

Nº 179057

Challenges of encoding digital data into DNA sequences and decoding it back to digital

Andre Guilherme da Costa Martins

*Palestra apresentada na
disciplina de Biologia Molecular
do Programa de Pós-graduação
Multicêntrico em Bioquímica e
Biologia Molecular, do Campus
Centro-Oeste Dona Lindu-UFSJ*

A série “Comunicação Técnica” compreende trabalhos elaborados por técnicos do IPT, apresentados em eventos, publicados em revistas especializadas ou quando seu conteúdo apresentar relevância pública.

PROIBIDO REPRODUÇÃO

Challenges of encoding digital data into DNA sequences and decoding it back to digital

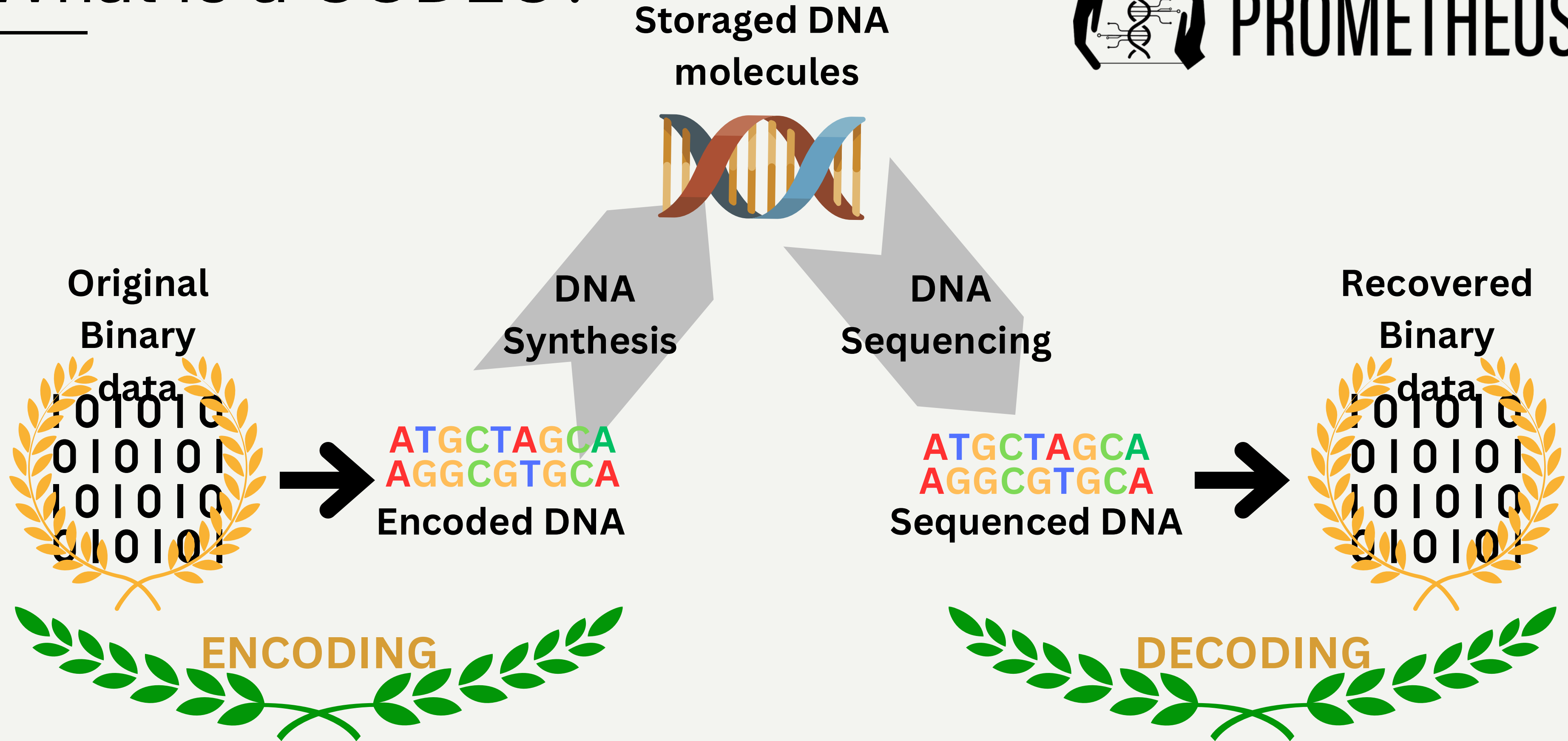


André Guilherme da Costa Martins, PhD Biomed. Sci.
Bioinformatics Researcher

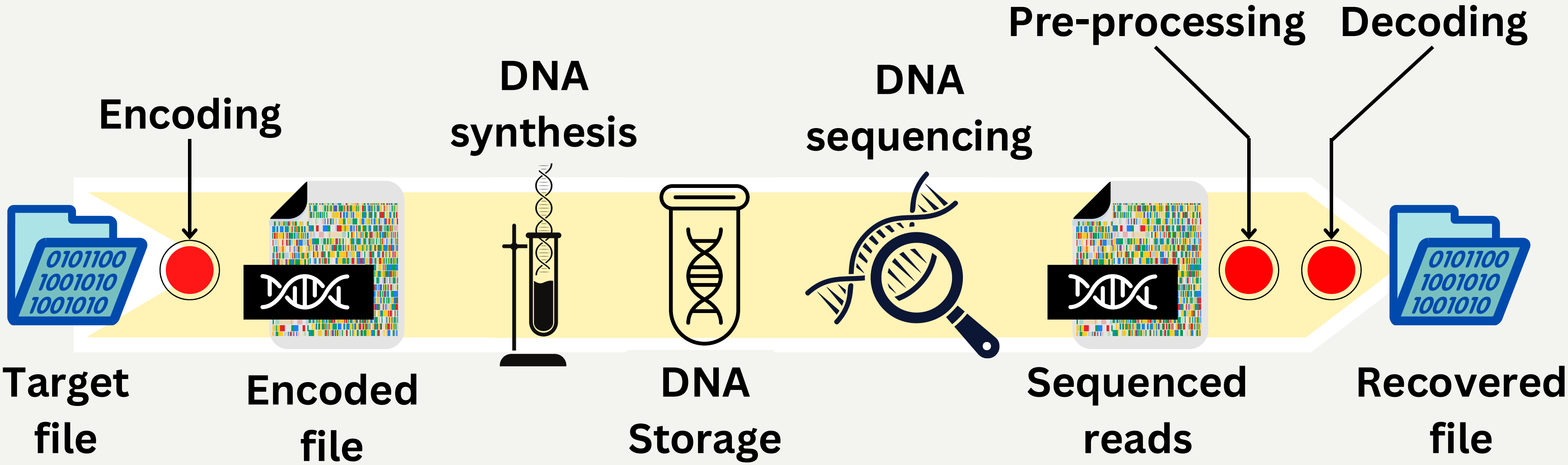
Institute for Technological Research - IPT, Brazil

andremartins@ipt.br

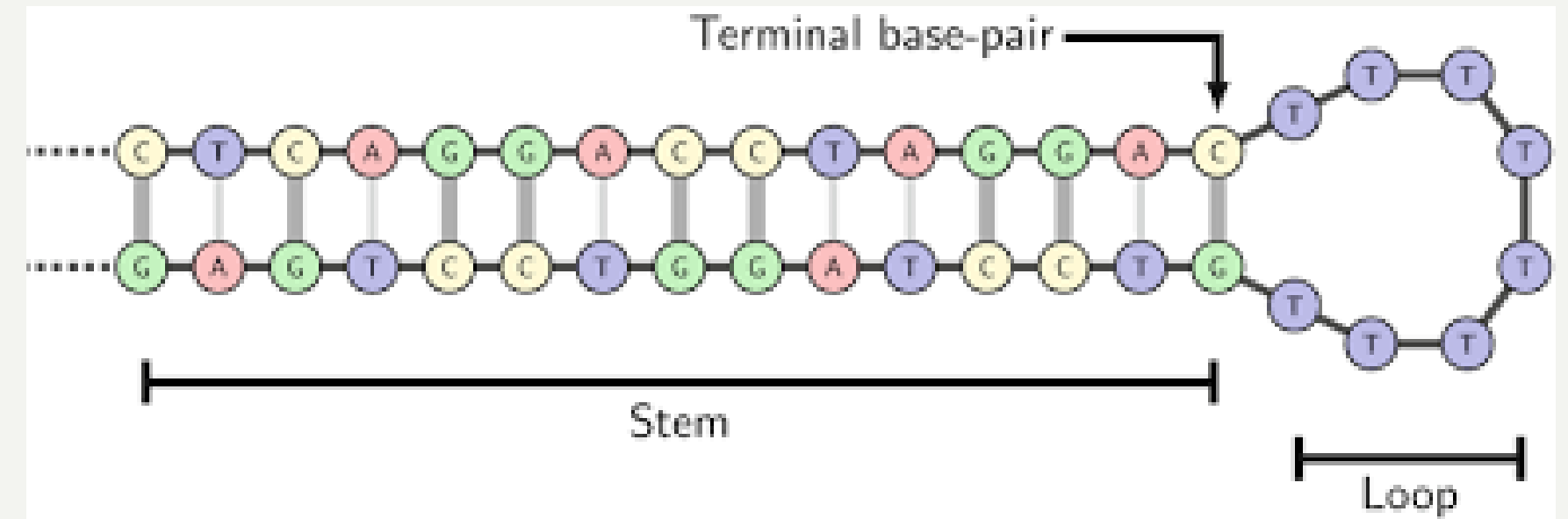
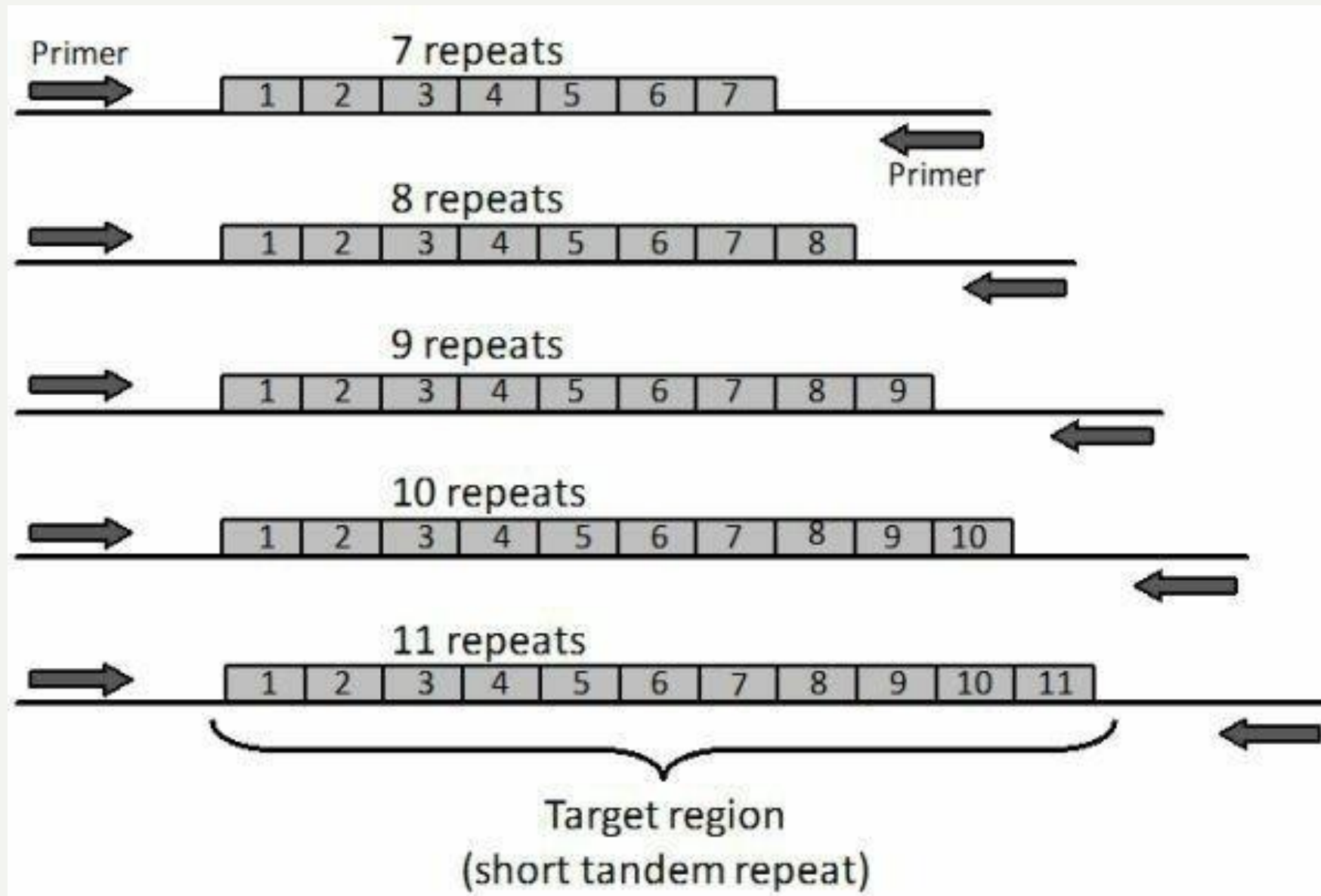
What is a CODEC?



The workflow



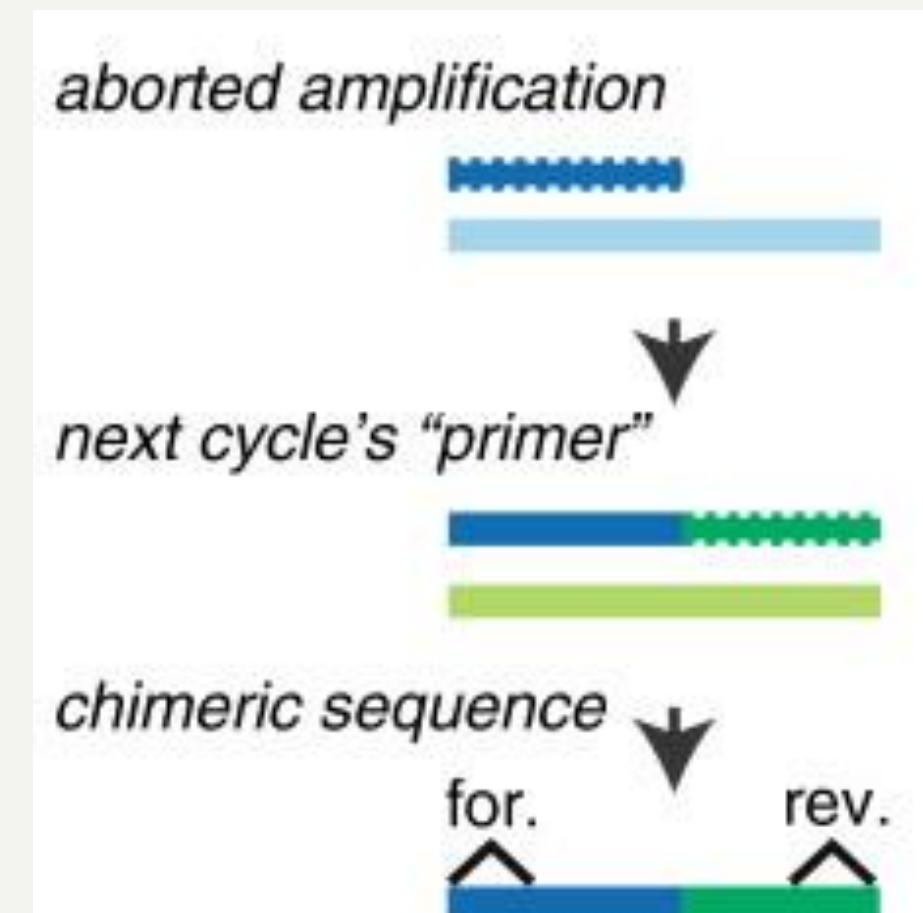
Challenges of encoding



DNA hairpins

- 1 2 3 ...
- 2-nucleotide repeat unit : (CA)(CA)(CA).....
 - 3 -nucleotide repeat unit : (GCC)(GCC)(GCC)
 - 4 -nucleotide repeat unit : (AATG)(AATG)(AATG)
 - 5 -nucleotide repeat unit : (AGAAA)(AGAAA)

Repetitive segments



Interspecific similarity

Challenges of encoding

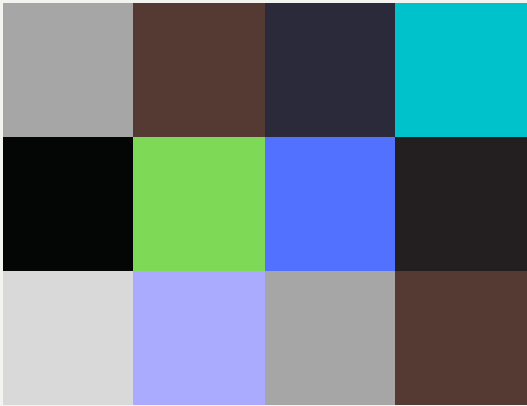
IPT - Main gate



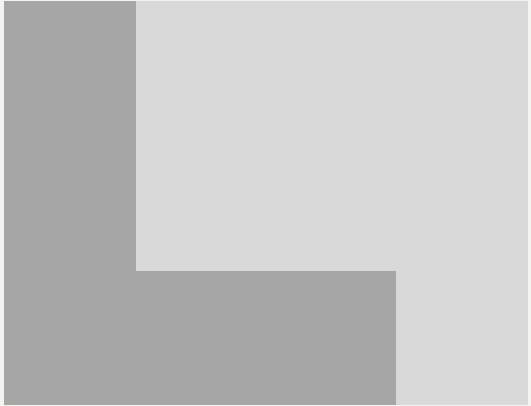
Bitmap image



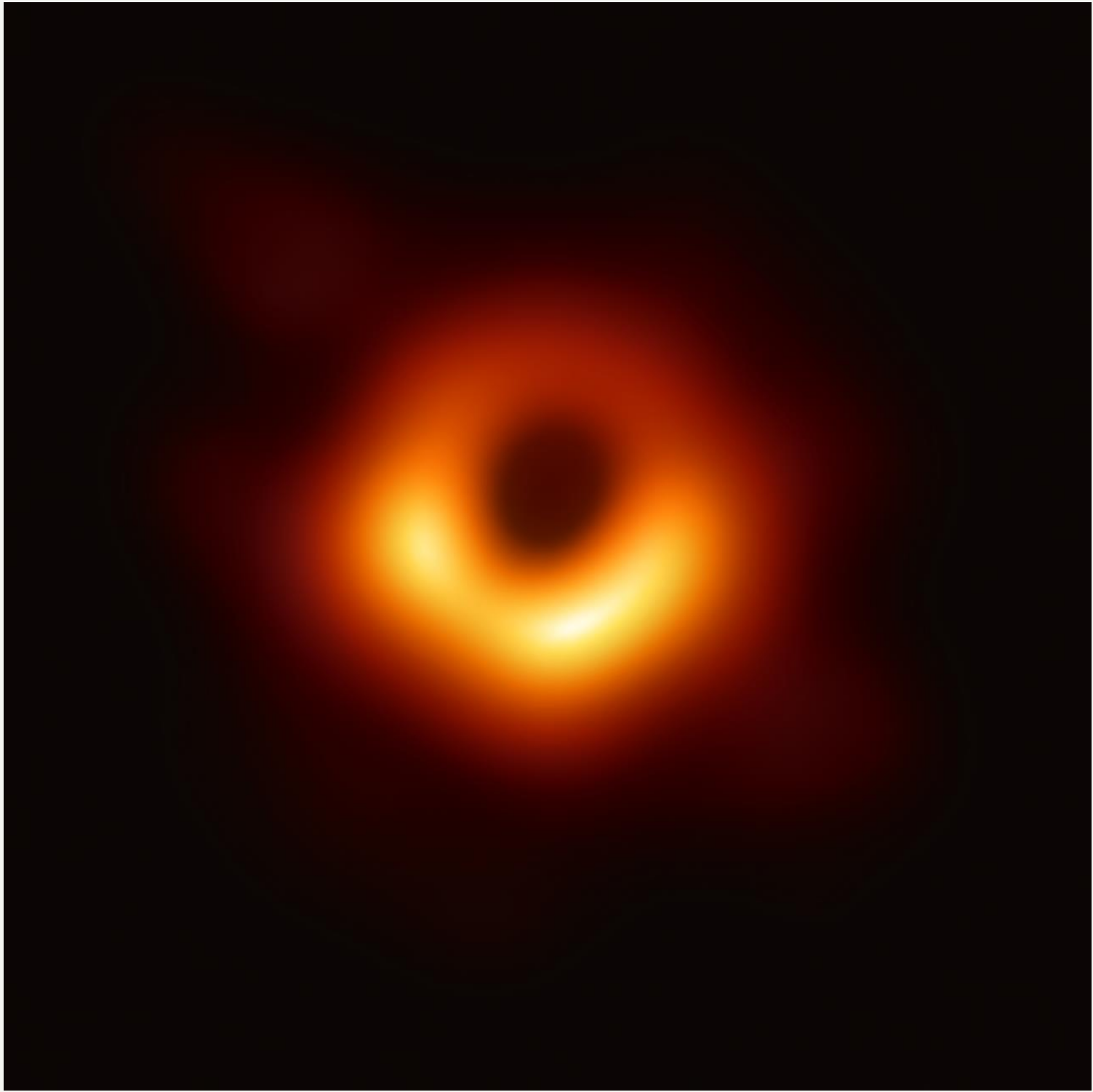
High complexity regions



Low complexity regions

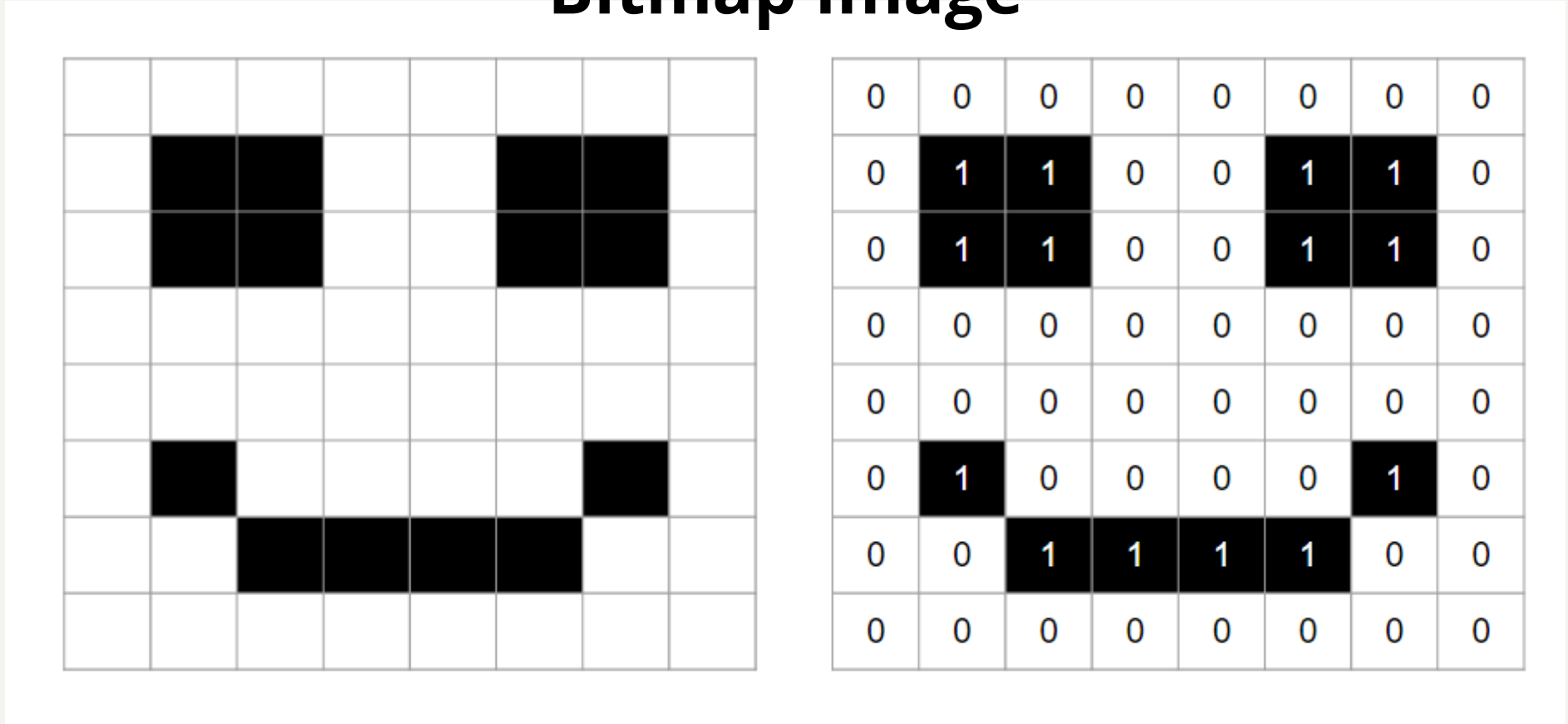


Challenges of encoding

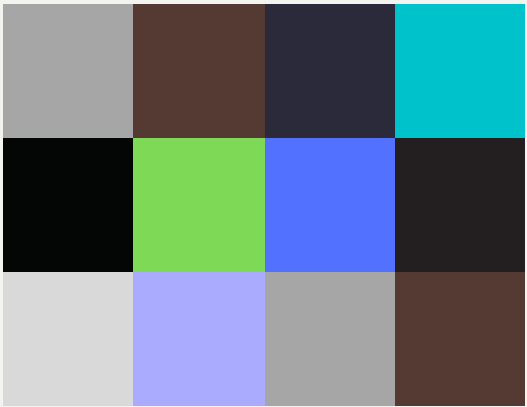


M87* black hole, taken on 11 April 2017

Bitmap image



High complexity regions

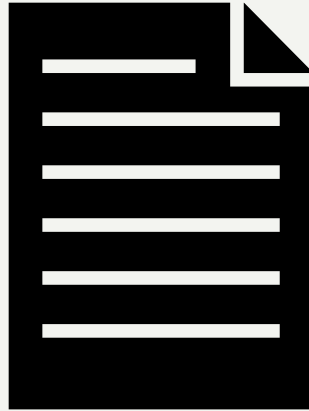


Low complexity regions

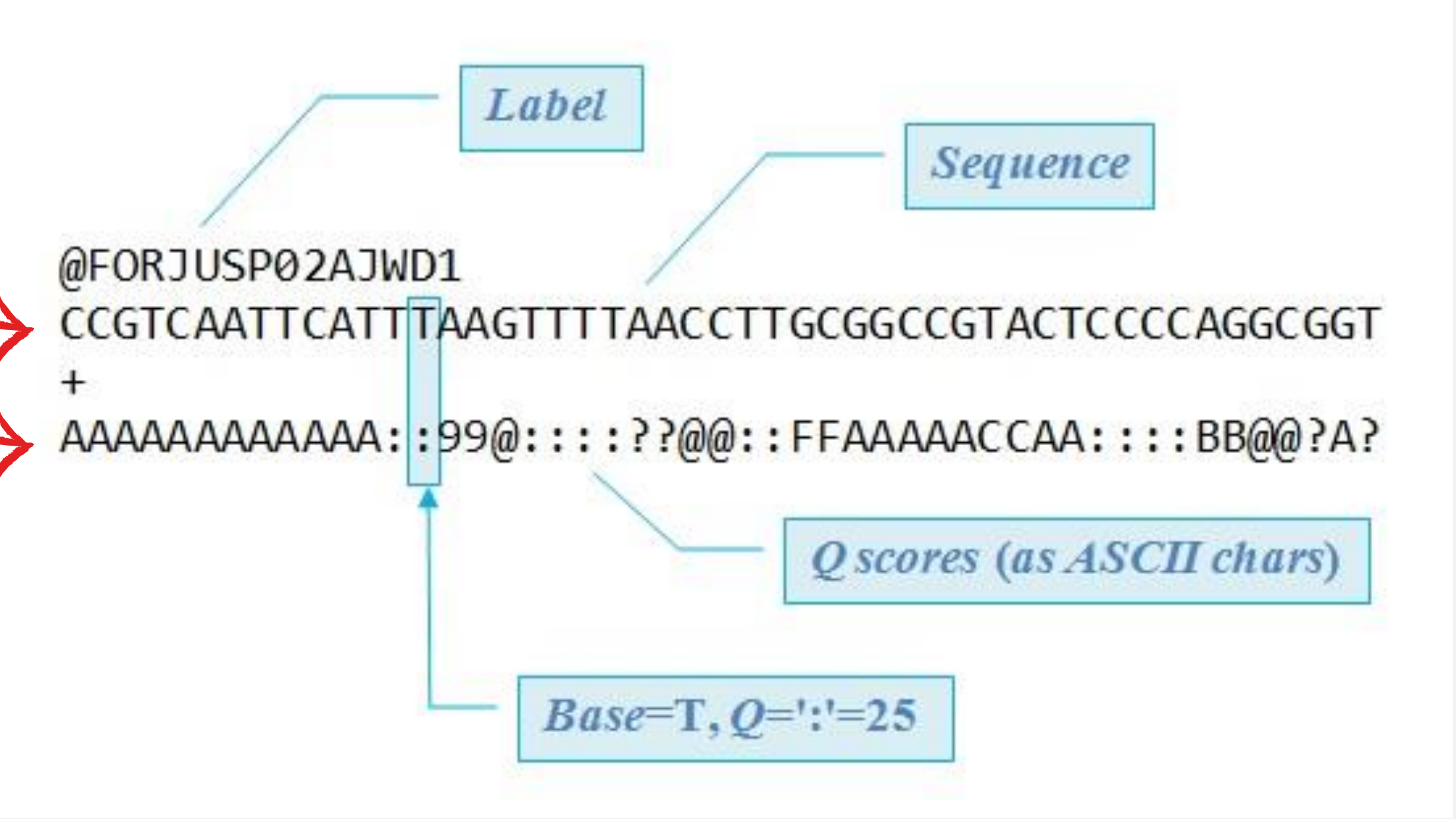


Challenges of processing NGS reads

Fastq



DNA
Sequence
quality



Challenges of processing NGS reads



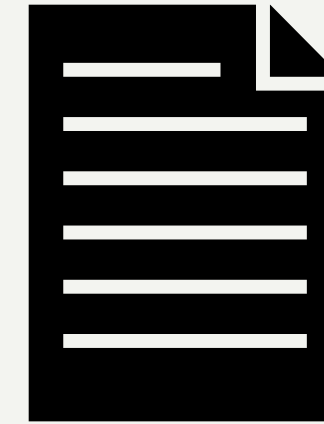
Encoding complexity

Fasta



Decoding complexity

Fastq



Challenges of processing NGS reads

 Analogy!

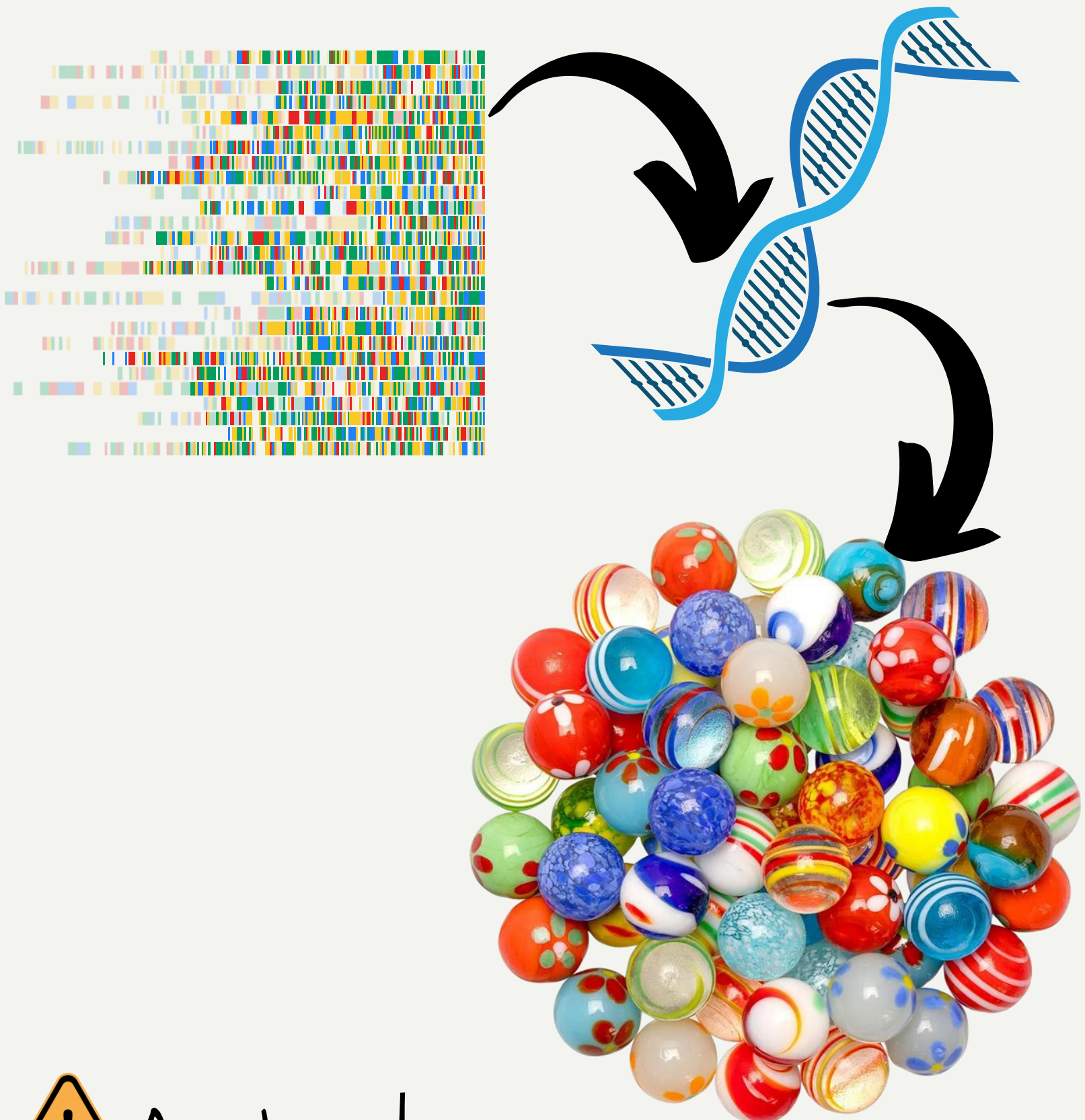

Fastq



Challenges of processing NGS reads

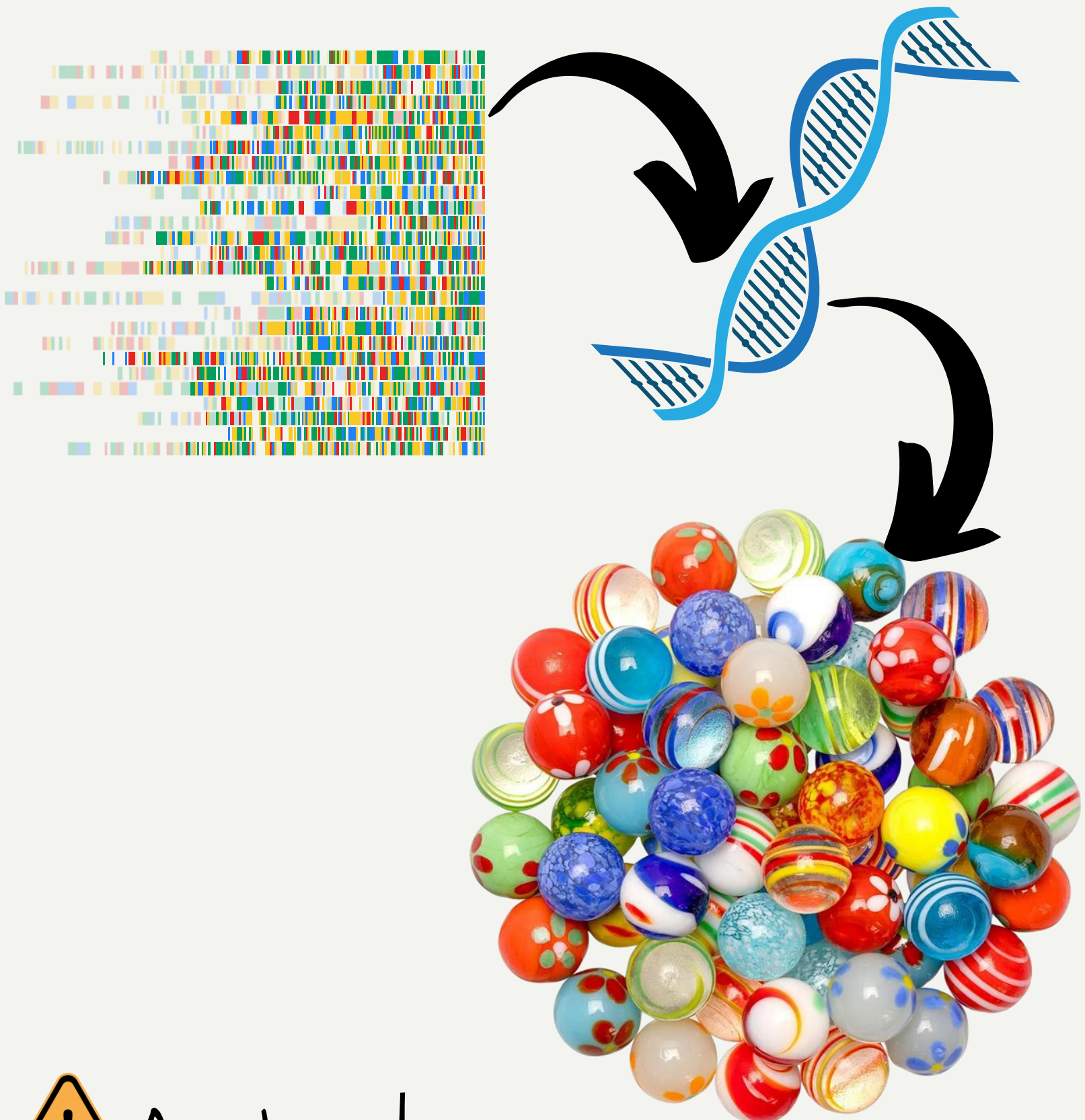
```
*Ep1035 20-3/1 -ACGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGTACACAATGGTAGAGCAG
Ep1035-1/1 --CGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035-3/1 ----CGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035-5/1 -ACGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 10-1/1 --CGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 10-3/1 -ACGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 10-5/1 -ACGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-AACAATGGTAGAGCAG
Ep1035 100-1/1 CACGACGTTGT--AACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 100-3/1 -ACGACGTTGT--AACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 100-5/1 --CGACGTTGT--AACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 13-1/1 CACGACGTTGTAAAACGACAAAAG-CATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 13-3/1 ---GACGTTGTAAAACGACAAAAG-CATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 13-5/1 ----ACGTTGTAAAACGACAAAAG-CATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 28-3/1 CACGACGTTGTAAAACGACAAAAGCCATTTCAT-CCCAGGT-CACAATGGTAGAGCAG
Ep1035 28-5/1 CACGACGTTGTAAAACGACAAAAGCCATTTCAT-CCCAGGT-CACAATGGTAGAGCAG
Ep1035 29-1/1 CACGA-GTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 29-3/1 -ACGA-GTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 29-5/1 CACGA-GTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 3-1/1 ---GACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 3-3/1 -ACGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 3-5/1 ----CGTTGTAAAACGACTAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 30-1/1 CACGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 30-3/1 -ACGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 30-5/1 CACGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
Ep1035 32-1/1 -ACGACGTTGTAAAACGACAAAAGCCATTTCATCCCCAGGT-CACAATGGTAGAGCAG
```


Challenges of decoding

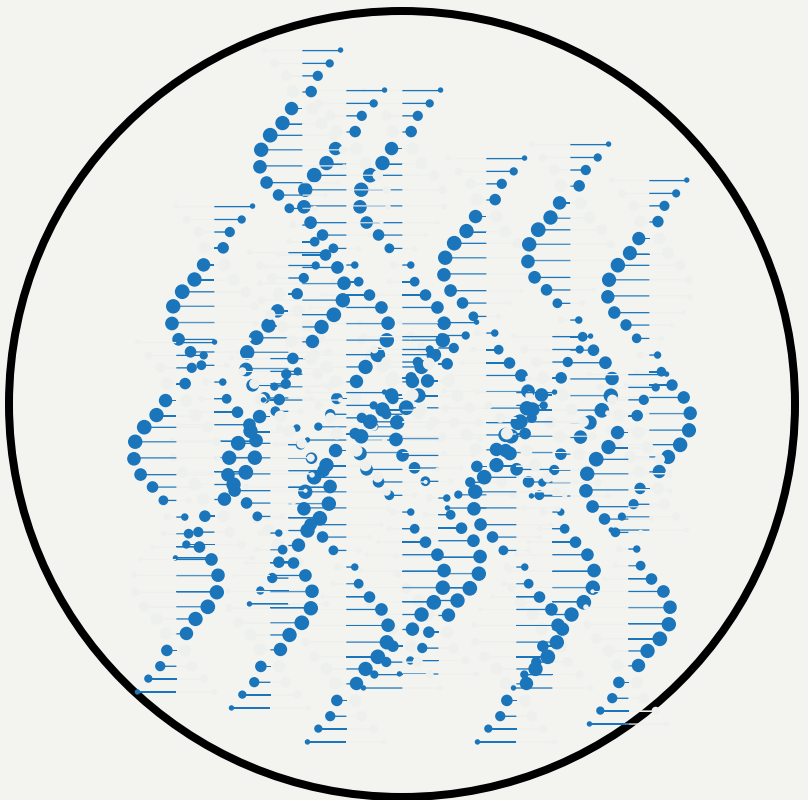


 *Analogy!*

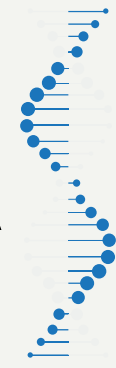
Challenges of decoding



Each DNA synthesis site



10^n copies of same ssDNA



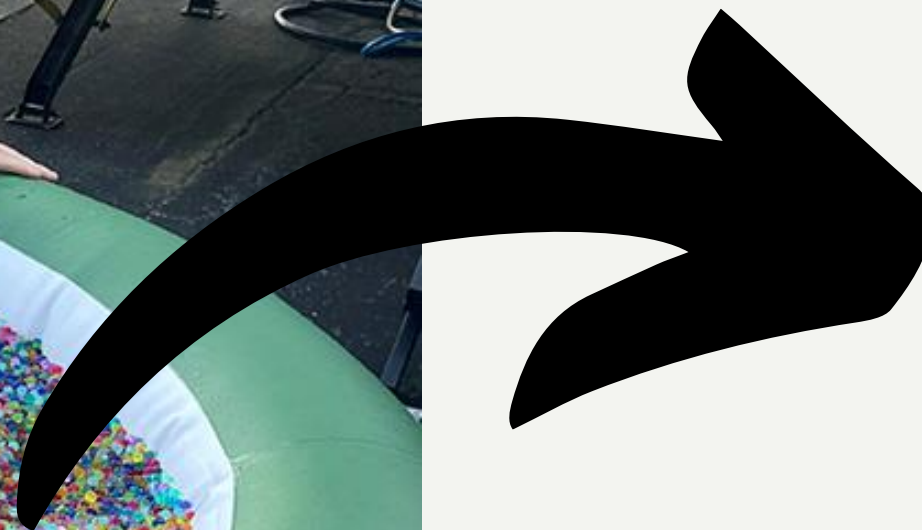
PCR amplification
dsDNA



 **Analogy!**

Challenges of decoding

 Analogy!



Challenges of decoding



Library preparation for NGS



Illumina Nextseq 2000

 Analogy!

Challenges of decoding



Expectations... :)

Challenges of decoding



Reality... :(

Developed tools




 PROMETHEUS



APOLLO 



ARTEMIS 



CHIRON 



HERMES 

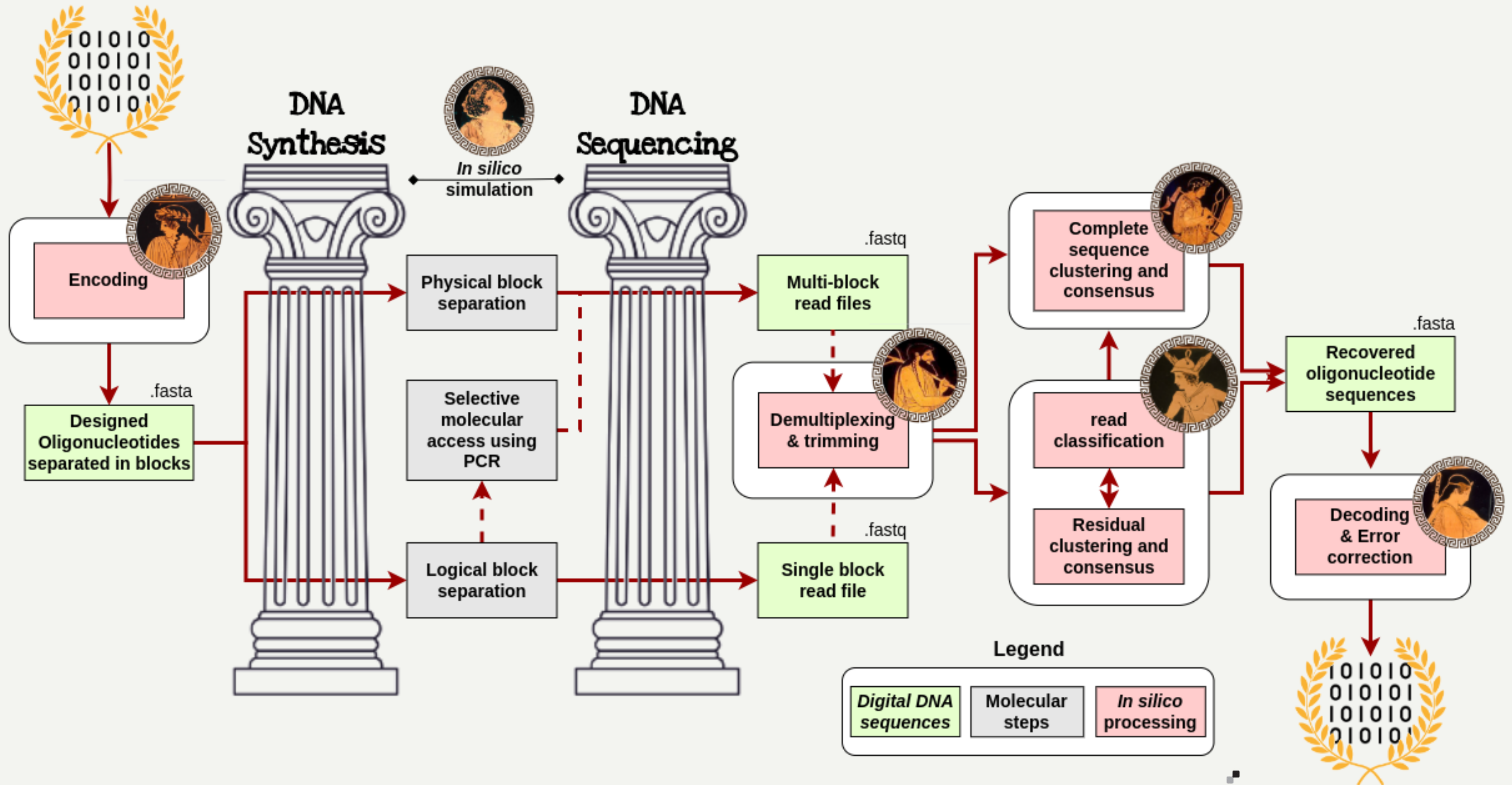


HEPHAESTUS 



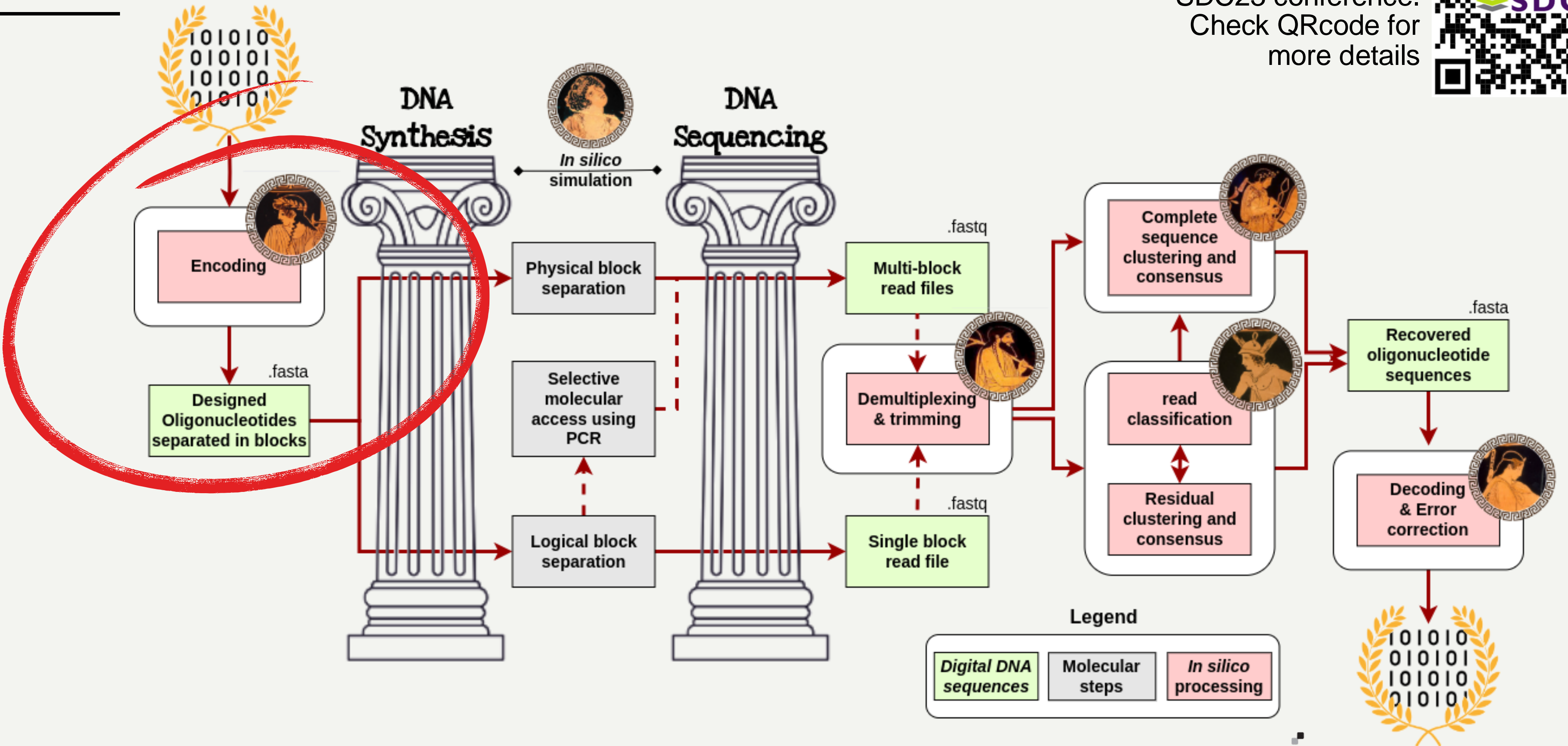
GAIA 

DDS Workflow

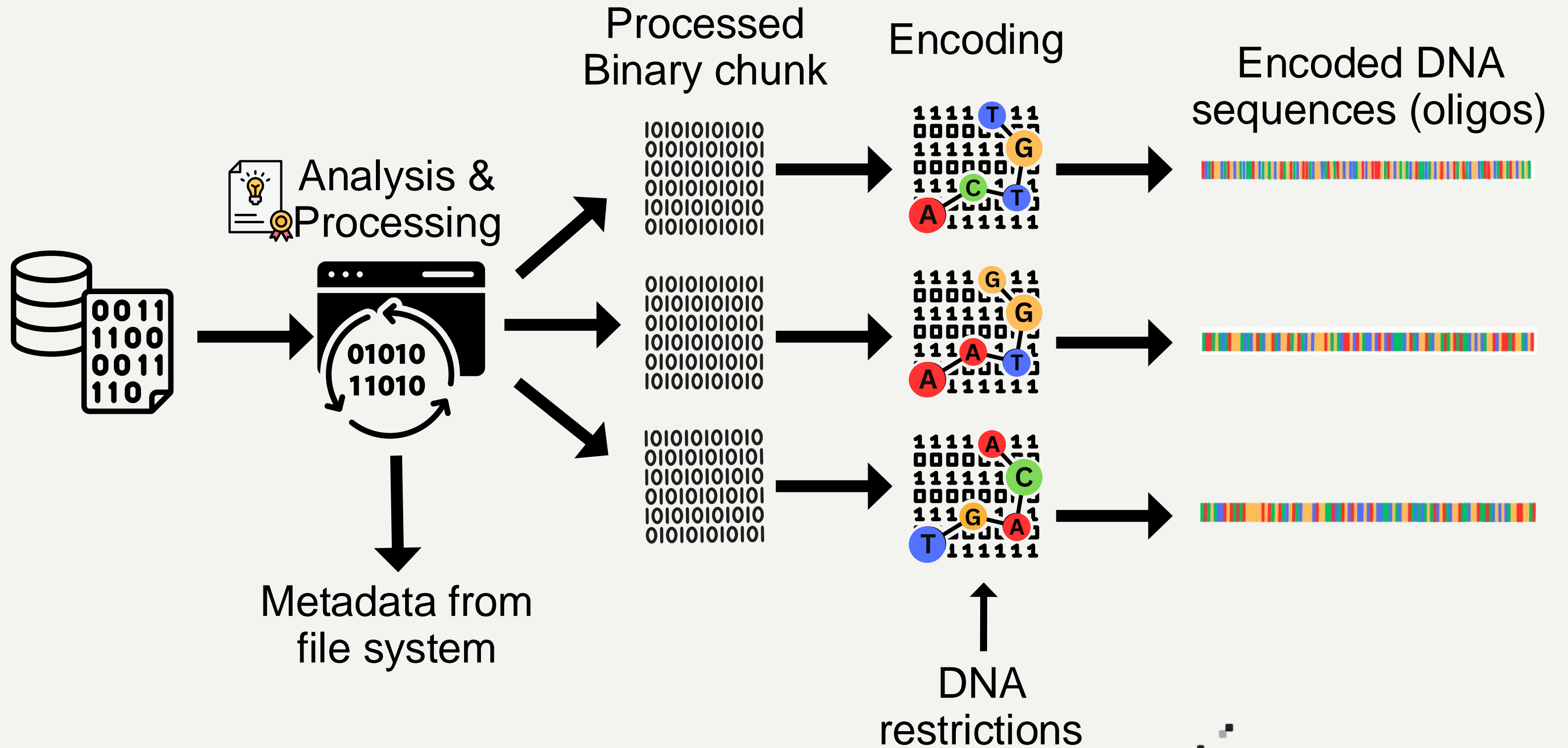
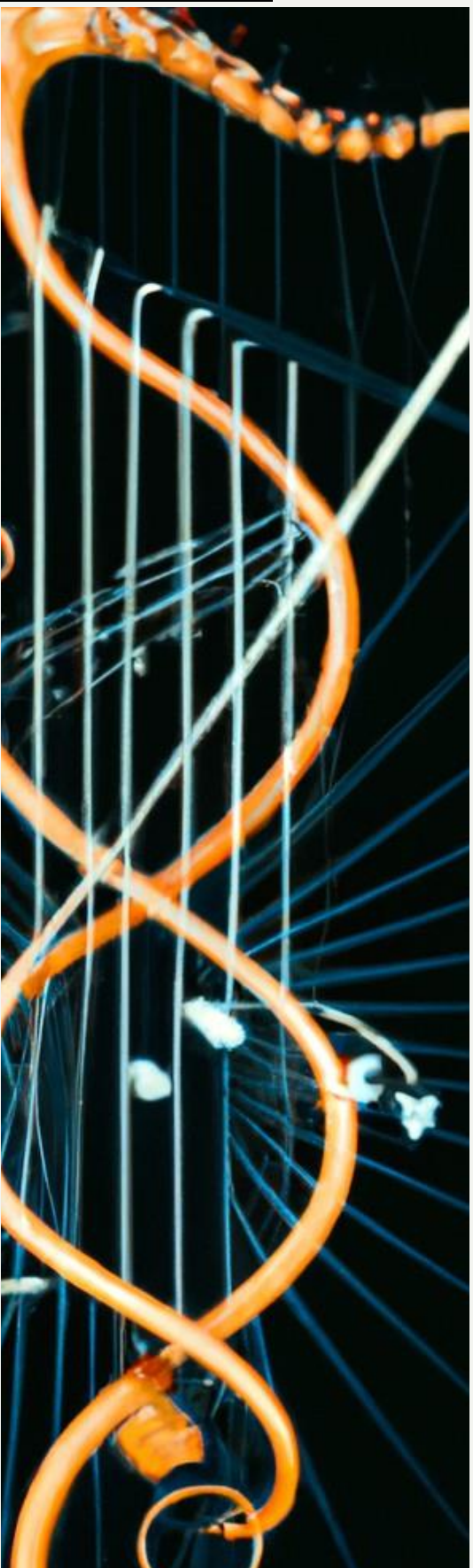


DDS Workflow

We have presented all Pantheon DNA Gods at SDC23 conference. Check QRcode for more details

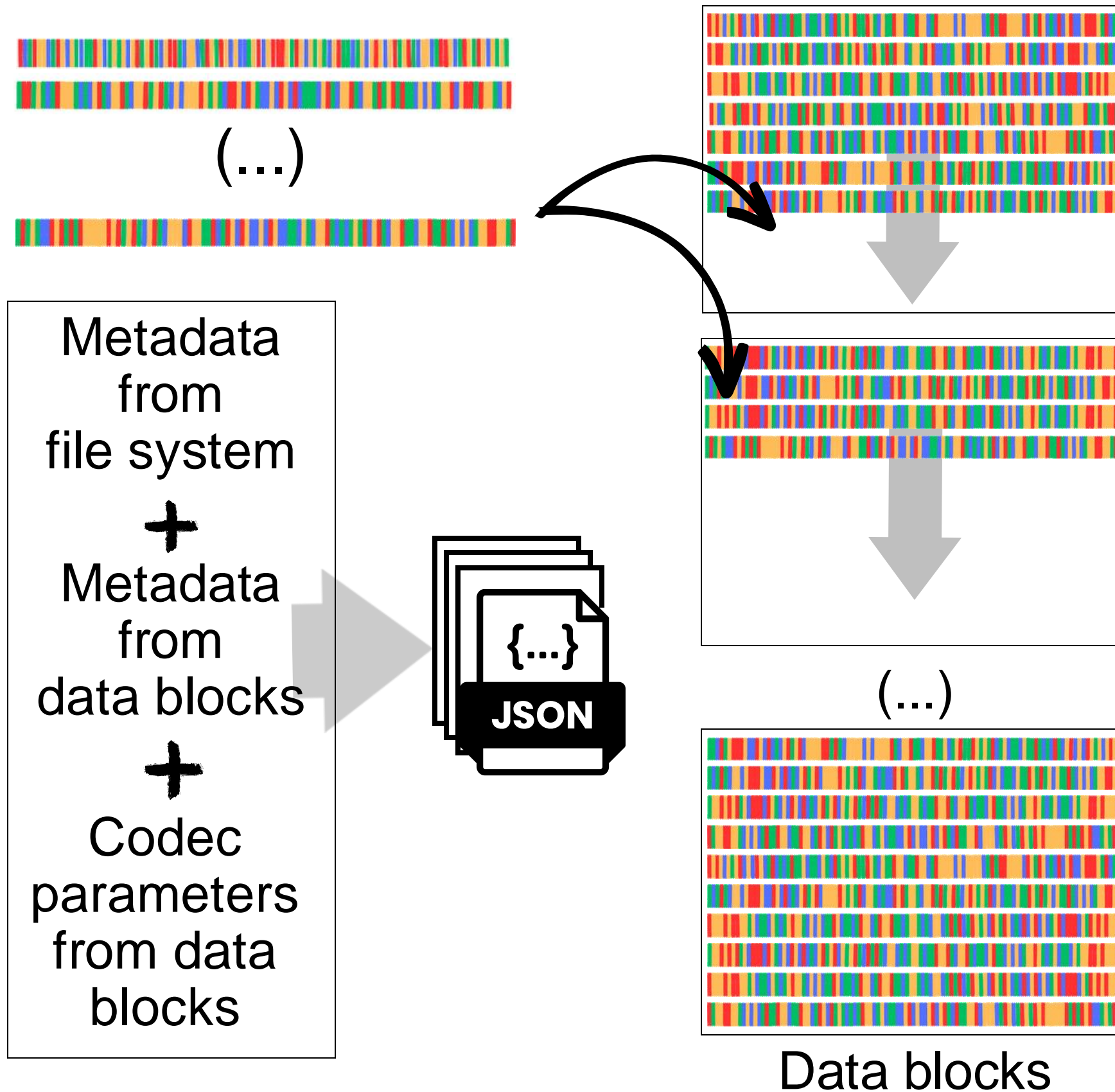


How data is organized into DNA?



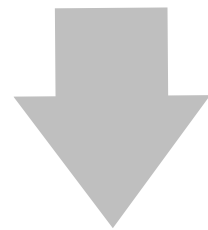
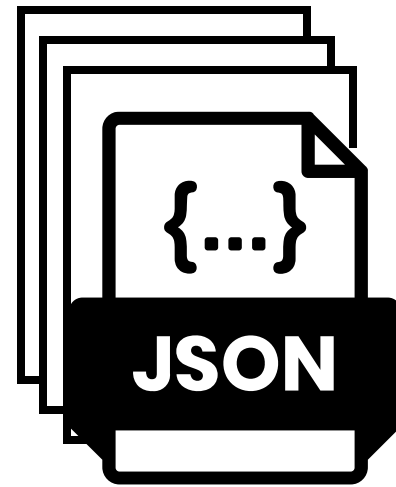
Data Block structure

-  Cytosine
-  Guanine
-  Adenine
-  Thymine



“Disk sectors”
 Analogy!

Data Block structure



Archive Metadata Block (AMB)



Outer
ECC



Outer
ECC



(...)

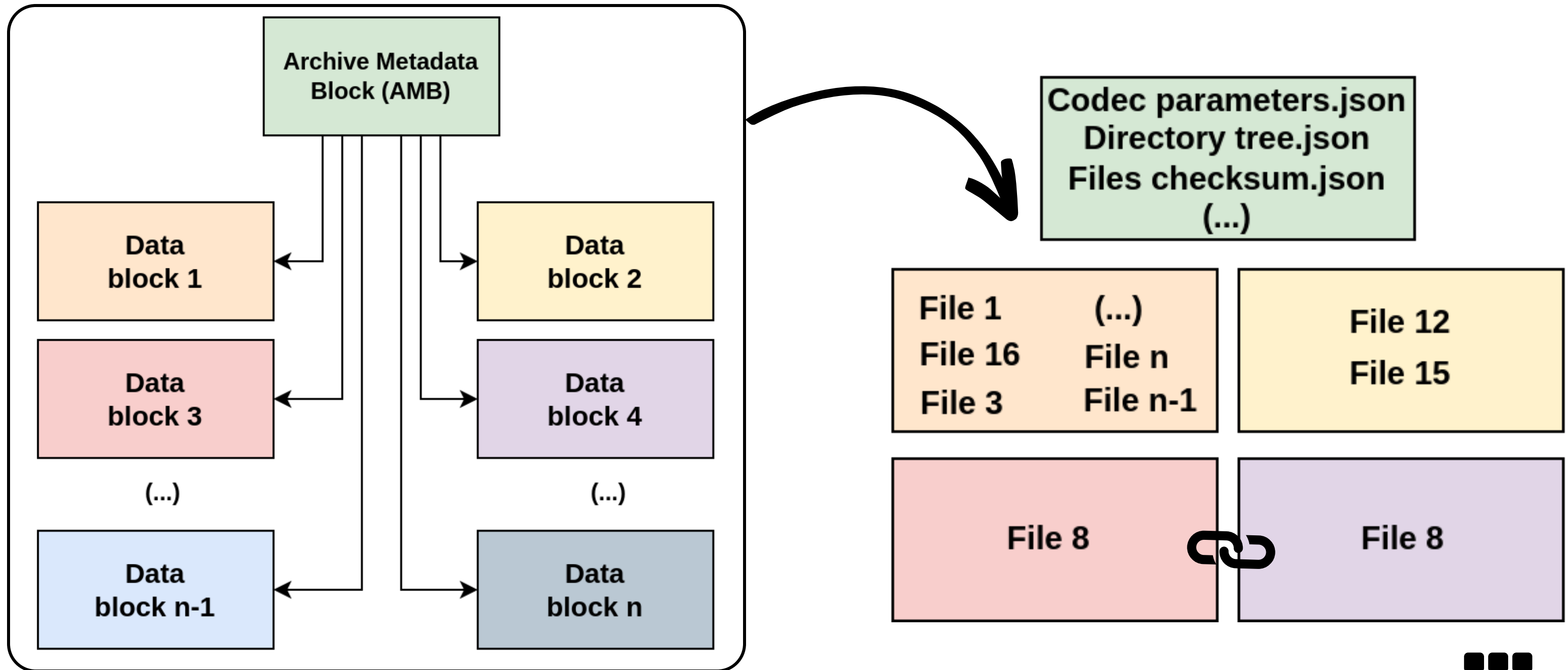


Outer
ECC



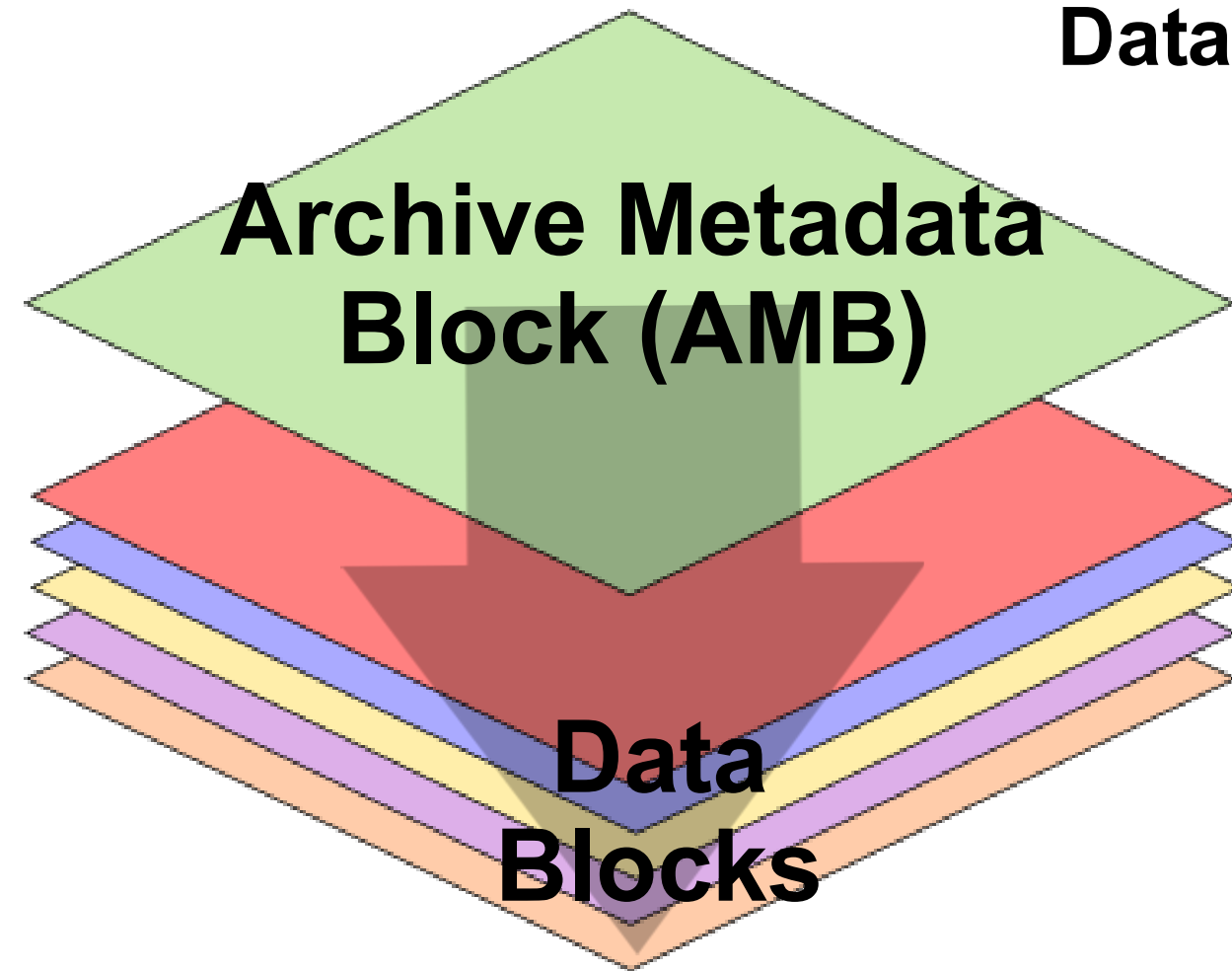
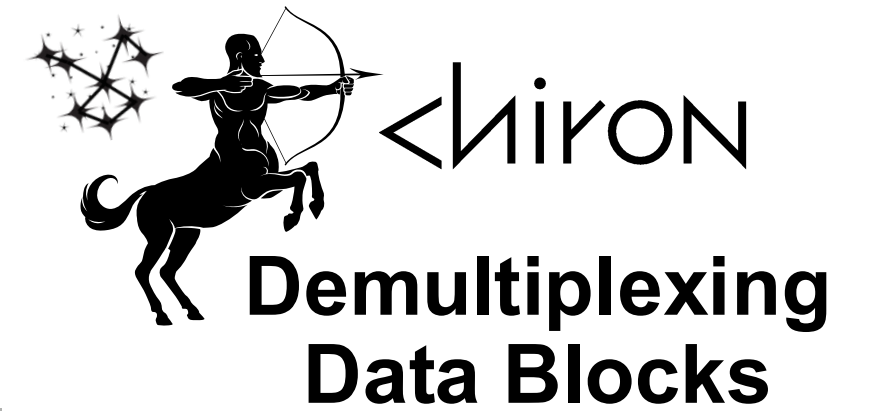
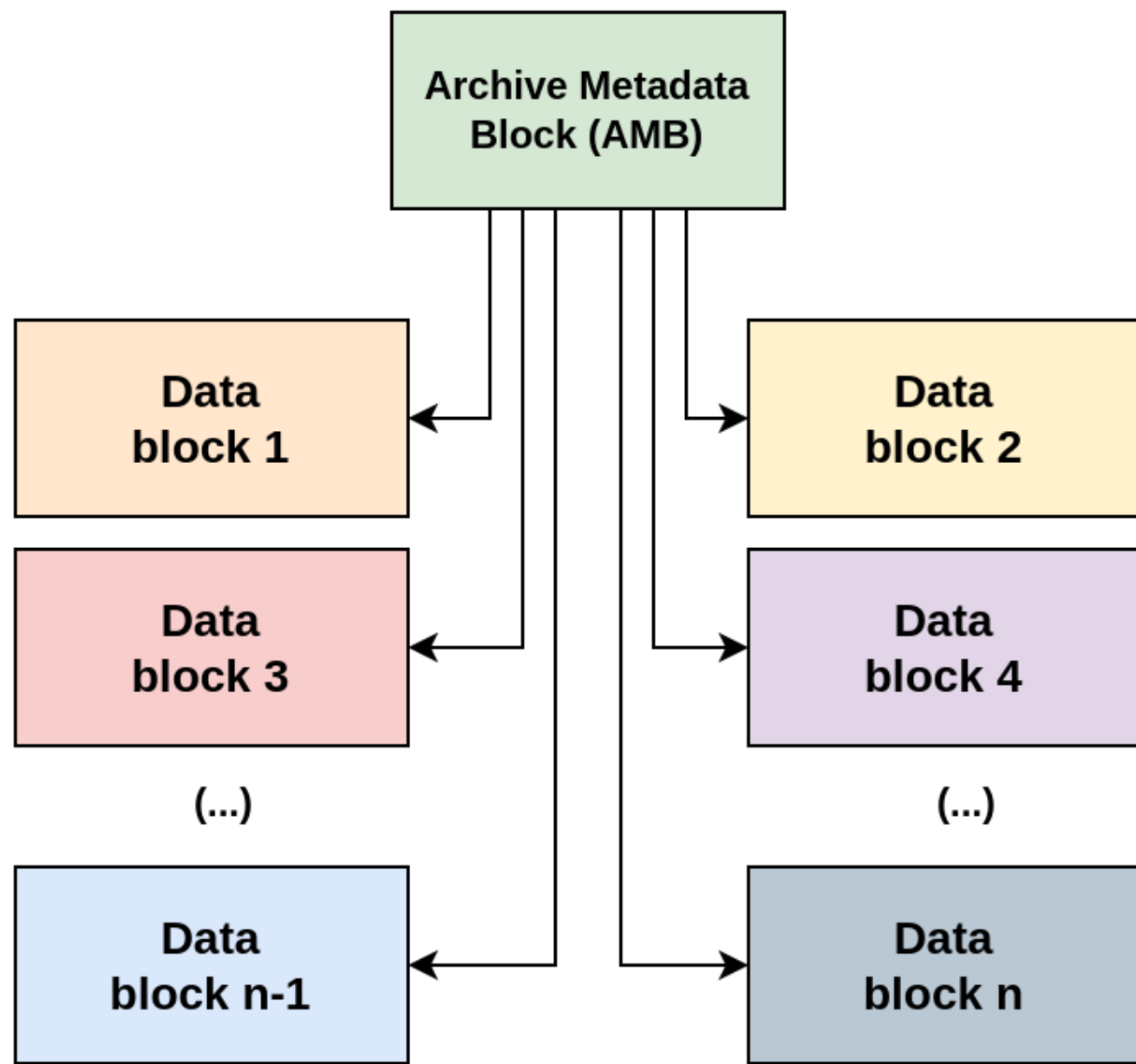
Data blocks

How data are organized into DNA?



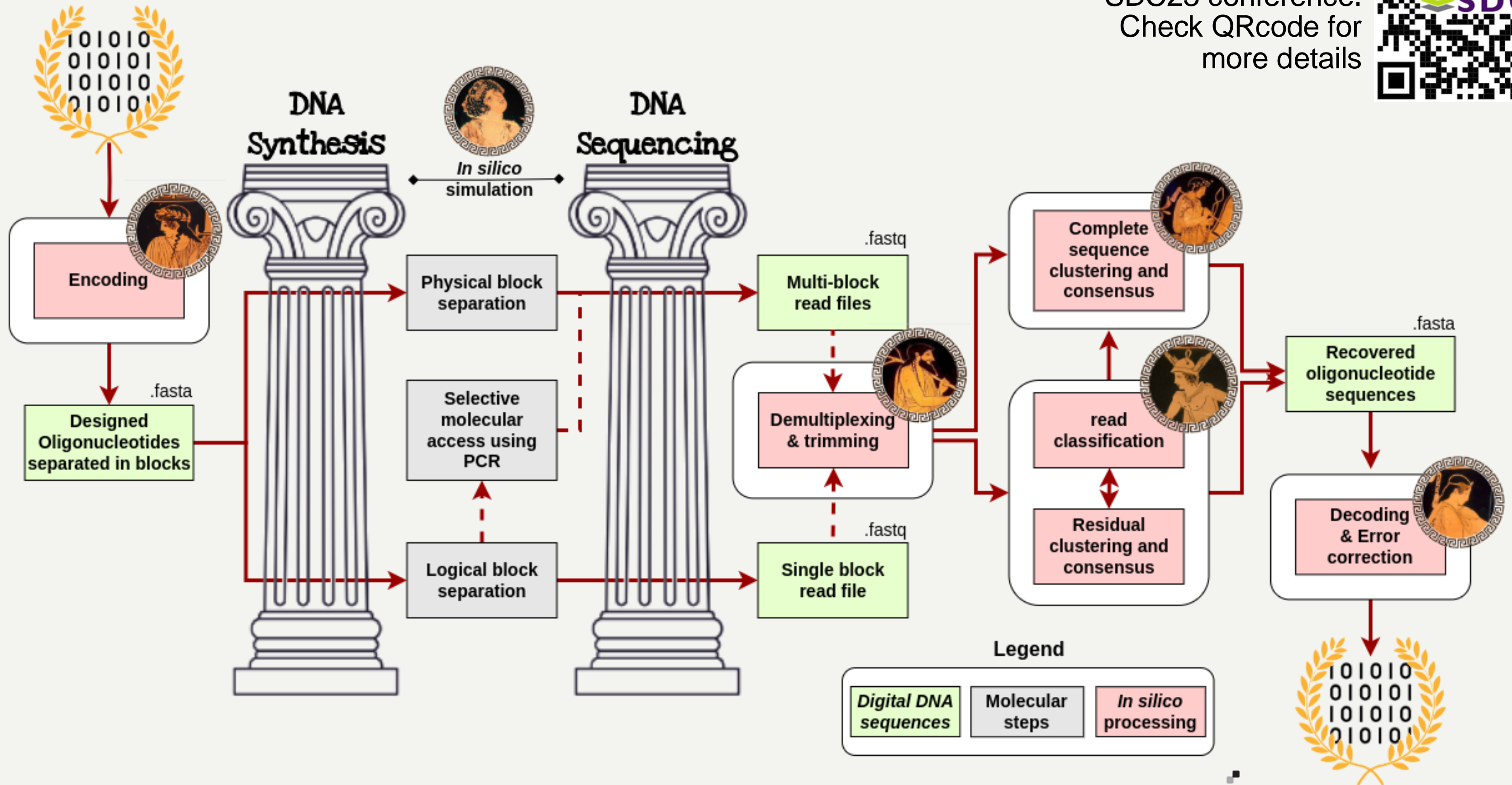
Filling the block is like a Tetris game... 

How to read it back?



DDS Workflow

We have presented all Pantheon DNA Gods at SDC23 conference. Check QRcode for more details



Pre-processing NGS reads

Pre-processing steps:

- Adapters/Primer trimming
- Demultiplexing coding blocks
- Merge read pair (paired-end strategy)
- Reorient DNA sequences
- Discard low-quality reads and contaminants

Block specific prime pairs



Planned data blocks

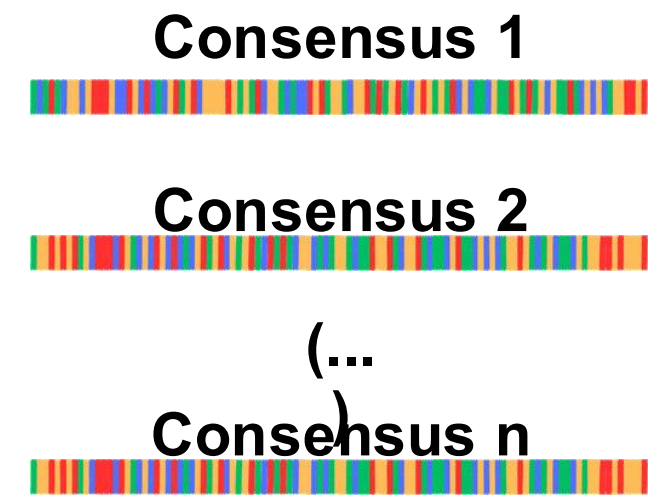
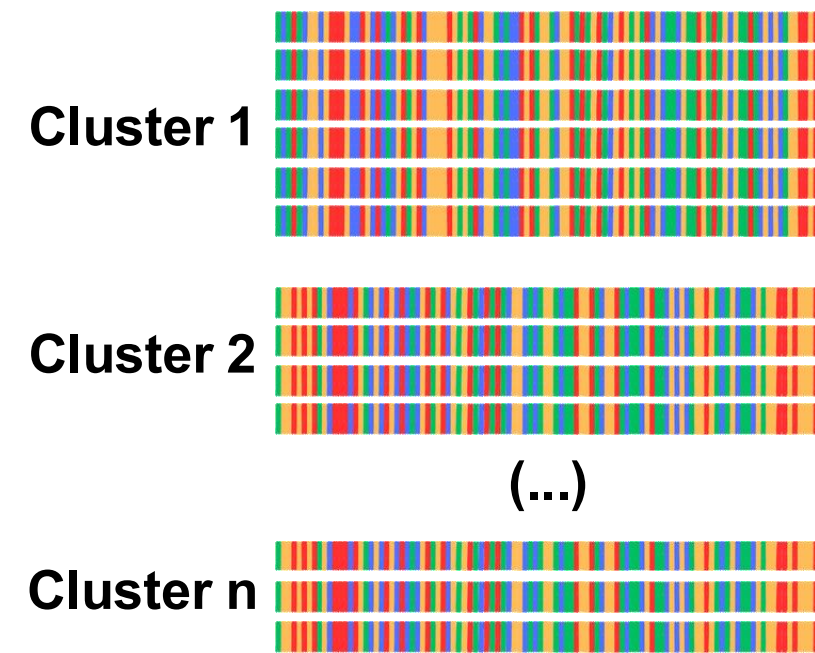
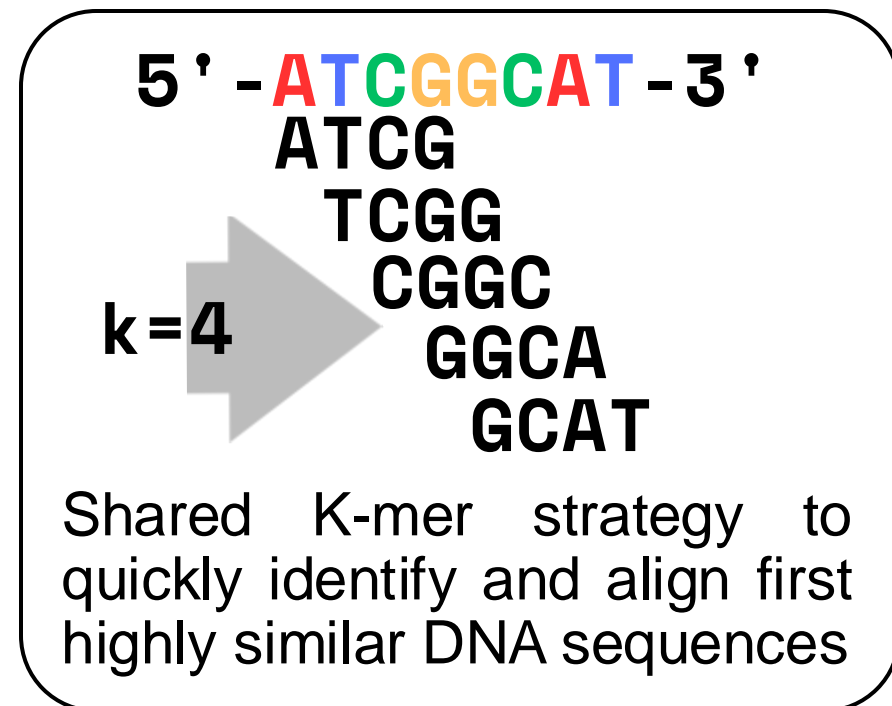
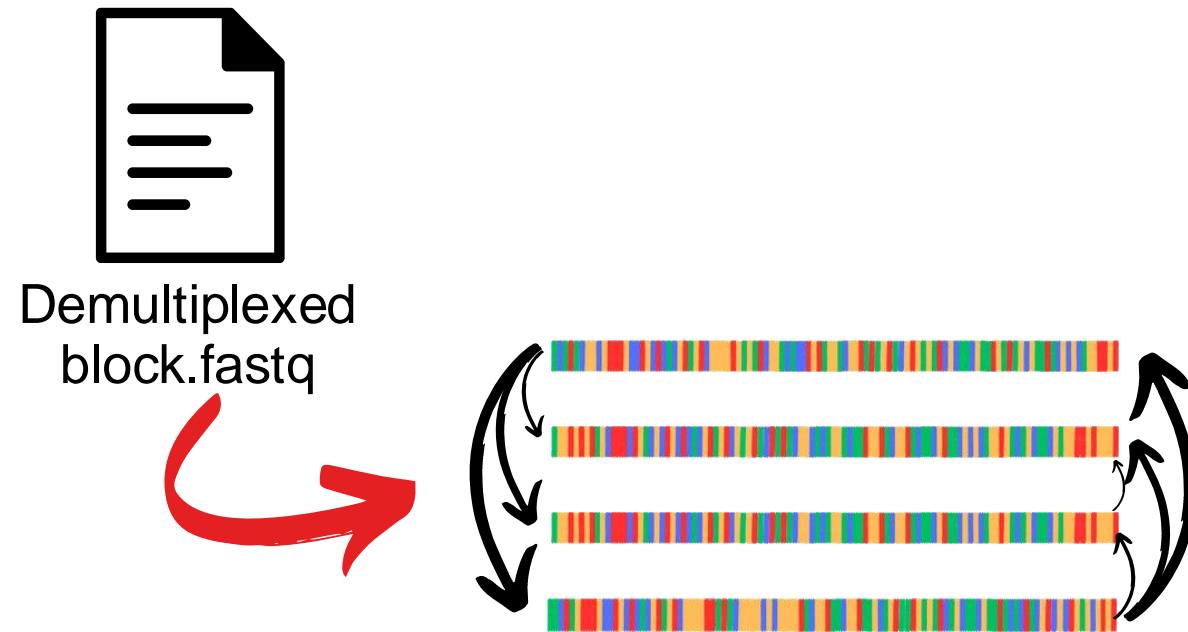


Sequenced DNA reads

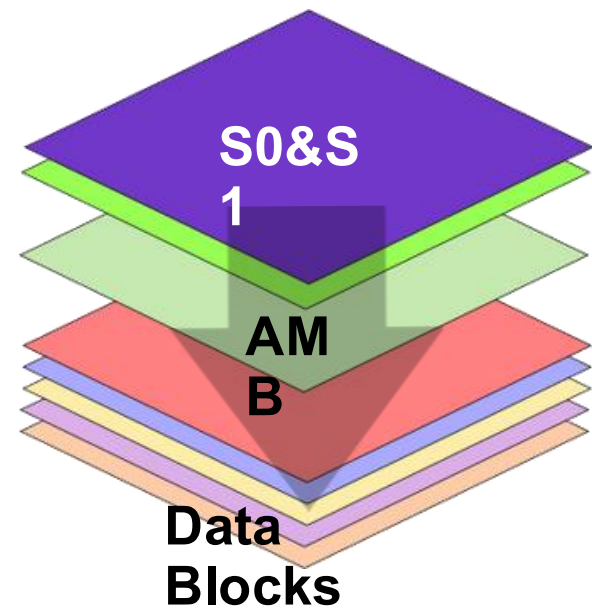
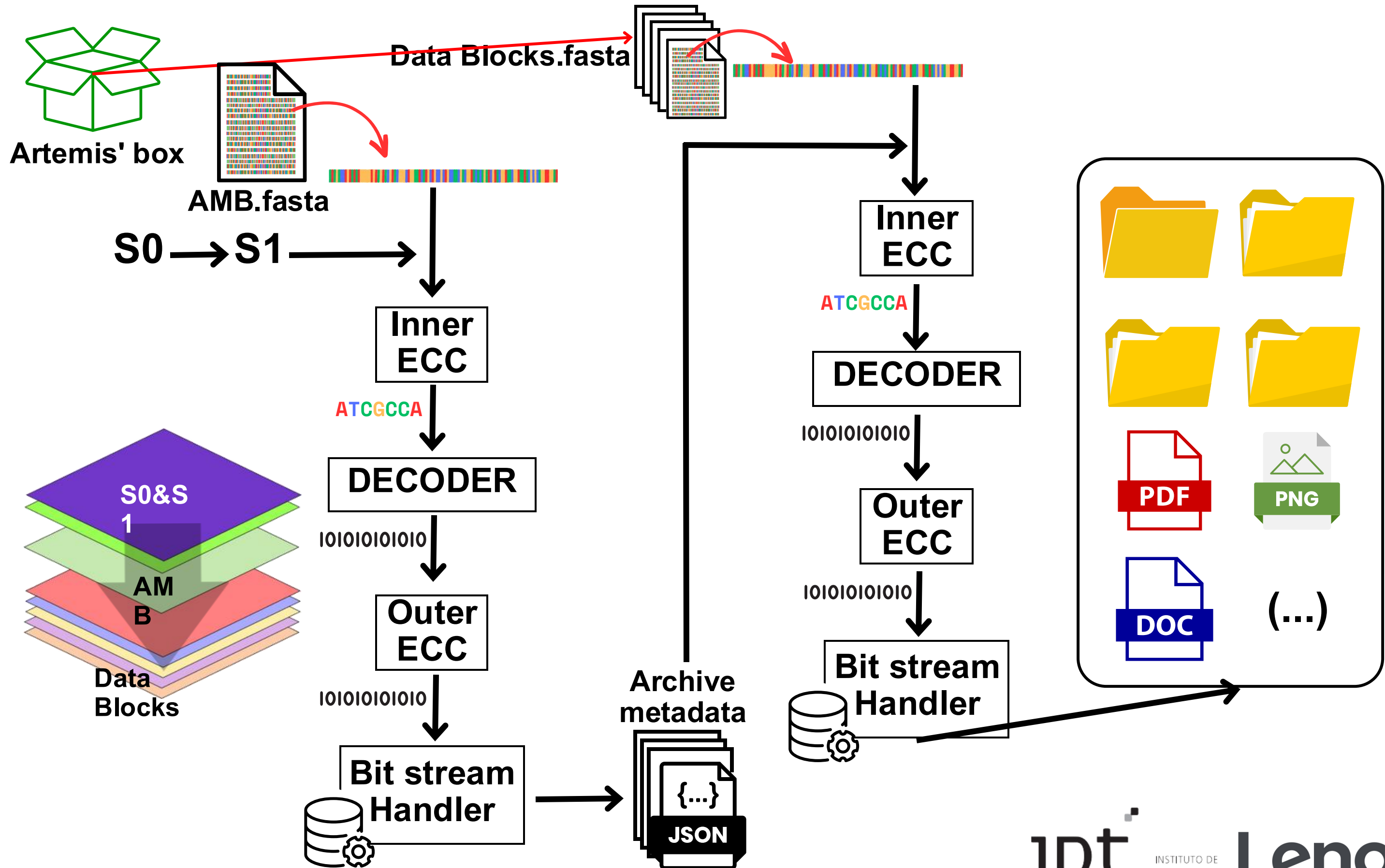
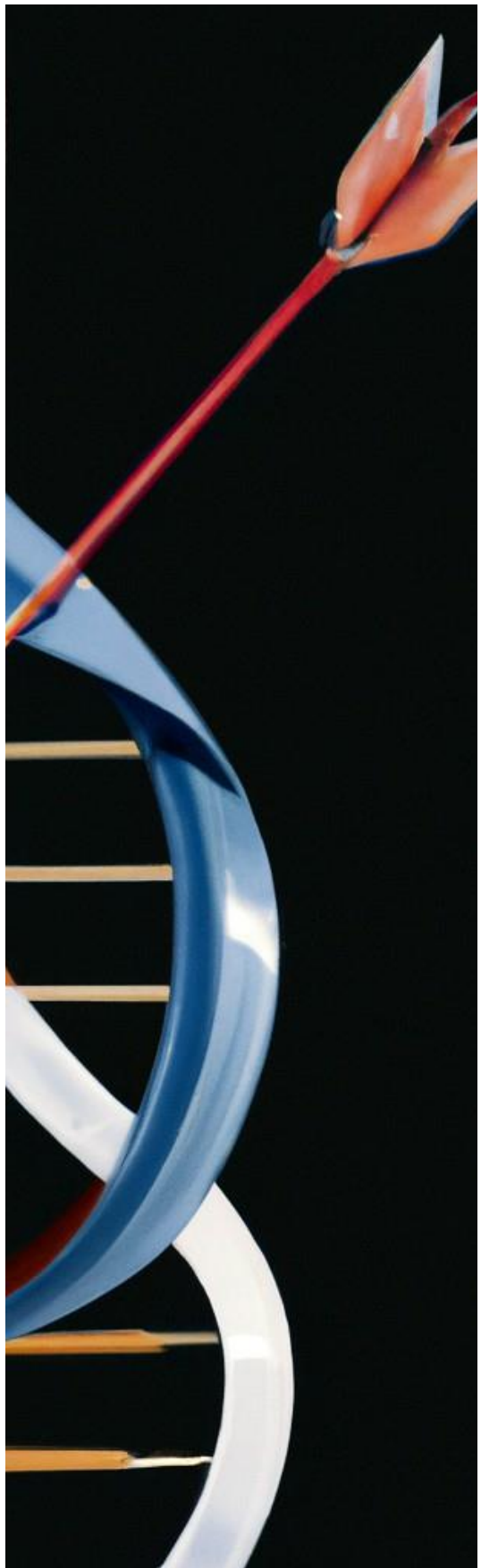


Sequenced & Demultiplexed into data blocks

Pre-processing NGS reads



Decoding



Simulation tool



- What Gaia does:

- Simulate different sequencing strategies



- Single or Paired-ends
- Library preparation
- Coverage variation
- Sequencing platforms

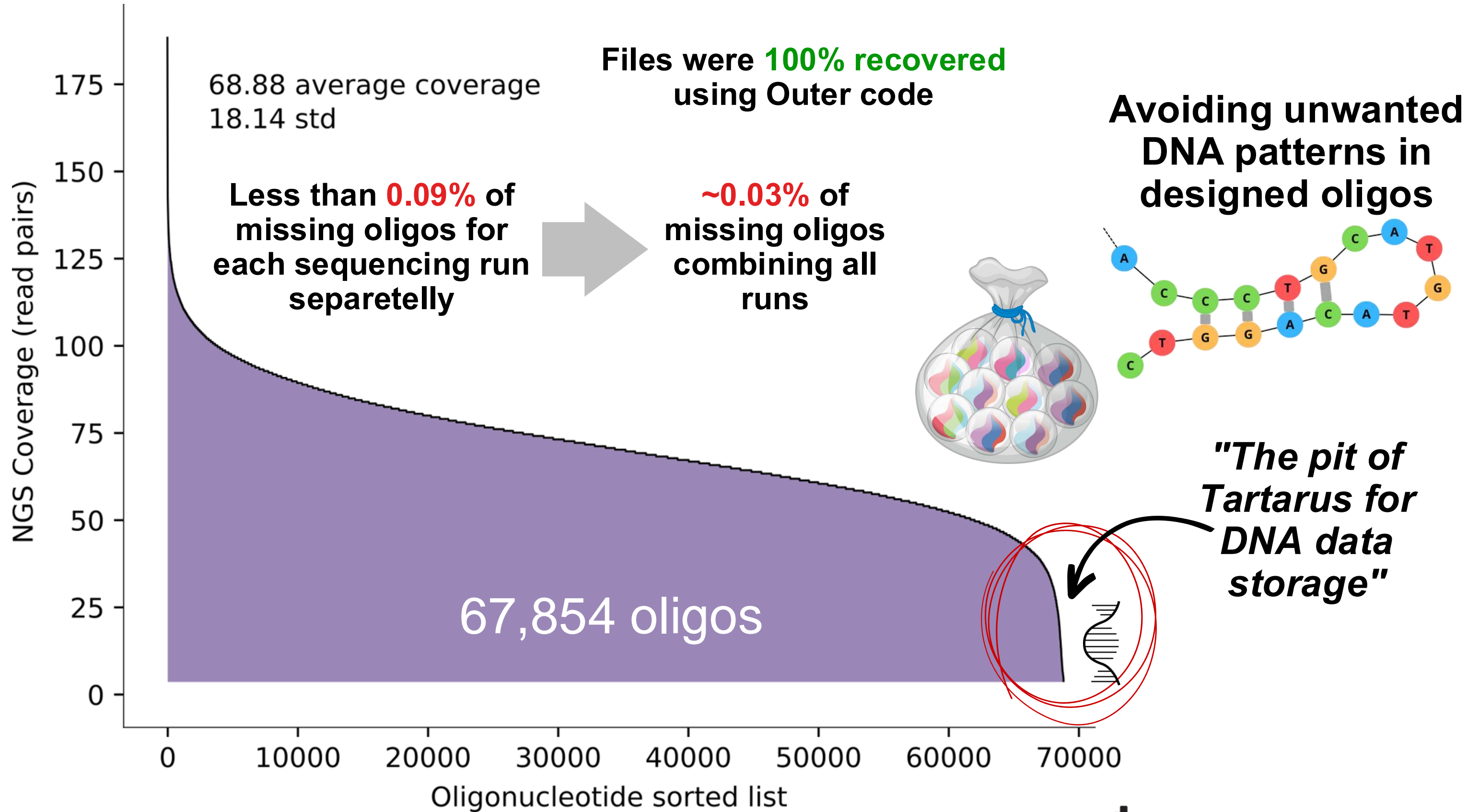
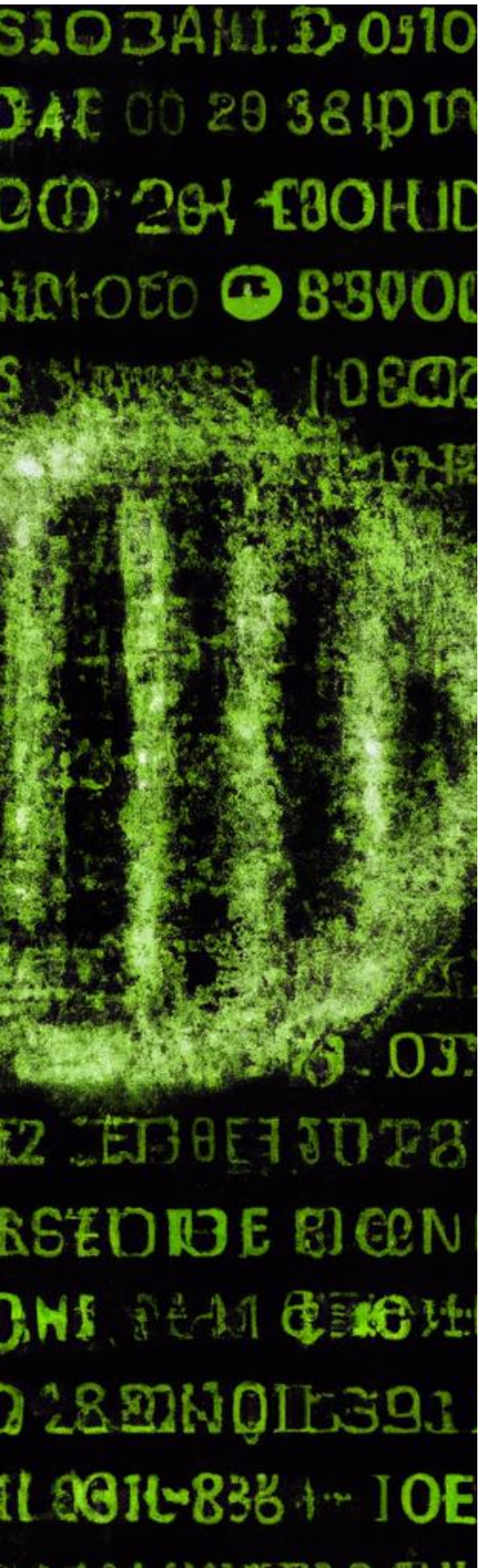
- Simulate different synthesis strategies and biases

- Pandora's box of bias models

ΜΑΝΤΙΚΟΡ



Testing CODEC with real data




Thank you!



**Bionanomanufacturing Unity-
IPT, São Paulo - SP, Brazil**

 bionano@ipt.br

 +55 11 3767 4100

