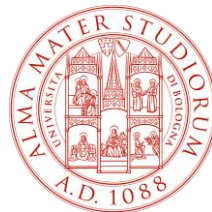# GOLAM: A framework for Analyzing Genomic Data

**Lorenzo Baldacci**, Matteo Golfarelli, Simone Graziani, and Stefano Rizzi

**DISI – University of Bologna, Italy**
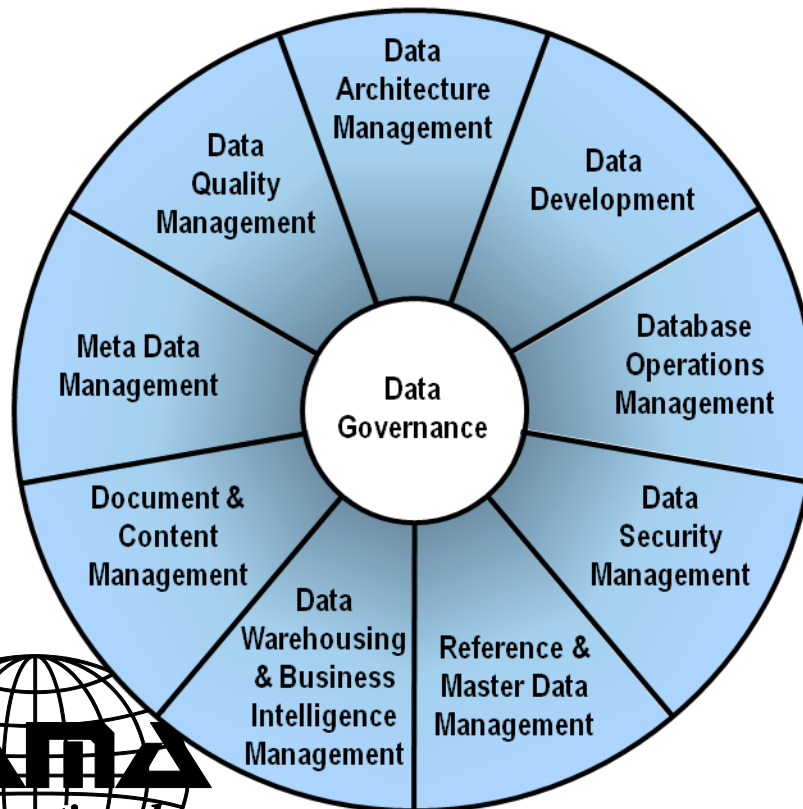
# Contents

- ➢ **Genomic Data**

- ➢ **Multidimensional Analytics applied to Genomic Data**

- ➢ **Discussion and Current Research Directions**

# Dealing with Genomic Data

- Genomic challenges: from sequencing to data management
  - ✓ Bioinformatics challenges are moving towards storage, retrieval, security, and presentation of genomic information
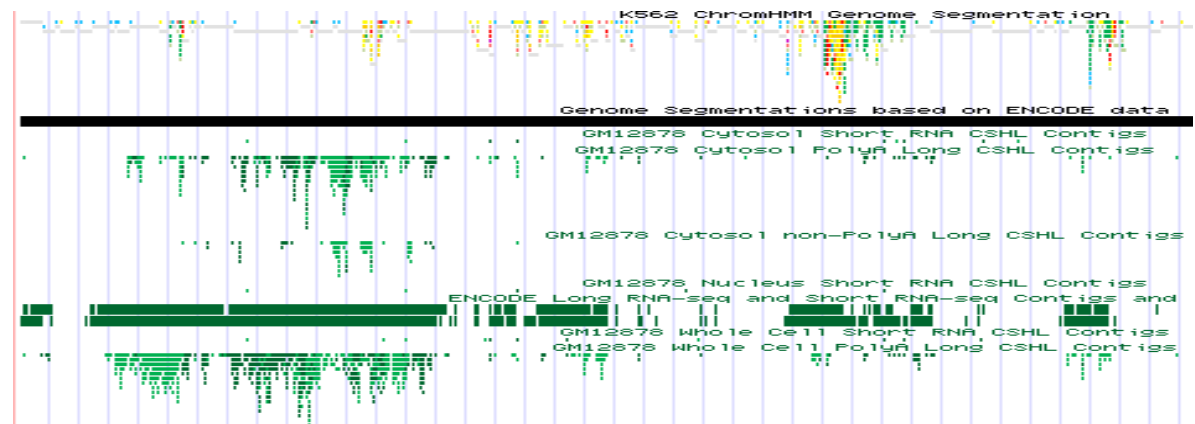


**Sequencing is not enough!**

# Genomic Analytics Challenges

- Ad hoc (vertical) applications supporting specific biological questions

- Laboratories and Biologists frequently undertake analytical tools development in-house

- Genomic browsers are effective when analyzing detailed data but fall short when analyzing aggregated data
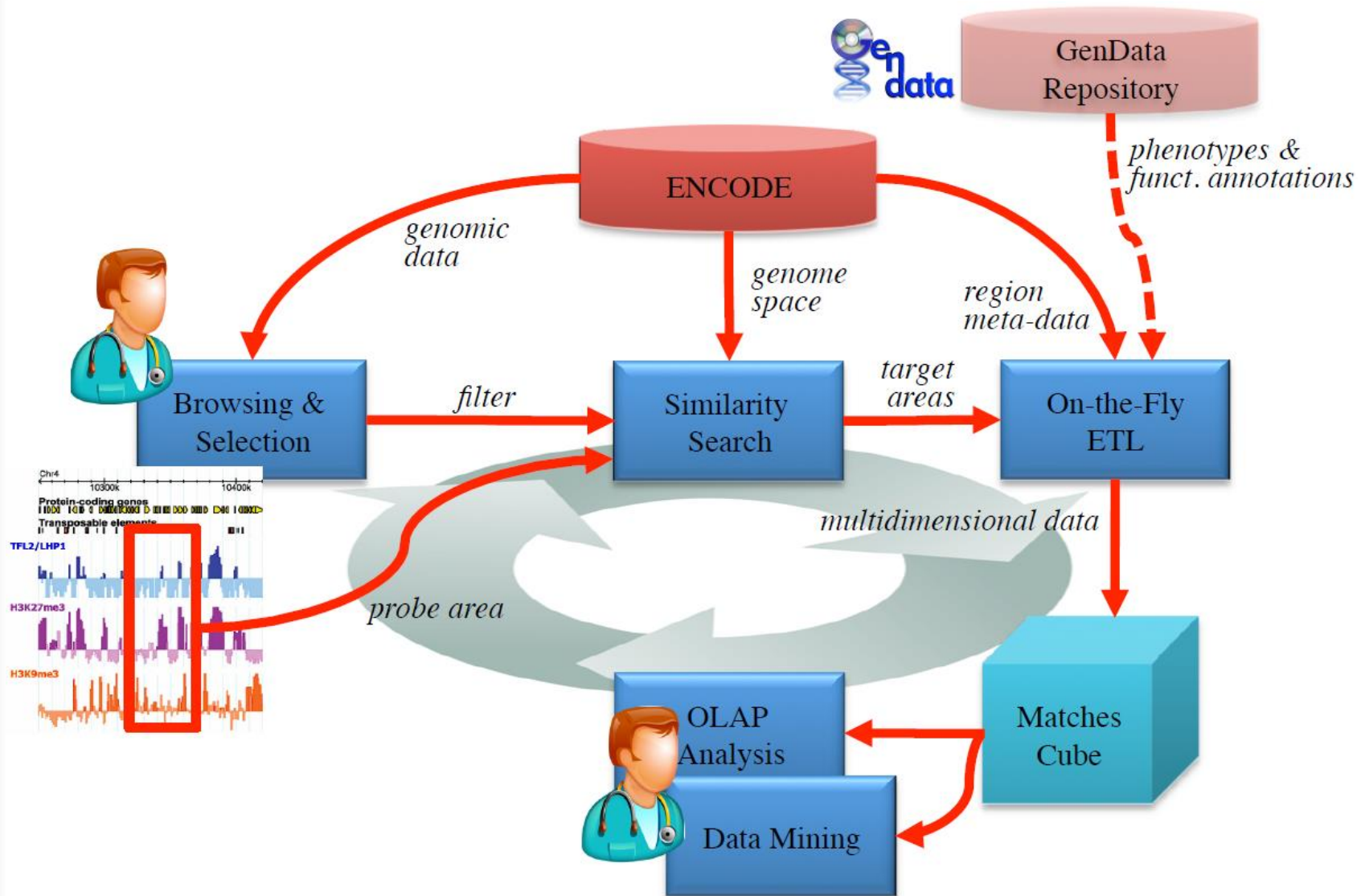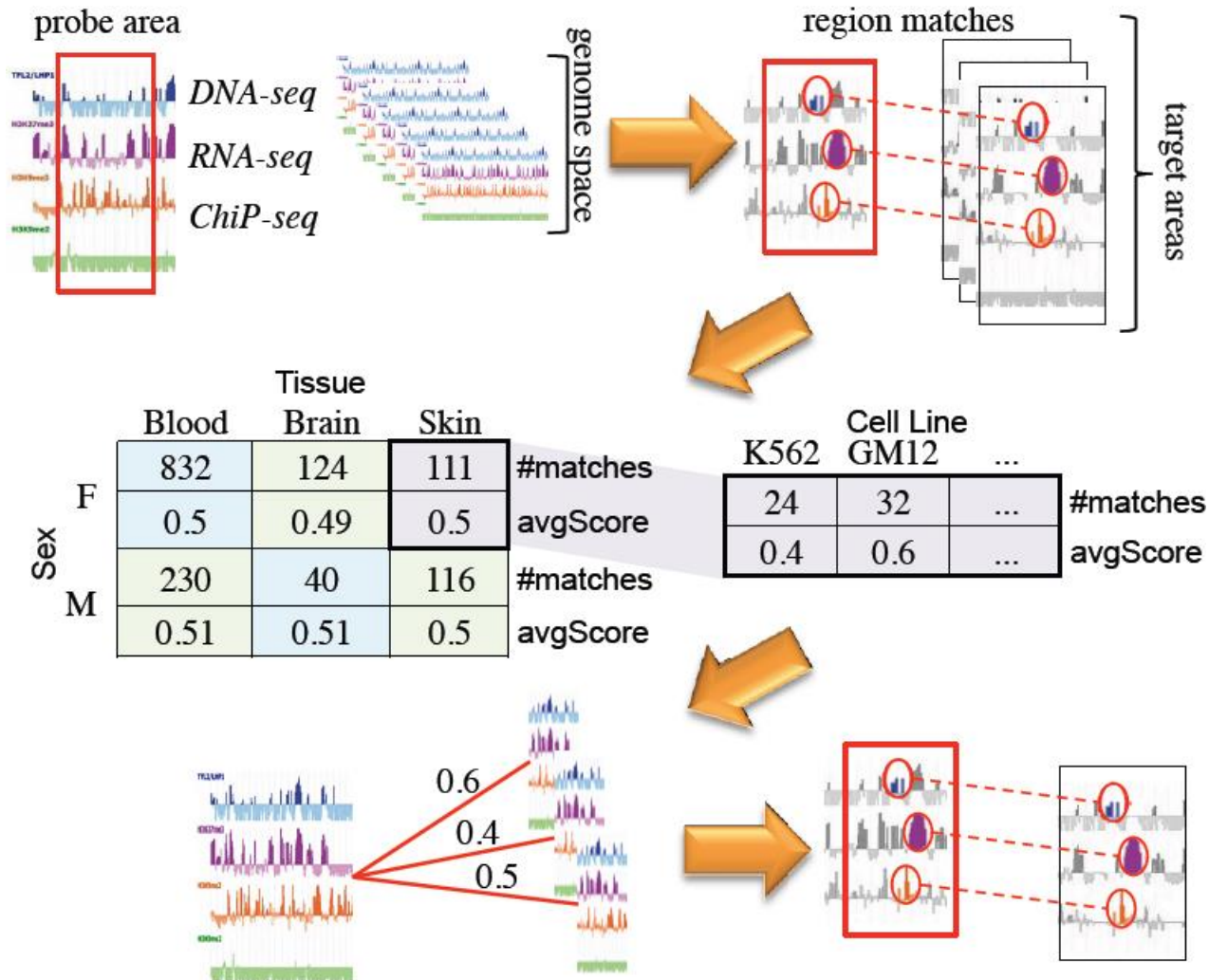
# GOLAM: Genomic OLAM

- is a framework for OLAP analysis and mining

- helps biologists in overcoming the rigidity of genome analysis methods

- automates and speeds up analysis sessions and introduces a multi-resolution view of the data

- addresses the issue of loading data to the cube "On-the-Fly"

# The GOLAM Framework

# Genomic OLAP

# Genomic OLAP

# Research Challenges

- Genomic data modeling
- Schemaless data structure to deal with
  - ✓ Metastars
- Non-traditional data sources
  - ✓ Files
  - ✓ Ontologies
  - ✓ Big Data

# ENCODE data model



**ENCODE** is the *Encyclopedia of DNA Elements* publicly available and recognized to be the standard repository for genomic data and functional elements.

- **Cell Line**: is the biological sample on which an experiment is carried out.
- **Experiment**: is a single analysis on a cell in a laboratory.
- **Region:** is a segment of an experiment. It holds biological functional information.

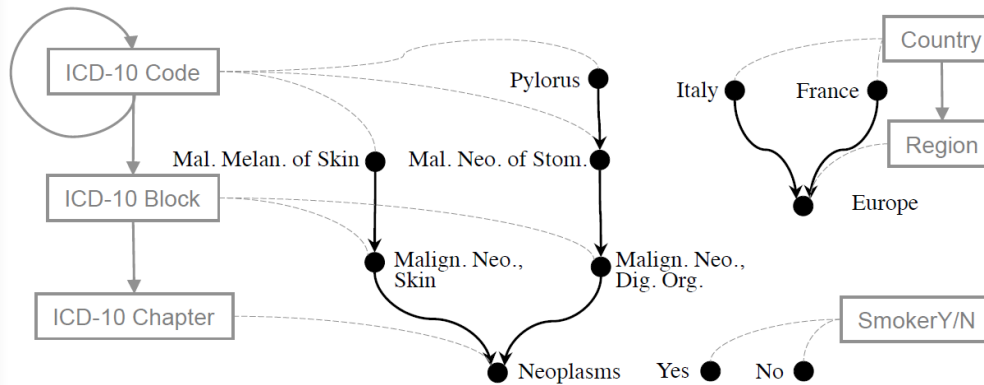# Multidimensional Schema

# Metastars



**Metastars** use meta-modelling coupled with traditional dimension tables to support non-onto, non-covering, and non-strict hierarchies.

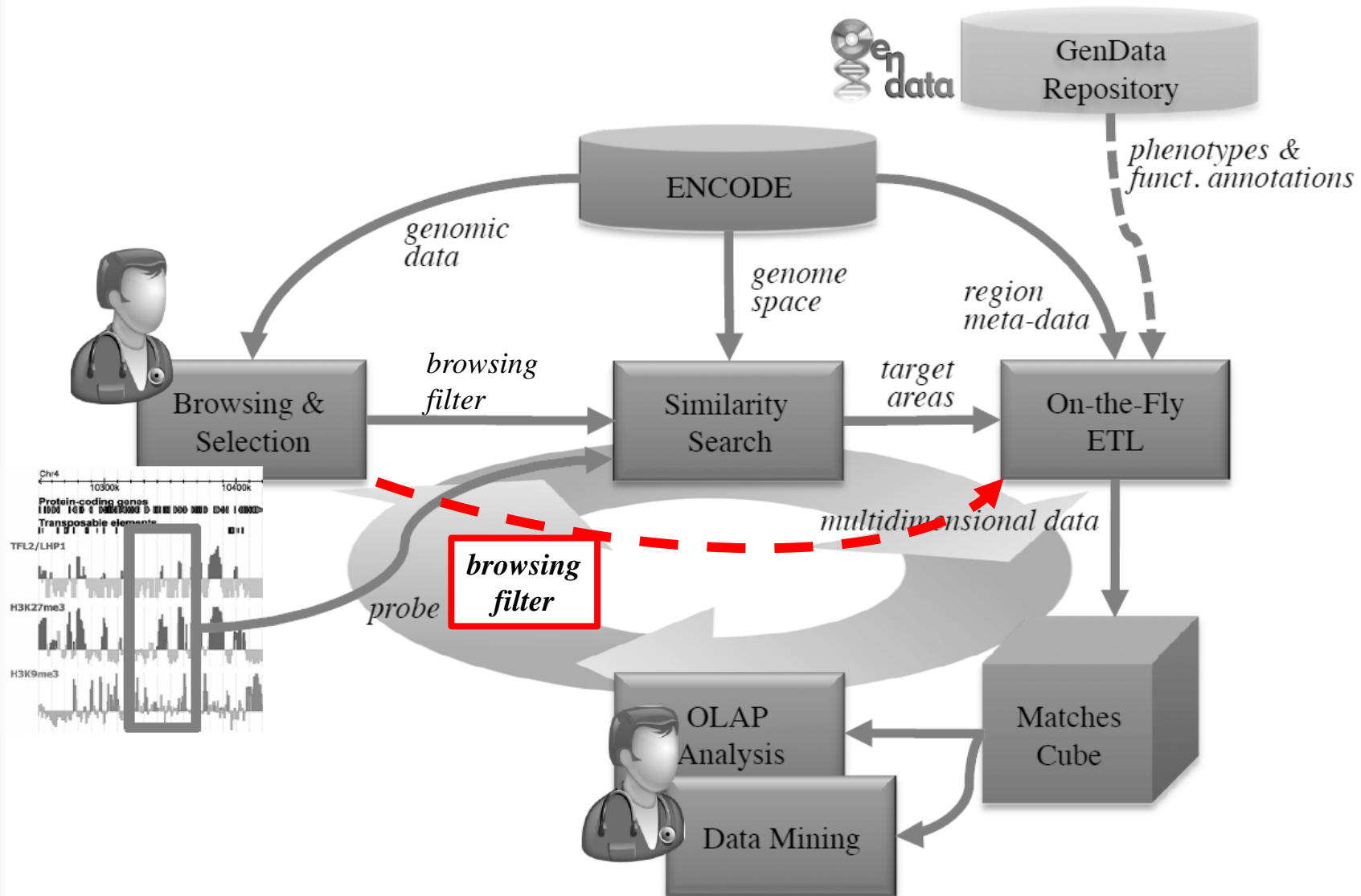By employing Metastars we are able to **model hierarchies** whenever the schema is **dynamic or missing**.

### DATA_T

| IdT | Value | Level |
|-----|-------|-------|
| 1 | Yes | SmokerY/N |
| 2 | No | SmokerY/N |
| 3 | Italy | Country |
| 4 | France | Country |
| 5 | Europe | Region |
| 6 | Neoplasms | ICD-10 Chapter |
| 7 | Mal. Neo., Skin | ICD-10 Block |
| 8 | Mal. Neo., Dig. Org. | ICD-10 Block |
| 9 | Mal. Melan. of Skin | ICD-10 Code |
| 10 | Mal. Neo. of Stom. | ICD-10 Code |
| 11 | Pylorus | ICD-10 Code |

### ROLLUP_T

| ChildId | FatherId |
|---------|----------|
| 1 | 1 |
| 2 | 2 |
| ... | ... |
| 3 | 5 |
| 4 | 5 |
| 7 | 6 |
| 8 | 6 |
| 9 | 7 |
| 10 | 8 |
| 11 | 10 |
| 9 | 6 |
| 10 | 6 |
| 11 | 8 |
| 11 | 6 |

# On-The-Fly ETL

# ETL Tests

- ENCODE counts a number of experiment files in the range of 25K, leading to **over a billion instances in the region dimension**

- The number of matches might count **over 300K events for each analysis session**

- Our preliminary tests aim at evaluating the On-the-Fly ETL from the efficiency point of view
  - ✓ The "eager" approach is compared vs. the approach that loads the matching genome space only
  - ✓ Tractability threshold is set afterwards

# ETL Tests

- A probe area composed of 50 regions and 3 different genome spaces

  - ✓ Test 1 ($T1$) consisting of 740K regions
  - ✓ Test 2 ($T2$) consisting of 4.42M regions
  - ✓ Test 3 ($T3$) consisting of 54.5M regions

Table 1: Genome spaces and matching genome spaces for tests

|  | Genome space | | Matching genome space | |
|---|---|---|---|---|
|  | ♯ files | ♯ regions | ♯ files | ♯ regions |
| $T1$ | 9 | $\approx 7.4 \times 10^5$ | 3 | $\approx 1.4 \times 10^5$ |
| $T2$ | 47 | $\approx 4.4 \times 10^6$ | 12 | $\approx 9.7 \times 10^5$ |
| $T3$ | 783 | $\approx 5.4 \times 10^7$ | 69 | $\approx 5.8 \times 10^6$ |

Table 2: ETL performance (in seconds)

|  | Genome space | Matching genome space |
|---|---|---|
| $T1$ | 185 | 37 |
| $T2$ | 1121 | 246 |
| $T3$ | 13128 | 1492 |

- The tractability threshold can be reasonably set to 50 ENCODE files

# Discussion

- With GOLAM we took a first step towards overcoming current limitations of genome analysis methods
  - ✓ Analysis session has been automated and speeded up
  - ✓ More analysis flexibility has been introduced

- We proved that ETL processes can be integrated within the analysis session in order to improve efficiency in those DW applications employed in non-traditional domains
  - ✓ big data
  - ✓ scientific data storage
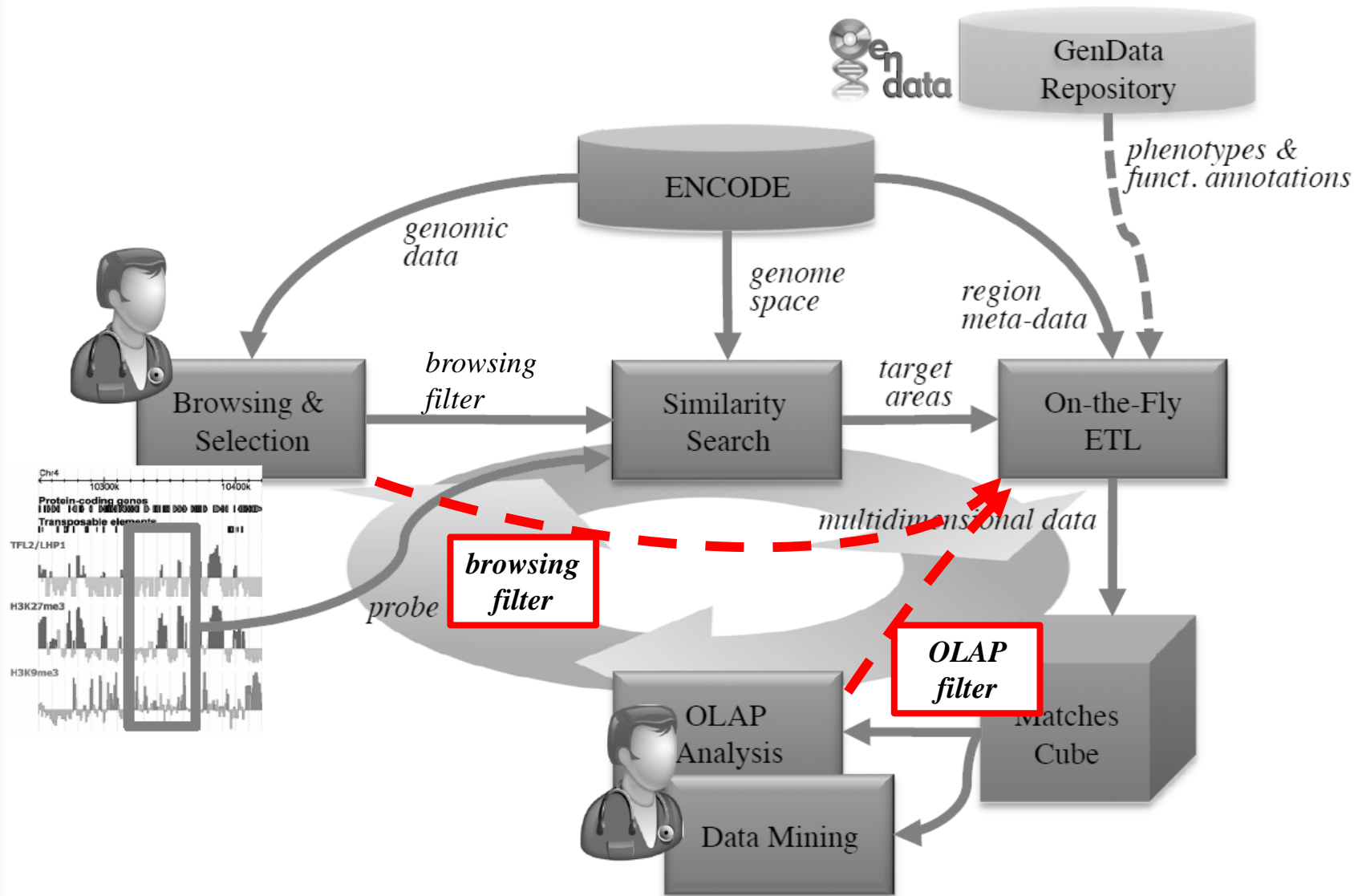  - ✓ open linked data
  - ✓ etc.

# (Not so) Future work

- On-the-Fly ETL can be further pushed so to be framed within the OLAP session itself. In this manner interesting data (according to the user) might be loaded into the cube
  - ✓ Multidimensional indexing and new dice operators must be employed by the ETL processes
  - ✓ Source data extraction must be done according to a cost function that considers many facets of optimization (e.g. time, cost, etc.)

- Approximate results could be provided in order to improve the overall session's responsiveness
  - ✓ ETL should be integrated and designed so that it gather data as a stream and exposes partial results to the user

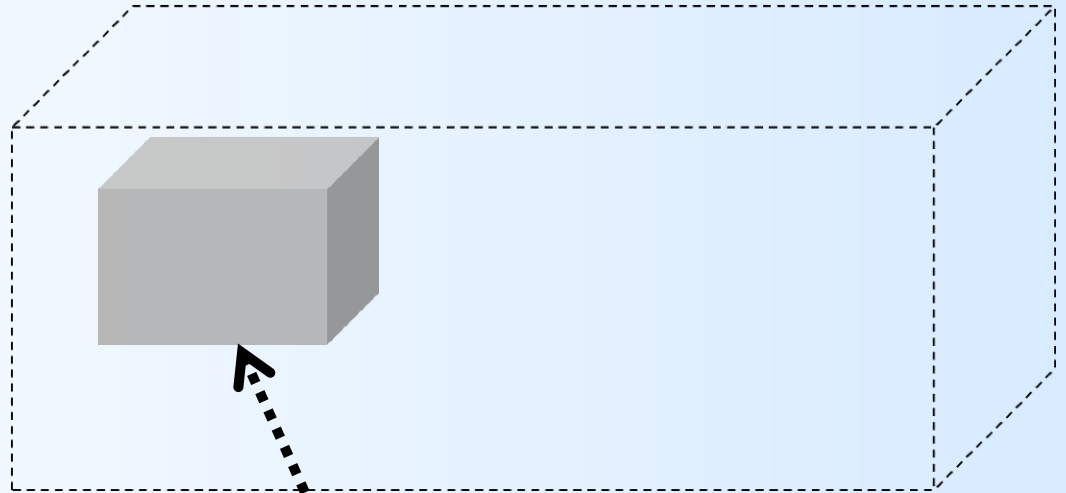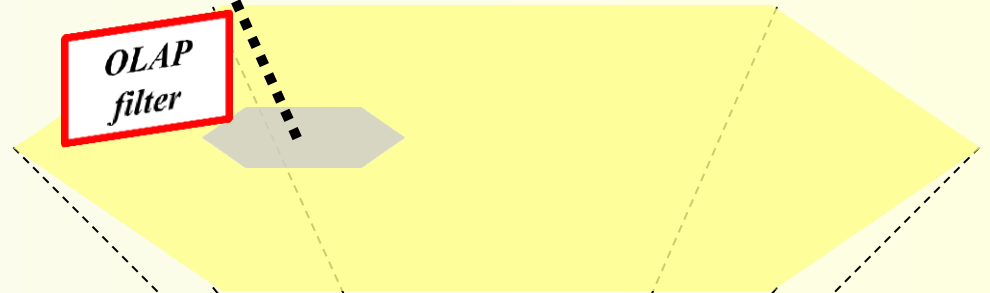# (Not so) Future work

# (Not so) Future work

*Genomic Space*

# (Not so) Future work



*Virtual Data Space*

*Genomic Space*

browsing filter

# (Not so) Future work

# (Not so) Future work



Multidimensional Space

Virtual Data Space

OLAP filter

Genomic Space

browsing filter

# (Not so) Future work

...questions time...