# Efficient Storage of Whole Human Genomes

## Motivation



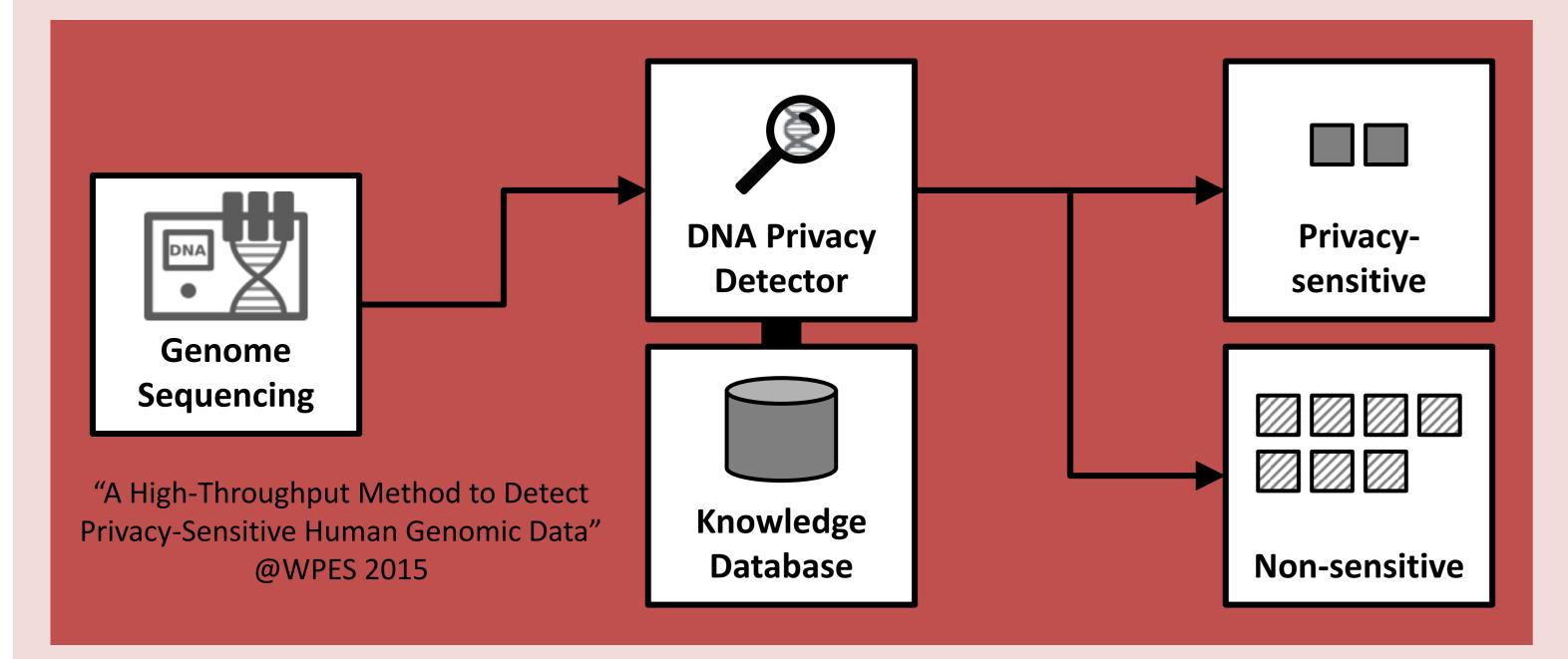
The **criticality** and **million-scale** size of sets of

human genomes require **systematic solutions** to store and share this data **efficiently** 

#### **Privacy-Awareness**

Some data are more important than others to our privacy Goals:

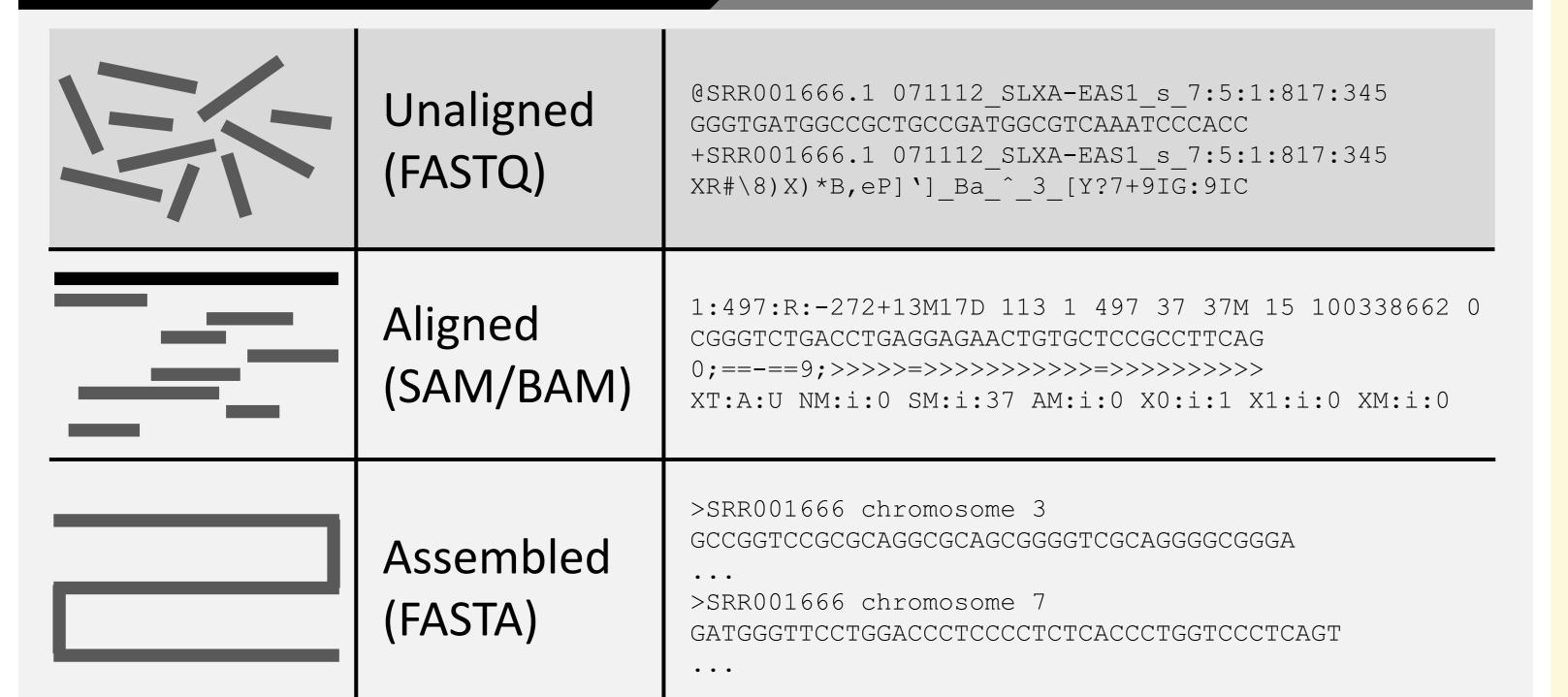
- Systematically detect the privacy-sensitive portions
- Create a database of known privacy-sensitive sequences
- Differentiate them according to their relevance to privacy



#### Challenges:

- ☐ High-throughput detection
- ☐ Completeness of the knowledge database
- Adapting subsequent components

# **Sequencing Genomes**



## <u>Vinicius Vielmo Cogo</u>, Alysson Neves Bessani

{vvcogo,anbessani}@fc.ul.pt

LaSIGE, Faculdade de Ciências, Universidade de Lisboa, Portugal







## State-of-the-Art

No automatic detection of privacy-sensitive sequences

High (lossy) compression ratio, slow (de)compression No deduplication for genomic data

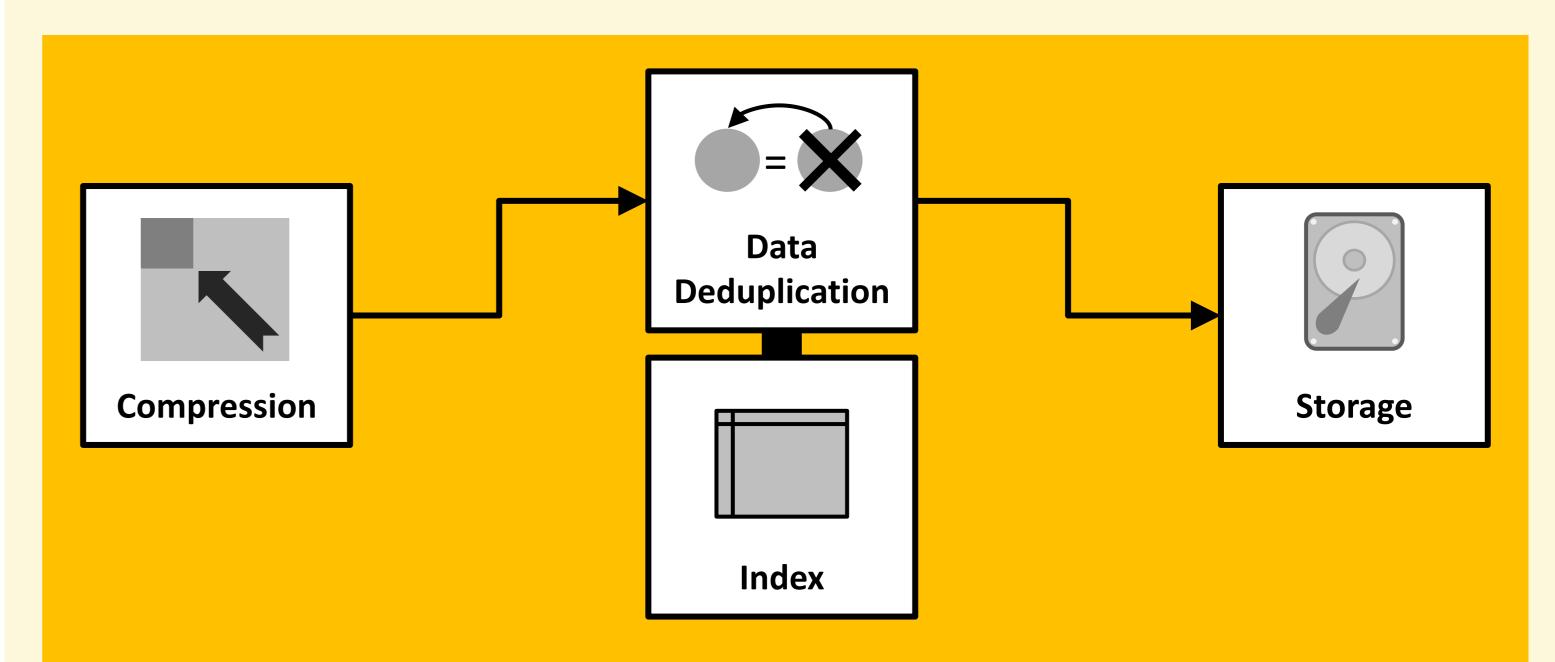
## Objective

Enhance the storage ecosystem with privacy-awareness and data reduction to enable the efficient storage of human genomes

#### Data Reduction

**Human genomes** are more than **99.5**% similar to each other Goals:

- Integrate lossless compression and deduplication
- Focus on compression ratio and decompression time
- Deduplicate similar portions



# Challenges:

- ☐ High compression ratio in lossless algorithms
- ☐ Fast decompression with few resources
- ☐ Grouping similar sequences and storing their differences
- ☐ Developing application-specific deduplication
- ☐ Storing a human genome with less than \$1 per year

Reduction	Gains	Hard Disk	Amazon S3	Amazon Glacier
0x Uncompressed	0%	195	99	25.2
5x	80%	39	19.8	5.04
200x	99.5%	0.975	0.495	0.126

Annual storage cost per genome (in \$)

- ☐ Adapt storage systems with application-specific knowledge
- ☐ Adapt data management to differentiated portions