

Nº 178455

DNA data storage transmuting digital information into DNA: an innovative odyssey

Adriano Galindo Leal

*Palestra apresentada no Palestra apresentada no InnScidSP,
24/07-04/08/2023. 54 slides.*

A série “Comunicação Técnica” compreende trabalhos elaborados por técnicos do IPT, palestras apresentadas, apresentados em eventos, publicados em revistas especializadas ou quando seu conteúdo apresentar relevância pública. **PROIBIDO A REPRODUÇÃO, APENAS PARA CONSULTA.**



DNA Data Storage

Transmuting Digital Information into DNA:
An Innovative Odyssey

Dr. Adriano Leal

*Senior Researcher
at IPT's Digital Technologies Business Unit
Artificial Intelligence and Analytics Section*

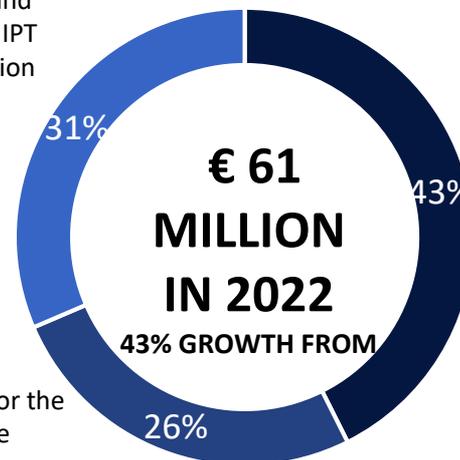
WHO ARE WE ?

IPT PROVIDES TECHNICAL SOLUTIONS FOR INDUSTRY, GOVERNMENTS AND SOCIETY, ENABLING THEM TO OVERCOME THE CHALLENGES OF OUR TIME

INCOMES

Sale of projects and services through IPT Support Foundation (FIPT)

São Paulo State Government basic funding



Sale of projects for the public and private sectors.

OUR NUMBERS (2022)



123 YEARS OF CONTRIBUTIONS TO SOCIETY



> 1,000 EMPLOYEES AND PARTNERS



41% REVENUE IN INNOVATION PROJECTS



> 1,830 CUSTOMERS SERVED



SATISFIED CUSTOMERS NPS 84 (LEVEL OF EXCELLENCE)



> 19,900 TECHNICAL DOCUMENTS ISSUED



> 2,000 TESTING AND ANALYSIS PROCEDURES IN THE PORTFOLIO

WHAT WE DO ?

RESEARCH,
DEVELOPMENT AND
INNOVATION

PRODUCTS AND
PROCESSES
SOFTWARES
FROM THE BENCH TO
THE PILOT
FUNDING
EMBRAPII

TESTS, TRIALS
AND ANALYSIS

TECHNICAL ANALYSIS OF
PRODUCTS AND
MATERIALS
PRODUCT EVALUATION
PRODUCT
CERTIFICATION

INSPECTION AND
MONITORING

CONSTRUCTION AND
STRUCTURES
MACHINERY AND
EQUIPMENT
ACCREDITED
INSPECTION BODY

METROLOGICAL
DEVELOPMENT,
MEASUREMENTS
AND CALIBRATIONS

PROFICIENCY PROGRAMS
STANDARDS
DEVELOPMENT
ADVANCED METROLOGY

CERTIFIED
REFERENCE
MATERIALS

METALS
CERAMIC
MINERAL
VISCOSITY
NORMAL SAND

TECNOLOGICAL
EDUCATION

PROFESSIONAL MASTER
EXTENSION COURSES
COURSES ON DEMAND



BUSINESS UNITS

BIONANOMANUFACTURING

Processes, Chemistry, PPEs, Biotech, Nanotech, Microfabrication

CITIES, INFRASTRUCTURE AND ENVIRONMENT

Territorial planning, Sustainability, Risks, Civil works

ENERGY

Generation, Infrastructure, Efficiency, Clean energy

BUILDING AND HOUSING

Confort, Performance, Safety, Materials, Sustainability

ADVANCED MATERIALS

Metallic, Polymeric, Composite, Cellulosic, Corrosion

DIGITAL TECHNOLOGIES

IoT, Embedded Systems, Intelligent Transport Systems, AI, Analytics

METROLOGICAL AND REGULATORY TECHNOLOGIES

Mechanics, Electrical, Flow Measurement, Aerodynamics, Chemistry



IPT'S HALLMARKS

#innovation



+ 120,000 square meters of laboratories
+ 1,000 qualified professionals
Countless ways to innovate

#quality



+ 2,000 tests and calibrations
+ 20,000 technical documents per year
Reference in quality services

#satisfaction



Level of Excellence in NPS
NPS 84
(Net Promoter Score)



O IPT opens its campus to the largest open innovation action in hardtech in Brazil, connecting distinct stakeholders of this ecosystem.

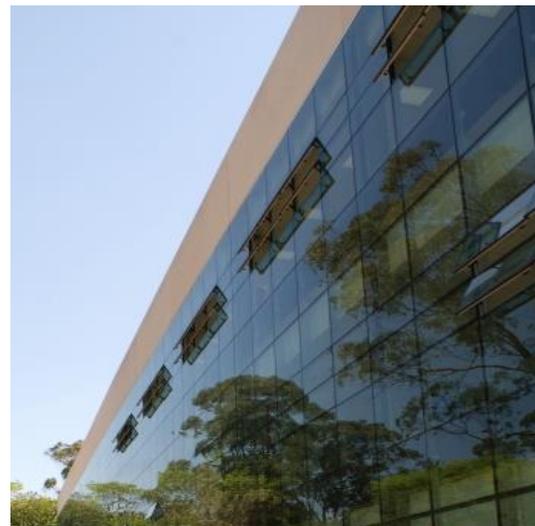
Cornerstone of the CITI Project – São Paulo State International Technology and Innovation Center.



Mode 1

Innovation hub

Become part of a unique and transformative ecosystem that brings together companies and startups that undertake together in the creation of technologies that drive new businesses.



Mode 2

Innovation center

Install your company's Technology Center within the IPT campus and leverage your development capacity.



PROF. DR. ADRIANO LEAL

Adriano Leal (IEEE Member 2006) received a B.Sc. degree in electrical engineering and an M.Sc. and Ph.D. from Polytechnic School at the University of São Paulo in 1996, 1999, and 2006, respectively.

For 11 years, he worked as an R & D engineer for the GAGTD research group in the Polytechnic School at the University of São Paulo. He was responsible for studying and developing automation and information systems for generation, transmission, and electricity distribution.

From April 2007 until November 2010, he has been a Researcher at Elucid Solutions, a consulting and TI company for several utility companies in Brazil.

Since December 2010, he has been a Senior Research Engineer at IPT – Institute for Technological Research. He was IEEE PES South Brazil Chapter President from 01/2015 until 04/2019.

From 05/2019 until 01/2020, he was Visiting Researcher at Technische Universität Berlin / TUB, at the Network Information Theory Group, under Dr. Renato Luís Garrido Cavalcante. Research and training activities were carried out on machine learning techniques to process data from the Internet of Things that comprise the Cyber-Physical Ecosystem of Smart Cities.

On his return, with the support of IPT's top management, he set up the artificial intelligence and analytics section, training researchers on various AI projects and helping assemble the IASMIN Platform.



Let's take flight on our
mental paraglides and
soar over the project's
landscape.

EMBARKING
ON AN
AERIAL
ODYSSEY:
LET'S
BEGIN...



Why DNA Data Storage



- DNA can be stable for thousands of years



- Potential to store Petabytes of data into 1 gram of DNA

Article | [Open Access](#) | [Published: 26 May 2022](#)

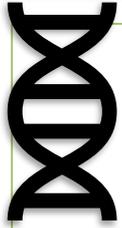
Bioarchaeological and palaeogenomic portrait of two Pompeians that died during the eruption of Vesuvius in 79 AD

[Gabriele Scorrano](#) , [Serena Viva](#), [Thomaz Pinotti](#), [Pier Francesco Fabbri](#) , [Olga Rickards](#) & [Fabio Macciardi](#) 

[Scientific Reports](#) **12**, Article number: 6468 (2022) | [Cite this article](#)

- DNA reading will never be outdated and will always be improved

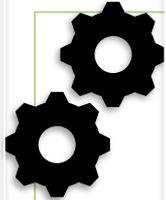
Scope



Enzymatic synthesis



Chemical synthesis



Process
miniaturization



Coding and
decoding



Sequencing
methods

DNA synthesis

	Enzymatic Synthesis	Chemical Synthesis
Benefits	<ul style="list-style-type: none">▪ Fast▪ Environmentally Friendly	<ul style="list-style-type: none">▪ Synthesis Control▪ Standardized Technique
Challenges	<ul style="list-style-type: none">▪ Synthesis Control	<ul style="list-style-type: none">▪ Hazardous chemicals



Working harder together to go further!

The CODEC Team



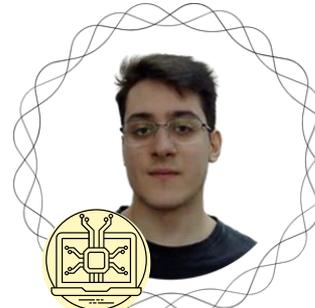
Diego



Alessandro



Adriano



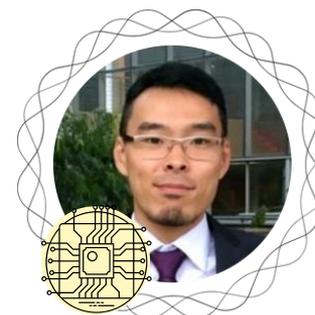
Matheus



Maria Cristina



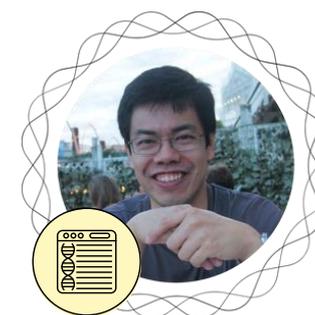
André Martins



Thiago



Takeo



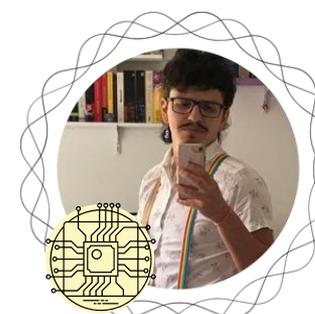
André Fujita



Cristina Maria



Leonardo



Allan



Nyanko

DNA data storage pipeline

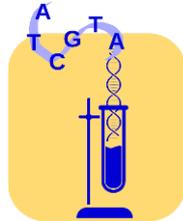
1. Converting information to binary code



2. Converting binary code to DNA code



3. DNA Synthesis



4. Storage



5. Recovery



6. DNA Sequencing



7. Decoding





DNA sequences containing encoded data

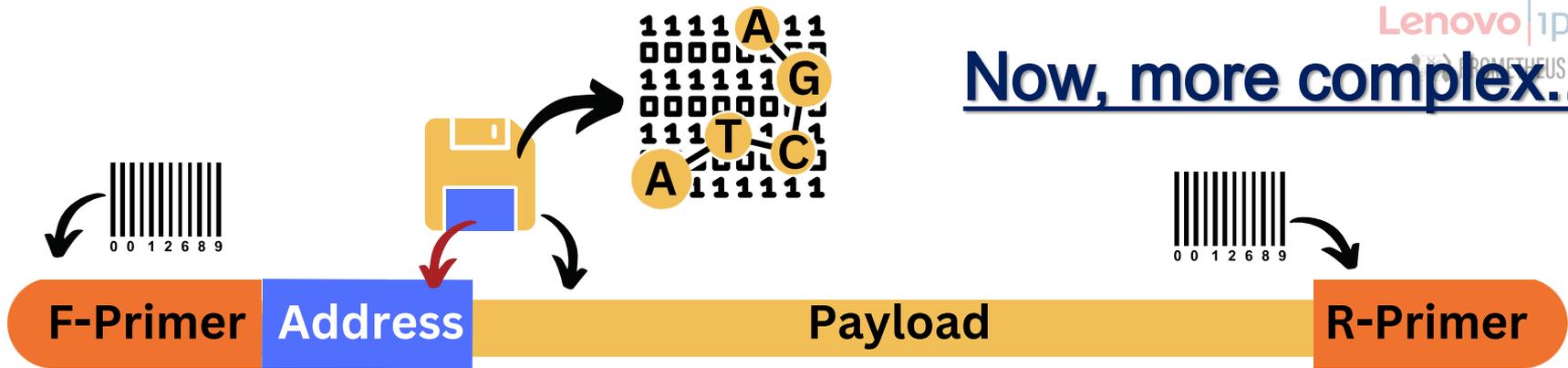
Forward Primer

5'CACGACGTTGTAAAACGACAGACAGGAGAAGCGTACTATATAAGGCCACAGA
CGATAAGGTGCTATCCGGTAGCATGCTGCACGACTATATCGTGTACGGTCACGC
TATATCGCATCACGGGACGCCGGGTCATAGCTGTTTCCTG

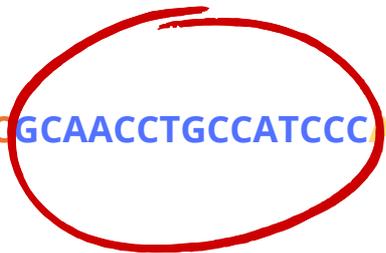
5'CACGACGTTGTAAAACGACTTCGTGGCAGATCAGTCCATAGCCGTCCAGACAA
GAACAGTACGGCCAAGAACATATCGTCCCAGATCCGACCATATCCTCTCTGATA
CGCCATATAACACAGAGGGATGGGTCATAGCTGTTTCCTG

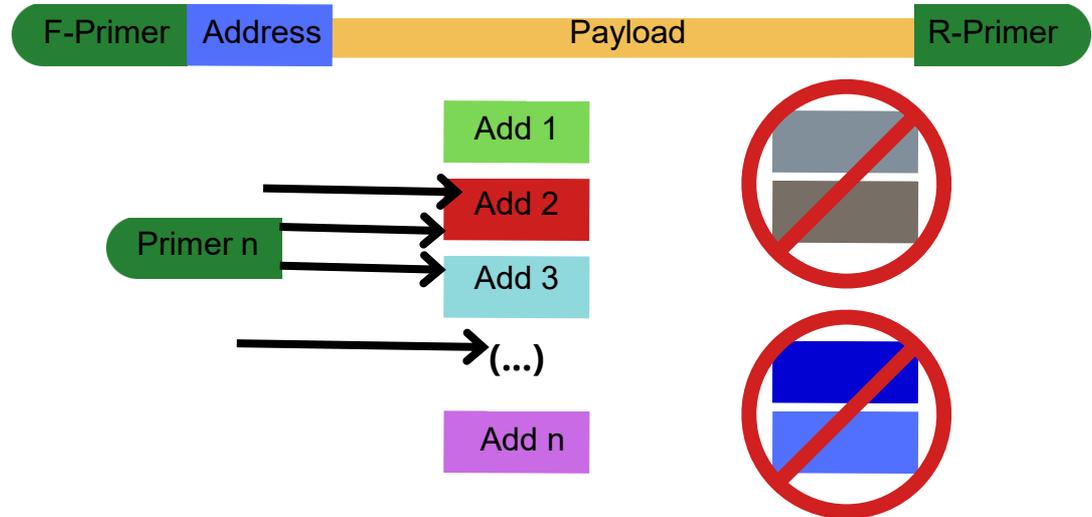
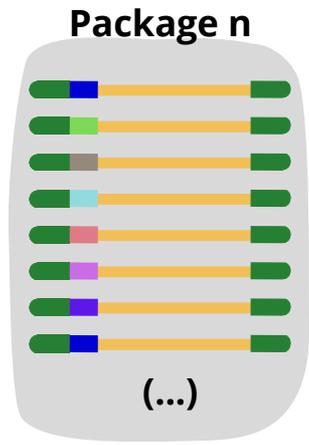
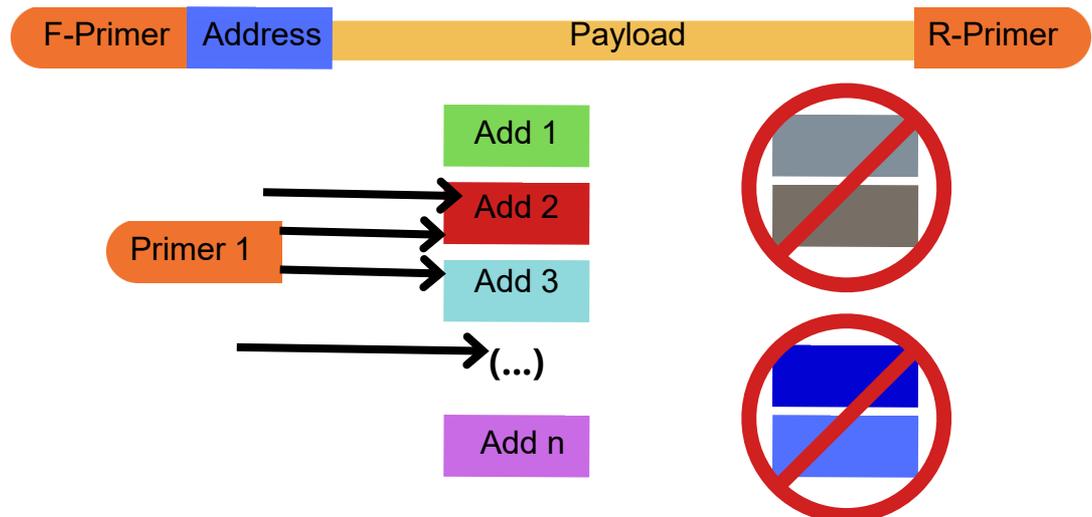
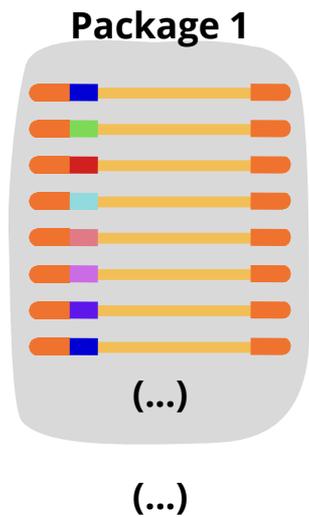
Reverse
Primer

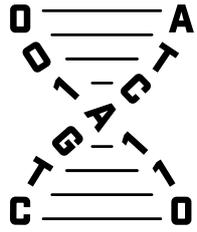
Now, more complex..



5'-(...)TACGCC**GCAACCTGCCATCC**AGCTAAAAGGGGAATCCCAATGGCAATCTGCACCGTCGGACGA(...)3'





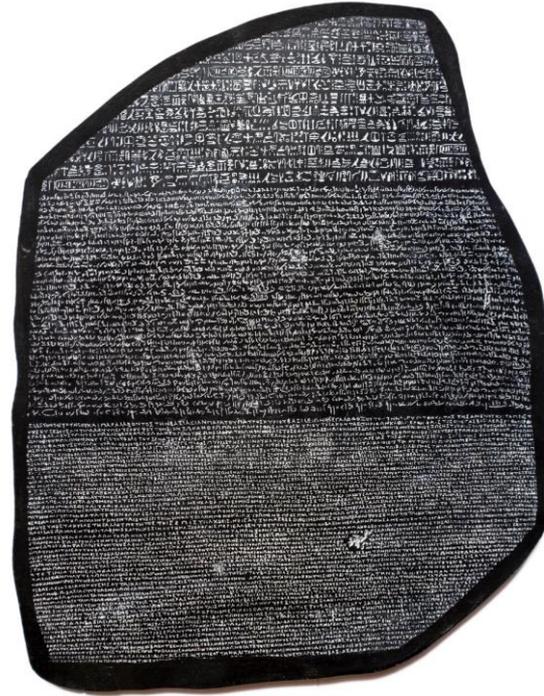


DNA DATA STORAGE ALLIANCE

A SNIA TECHNOLOGY AFFILIATE

Technical Working Group

DNA Archive Rosetta Stone - DARS



"The goal of our group is to define a simple and minimal reserved area in the archive that will function like a "rosetta-stone" and will give the reader/user an idea about how the rest of the archive is encoded."

DNA ZONE

Alliance Codec/Custom

Proprietary Codec

Free text

```
<Sector 1>  
<UID>  
<Description>  
<Codec Params>  
  <Compression>  
  <Encryption>  
  <Map File Primers>  
  <Oligo Size>  
  ...  
<Timestamp>  
<Vendor Details>  
<Hash>
```

```
<Bucket 1>  
<Primers>  
<Contents>  
<Desc>  
  
<Bucket 2>  
<Primers>  
<Contents>  
...
```

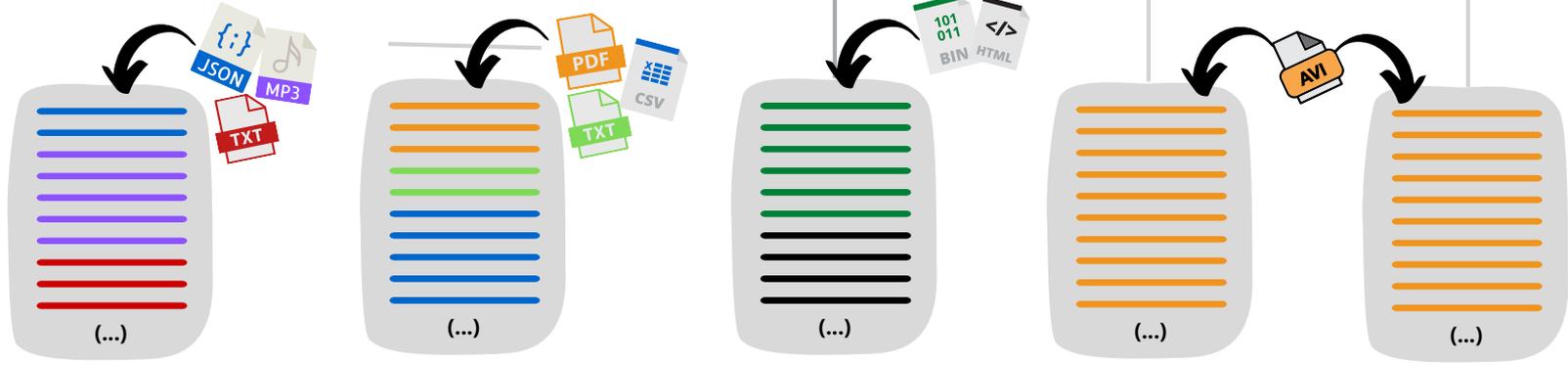
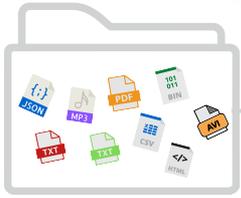
DNA Map file

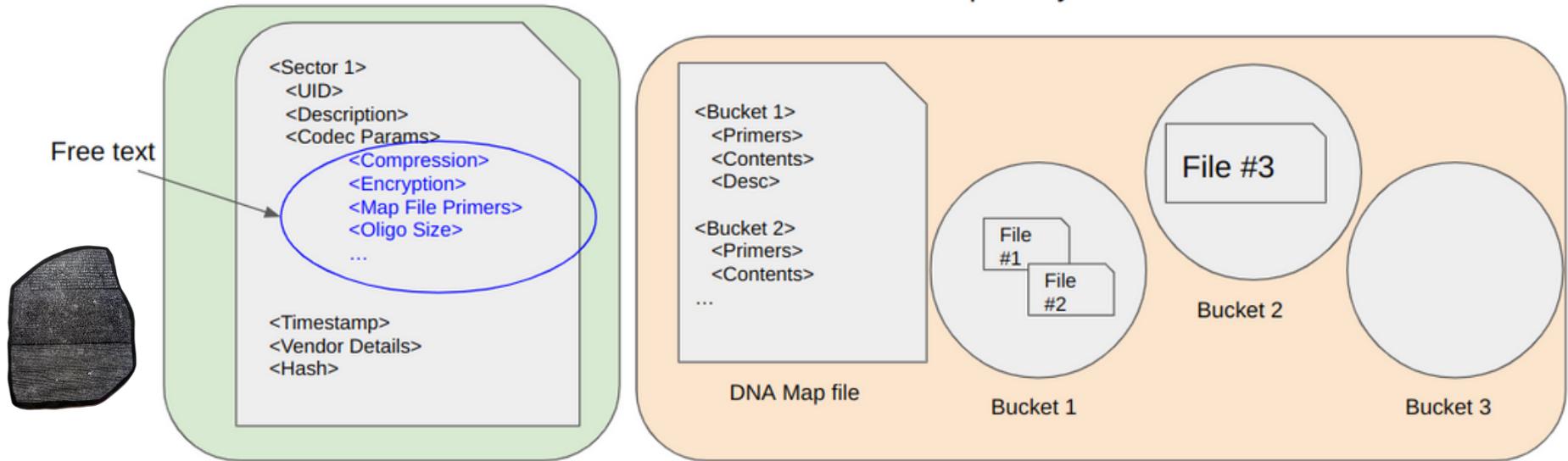
Bucket 1

File #3

Bucket 2

Bucket 3





Sector 0: DNA origin - Company name; codec used

Sector 1 (S1): File Metadata



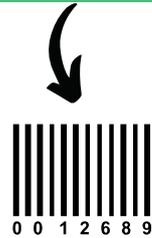
Sector 0: Company name; codec used

F-Primer

vendor

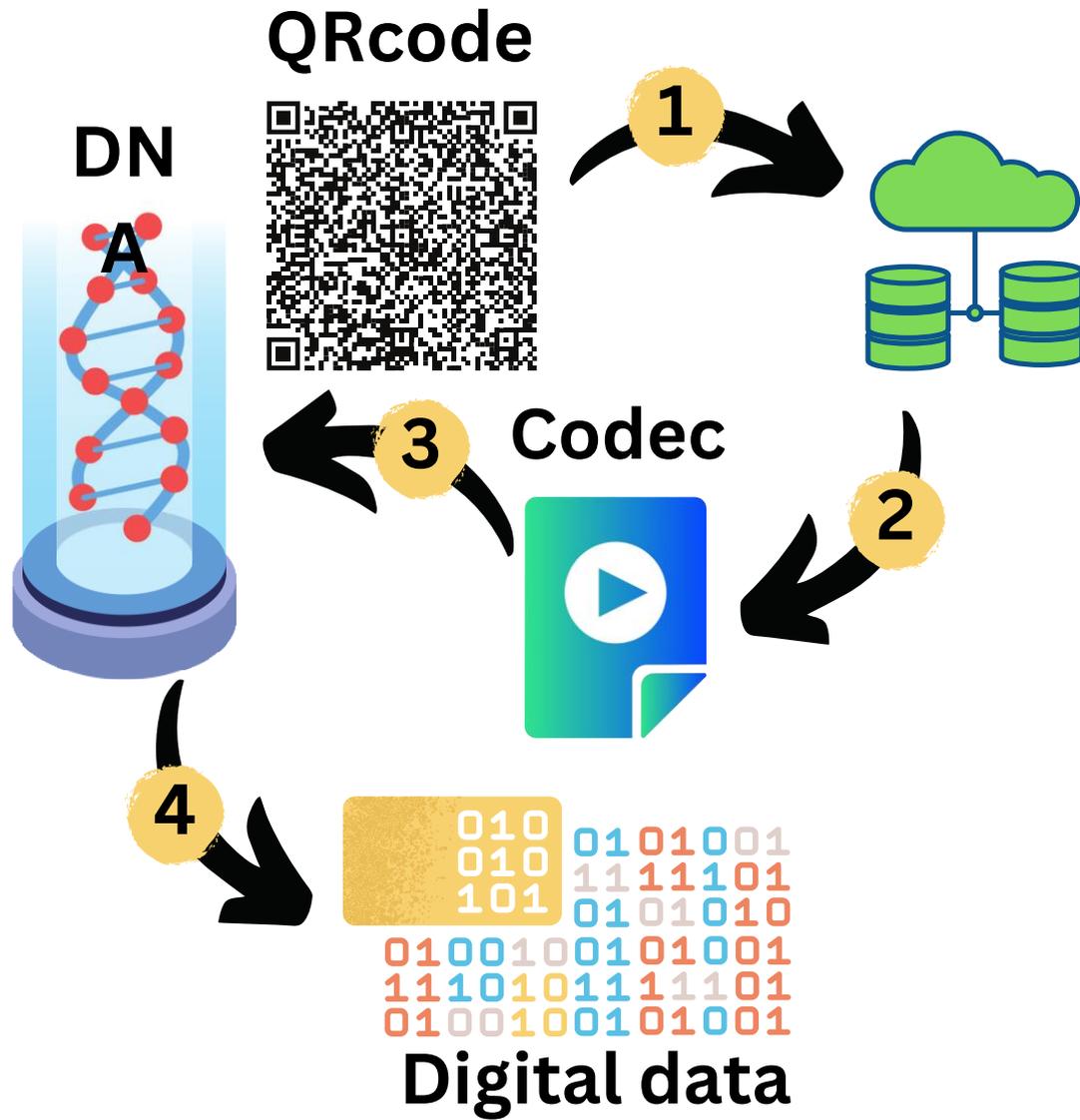
Codec

R-Primer

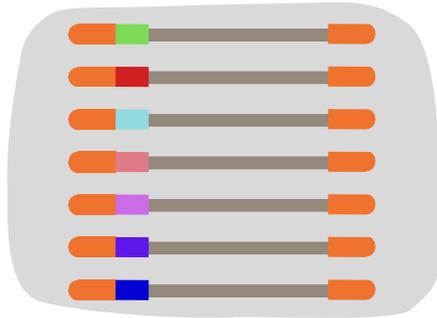


ID	Vendor Name	Contact Details	Notes
ATGCGATGCGATGCGATGC GATGCGATGCGATGCG	Twist Bioscience	twistbioscience.com/get_codec +14155307827	

ID	Vendor Name	Codec Details	Notes
ATGCGATGCGATGCGATGC GATGCGATGCGATGCG		Version 1.4	

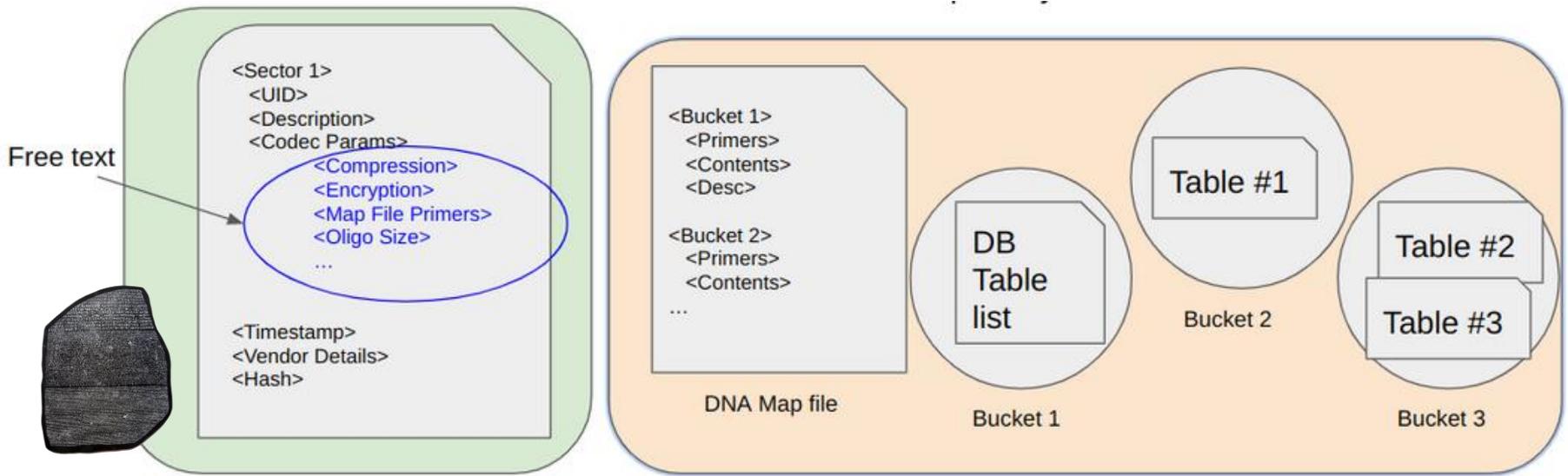


Sector 1 (S1): File Metadata

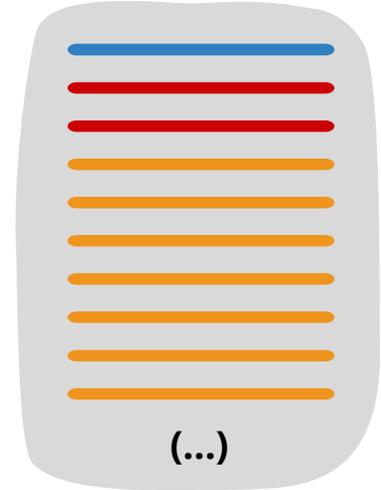
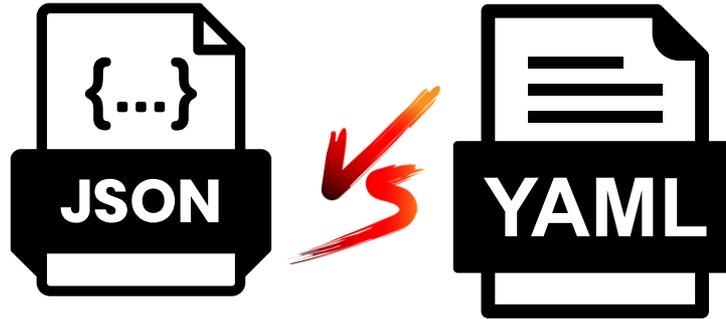
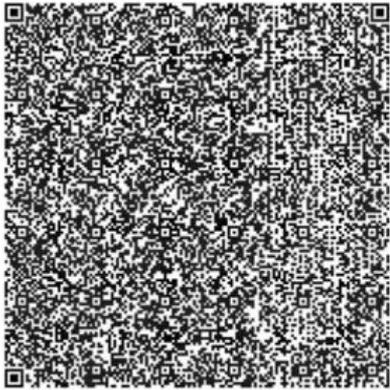


Emergency rescue

- File data
 - List os files
 - file sizes
 - Creation date
 - MD5 Checksum
 - User owner
 - File location (package)
- CODEC
 - Encoding parameters
 - Primer data
 - Number of oligos



Sector 1 (S1): File Metadata

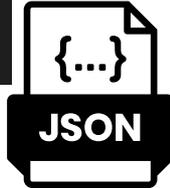


QR codes can hold a maximum of 3KB of data (177x177 nodes) according to <https://dhiway.com/a-definitive-guide-to-qr-code/#:~:text=A%20QR%20code%20has%20a,characters%20or%204%2C269%20alphanumeric%20ones.>

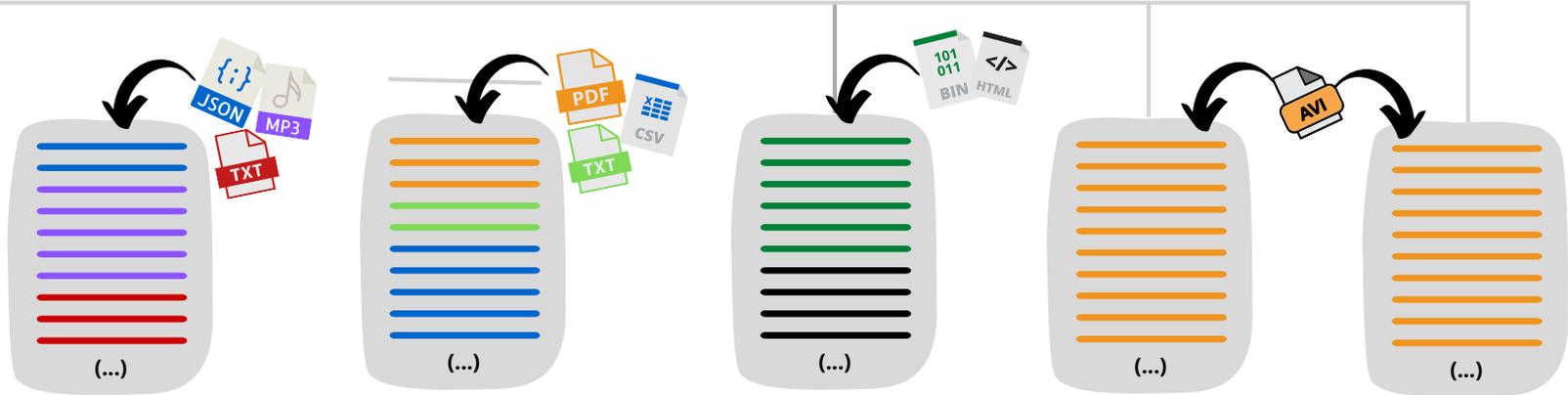
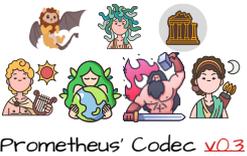
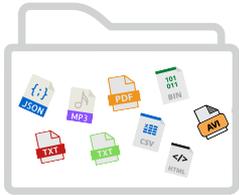
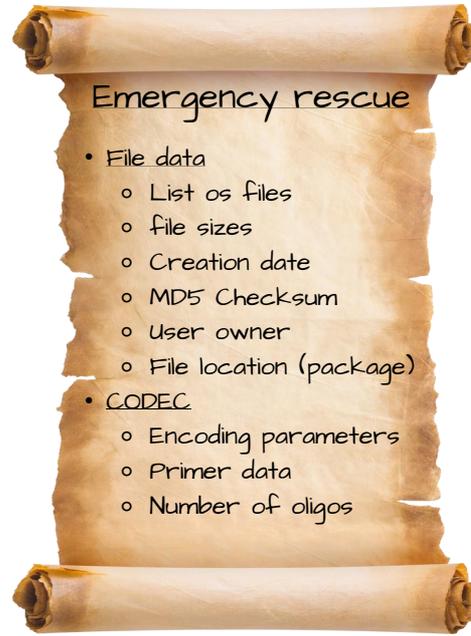
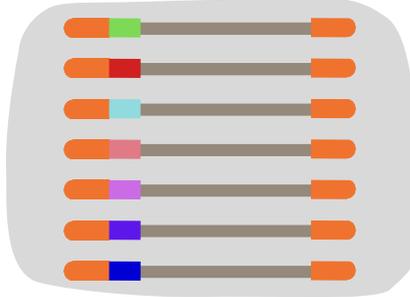
```

codecsn X
C: > Users > thiagoayagi > Desktop > codecsn > ...
1
2 "address_method": "baselevel_ordered",
3 "inner_ecc_method": "ReedSolomon",
4 "interleaving": true,
5 "mapped_oligo_size": 100,
6 "mapping_method": "direct",
7 "max_block_size": 5100,
8 "outer_ecc_method": "ReedSolomon",
9 "primer_cut": 3,
10 "primer_list": [
11 "CACGACGTTGTAGGACGAC",
12 "CAGGGAGCAGCTATGACC",
13 "GGTTGGCCCACTCAGCAG",
14 "GGAGGCAGCTATGACCATG",
15 "TGTAAGACGACGGCCAGTG",
16 "AGCGGATAACAATGGCACC",
17 "TCCCGGACTCACTATAGG",
18 "ATGTCGGGCTAACCCGCC",
19 "GCTAGTTATTGCTCAGCGG",
20 "TAGTTATTGCTCAGCGGTG",
21 ],
22 "rand_type": "None",
23 "segment_size": 216,
24 "xor_kernel": "kbf1f2"
25

```

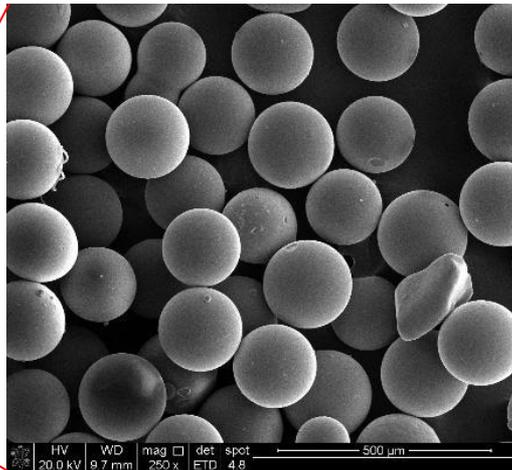
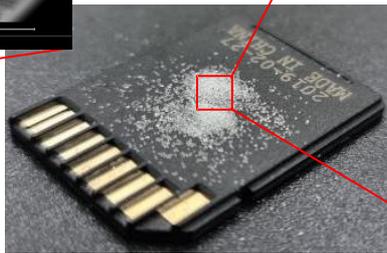
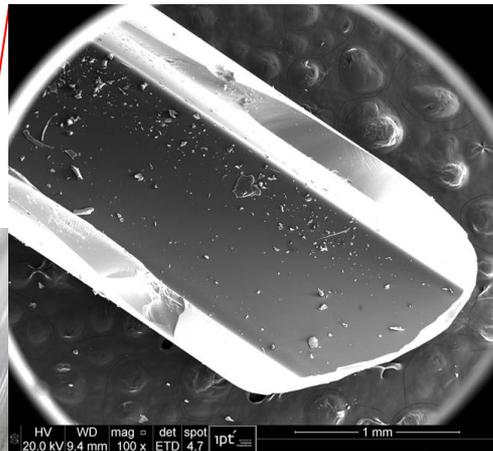
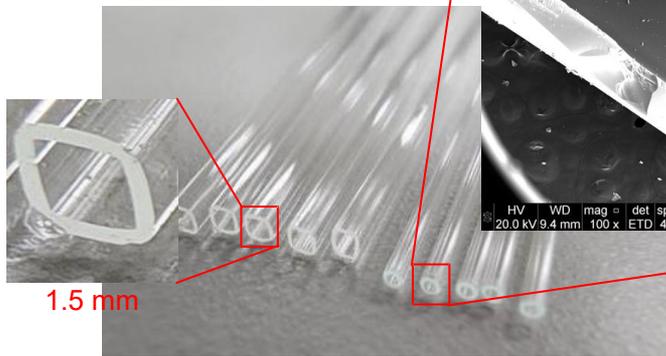


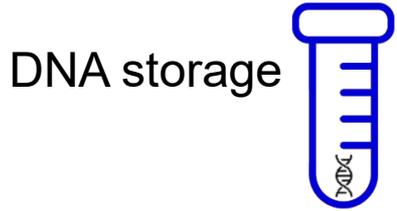
Package 0 Metadata



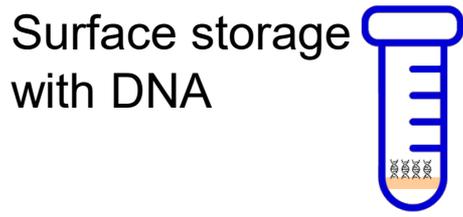


Chemical DNA synthesis on different surfaces





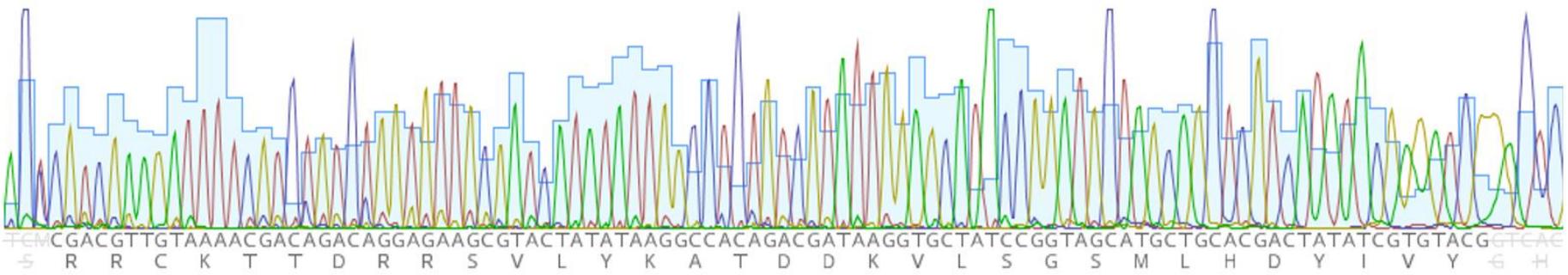
Creating copies by PCR methodology
PCR = Polymerase Chain Reactions



Sequencing



Sequencing of synthesized DNA



5'CACGACGTTGTAAAACGACAGACAGGAGAAGCGTACTATATAAGGCCACAGACGATAAGGTGCT
ATCCGGTAGCATGCTGCACGACTATATCGTGTACGGTCACGCTATATCGCATCACGGGACGCCGG
GTCATAGCTGTTTCCTG

5'CACGACGTTGTAAAACGACTCGTGGCAGATCAGTCCATAGCCGTCCAGACAAGAACAGTACGG
CCAAGAACATATCGTCCCAGATCCGACCATATCCTCTCTGATACGCCATATAACACAGAGGGATGG
GTCATAGCTGTTTCCTG





 ATCTGCA

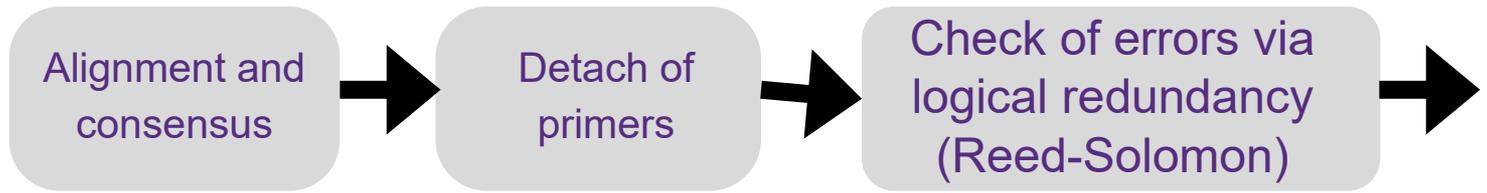
 ATCTGCG

 ATCAGCA



```

ATCTGCA
ATCTGCG
ATCAGCA
ATCTGCA
ATCTGTA
ATATGCA
  
```



```

primer  |-----|
read1  TC*CGACGTTGTA AAAACGACAGACAGGAGAAGCGTACTATATAAGGCCACAGACGATAAGGTGCTATCCGGTAGCATGCTGCAGACTATATCGTGTACGGTCACGC*--**C*---*ACG-----|-----|
read2  -----C*A-----*****TAAGGC-*CAGACGATAAGGTGCTATCCGGTAGCATGCTGCAGACTATATCGTGTACGGTCACGCTATATCGCATACAGGGACGCCGGGTCATAGCTGTTTC**GA
consensus C*CGACGTTGTA AAAACGACAGACAGGAGAAGCGTACTATATAAGGCCACAGACGATAAGGTGCTATCCGGTAGCATGCTGCAGACTATATCGTGTACGGTCACGCTATATCGCATACAGGGACGCCGGGTCATAGCTGTTTC**G
original CACGACGTTGTA AAAACGACAGACAGGAGAAGCGTACTATATAAGGCCACAGACGATAAGGTGCTATCCGGTAGCATGCTGCAGACTATATCGTGTACGGTCACGCTATATCGCATACAGGGACGCCGGGTCATAGCTGTTTCCTG
errors  x                                                                                                                                           xx
  
```

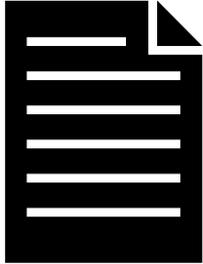
```

primer  |-----|
read1  TC*CGACGTTGTA AAAACGACTCGTGGCAGATCAGTCCATAGCCGTCCAGACAAGAACAGTACGGCCAAGAACATATCGTCCCAGATCCGACCAT*-CCTCTC-*A**--GC*****A-----|-----|
read2  -----*****G---T---CGT-CAGAC-AGAACAGTACGGCCAAGAACATATCGTCCCAGATCCGACCATATCCTCTCTGATACGCCATATAACACA*AGGGATGGGTCATAG*TGTTTC**GA
consensus C*CGACGTTGTA AAAACGACTCGTGGCAGATCAGTCCATAGCCGTCCAGACAAGAACAGTACGGCCAAGAACATATCGTCCCAGATCCGACCATATCCTCTCTGATACGCCATATAACACA*AGGGATGGGTCATAG*TGTTTC**G
original CACGACGTTGTA AAAACGACTCGTGGCAGATCAGTCCATAGCCGTCCAGACAAGAACAGTACGGCCAAGAACATATCGTCCCAGATCCGACCATATCCTCTCTGATACGCCATATAACACAAGGGATGGGTCATAGCTGTTTCCTG
errors  x                                                                                                                                           x                                                                                                                                           xx
  
```

Fasta



Fastq

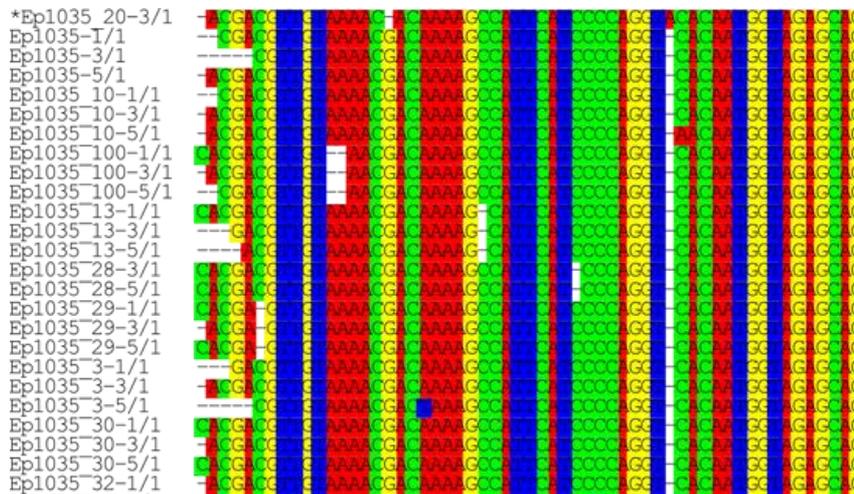
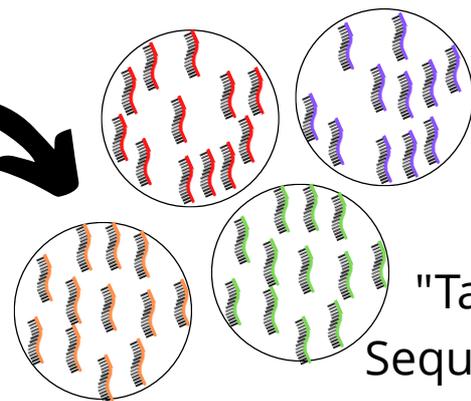
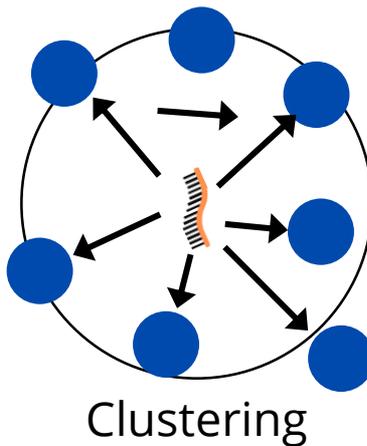
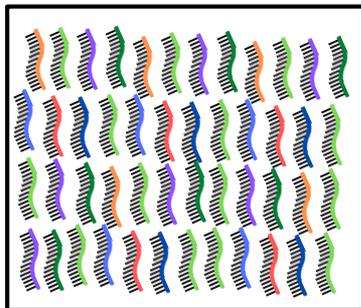


 Analogy!



Review of Post-Sequencing Processing

Fastq
(reads)



Alignment

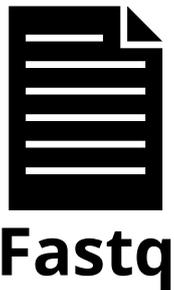


Fastq



Pre-processing
sequenced data 

Sequence clustering



Let's check the pile
for the page 15



 **Warning!** Analogy!

• Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam finibus vitae ____ ac blandit. Cras eget bibendum tortor, scelerisque pharetra leo. Suspendisse vestibulum odio sapien, id ultrices nibh ultricies finibus. Integer fermentum mauris **et** ultrices commod**a**.

• Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam finibus vitae elit ac blandit. Cras eget bibendum **r**tortor, scelerisque p_ aretra leo. Suspendisse vestibulum odio sapien, id ultrices nibh ultricies finibus. Integer fermentum mauris **ut** ultrices commod**o**.

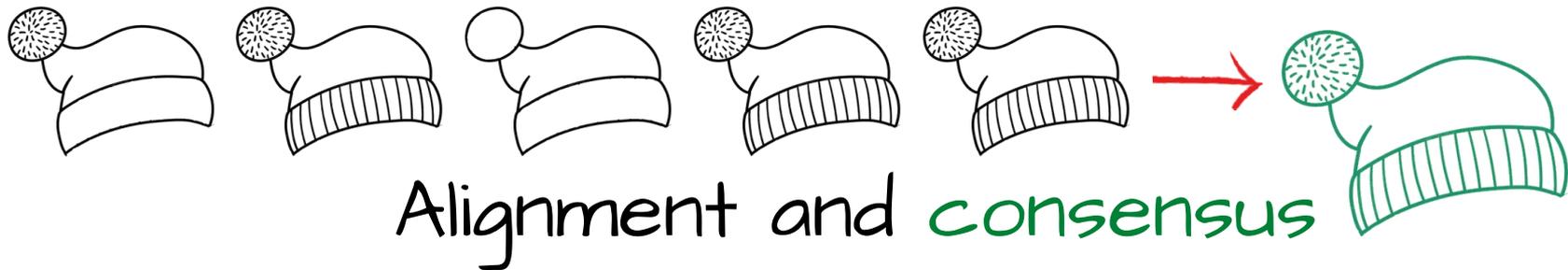
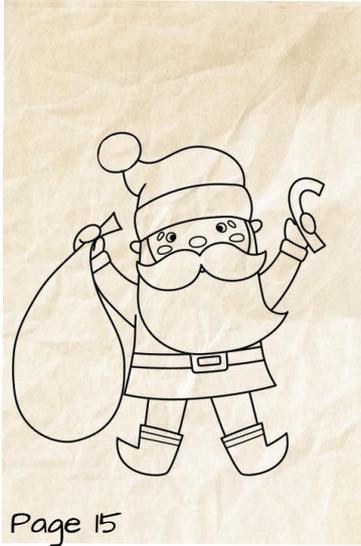
• Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam finibus vitae elit ac blandit. Cras eget bibendum **t**tortor, scelerisque pharetra **l??**. Suspendisse vestibulum odio sapien, id ultrices nibh ultricies finibus. Integer fermentum mauris **at** ultrices commod**u**.

- Identify similar pages
- group them
- make a consensus
- repeat n times

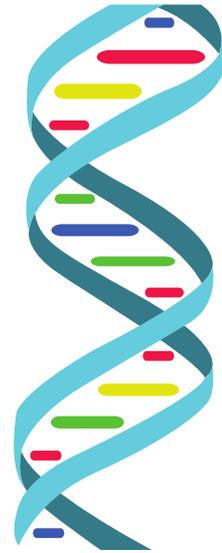


! Analogy!

Which one is the correct version of Santa Claus?



For DNA storage, the challenge is similar, find and group the right sequence and then, get consensus, but repeating these millions of times

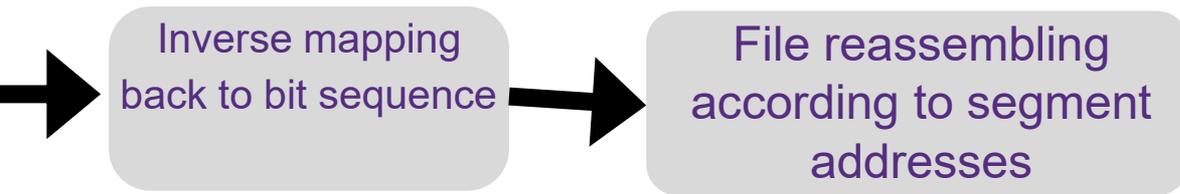


If any error remains in the consensus, Artemis will attempt to fix it using error correction methods such as Reed-Solomon or Hamming code.



>>> DECODING >>>





0101
1001
0110

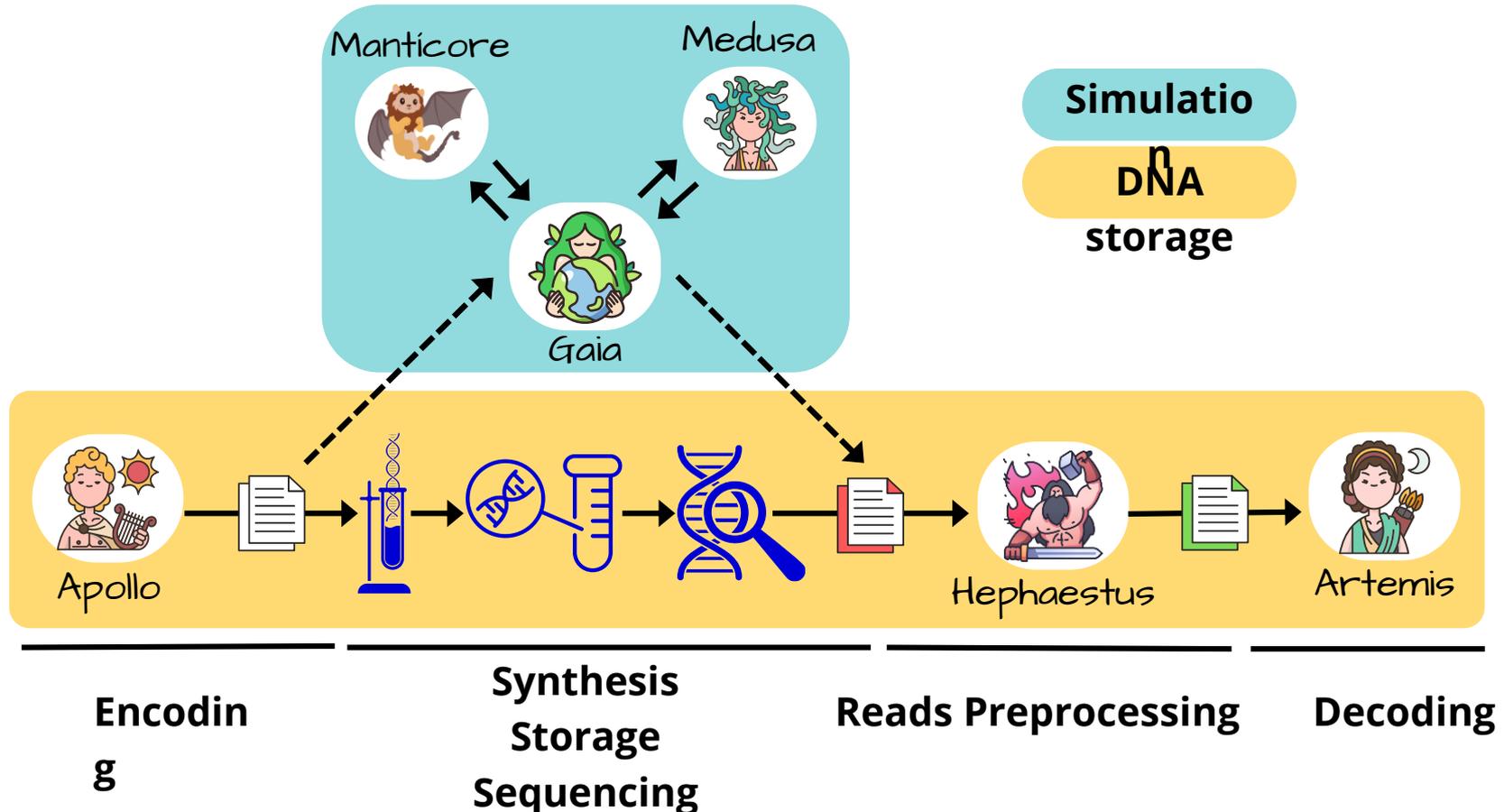
010
010
101

0101001
1111101
0101010
0100100101001
1110101111101
0100100101001

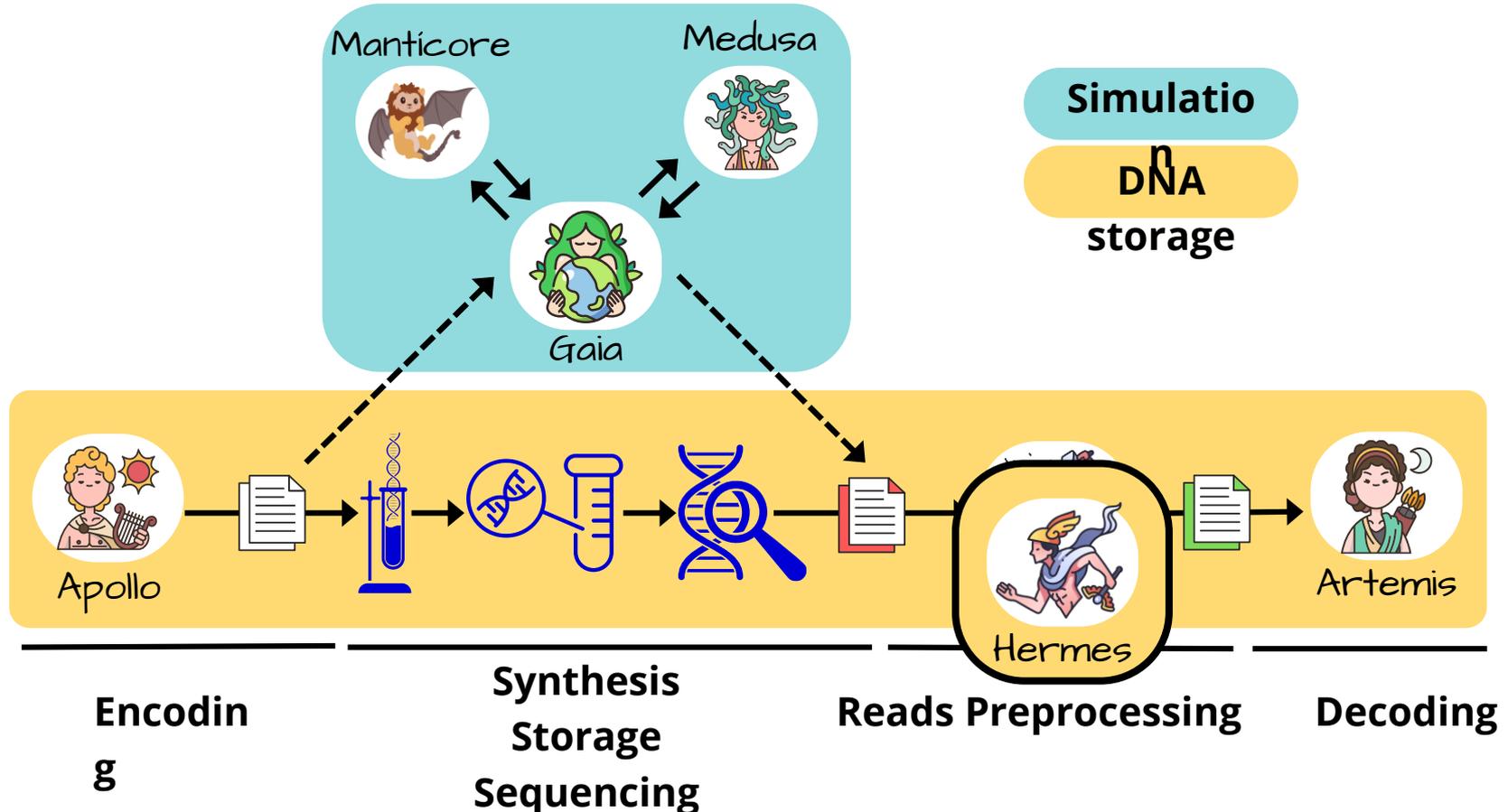


But the Process was Slow and we Needed
to Sophisticate...

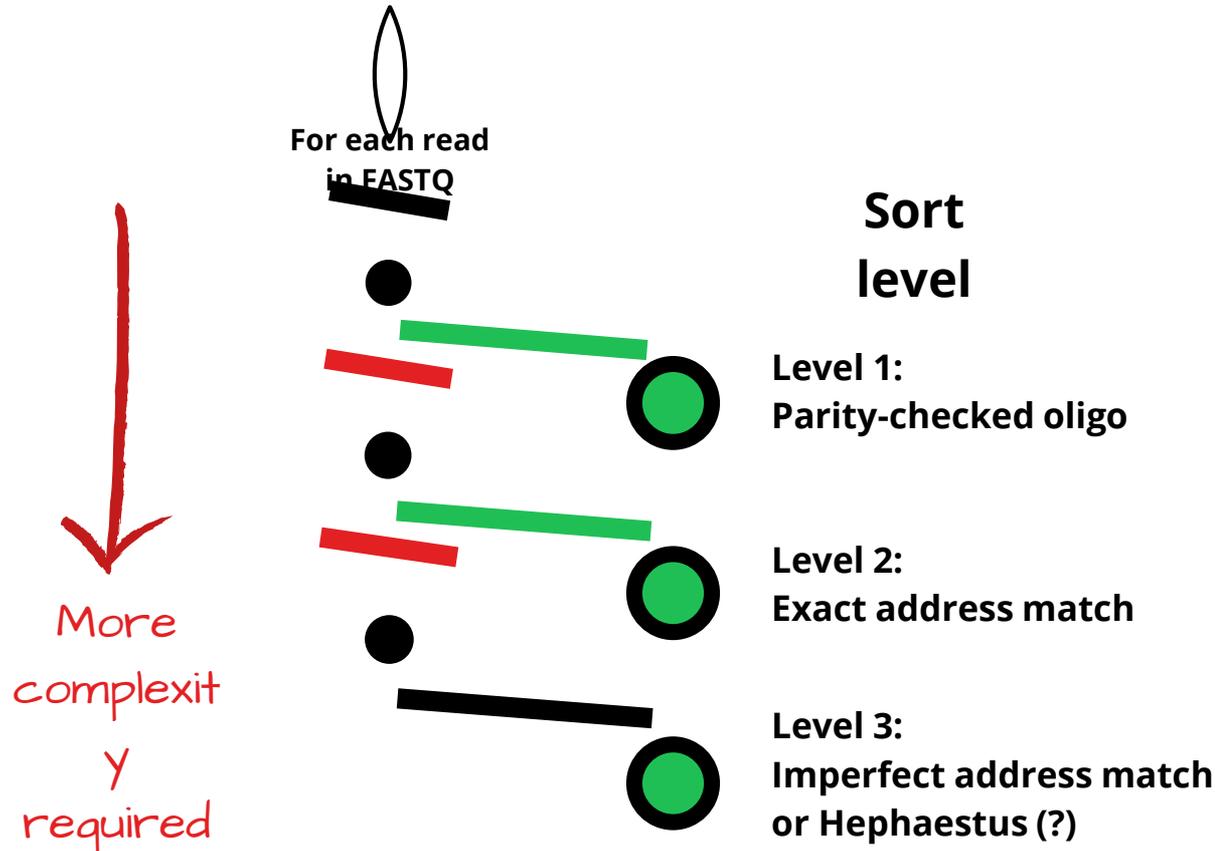
CODEC v0.4 Update



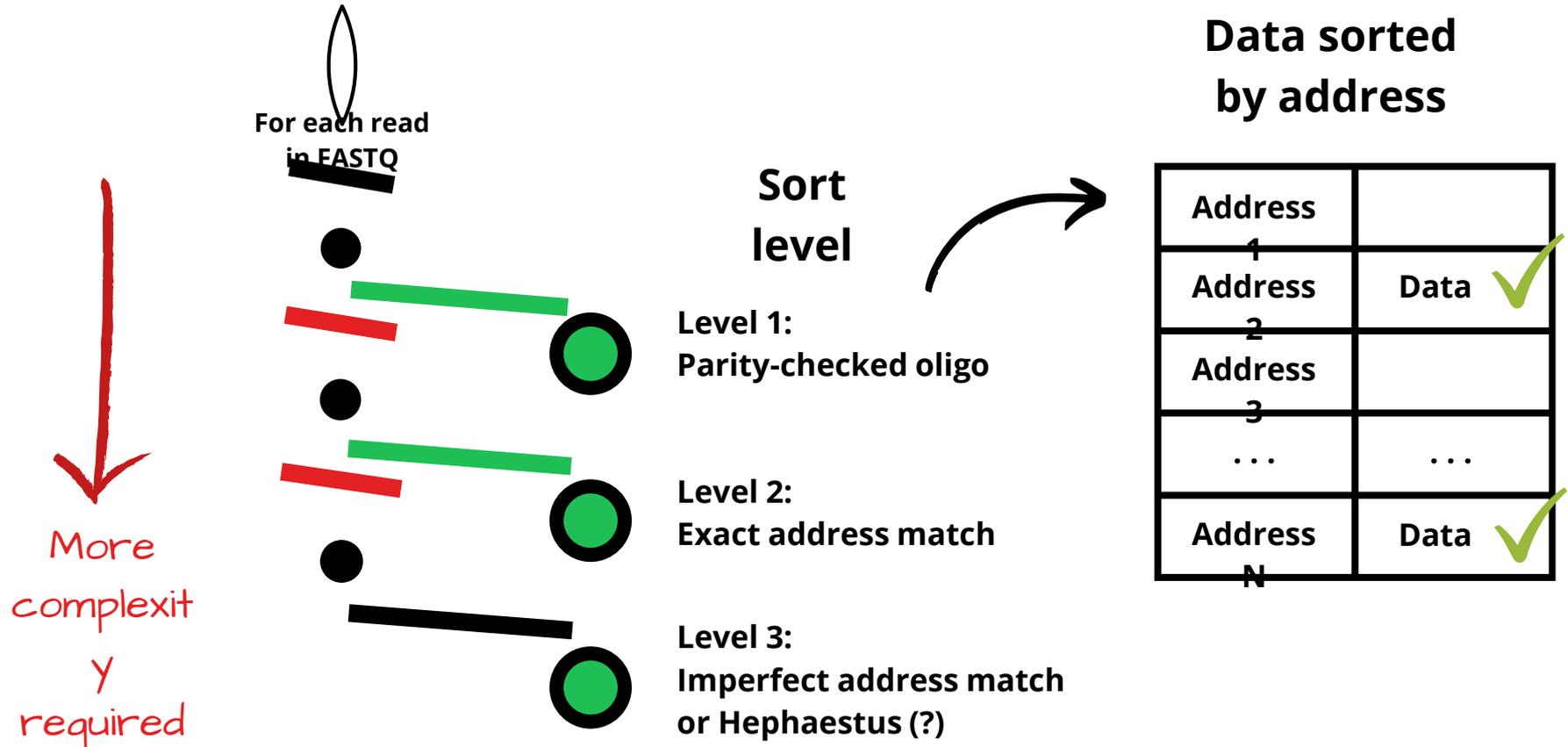
CODEC v0.4 Update



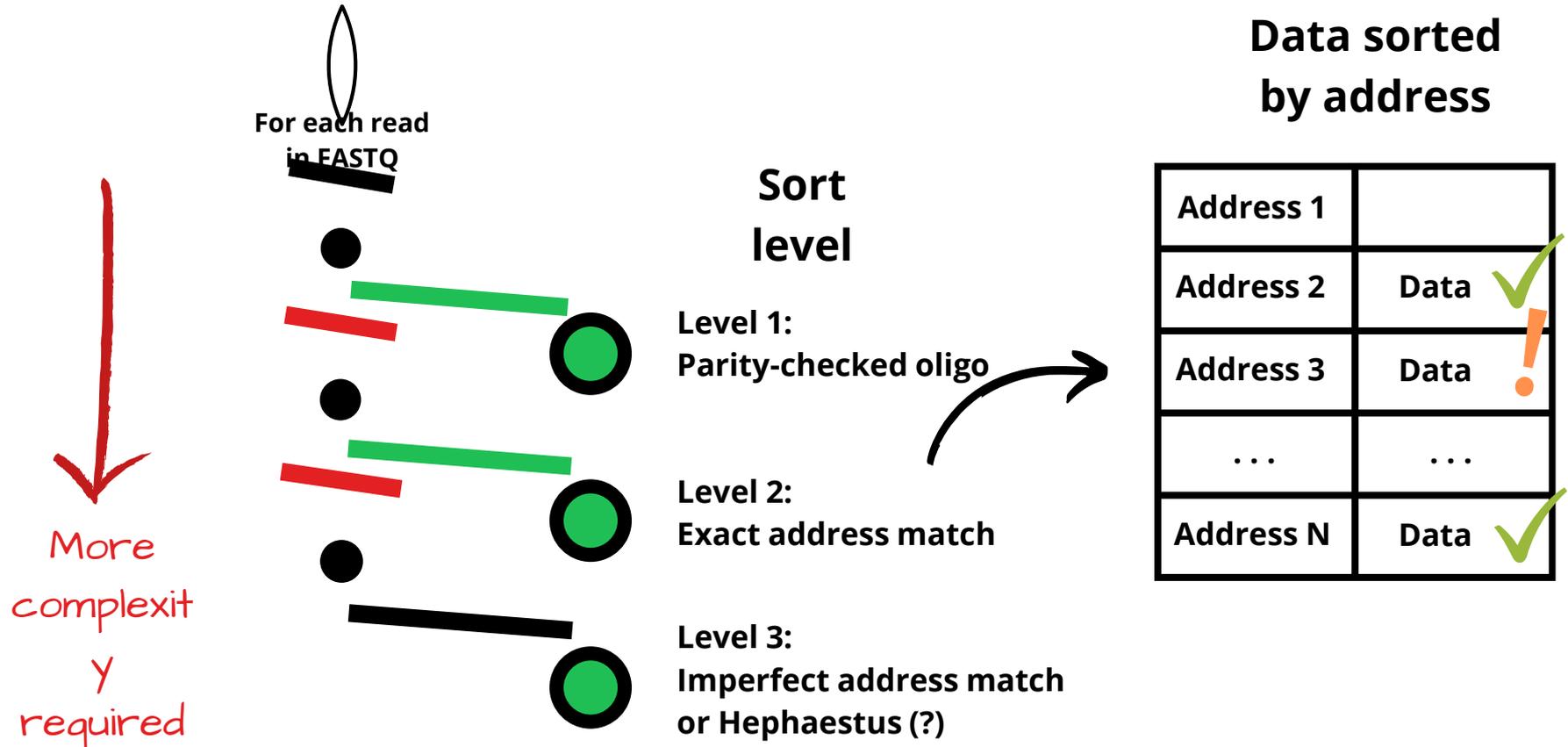
Post-sequencing processing pipeline (under discussions)

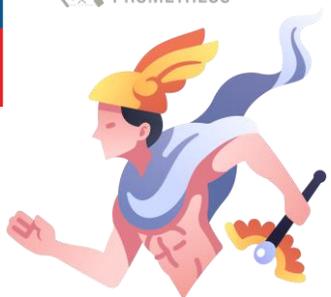


Post-sequencing processing pipeline (under discussions)



Post-sequencing processing pipeline (under discussions)





Fastq



Address

ECC sequence

5'-(...)**TACGCC****GCAACCTGCCATCC****AAGCTAAAGGGGAATCCCAATGGCA****ATCTGCACCGT****GGACGA**(...)-3'



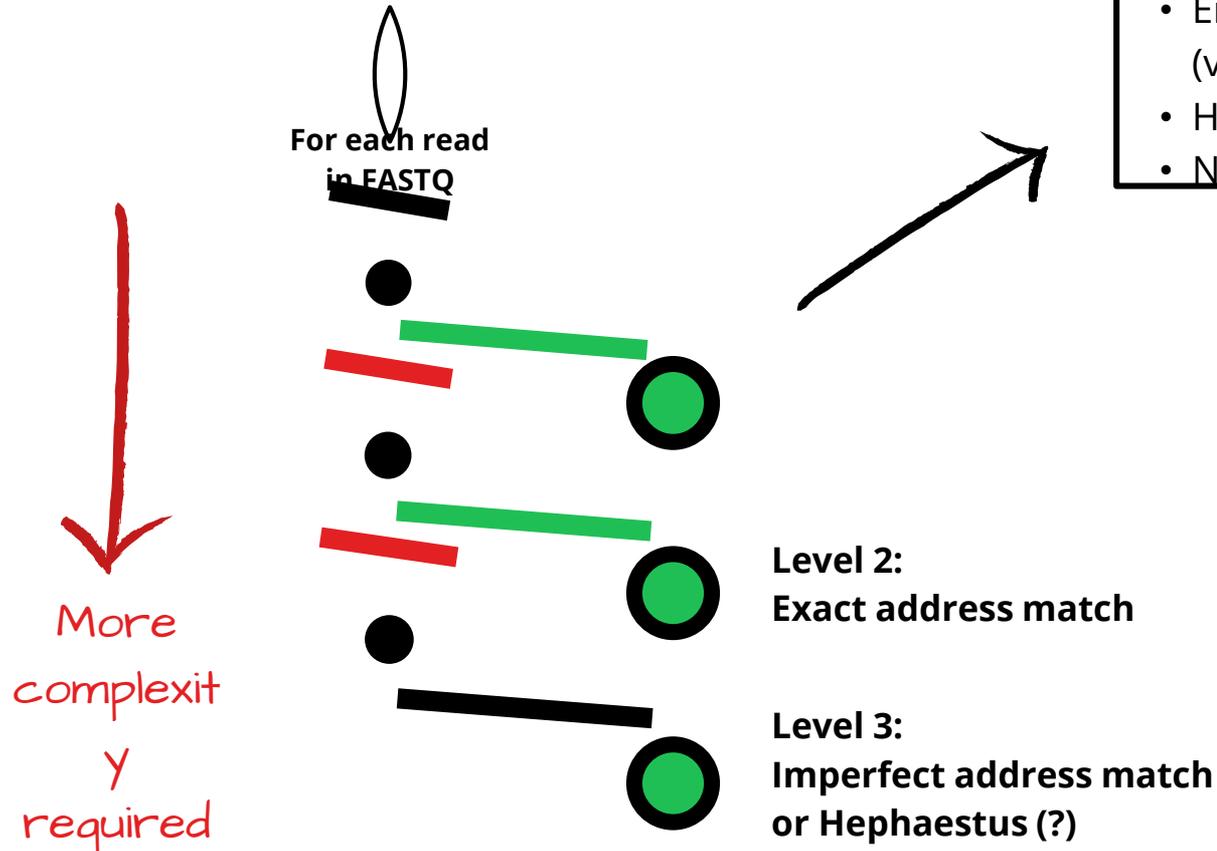
- Integrity check (sequence length and ECC)
- Try to identify the address
- Pass integrity and identify the address
- Fails integrity, but identify the address
- Fails integrity, and fails to get address



This is a perfect sequence, hold to send to Artemis for decoding



Post-sequencing processing pipeline (under discussions)





Fastq



Address

5'-(...)TACGCC**GCAACCTGCCATCC**AAGCTAAAGGGGAATCCCAATGGCA**ATCTGCACCGT**GGACGA(...)3'

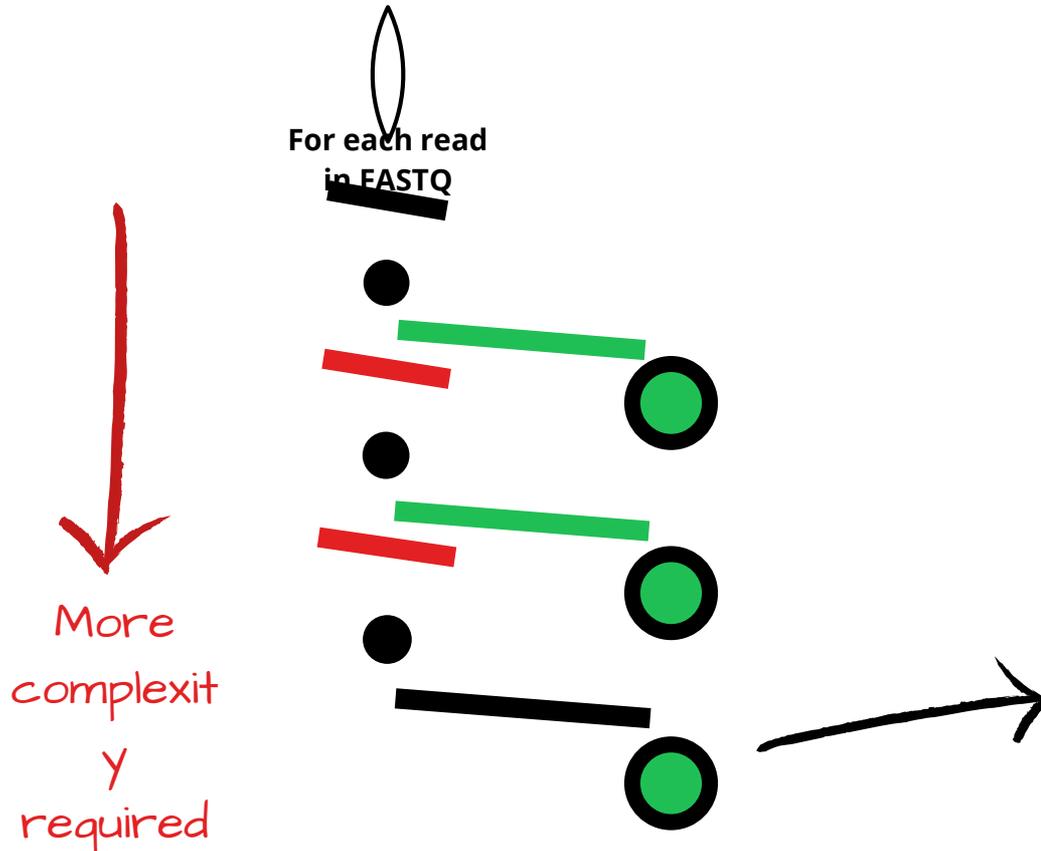
ECC sequence

- Integrity check (sequence length and ECC)
- Try to identify the address
- Pass integrity and identify the address
- Fails integrity, but identify the address
- Fails integrity, and fails to get address

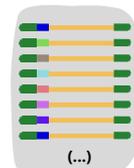
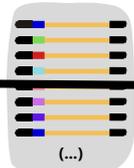
Put this sequence in the pile for that address. If we don't find a perfect one, do MSA and consensus using the pile.



Post-sequencing processing pipeline (under discussions)



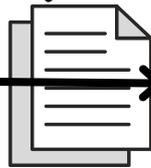
Packages



Illumina



FASTq



Parallel clustering



FASTA



Consensus sequences



Finally, a small Sample

Input file: photo from workshop
(1 data block)



- 5863 original sequences (targets)
- 64493 simulated sequences (reads)
- 100% of sequences have 2 substitution errors

Thank you!

Dr. Adriano Leal
leal@ipt.br



[linkedin.com/school/iptsp/](https://www.linkedin.com/school/iptsp/)



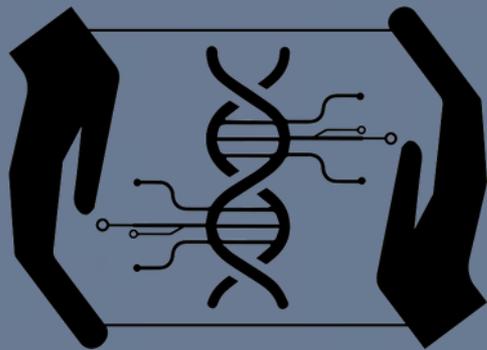
[instagram.com/ipt_oficial/](https://www.instagram.com/ipt_oficial/)



[youtube.com/@IPTbr/](https://www.youtube.com/@IPTbr/)

www.ipt.br





Lenovo™

ipt

