

**IDA  
2022**

**Symposium on Intelligent Data Analysis 2022**  
**April 20–22, 2022, Rennes, France**

# Data Storage on DNA

Dominique Lavenier

GenScale

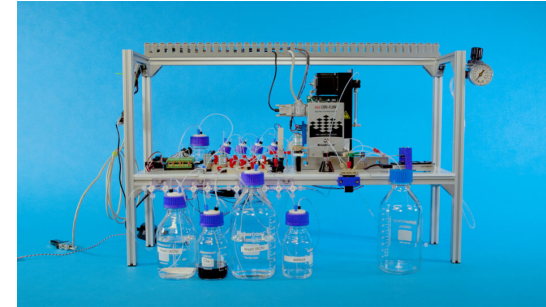
*Univ. Rennes, IRISA/CNRS, Inria*



# Agenda

## 1. Introduction

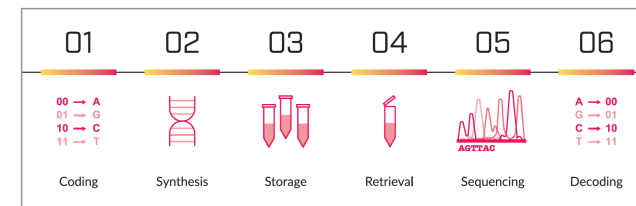
- Why new storage medium ?
- Why DNA ?



Source: <https://news.microsoft.com/innovation-stories/hello-data-dna-storage/>

## 2. DNA storage principle

- How does it work ?
- What are the main challenges ?



Source: *An Introduction to DNA Data Storage*, DNA storage alliance, June 21

## 3. dnarXiv project

- What's going on in Rennes ?



# Exponential Data Growth

According to IDC\*

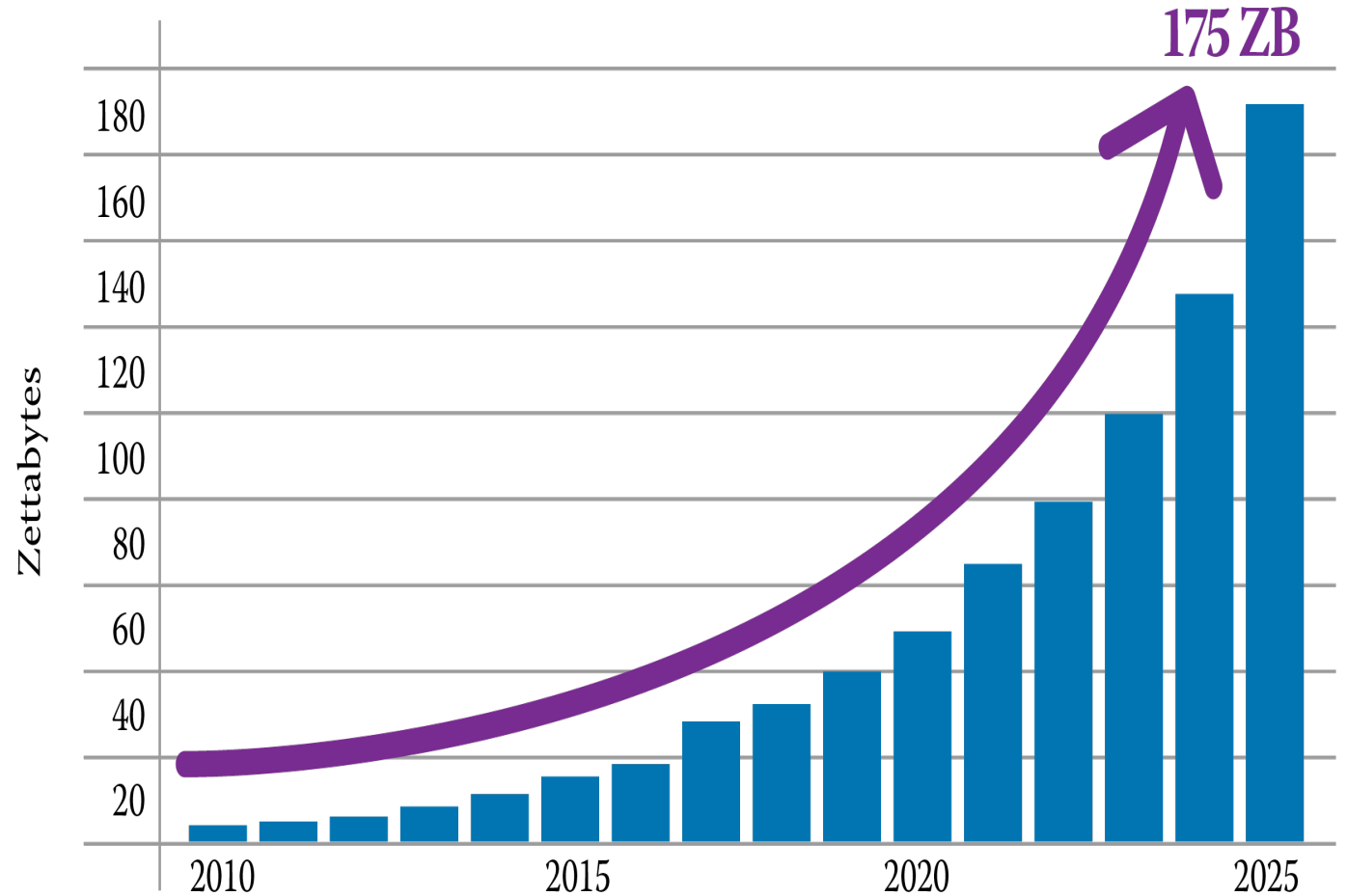
- annual data growth rate:  
→ 23% [2020-2025]
- should reach 175 ZB in 2025

1 Zetta Bytes (ZB) =  $10^{21}$  Bytes

- 1 billion of TB

IoT devices are expected to create over 90 ZB of data in 2025

IDC: International Data Corporation

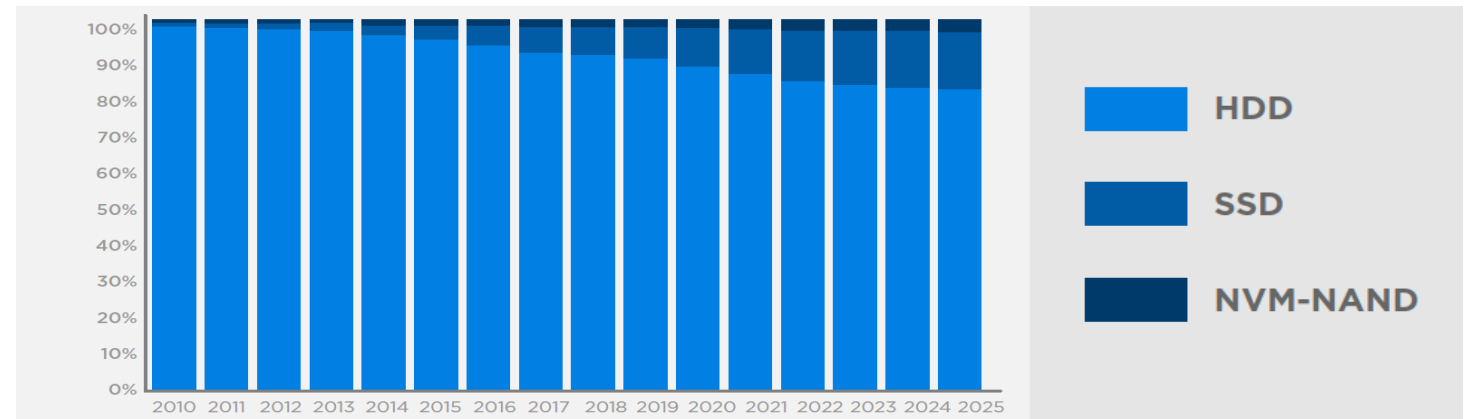
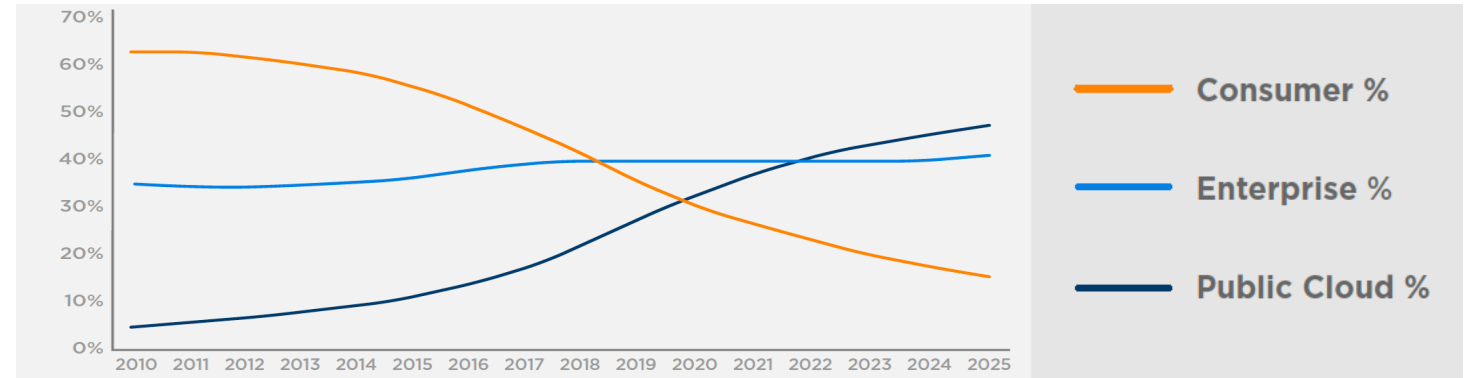


Source: IDC White paper, Data Age 2025, nov. 2018, The Digitization of the world, From Edge to Core

# Where the data will be stored ?

By 2025:

- need to store **22 ZB** of data
- **85%** of all data worldwide will reside in public or private cloud environment
- **80%** of cloud data will be stored on HDD



Source: IDC White paper, Data Age 2025, nov. 2018, The Digitization of the world, From Edge to Core

# Challenges for archival storage technologies

## Storage maintenance and replacement costs

- HDD / SSD devices life span : 10 years (max)



Device storage have to be periodically changed

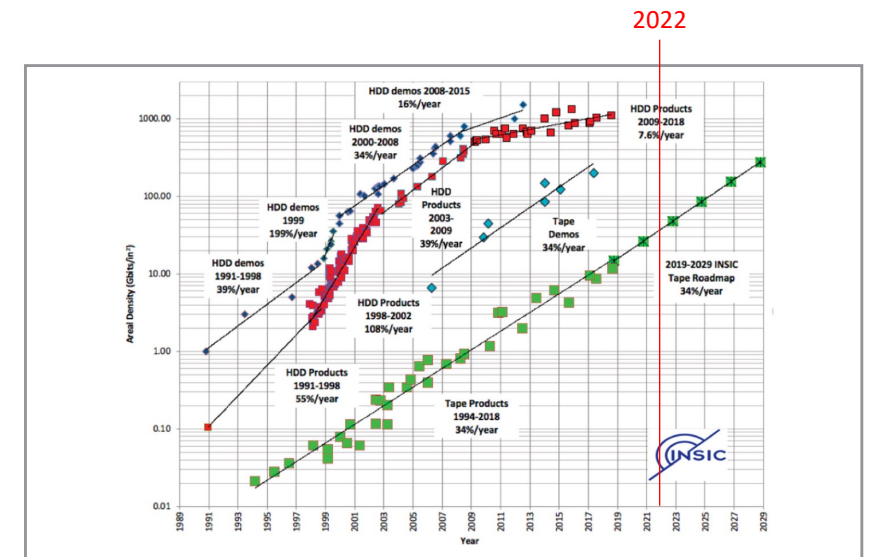
## Density limitations

- Areal density for magnetic media is slowing



## Energy and sustainability concerns

- Increasing electric consumption of data centers
- HDD relies on rare earth metals



# Storage hierarchy

## HOT DATA

- data accessed frequently
- storage: high performance devices: SSD

## WARM DATA

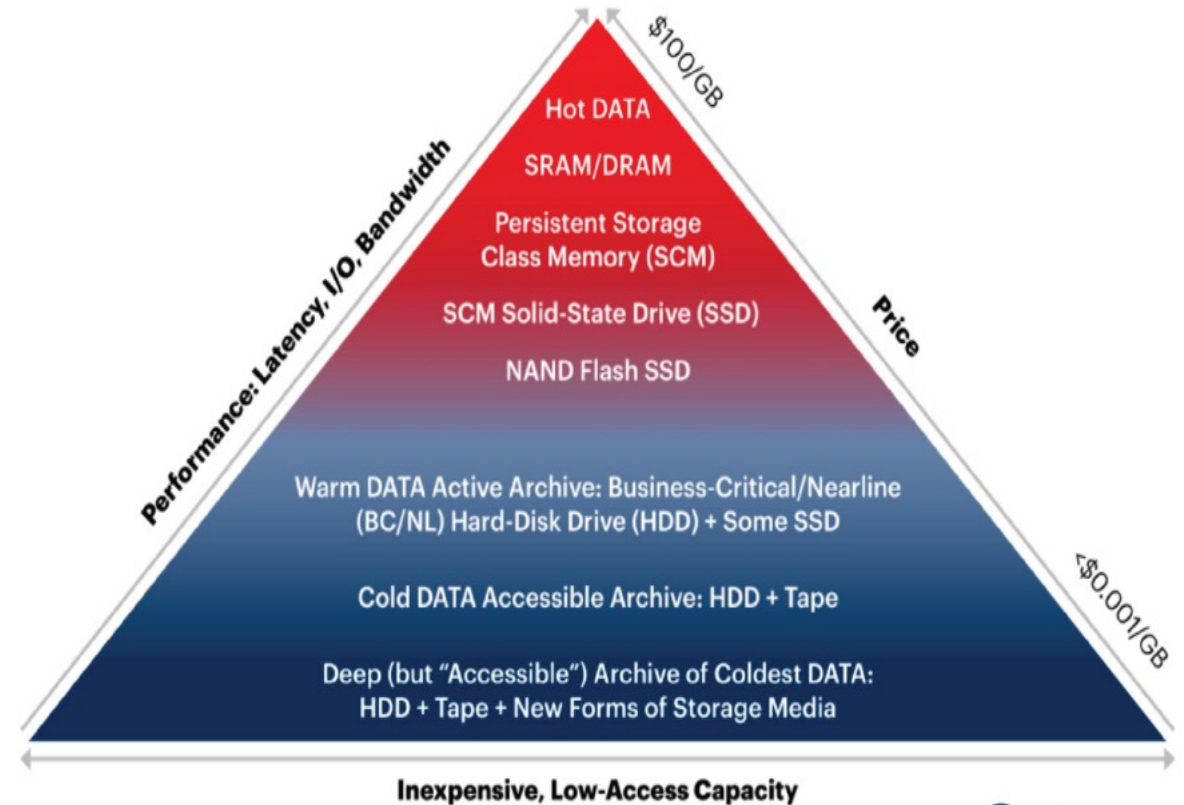
- data accessed less frequently
- storage: HDD

## COLD DATA

- data accessed unfrequently
- storage: tapes

## COLDEST DATA

- data never accessed (or with low probability)
- storage: tapes or new storage media



Gartner

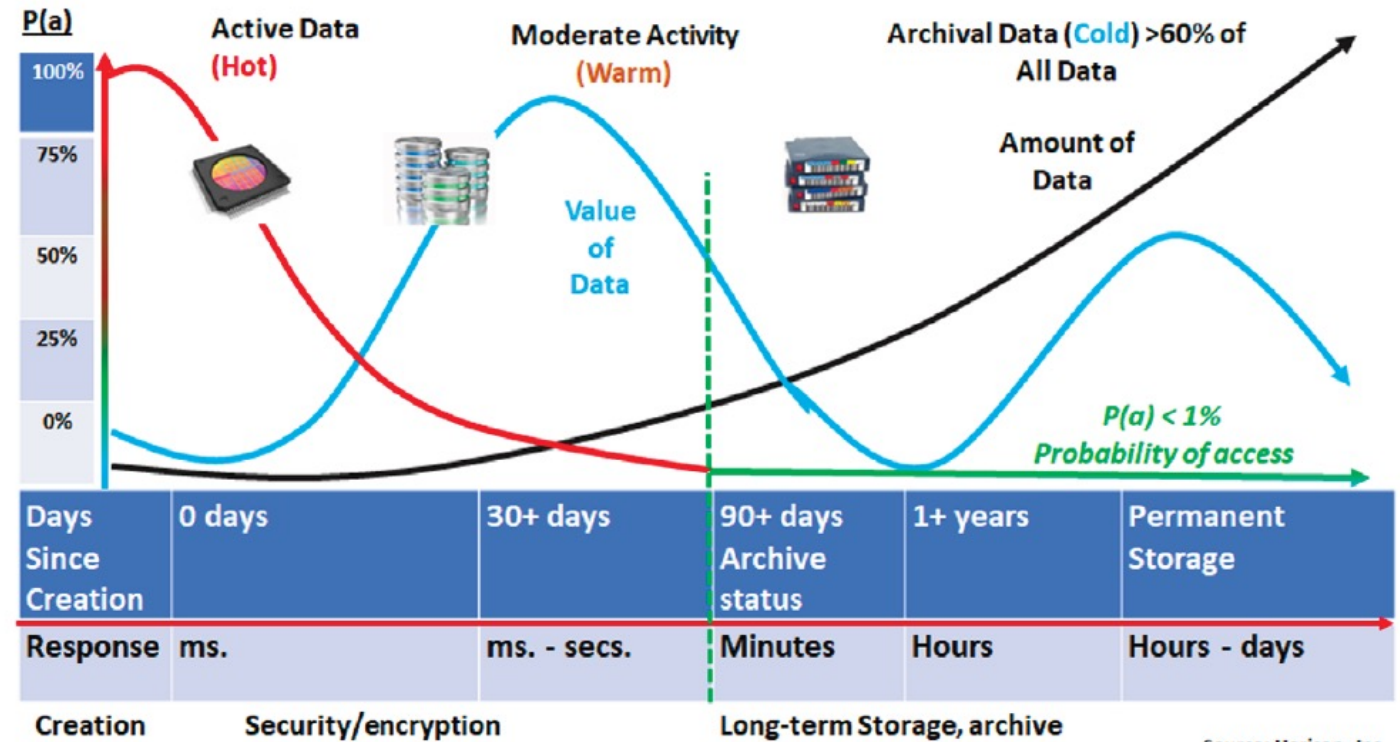
Source: Source: Gartner Product Manager Insight: Tape to the Future? Oct. 2020

# Data life span

Most data become archival after 90 days

Over 60% of all data is archival

Retention of 50-100+ years is common



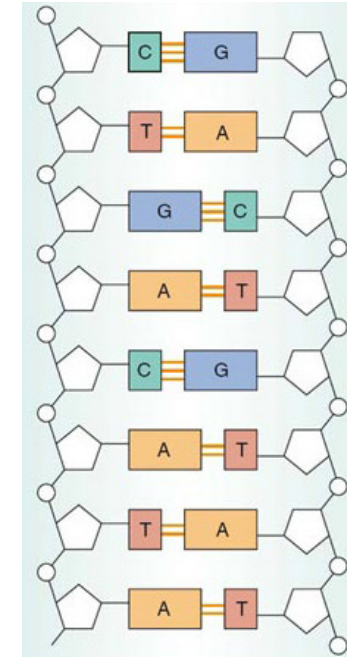
Source: Horison, Inc

# DNA as storage media

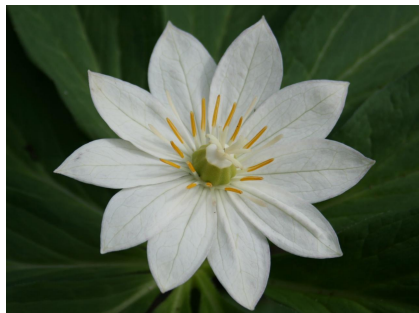
## Advantages

- density
- durability
- medium immutability
- low TCO (Total Cost of Ownership)
- Energy efficiency and sustainability

DNA molecule



4 bases: A C G T  
| | | | base pairs (bp)  
T G C A



	<i>Genome size</i>
Virus	$3 \times 10^4$ bp
Bacteria	$3 \times 10^6$ bp
<b>Human</b>	$3 \times 10^9$ bp
Salamander	$30 \times 10^9$ bp
Paris Japonica	$150 \times 10^9$ bp

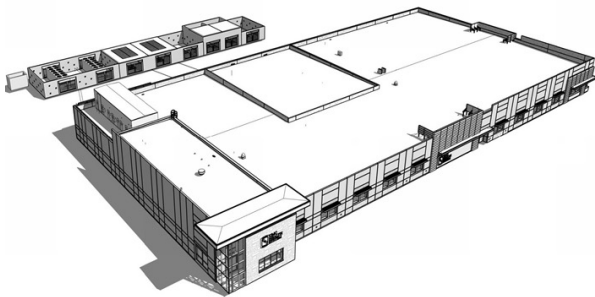




# DNA as storage medium

## Advantages

- density
- durability
- medium immutability
- low TCO (Total Cost of Ownership)
- Energy efficiency and sustainability

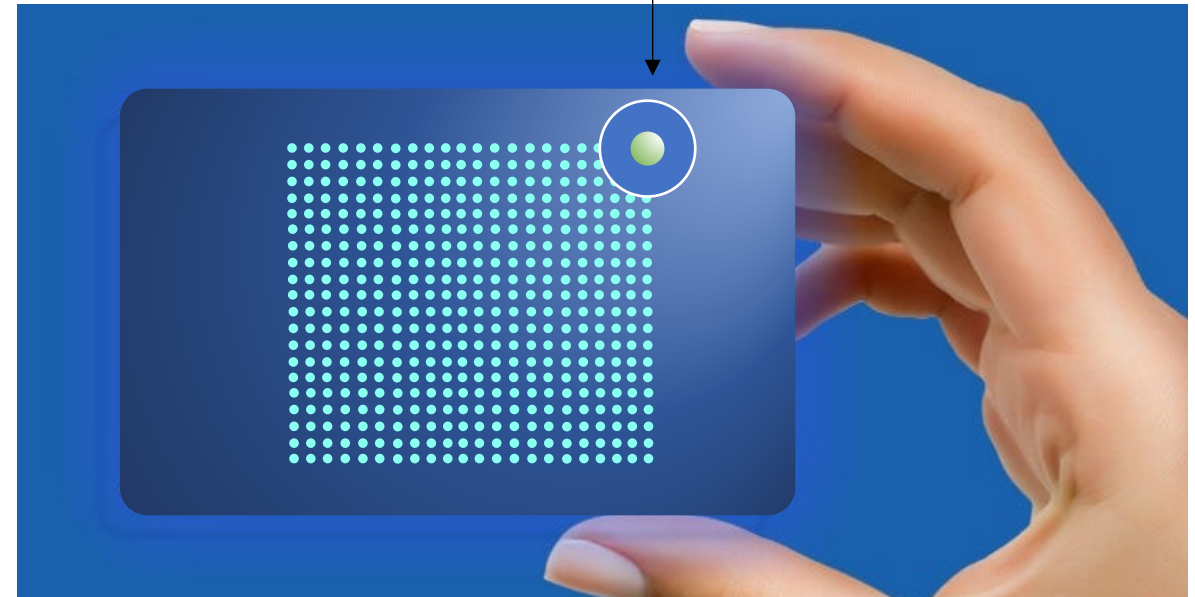


1 ExaBytes of cold data



2,500 credit cards

1 mm dot → 1 TB of data



400 hard drives on a credit card

# DNA as storage medium

## Advantages

- density
- **durability**
- medium immutability
- low TCO (Total Cost of Ownership)
- Energy efficiency and sustainability

Stable at room temperature for a very long period

No need to replicate over time



**nature**

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾ [Subscribe](#)

[nature](#) > [articles](#) > [article](#)

Article | [Published: 17 February 2021](#)

**Million-year-old DNA sheds light on the genomic history of mammoths**

# DNA as storage medium

## Advantages

- density
- durability
- **medium immutability**
- low TCO (Total Cost of Ownership)
- Energy efficiency and sustainability



There will always be DNA readers

In the future, DNA readers will be smaller, cheaper, faster

How to read a floppy disk today ?



# DNA as storage medium

## Advantages

- density
- durability
- medium immutability
- **low TCO (Total Cost of Ownership)**
- Energy efficiency and sustainability



TCO: the money put on the resources to preserve data for a certain amount of years

No need to periodically replace SSD, HDD, tapes, ...

Low electric consumption

- no electronic devices
- no air conditioning

Cheap data replication

# DNA as storage medium

## Advantages

- density
- durability
- medium immutability
- low TCO (Total Cost of Ownership)
- **Energy efficiency and sustainability**

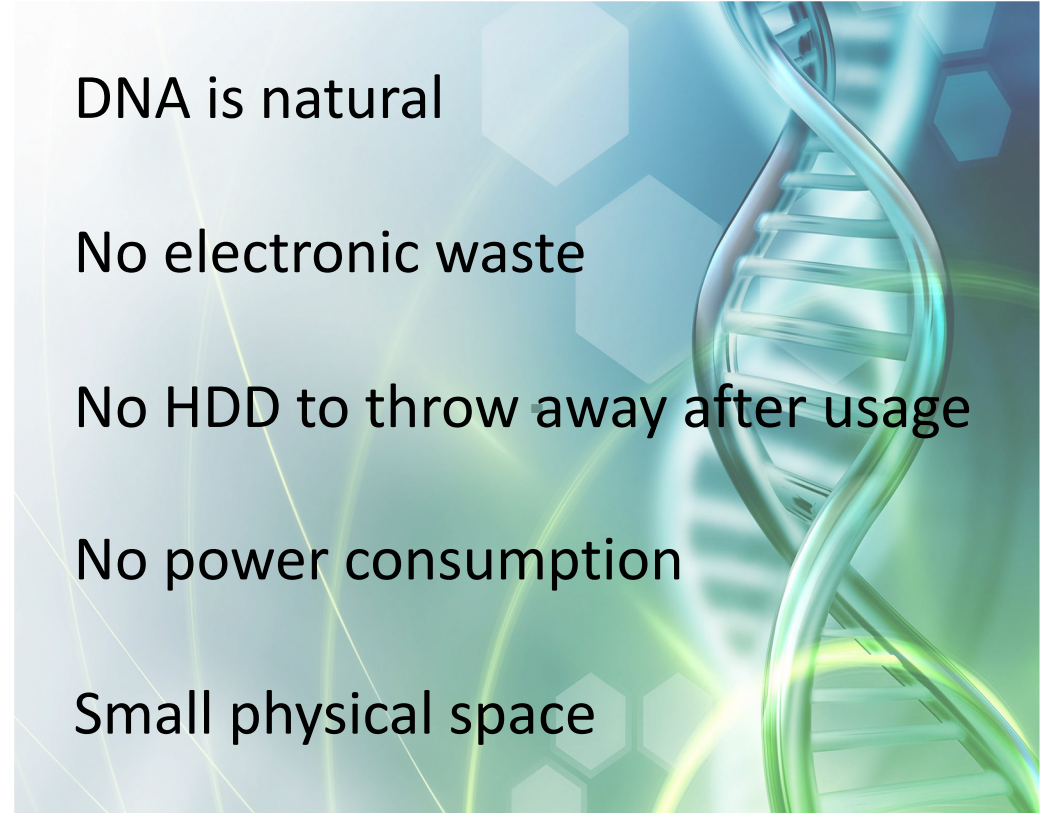
DNA is natural

No electronic waste

No HDD to throw away after usage

No power consumption

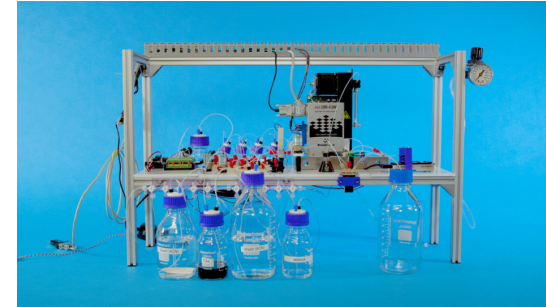
Small physical space



# Agenda

## 1. Introduction

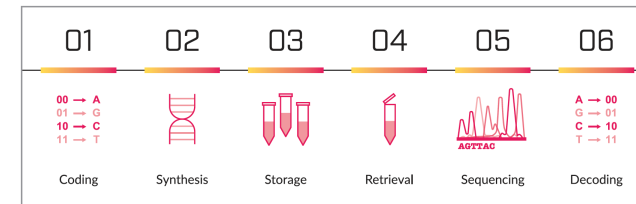
- Why new storage medium ?
- Why DNA ?



Source: <https://news.microsoft.com/innovation-stories/hello-data-dna-storage/>

## 2. DNA storage principle

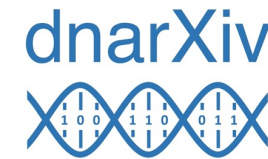
- How does it work ?
- What are the main challenges ?



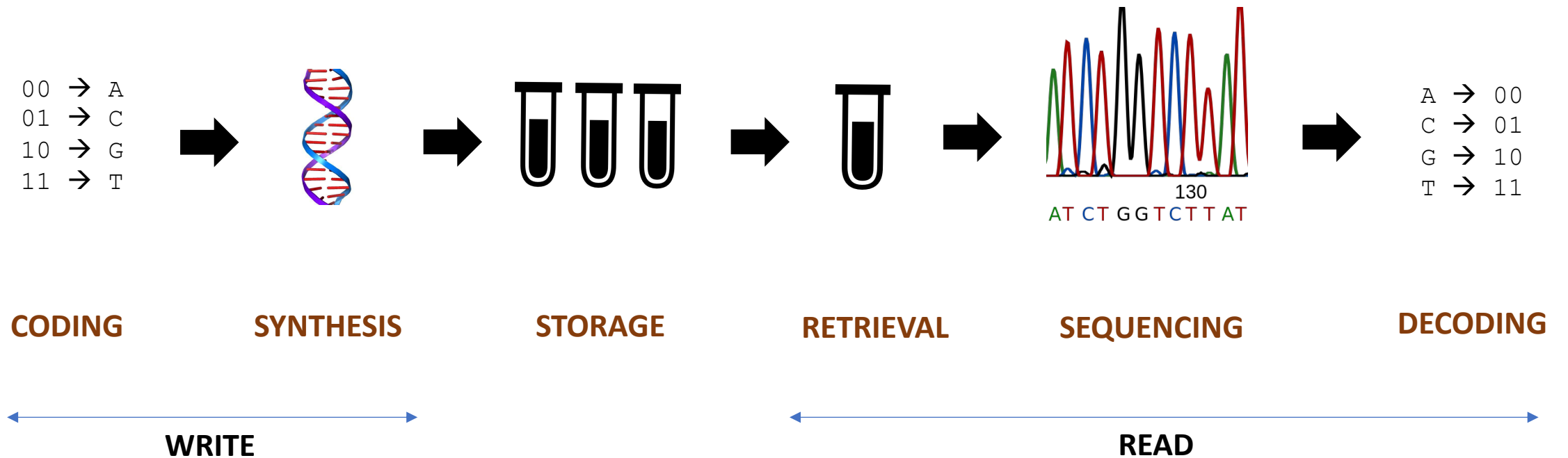
Source: *An Introduction to DNA Data Storage*, DNA storage alliance, June 21

## 3. dnarXiv project

- What's going on in Rennes ?

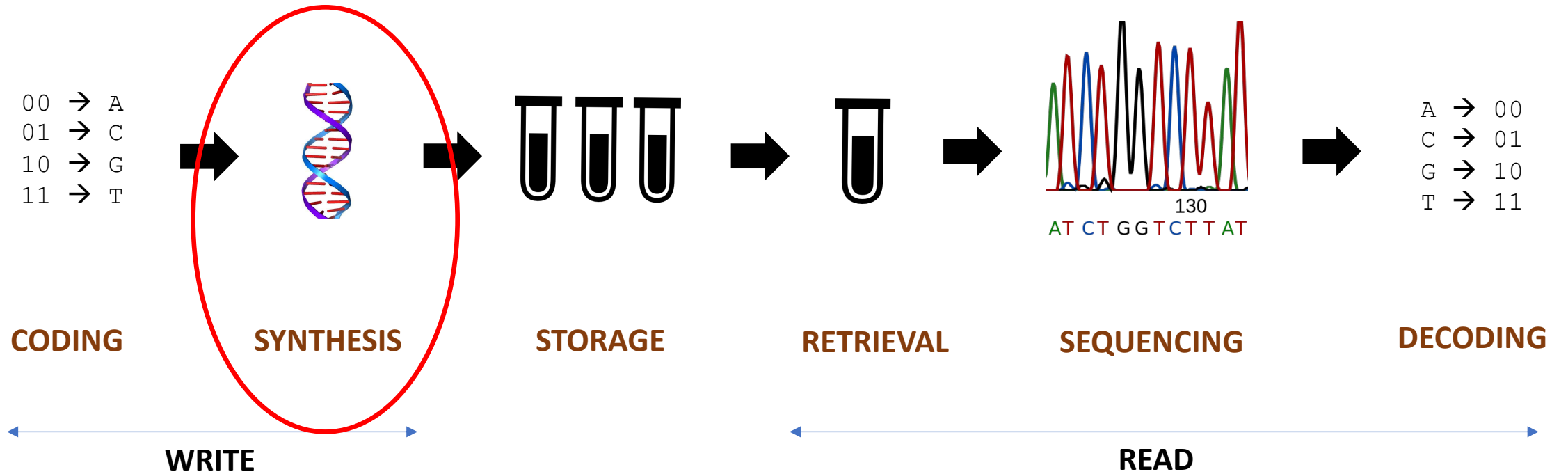


# Principle of DNA data storage





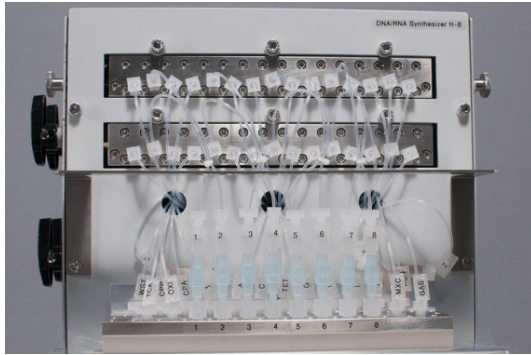
# Principle of DNA data storage



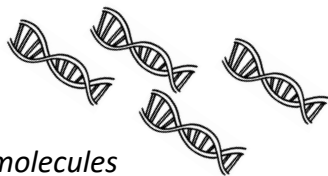
# Synthesis

**ATGGACCGGTGACACCCGGTTGGACCGTGA**

*1 text sequence*

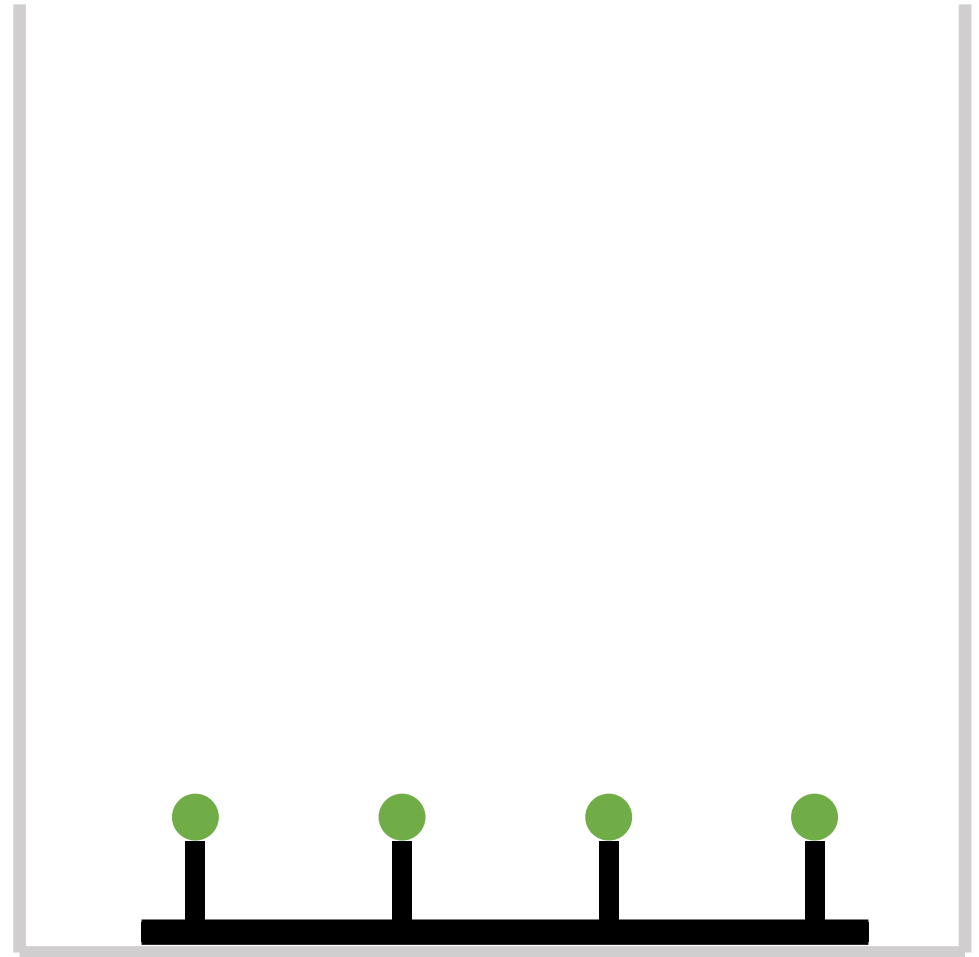


DNA  
Synthesizer



*Million/Billions of molecules*

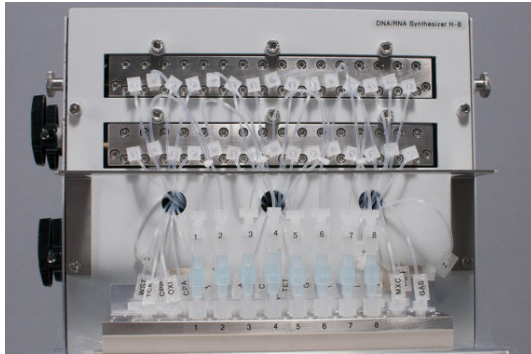
AGTCGTAC...



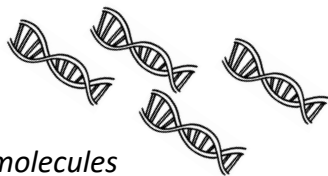
# Synthesis

**ATGGACCGGTGACACCCGGTTGGACCGTGA**

*1 text sequence*

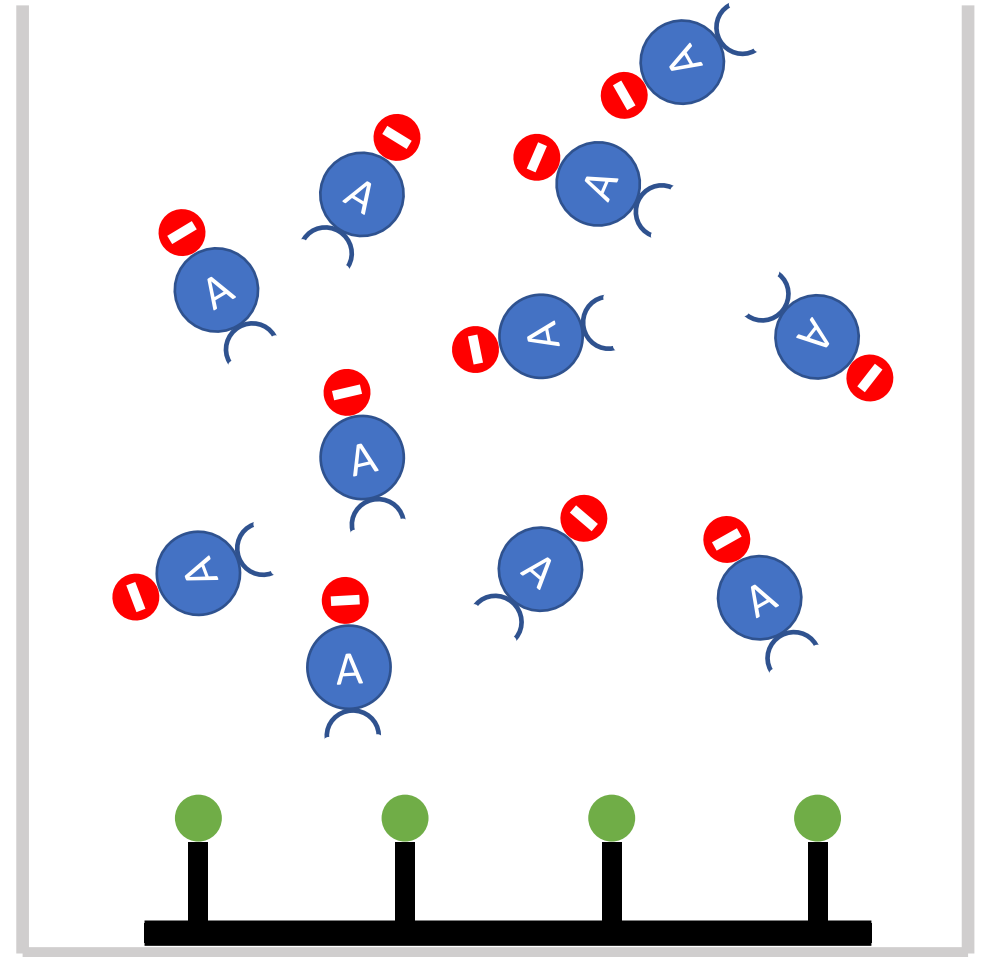


DNA  
Synthesizer



*Million/Billions of molecules*

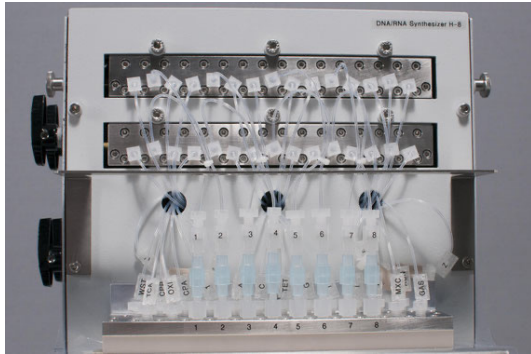
AGTCGTAC...



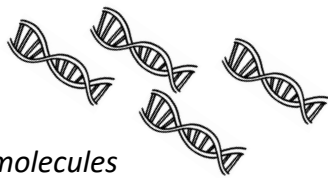
# Synthesis

**ATGGACCGGTGACACCCGGTTGGACCGTGA**

*1 text sequence*

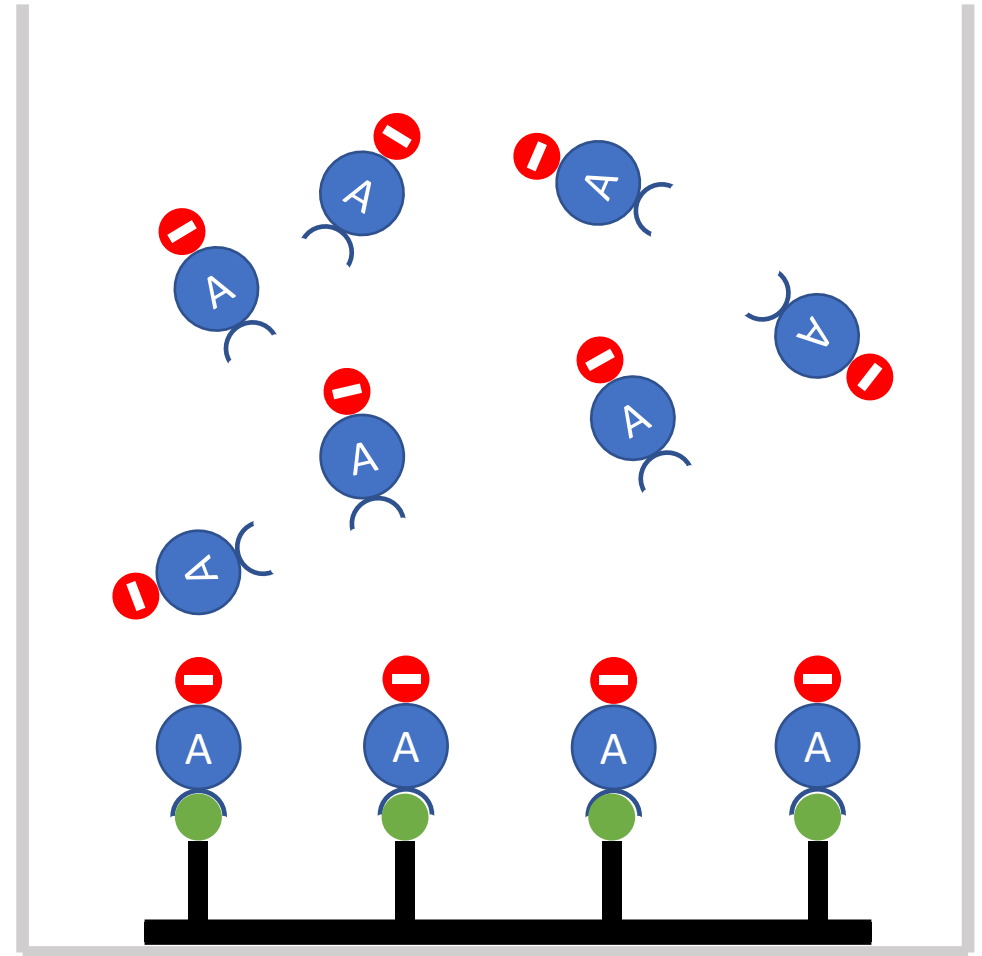


DNA  
Synthesizer



*Million/Billions of molecules*

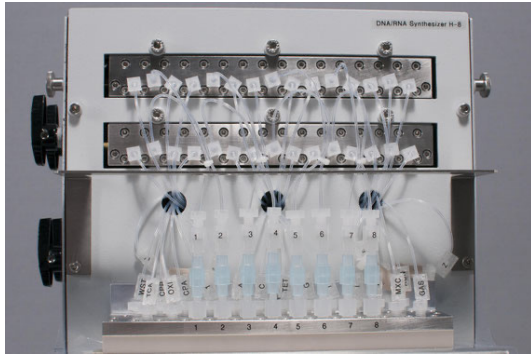
AGTCGTAC...



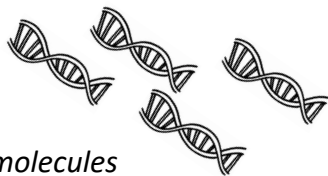
# Synthesis

**ATGGACCGGTGACACCCGGTTGGACCGTGA**

*1 text sequence*

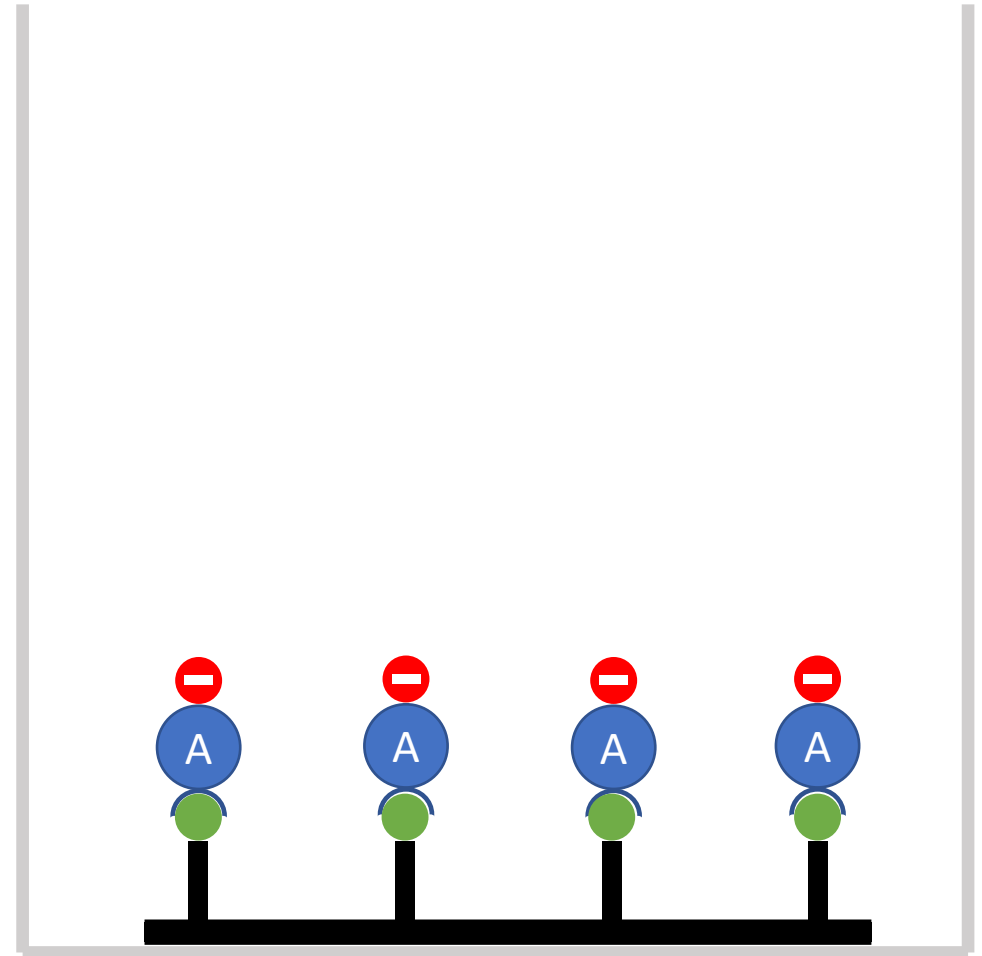


DNA  
Synthesizer



*Million/Billions of molecules*

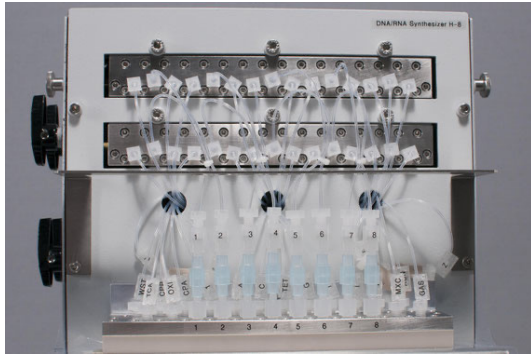
AGTCGTAC...



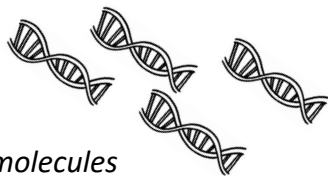
# Synthesis

**ATGGACCGGTGACACCCGGTTGGACCGTGA**

*1 text sequence*

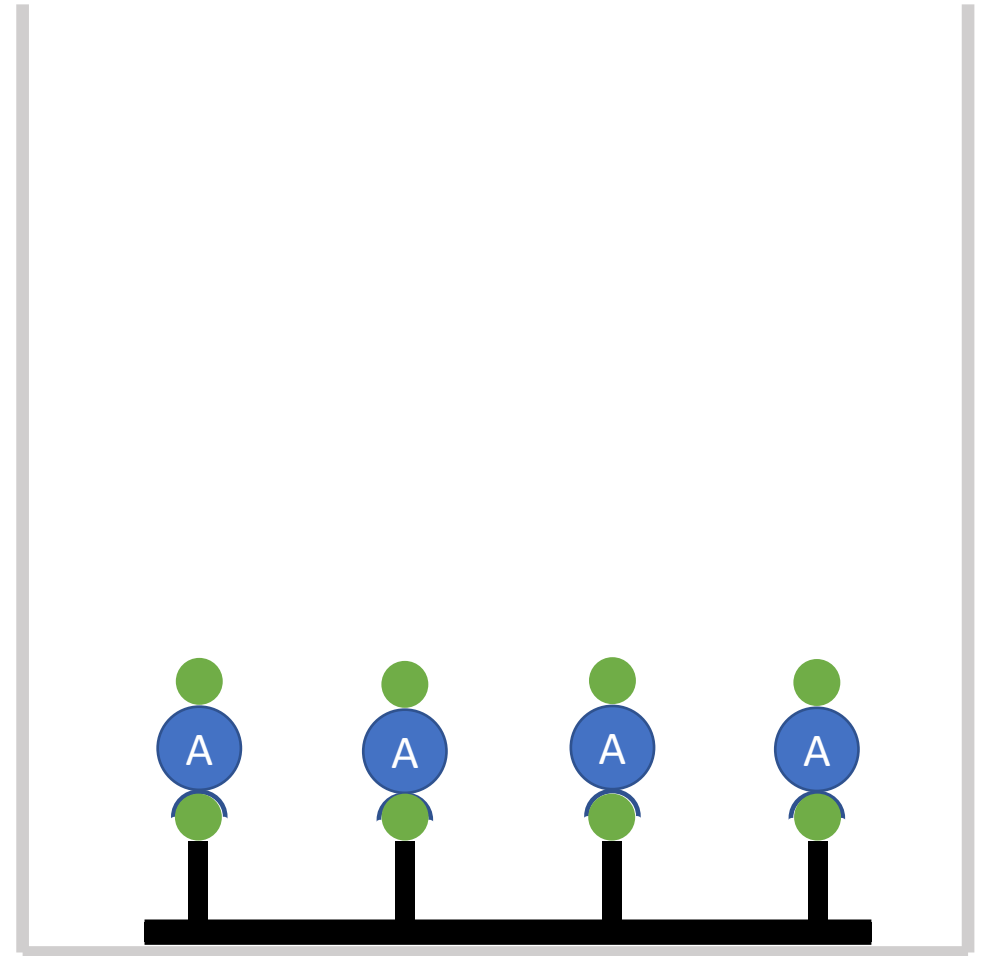


DNA  
Synthesizer



*Million/Billions of molecules*

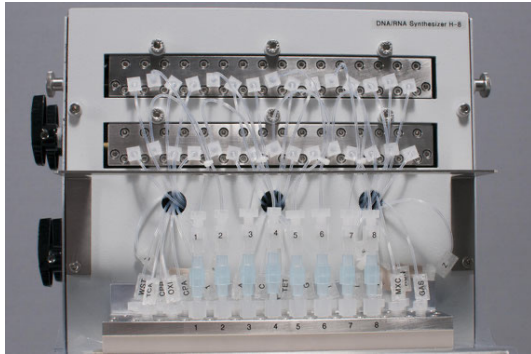
AGTCGTAC...



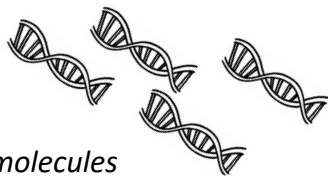
# Synthesis

**ATGGACCGGTGACACCCGGTTGGACCGTGA**

*1 text sequence*

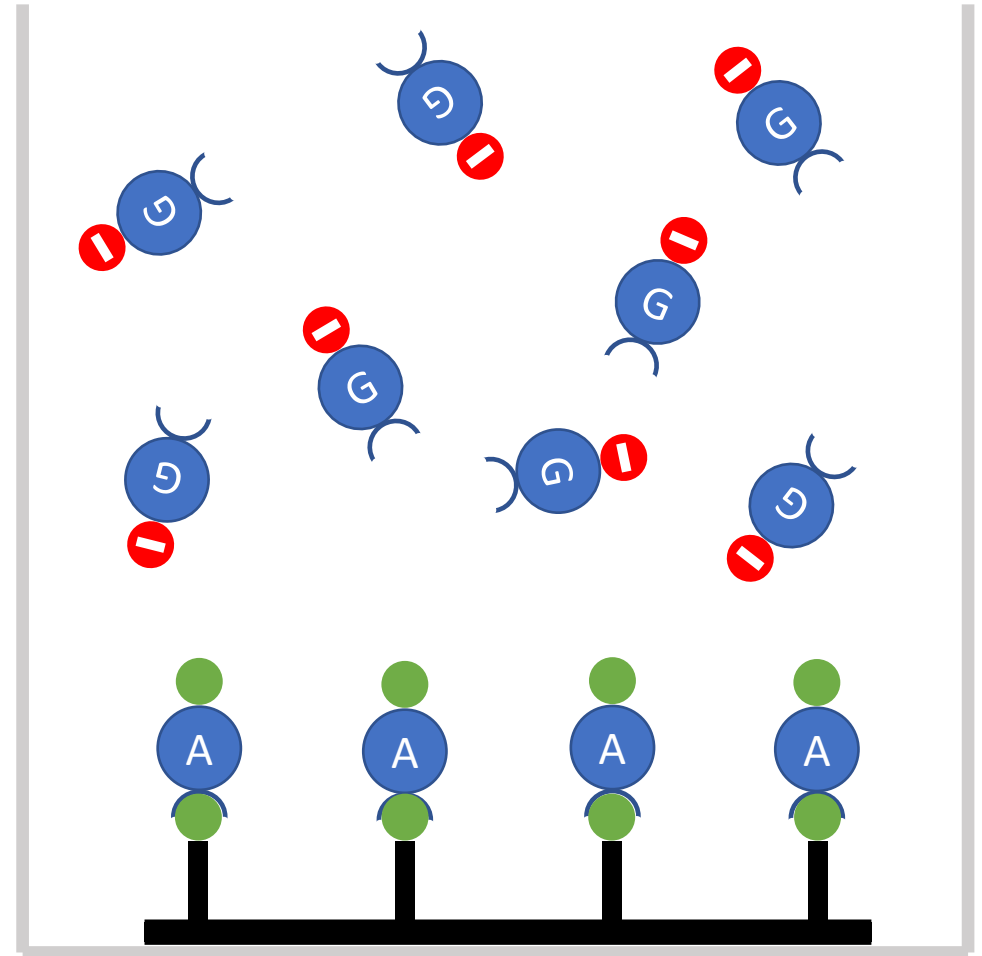


DNA  
Synthesizer



*Million/Billions of molecules*

AGTCGTAC...



# Synthesis limitation

Synthesizers can only generate short molecules (oligonucleotides)

- oligonucleotide length < 200 nucleotides

Very slow process

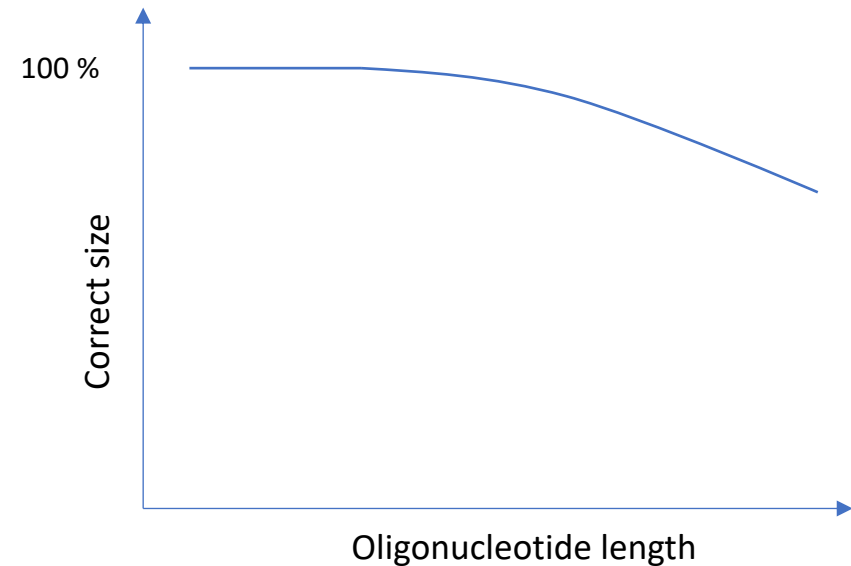
- 1 nucleotide / minute

High cost

- 100 € / Kbyte

Errors

- variation on oligonucleotide length





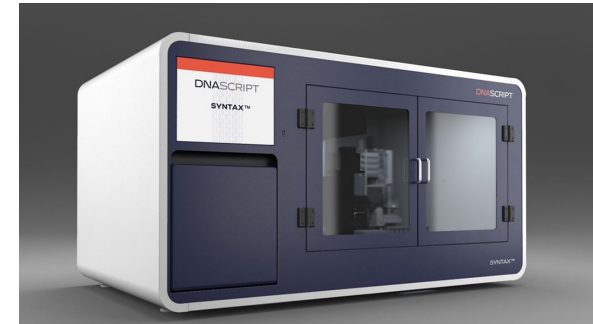
# Synthesis perspectives

## New technologies

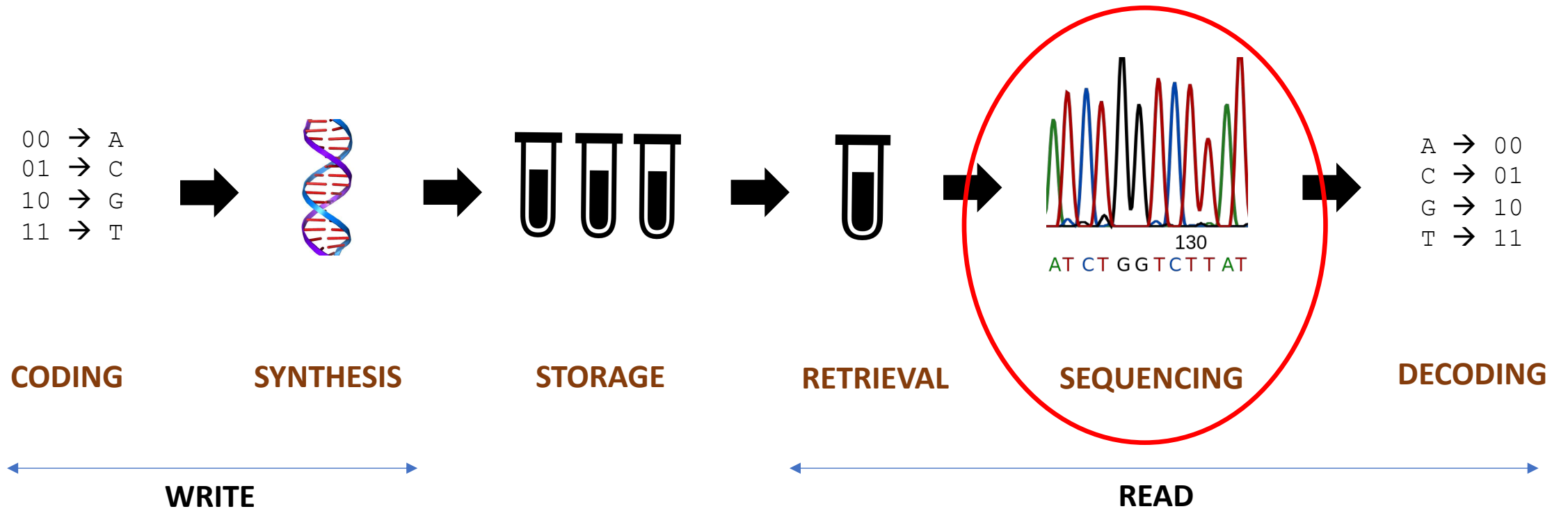
- Enzymatic DNA synthesis
  - exploit the capabilities of the TdT enzyme
  - longer oligonucleotides
  - faster

## Miniaturization & Parallelization

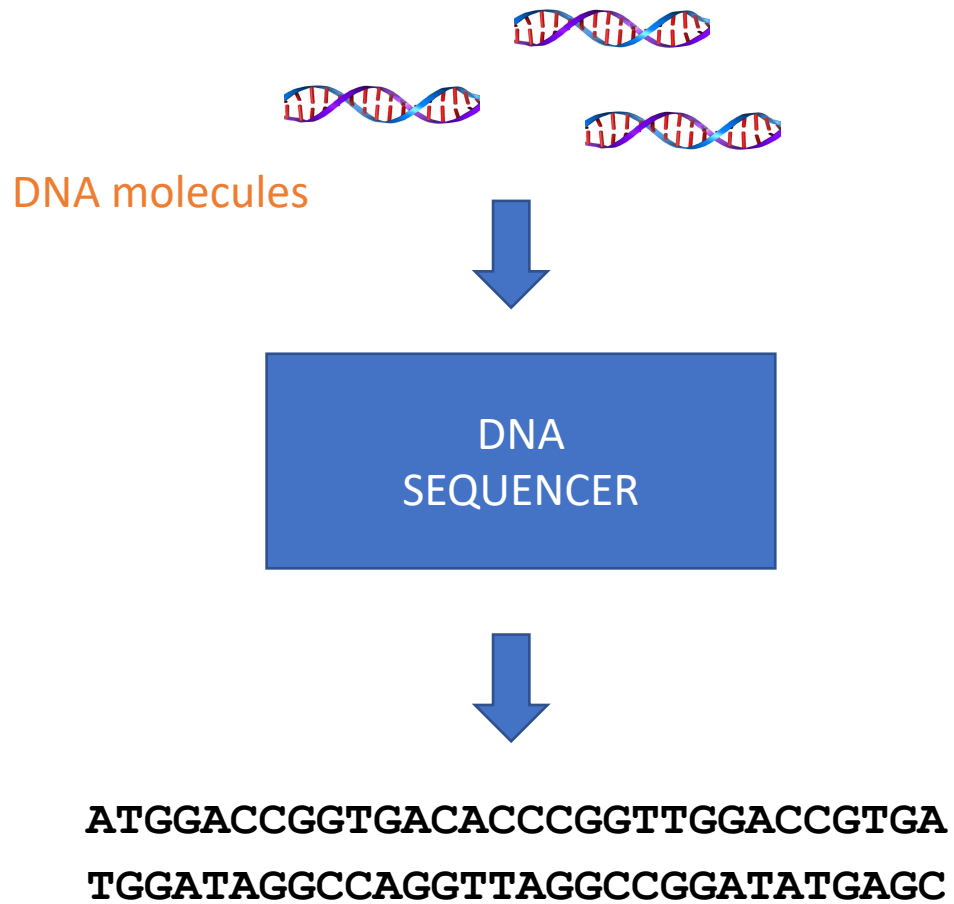
- Roadmap 2025
  - 1 Tbytes
  - 1 day
  - 1000 \$



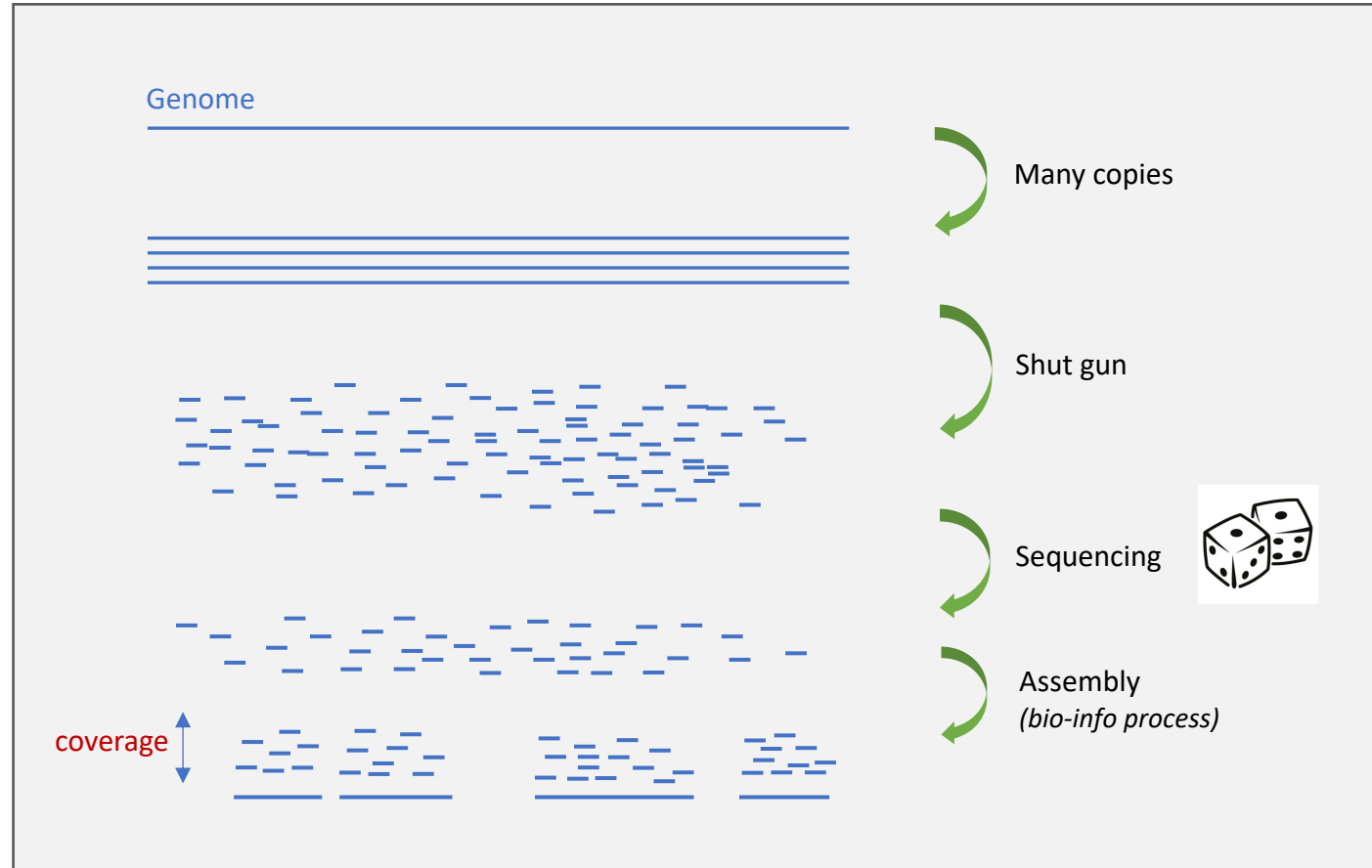
# Principle of DNA data storage



# Sequencing



Reading the text of a genome





# Sequencing cost & throughput

1 Human genome = 3.2 Gbp (~ 800 Mbytes)

## Cost

1 Human Genome < 1000 \$  
→ ~ 1 MBytes / \$

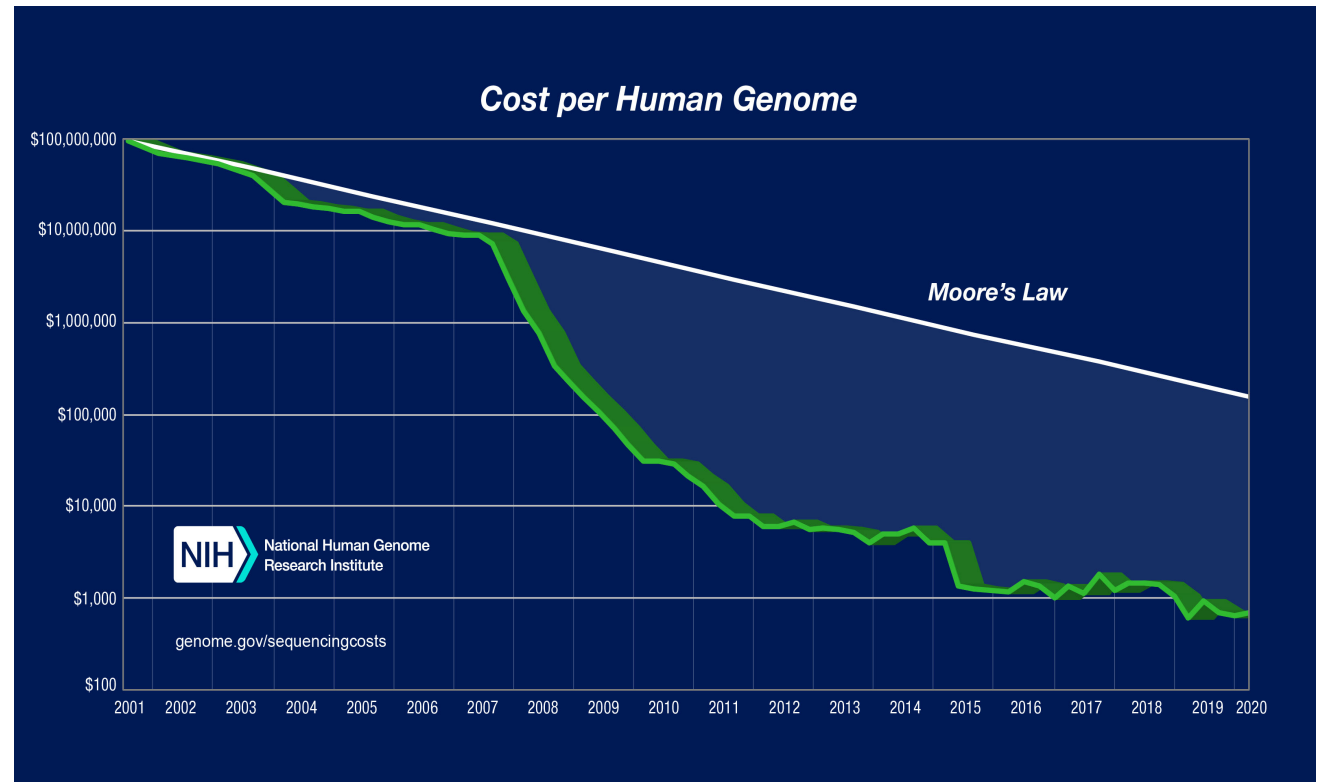
## Throughput

Illumina (NovaSeq 6000)

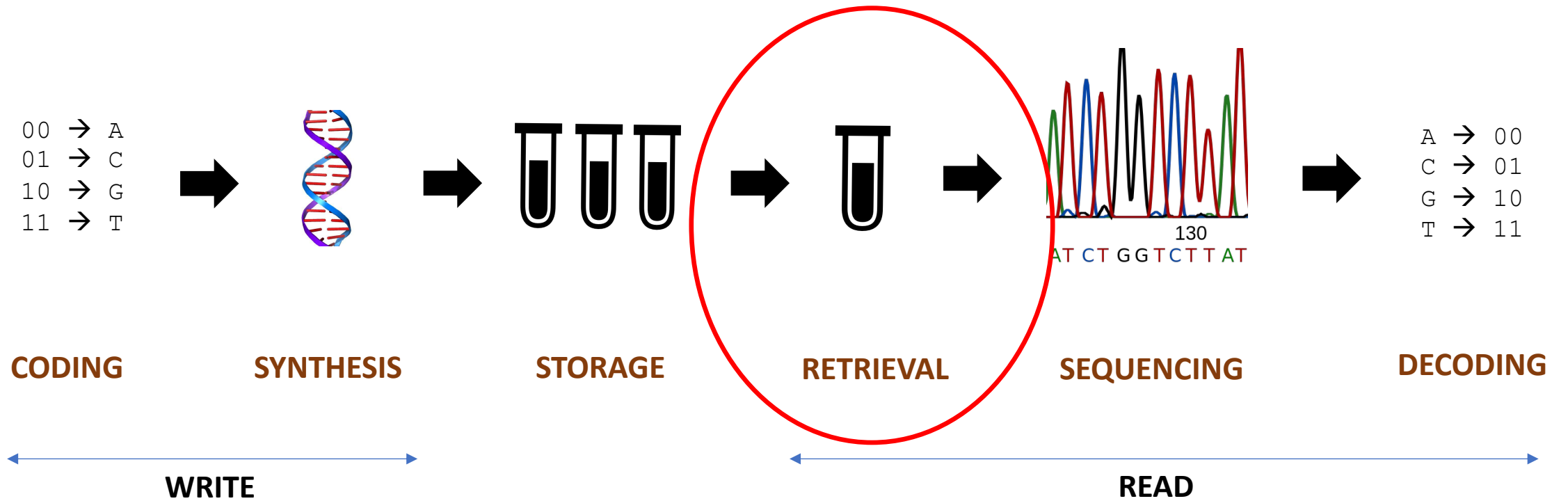
70 x 10<sup>3</sup> nucleotides / sec

ONT (Promethion, 48 flow cells)

20 x 10<sup>3</sup> nucleotides / sec



# Principle of DNA data storage



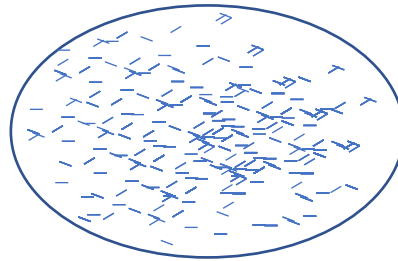
# Retrieval



Encoding  
+  
Synthesis



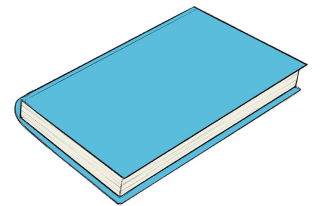
container



*Trillions of DNA molecules*



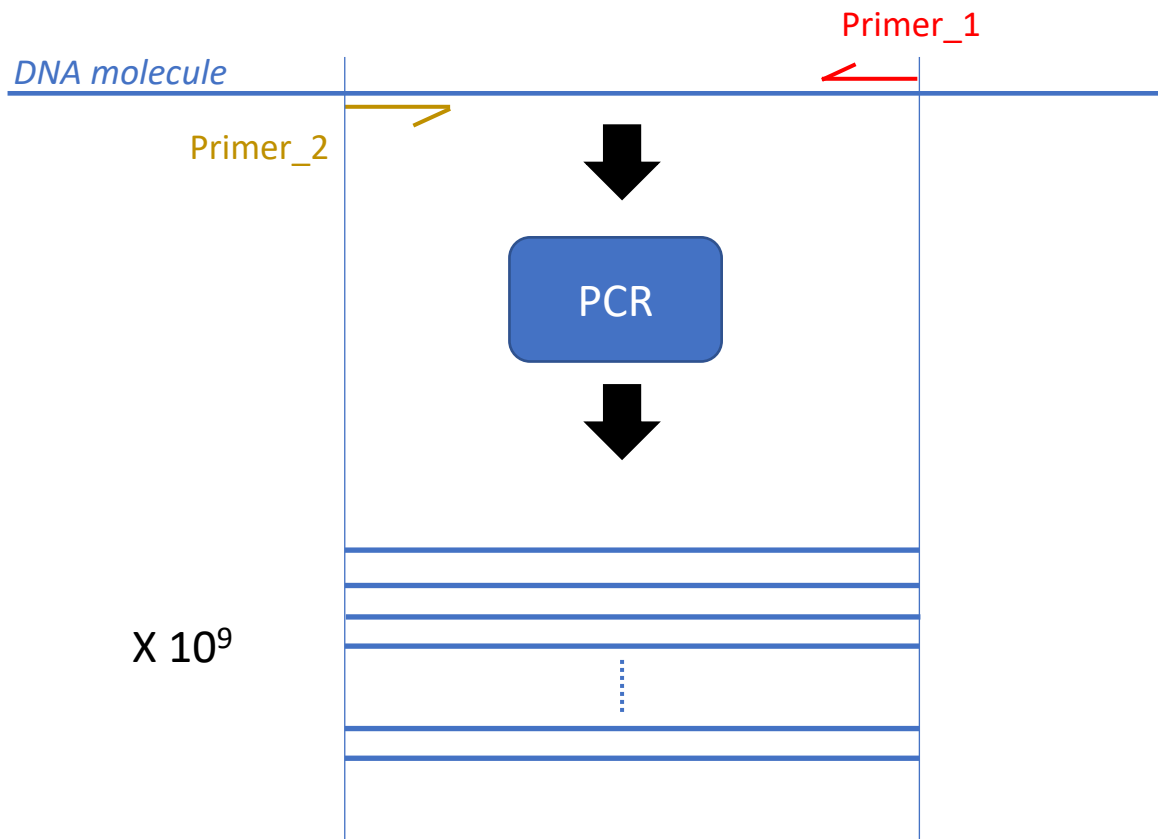
?



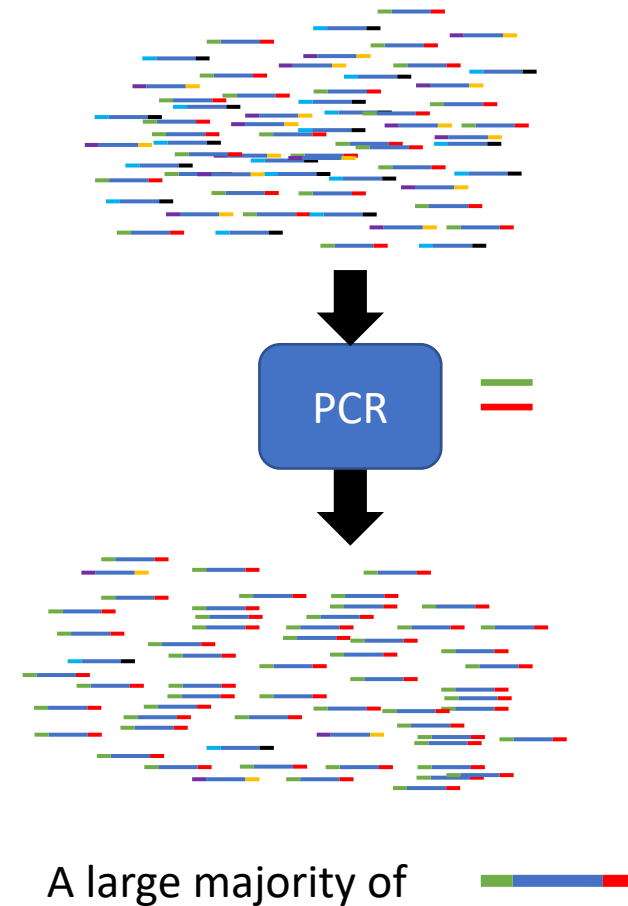
How avoiding to sequence  
all molecules to retrieve a  
specific document ?

# Selection by PCR

PCR: amplification of DNA molecules

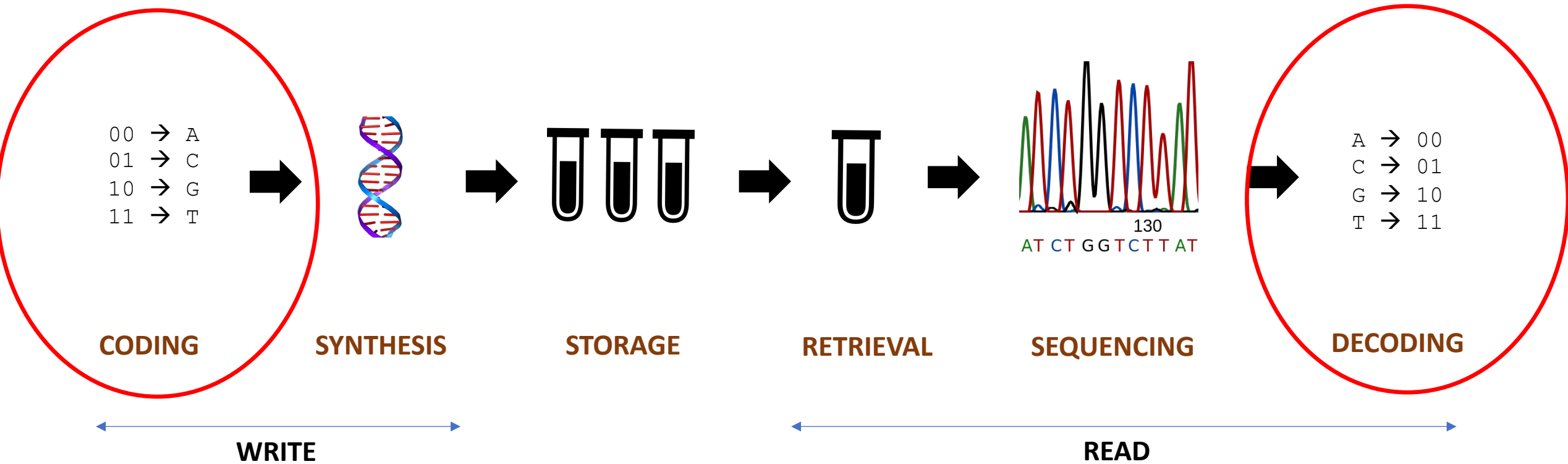


All DNA molecules from the same document will start and end by the same couple of primers

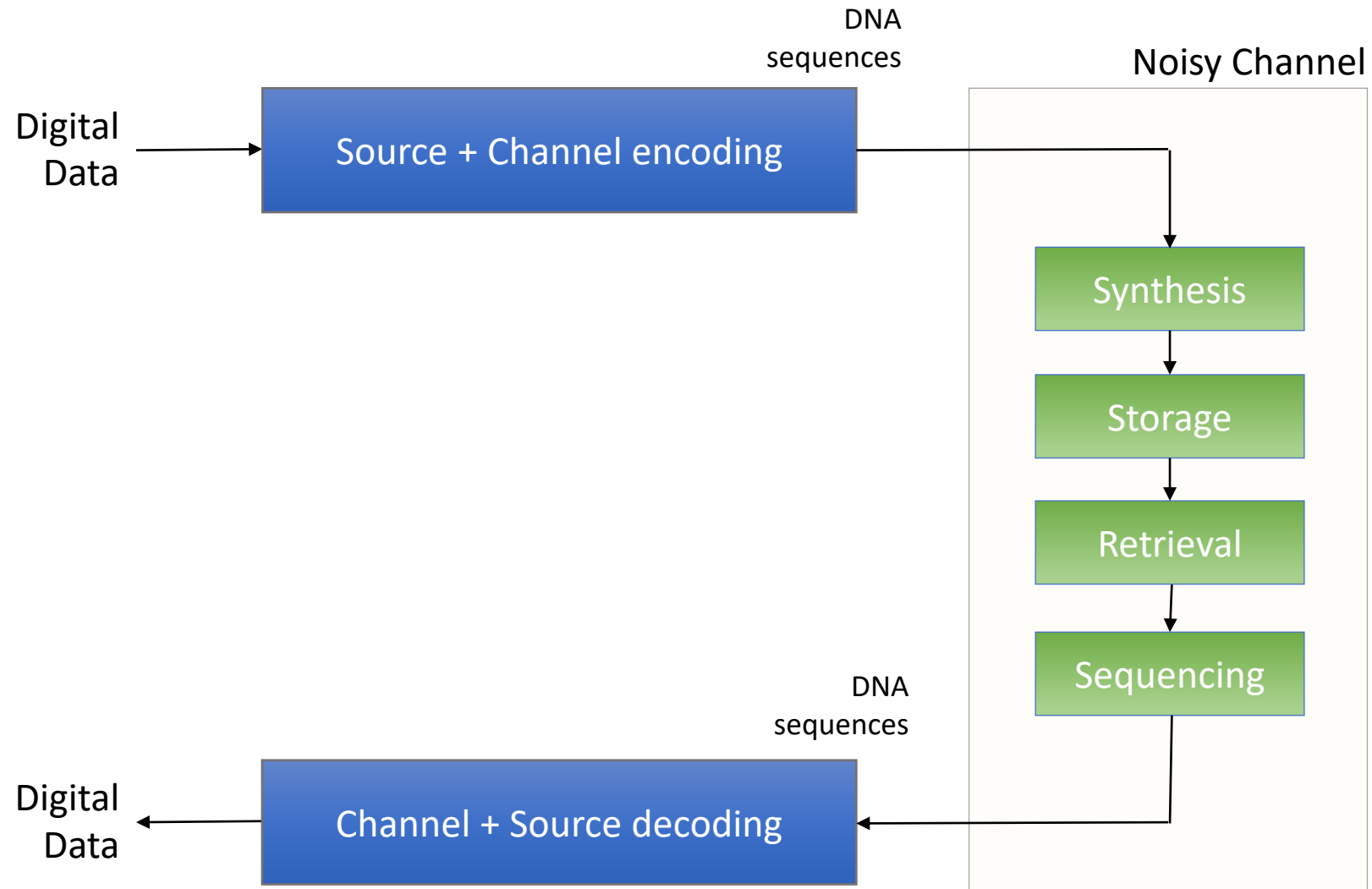




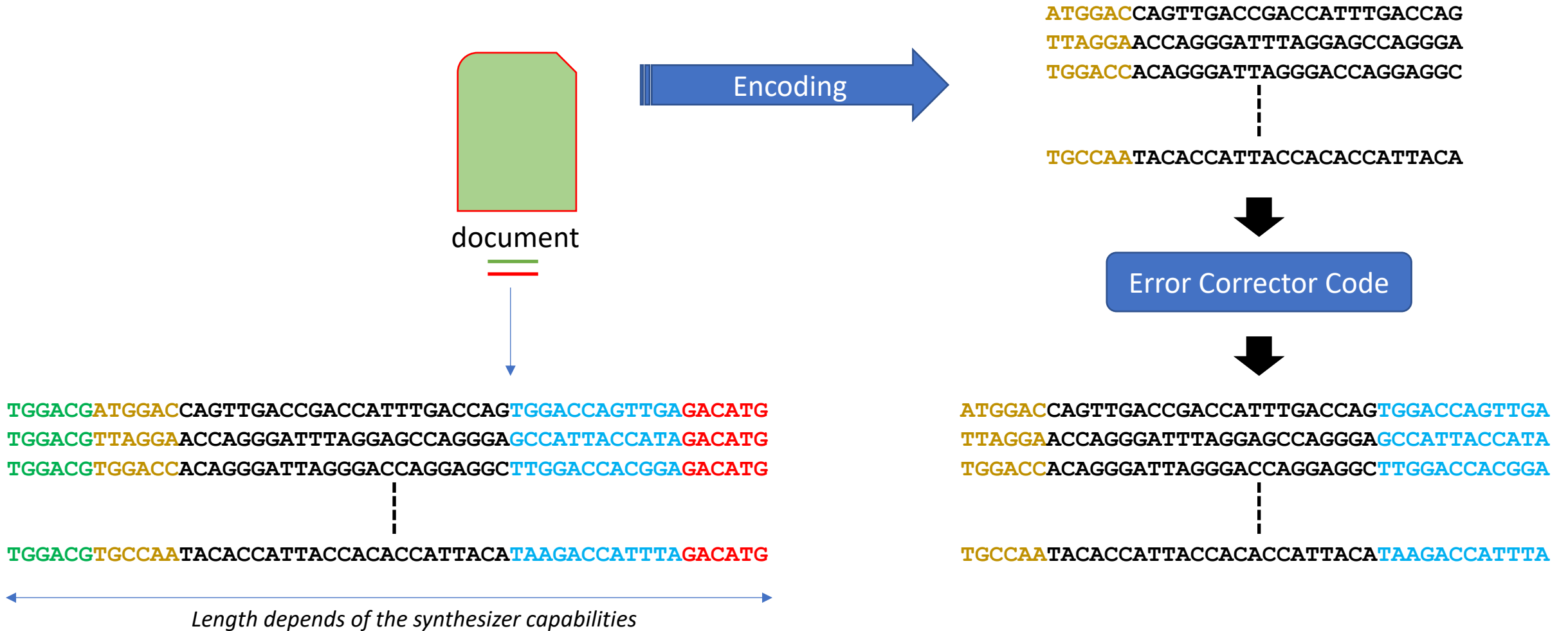
# Principle of DNA data storage



# Coding / Decoding



# Coding scheme



# Coding Challenges

Maximize the ratio: 
$$\frac{\text{Useful information}}{\text{Indexing} + \text{CCR}}$$

TACGAGACCAGTTGACATTTGACCAGTCAGTTGAGATG

## Constraints on sequences

- no (long) homopolymers: ATTAGACTTTTTCGAGTA
- no specific motifs
- AT / GC ratio between 0.45 & 0.55

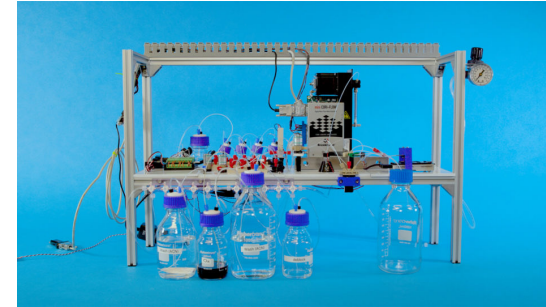
## Error Corrector Codes

- have to consider the following errors: insertion, deletion, substitution

# Agenda

## 1. Introduction

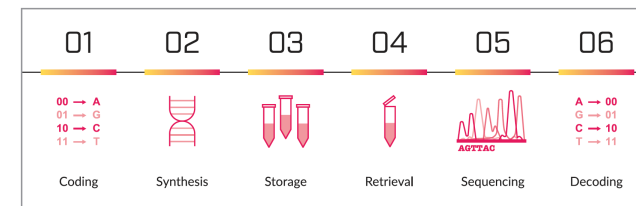
- Why new storage medium ?
- Why DNA ?
- State-of-the-art



Source: <https://news.microsoft.com/innovation-stories/hello-data-dna-storage/>

## 2. DNA storage principle

- How does it work ?
- What are the main challenges ?



Source: *An Introduction to DNA Data Storage*, DNA storage alliance, June 21

## 3. **dnarXiv** project

- What's going on in Rennes ?



# dnarXiv project



Starting date: October 2020

Objective: Investigate information archiving on DNA

Multidisciplinary project

- Signal processing / computer science
- Bio-informatics
- Bio-technology: synthesis & sequencing



# dnarXiv research axes

## Coding/decoding

Design codes dedicated to DNA storage medium

## Bio-technology

Synthesis: store long DNA molecules

Sequencing: 3<sup>rd</sup> generation sequencing technology

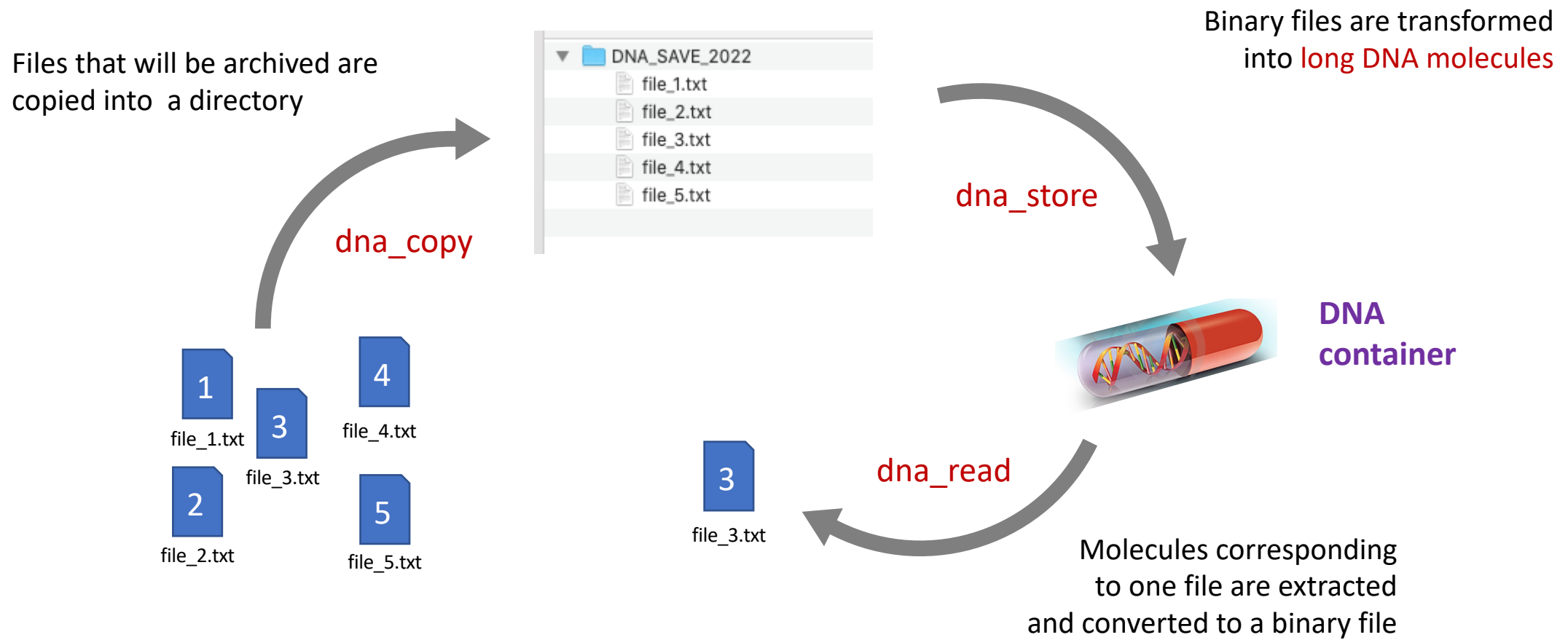
## Security

Investigate the possibilities to mixed digital technologies and bio-tech properties

## Platform

- To integrate codes & Software developed in dnarXiv
- To connect IT and bio-tech parts
- To conduct experiments
  - wet-lab
  - in-silico

# dnarXiv approach





# dnarXiv bio-technology

Objective: optimize the ratio:  $\frac{\text{Useful information}}{\text{Indexing}}$

AGTTGACACCAGTTGACATTTGACCAG



Design long molecules from short oligonucleotides

AGATGCCAGTTGACTTGGGACCAGGATTTAGGGGACACCAGGGGATTGTAGGACCCACGGATTGAATTTGACCAGTTGGACCAGGTTCCATTA

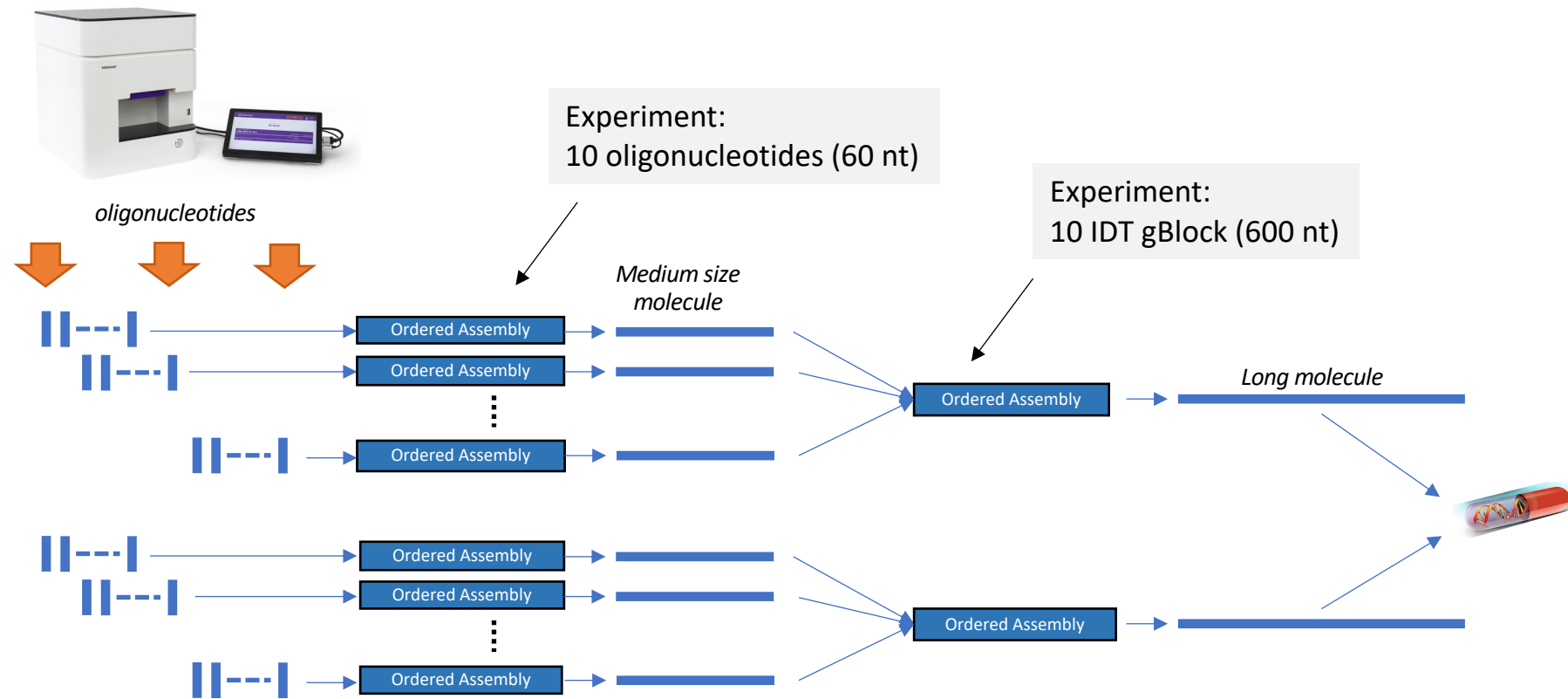


Use 3<sup>rd</sup> generation sequencing technology

- long reads
- Oxford nanopore technology



# Oligos → long molecules



# Ordered assembly

AGGACCAGGGATTTAGGCCAGATATGAGGACCGGATTAGGCCCGGGTATATATGATCCATGGGAC...

AGGACCAGGGATTTAG **GCCAGATATGAGGACCGGAT** **TAGGCCCGGGTATATATGAT** **CCATGGGAC**...

AGGACCAGGGATTTAG

**GCCAGATATGAGGACCGGAT**

**TAGGCCCGGGTATATATGAT**

|||||

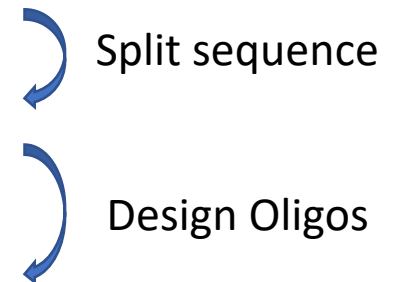
|||||

|||||

TCCTGGTCCCTAAAT**CGGT**

CTATACTCCTGGCCTA**ATCC**

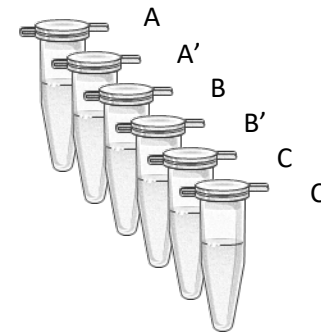
GGCCCATATATACTA**GGTA**



A : AGGACCAGGGATTTAG  
A' : **GCCAGATATGAGGACCGGAT**  
B : **TAGGCCCGGGTATATATGAT**  
B' : TCCTGGTCCCTAAAT**CGGT**  
C : CTATACTCCTGGCCTA**ATCC**  
C' : GGGCCATATATACTA**GGTA**  
...

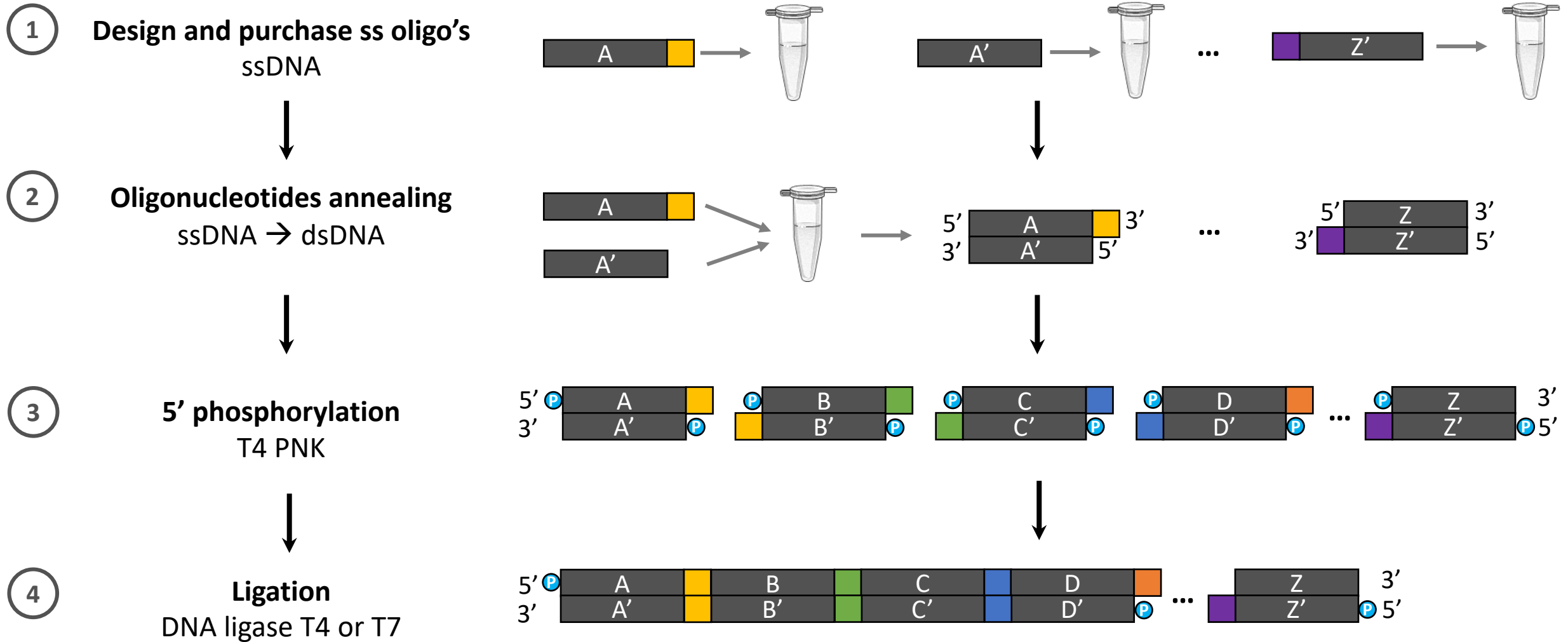


Oligonucleotides  
Synthesis

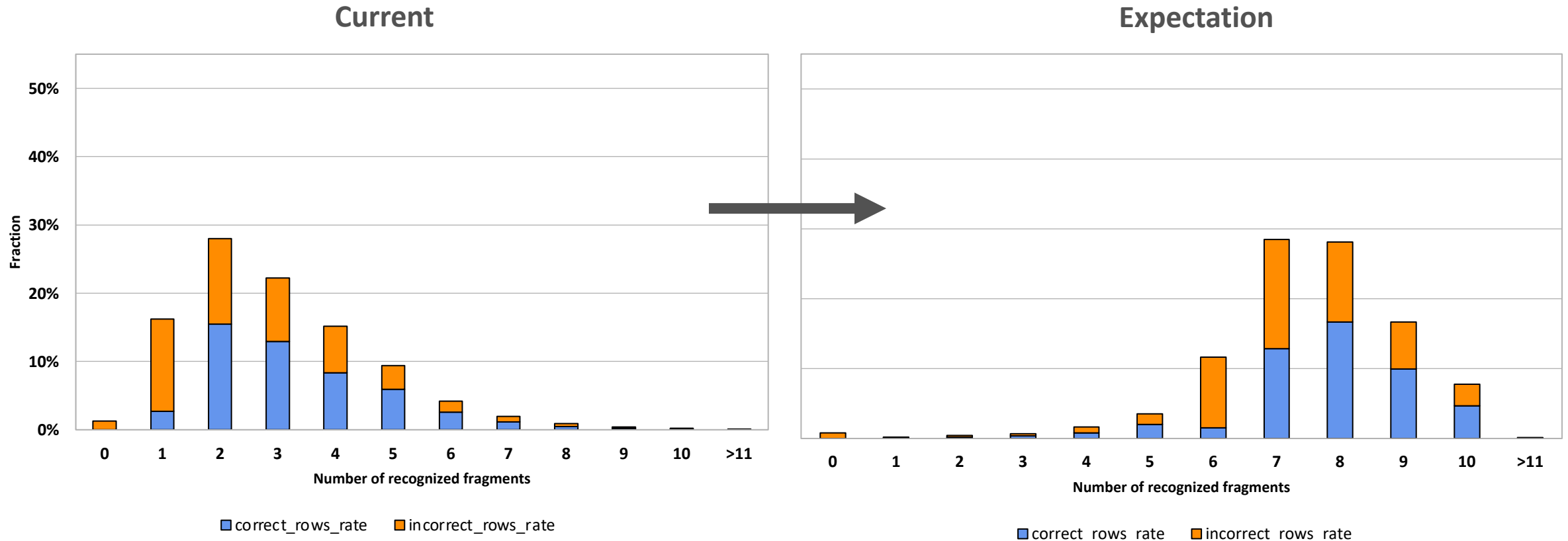


List of oligos

# Synthetic DNA assemblies workflow

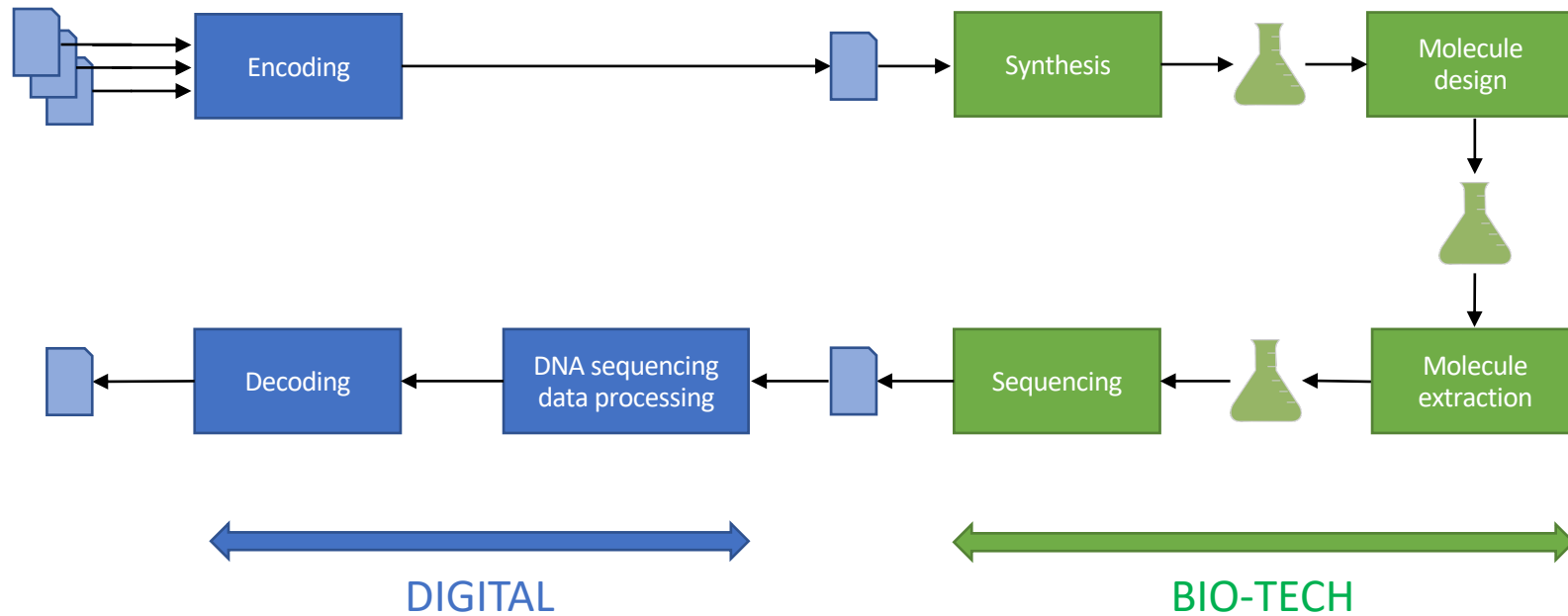


# Ordered assembly experiments

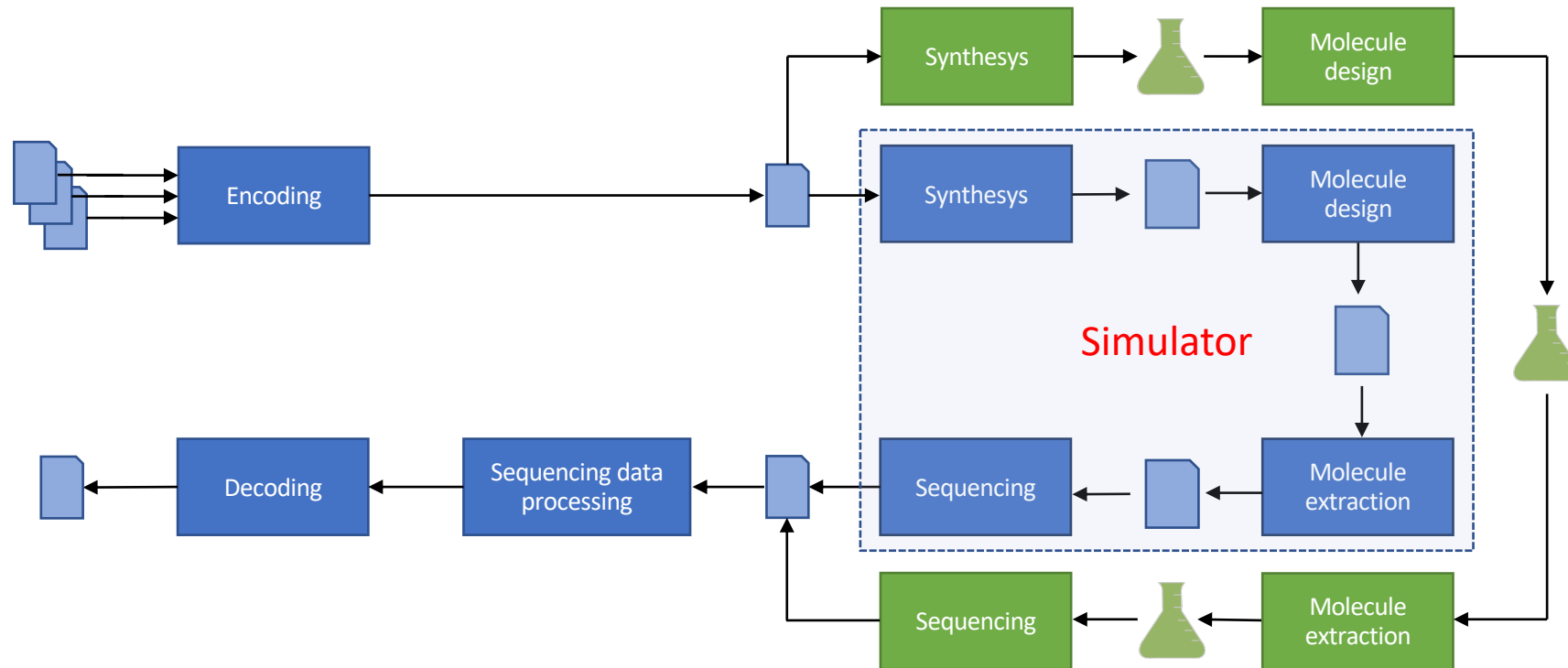


**Work on the limiting steps of the process:** inhibitor reaction products, ATP degradation, enzymes stability ...  
**Next:** try to assemble more fragments and increase correct row rates.

# dnarXiv platform



# in-silico alternative



# Take home message

## Data explosion

- Need for new storage medium

## DNA molecules

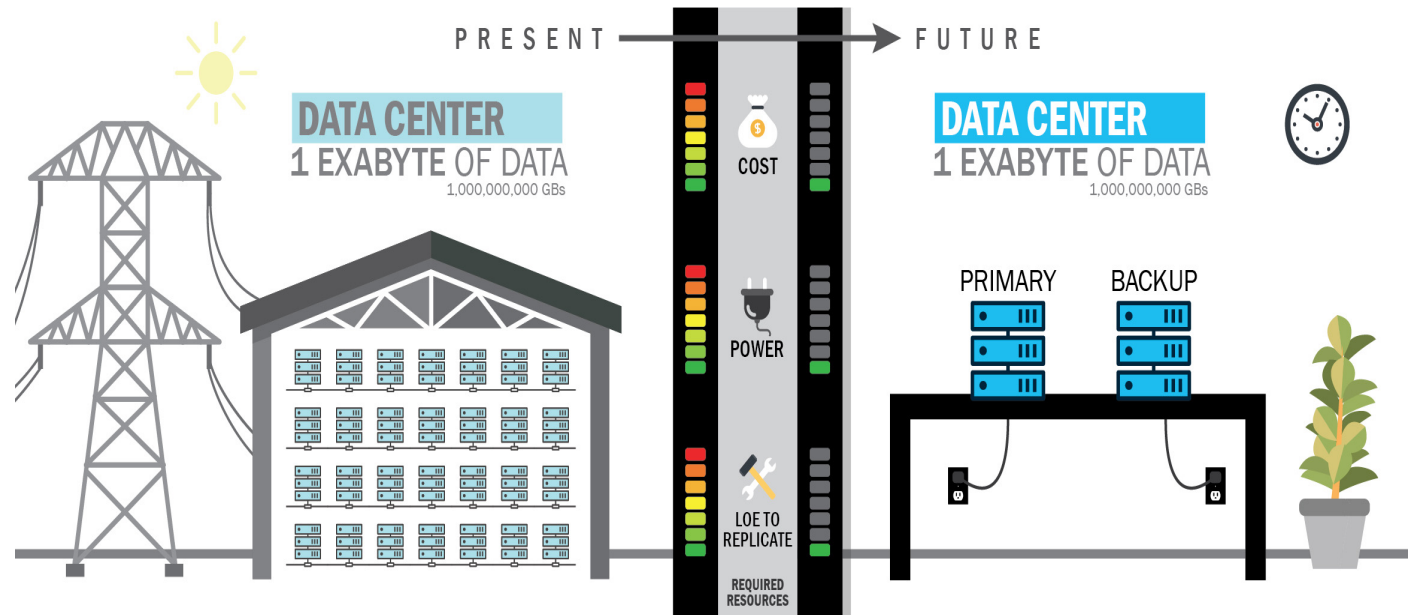
- density / durability
- cold data (archiving)

## Today bottleneck = Synthesis

- slow & costly
- keys: miniaturization, parallelization, automation

## New active research field

- international initiatives



<https://www.iarpa.gov/research-programs/mist>



# International initiatives

## MIST Molecular Information Storage (IARPA)

- 2020 – 2023 (4 years)
- Roadmap 2024: write 1TB/read 10TB, 1 day, 1000 \$



<https://www.iarpa.gov/research-programs/mist>

## DNA Data Storage Alliance

- formed in October 2020
- founders: Illumina, Microsoft, Twist Bioscience and Western Digital
- Today : 50 members



<https://dnastoragealliance.org/>

## EIC Pathfinder Challenge: DNA-based digital data storage

- Deadline: oct 2022



## PEPR MolecularArXiv, France

- start year: 2022



**IDA  
2022**

**Symposium on Intelligent Data Analysis 2022**  
**April 20–22, 2022, Rennes, France**

**Thank you for your attention**

*dominique.lavenier@irisa.fr*



# State of the art

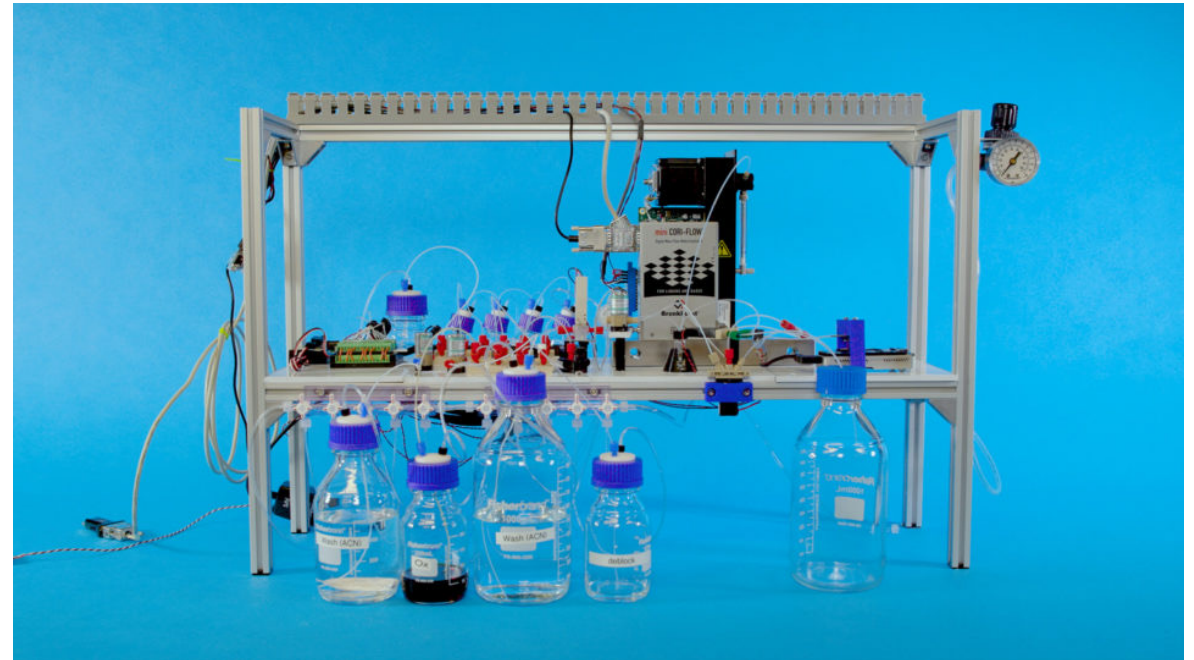
MICROSOFT & University of Washington  
March 2019

First fully automated system to store and retrieve data in manufactured DNA

Write & read automatically the word “**Hello**”

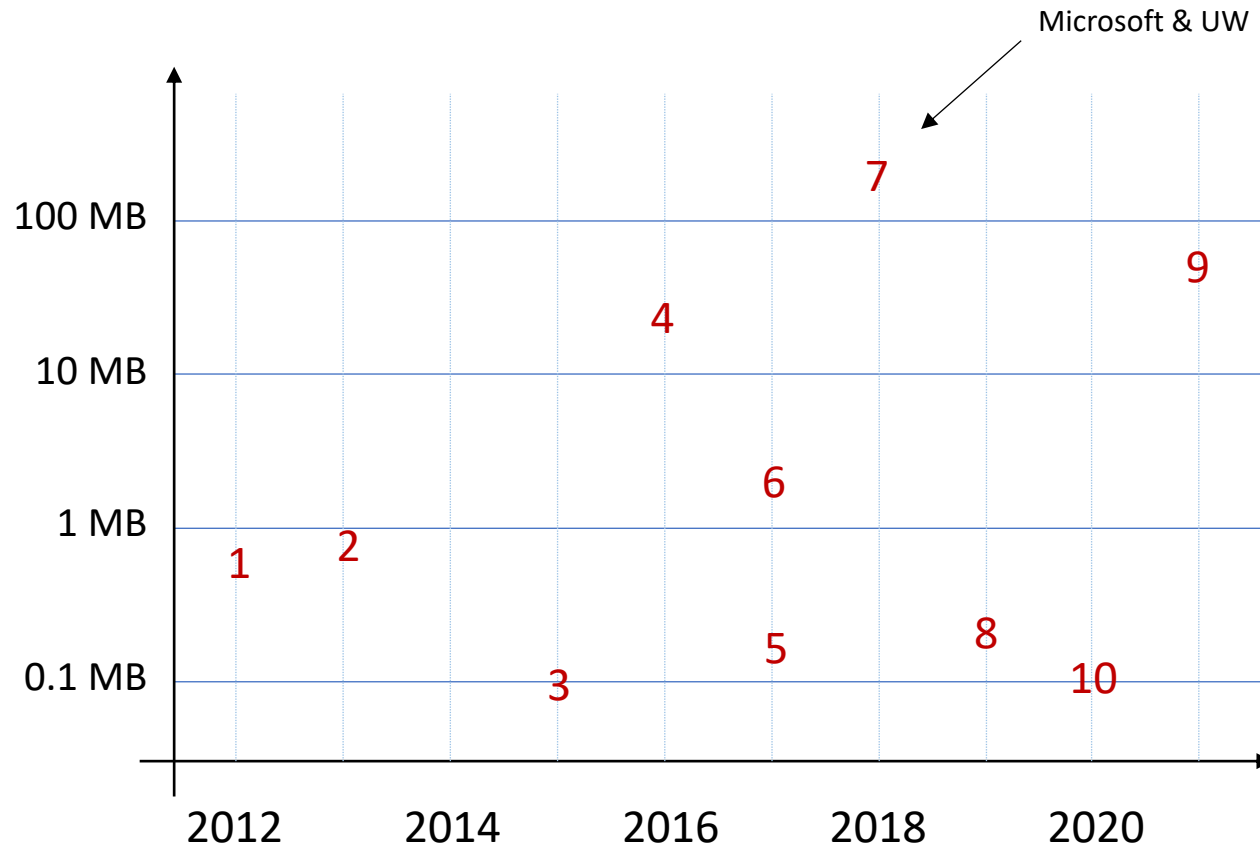
Process time : 21 hours

Takahashi, C.N., Nguyen, B.H., Strauss, K. *et al.* Demonstration of End-to-End Automation of DNA Data Storage. *Sci Rep* 9, 4998 (2019)



Source: <https://news.microsoft.com/innovation-stories/hello-data-dna-storage/>

# DNA Storage Experiments



	length of seqs.	number of seqs	data stored (in MB)	error correction	
1	Church, Gao, and Kosuri [27]	115	54,898	0.65	None
2	Goldman, Bertone, Chen, <i>et al.</i> [38]	117	153,335	0.75	Repetition
3	Grass, Heckel, Puddu, <i>et al.</i> [39]	117	4,991	0.08	RS
4	Blawat, Gaedke, Hütter, <i>et al.</i> [12]	190	900,000	22	RS
5	Bornholt, Lopez, Carmean, <i>et al.</i> [14]	120	45,652	0.15	RS
6	Erlich and Zielinski [33]	152	72,000	2.14	Fountain
7	Organick, Ang, Chen, <i>et al.</i> [80]	150	$13.4 \cdot 10^9$	200.2	RS
8	Chandak, Tatwawadi, Lau, <i>et al.</i> [17]	150	13,716	0.192	LDPC
9	Heckel and Grass [46]	105	$3.88 \cdot 10^9$	63.1	RS
10	Antkowiak, Lietard, Darestani, <i>et al.</i> [7]	60	16,383	0.1	RS

I. Shomorony, R. Heckel (2022), "Information-Theoretic Foundations of DNA Data Storage", *Foundations and Trends® in Communications and Information Theory*: Vol. 19: No. 1, pp 1-106  
<http://dx.doi.org/10.1561/0100000117>