# Pantheon DNA Data Storage CODEC

Experiences, Challenges, and Innovations

André Guilherme da Costa Martins, PhD Biomed. Sci.

Bioinformatics researcher

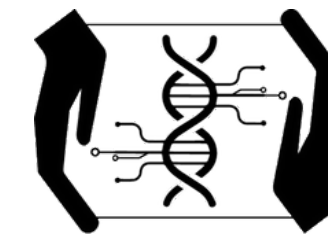Institute for Technological Research - IPT, Brazil

andremartins@ipt.br

# Who are we?

The Institute for Technological Research - IPT has been contributing actively for 124 years to science and technical advances.
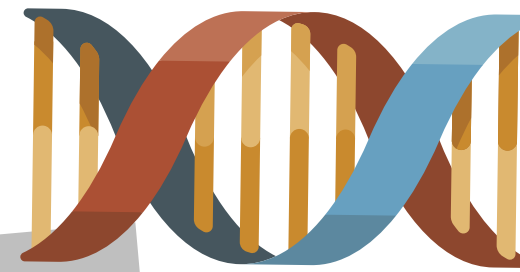
We provide technical solutions for industry, governments, and society, enabling them to overcome the challenges of our time.
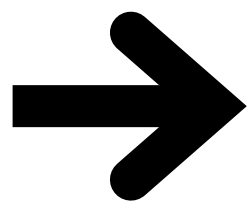
Storaged DNA molecules
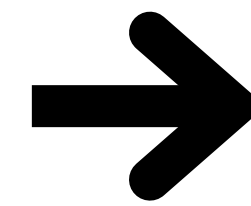
PROMETHEUS

Original Binary data

ATGCTAGCA
AGGCGTGCA
Encoded DNA

DNA Synthesis

DNA Sequencing

ATGCTAGCA
AGGCGTGCA
Sequenced DNA

Recovered Binary data

ENCODING

DECODING

# The Pantheon CODEC

- **A versatile CODEC:**
  - Robust DNA data architecture
    - Binary data pre-processing
    - Multiple choices for mapping algorithms
    - Multi-layer ECC strategy
  - Supports SNIA's sectors (S0 & S1)
  - Includes NGS processing algorithms
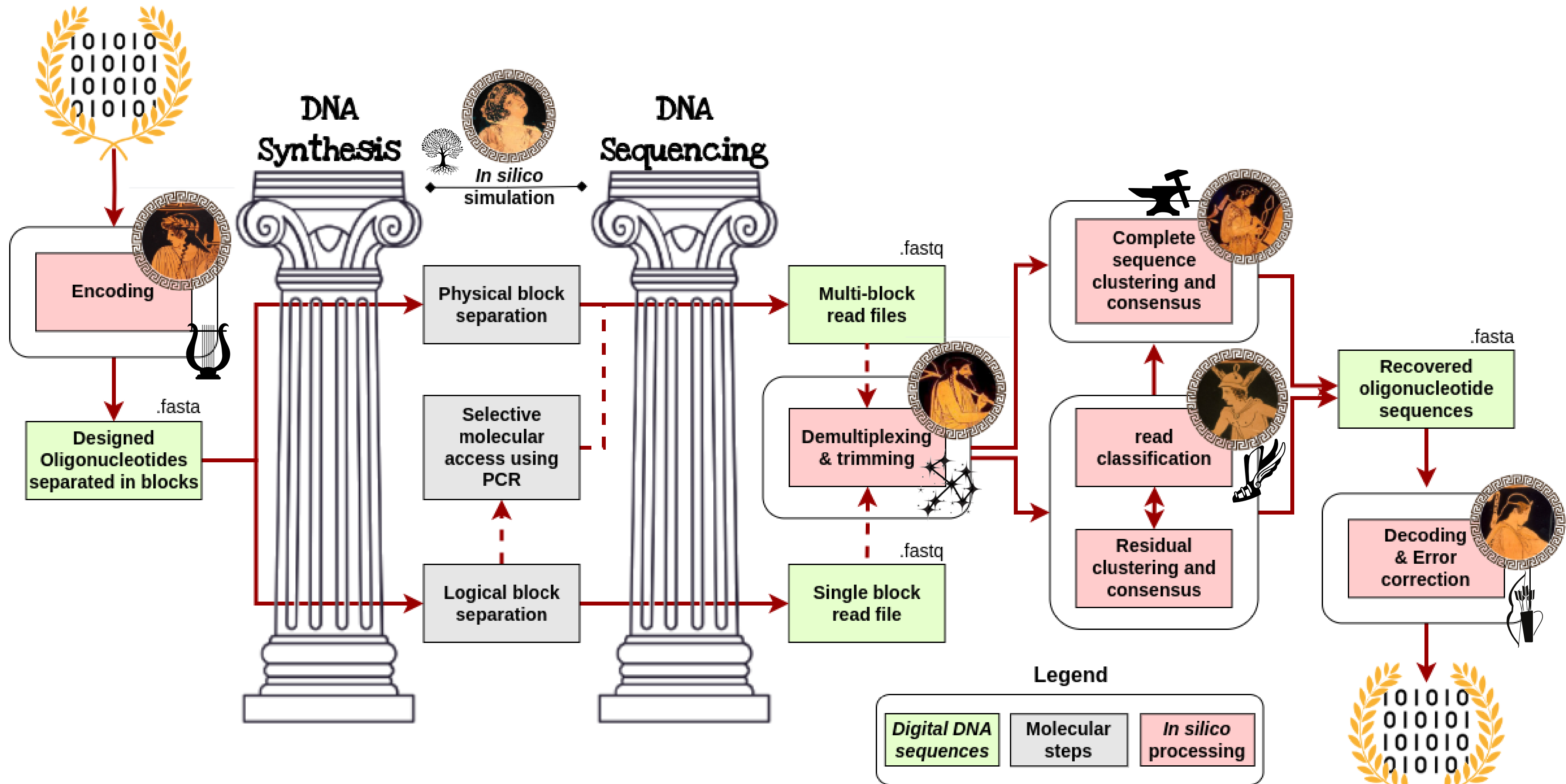    - Compatible with multiple sequencing and storage strategies
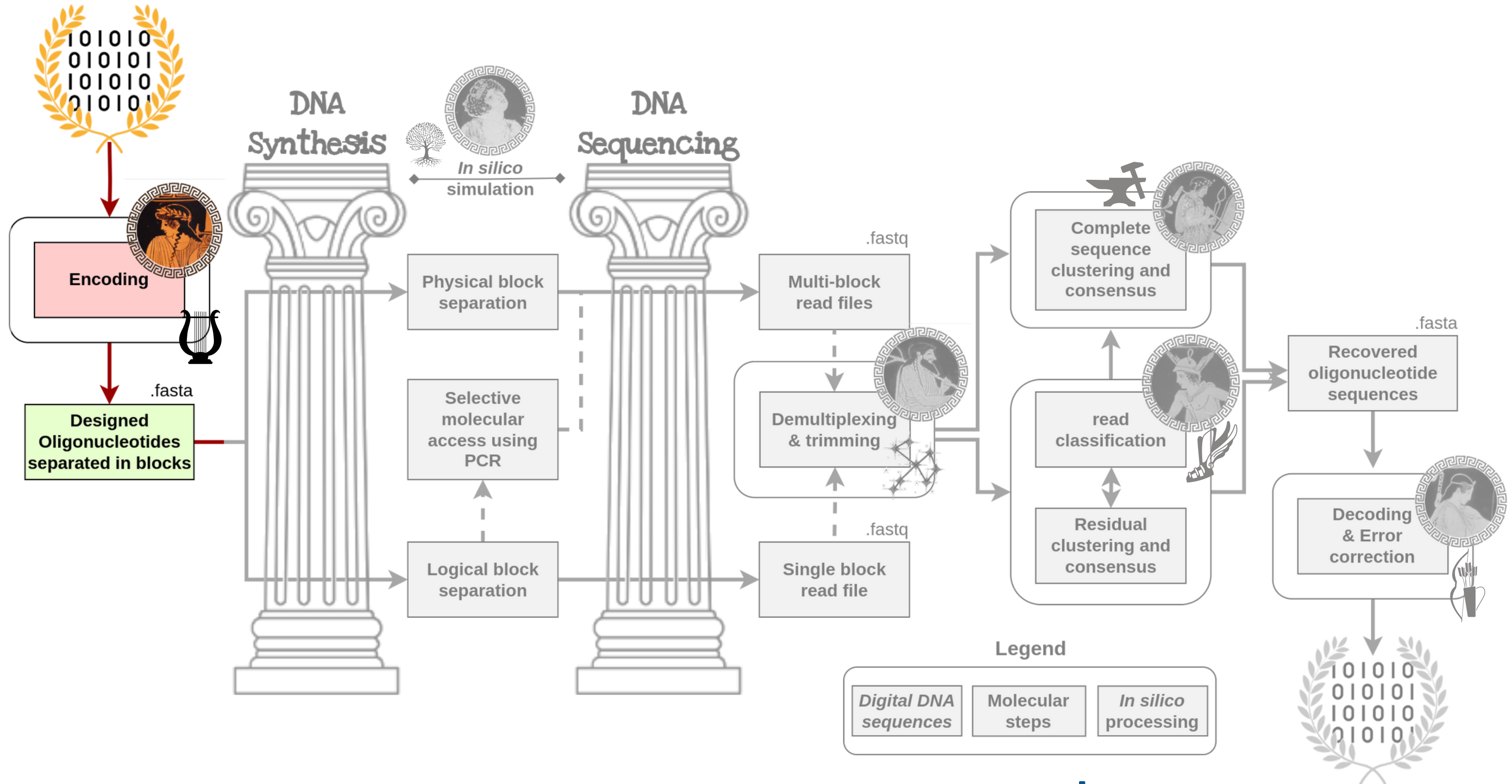
Apollo

Artemis

Chiron

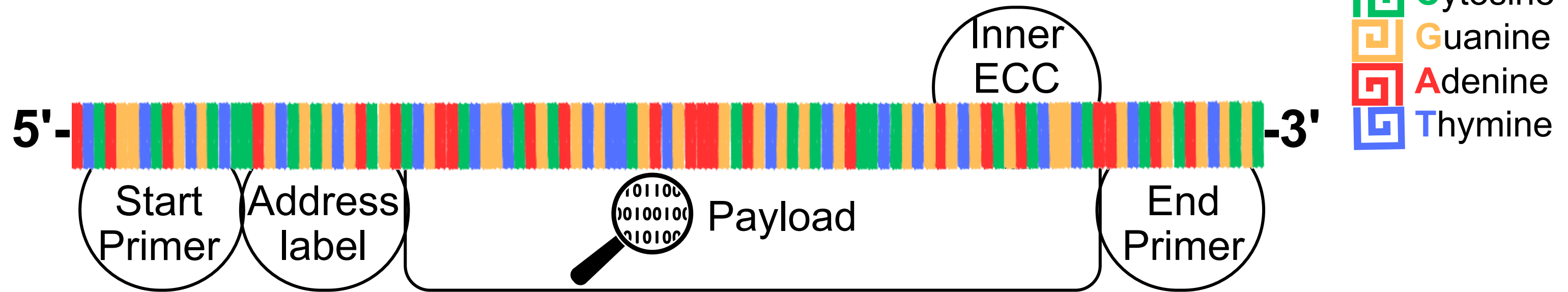Hermes

Hephaestus

Gaia

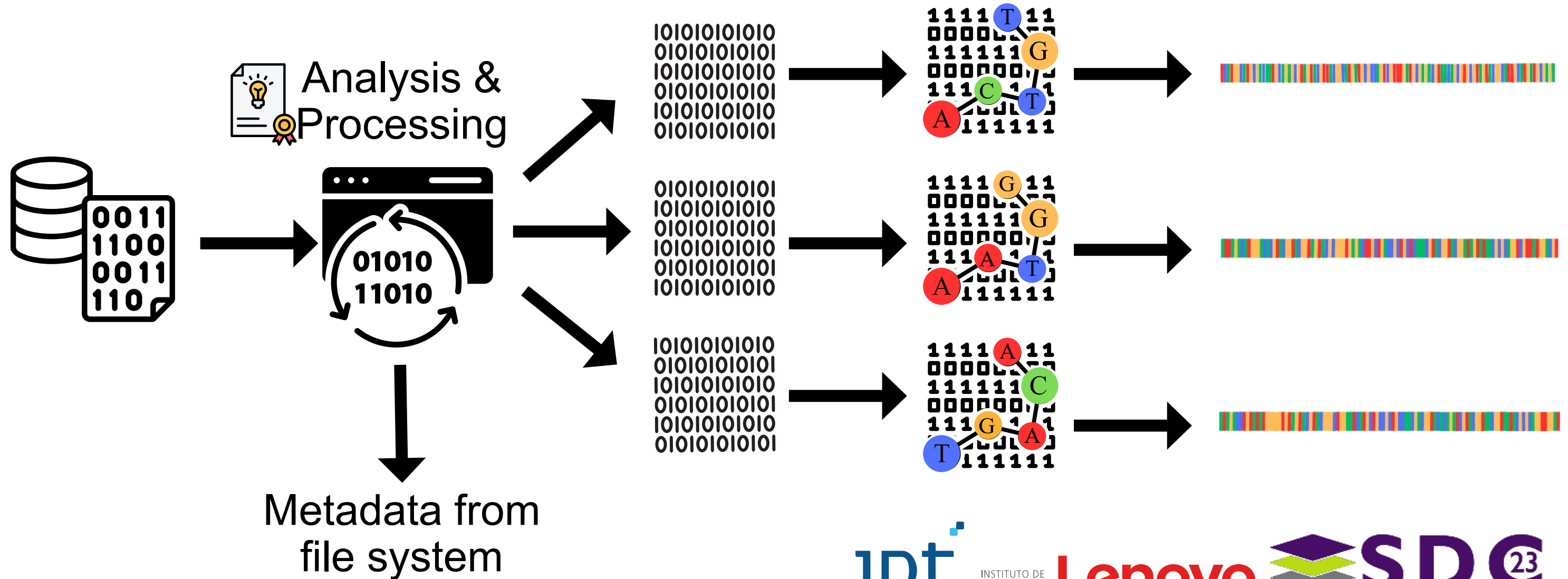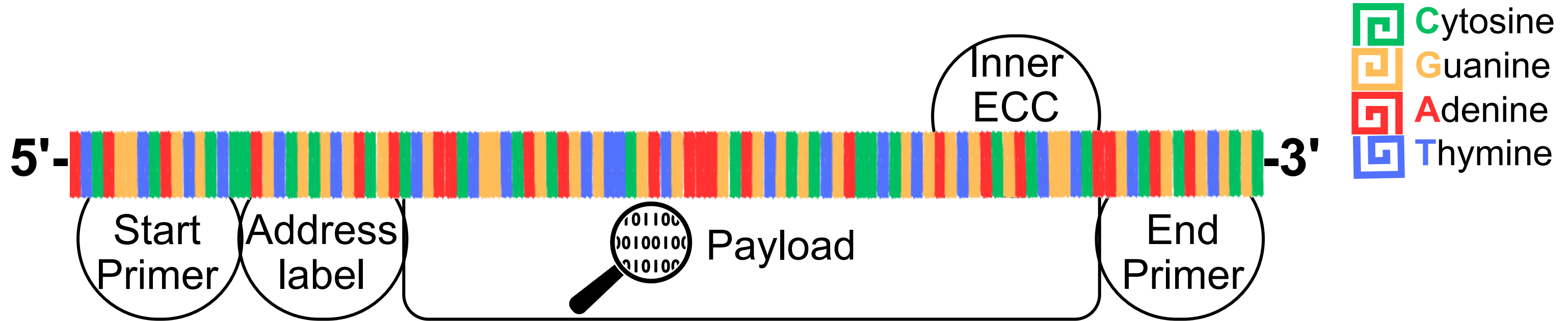PROMETHEUS

# Apollo, the encoding module



Apollo and the Muses - Michel Dorigny, early 1640s

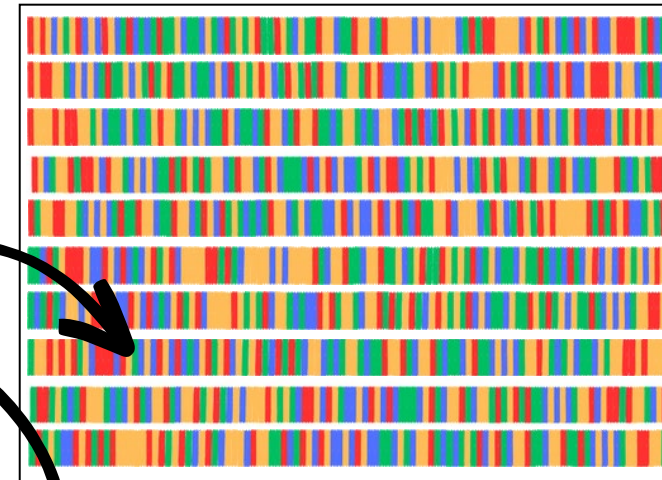# Apollo: oligonucleotide structure

# Apollo: encoding process

# Apollo: metadata & Outer ECC



Metadata
from
file system

**+**

Metadata
from
data blocks

**+**

Codec
parameters
from data
blocks

{...}
JSON

Data blocks

PROMETHEUS

# Apollo: metadata & Outer ECC



Archive Metadata Block (AMB)

Data blocks

Outer ECC

Outer ECC

Outer ECC

# Apollo: primers architecture

# Apollo: DNA blocks architecture

# Apollo: DNA blocks architecture



- **System directory tree**
- **Files checksum**
- **Files coordenates within blocks**
- **CODEC parameters for blocks**
- **Primers data**
- CODEC manual
- Other types of data to assist data recover

# Apollo: DNA blocks architecture

# Apollo: Integration with SNIA's specifications for sectors S0 & S1

**S1 - CODEC parameters**

- General CODEC & ECC parameters
- AMB especific ECC & CODEC parameters
- AMB files checksum



S0 & S1

AMB

Data Blocks

# Chiron, the NGS pre-processing module



The Education of Achilles - Bénigne Gagneraux, 1785

# Chiron: Pre-processing NGS reads

Pre-processing steps:
- Merge read pair (paired-end strategy)
- Adapters/Primer trimming
- Demultiplexing coding blocks
- Reorient DNA sequences
- Discard low-quality reads

PROMETHEUS

IPt INSTITUTO DE PESQUISAS TECNOLÓGICAS  Lenovo.  SD C 23

# Chiron: Pre-processing NGS reads

Pre-processing steps:

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

**C**ytosine
**G**uanine
**A**denine
**T**hymine

# Chiron: Pre-processing NGS reads

Pre-processing steps:

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

**R1**

**R2**



Cytosine
Guanine
Adenine
Thymine

23

# Chiron: Pre-processing NGS reads

Pre-processing steps:
- Merge read pair (paired-end strategy)
- Adapters/Primer trimming
- Demultiplexing coding blocks
- Reorient DNA sequences
- Discard low-quality reads

Block specific prime pairs

Planned data blocks

Sequenced DNA reads

Sequenced & Demultiplexed into data blocks
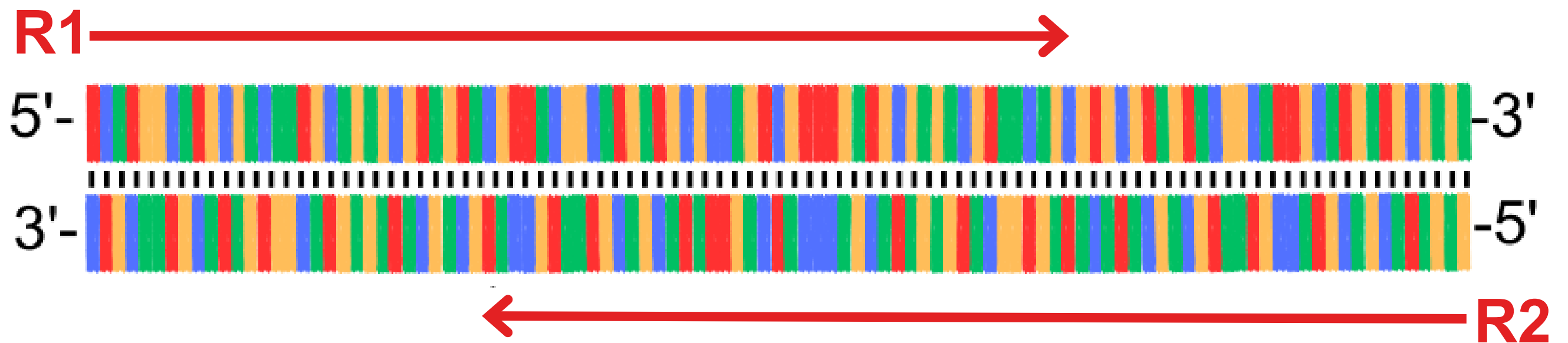
PROMETHEUS

# Hermes, address-oriented module



Mercury (Hermes) - Buti Lodovico, 16th century

# Hermes: address library

# Hermes: parity check



Inner ECC

**C**ytosine
**G**uanine
**A**denine
**T**hymine

5'- ... -3'

Payload

Data block

| Addr. | Mapped data | Inner ECC |

Screenshot of the first 35 DNA sequences from a planned data block in a multi-fasta format. The payload was mapped using a G&A rich scheme.

# Hermes: basic algorithm, part 1

# Hermes: basic algorithm, part 2



e.g. read pile
for addr. #2

FALSE / TRUE — minimum coverage

Multiple sequence alignment (MSA)

consensus sequence

FALSE / TRUE — Length & parity check

Hephaestus' box

Artemis' box

| | Length | Parity |
|---|---|---|
| 5' - ( . . . ) ATCGGCAT- ACTAC( . . . ) - 3' | ✗ | ✗ |
| 5' - ( . . . ) ATCGGCAT- AC- AC( . . . ) - 3' | ✗ | ✗ |
| 5' - ( . . . ) ATC- GCAG- ACTAC( . . . ) - 3' | ✗ | ✗ |
| 5' - ( . . . ) ACCGCCATTACTAC( . . . ) - 3' | ✓ | ✗ |
| 5' - ( . . . ) AACGGCAT- ACTAC( . . . ) - 3' | ✗ | ✗ |

| | Length | Parity |
|---|---|---|
| 5' - ( . . . ) ATCGCCATTACTAC( . . . ) - 3' | ✓ | ✓ |

PROMETHEUS

IPT INSTITUTO DE PESQUISAS TECNOLÓGICAS   Lenovo   SDC 23

# Hephaestus, full-length clustering module



Thetis Receiving the Weapons of Achilles from Hephaestus
Anthony van Dyck, 1632

# Hephaestus: pair-wise clustering and consensus



Demultiplexed
block.fastq

5' - ATCGGCAT - 3'

ATCG
TCGG
k=4 → CGGC
GGCA
GCAT

Shared K-mer strategy to quickly identify and align first highly similar DNA sequences

Cluster 1

Cluster 2

(...)

Cluster n

Demultiplexing speeds up clustering

Consensus 1

Consensus 2

(...

Consensus n

PROMETHEUS

# Hephaestus: pair-wise clustering and consensus



Demultiplexed block.fastq

5' - ATCGGCAT- 3'

k=4

ATCG
TCGG
CGGC
GGCA
GCAT

Shared K-mer strategy to quickly identify and align first highly similar DNA sequences

Cluster 1

Cluster 2

(...)

Cluster n

Strategy adapted from metagenomics studies

Consensus 1

Consensus 2

(...

Consensus n

PROMETHEUS

ipt INSTITUTO DE PESQUISAS TECNOLÓGICAS

Lenovo.

SD©

23

# Hephaestus: basic clustering algorithm

# Artemis, the decoding module



Artemis returning from the hunt - Colombel, 1697

# Artemis: decoding the archive metadata block
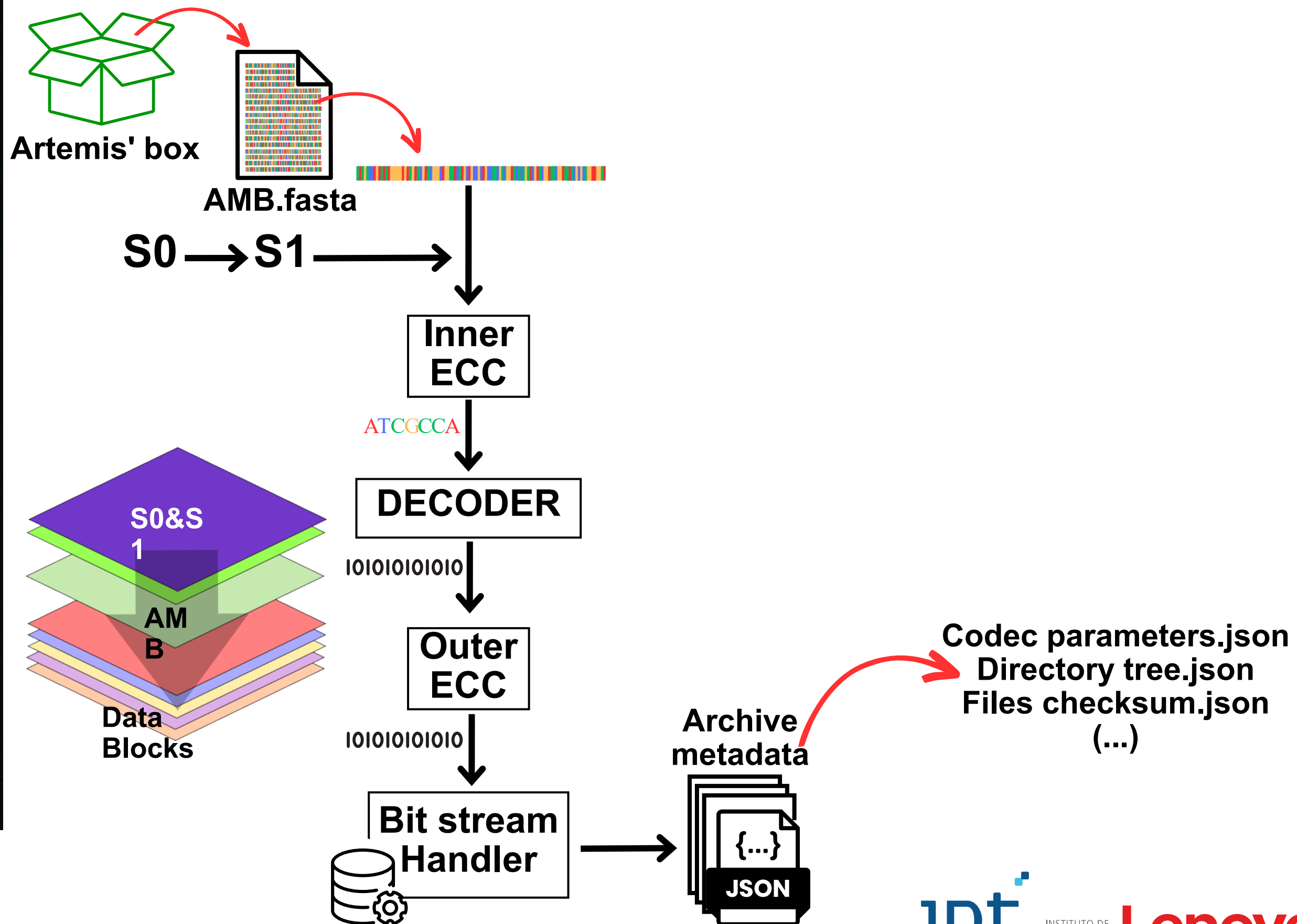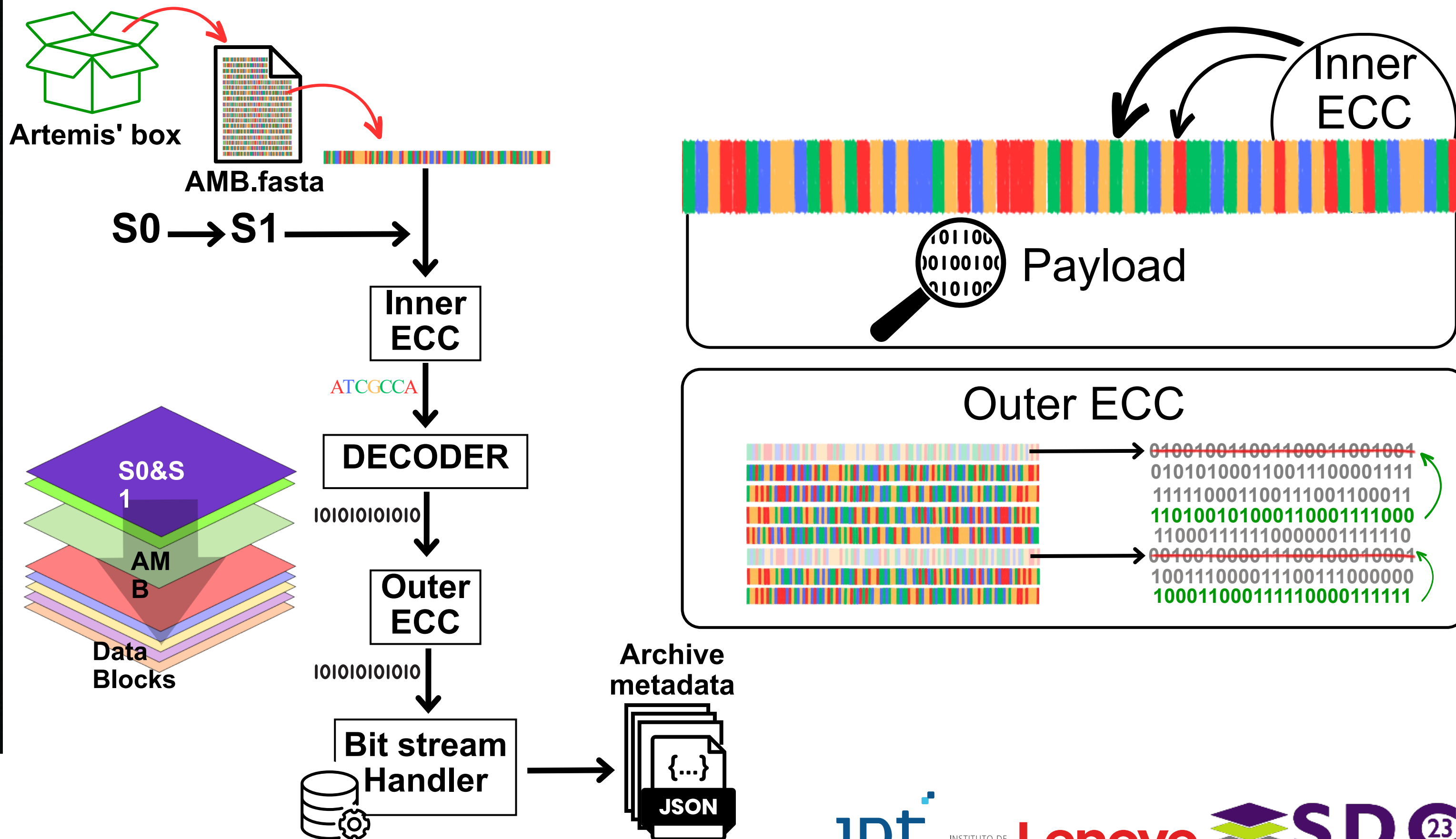
# Artemis: applying ECC to correct errors

# Artemis: Decoding the data blocks

# Gaia, the development module



Tellus Mater (Mother Gaia) panel,
Ara Pacis Augustae - Rome, Italy

# Gaia: a sandbox module to support development

- ## 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

# Gaia: a sandbox module to support development

- ## What Gaia does:
  - Simulate different sequencing strategies
    - Single or Paired-ends
    - Library preparation
    - Coverage variation
    - Sequencing platforms
  - Simulate different synthesis strategies and biases
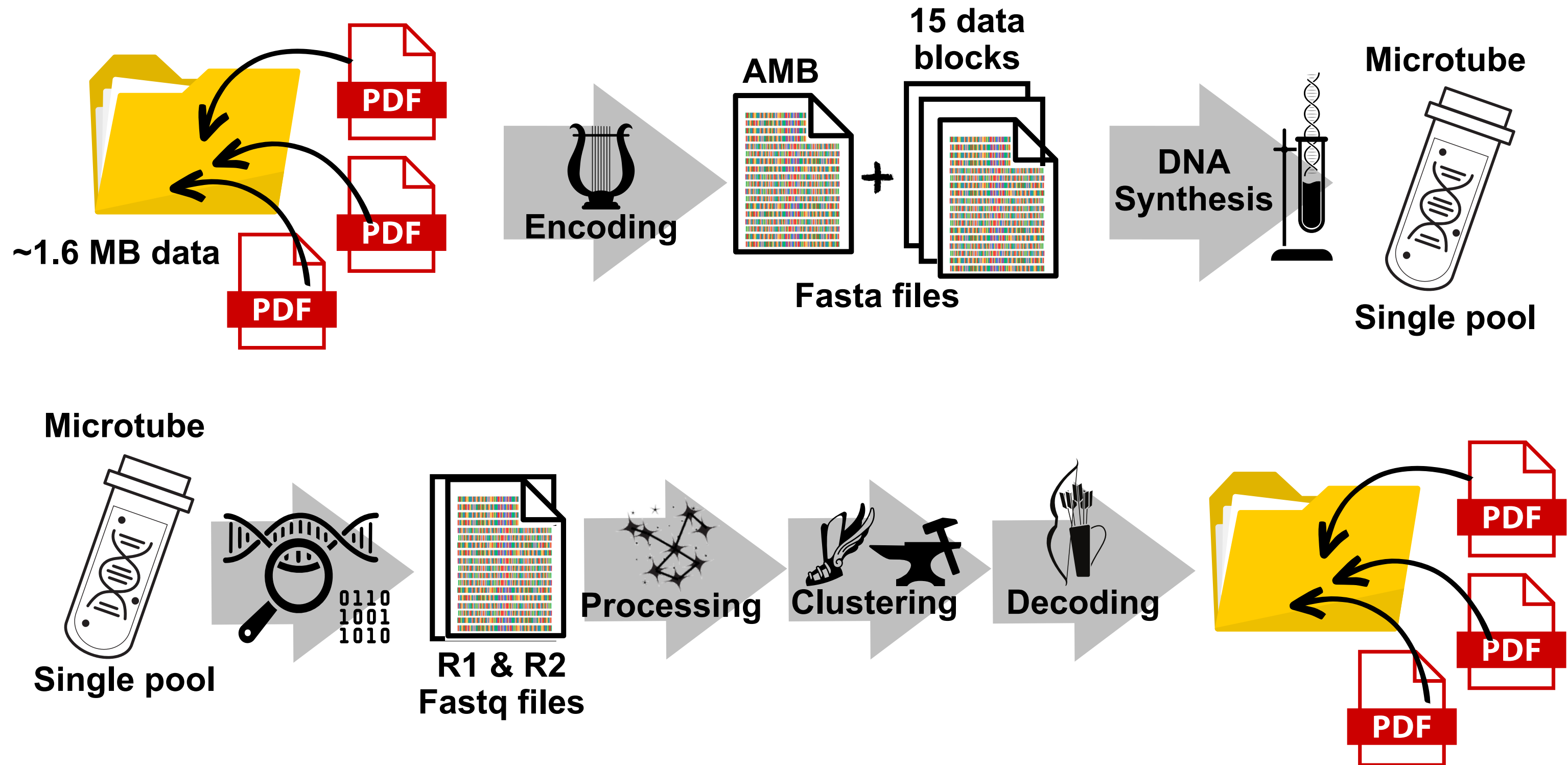    - Pandora's box of bias models

Medusa

Manticore

PROMETHEUS

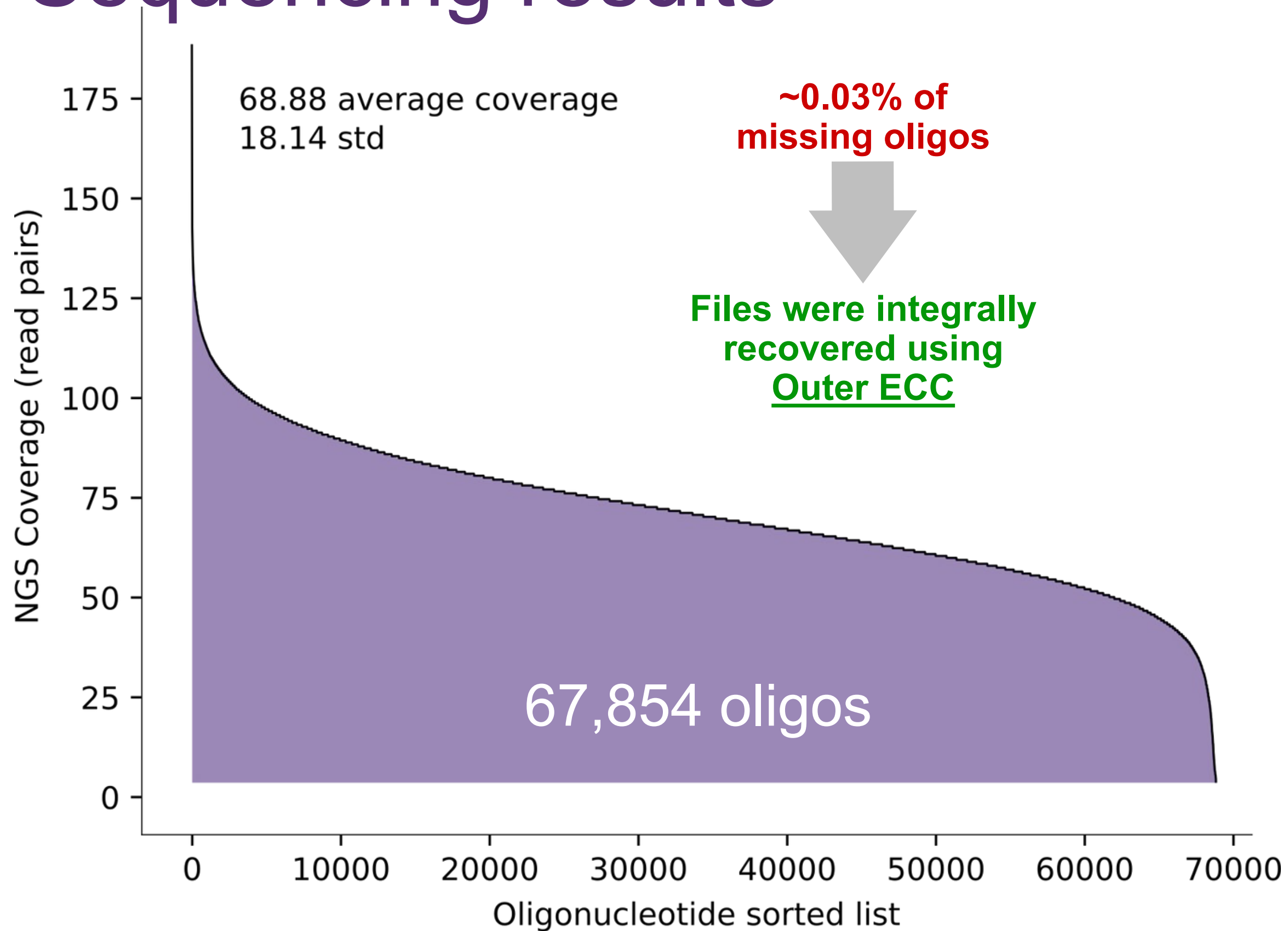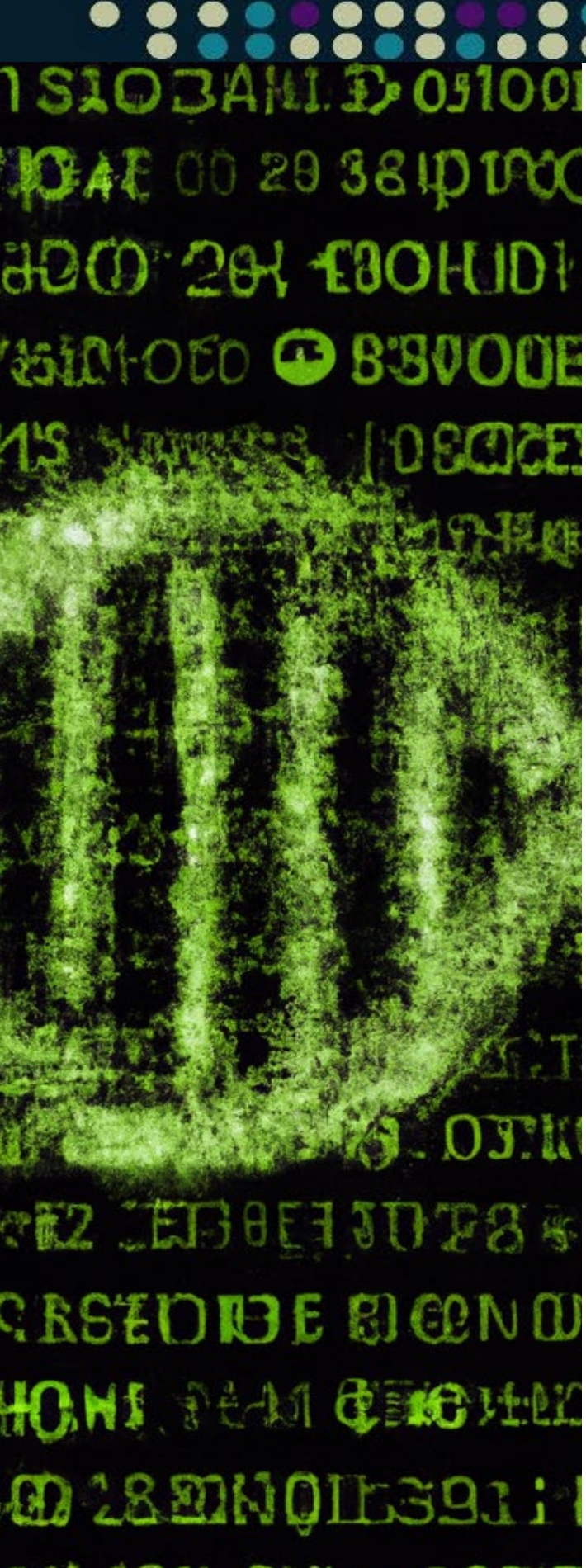ipt INSTITUTO DE PESQUISAS TECNOLÓGICAS  Lenovo  SDC 23

# Testing Phanteon CODEC with real data

# Experiment workflow

# Sequencing results



68.88 average coverage
18.14 std

**~0.03% of missing oligos**

**Files were integrally recovered using Outer ECC**

67,854 oligos

NGS Coverage (read pairs)

Oligonucleotide sorted list

PROMETHEUS

# Sequencing results



68.88 average coverage
18.14 std

**~0.03% of missing oligos**

**Files were integrally recovered using Outer ECC**

67,854 oligos

**162 oligos were low coveraged (≤ 5)**

**22 oligos were missing (coverage = 0)**

NGS Coverage (read pairs)

Oligonucleotide sorted list

PROMETHEUS

# Sequencing results



68.88 average coverage
18.14 std

~0.03% of
missing oligos

Files were integrally
recovered using
Outer ECC

67,854 oligos

"The pit of
Tartarus for
DNA data
storage"

# Sequencing results

68.88 average coverage
18.14 std

**~0.03% of missing oligos**

**Files were integrally recovered using Outer ECC**

**Avoiding unwanted DNA patterns in designed oligos**

"The pit of Tartarus for DNA data storage"

67,854 oligos

NGS Coverage (read pairs)

Oligonucleotide sorted list

PROMETHEUS

# CODEC performance

**00:00:31.490 (AMB + DB)**

Apollo

**00:00:00.200 AMB**

Hermes

**00:01:30.730 AMB**

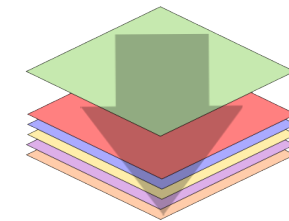Chiron

**00:00:07.800 AMB**

Artemis

**00:00:01.700 AMB**

Hephaestus

4 fastq files with 9.3M reads in total
A single CPU core was used
AMB (Archive metadata block)

PROMETHEUS

# CODEC performance

**00:00:31.490 (AMB + DB)**

Apollo

**00:00:00.590 DB**

Hermes

**00:15:42.000 DB**
**00:00:58.350 DB (fast mode)***

*Increases the sequence loss

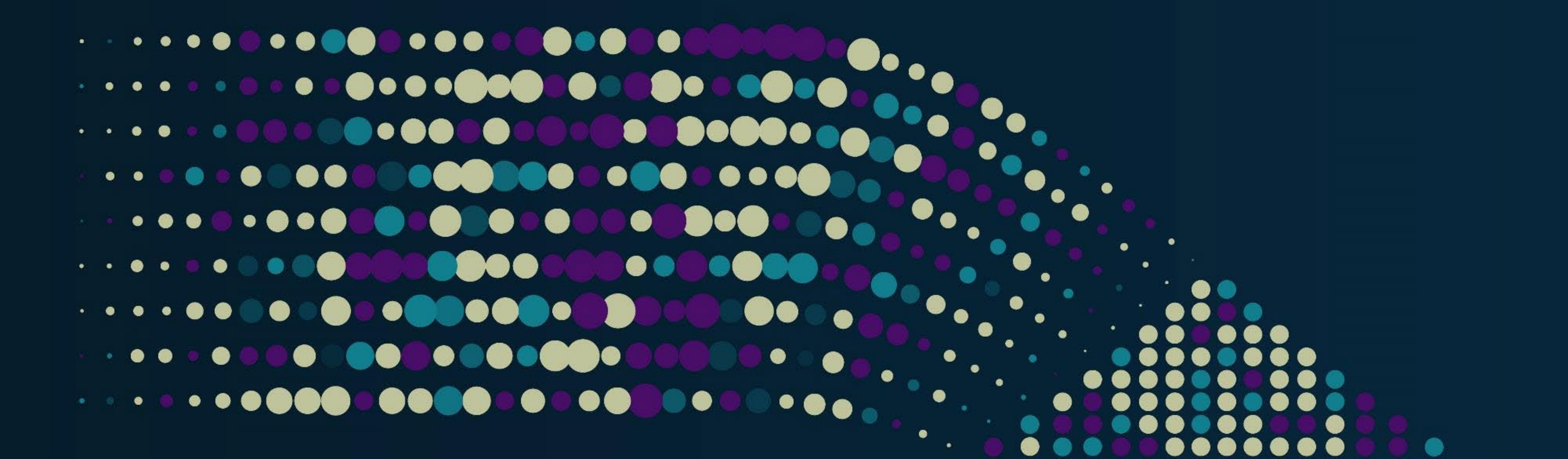Chiron

**00:00:21.040 DB**

Artemis

**00:01:00.000 DB**
**09:11:28.000 DB****

**no prior demultiplexing step

Hephaestus

4 fastq files with 9.3M reads in total
A single CPU core was used
DB (Data Blocks)

PROMETHEUS

# Please take a moment to rate this session.

Your feedback is important to us.

PROMETHEUS

IPT INSTITUTO DE PESQUISAS TECNOLÓGICAS

Lenovo.

SDC 23