



**HAL**  
open science

# Secure data storage into DNA molecules compliant with biological constraints: Ensuring the confidentiality of data stored into DNA molecules

Chloé Berton, Gouenou Coatrieux, Dominique Lavenier

## ► To cite this version:

Chloé Berton, Gouenou Coatrieux, Dominique Lavenier. Secure data storage into DNA molecules compliant with biological constraints: Ensuring the confidentiality of data stored into DNA molecules. RITS 2022 - Recherche en Imagerie et Technologies pour la Santé, May 2022, Brest, France. pp.1-1. hal-03817968

**HAL Id: hal-03817968**

**<https://hal.science/hal-03817968>**

Submitted on 17 Oct 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

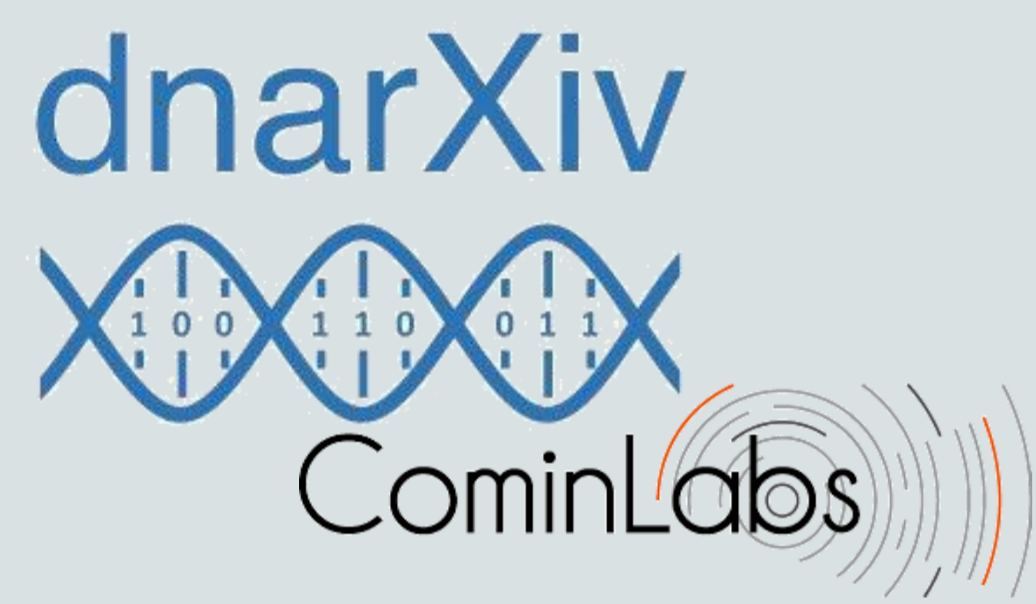
# Secure data storage into DNA molecules compliant with biological constraints

## Ensuring the confidentiality of data stored into DNA molecules

### Auteurs

Chloé Berton  
Gouenou Coatrieux  
Dominique Lavenier

### Projet

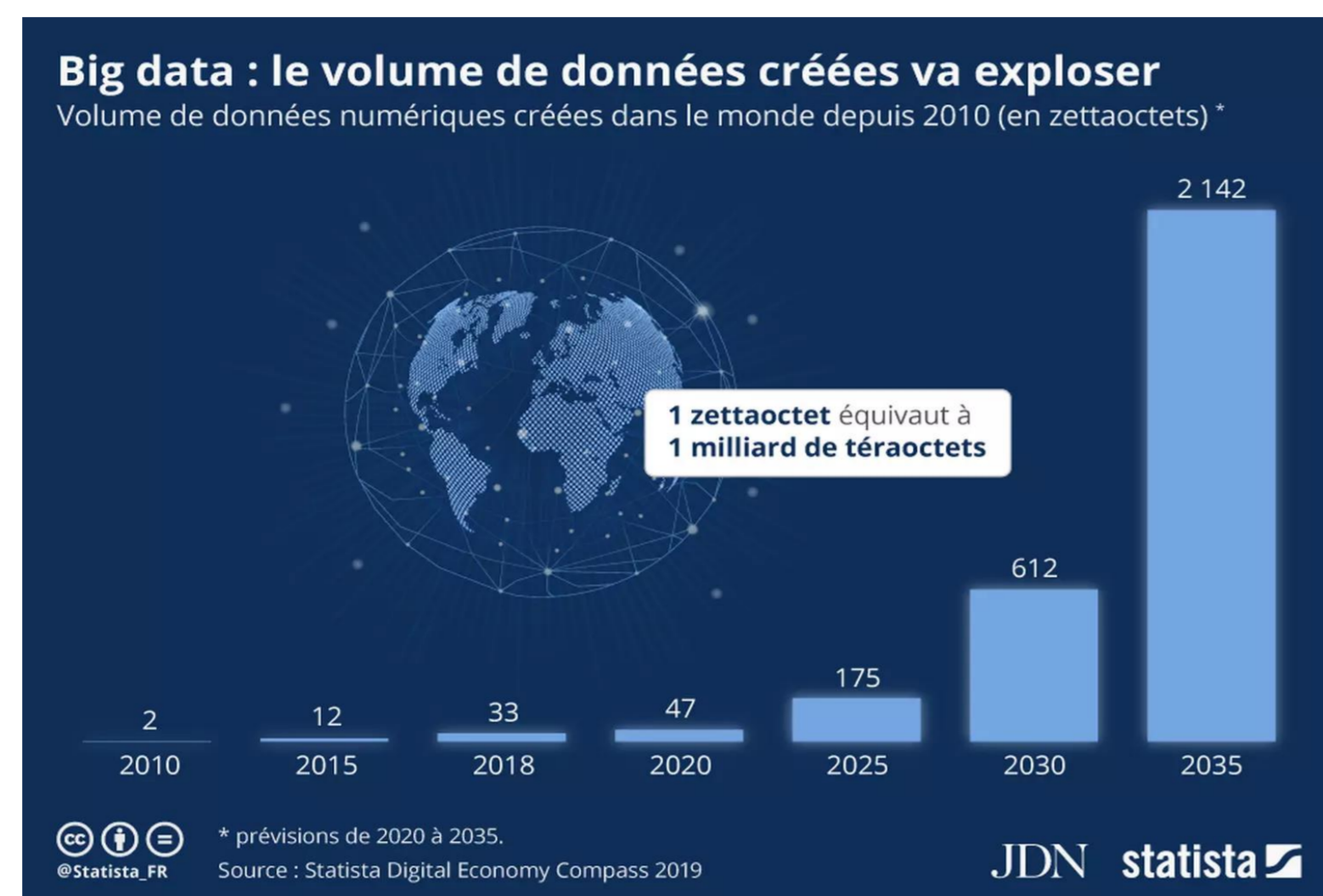


### References

[1] Rydning, D. R. J. G. J. (2018). The digitization of the world from edge to core. Framingham: International Data Corporation, 16.

[2] De Silva, P. Y., & Ganegoda, G. U. (2016). New trends of digital data storage in DNA. BioMed research international, 2016.

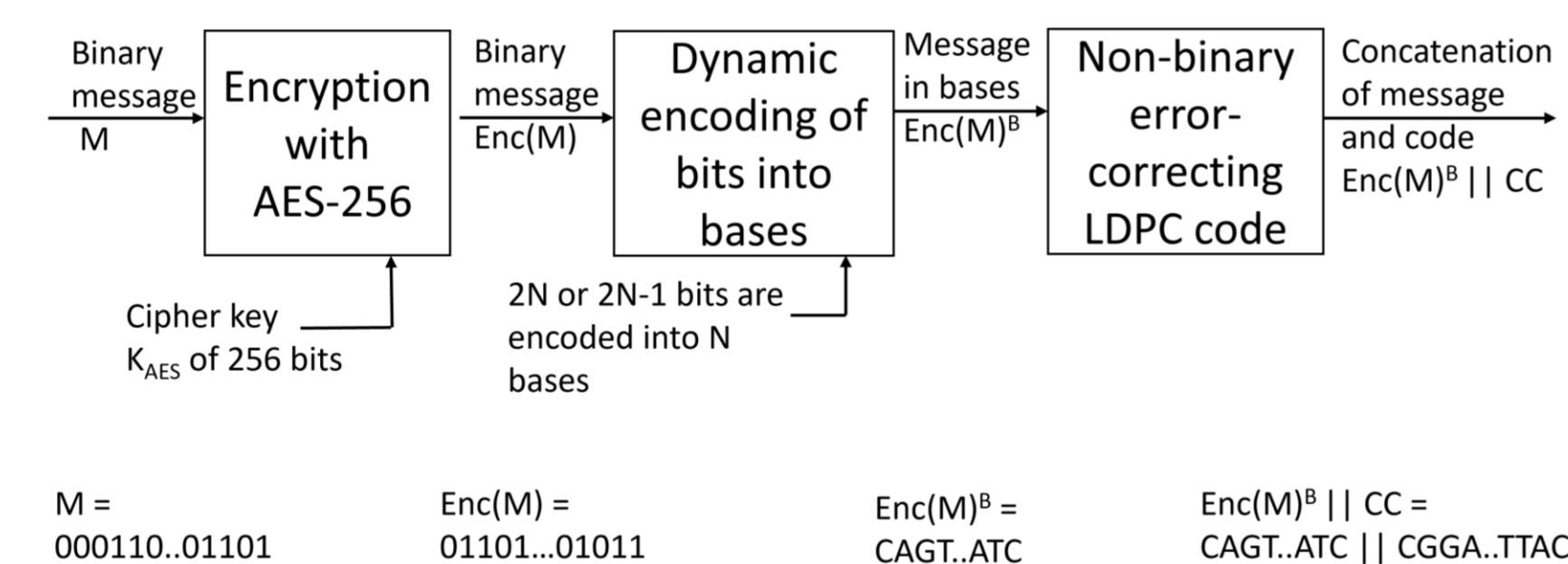
[3] Hamoum, B., Dupraz, E., Conde-Canencia, L., & Lavenier, D. (2021, August). Channel Model with Memory for DNA Data Storage with Nanopore Sequencing. In 2021 11th International Symposium on Topics in Coding (ISTC) (pp. 1-5). IEEE.



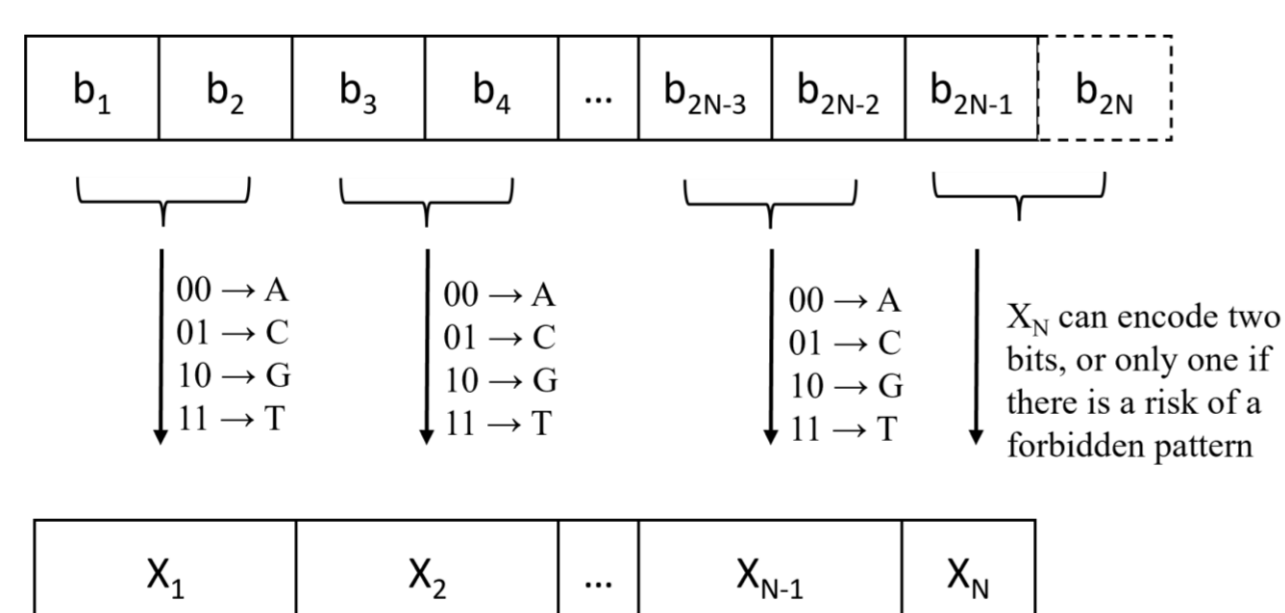
Volume of data created or replicated in the world, projection for 2020 to 2035 - Source : Statista

## DNA data storage - A new storage medium

- **Principle** – [WRITING] Encode binary data into 4-base sequences following the DNA structure, transfer this data into synthetic DNA molecules. [READING] Amplify encoded sequences of interest and get several reads of them with a sequencing device. Reads are then processed and decoded back to binary data
- **Constraints** – i) Biological DNA synthesis and sequencing are imperfect and introduce errors. ii) devices have structural DNA requirements when generating 4-base sequences
- **Vulnerabilities** – This chain is notably vulnerable to: theft or cloning of molecules; spying attacks on the sequencing or synthesis devices; DDoS attack by adding fake DNA sequence to confuse sequencing



Encoding solution to ensure data confidentiality in the entire DNA data storage channel



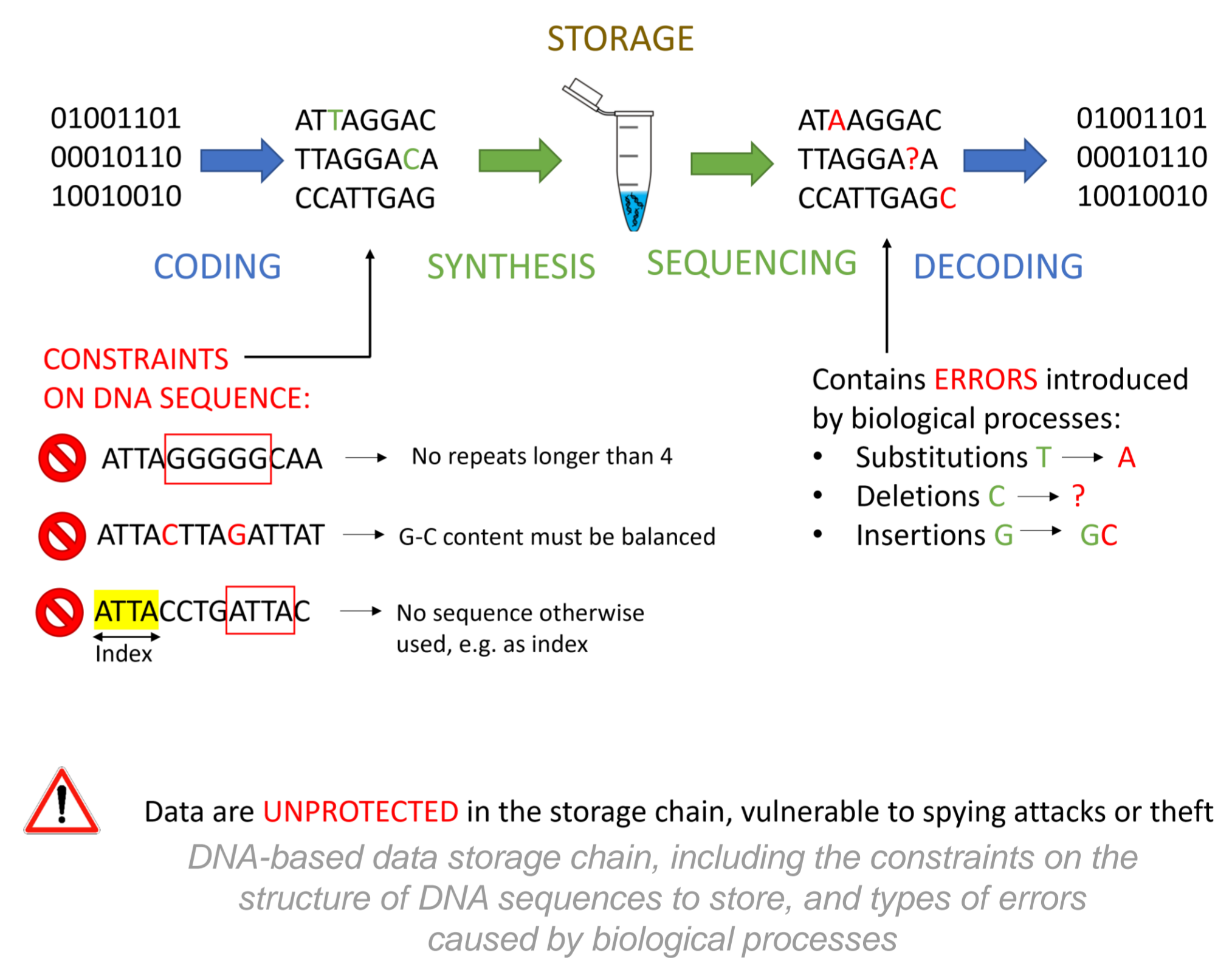
Dynamic information encoding into bases. In a block of N bases, each base  $X_i$  encodes two bits:  $b_{2i-1}$  and  $b_{2i}$ , except  $X_N$ , which encodes either  $b_{2N-1}$  only, or  $b_{2N-1}$  and  $b_{2N}$

## Experimental results

- **Simulation** of the biological processes using the simulator from [3]
- **Results:** no homopolymers longer than  $N$ , G-C content of 43-57%, data recovery without errors. For  $N=4$ , information rate of 1,875 bits per base.

## Motivation - Considering the security of a promising storage medium

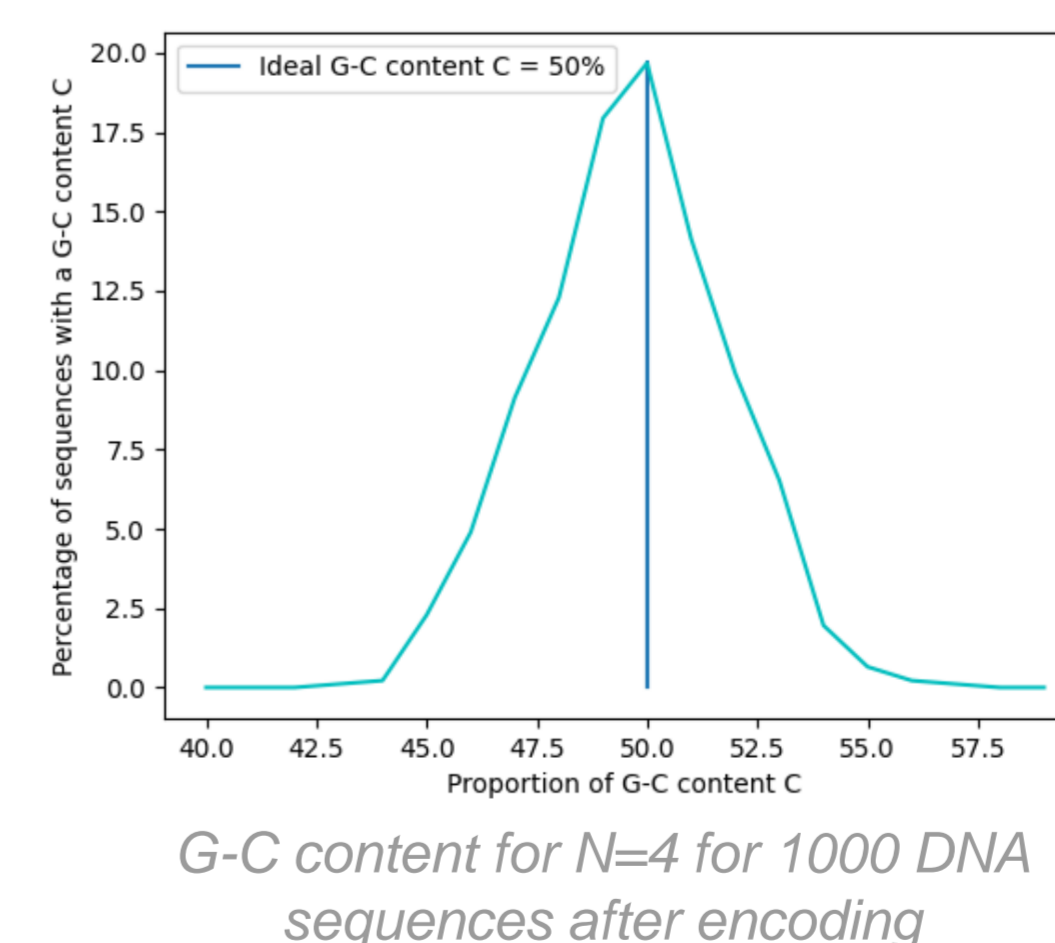
- **Context** – Actual storage technologies (flash memory, hard drives, magnetic tapes,..) are outpaced by the exponential rise of digital data production [1]
- **Advantages of DNA storage [2]** – Density of  $10^{21}$  bytes in one gram ( $10^6$  times more compact than hard disks), durability for centuries, energy cost close to zero (molecules kept at room temperature with no maintenance)
- **Motivation** – Introducing security to ensure the confidentiality of the data stored into DNA molecules



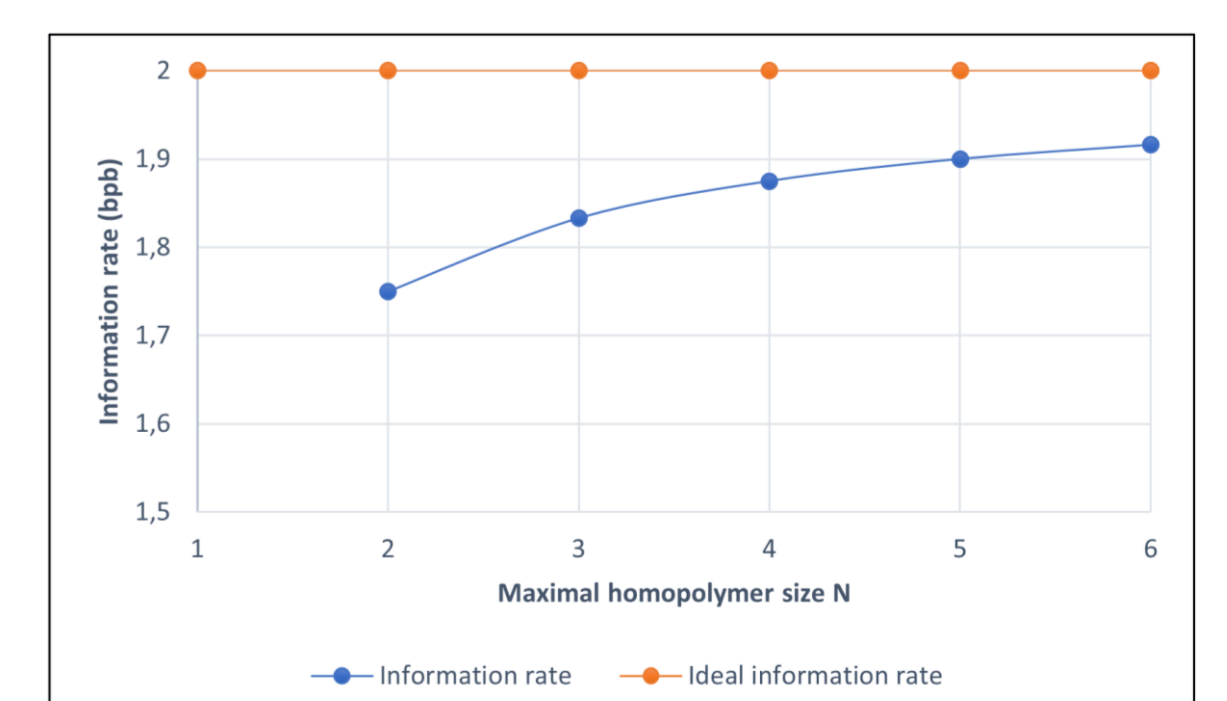
Data are UNPROTECTED in the storage chain, vulnerable to spying attacks or theft DNA-based data storage chain, including the constraints on the structure of DNA sequences to store, and types of errors caused by biological processes

## Dynamic encoding - a solution for confidentiality

- **Challenge** - Ensure confidentiality under biological constraints while approaching the ideal rate of 2 bits of information per base
- **Solution** – A three step coding process that includes encryption, dynamic data encoding and error-correction code
  - **Step 1:** Encryption with AES-256 to ensure confidentiality and to regulate the G-C base rate and homopolymers
  - **Step 2:** Dynamic encoding to manage unwanted base patterns; encoding based on the addition or not of one bit of data every  $N$  bases to avoid homopolymers longer than  $N$
  - **Step 3:** Non-binary LDPC error-correction code to correct any base substitutions, deletions or insertions



G-C content for  $N=4$  for 1000 DNA sequences after encoding



Information rate (bits per base) depending on  $N$ , the maximal homopolymer size

## Conclusion and future work

- **Confidentiality** in the entire storage chain that takes into account **biological constraints**
- Encoding solution **independent** from encryption algorithm and error-correction code, **adaptable** to the size of unwanted patterns
- Extend the approach to other synthesis and sequencing technologies

Contact: chloe.berton@imt-atlantique.fr, gouenou.coatrieux@imt-atlantique.fr