IMPERIAL

DNA data storage: A generative tool for Motif-based DNA storage

Samira Brunmayr, Omer S. Sella, Thomas Heinis

23rd USENIX Conference on File and Storage Technologies February 27, 2025

Data Storage

- High demand
- Storage devices deteriorate over time
- Require considerable space and energy



DNA as a Storage Medium

- 10¹⁸ bytes per mm³
- Store data for 2000 years
- Store Wikipedia in test tube



How to store data in DNA?

GTAGGTACAG

Bases	Binary		
Α	00		
Т	01		
С	10		
G	11		

Encoding example

How to store data in DNA?

- 1. Encoding
- 2. Synthesis: construct DNA strands
- 3. Storage of DNA
- 4. Sequencing: access data in DNA
- 5. Decoding



Current Technological Limitations

- Strand length: <150-200 bases
- Synthesis: \$400 million/Terabyte



Solution: Motif-based DNA Data Storage

• Motif: Short DNA sequence



















Motif Structure

- Payload: Actual Information
- Key:
 - Used for Assembly
 - Access Data



GTAGGTACAG

• Insertion



- Insertion
- Deletion



- Insertion
- Deletion
- Substitution



- Homopolymer
- Hairpins
- GC-Content

- Homopolymer
- Hairpins
- GC-Content



- Homopolymer
- Hairpins
- GC-Content





- Homopolymer
- Hairpins
- GC-Content





What we want:

- Cost-efficient: Motif-based storage
- **Reliable:** conform to biological constraints





The Motif Generation Tool

Motif Generation Tool

- Input:
 - Information on Motifs: payload length, ...
 - Constraints on:
 - Homopolymer
 - Hairpins
 - GC-Content
- Output:
 - Set of motifs conforming to the constraints





Transition Probabilities

 $score(i) = score_{hom}(i) \times score_{hairpin}(i) \times score_{gc}(i) \times score_{noKeyInPayload}(i)$

$$score(i) = e^{logScore(i)}$$

$$p_i = \frac{score(i)}{\sum_{j \in \{A,T,C,G\}} score(j)}$$

Transition Probabilities

 $score(i) = score_{hom}(i) \times score_{hairpin}(i) \times score_{gc}(i) \times score_{noKeyInPayload}(i)$

$$score(i) = e^{logScore(i)}$$

$$p_i = \frac{score(i)}{\sum_{j \in \{A,T,C,G\}} score(j)}$$





$$ls_{hom}(i) = -(h_{hom})^{homLen/maxHom} + 1$$



 $ls_{hom}(i) = -(h_{hom})^{homLen/maxHom} + 1$



 $ls_{hom}(i) = -(h_{hom})^{homLen/maxHom} + 1$

• maxHom : Constraint



 $ls_{hom}(i) = -(\mathbf{h}_{hom})^{homLen/maxHom} + 1$ Homopolymer Log Score • h_{hom} : controls shape of log score -10 Log Score) h_{hom} = 5 h_{hom} = 15
h_{hom} = 30 -20 --30 2 3 0 4 5 36

Evaluation

Evaluation Setup

Other tools for comparison:

- DNA Fountain, Erlich et al.
- Finite State Machine Encoding, Sella et al.
- Shortmer combinatorial encoding scheme, Preuss et al.
- Randomly generated DNA sequences

Evaluation:

- For each constraint separately
- Case Study

Single Constraint



Single Constraint



Case Study

• Currently applicable constraints

Parameters		Values	
Motifs	keySize	20	
	keyNum	8	
	payloadSize	60	
	payloadNum	15	
Homopolymers	maxHom	5	
Hairpins	maxHairpin	1	
	minLoopSize	6	
	maxLoopSize	7	
GC-Content	minGC	25	
	maxGC	65	

Case Study

	Motif Generation Tool	DNA Fountain [11]	Euclid [21]	Preuss et al. [17]	Randomly generated
Time taken to generate a set of motifs conforming to the constraints in Table 1 (seconds)	2.54 sec	>5min	>5min	>5min	>5min

Summary

- Motif-based DNA data storage cost efficient
- Motif Generation Tool for automated reliable motif generation
- Outperforms other tools based on currently industry standard biological constraints

Summary

- Motif-based DNA data storage cost efficient
- Motif Generation Tool for automated reliable motif generation
- Outperforms other tools based on currently industry standard biological constraints

Thank you for your attention!

Q&A

Sponsors: DNAMIC and NEO