



# Efficient storage of and in DNA: genomic data compression & DNA based storage

Shubham Chandak

PhD '21, Electrical Engineering, Stanford University

Currently Applied Scientist, S3, Amazon Web Services

-

*Biochemical Engineering and Biotechnology Department Seminar*

*IIT Delhi*

*Apr 28, 2022*

# Outline

- Introduction to genomic sequencing technologies
- Genomic data compression: SPRING
- Using DNA as a storage medium

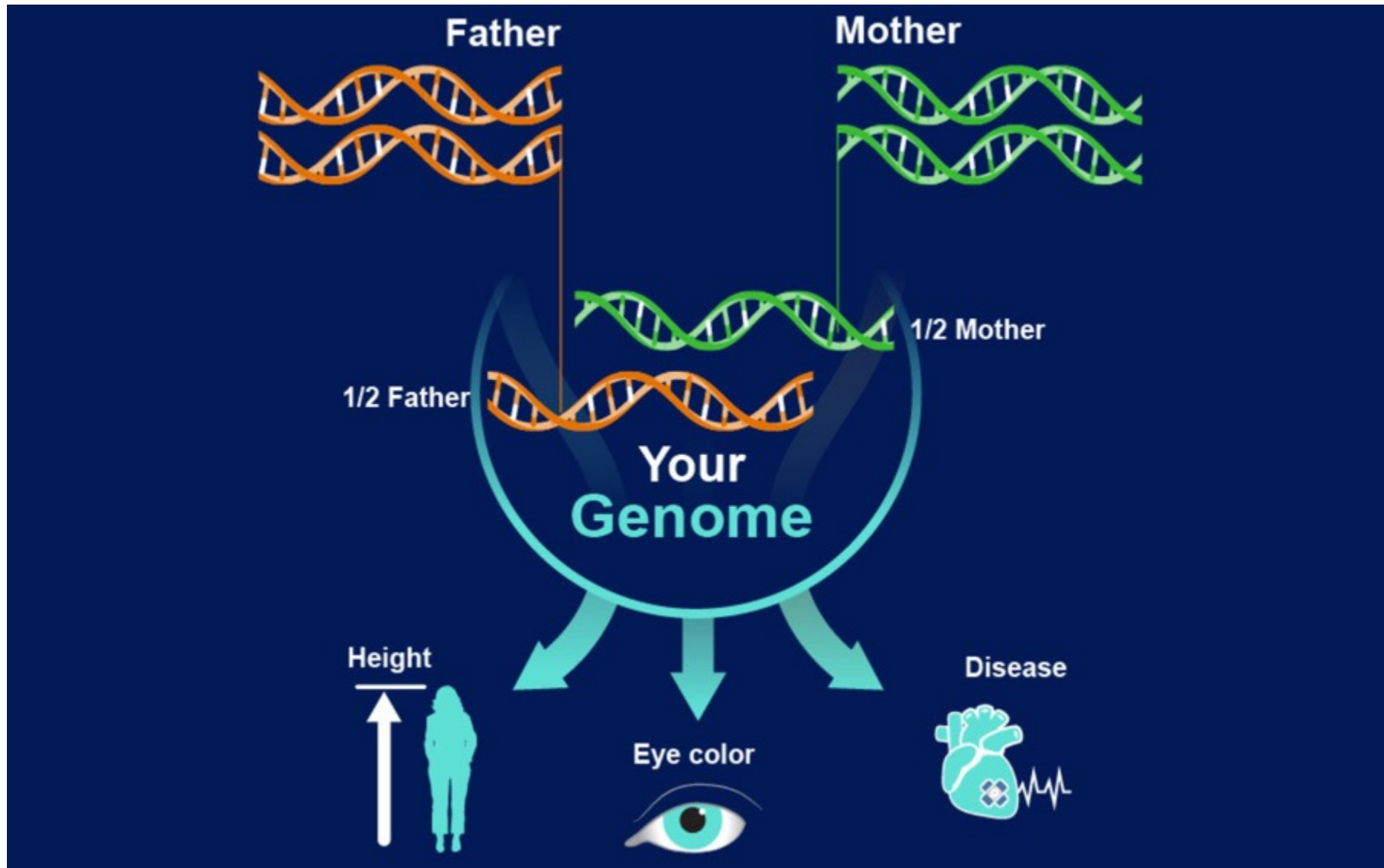
# Introduction to genomic sequencing technologies

What is the genome?

What is genome sequencing?

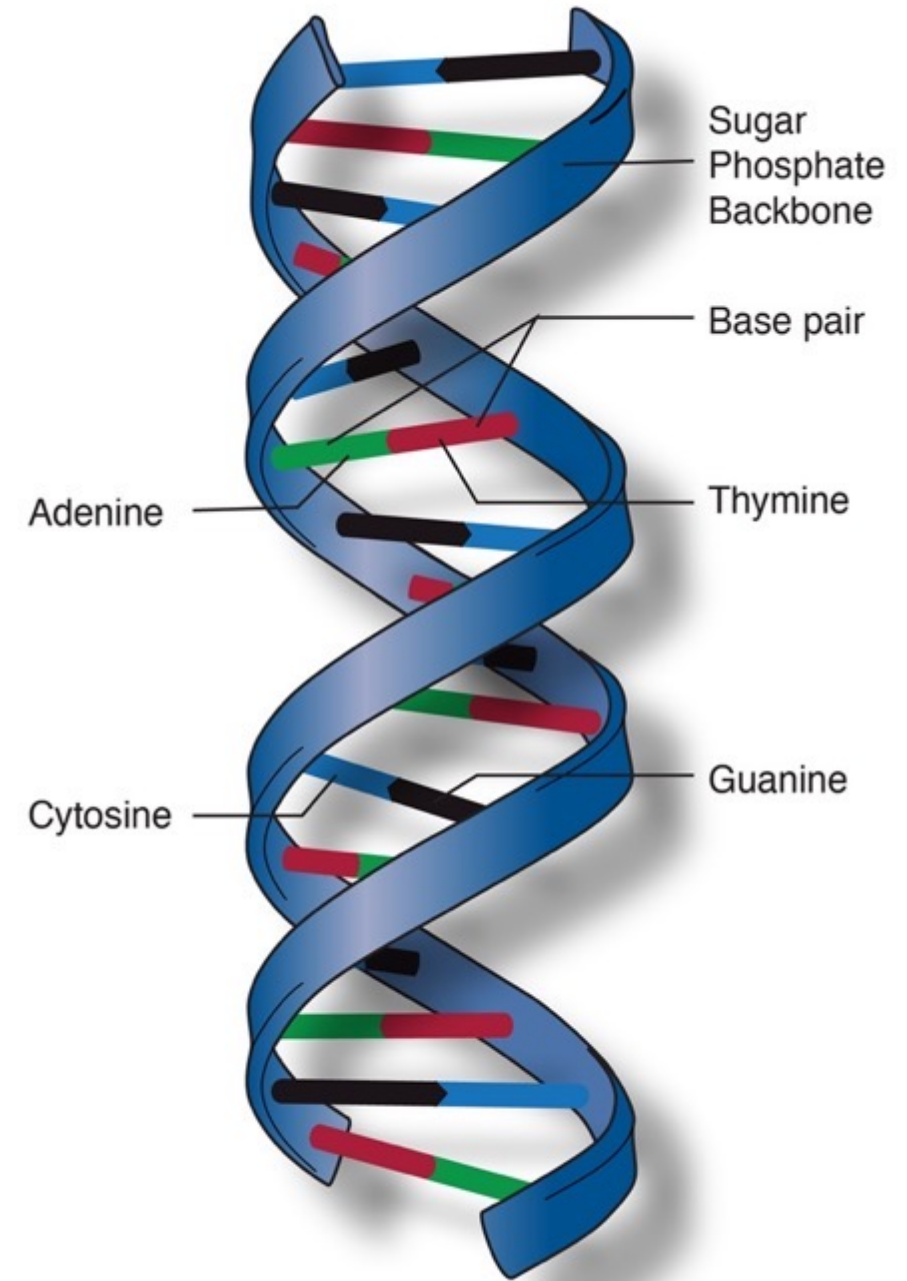
Why compression?

Raw data and downstream analysis

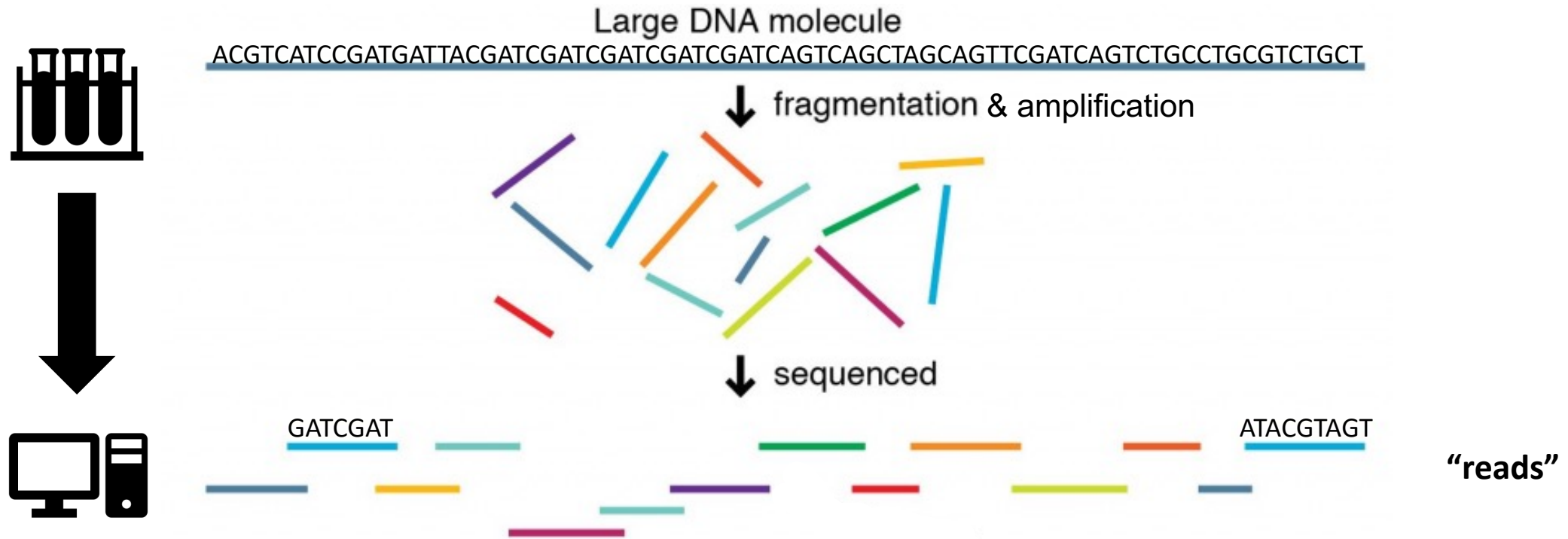


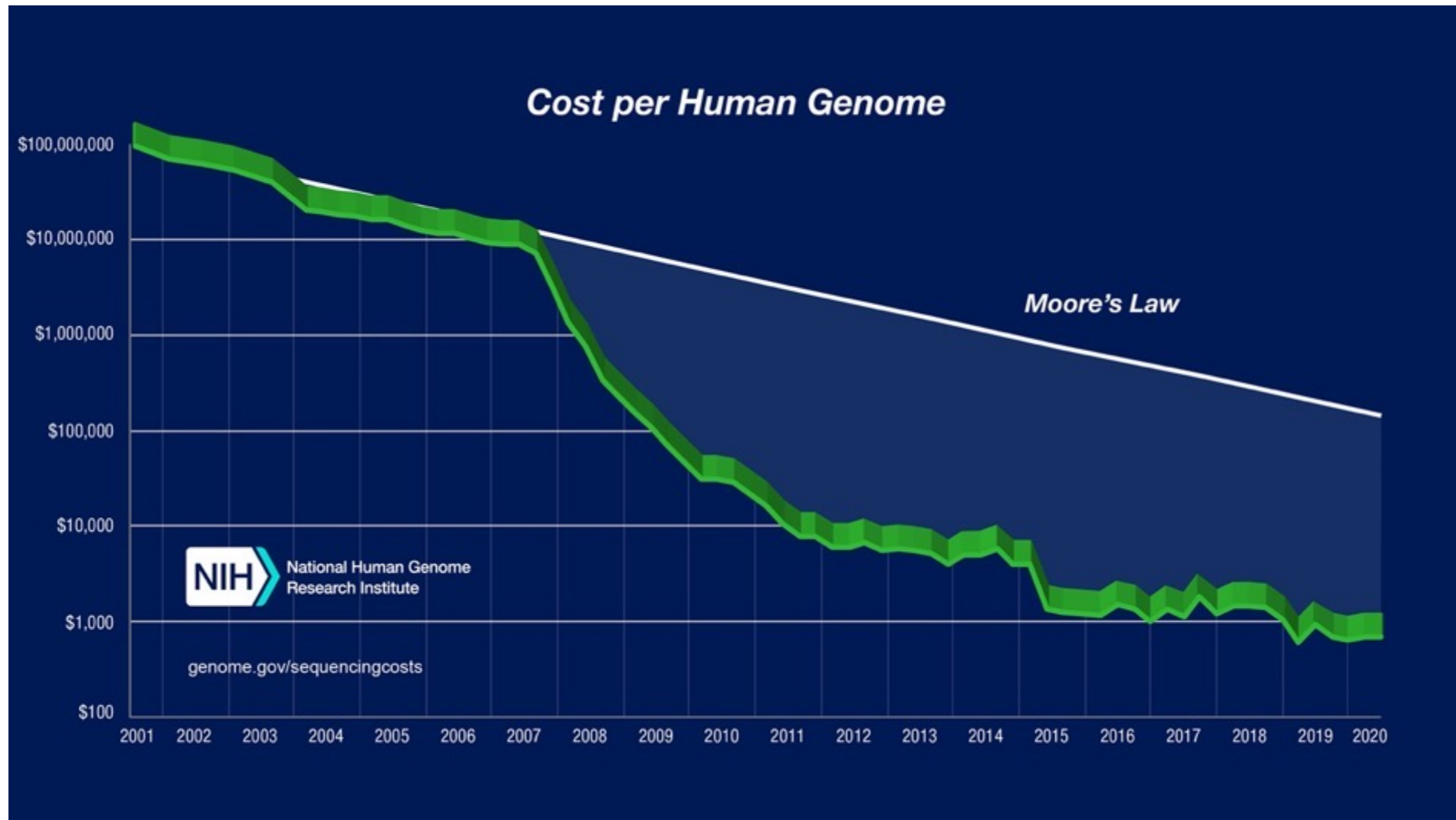
# What is the genome?

- Sequence of DNA bases in {A, C, G, T}
- Two complementary strands
- For humans:
  - 3 billion bases (x2)
  - Across 23 (x2) chromosomes



# Genome sequencing

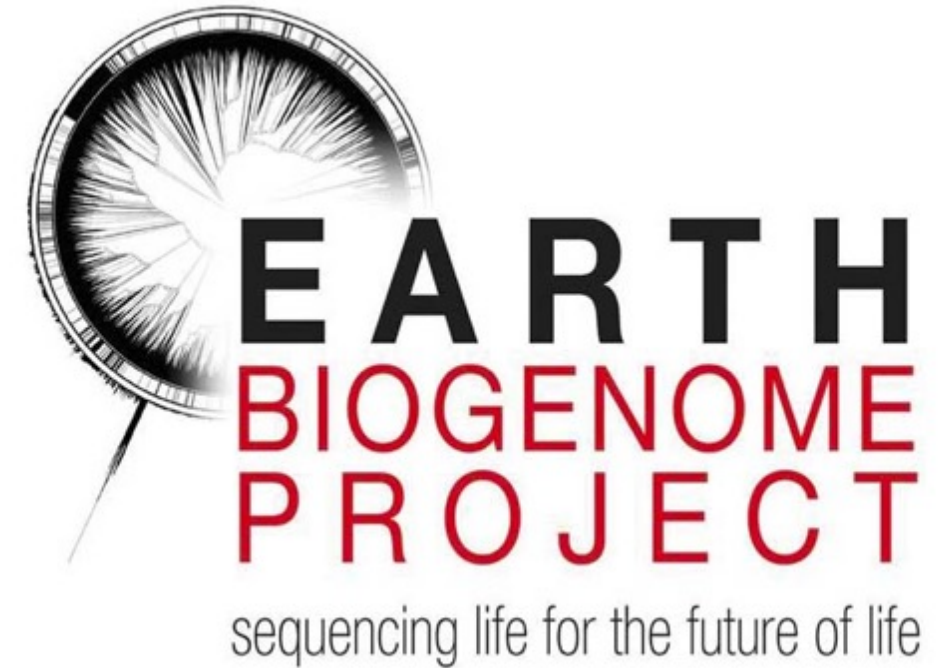








500K human genomes



~1.5M eukaryote species

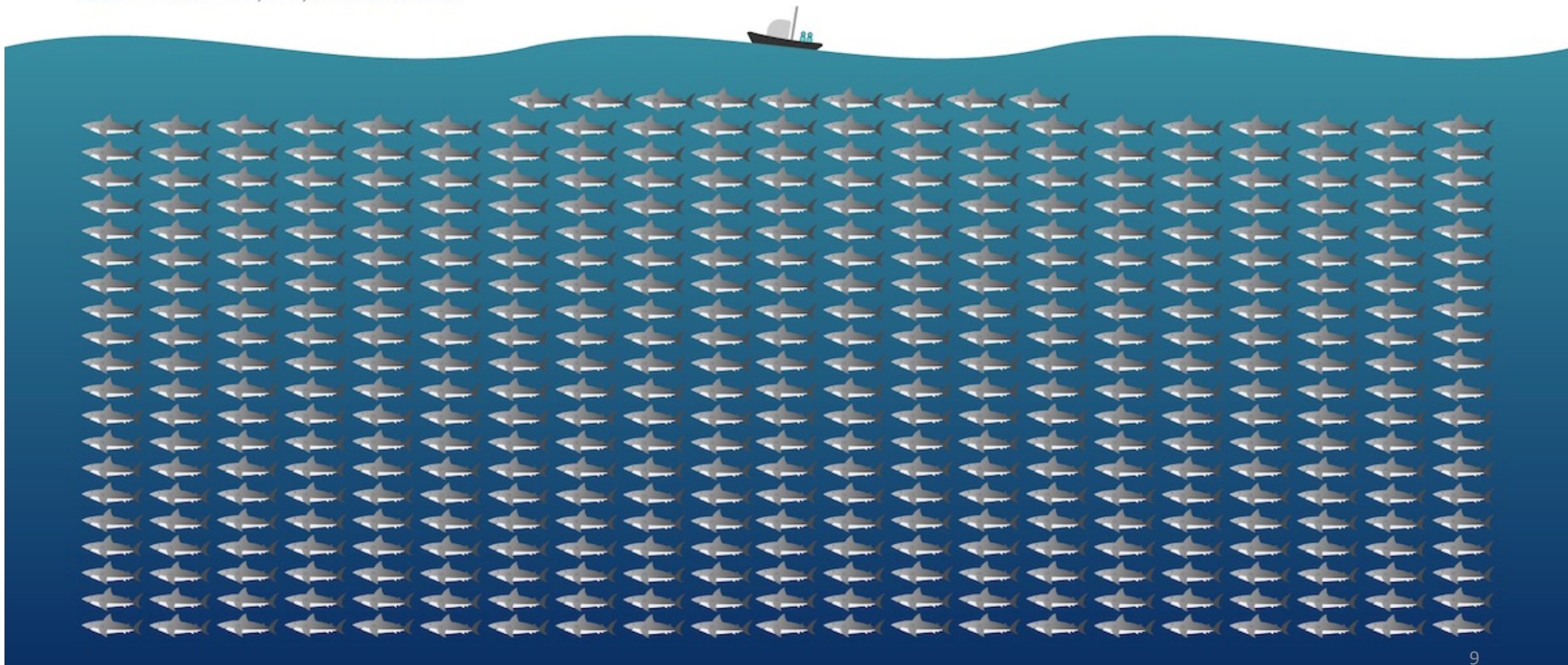




# How big is 40 exabytes?

Genomics projects will generate 40 exabytes of data in the next decade.

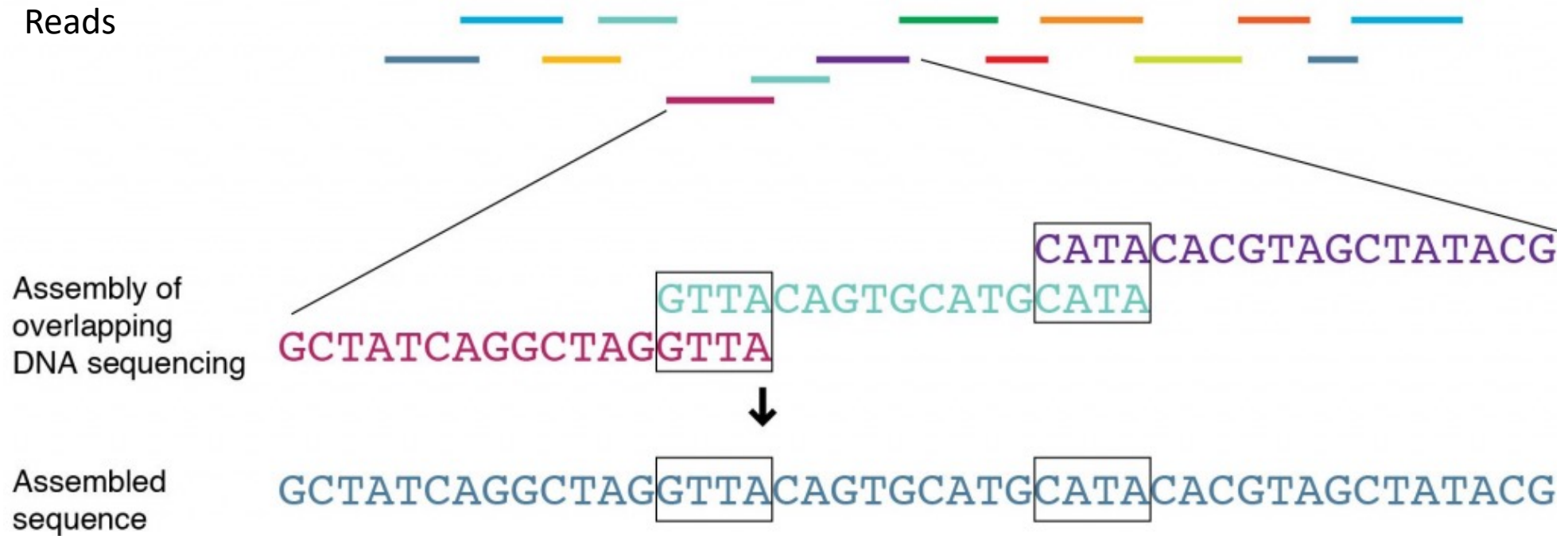
*Each shark = 100,000,000 GB of data*



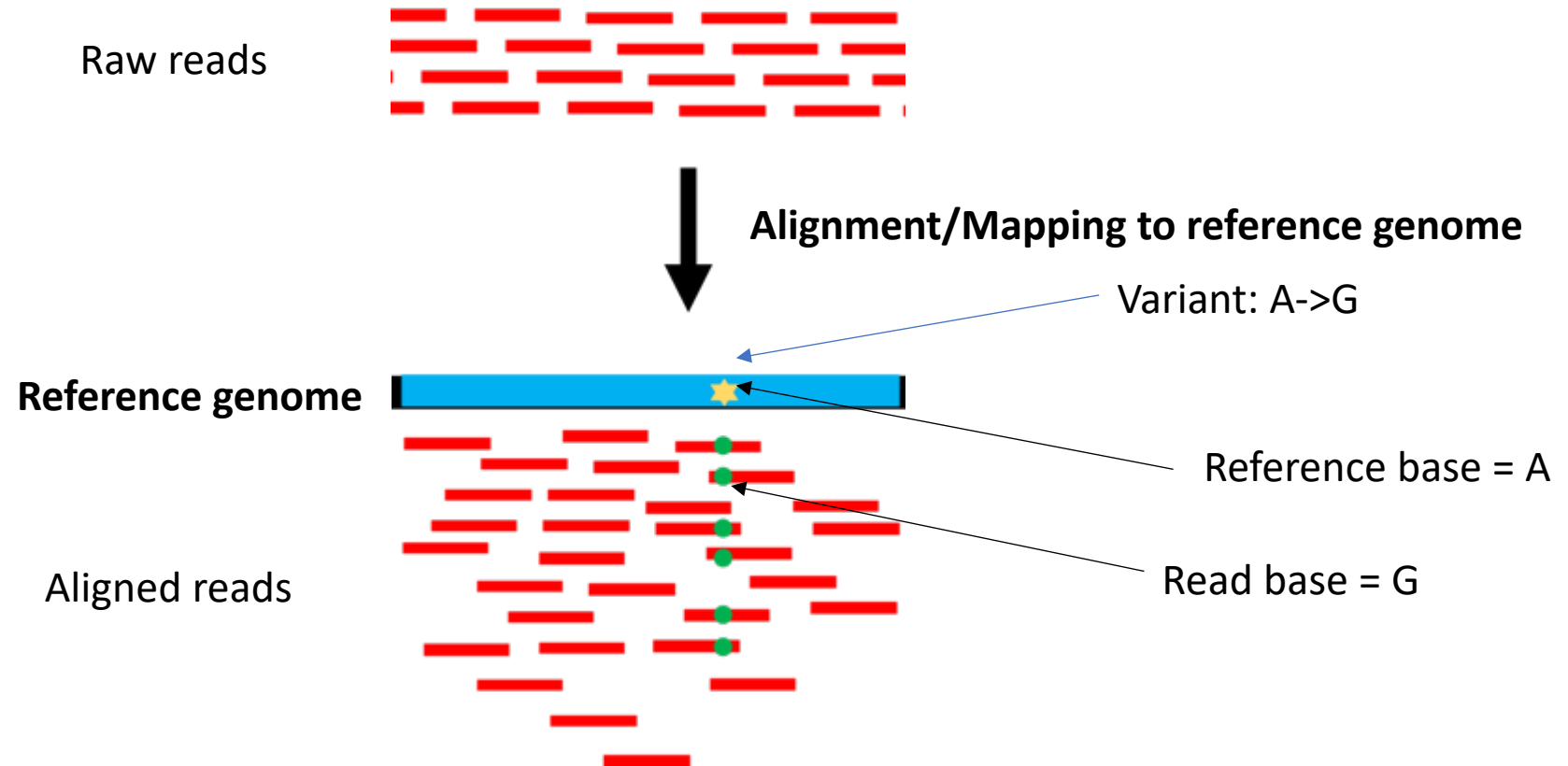
# Sequencing & downstream analysis

- Aim: learn about the genome from the sequenced reads
- Two major analysis pipelines:
  - Assembly
  - Alignment + Variant Calling

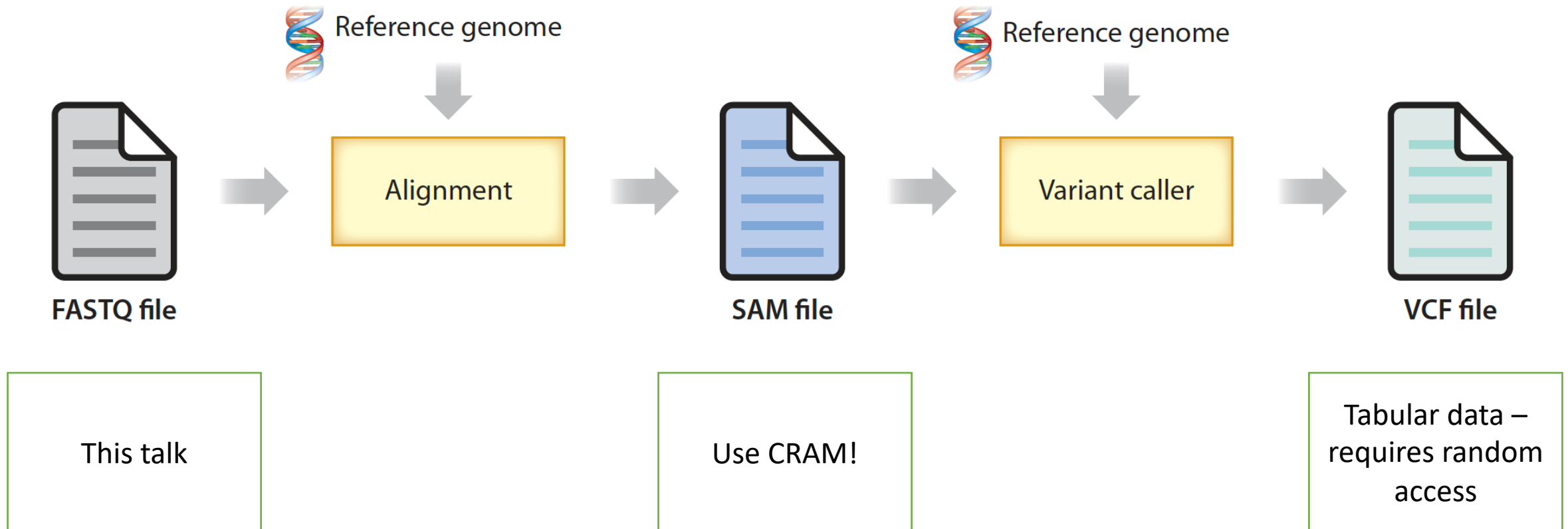
# Genome assembly



# Alignment and Variant Calling



# File formats in the pipeline



# Sequencing & downstream analysis

- Aim: learn about the genome from the sequenced reads
- Two major analysis pipelines:
  - Assembly
  - Alignment + Variant Calling
- Several sequencing methods with different features
  - We focus on two leading technologies



# Sequencing technologies



Illumina NextSeq 550

- High throughput
- Short reads
- Low error rate



Oxford Nanopore MinION

- Portable and real-time
- Long reads
- Native DNA & direct RNA sequencing

Image source:

<https://www.genengnews.com/uncategorized/first-nanopore-sequencing-of-human-genome/>

<https://www.illumina.com/systems/sequencing-platforms/nextseq.html>



# Outline

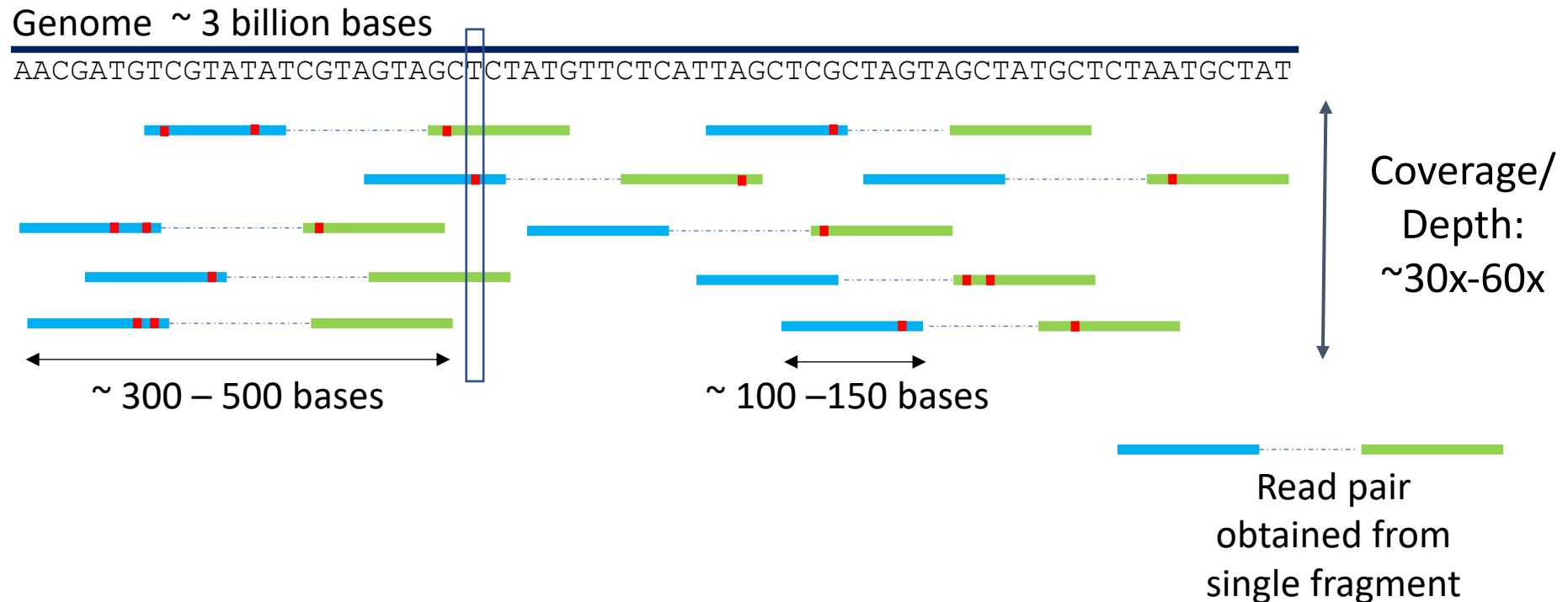
- Introduction to genomic sequencing technologies
- **Genomic data compression: SPRING**
- Using DNA as a storage medium

Chandak, Shubham, et al. "SPRING: a next-generation compressor for FASTQ data." *Bioinformatics* 35.15 (2019): 2674-2676.

*Joint work with Kedar Tatwawadi, Idoia Ochoa, Mikel Hernaez, Tsachy Weissman*

# Paired-end genome sequencing

- Genome: long string of bases {A, C, G, T}
- Sequenced as noisy paired substrings (*reads*):



# Why store raw reads?

- Pipelines improve with time - need raw data for reanalysis
- For temporary storage or regulatory requirements
- When reference genome not available – e.g., de novo assembly or metagenomics

# FASTQ format

## File 1

```
@ERR174324.1 HSQ1009_86:1:1101:1192:2116/1
ATTCNGTCACTTCTCACCAGGCCCTCATTCAACACTGGGAATTAAAATTCGAC...
+
CCCF#2ADHHHHHJJJIJJJIJJJJJJJGIJJJJJJJIJJJIJJJJJGIJJ...
:
```

## Read

## Quality scores

## File 2

@ERR174324.2 HSQ1009\_86:1:1101:1192:2116/2  
CAGANAGAGACTCTGTCTCAAAAAACAAACAAACAAACAAAAGTCTTA...  
+  
CCCF#2ADHFHHJJJJJJJJJJJJJJJJJJJJHIIJJJJJJJIIIJJ...  
⋮

## Read identifier

We'll mostly focus on **reads** in this talk.

# Read compression

- For a typical 25x human dataset:
  - Uncompressed: 79 GB (1 byte/base)
  - Gzip: ~20 GB (2 bits/base) – still far from optimal

# Read compression results

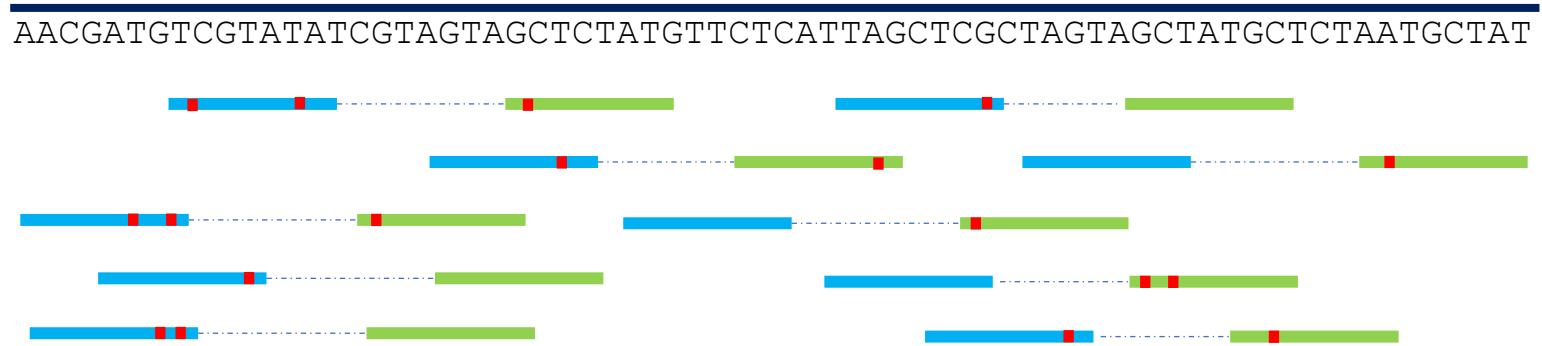
Compressor	25x human
Uncompressed	79 GB
Gzip	~20 GB
<b>SPRING</b>	<b>3 GB</b>

# Read compression results

Compressor	25x human	100x human
Uncompressed	79 GB	319 GB
Gzip	~20 GB	~80 GB
<b>SPRING</b>	<b>3 GB</b>	<b>10 GB</b>

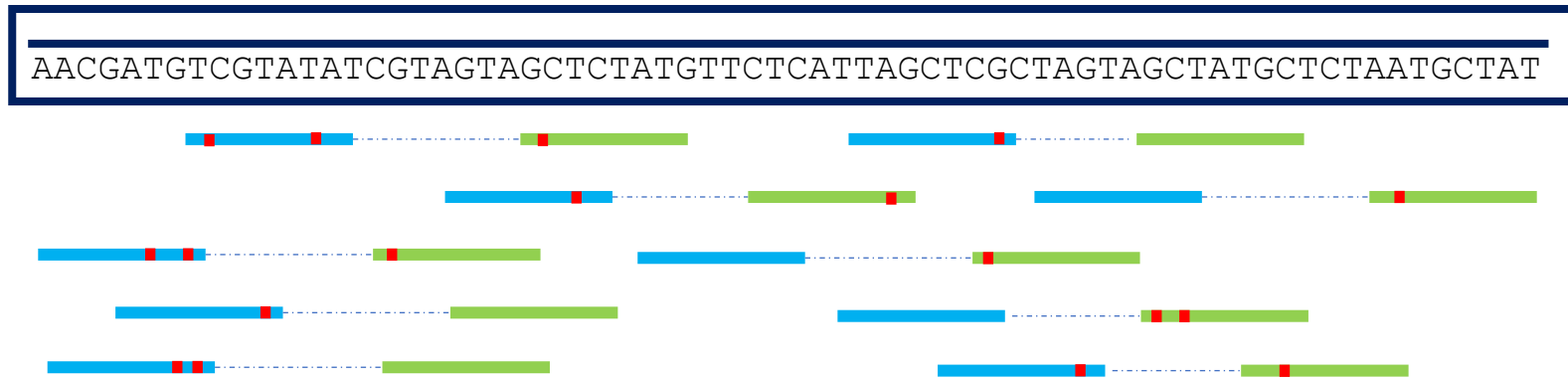


# Key idea



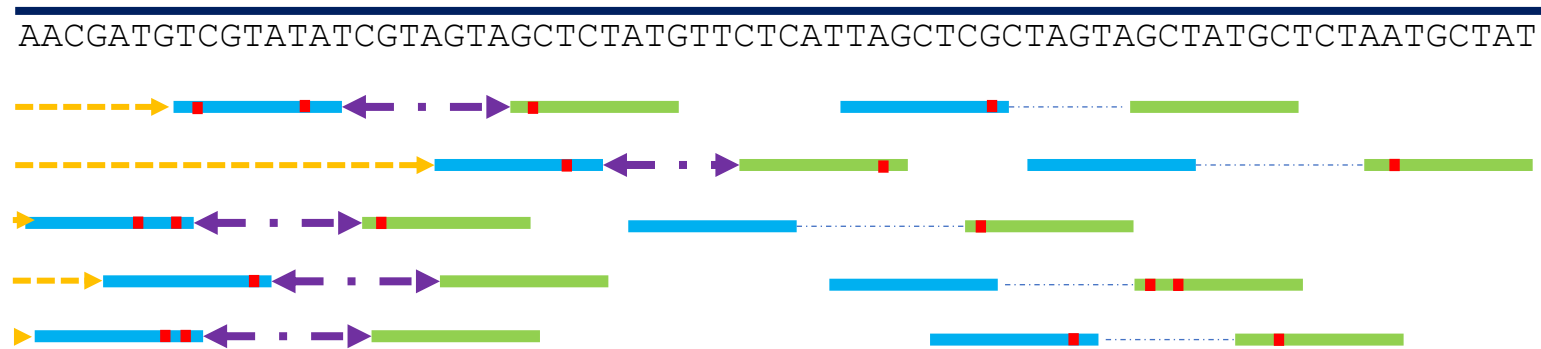
- Storing reads equivalent to

# Key idea



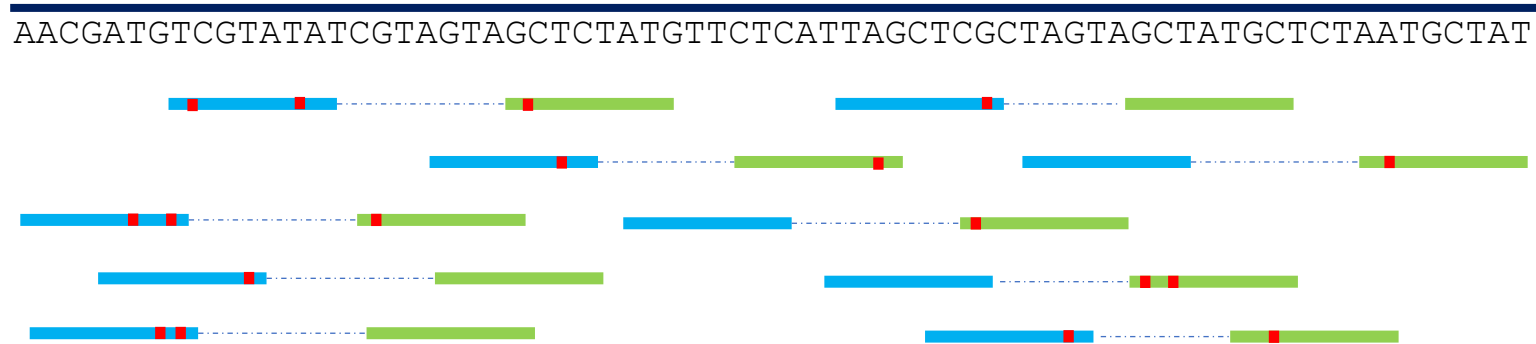
- Storing reads equivalent to
  - Store genome

# Key idea



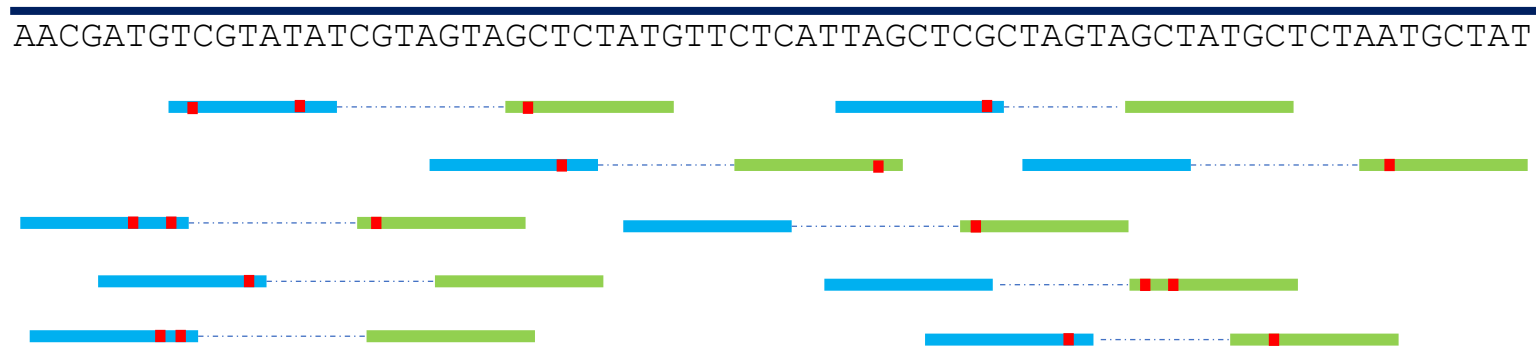
- Storing reads equivalent to
  - Store genome
  - Store read positions in genome (+ gap between paired reads)

# Key idea



- Storing reads equivalent to
  - Store genome
  - Store read positions in genome (+ gap between paired reads)
  - Store noise in reads

# Key idea

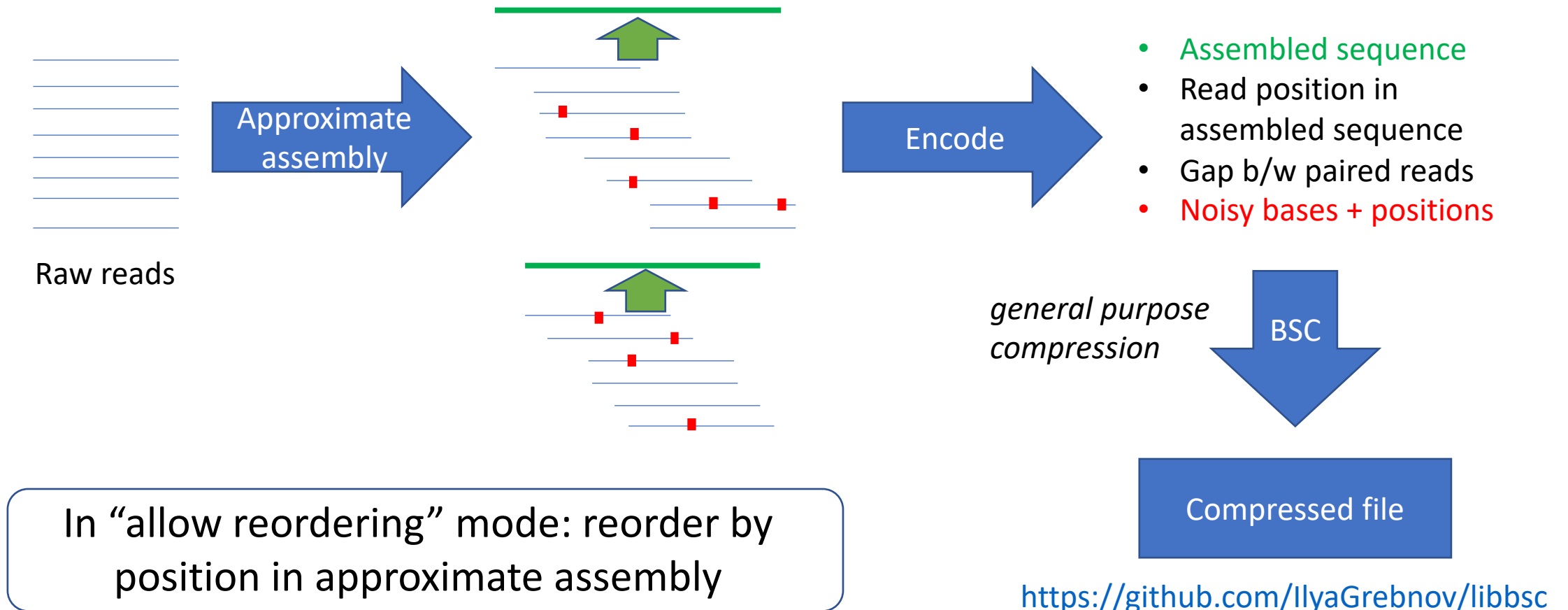


- Storing reads equivalent to
  - Store genome
  - Store read positions in genome (+ gap between paired reads)
  - Store noise in reads
- Theoretical calculations show this outperforms previous compressors

# Key idea

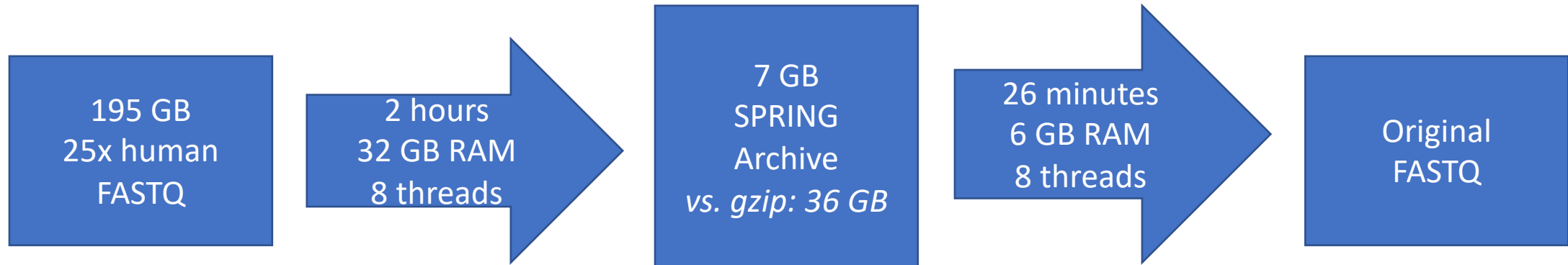
- But... How to get the genome from the reads?
- Genome assembly too expensive - big challenges:
  - resolve repeats
  - get very long pieces of genome from shorter assemblies
- Solution: Don't need perfect assembly for compression!

# SPRING workflow





# SPRING as a practical tool



- Easy to use with support for:
  - Lossless and lossy modes
  - Variable length reads, long reads, etc.
  - Compressed in blocks to allow partial/streaming decompression
  - Scalable to large datasets
  - Gzipped I/O
- GitHub: <https://github.com/shubhamchandak94/SPRING/>

# Future directions

- Another paradigm: reference-based FASTQ compression
  - Illumina ORA/Enancio, Petagene
- More recent work on compression for long read data
  - Meng, Q., Chandak, S., Zhu, Y., & Weissman, T. (2021). NanoSpring: reference-free lossless compression of nanopore sequencing reads using an approximate assembly approach. *bioRxiv*.
  - Shubham Chandak, Kedar Tatwawadi, Srivatsan Sridhar, Tsachy Weissman, Impact of lossy compression of nanopore raw signal data on basecalling and consensus accuracy, *Bioinformatics*, Volume 36, Issue 22-23, 1 December 2020, Pages 5313–5321.

# Outline

- Introduction to genomic sequencing technologies
- Genomic data compression: SPRING
- Using DNA as a storage medium

Lau, Billy T., **Chandak S.**, et al. "Magnetic DNA random access memory with nanopore readouts and exponentially-scaled combinatorial addressing." *bioRxiv* (2021).

**S. Chandak et al.**; "Overcoming high nanopore basecaller error rates for DNA storage via basecaller-decoder integration and convolutional codes," *ICASSP 2020*.

**S. Chandak et al.**; "Improved read/write cost tradeoff in DNA-based data storage using LDPC codes," *Allerton 2019*.

# Team and funding



Shubham  
Chandak



Joachim  
Neu



Jay  
Mardia



Billy  
Lau



Matt  
Kubit



Reyna  
Hulett



Peter  
Griffin



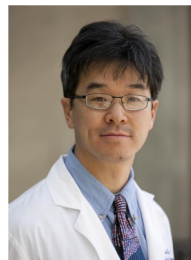
Sharmili  
Roy



Tsachy  
Weissman



Mary  
Wootters



Hanlee  
Ji



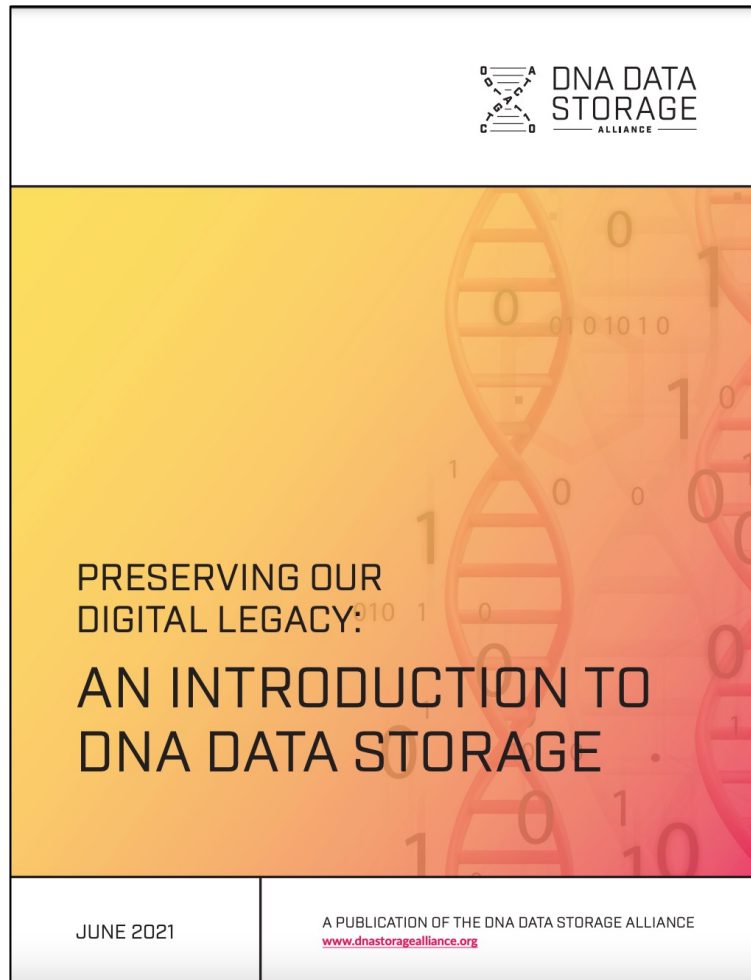
**SemiSynBio: Highly scalable random  
access DNA data storage with  
nanopore-based reading**

*Beckman Center Innovative Technology Seed Grant*

**Scalable Long-Term DNA Storage with Error Correction and  
Random-Access Retrieval**



National Institutes  
of Health



## FOUNDERS

illumina®

Illumina

Microsoft

Microsoft

T W I S T  
BIOSCIENCE

Twist Bioscience

Western Digital.

Western Digital

**BBC** | [Sign in](#) | [Home](#) | [News](#) | [Sport](#) | [Reel](#) | [Worklife](#) | [Travel](#)

# NEWS

[Home](#) | [War in Ukraine](#) | [Coronavirus](#) | [Climate](#) | [Video](#) | [World](#) | [Asia](#) | [UK](#) | [Business](#) | [Tech](#) | [Science](#)

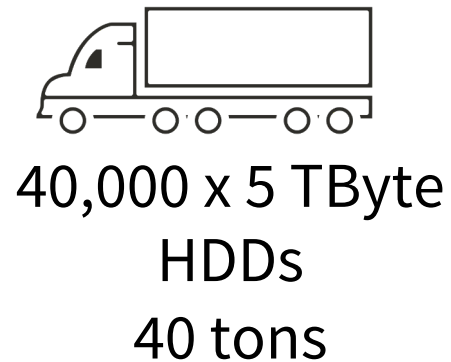
Science

## Scientists claim big advance in using DNA to store data

By Paul Rincon  
Science editor, BBC News website

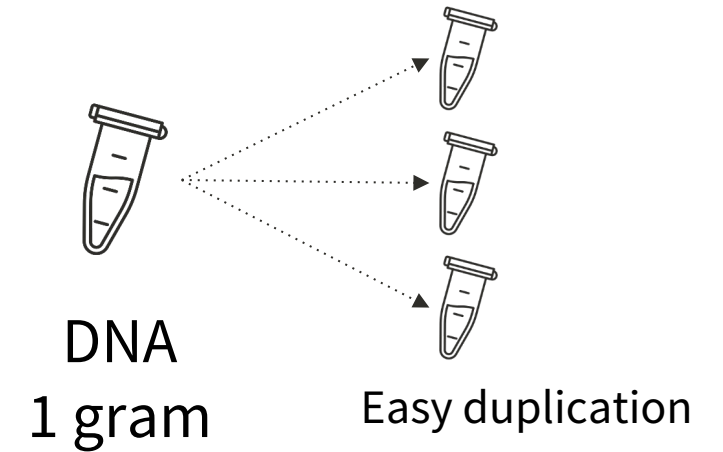
🕒 1 December 2021

# Why DNA-based Storage?



10's of years

**200  
Petabyte**



1000's of years

# Building Blocks

- Ability to “**read/sequence**” the DNA from the solution.



- Ability to “**write/synthesize**” artificial DNA (sequence of {A,C,G,T})



**Agilent Technologies**

**Current ability:** short DNA oligo sequences (~**150** length) at scale  
(Array Synthesis)



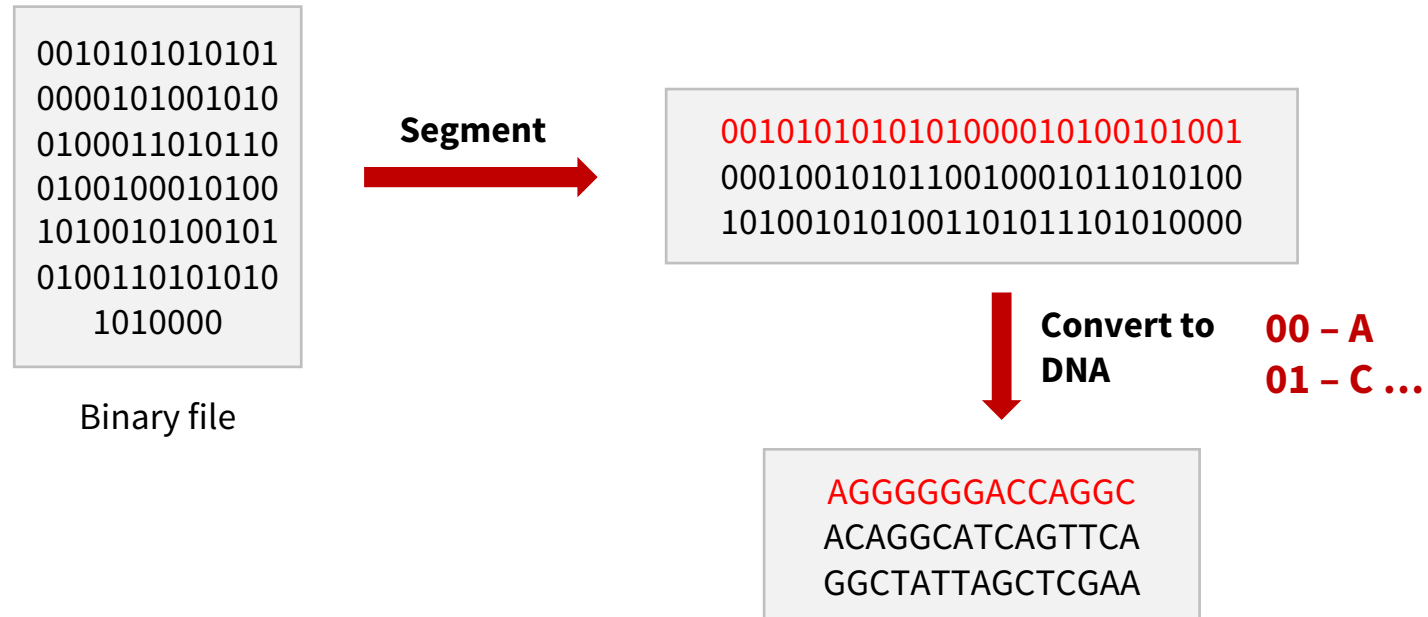
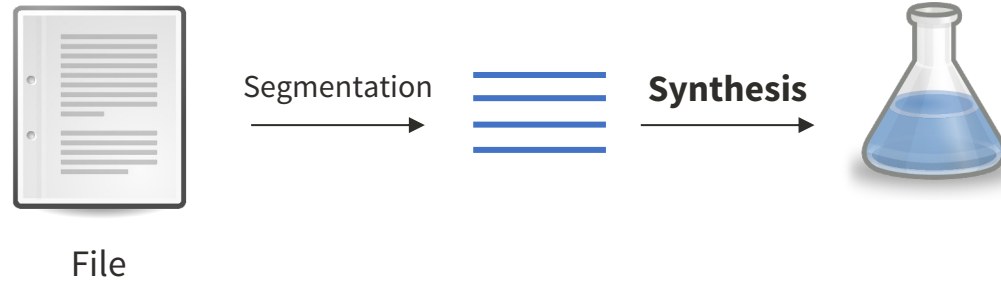
# Typical DNA Storage System



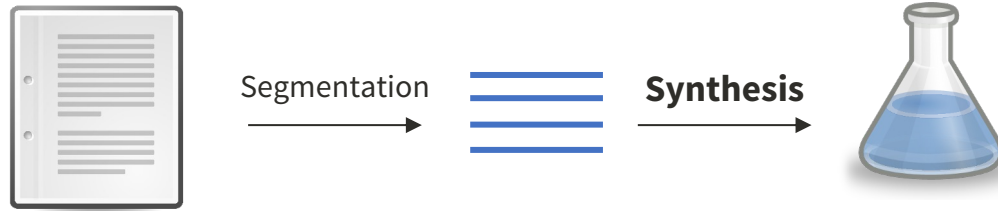
File

**Can only synthesize short DNA  
oligo sequences ~150 bases**

# Typical DNA Storage System

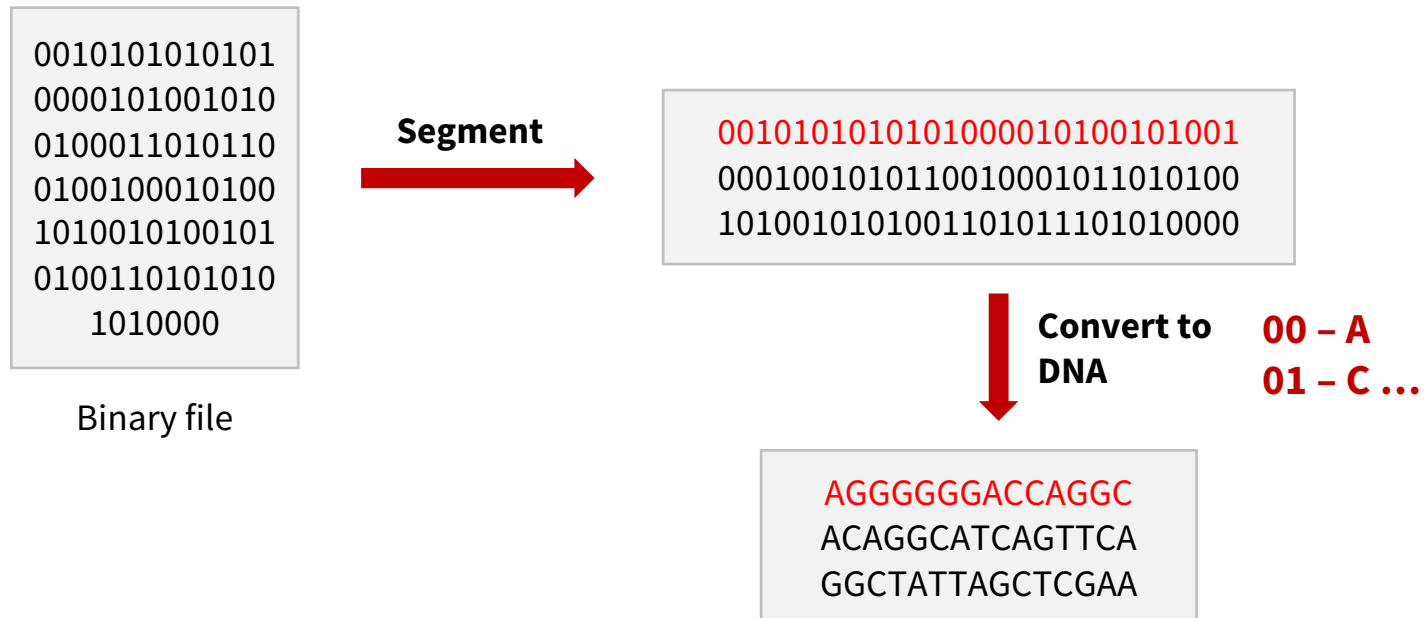


# Typical DNA Storage System



File

**Order of sequences lost  
in the solution!**



# Typical DNA Storage System

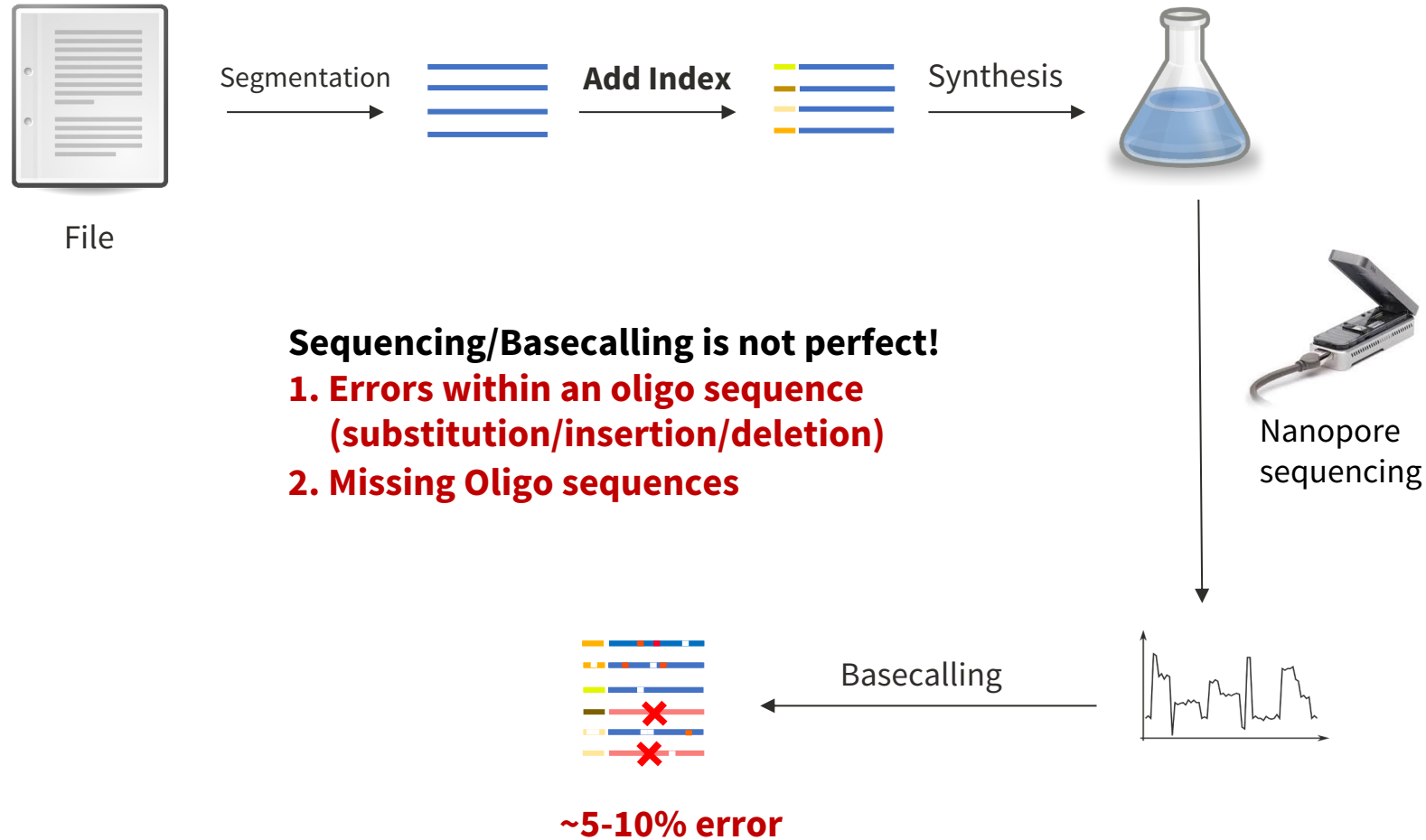


**Order of sequences lost in the solution! – Add Index**

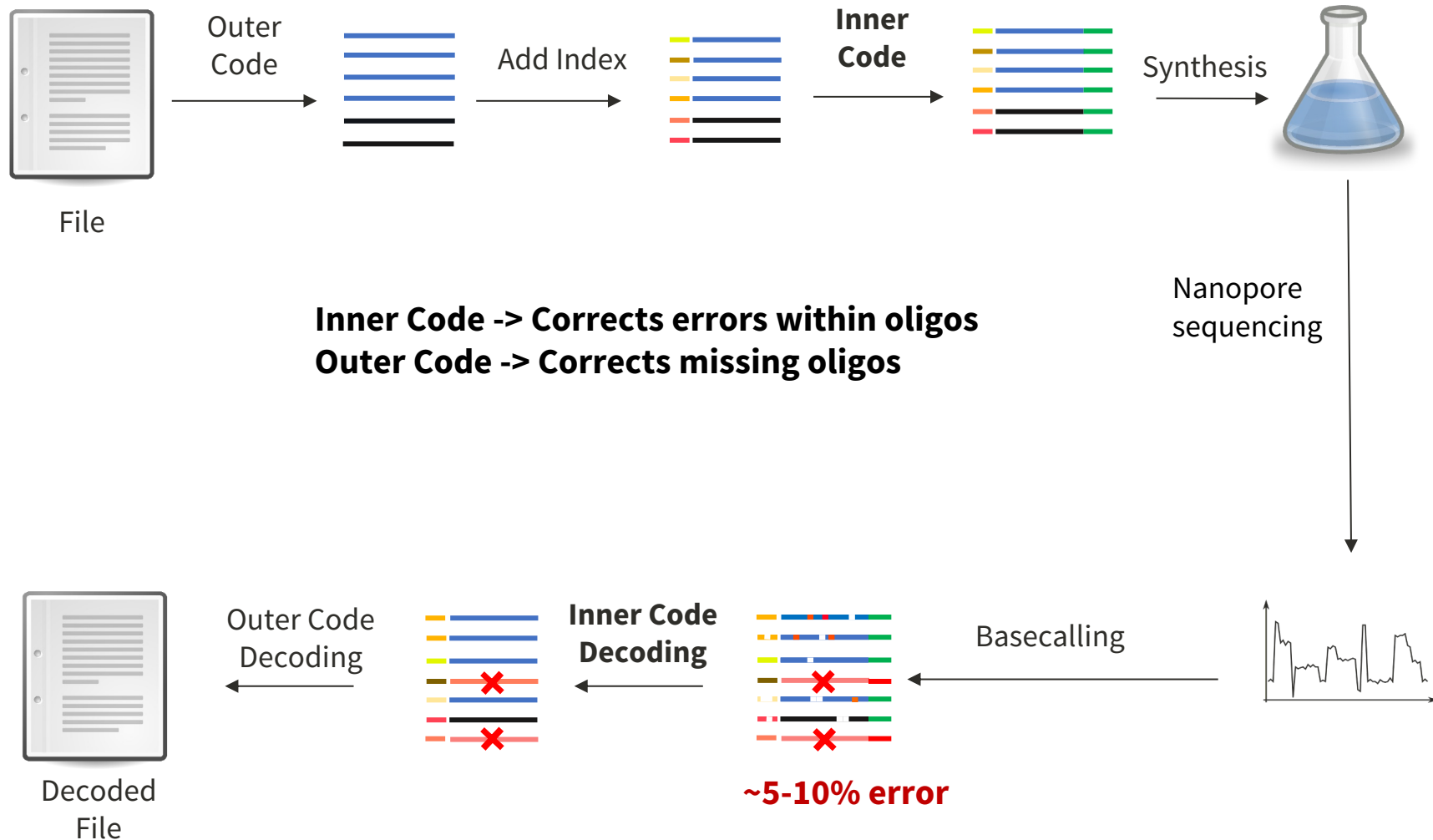
```
000010101010101000010100101001  
000101001010010010001010010100  
001010010101001001010101010000
```

Length of index in binary segment at least  $\log_2(\text{number of segments})$

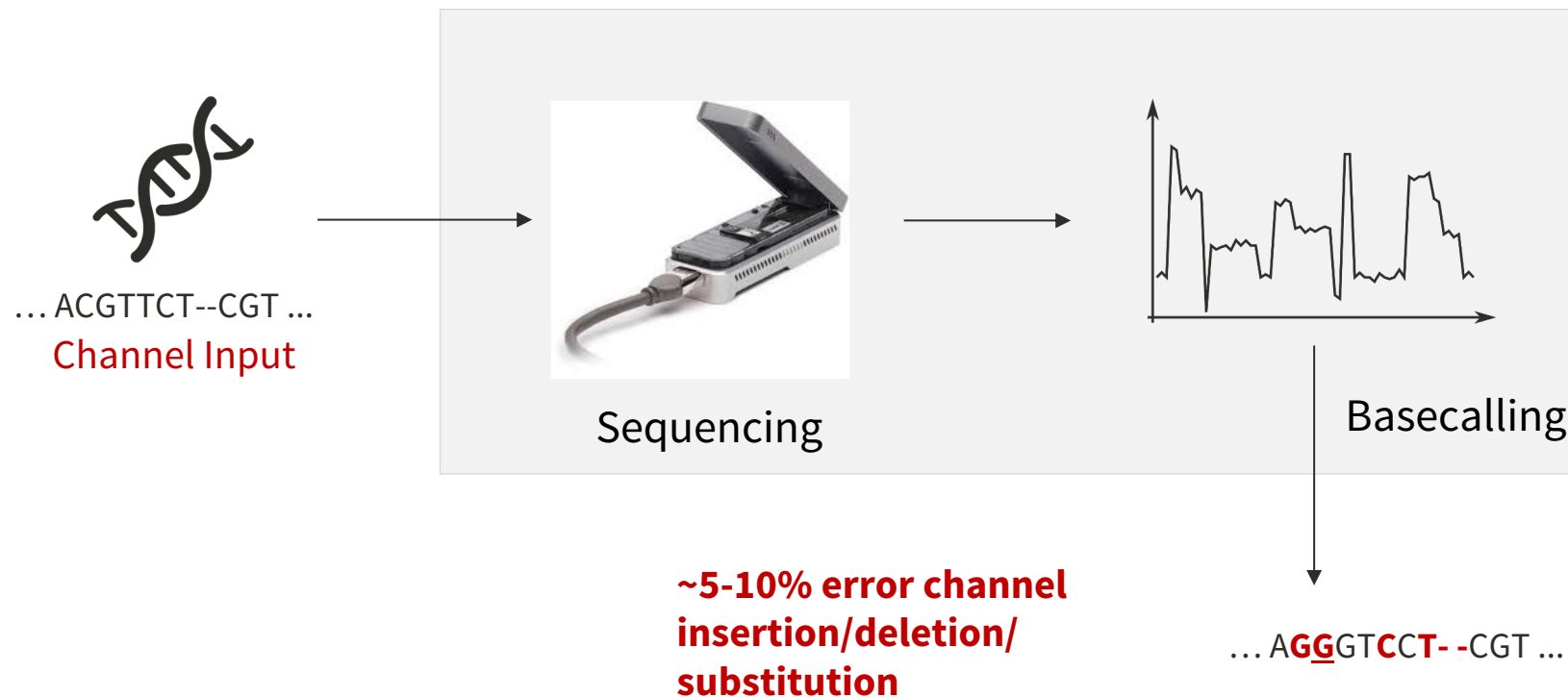
# Typical DNA Storage System



# Typical DNA Storage System

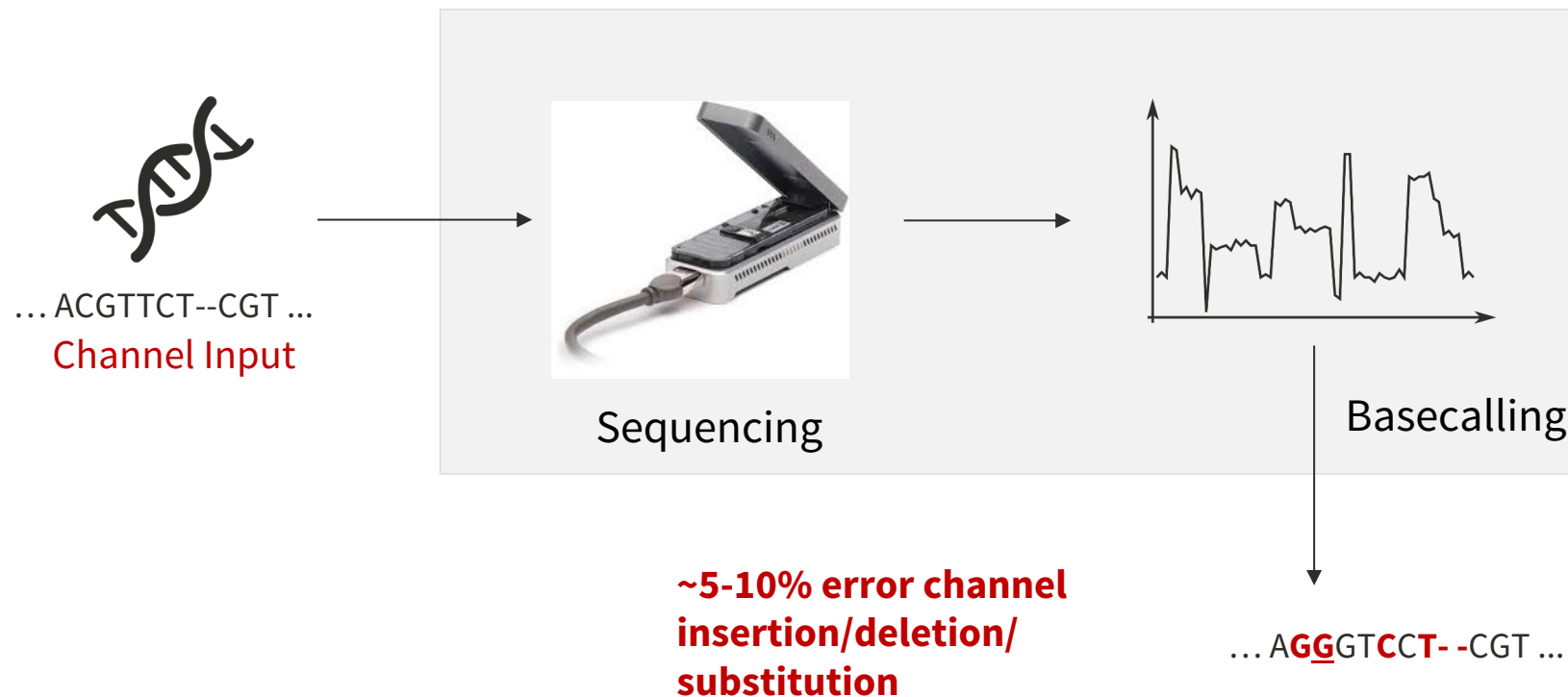


# Channel Model – Insertion/Deletion Channel



# Channel Model – Insertion/Deletion Channel

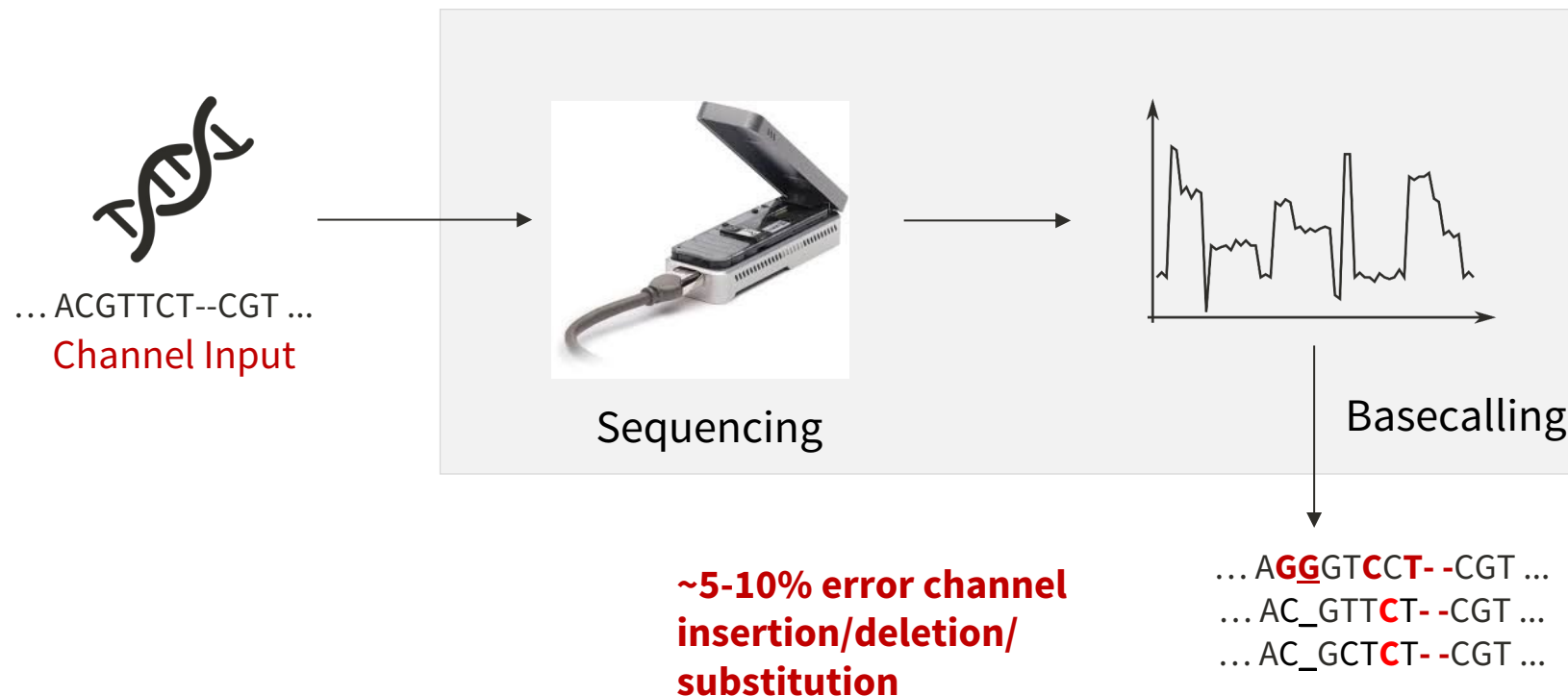
- **Basecalling Error:** No good practical error correction code for 5-10% Insertion/Deletions



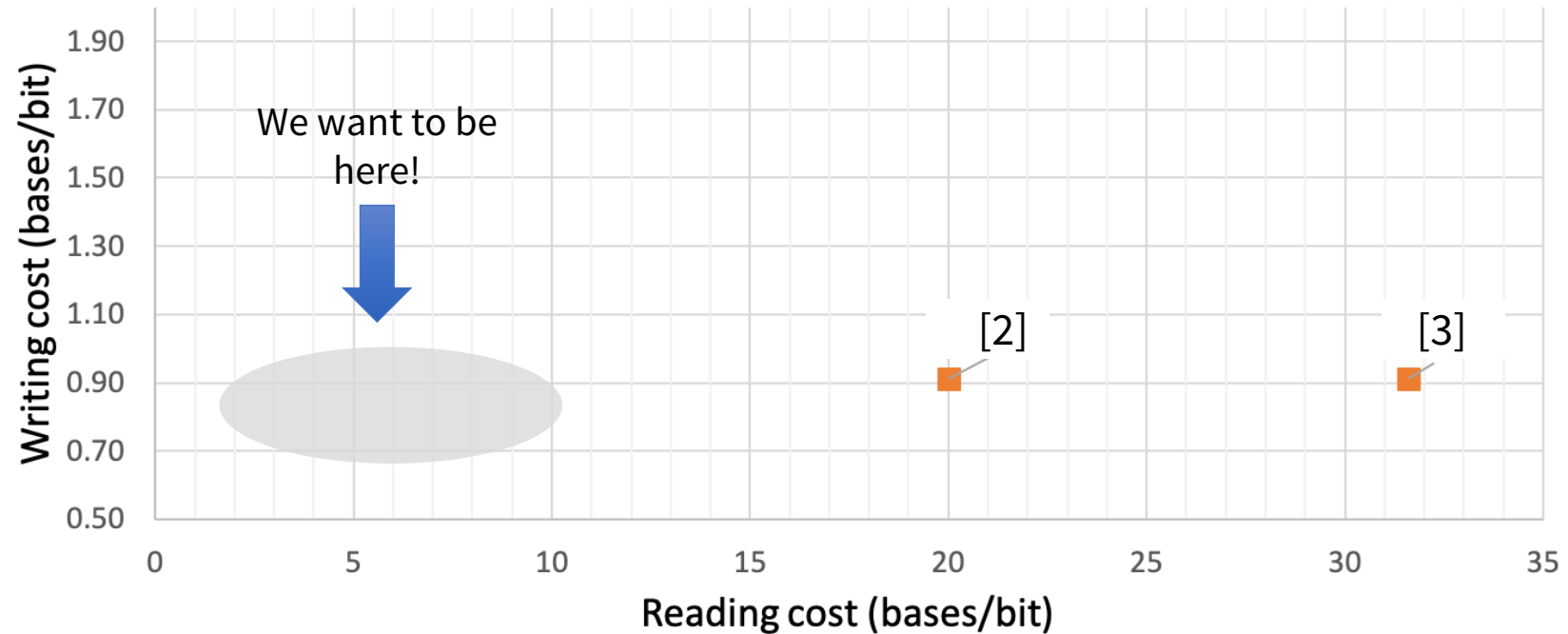


# Channel Model – Insertion/Deletion Channel

- **Basecalling Error:** No good practical error correction code for 5-10% Insertion/Deletions
- **Common Idea:** Sequence the input lot of times (~30-40x)
  - cluster *index*, and perform “averaging” to reduce the error



# Previous Works

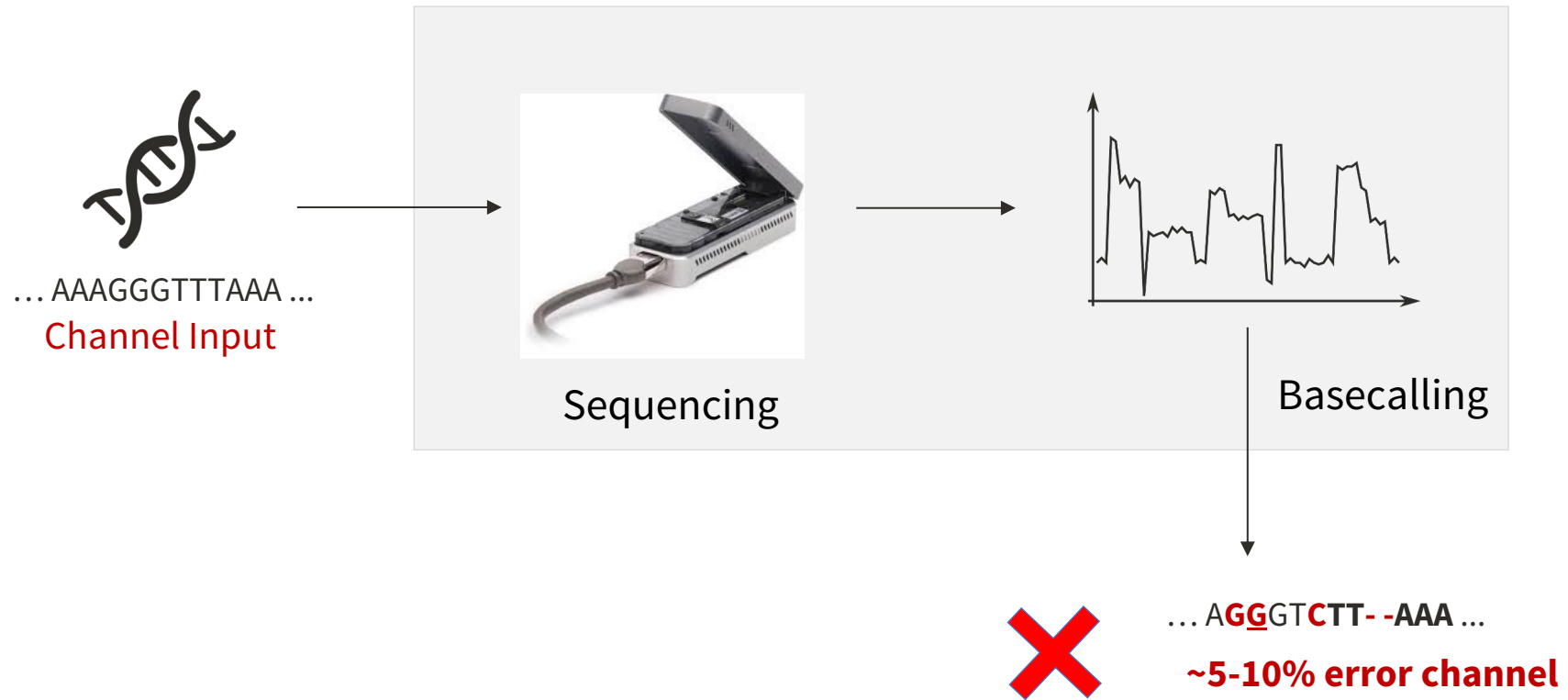


## Tradeoff between reading and writing costs

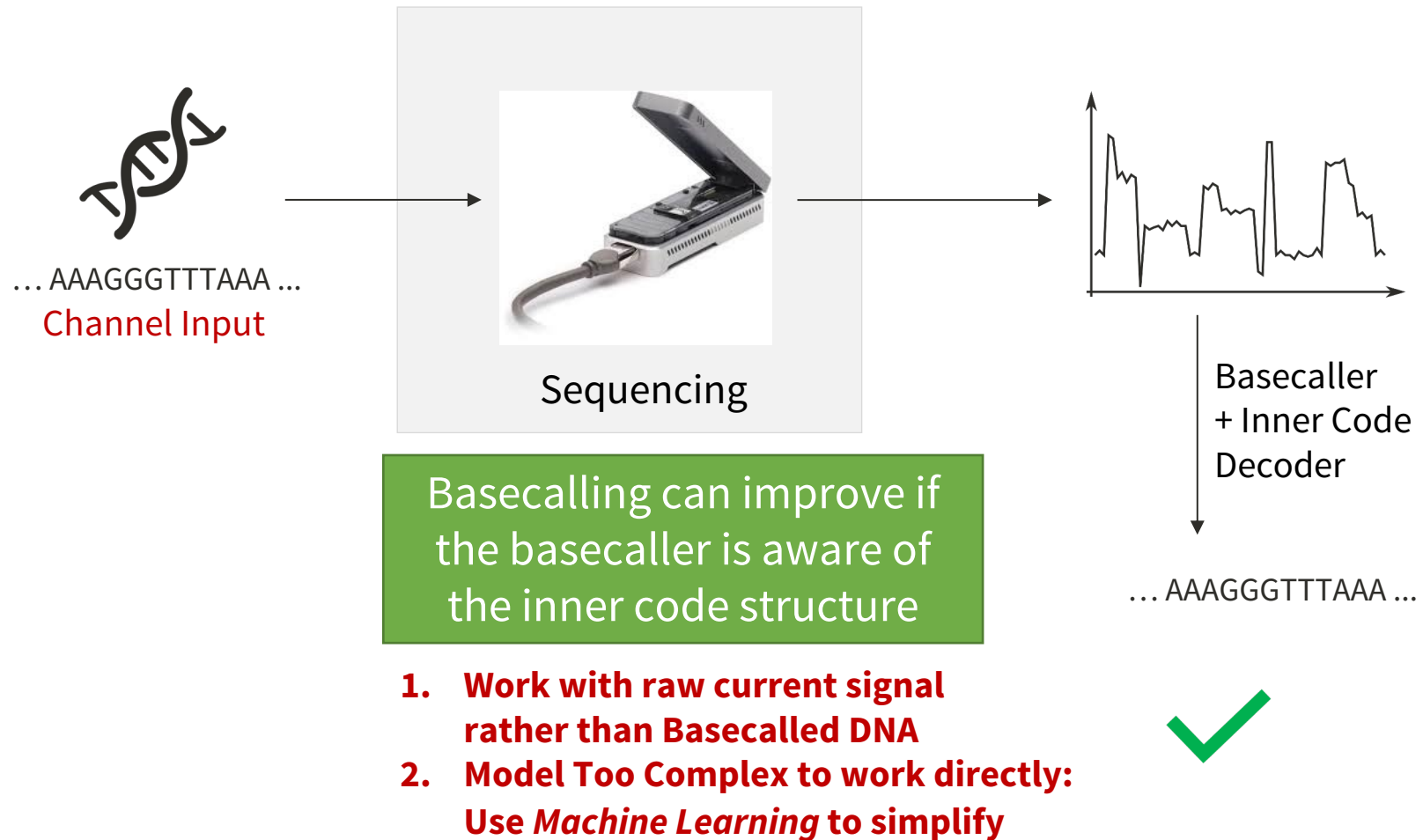
[2] L. Organick *et al.*, "Random access in large-scale DNA data storage," *Nature biotechnology*, vol. 36, no. 3, p. 242, 2018.

[3] Randolph Lopez *et al.*, "DNA assembly for nanopore data storage readout," *Nature communications*, vol. 10, no. 1, pp. 2933, 2019.

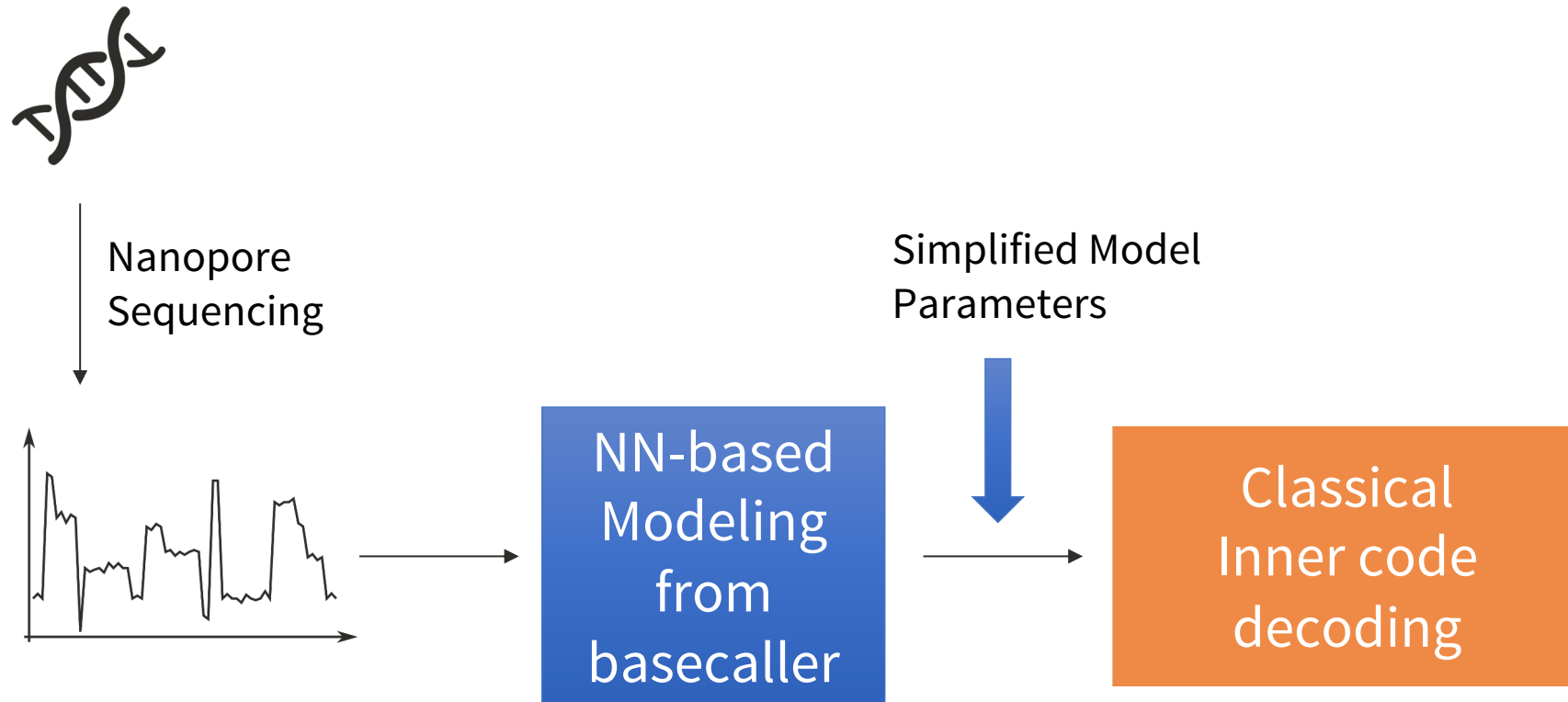
# Nanopore Error Model



# Key Insight!

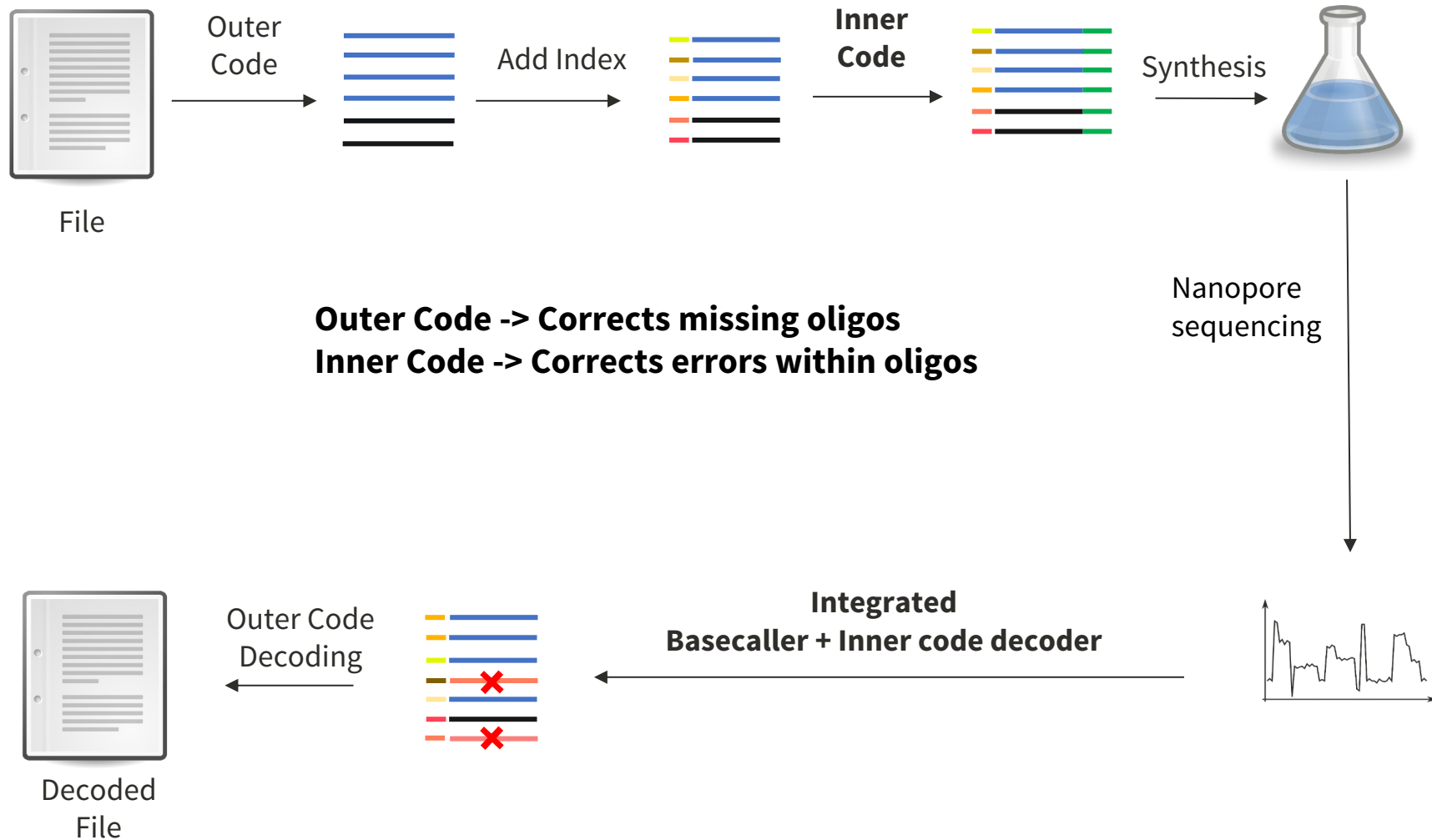


# Inner code decoding

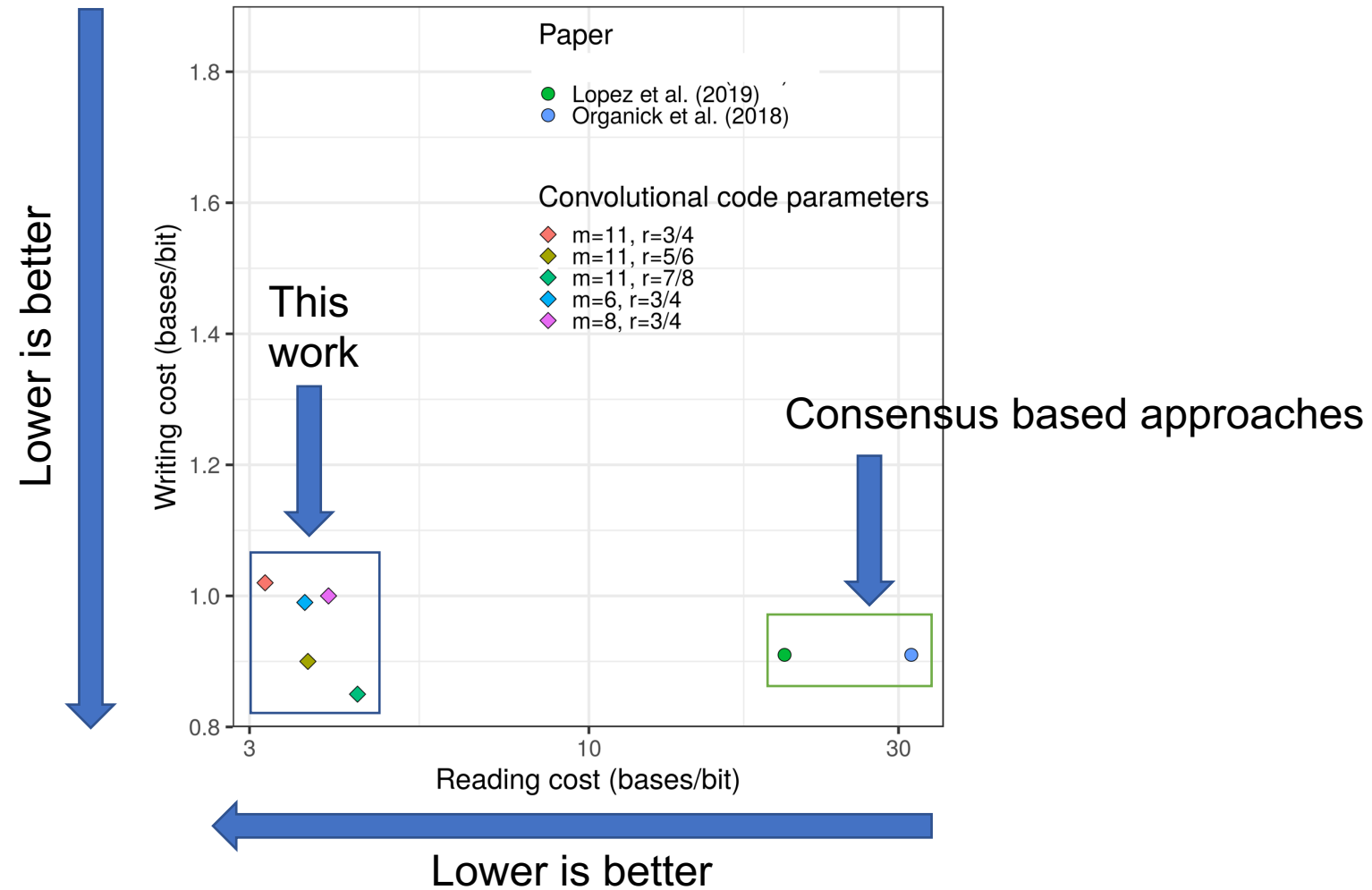


NN-based model simplifies the complex model into a simpler Markov model!  
It repurposes the basecaller's NN model which is optimized based on large amounts of data.

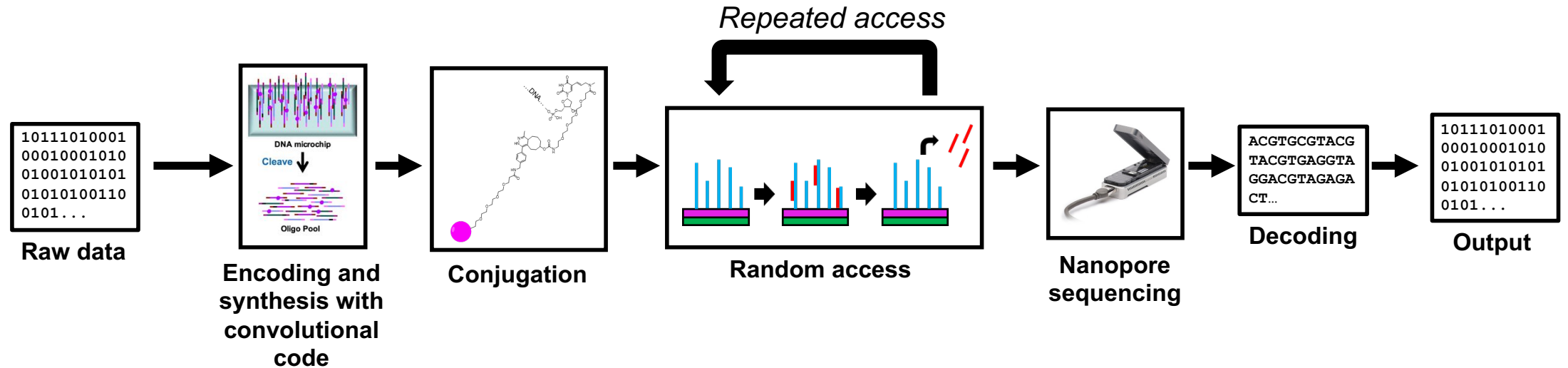
# Integrated Decoding System



# Results



# Overall approach





# Future directions

- Automation of DNA data retrieval with liquid handling robots
- Possibility of real-time data decoding with nanopore sequencers
- Design cheaper, possibly more error prone synthesis platforms

Thank you!