

DNA Storage Toolkit:

A Modular End-to-End DNA Data Storage Codec and Simulator



Puru
Sharma



Gary Goh
Yipeng



Bin
Gao



Longshen
Ou



Dehui
Lin



Deepak
Sharma



Djordje
Jevdjic



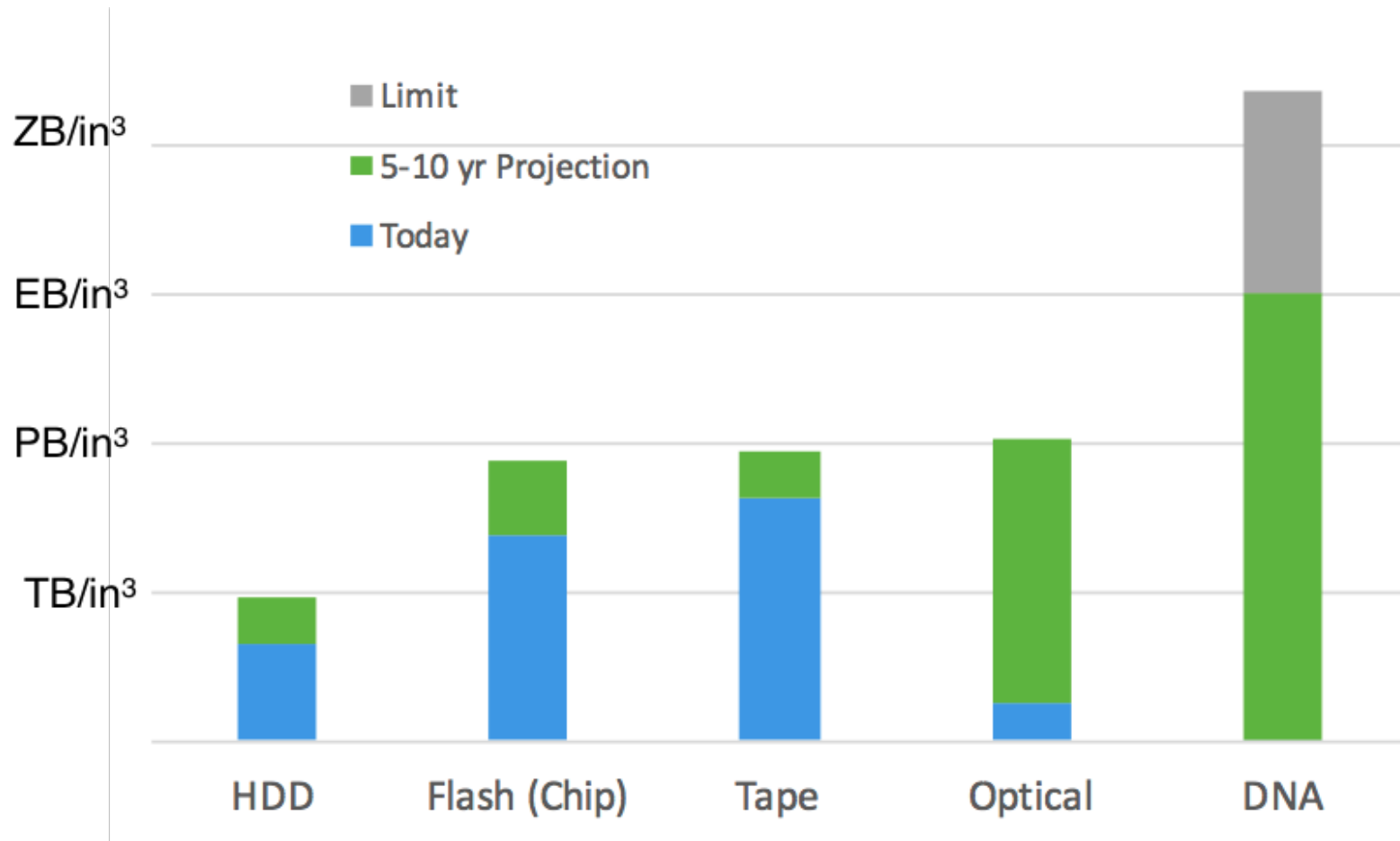
NUS
National University
of Singapore

Why storing data in DNA molecules?

1. Incredible density
 - 6-7 orders of magnitude ahead of best alternatives!
2. Unmatched durability
 - Thousands/millions/billions of years (vs. 3-5 years for disks/flash)
3. Never obsolete: R/W interfaces will only improve with time
4. Efficient random access
5. Convenient for many data-parallel & near-data computations



Storage Density Projections*



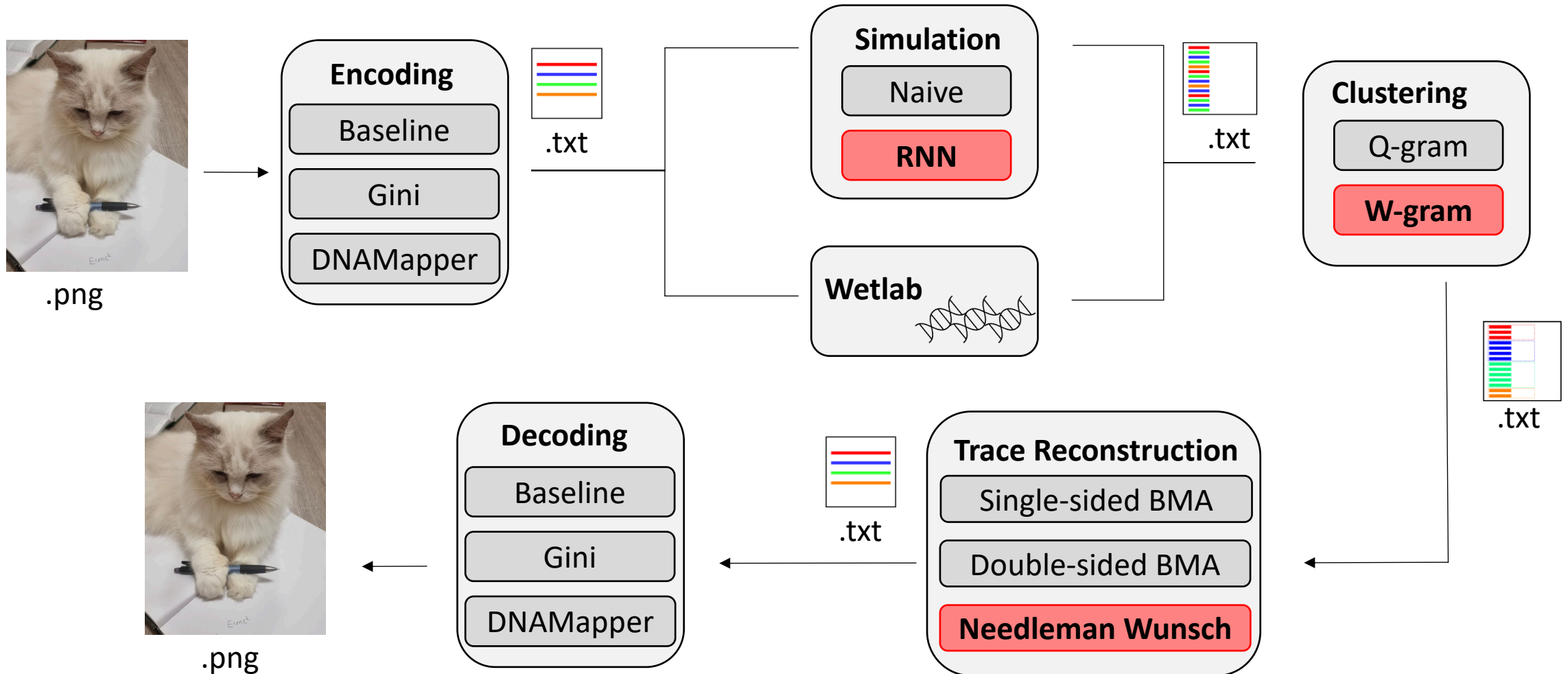
*Credit: Luis Ceze & Karin Strauss, Microsoft

Key Problems with DNA Storage

1. Expensive R/W interfaces
 - Writing cost: \$1K - \$10K/MiB, reading cost: \$10 - \$10K/MiB
→ High cost also makes research very challenging
2. High read/write latency
 - Takes days to write DNA, hours (to days) to read (OK for archival storage)
→ But this delay slows down research
3. Extremely error-prone interfaces
 - Errors are very peculiar and hard to simulate
4. Requires expensive equipment and wetlab expertise
5. No complete open-source codec available



Our DNA Storage Pipeline

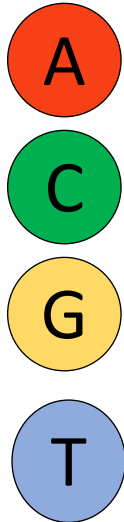


Outline

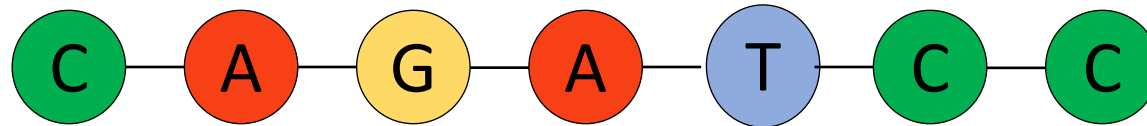
- Introduction
- **DNA Storage Basics**
- Our Toolkit
 - Encoding
 - Simulation
 - Clustering
 - Trace Reconstruction
 - Decoding
- Conclusion

DNA Molecules

4 nucleotides



Synthetic DNA molecule

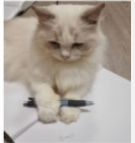


- Artificially created string of nucleotides
- No biological meaning

$$\log_2 |\{A, C, G, T\}| = 2 \text{ bits of data per nucleotide}$$

DNA Storage Pipeline

Encoding

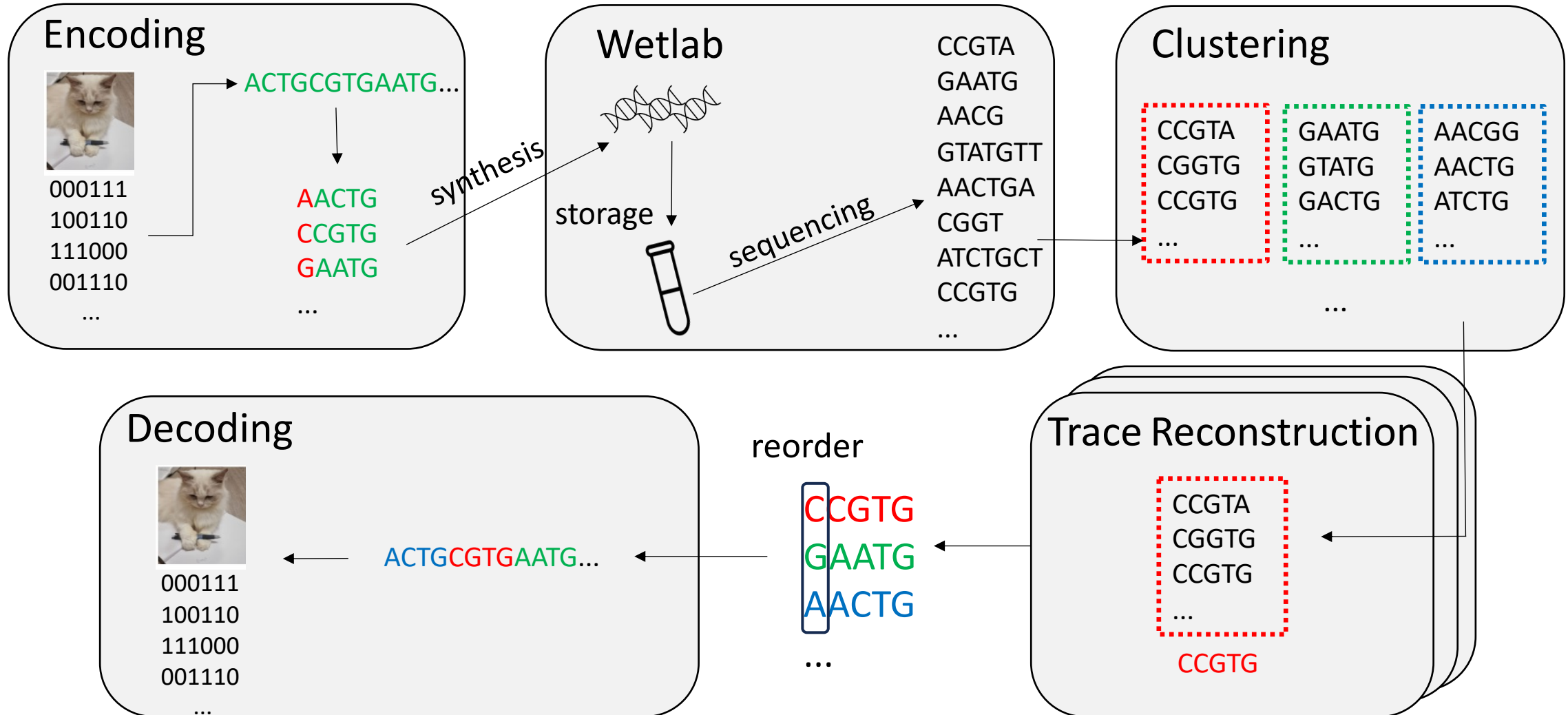


000111
100110
111000
001110
...

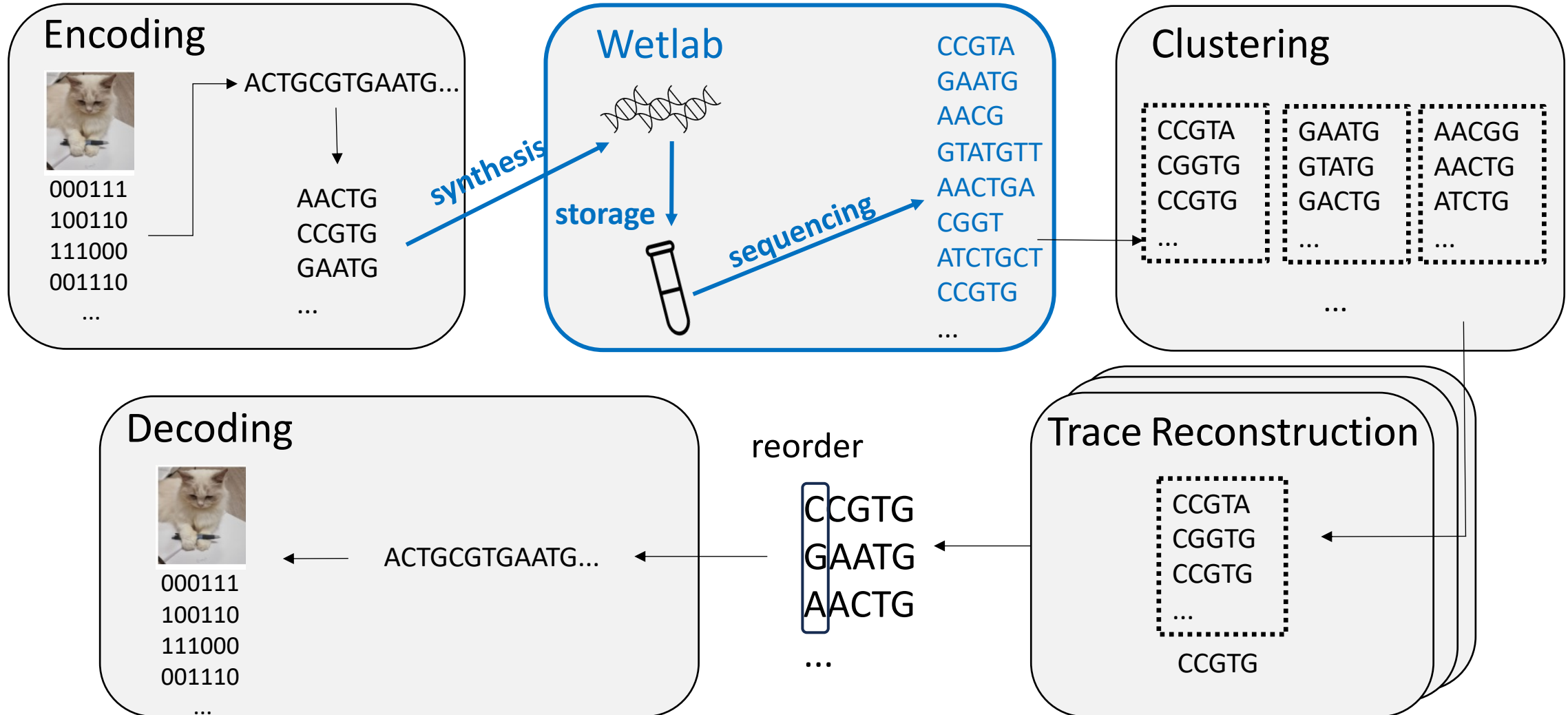
→ ACTGCGTGAATG...

↓
A**ACTG**
C**CGTG**
G**AATG**
...

DNA Storage Pipeline



DNA Storage Pipeline – wetlab steps



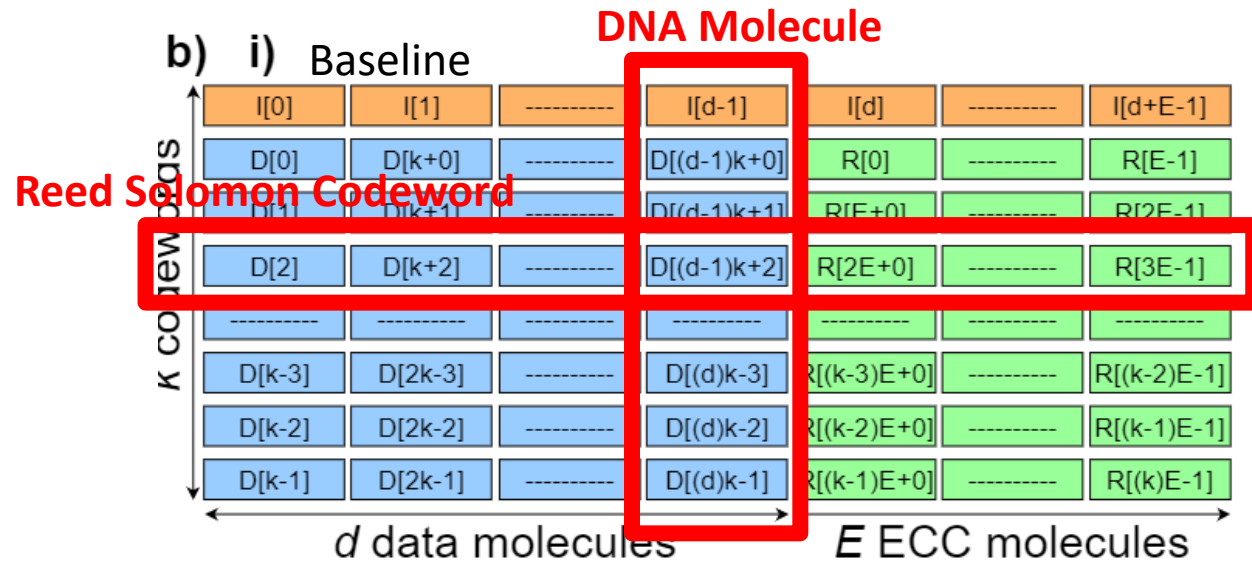
Outline

- Introduction
- DNA Storage Basics
- **Our Toolkit**
 - Encoding
 - Simulation
 - Clustering
 - Trace Reconstruction
 - Decoding
- Conclusion

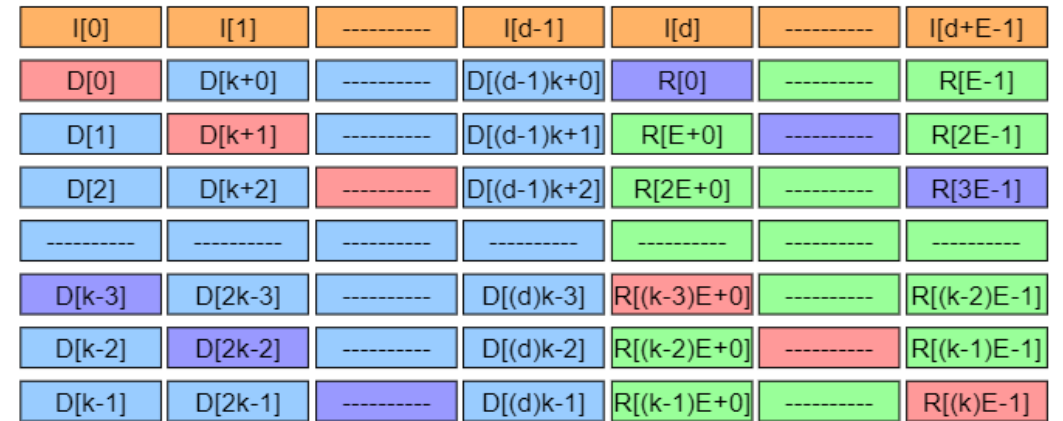
Encoding

3 encoding schemes provided.

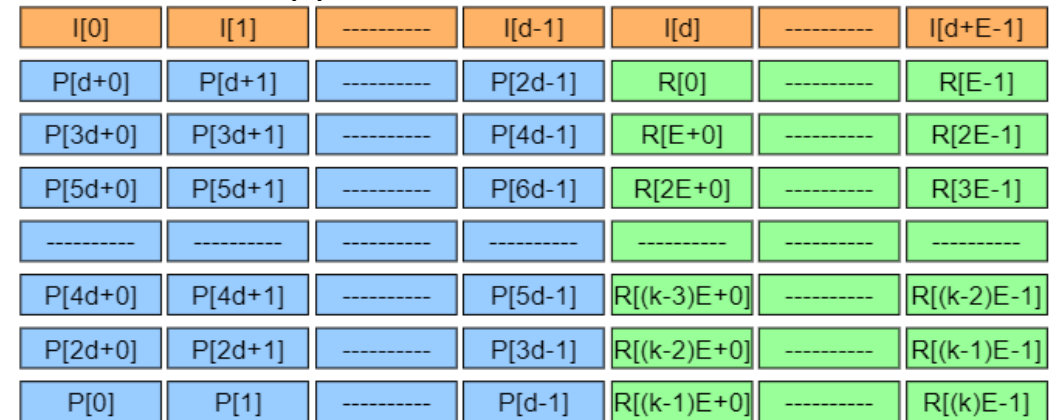
- i) Baseline: Introduced by Organick et al. [1]
- ii) Gini: Introduced by Lin et al. [2]
- iii) DNAMapper: Introduced by Lin et al. [2]



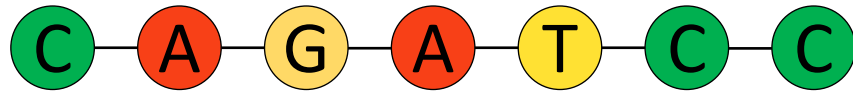
ii) Gini



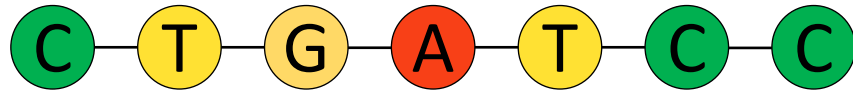
iii) DNA Mapper



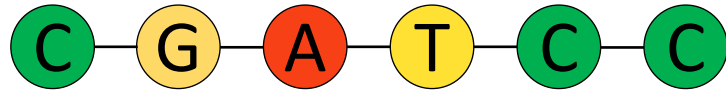
Wetlab



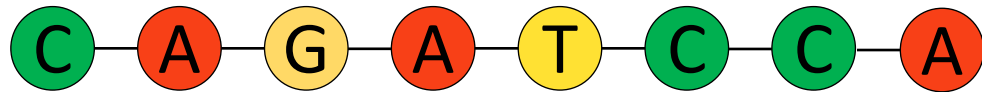
Original



Substitution (A→T)



Deletion (of A)



Insertion (of A)

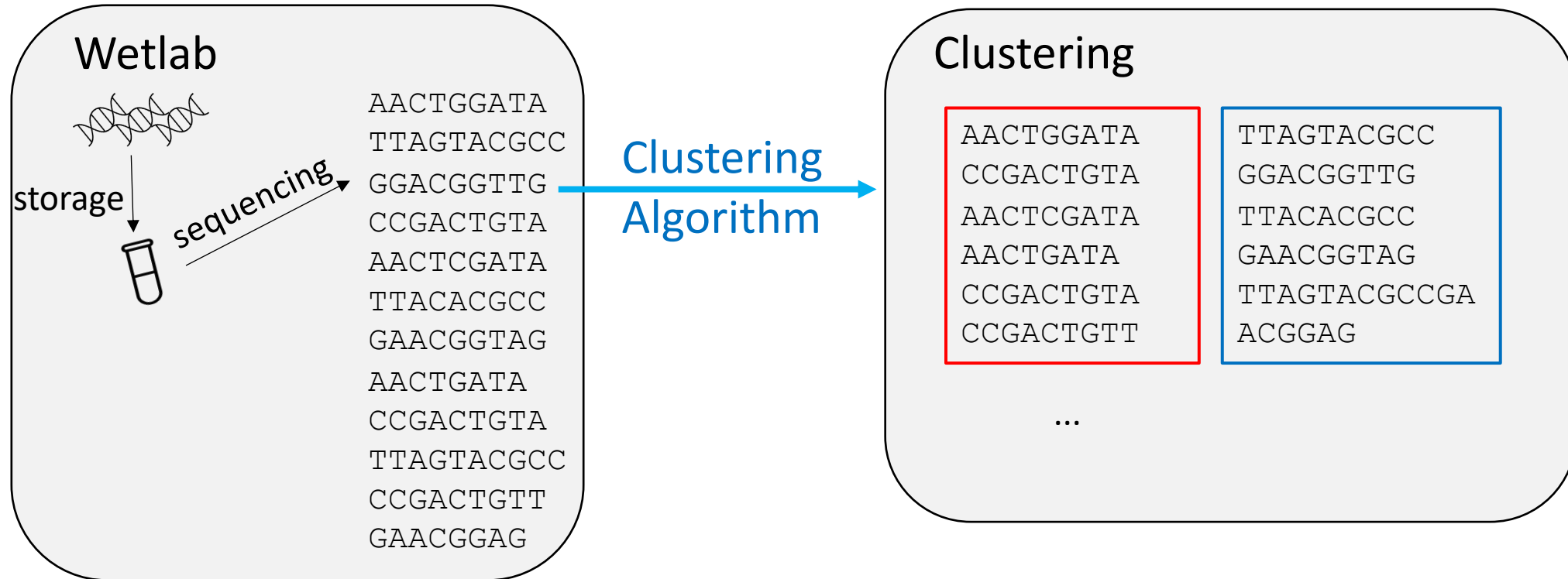
Sequencing produces many (buggy) copies of each molecule:

synthesized
CAGATCC

Simulator? 

sequenced
CAGATCC
CAGATC
AAGATCCA
AGATTCC

Clustering



Clustering

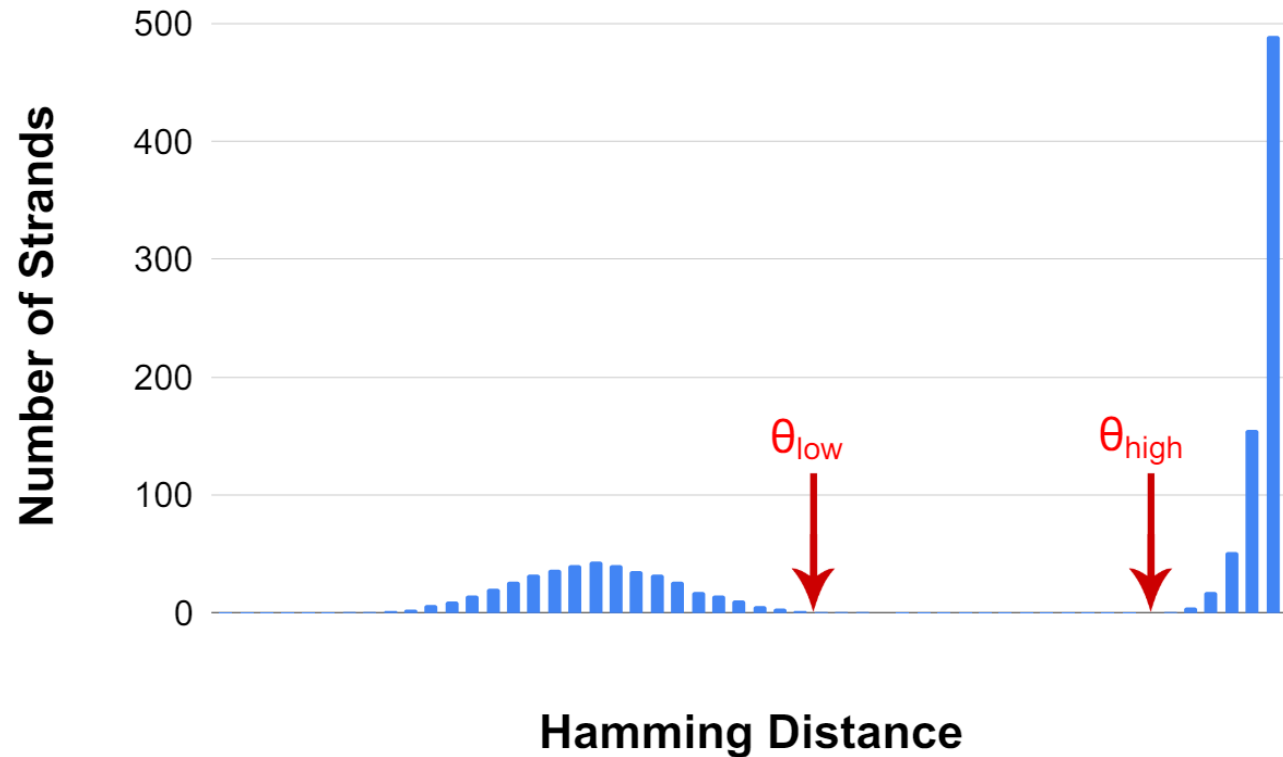
- Using edit distance for clustering is too slow.
- We can approximate using q-gram binary signatures and Hamming distance.

	CT	GA	GT	
ACGT GA AC	0	1	1	} Low Hamming distance
ACGTCC GA AC	0	1	1	
ACGT GA AC	0	1	1	} High Hamming distance
TAC CT ATTCC	1	0	0	

Autotuning the threshold parameters

The threshold parameters for the binary signature Hamming distance need to be tuned based on dataset.

Take a tiny sample and plot the Hamming distance for them:



Trace Reconstruction

GTACCAGTCGAGTAAAGC
GCCGTGCGTAAGCT
GTACAATGTCGTGTAAC
GTCATGGTCAGTAAGC
GTACAGTCCGTAAAGC
TACGTGTATAGC
GATACAGCACGTGAAGC

What is the consensus?

Trace Reconstruction

GTACCAGTCGAGTAAAGC

GCCGTGCGTAAGCT

GTACAATGTCGTGTAAC

GTCATGGTCAGTAAGC

GTACAGTCCGTAAAGC

TACGTGTATAGC

GAATACAGCACGTGAAGC

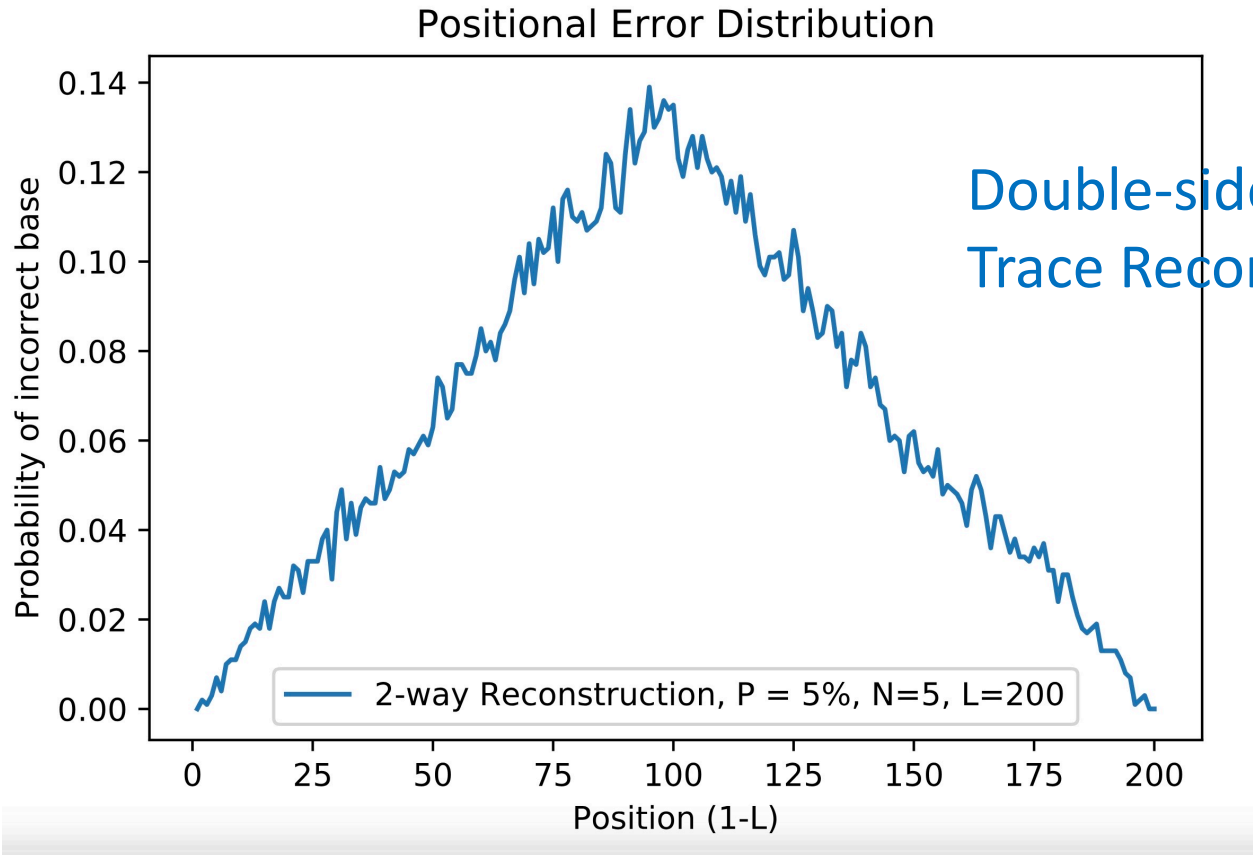
G

Trace Reconstruction

G**T**ACCAGTCGAGTAAAGC
G**C**CGTGCGTAAGCT
G**T**ACAATGTCGTGTAAC
G**T**CATGGTCAGTAAGC
G**T**ACAGTCCGTAAAGC
-**T**ACGTGTATAGC
G**A**TACAGCACGTGAAGC

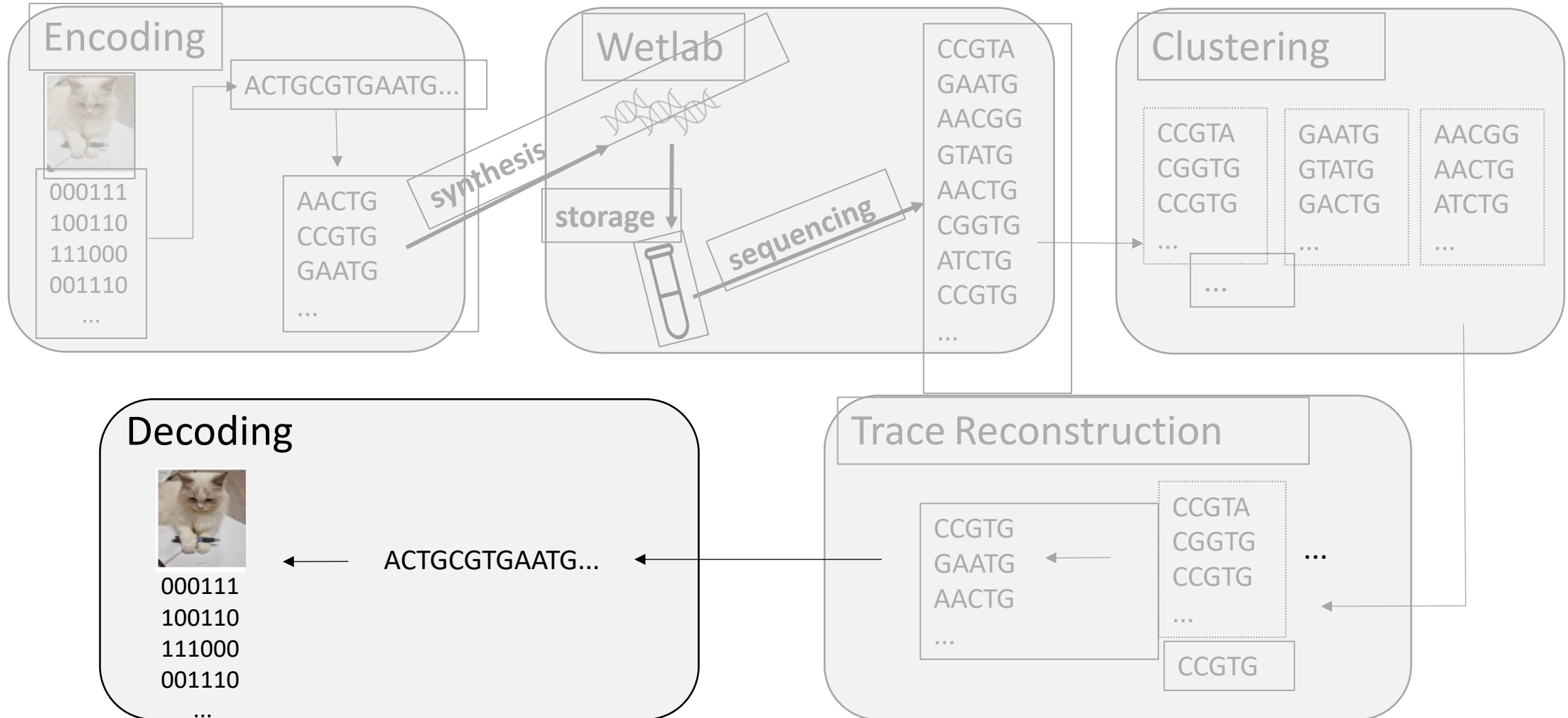
GT

Trace Reconstruction



Double-sided Bitwise Majority Trace Reconstruction

Decoding



Wetlab Simulation

Sequencing produces many (buggy) copies of each molecule:

synthesized
CAGATCC

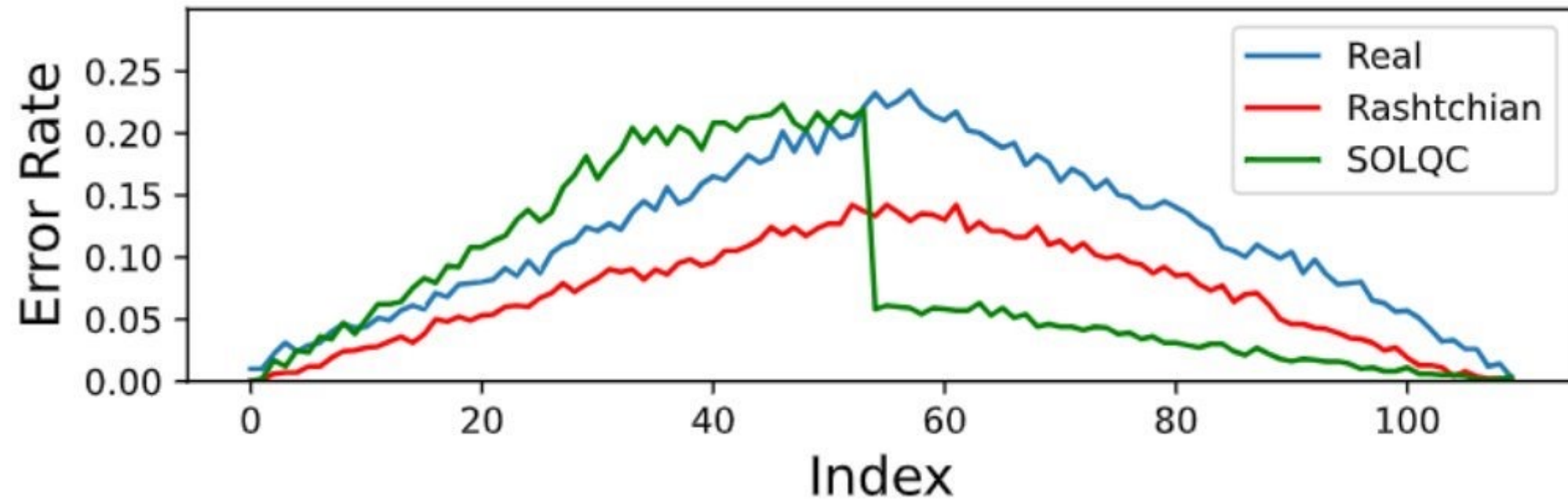
Simulator?



sequenced
CAGATCC
CAGATC
AAGATCCA
AGATTCC

Wetlab Simulation

We evaluate the simulation by comparing how the Trace Reconstruction module performs on the simulated data.



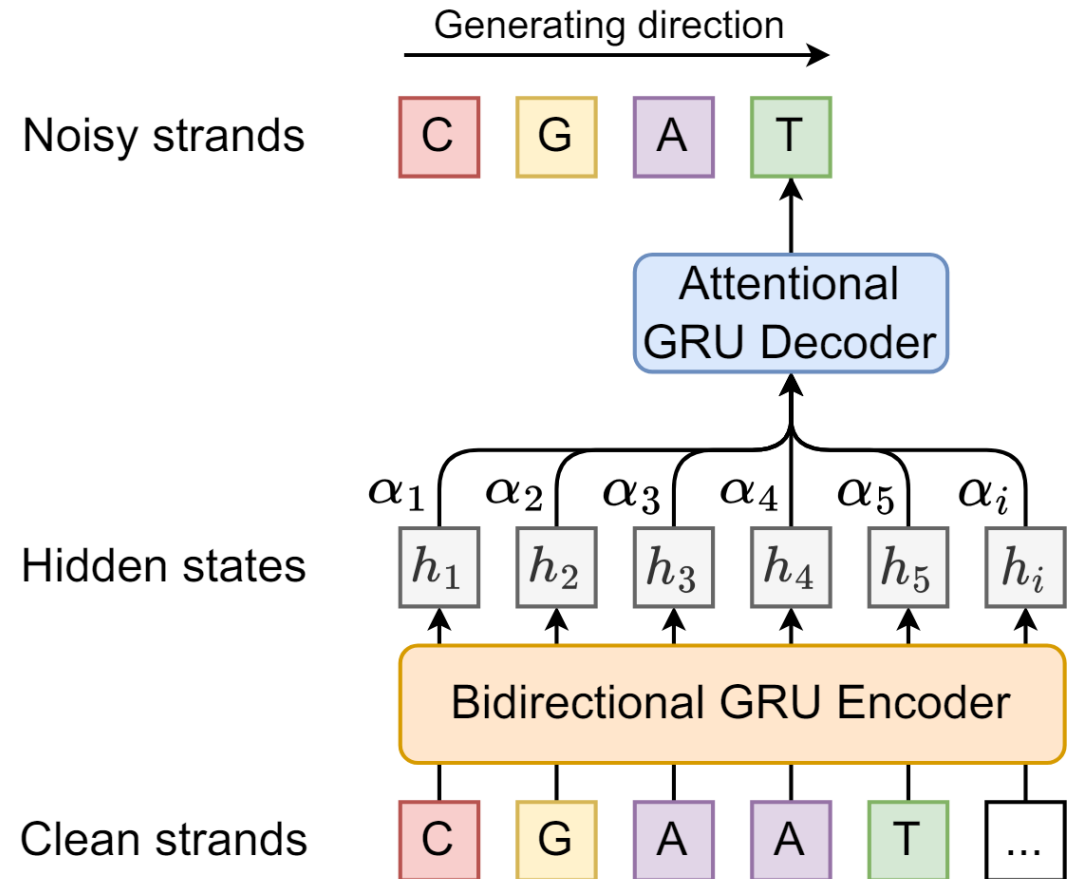
Our Simulation

Sequence-to-sequence problem.

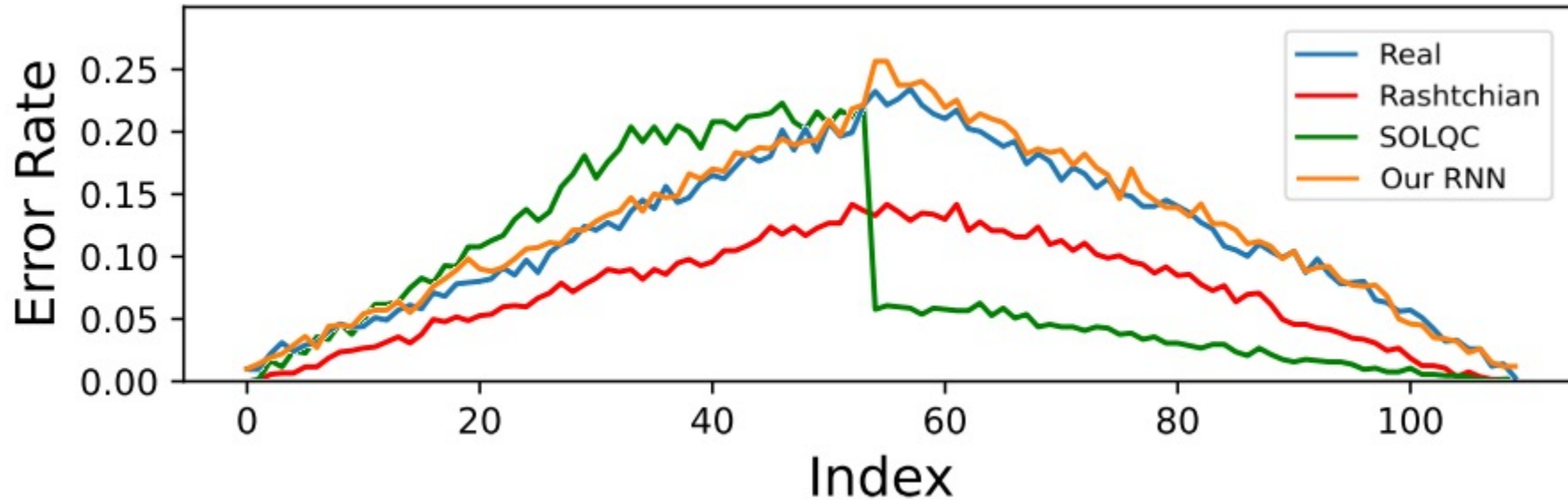
Model based on basic NMT model [3].

Attention based encoder-decoder.

[Hidden layer size = 128, Greedy sampling.]



Our Simulation

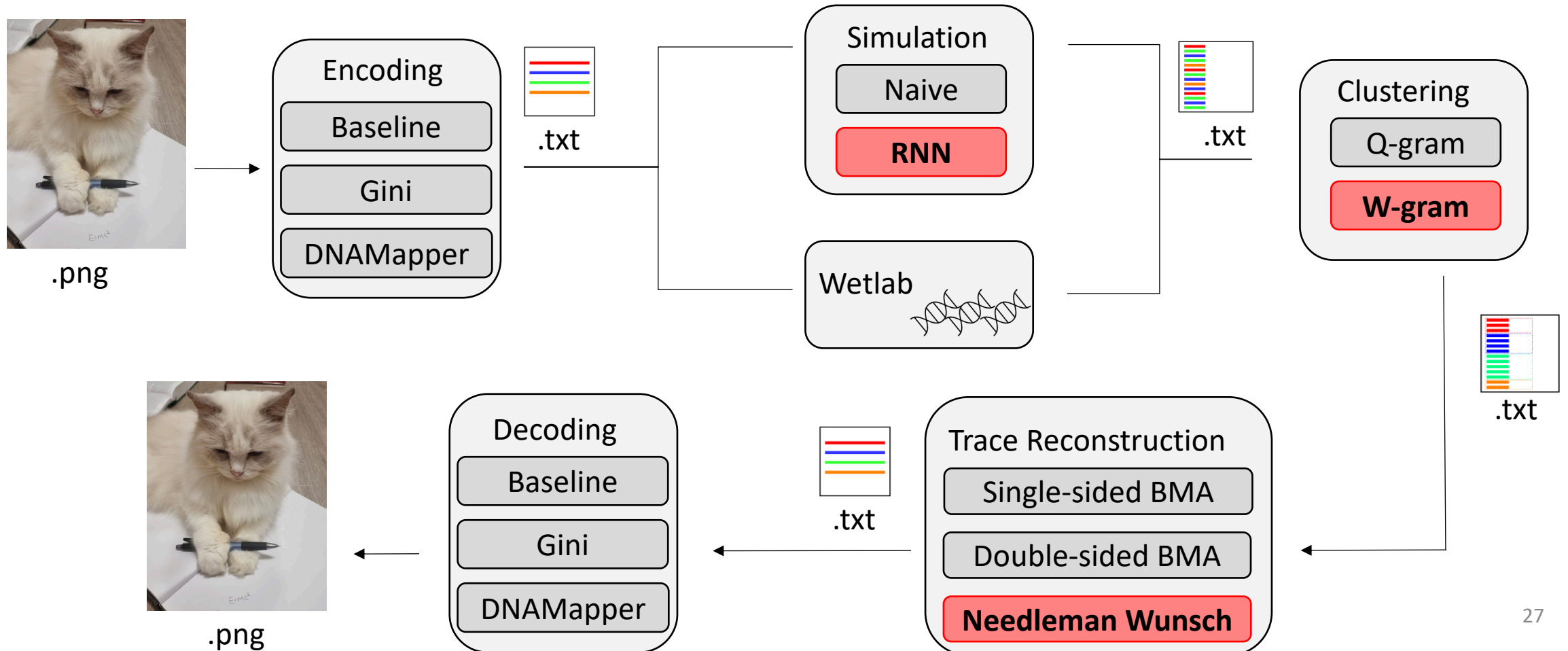


The performance of the reconstruction module on **our simulated data** is very similar to **real data**.

We can evaluate later modules of this pipeline using our simulated data!

Toolkit

<https://github.com/prongs1996/DNAStorageToolkit/>



Conclusion

- First open end-to-end DNA storage toolkit
- Very accurate simulator for wetlab steps



Repository: <https://github.com/prongs1996/DNAStorageToolkit/>

Thank you! Questions?



Puru
Sharma



Gary Goh
Yipeng



Bin
Gao



Longshen
Ou



Dehui
Lin



Deepak
Sharma



Djordje
Jevdjic



NUS
National University
of Singapore