strings from their composition multisets.

[GPM20] R. Gabrys, S. Pattabiraman, and O. Milenkovic, “Mass Error-Correction Codes for Polymer-Based Data Storage”, in *2020 IEEE International Symposium on Information Theory (ISIT)*, Los Angeles, CA, USA: IEEE, Jun. 2020, pp. 25–30

[GPM20a] R. Gabrys, S. Pattabiraman, and O. Milenkovic, *Reconstructing Mixtures of Coded Strings from Prefix and Suffix Compositions*, Number: [arXiv:2010.11116](#) [arXiv:2010.11116 \[cs, math\]](#), Oct. 2020

[GM22] U. Gupta and H. MahdaviFar, “A New Algebraic Approach for String Reconstruction from Substring Compositions”, [arXiv:2201.09955 \[cs, math\]](#), Jan. 2022, [arXiv: 2201.09955](#)

Introduction

Problem Statement & Results

Conclusion

Problem Statement

1. Are composition-insertion-correcting codes equivalent to composition-deletion-correcting codes?

Problem Statement



1. Are composition-insertion-correcting codes equivalent to composition-deletion-correcting codes?

→ Yes!

Problem Statement

1. Are composition-insertion-correcting codes equivalent to composition-deletion-correcting codes?
→ Yes!
2. How robust is $\mathcal{S}_R(n)$ to insertions and deletions in composition multisets?

Problem Statement

1. Are composition-insertion-correcting codes equivalent to composition-deletion-correcting codes?
→ Yes!
2. How robust is $\mathcal{S}_R(n)$ to insertions and deletions in composition multisets?
→ Can correct deletion of at most one complete multiset.

Problem Statement

1. Are composition-insertion-correcting codes equivalent to composition-deletion-correcting codes?
→ Yes!
2. How robust is $\mathcal{S}_R(n)$ to insertions and deletions in composition multisets?
→ Can correct deletion of at most one complete multiset.
→ Similar construction exists to correct deletions of more multisets.

Insertion-Deletion Equivalence



Lemma

A code can correct the deletion of t multisets, if and only if it can correct composition insertion errors in those t multisets.

Insertion-Deletion Equivalence

Lemma

A code can correct the deletion of t multisets, if and only if it can correct composition insertion errors in those t multisets.

Proof idea:

- Define **t -multiset deletion ball** of \mathbf{s} as

$$D_t(\mathbf{s}) = \{C'(\mathbf{s}) = \bigcup_{i \in \{1, \dots, n\} \setminus \mathcal{I}} C_i(\mathbf{s}) : \mathcal{I} \subseteq \{1, \dots, n\}, |\mathcal{I}| \leq t\}$$

Insertion-Deletion Equivalence

Lemma

A code can correct the deletion of t multisets, if and only if it can correct composition insertion errors in those t multisets.

Proof idea:

- Define t -multiset deletion ball of \mathbf{s} as

$$D_t(\mathbf{s}) = \{C'(\mathbf{s}) = \bigcup_{i \in \{1, \dots, n\} \setminus \mathcal{I}} C_i(\mathbf{s}) : \mathcal{I} \subseteq \{1, \dots, n\}, |\mathcal{I}| \leq t\}$$

- Consider a \mathbf{v} such that $D_t(\mathbf{s}) \cap D_t(\mathbf{v}) \neq \emptyset$. So, for some $\mathcal{I} \subseteq \{1, \dots, n\}$ where $|\mathcal{I}| \leq t$,

$$C_j(\mathbf{s}) \neq C_j(\mathbf{v}) \quad \forall j \in \mathcal{I}$$

Insertion-Deletion Equivalence

Lemma

A code can correct the deletion of t multisets, if and only if it can correct composition insertion errors in those t multisets.

Proof idea:

- Define t -multiset deletion ball of \mathbf{s} as

$$D_t(\mathbf{s}) = \{C'(\mathbf{s}) = \bigcup_{i \in \{1, \dots, n\} \setminus \mathcal{I}} C_i(\mathbf{s}) : \mathcal{I} \subseteq \{1, \dots, n\}, |\mathcal{I}| \leq t\}$$

- Consider a \mathbf{v} such that $D_t(\mathbf{s}) \cap D_t(\mathbf{v}) \neq \emptyset$. So, for some $\mathcal{I} \subseteq \{1, \dots, n\}$ where $|\mathcal{I}| \leq t$,

$$C_j(\mathbf{s}) \neq C_j(\mathbf{v}) \quad \forall j \in \mathcal{I}$$

$$C_j(\mathbf{s}) = C_j(\mathbf{v}) \quad \forall j \in \{1, \dots, n\} \setminus \mathcal{I}$$

Insertion-Deletion Equivalence

- Let $C'_k(\mathbf{s}) = C_k(\mathbf{s}) \cup C_k(\mathbf{v}) = C'_k(\mathbf{v})$.

$C_k(\mathbf{s})$ corrupted by insertion of compositions $C_k(\mathbf{v}) \setminus C_k(\mathbf{s})$.

Insertion-Deletion Equivalence

- Let $C'_k(\mathbf{s}) = C_k(\mathbf{s}) \cup C_k(\mathbf{v}) = C'_k(\mathbf{v})$.

$C_k(\mathbf{s})$ corrupted by insertion of compositions $C_k(\mathbf{v}) \setminus C_k(\mathbf{s})$.

- Define **t -multiset insertion ball** as

$$I_t(\mathbf{s}) = \{C'(\mathbf{s}) = \bigcup_{i \in \{1, \dots, n\}} C'_i(\mathbf{s}) : \mathcal{I} \subseteq \{1, \dots, n\}, |\mathcal{I}| \leq t, \\ C_i(\mathbf{s}) \subset C'_i(\mathbf{s}) \ \forall i \in \mathcal{I}, C'_i(\mathbf{s}) = C_i(\mathbf{s}) \ \forall i \in \{1, \dots, n\} \setminus \mathcal{I}\}.$$

Insertion-Deletion Equivalence

- Let $C'_k(\mathbf{s}) = C_k(\mathbf{s}) \cup C_k(\mathbf{v}) = C'_k(\mathbf{v})$.

$C_k(\mathbf{s})$ corrupted by insertion of compositions $C_k(\mathbf{v}) \setminus C_k(\mathbf{s})$.

- Define t -multiset insertion ball as

$$I_t(\mathbf{s}) = \{C'(\mathbf{s}) = \bigcup_{i \in \{1, \dots, n\}} C'_i(\mathbf{s}) : \mathcal{I} \subseteq \{1, \dots, n\}, |\mathcal{I}| \leq t, \\ C_i(\mathbf{s}) \subset C'_i(\mathbf{s}) \ \forall i \in \mathcal{I}, C'_i(\mathbf{s}) = C_i(\mathbf{s}) \ \forall i \in \{1, \dots, n\} \setminus \mathcal{I}\}.$$

- Thus

$$\bigcup_{k \in \mathcal{I}} C'_k(\mathbf{s}) \cup \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{s}) = \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{v}) \cup \bigcup_{k \in \mathcal{I}} C'_k(\mathbf{v}).$$

Insertion-Deletion Equivalence

- Let $C'_k(\mathbf{s}) = C_k(\mathbf{s}) \cup C_k(\mathbf{v}) = C'_k(\mathbf{v})$.

$C_k(\mathbf{s})$ corrupted by insertion of compositions $C_k(\mathbf{v}) \setminus C_k(\mathbf{s})$.

- Define t -multiset insertion ball as

$$I_t(\mathbf{s}) = \{C'(\mathbf{s}) = \bigcup_{i \in \{1, \dots, n\}} C'_i(\mathbf{s}) : \mathcal{I} \subseteq \{1, \dots, n\}, |\mathcal{I}| \leq t, \\ C_i(\mathbf{s}) \subset C'_i(\mathbf{s}) \ \forall i \in \mathcal{I}, C'_i(\mathbf{s}) = C_i(\mathbf{s}) \ \forall i \in \{1, \dots, n\} \setminus \mathcal{I}\}.$$

- Thus

$$\bigcup_{k \in \mathcal{I}} C'_k(\mathbf{s}) \cup \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{s}) = \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{v}) \cup \bigcup_{k \in \mathcal{I}} C'_k(\mathbf{v}).$$

→ Equal by definition of $D_t(\mathbf{s})$.

Insertion-Deletion Equivalence

- Let $C'_k(\mathbf{s}) = C_k(\mathbf{s}) \cup C_k(\mathbf{v}) = C'_k(\mathbf{v})$.

$C_k(\mathbf{s})$ corrupted by insertion of compositions $C_k(\mathbf{v}) \setminus C_k(\mathbf{s})$.

- Define t -multiset insertion ball as

$$I_t(\mathbf{s}) = \{C'(\mathbf{s}) = \bigcup_{i \in \{1, \dots, n\}} C'_i(\mathbf{s}) : \mathcal{I} \subseteq \{1, \dots, n\}, |\mathcal{I}| \leq t, \\ C_i(\mathbf{s}) \subset C'_i(\mathbf{s}) \ \forall i \in \mathcal{I}, C'_i(\mathbf{s}) = C_i(\mathbf{s}) \ \forall i \in \{1, \dots, n\} \setminus \mathcal{I}\}.$$

- Thus

$$\bigcup_{k \in \mathcal{I}} C'_k(\mathbf{s}) \cup \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{s}) = \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{v}) \cup \bigcup_{k \in \mathcal{I}} C'_k(\mathbf{v}).$$

→ Equal since $C'_k(\mathbf{s}) = C'_k(\mathbf{v})$.

Insertion-Deletion Equivalence

- Let $C'_k(\mathbf{s}) = C_k(\mathbf{s}) \cup C_k(\mathbf{v}) = C'_k(\mathbf{v})$.

$C_k(\mathbf{s})$ corrupted by insertion of compositions $C_k(\mathbf{v}) \setminus C_k(\mathbf{s})$.

- Define t -multiset insertion ball as

$$I_t(\mathbf{s}) = \{C'(\mathbf{s}) = \bigcup_{i \in \{1, \dots, n\}} C'_i(\mathbf{s}) : \mathcal{I} \subseteq \{1, \dots, n\}, |\mathcal{I}| \leq t, \\ C_i(\mathbf{s}) \subset C'_i(\mathbf{s}) \ \forall i \in \mathcal{I}, C'_i(\mathbf{s}) = C_i(\mathbf{s}) \ \forall i \in \{1, \dots, n\} \setminus \mathcal{I}\}.$$

- Thus

$$\bigcup_{k \in \mathcal{I}} C'_k(\mathbf{s}) \cup \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{s}) = \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{v}) \cup \bigcup_{k \in \mathcal{I}} C'_k(\mathbf{v}). \\ \in I_t(\mathbf{s}) \cap I_t(\mathbf{v})$$

Insertion-Deletion Equivalence

- Let $C'_k(\mathbf{s}) = C_k(\mathbf{s}) \cup C_k(\mathbf{v}) = C'_k(\mathbf{v})$.

$C_k(\mathbf{s})$ corrupted by insertion of compositions $C_k(\mathbf{v}) \setminus C_k(\mathbf{s})$.

- Define t -multiset insertion ball as

$$I_t(\mathbf{s}) = \{C'(\mathbf{s}) = \bigcup_{i \in \{1, \dots, n\}} C'_i(\mathbf{s}) : \mathcal{I} \subseteq \{1, \dots, n\}, |\mathcal{I}| \leq t, \\ C_i(\mathbf{s}) \subset C'_i(\mathbf{s}) \ \forall i \in \mathcal{I}, C'_i(\mathbf{s}) = C_i(\mathbf{s}) \ \forall i \in \{1, \dots, n\} \setminus \mathcal{I}\}.$$

- Thus

$$\bigcup_{k \in \mathcal{I}} C'_k(\mathbf{s}) \cup \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{s}) = \bigcup_{j \in \{1, \dots, n\} \setminus \mathcal{I}} C_j(\mathbf{v}) \cup \bigcup_{k \in \mathcal{I}} C'_k(\mathbf{v}). \\ \in I_t(\mathbf{s}) \cap I_t(\mathbf{v})$$

- $D_t(\mathbf{s}) \cap D_t(\mathbf{v}) \neq \emptyset \iff I_t(\mathbf{s}) \cap I_t(\mathbf{v}) \neq \emptyset$

Multiset Deletions

Definition

An **asymmetric multiset deletion** occurs in the composition multiset $C(\mathbf{s})$ of a string $\mathbf{s} \in \{0, 1\}^n$, if for some $i \in [n]$, the multiset $C_i(\mathbf{s})$ is entirely missing, while $C_{n-i+1}(\mathbf{s})$ is not corrupted.

Multiset Deletions

Definition

An **asymmetric multiset deletion** occurs in the composition multiset $C(\mathbf{s})$ of a string $\mathbf{s} \in \{0, 1\}^n$, if for some $i \in [n]$, the multiset $C_i(\mathbf{s})$ is entirely missing, while $C_{n-i+1}(\mathbf{s})$ is not corrupted.

Lemma

$S_R(n)$ can correct a single asymmetric multiset deletion.

Asymmetric Multiset Deletions

Construction

$$\begin{aligned} \mathcal{S}_{DA}^{(t)}(n) = & \{ \mathbf{s} \in \{0, 1\}^n : s_1 = 0, s_n = 1, \text{ and} \\ & \exists I \subset \{2, \dots, \frac{n}{2}\}, |I| \geq t, \text{ such that} \\ & \forall i \in I, s_i \neq s_{n+1-i}, \text{ and } \forall i \notin I, s_i = s_{n+1-i}, \\ & \mathbf{s}_{[n/2] \cap I} \text{ is a string wherein each} \\ & \text{prefix has at least } t \text{ more 0s than 1s} \}. \end{aligned}$$

Asymmetric Multiset Deletions

Construction

$$\begin{aligned} \mathcal{S}_{DA}^{(t)}(n) = \{ & \mathbf{s} \in \{0, 1\}^n : \mathbf{s}_1 = 0, \mathbf{s}_n = 1, \text{ and} \\ & \exists I \subset \{2, \dots, \frac{n}{2}\}, |I| \geq t, \text{ such that} \\ & \forall i \in I, \mathbf{s}_i \neq \mathbf{s}_{n+1-i}, \text{ and } \forall i \notin I, \mathbf{s}_i = \mathbf{s}_{n+1-i}, \\ & \mathbf{s}_{[n/2] \cap I} \text{ is a string wherein each} \\ & \text{prefix has at least } t \text{ more 0s than 1s} \}. \end{aligned}$$

- Is an extension of $\mathcal{S}_R(n)$ [PGM19] .
- Involves at most $0.5 \log(n - 2t) + 2t + 3$ redundant bits.

[PGM19] S. Pattabiraman, R. Gabrys, and O. Milenkovic, "Reconstruction and Error-Correction Codes for Polymer-Based Data Storage", in *2019 IEEE Information Theory Workshop (ITW)*, Visby, Sweden: IEEE, Aug. 2019, pp. 1–5

Asymmetric Multiset Deletions

Theorem

$\mathcal{S}_{DA}^{(t)}(n)$ can correct up to t -asymmetric multiset deletions.

- Proof idea [PGM21]:

[PGM21] S. Pattabiraman, R. Gabrys, and O. Milenkovic, “Coding for Polymer-Based Data Storage”, [arXiv:2003.02121](#) [cs, math], Jun. 2021, arXiv: 2003.02121

Asymmetric Multiset Deletions

Theorem

$\mathcal{S}_{DA}^{(t)}(n)$ can correct up to t -asymmetric multiset deletions.

- Proof idea [PGM21]:

1. For $\mathbf{s} \in \mathcal{S}_{DA}^{(t)}(n)$, build a set of all strings $\mathbf{v} \in \mathcal{S}_{DA}^{(t)}(n)$ such that $D_{t+1}(\mathbf{s}) \cap D_{t+1}(\mathbf{v}) \neq \emptyset$.

[PGM21] S. Pattabiraman, R. Gabrys, and O. Milenkovic, "Coding for Polymer-Based Data Storage", [arXiv:2003.02121](https://arxiv.org/abs/2003.02121) [cs, math], Jun. 2021, arXiv: 2003.02121

Asymmetric Multiset Deletions

Theorem

$\mathcal{S}_{DA}^{(t)}(n)$ can correct up to t -asymmetric multiset deletions.

- Proof idea [PGM21]:

1. For $\mathbf{s} \in \mathcal{S}_{DA}^{(t)}(n)$, build a set of all strings $\mathbf{v} \in \mathcal{S}_{DA}^{(t)}(n)$ such that $D_{t+1}(\mathbf{s}) \cap D_{t+1}(\mathbf{v}) \neq \emptyset$.
 -Includes all \mathbf{v} such that if $(\mathbf{s}_1^k, \mathbf{s}_{n-k+1}^n) = (\mathbf{v}_1^k, \mathbf{v}_{n-k+1}^n)$ and $s_{k+1} \neq v_{k+1}$, then for $i \in \{1, \dots, t+1\}$

$$|C_{n-k-i}(\mathbf{s}) \setminus C_{n-k-i}(\mathbf{v})| \leq 2.$$

[PGM21] S. Pattabiraman, R. Gabrys, and O. Milenkovic, "Coding for Polymer-Based Data Storage", [arXiv:2003.02121](https://arxiv.org/abs/2003.02121) [cs, math], Jun. 2021, arXiv: 2003.02121

Asymmetric Multiset Deletions

Theorem

$\mathcal{S}_{DA}^{(t)}(n)$ can correct up to t -asymmetric multiset deletions.

- Proof idea [PGM21]:

1. For $\mathbf{s} \in \mathcal{S}_{DA}^{(t)}(n)$, build a set of all strings $\mathbf{v} \in \mathcal{S}_{DA}^{(t)}(n)$ such that $D_{t+1}(\mathbf{s}) \cap D_{t+1}(\mathbf{v}) \neq \emptyset$.
-Includes all \mathbf{v} such that if $(\mathbf{s}_1^k, \mathbf{s}_{n-k+1}^n) = (\mathbf{v}_1^k, \mathbf{v}_{n-k+1}^n)$ and $s_{k+1} \neq v_{k+1}$, then for $i \in \{1, \dots, t+1\}$

$$|C_{n-k-i}(\mathbf{s}) \setminus C_{n-k-i}(\mathbf{v})| \leq 2.$$

2. We infer that no \mathbf{v} simultaneously satisfies

$$C_{n-k-t-1}(\mathbf{s}) = C_{n-k-t-1}(\mathbf{v}).$$

[PGM21] S. Pattabiraman, R. Gabrys, and O. Milenkovic, "Coding for Polymer-Based Data Storage", [arXiv:2003.02121](https://arxiv.org/abs/2003.02121) [cs, math], Jun. 2021, arXiv: 2003.02121

Asymmetric Multiset Deletions

- Thus, if $C_{n-k-1}(\mathbf{s}), \dots, C_{n-k-t}(\mathbf{s})$ are deleted, each $\mathbf{v} \in S_{DA}^{(t)}(n)$ upholds

$$C_{n-k-t-1}(\mathbf{s}) \neq C_{n-k-t-1}(\mathbf{v})$$

Asymmetric Multiset Deletions

- Thus, if $C_{n-k-1}(\mathbf{s}), \dots, C_{n-k-t}(\mathbf{s})$ are deleted, each $\mathbf{v} \in S_{DA}^{(t)}(n)$ upholds

$$C_{n-k-t-1}(\mathbf{s}) \neq C_{n-k-t-1}(\mathbf{v})$$

$$\implies D_t(\mathbf{s}) \cap D_t(\mathbf{v}) = \emptyset.$$

Asymmetric Multiset Deletions

- Thus, if $C_{n-k-1}(\mathbf{s}), \dots, C_{n-k-t}(\mathbf{s})$ are deleted, each $\mathbf{v} \in \mathcal{S}_{DA}^{(t)}(n)$ upholds

$$C_{n-k-t-1}(\mathbf{s}) \neq C_{n-k-t-1}(\mathbf{v})$$

$$\implies D_t(\mathbf{s}) \cap D_t(\mathbf{v}) = \emptyset.$$

→ Thus, \mathbf{s} can be recovered uniquely.

Asymmetric Multiset Deletions

- Thus, if $C_{n-k-1}(\mathbf{s}), \dots, C_{n-k-t}(\mathbf{s})$ are deleted, each $\mathbf{v} \in S_{DA}^{(t)}(n)$ upholds

$$C_{n-k-t-1}(\mathbf{s}) \neq C_{n-k-t-1}(\mathbf{v})$$

$$\implies D_t(\mathbf{s}) \cap D_t(\mathbf{v}) = \emptyset.$$

→ Thus, \mathbf{s} can be recovered uniquely.

- When deleted multisets are nonconsecutive, proof follows similarly.

Introduction

Problem Statement & Results

Conclusion

Summary of Results

- A code that can correct t multiset deletions, iff correct composition insertions in t multisets.

Summary of Results

- A code that can correct t multiset deletions, iff correct composition insertions in t multisets.
- $\mathcal{S}_R(n)$ can correct the deletion of a single multiset.

Summary of Results

- A code that can correct t multiset deletions, iff correct composition insertions in t multisets.
- $\mathcal{S}_R(n)$ can correct the deletion of a single multiset.
- To correct the deletion of t asymmetric multisets, $0.5 \log(n - 2t) + 2t + 3$ redundant bits suffice.

Summary of Results

- A code that can correct t multiset deletions, iff correct composition insertions in t multisets.
- $\mathcal{S}_R(n)$ can correct the deletion of a single multiset.
- To correct the deletion of t asymmetric multisets, $0.5 \log(n - 2t) + 2t + 3$ redundant bits suffice.
- For any $1 \leq i \leq n$, $\mathcal{S}_R(n)$ can correct the deletion of multisets $C_i(\mathbf{s})$ and $C_{n-i+1}(\mathbf{s})$.

Future Work

- **Upper bounds** on codes for unique reconstruction from composition multisets.

Future Work

- **Upper bounds** on codes for unique reconstruction from composition multisets.
- Extension to **non-binary** alphabets.

Future Work

- **Upper bounds** on codes for unique reconstruction from composition multisets.
- Extension to **non-binary** alphabets.
- Extension to **circular polymers**, i.e. bits on a ring.

Thank you!

References



A. Al Ouahabi, J.-A. Amalian, L. Charles, and J.-F. Lutz, “Mass spectrometry sequencing of long digital polymers facilitated by programmed inter-byte fragmentation”, *Nature Communications*, vol. 8, no. 1, p. 967, Dec. 2017.



J. Acharya, H. Das, O. Milenkovic, A. Orlitsky, and S. Pan, “String Reconstruction from Substring Compositions”, *SIAM Journal on Discrete Mathematics*, vol. 29, no. 3, pp. 1340–1371, Jan. 2015.



S. Pattabiraman, R. Gabrys, and O. Milenkovic, “Reconstruction and Error-Correction Codes for Polymer-Based Data Storage”, in *2019 IEEE Information Theory Workshop (ITW)*, Visby, Sweden: IEEE, Aug. 2019, pp. 1–5.



R. Gabrys, S. Pattabiraman, and O. Milenkovic, “Mass Error-Correction Codes for Polymer-Based Data Storage”, in *2020 IEEE International Symposium on Information Theory (ISIT)*, Los Angeles, CA, USA: IEEE, Jun. 2020, pp. 25–30.



—, *Reconstructing Mixtures of Coded Strings from Prefix and Suffix Compositions*, Number: arXiv:2010.11116 arXiv:2010.11116 [cs, math], Oct. 2020.



U. Gupta and H. MahdaviFar, “A New Algebraic Approach for String Reconstruction from Substring Compositions”, *arXiv:2201.09955 [cs, math]*, Jan. 2022, arXiv: 2201.09955.



S. Pattabiraman, R. Gabrys, and O. Milenkovic, “Coding for Polymer-Based Data Storage”, *arXiv:2003.02121 [cs, math]*, Jun. 2021, arXiv: 2003.02121.