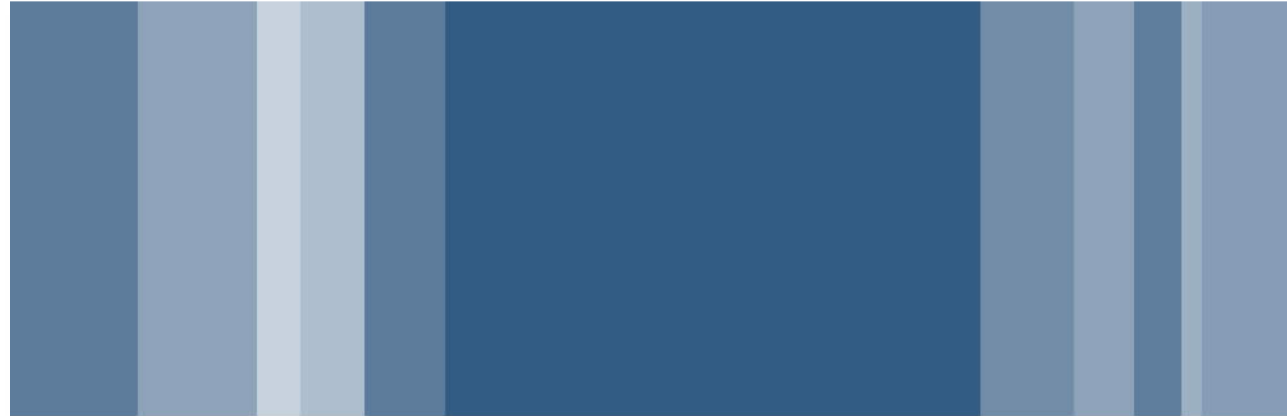**Genomic Computing**
*Politecnico di Milano*

POLITECNICO
DI MILANO

Dipartimento
di Elettronica, Informazione e
Bioingegneria

**Genomic Data Model and GenoMetric Query Language as research enabler to discover genome properties**

*Marco Masseroli and Stefano Ceri*

*(joint work with several PhD students)*

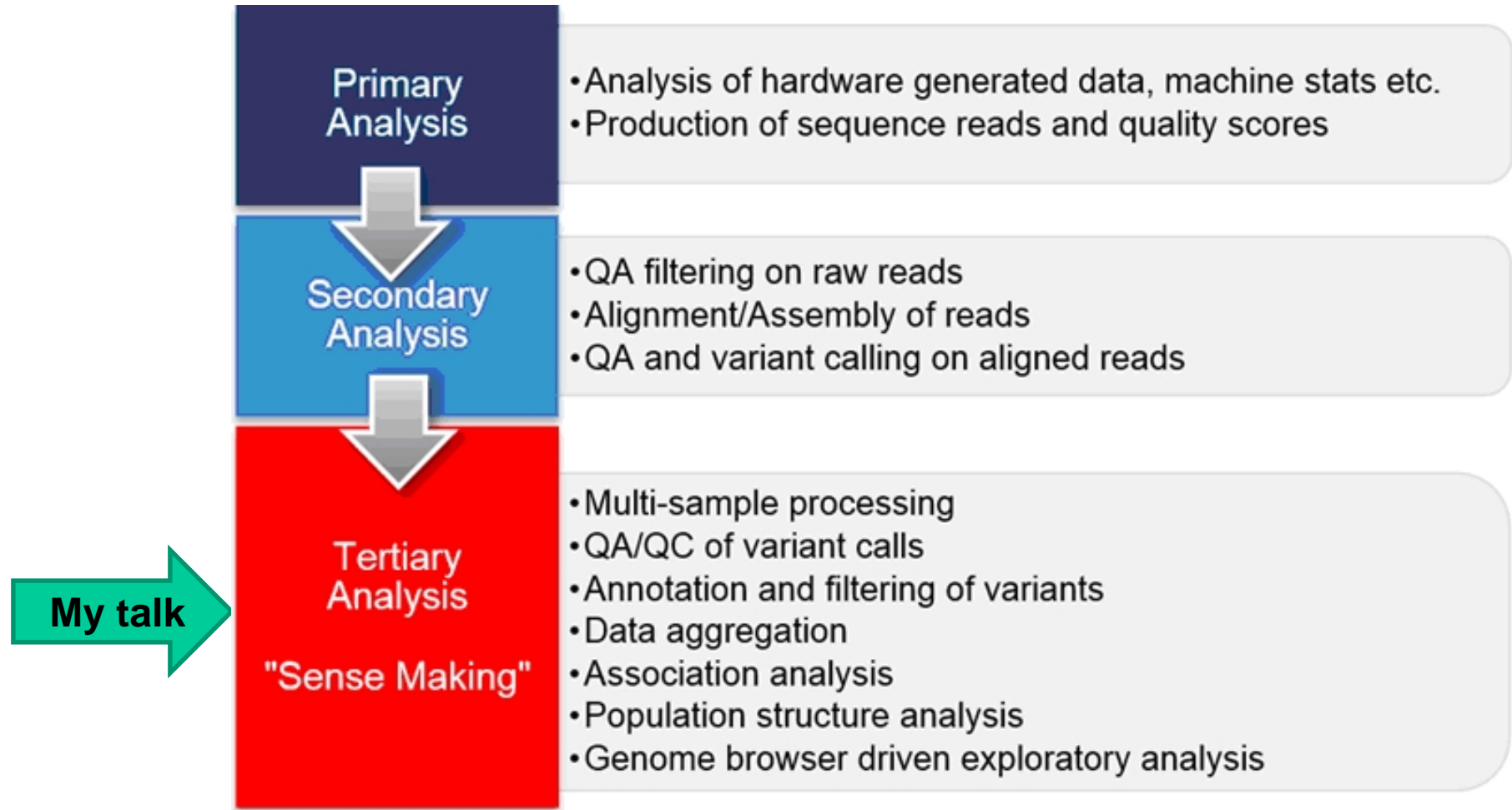*Politecnico di Milano, BioInformatics Group*

- <u>Next Generation Sequencing</u> technology is about to provide affordable (in time and cost) and precise determinations of genome wide:
  - DNA sequence / variations (DNA-seq)
  - gene subregions' activity (RNA-seq) [all gene test]
  - protein-DNA interaction regions (ChIP-seq)
  - open chromatin (DNase-seq)

Goal of $1,000 full genome sequencing in under an hour has just met



- Very many DNA-interacting proteins / subjects / conditions will be soon evaluated
  - Personalized medicine (diagnosis and treatment)
  - Each NGS test can generate 0.4TB -> **Big Data** scenario

Source: http://blog.goldenhelix.com/grudy/a-hitchhiker%E2%80%99s-guide-to-next-generation-sequencing-part-2/

Medical Literature

Biologist

Clinician

Clinical Protocols

Ontological Knowledge

Genome Browser

Data Analytics

Genomic Data Integration

Publishing / Crawling / Searching

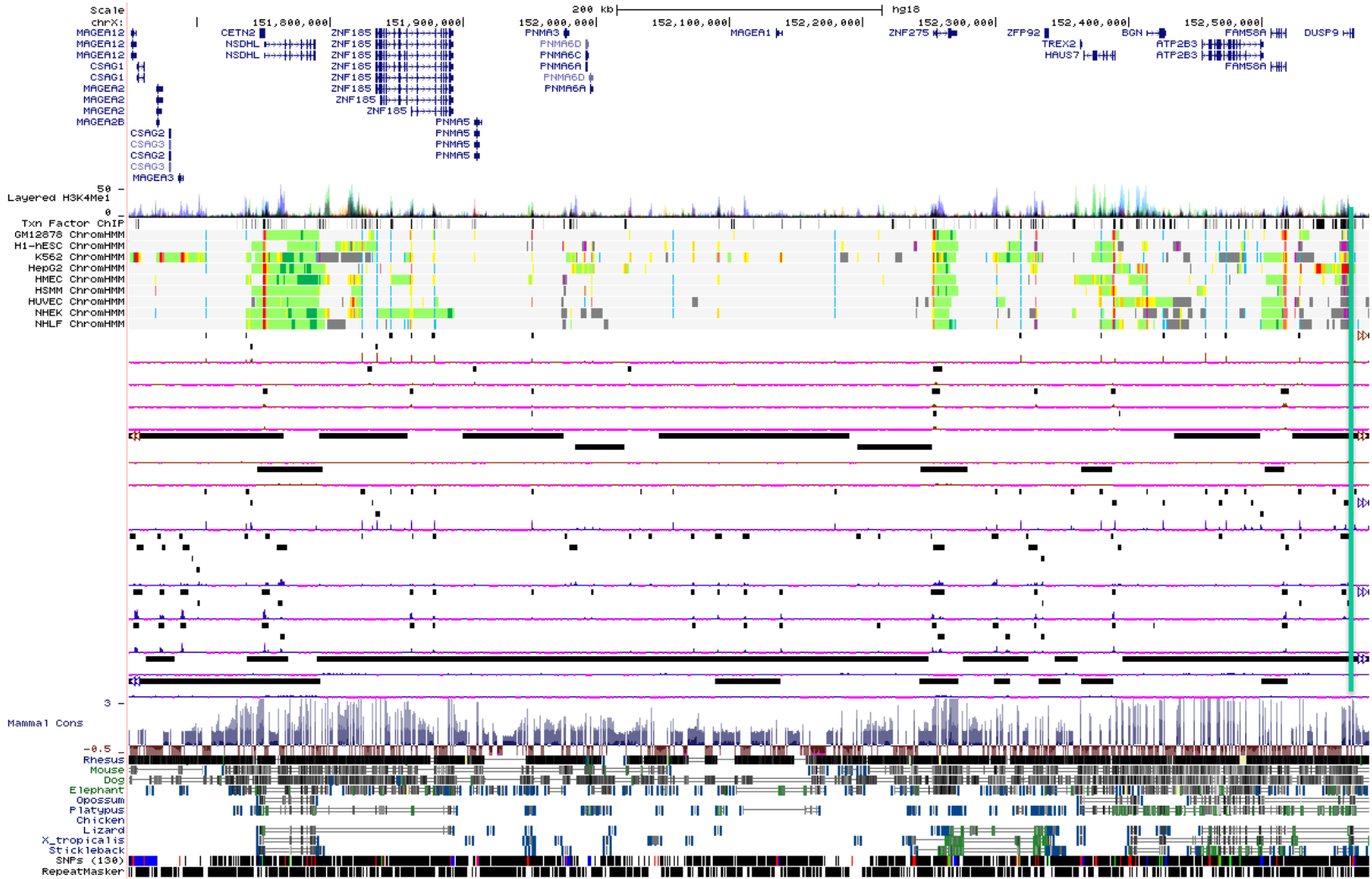Clinical Data Integration

Heterogeneous Genomic Data Sources

**Personal Patient Data**

Heterogeneous Clinic Data Sources
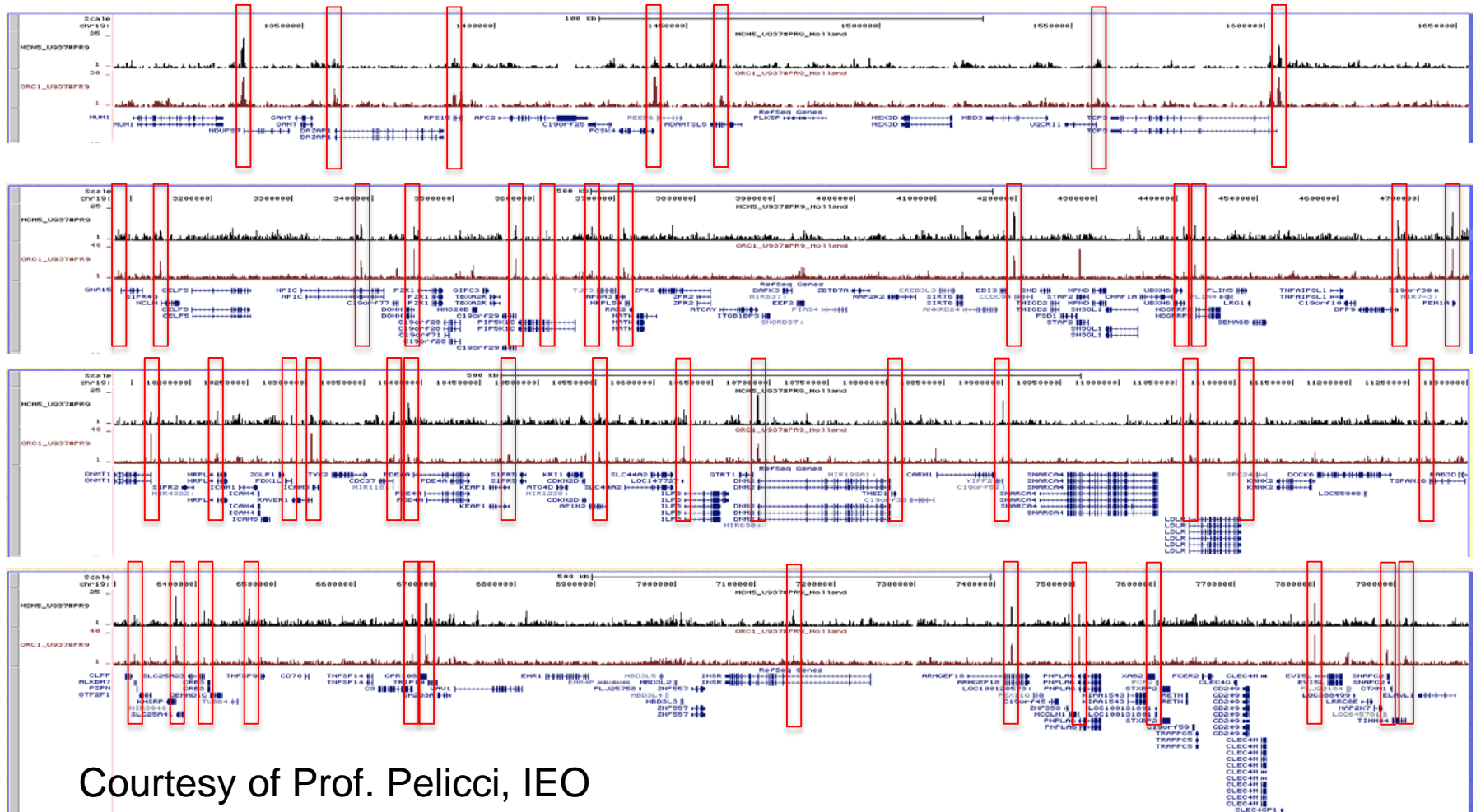
A number of genomic features (tracks)

Data tracks

One macro genomic region

- Working together with biologists for giving answers to the problems behind the «courtesy» slide



Courtesy of Prof. Pelicci, IEO

- **(Epi)genotype-phenotype relationship discovery**: understanding genomic regions, genome variations and their associations with different phenotypes
  - highly heterogeneous scenario

- It requires evaluating, in several different conditions and types of individuals:
  - genome (DNA) sequence variations
  - gene activity & its regulation
  - occurring interactions

Scientist's typical questions

(from our interaction with IEO - European Oncology Institute and IIT - Italian Institute of Technology)

- *Can <u>interesting DNA regions</u> and their relationships be discovered using genome-wide queries?*

- *Can <u>genomic data of patients</u> be grouped according to <u>clinical phenotype</u> and compared?*

- *Can the genomic features of all the genes involved in the same biological process be extracted and then analyzed?*

- *Can we retrieve portions of the genome of given patients, extracting them from remote servers and comparing them?*

- *Can interesting DNA regions and their relationships be discovered using genome-wide queries?*

    **Genometric query language**

- *Can genomic data of patients be grouped according to clinical phenotype and compared?*

    *Genometric query language +* **clustering**

- *Can all the features of the genes involved in the same biological process be extracted and then analyzed?*

    *Genometric query language +* **data analysis**

- *Can we retrieve portions of the genome of given patients, extracting them from remote servers and comparing them?*

    *Genometric query language +* **indexing & search**

- **Data model**: design a <u>simple</u> and <u>format-independent</u> data model for describing datasets with both <u>genomic regions</u> and general <u>provenance</u> <u>information</u> (including phenotype)

- **Query language**: design a query language where both <u>genometric aspects</u> (about the placement of regions on the genome) and <u>provenance</u> can be <u>queried</u> at a high level of data independence and transparency

- **Integrative data analysis**: translating query results into a <u>genome space</u> which is the ideal start point for <u>correlation</u> and <u>network analysis</u>

- **Data search**: design protocols for <u>data crawling</u> and <u>indexing</u> based on the data model
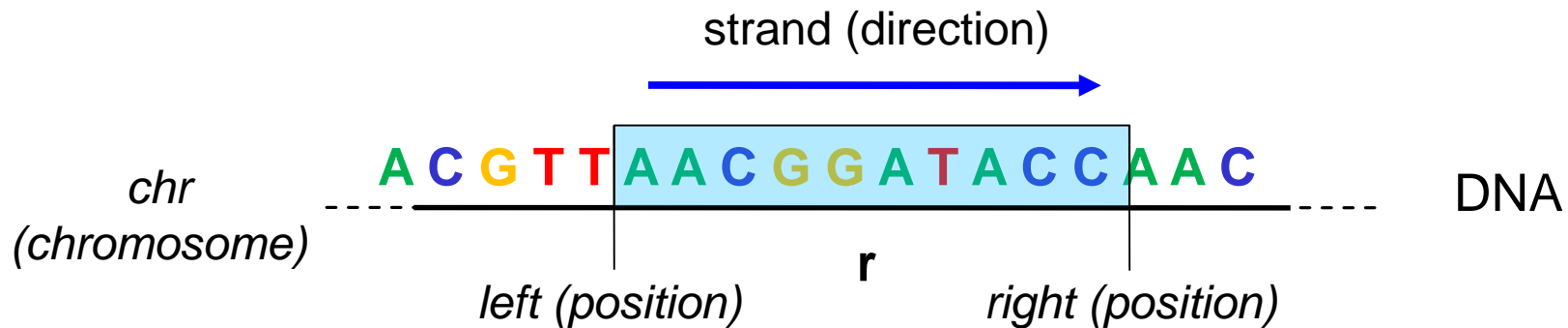
# Genomic Data Model

Within the same sample, <u>two</u> kinds of data:

- **Region values** aligned w.r.t. a given reference, with specific left-right ends within a chromosome, and with several <u>associated attributes</u> (e.g. p-value of region significance)



- **Metadata,** with free-format attribute-value pairs, storing all the knowledge about the sample

- <u>Regions</u> of the model are **data format independent** and provide an <u>interoperability</u> framework for <u>comparing</u> data on <u>mutations</u>, <u>expression</u> or <u>regulation</u> using regions as common ground

- Metadata attribute-value pairs of the model are **info-system independent** and provide an <u>interoperability</u> framework for <u>comparing samples</u> based upon their <u>biological aspects</u>
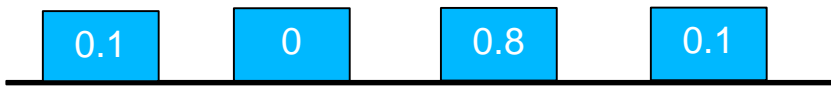
Sample 1

Tumor_type = brca
Patient_age = 75

Sample 2

Tumor_type = brca
Patient_age = 63
Sex = Female

Sample 3

Tumor_type = brca
Patient_age = 58

- **Region values**: {*expID, region:(chr, left, right, strand), p-value*}

```
1    (3, 3245,4535, +)       0.0000000024
1    (3, 5443, 6553, +)      0.00000044
1    (1, 59873, 85443, *)    0.0000000035
1    (4, 653, 899, -)        0.0000000043
1    (15, 9874, 32345, +)    0.000000026
2    (2, 586, 910, *)        0.000000051
2    (5, 1274, 2421, -)      0.0000000016
2    (20, 35742, 39145, +)   0.00000057
......
```

- **Metadata**: {*expID, attribute, value*}

```
1       taxonomy        Homo sapiens
1       tissue          Brain
1       type            ChIP-seq
1       antibody        cMyc
2       taxonomy        Homo sapiens
2       tissue          Breast
2       type            ChIP-seq
......
```
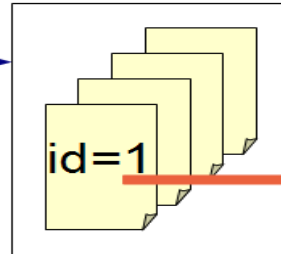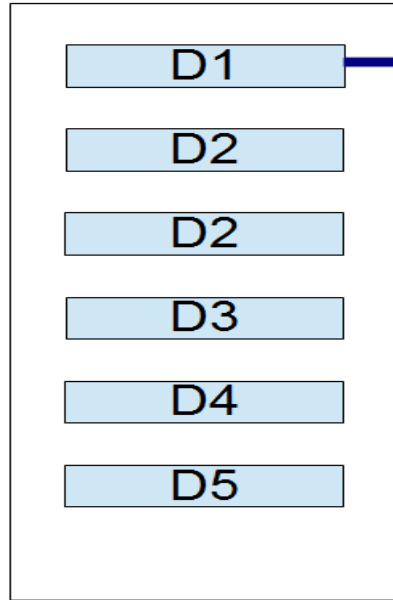
## Samples and datasets

- Every **sample** corresponds to an «experiment», with an ID
- Every **dataset** is a named collection of samples with the same region data schema

Data **format independent**; **interoperability** framework for comparing data samples based upon their biological aspects

**ENCODE NARROW** (or point source) **PEAK** format: It is used for called regions of signal enrichment based on pooled, normalized (interpreted) data, which usually represent genomic features.

```
chrom    chromStart chromEnd  name score  strand  signalValue   pValue   qValue peak
chr1     9356548    9356648   .    0      .       182           5.0945   -1     50
chr1     9358722    9358822   .    0      .       91            4.6052   -1     40
```

```xml
<?xml version = "1.0" encoding = "UTF-8" standalone = "yes"?>
<gdmSchemaCollection xmlns = "http://www.bioinformatics.deib.polimi.it/GDM/"
     name = "global_schemas">
  <gdmSchema type = "NARROWPEAK">
     <field type = "string">chr</field>   // Name of reference sequence chromosome or scaffold
     <field type = "long">left</field>   // Starting position of the feature in the chromosome or scaffold
     <field type = "long">right</field>   // Ending position of the feature in the chromosome or scaffold
     <field type = "string">name</field>   // Feature / region name ('.' if not assigned)
     <field type = "int">score</field>   // Feature score (how dark the region is shown in a genome
                                                browser (0-1000))
     <field type = "char">strand</field>   // Chromosome strand
     <field type = "double">signalvalue</field>   // Overall (usually, average) enrichment for the region
     <field type = "double">pvalue</field>   // Statistical significance (-log10) for the region (-1 if not
                                                assigned)
     <field type = "double">qvalue</field>   // Statistical significance using false discovery rate (-log10)
                                                for the region (-1 if not assigned)
     <field type = "int">peak</field>   // Point-source called for the region; 0-based offset from region left
                                                end (-1 if not assigned)
  </gdmSchema>
</gdmSchemaCollection>
```

```
(id,(chr,left,right,strand),(name,score,signalvalue,pvalue,qvalue,peak))
(1,("chr1",9356548,9356648,'*'),(".",0,182,5.0945,-1,50))
(1,("chr1",9358722,9358822,'*'),(".",0,91,4.6052,-1,40))
```

**VCF** (Variant Call Format) format: It is a flexible and extendable line-oriented text format developed by the 1000 Genomes Project for releases of single nucleotide variants, indels, copy number variants and structural variants.

```
CHROM    POS         ID              REF   ALT   QUAL     FILTER
22       16050075    .               A     G     100      PASS
22       16050678    rs139377059     C     T     100      PASS
```

```xml
<?xml version = "1.0" encoding = "UTF-8" standalone = "yes"?>
<gdmSchemaCollection xmlns = "http://www.bioinformatics.deib.polimi.it/GDM/"
    name = "global_schemas">
  <gdmSchema type = "VCF">
    <field type = "string">chr</field>   // Name of reference sequence (e.g., a chromosome) on which
                                          //       the variation is being called
    <field type = "long">left</field>    // Starting position of the variation on the given reference sequence
    <field type = "string">id</field>    // The identifier of the variation (e.g., a dbSNP rs identifier or "." if
                                          //       unknown)
    <field type = "string">ref</field>   // The reference base (or bases in the case of an InDel) at the
                                          //       given position on the given reference sequence
    <field type = "string">alt</field>   // The list of alternative alleles at the given position
    <field type = "int">qual</field>     // A quality score associated with the inference of the given alleles
    <field type = "string">filter</field>   // A flag indicating which of a given set of filters the variation has
                                             //       passed
  </gdmSchema>'
</gdmSchemaCollection>
```

```
(id,(chr,left,right,strand),(id,ref,alt,qual,filter))
(1,("22",16050075,16050075,'*'),(".","A","G",100,"PASS"))
(1,("22",16050678,16050678,'*'),("rs139377059","C","T",100,"PASS"))
```

**CpG Islands (UCSC)** annotations: They are regions of DNA where a cytosine nucleotide is followed by a guanine nucleotide in the linear sequence of bases along the 5' -> 3' direction, provided in a kind of **BED** (Browser Extensible Data) format.

```
Chrom  chromStart  chromEnd  name       length  cpgNum  gcNum  perCpg  perGc  obsExp
chr1   28735       29810     CpG: 116   1075    116     787    21.6    73.2   0.83
chr1   135124      135563    CpG: 30    439     30      439    13.7    67.2   0.64
```

```xml
<?xml version = "1.0" encoding = "UTF-8" standalone = "yes"?>
<gdmSchemaCollection xmlns = "http://www.bioinformatics.deib.polimi.it/GDM/"
     name = "global_schemas">
  <gdmSchema type = "CpG">
    <field type = "string">chr</field>    // Name of reference sequence chromosome or scaffold
    <field type = "long">left</field>    // Starting position of the feature in the chromosome or scaffold
    <field type = "long">right</field>    // Ending position of the feature in the chromosome or scaffold
    <field type = "string">name</field>    // CpG island name
    <field type = "long">length</field>    // Island length (right - left)
    <field type = "long">cpgNum</field>    // Number of CpGs in island
    <field type = "long">gcNum</field>    // Number of C and G in island
    <field type = "double">perCpG</field>    // Percentage of island that is CpG
    <field type = "double">perGC</field>     // Percentage of island that is C or G
    <field type = "double">obsExp</field>     // Ratio of observed (cpgNum) to expected (numC * numG /
                                                  length) CpGs in island
  </gdmSchema>
</gdmSchemaCollection>
```

```
(id,(chr,left,right,strand),(name,length,cpgNum,gcNum,perCpG,perGG,obsExp))
(1,("chr1",28735,29810,'*'),("CpG: 116",1075,116,787,21.6,73.2,0.83))
(1,("chr1",135124,135563,'*'),("CpG: 30",439,30,439,13.7,67.2,0.64))
```

**DNA-seq (mutations)**

(id, ('chr,start,stop,strand), (A,G,C,T,del,ins,inserted,ambig,Max,Error,A2T,A2C,A2G,C2A,C2G,C2T))

(1, (chr1,  917179, 917180,*), (0,0,0,0,1,0,'.','.',0,0,0,0,0,0,0,0))

(1, (chr1,  917179, 917179,*), (0,0,0,0,0,1,G,'.',0,0,0,0,0,0,0,0))


**RNA-seq (gene expression)**

(id, ((chr,start,stop,strand), (source,type,score,frame,geneID,transcriptID,RPKM1,RPKM2,iIDR))

(1, (chr8, 101960824, 101964847,-),   ('GencodeV10', 'transcript', 0.026615, NULL, 'ENSG00000164924.11', 'ENST00000418997.1', 0.209968, 0.193078, 0.058))


**Annotations**

(id, (chr,start,stop,strand), (proteinID,alignID,type))

(1, (chr1, 11873, 11873, +), ('uc001aaa.3', 'uc001aaa.3', 'cds'))
(1, (chr1, 11873, 12227, +), ('uc001aaa.3', 'uc001aaa.3', 'exon'))
(1, (chr1, 12612, 12721, +), ('uc001aaa.3', 'uc001aaa.3', 'exon'))
(1, (chr1, 13220, 14409, +), ('uc001aaa.3', 'uc001aaa.3', 'exon'))


**ChIA-PET** (denoting 3D genomic loops, head is assembled with coordinates, tail is in the schema)

(id,(chr,headstart,headstop,strand), (loopType, tailChr, tailStart, tailStop, PETcount, pValue, qValue))

(1, (chr1,7385626,7389841,*), ('Inter-Chromosome', chr17, 3081653, 3084755, 50, 0.0, 0.0)

# Query Language

(Motivational example and detailed description)

The language allows for queries on the genome involving <u>large datasets</u> describing:

- <u>Genomic signals</u>   (i.e. experiment dataset regions)
- <u>Reference regions</u>  (e.g. TSS, genes, promoters, enhancers)
- <u>Distance rules</u>     (e.g. the nearest enhancer that stands
  at least at 100 kb from the nearest gene)



Reference DNA

Gene

Enhancer

Promoter

Experimental dataset 1

Experimental dataset 2

Experimental dataset 3

Distance pattern

Reference regions

## Identification of distal bindings in transcription regulatory regions

Find all CTCF transcription factor (TF) binding regions of ChIP-seq data regarding human cancer cell line HeLa-S3, which are farther than x kb (e.g. 1000 kb) from the transcription start site (TSS) of the nearest gene.

Then, for the same cell line find all H3K4me1 histone modification (HM) regions that are the nearest regions farther than x kb from a TSS.

Finally, consider known enhancer (EN) regions and return a list of EN-HM-overlapping TF regions.



**HM:** Histone mark experiment
**TF:** Transcription factor experiment
**REF:** Reference DNA regions
**EN:** Enhancer
**x:** Threshold distance

DNA region

Not respecting the distance threshold

Not EN or HM overlapping

```
HM = SELECT(dataType == 'ChipSeq' AND cell == 'HeLa-S3'
                                   AND antibody == 'H3K4me1') PEAK;
TF = SELECT(dataType == 'ChipSeq' AND cell == 'HeLa-S3'
                                   AND antibody == 'CTCF') PEAK;
TSS = SELECT(type == 'TSS') ANNOTATION;
EN = SELECT(type == 'enhancer') ANNOTATION;

HMa = JOIN(distance > 1000000, minDistance(5); output: right) TSS HM;
TFa = JOIN(distance > 1000000, minDistance(5); output: right) TSS TF;
HMb = JOIN(distance < 0; output: int) EN HMa;
HMc = MERGE() HMb;
TF_res = JOIN(distance < 0; output: right) HMc TFa;
```



**HM:** Histone mark experiment

**TF:** Transcription factor experiment

**REF:** Reference DNA regions

**EN:** Enhancer

**x:** Threshold distance

DNA region

Not respecting the distance threshold

Not EN or HM overlapping

GenoMetric Query Language (GMQL) is defined as a sequence of algebraic operations following the structure:

**< variable > = < operator > (< parameters >) < variable >**

– Every variable is a dataset including many samples

– Offers high-level, declarative operations which operate both on regions and meta-data -> thus, each operation progressively builds the regions and meta-data of its result

– Inspired by SQL and *Pig Latin*

– Targeted towards cloud computing

## Classic relational operations – with genomic extensions

- SELECT, PROJECT, EXTEND, ORDER, GROUP, MERGE, UNION, DIFFERENCE

## Domain-specific genomic operations:

- COVER, (GENOMETRIC) JOIN, MAP

## Utilities:

- MATERIALIZE

Selection of the samples where a selection predicate p is true (e.g. select patients younger than 70 years)

```
S2 = SELECT(p) S1;
```

Example:   `S2 = SELECT(Patient_age < 70) S1;`

| 0.1 | 0.6 | Tumor_type = brca<br>Patient_age = 75 |

| 0.5 | 0 | Tumor_type = brca<br>Patient_age = 63<br>Gender = Female |

| 0.1 | 0 | 0.8 | 0.1 | Tumor_type = brca<br>Patient_age = 58 |

Selection of the regions where a selection predicate p is true (e.g. select those regions which have a score greater than 0.5)

$$S2 = SELECT(region: p) \ S1;$$

Example:   $S2 = SELECT(region: score > 0.5) \ S1;$



Tumor_type = brca
Patient_age = 75

Tumor_type = brca
Patient_age = 63
Gender = Female

Tumor_type = brca
Patient_age = 58

Projection of the regions: for each gene in a set, take its promoter (e.g. from -2kbp, to +1kbp from the TSS)

```
S2 = PROJECT(p) S1;
```

Example:   `S2 = PROJECT(region_update:`
`start = start – 2000, stop = start + 1000) S1;`



S1    Tumor_type = brca
      Patient_age = 75

S2    Tumor_type = brca
      Patient_age = 75

Count the regions in each sample and store it in metadata

```
S2 = EXTEND(p) S1;
```

Example: `S2 = EXTEND(Region_count AS COUNT()) S1;`



Tumor_type = brca
Patient_age = 75
**Region_count = 3**

Tumor_type = esca
Patient_age = 78
**Region_count = 5**

Tumor_type = chol
Patient_age = 85
**Region_count = 2**

POLITECNICO
DI MILANO

Order by region_count metadata and take the top 2 samples

$$S2 = ORDER(Ai; [TOP: k]) S1;$$

Example: `S2 = ORDER(`Region_count; `TOP: 2) S1;`

| 2 | 2 | | 5 | 10 | 3 | |

Tumor_type = esca
Patient_age = 78
Region_count = 5
**Order = 1**

| 5 | 1 | | 0 | |

Tumor_type = brca
Patient_age = 75
Region_count = 3
**Order = 2**

| 5 | | 3 | |

Tumor_type = chol
Patient_age = 85
Region_count = 2
**Order = 3**

Group samples according to the value of tumor and compute the region minimum score of each group



Tumor_type = brca
Patient_age = 75
**Group = 1**
**Min = 0**

Tumor_type = esca
Patient_age = 78
**Group = 2**
**Min = 1**

Tumor_type = esca
Patient_age = 78
**Group = 2**
**Min = 1**

Tumor_type = chol
Patient_age = 87
**Group = 3**
**Min = 3**

Collapse a bunch of samples (both region and metadata) into an unique one        `S2 = MERGE() S1;`



| S1.s1 | Type = ChipSeq<br>Antibody = CTCF<br>Replicate = 1 |
| S1.s2 | Type = ChipSeq<br>Antibody = CTCF<br>Replicate = 2 |
| S1.s3 | Type = ChipSeq<br>Antibody = CTCF<br>Replicate = 3 |
| S2 | Type = ChipSeq<br>Antibody = CTCF<br>Replicate = 1<br>Replicate = 2<br>Replicate = 3 |

Return a single dataset with all the samples in two input datasets, merging their region attributes if different

```
S3 = UNION() S1 S2;
```



S1 — Tumor_type = brca
Experiment = rnaseq

S2 — Tumor_type = brca
Experiment = mirna

S3 — Tumor_type = brca
Experiment = rnaseq

Tumor_type = brca
Experiment = mirna

Return all the regions in the first dataset that do not overlap any region in the second one

$$S3 = DIFFERENCE() \ S1 \ S2;$$

S1

Tumor_type = brca
Experiment = rnaseq

S2

Tumor_type = brca
Experiment = mirna

S3

S1.Tumor_type = brca
S1.Experiment = rnaseq
S2.Tumor_type = brca
S2.Experiment = mirna

- Produces new regions where there are between MIN and MAX regions of the operand dataset

$$\texttt{S2 = COVER(min, max) S1;}$$

COVER(ALL,ALL) AND     COVER(1,ANY) OR     COVER(2,ANY)

- ALL: number of samples in the dataset
- Jaccard indexes can be used instead of min-max
- An aggregate function *f* can be computed for regions forming the cover

COVER(2, ANY): find portions of the genome that are covered by at least two regions

S2 = COVER(2, ANY) S1;

S1.s1

Tumor_type = brca
Tumor_grade = g3

S1.s2

Tumor_type = brca
Tumor_grade = g2

S1.s3

Tumor_type = brca
Tumor_grade = g2

S2

2    3 2    1  2  1      1

Tumor_type = brca
Tumor_grade = g2
Tumor_grade = g3

- Given two sets of samples, JOIN builds the pairs of regions and metadata where a join predicate p is true.
- Region of results are composed from regions of the operands

```
S3 = JOIN(p, comp-op) S1 S2;
```

- Functions *minDistance* and *distance* can be used in the predicate

Metadata join: select pairs of matching samples (e.g. with the same "Type")

Type = **uvm**
Patient = 123
Gender = M

Type = **brca**
Patient = 211

Type = **brca**
Patient = 10
Age = 88

Type = **sarc**
Patient = 12

Type = **brca**
Patient = 333
Grade = g3

Type = **sarc**
Patient = 444
Age = 88

Type = **sarc**
Patient = 12

Join at min-distance: associate each region in the former dataset with the closest in the latter

```
S3 = JOIN(MINDISTANCE(1); output: RIGHT) S1 S2;
```



S1 — feature = transcripts

S2 — feature = TFBSs

S3 — S1.feature = transcripts
S2.feature = TFBSs

```
S2 = JOIN(DISTANCE < 1000; output: CAT) R, S1;
```



R

$d = 2500$

$d = 1100$

$d = 400$

$d = 0$

All pairs

3    5

e1

R

Matching pairs and region composition

3    5

e1

3    5

e2

- Computes aggregate functions over samples of S1 which intersect with the regions of R

```
S2 = MAP(newAttr AS MIN(attr)) R S1;
```

Compute an aggregate function (e.g. AVERAGE) on all the regions intersecting the reference

```
S2 = MAP(average_score AS AVG(score)) R S1;
```

COUNT is computed by default

```
S2 = MAP() R S1;
```



R — annotation = genes / provider = RefSeq

S1 — feature = SNP

S2 — R.annotation = genes / R.provider = RefSeq / S1.features = SNP

- **MAP** operations, through <u>reference regions</u> **R**, <u>extract</u> and <u>standardize</u> **genomic features** expressed in distinct datasets

GMQL MAP



| | $R_1$ | | | $R_2$ | | $R_3$ | |
|---|---|---|---|---|---|---|---|
| DHS | | | | | | | |
| RNAPII | | | | | | | |
| H3K4me1 | | | | | | | |
| ... | | | | | | | |
| Gene | | | | A | | | |
| TSS | | | | | | | |
| Enhancer | E | | | | | | |

**Genomic Space**

| | $E_1$ | $E_2$ | $E_3$ | | $E_m$ | $E_{m+1}$ | $E_{m+2}$ | $E_{m+3}$ |
|---|---|---|---|---|---|---|---|---|
| $R_1$ | 10 | 3 | 2 | | 12 | 0 | 1 | 0 |
| $R_2$ | 1 | 48 | 12 | | 3 | 0 | 0 | 0 |
| $R_3$ | 11 | 9 | 10 | ... | 0 | 1 | 0 | 1 |
| ... | | | | | ... | | | |
| $R_{n-1}$ | 56 | 47 | 1 | | 6 | 1 | 0 | 1 |
| $R_n$ | 46 | 3 | 21 | | 13 | 1 | 0 | 0 |

- **Genome Space**: simplified structured outcome, ideal format for data analysis

GMQL MAP



| | | R₁ | | R₂ | | R₃ | |

- **Genometric spaces** represent <u>adjacency matrices</u>, i.e. networks
  - <u>Network analysis methods</u> (e.g. page rank, hub/authority, community detection, …)

```
Res = MAP(count_name: mutCount) Genes Dataset;
```



Dataset

Genes

Res

It requires:
- Partitioning by experiment classes
- Adding names to regions and to experiments (from metadata)
- Adding colors

**GenoMetric Space Explorer**: http://www.bioinformatics.deib.polimi.it/GeMSE/

| Attribute ▲ | Value | P_0 | P_1 | P_2 | P_3 | P_4 | P_5 | P_6 | P_7 |
|---|---|---|---|---|---|---|---|---|---|
| Disease | Burkitt lymphoma | 5 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| Disease | Chronic lymphocytic leukemia | 5 | 3 | 2 | 2 | 2 | 1 | 1 | 0 |
| Drug | Tetracycline | 10 | 5 | 3 | 3 | 3 | 2 | 2 | 1 |
| Exposure Time | 30min | 7 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| Exposure Time | 6h | 3 | 4 | 3 | 3 | 2 | 2 | 2 | 1 |
| Transcription Factor | GATA1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| Transcription Factor | GATA2 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Transcription Factor | EZH2 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| Transcription Factor | NANOG | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| Transcription Factor | RUNX2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Transcription Factor | EP300 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| Transcription Factor | IRF4 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| Transcription Factor | c-Myc | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Transcription Factor | RAD21 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Transcription Factor | c-Jun | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| Transcription Factor | P53 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| Transcription Factor | c-Fos | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| Transcription Factor | CNTN2 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| Transcription Factor | SRC | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| Transcription Factor | MAX | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 |

For biological/clinical interpretation of genomic data processing, and data stratification based on of biological/clinical metadata values and/or patterns of different genomic feature regions

# Implementation

(Ver. 1 & Ver. 2)

- GMQL similar to *Pig Latin* (by *Yahoo! Research*)
  - Algebraic language for data-intensive applications on *Apache Hadoop, a* framework for parallel computing which executes *Google MapReduce* programs
- Implementation strategy: develop a **translator** to *Pig Latin*
  - Easier development and maintenance
  - Big company involvement ensures development
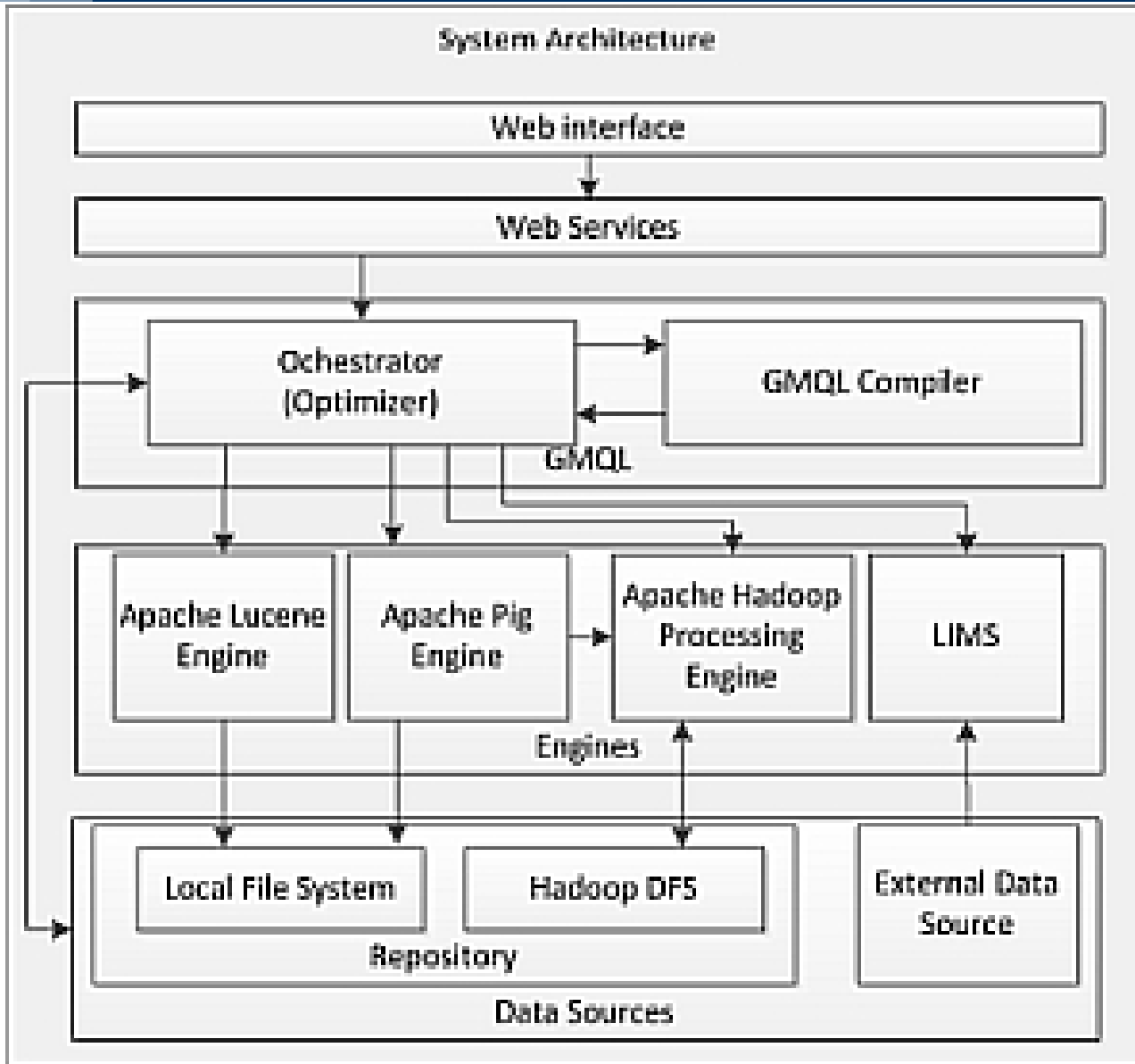  - Use cloud computing power to obtain efficiency and scalability

Repository Architecture

Main Server Local System

Public — Same control data as private users

Private — User 1: GMQL history, Schema, Metadata, MetaData Index, DataSet XML

Hadoop Distributed File System (HDFS)

Public — Meta Data: BED, Nar; Samples

Private — user1: Samples, Meta data; User#: Samples, Meta data

**GMQL query**          **Translator**          **Pig over Hadoop**

```
1  JUN_2 = SELECT (antibody == 'c-Jun' AND
2                  dataType = 'ChipSeq' AND
3                  treatment == 'None') ENCODE_PEAK;
4  |
5  GENES = PROJECT (feature == 'gene') GENE_NN;
6
7  mutations_number = MAP (mutCount: count) exon_regions
8                                           dna_datasets
9
```

Motivation:
- Clear & compact user code
- User-transparent optimization

```
HM = SELECT (dataType=='ChipSeq' AND cell=='HeLa-s3'
                AND antibody=='H3K4me1') ENCODE_BROAD;

TF = SELECT (dataType=='ChipSeq' AND cell=='HeLa-s3'
                AND antibody=='CTCF') ENCODE_BROAD;

TSS = SELECT (type=='TSS') ANNOTATION;

EN = SELECT(type=='enhancer') ANNOTATION;

HMa = JOIN (minDistance AND distance > 1000, right) TSS HM;

TFa2 = JOIN (distance < 0, right) EN TF;

TFb2 = JOIN (minDistance AND distance>1000,right) TSS TFa2;

TF_res = JOIN (distance < 0, right) HMa TFb2;

MATERIALIZE TF_res;
```

```
TSS_meta_group = group TSS_meta by id parallel 8;

HM_meta_group = group HM_meta by id parallel 8;

TSS_HM_exp_cross = foreach (cross TSS_exp_group ,
        HM_exp_group parallel 8)
        generate ($0, $2), ($1, $3);

TSS_HM_meta_cross = foreach (cross TSS_meta_group ,
        HM_meta_group parallel 8)
        generate ($0, $2), ($1, $3);

define HMa_joiner
        GenometricPig.MinDistancePlusDistanceJoin
        ('1000 1000 r2');

TSS_HM_meta_cross_flattened =
        foreach TSS_HM_meta_cross generate $0,
        flatten($1);

HMa_meta = foreach (union
        (foreach TSS_HM_meta_cross_flattened
        generate GenometricPig.NewId($0.$0, $0.$1),
        flatten($1)), (foreach TSS_HM_meta_cross_flattened
        generate GenometricPig.NewId($0.$0, $0.$1),
        flatten($2))) generate $0 as id, $2
        as attribute, $3 as value;

HMa_exp = foreach (foreach
        (foreach TSS_HM_exp_cross generate HMa_joiner($1))
        generate flatten($0)) generate $0 as id, $1
                as region, $2 as value;
```

- 1 statement => 25 *Pig Latin* lines of code + auxiliary Java function
- The translator takes also care of updating the variable schema
- Error handling

| GQL Operator | *PigLatin* Translation |
|---|---|
| S2 = SELECT (p) S1 | s2_pred = group S1_META by id;<br>s2_ids = foreach (filter S1_pred::$1 by $\tau(p)$) generate id;<br>S2_META = foreach(join S1_META by id, s2_ids by id)<br>    generate S1_META.id, S1_META.attribute, S1_META.value;<br>S2_EXP = foreach (join S1_EXP by id, s2_ids by id) generate id, region, value; |
| S2 = PROJECT(p,[f1,f2]) S1 | S2_EXP = foreach (filter S1_EXP by $\tau(p)$) generate id,(chr,$\tau$(f1),$\tau$(f2),strand),value; |
| S3 = JOIN ($P_R \wedge P_M$,o) S1, S2 | s3temp = foreach (filter (cross S1_EXP, S1_EXP) by $\tau(P_R)$)<br>    generate (S1_EXP::id as id1, S2_EXP::id as id2) as ids,<br>    ($\tau$(o)(S1_EXP::region, S2_EXP::region)), (S1_EXP::value, S2_EXP::VALUE);<br>s1pred = group S1_META by id;<br>s2pred = group S2_META by id;<br>s3pair = foreach (filter (cross s1_pred, s2_pred) by $\tau(P_M)$)<br>    generate (s1pred::id as id3, s2pred::id as id4) as ids;<br>s3_quad = foreach (join s3temp by ids, s3pair by ids)<br>    generate $\tau(new(ids))$ as id, s3temp.ids, s3_pair.id1 as id1, s3_pair.id2 as id2;<br>s3_EXP = foreach (join s3temp by ids, s3quad by ids) generate s3quad.id, s3temp.$1, s3temp.$2;<br>s3_META = union<br>    (foreach (join S1_META by id, s3quad by id1)<br>     generate s3quad.id, S1_META.attribute, S1_META.value)<br>    (foreach (join S2_META by id, s3_quad by id1)<br>     generate s3quad.id, S2_META.attribute, S2_META.value); |
| S2 = MAP (F1:f1,...,Fn:fn) R, S1 | space = foreach (cross (distinct (foreach S1_EXP generate id)),R) generate S1::id, R::region;<br>not_null = foreach (filter (cross R, S1_EXP) by $\tau$(intersect)(R::region, S1_EXP::region))<br>    generate S1_EXP::id, R::region, S1_EXP::value;<br>S2_EXP = foreach (group (join space by region left, not_null by region) by region)<br>    generate $0::id, $0::region, f1($1) as F1,....,fn($1) as Fn; |
| R = EXTRACT S | R = foreach S_EXP generate region; |
| R = EXTRACT (<br>    cover (N [,K]) ) S | *Translation is for* $N \leq 2$, *COVER code can be automatically generated for arbitrary N.*<br>R1 = foreach S generate region;<br>f = filter (cross R1, R1) by $0::chr == $1::chr<br>    and ($0::left < $1::right or ($0::left == $1::left and $0::right < $1::right));<br>i = foreach f generate {<br>    ($0, ($1<$2?$1:$2), 1), ($2, ($1<$3?$1:$3), ($1>$2?2:0)),<br>    (($1>$2?($1>$3?$3:$1):$2), ($1>$3?$1:$3),1)};<br>cover_count = foreach i generate flatten($0);<br>R = filter cover_count by $3 >= N [and $3 <= K]; |
| RN = UNION R1,R2 [,Ri...] | RN = union R1,R2 [,Ri...]; |
| SN = UNION S1, S2 [,Si...] | SN_EXP = union S1_EXP, S2_EXP [, Si_EXP];<br>SN_META = union S1_META, S2_META [, Si_META]; |

1. Parallelism by splitting computations:
   - By chromosome
   - By experiment
2. Join and Map have a translation which avoids cross products, based on sequential scan of regions

   **Pig Latin** shows its ability to scale on hundreds or thousands of experiments and multi-node systems
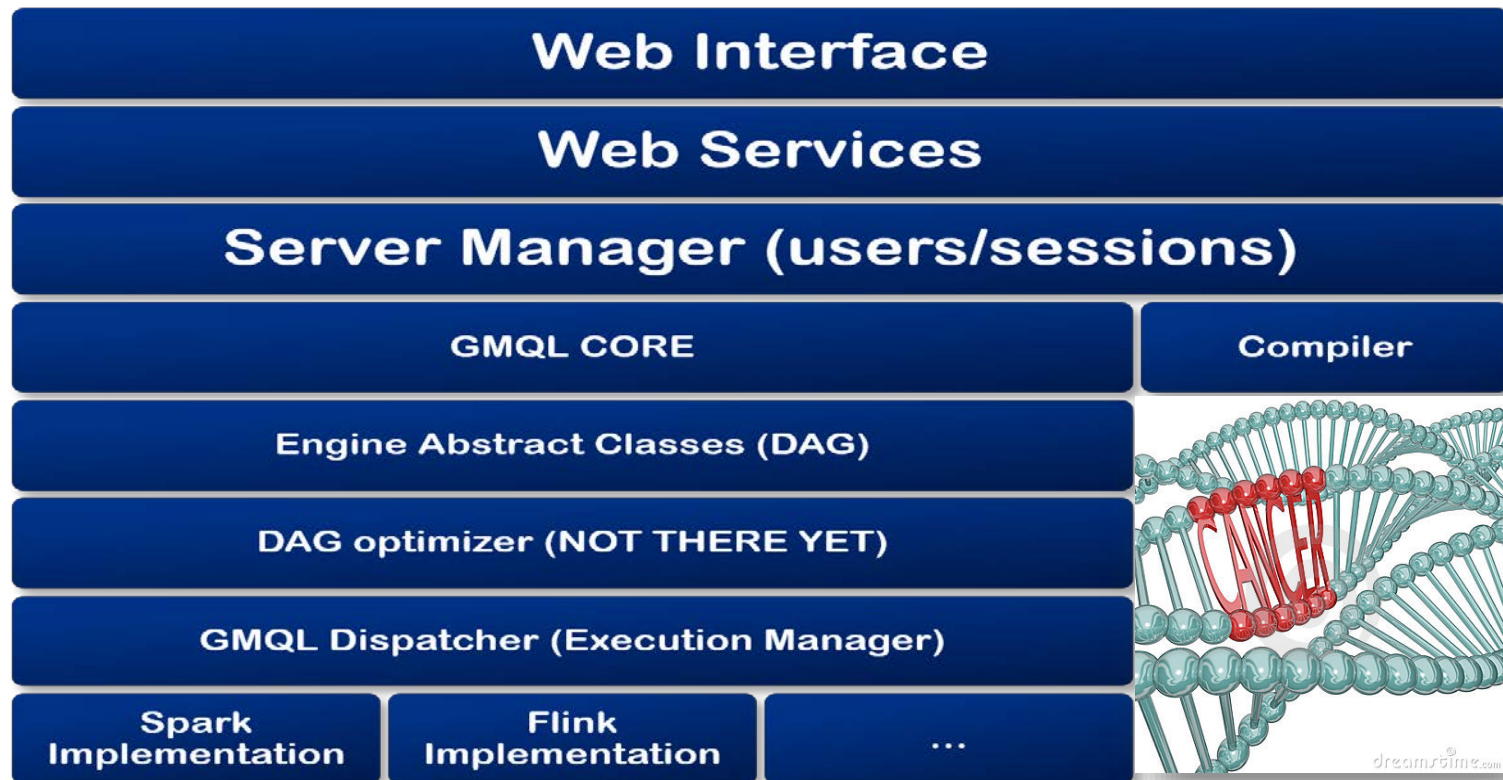
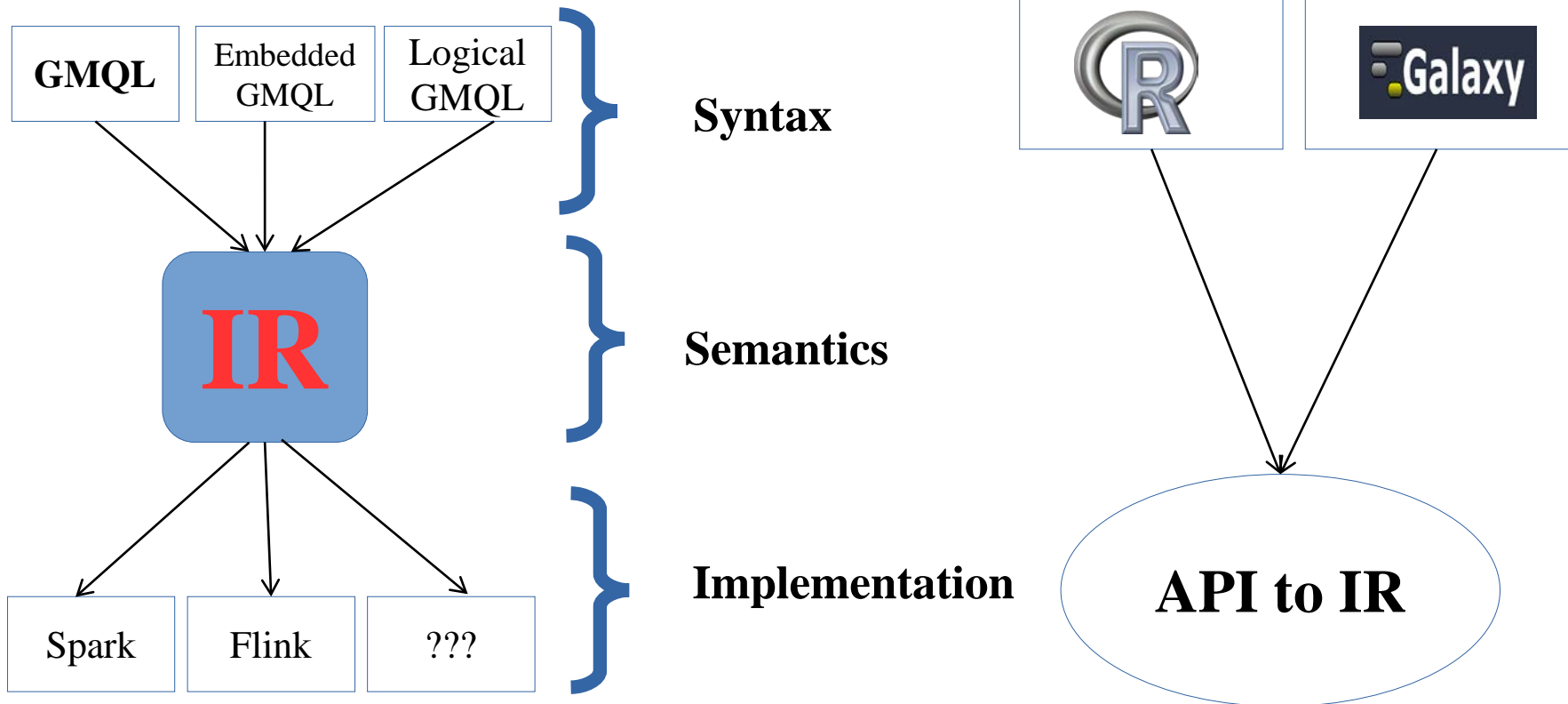- Holistic data management system for genomics
- Uses cloud-based computing for querying thousands of heterogeneous datasets

- A different approach, with language-independent intermediate representation

- Targeting also usability from within R and Galaxy

```scala
import it.polimi.genomics.GMQLServer.GmqlServer
import it.polimi.genomics.core.DataStructures.CoverParameters.{CoverFlag, N}
import it.polimi.genomics.spark.implementation.GMQLSparkExecutor
import it.polimi.genomics.spark.implementation.loaders.test3Parser
import org.apache.spark.{SparkContext, SparkConf}

object Cover {

  def main(args : Array[String]) {

    val conf = new SparkConf()
    val sc:SparkContext =new SparkContext(conf)

    val server = new GmqlServer(new GMQLSparkExecutor(sc=sc))

    val ex_data_path = "/home/abdulrahman/Desktop/datasets/coverData/"
    val output_path = "/home/abdulrahman/testCover/res/"


    val dataAsTheyAre = server READ ex_data_path USING test3Parser()

    val cover = dataAsTheyAre.COVER(CoverFlag.COVER, N(2), N(3), List(), None )

    server setOutputPath output_path MATERIALIZE cover

    server.run()

  }

}
```

## *Example*

- New optimization options

```
┌─────────┐      ┌─────────┐      ┌───────────┐      ┌──────────────────┐
│  GMQL   │  ⟹  │   IR    │  ⟹  │ Optimized │  ⟹  │  Implementation  │
│         │      │         │      │    IR     │      │                  │
└─────────┘      └─────────┘      └───────────┘      └──────────────────┘
                                        ⬆                     ⬆
                                 ┌───────────┐         ┌──────────────┐
                                 │   Query   │         │  Low Level   │
                                 │ Optimizer │         │  Optimizer   │
                                 └───────────┘         └──────────────┘
```

1) Node reordering / deletion
2) Select condition refinement

1) Alternative algorithms
2) Parallelism tuning
3) Data partitioning
4) Caching

**Idea**:

- Let Flink/Spark/… engines implement common and well known optimization

- Exploit the intermediate representation in order to implement optimizations which are driven by the semantics of GMQL

  - Meta-first optimization

  - Operator swapping optimization

- Other optimizations based on algorithms for <u>parallel execution</u> on the cloud

Under certain conditions (meta-separability), it is possible to compute the metadata side of the query strictly before the region data side.

GMQL queries are always meta-separable, except for the ones which use the EXTEND operator

(EXTEND operator computes and aggregates on the region data and stores the result in the metadata)

- Compute metadata side of the query
- Retrieve the IDs from the metadata result
- Use the IDs to selectively load only the files that will appear in the output

Affected queries are the ones which contain one or more metadata selection (far from the Readings), metadata join and metadata group by; those operations cut the size of the output

Some reordering of the execution plan can not be inferred by lower level optimizer, since they are motivated by GMQL semantics

Bin1    Bin2    Bin3    Bin4    Bin5    Bin6    Bin7

Strategy for intersection:
1. Partition the genome in bins
2. Assign each region to all the bins it overlaps
3. Search for intersections within each bin

In the case of more complex operations, we change the way in which the regions are assigned to the bins

**Avoiding output duplicates**:



In order to avoid the duplicates production, when two regions overlap, an output is emitted if, and only if, **at least one of them begins in the considered bin**

- Bin 2: overlap => red region begins => **Output**
- Bin 3: overlap => no region begins => **Output not emitted!**

- **Smaller bins**: smaller search space, but higher number of replicates
- **Optimal binning size depends on**:
  - Number of regions and local density
  - Region length distribution
  - GMQL operation and parameters
  - System settings (e.g., number of nodes, amount of memory, …)

# Repository

Stores experimental datasets and annotations collected from external databases

- ENCODE (more than 4000 processed datasets for humans and mices, relevant to epigenomic research)

- Roadmap Epigenomics (about 1000 human epigenomic datasets for stem cells and ex-vivo tissues)

- TCGA (The Cancer Genome Atlas, providing more than 50,000 processed datasets for more than 30 cancer types, including mutations, copy number variations, gene and miRNA expressions, methylations)

Annotation data are also extracted from external references, based upon the needs of given research projects

- Genes (UCSC, RefSeq, Ensembl, GENCODE)

- Transcription Start Sites (SwitchGear)

- Transcription Factor Binding Sites (UCSC, ENCODE)

- CpG islands (UCSC)

- miRNA target sites (UCSC)

- Enhancers (Vista)

| Consortium | Imported datasets | # of samples | File size (MB) |
|---|---|---|---|
| **ENCODE** | GRCh38_ENCODE_BROAD | 850 | 6,869 |
| | GRCh38_ENCODE_NARROW | 11,573 | 128,316 |
| | HG19_ENCODE_BROAD | 844 | 18,382 |
| | HG19_ENCODE_NARROW | 10,342 | 111,925 |
| **ROADMAP EPIGENOMICS** | HG19_ROADMAP_EPIGENOMICS_BED | 156 | 968 |
| | HG19_ROADMAP_EPIGENOMICS_BROAD | 979 | 24,332 |
| | HG19_ROADMAP_EPIGENOMICS_DMR | 66 | 3,060 |
| | HG19_ROADMAP_EPIGENOMICS_GAPPED | 979 | 6,875 |
| | HG19_ROADMAP_EPIGENOMICS_NARROW | 1,032 | 11,788 |
| | HG19_ROADMAP_EPIGENOMICS_RNA_expression | 399 | 2,453 |
| **TCGA** | HG19_TCGA_cnv | 22,632 | 797 |
| | HG19_TCGA_dnamethylation | 12,860 | 247,742 |
| | HG19_TCGA_dnaseq | 6,914 | 286 |
| | HG19_TCGA_mirnaseq_isoform | 9,909 | 4,207 |
| | HG19_TCGA_mirnaseq_mirna | 9,909 | 746 |
| | HG19_TCGA_rnaseq_exon | 3,675 | 47,668 |
| | HG19_TCGA_rnaseq_gene | 3,675 | 5,327 |
| | HG19_TCGA_rnaseq_spljxn | 3,675 | 44,377 |
| | HG19_TCGA_rnaseqv2_exon | 9,825 | 124,343 |
| | HG19_TCGA_rnaseqv2_gene | 9,825 | 21,862 |
| | HG19_TCGA_rnaseqv2_isoform | 9,825 | 53,082 |
| | HG19_TCGA_rnaseqv2_spljxn | 9,825 | 115,088 |
| **GDC - TCGA** | GRCh38_TCGA_copy_number | 22,374 | 686 |
| | GRCh38_TCGA_copy_number_masked | 22,375 | 337 |
| | GRCh38_TCGA_gene_expression | 11,091 | 56,542 |
| | GRCh38_TCGA_methylation | 12,218 | 1,348,516 |
| | GRCh38_TCGA_miRNA_expression | 10,947 | 1,502 |
| | GRCh38_TCGA_miRNA_isoform_expression | 10,999 | 5,004 |
| | GRCh38_TCGA_somatic_mutation_masked | 10,188 | 2,280 |
| **GENCODE** | GRCh38_ANNOTATION_GENCODE | 24 | 1,798 |
| | HG19_ANNOTATION_GENCODE | 20 | 1,324 |
| **REFSEQ** | GRCh38_ANNOTATION_REFSEQ | 31 | 740 |
| | HG19_ANNOTATION_REFSEQ | 30 | 275 |
| **Grand total** | 33 datasets | 240,066 | 2,399,497 |

# User Interface

POLITECNICO DI MILANO

GMQL    GMQL-REST    Demo Video    Documentation

Hello Marco Masseroli    Logout

**Data sets**

+ Private
+ Public

⊕ Add    🗑 Delete    ⊕ Download    ☁ UCSC G. B.

**Query editor**    Select query ▼

```
1
```

**Query name**    fileName

**Output Format**    ⦿ GTF    ◯ Tab Delimited

⚙ Compile    ▶ Execute

☰ Show jobs

**Metadata browser**

```
DATA_SET_VAR = SELECT() HG19_BED_ANNOTATION;
```

➕ Add new condition    ▤ Test condition

**Sample Metadata**

**Schema**

**Schema type:** bed

| Field name | Field type | Heat map |
|------------|------------|----------|
| chr | STRING | |
| start | LONG | |

# Results are provided to user in GTF or Tab-delimited format

https://pygmql.readthedocs.io/en/latest/
https://bioconductor.org/packages/release/bioc/html/RGMQL.html

Integrated environments where the bioinformatician can:

- Run GMQL queries on **local** or **remote** data
- Integrate the results with **external libraries** of Python or R/Bioconductor
- Visualize the results

☁ Use FireCloud   </> API   📖 User Guide   📢 **Blog**   💬 Forum   📅 Events

← A PRODUCTIVE HACKATHON: MAKING DATA MORE...

🏠

FEATURED WORKSPACES | WHAT ARE THEY, WHA... →

## New Featured workspace showcasing the GenoMetric Query Language

Posted by Tiffany_at_Broad on 1 Jun 2018                                    💬 (0)

We are excited to introduce a new Featured workspace that demonstrates the GenoMetric Query Language (GMQL) created by a team from Politecnico di Milano in Italy. For some context on Featured workspaces, please read our previous blog post.

GMQL is a high-level, declarative language supporting queries over thousands of heterogeneous datasets and samples; as such, it enables genomic "big data" analysis. Based on Hadoop framework and the Apache Spark platform, GMQL is designed to be highly scalable, flexible, and simple to use. You can try the system here through its several interfaces, with documentation and biological query examples on ENCODE, TCGA and other public datasets or clone the Featured workspace and launch an example analysis.

The **GMQL 101 workspace** features three methods, each with increasing levels of complexity to give you a taste of how the query language works. One method shows how to join two datasets, and then extracts a third dataset based on a specific condition: pairs of regions that are less than 1000 bases a part. The second method takes a VCF and performs an epigenomic analysis using gene annotation and Chip-Seq results. It shows how you can select high confidence regions, use RefSeq annotations to find regions that overlap a gene, and count the mutations falling within the high confidence regions. Finally, the third method is a combination of GATK4's Mutect 2 pipeline and the second method, showing an epigenomic analysis from start (calling somatic variants) to finish (annotating variants). For any GMQL-specific questions or problems you can visit the GMQL GitHub page.

**Many thanks to Luca Nanni, Arif Canakoglu, Pietro Pinoli, and Stefano Ceri** for putting together this workspace. It takes a lot of thought and effort to create a valuable learning resource like this, and we are still figuring out the most successful way to do this. Please share your thoughts in the Comments section below on the effectiveness of this workspace and any other Featured workspaces you try out. If you are interested in featuring examples of your methods in this way, please tell us here, and we can talk to you about the process.

# Applications

**Source**: ENCODE ChIP-seq datasets for transcription factors (TF)

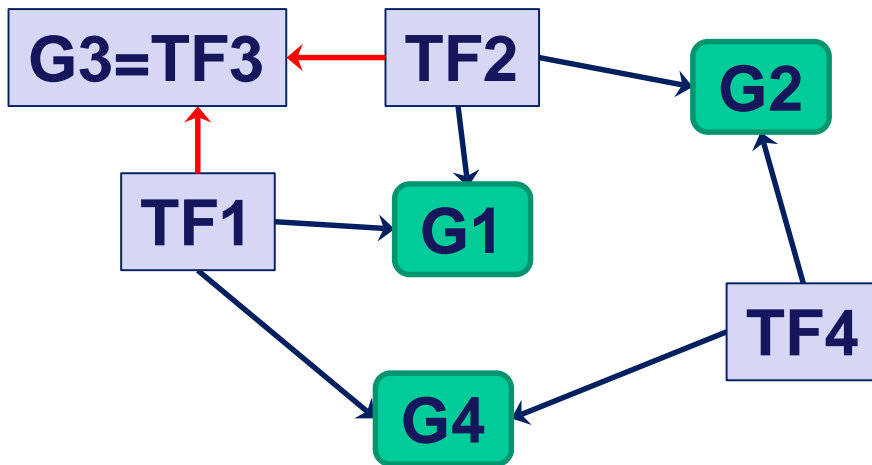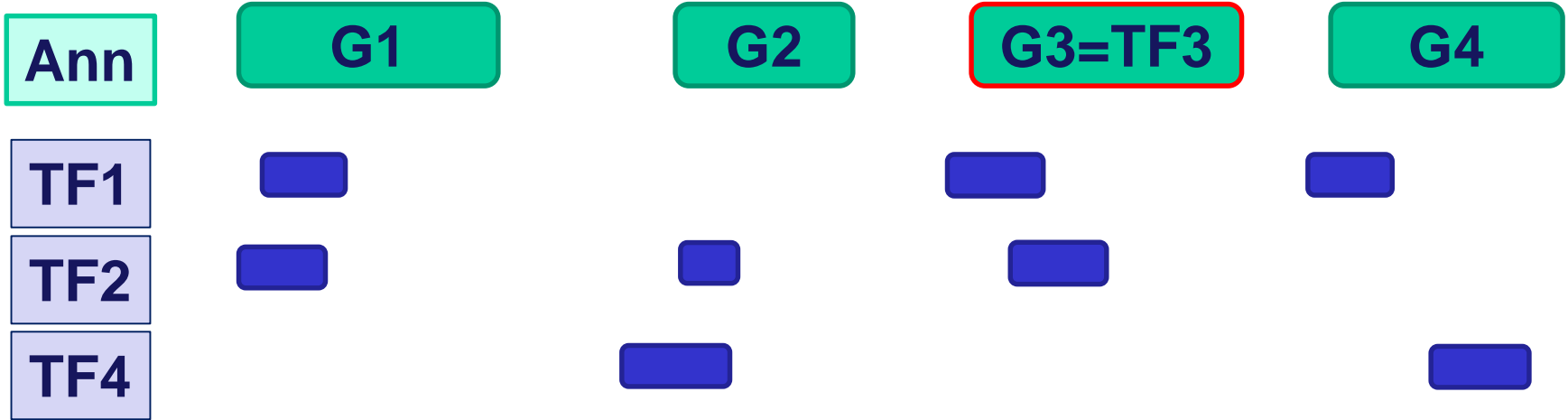**Goal**: Generation of a transcriptional network from ChIP-seq data

**Method**: Select TFs and genes, and derive TFs-Genes links. Build a TFxGenes matrix $M_{ij}$ such that $M_{ij}$ = number of binding sites of TF i in the gene region j, $M_{ij} = 0$ if no binding.

```
# Extract gene information, some of them tagged with TF-encoding
GENE_ANN = SELECT(type == 'gene' AND provider == 'RefSeq')
                                              ANNOTATION;
GENES = SELECT(region: feature == 'gene') GENE_ANN;
TF_GENES = SELECT(region: encode_TF == 'yes') GENES; # red in
                                                 next slide
# Collect TF samples (122)
TF = SELECT(dataType = 'ChipSeq' AND subType = 'TF' AND
        cell == 'k562' AND treatement == 'None') ENCODE_PEAK;
# Build TF Genomic Space
GS_TF = MAP(count_name: binding_num) GENES TF;
```

| Ann | G1 | G2 | G3=TF3 | G4 |
|-----|----|----|--------|----|
| TF1 | ■ | | ■ | ■ |
| TF2 | ■ | ■ | ■ | |
| TF4 | | ■ | | ■ |

Cell line: **K562**
(CML)
# TFs:      **122**
# Nodes:   **6240**
# Edges: **30587**

```
# Collect open chromatin samples

DHS_2 = SELECT(dataType == 'DnaseSeq' AND
               cell == 'k562') ENCODE_PEAK;     # 2 samples

FAIRE = SELECT(dataType == 'FaireSeq' AND
               cell == 'k562') ENCODE_PEAK;     # 1 sample


# Merge DHS replicates in one sample

DHS = FLAT(2, ANY; aggregate: pValue AS MIN(pValue))
                                               DHS_2;

# Merge open chromatin regions from DHS and FAIRE assays

DHS_FAIRE = UNION() DHS FAIRE;

OPEN = COVER(1, ANY; aggregate: pValue AS MIN(pValue))
                                               DHS_FAIRE;

# Extract TFs in open chromatin regions only (active DNA
# binding)

TF_OPEN = JOIN(distance < 0; output: left) TF OPEN;
```

\# TF Genomic Space

**GS_TF_0 = MAP(count_name: binding_num) GENES TF_OPEN;**

| | G1 | G2 | G3 TF3 | G4 | G5 | … | Gn |
|---|---|---|---|---|---|---|---|
| TF1 | 1 | 0 | 1 | 1 | 0 | … | 1 |
| TF2 | 1 | 0 | 1 | 0 | 0 | … | 1 |
| TF3 | 0 | 0 | 0 | 0 | 1 | … | 0 |
| … | 0 | 0 | 0 | 1 | 0 | … | 1 |
| TFn | 0 | 0 | 0 | 1 | 1 | … | 1 |

**GS_TF = SELECT(binding_num > 0) GS_TF_0;**

Cell line: **K562** (CML)

\# TFs:        **95**

\# Nodes: **1717**

\# Edges: **2367**

# Summary & Outlook

- **GDM**: a data <u>format-independent</u> genomic data model
  - For <u>genomic region</u> data and related <u>metadata</u>
  - Easing <u>integration</u> and processing of <u>heterogeneous</u> genomic data

- **GMQL**: a high-level <u>declarative</u> language
  - <u>Easing</u> the expression of even <u>complex queries</u> on <u>numerous data</u> of multiple <u>different types</u>
  - Running also on <u>cloud computing</u> environments
  - Supporting a <u>first processing</u> also of <u>big data</u>, to extract the relevant (usually smaller) ones for <u>further processing</u>

- **Several GDM & GMQL application examples**
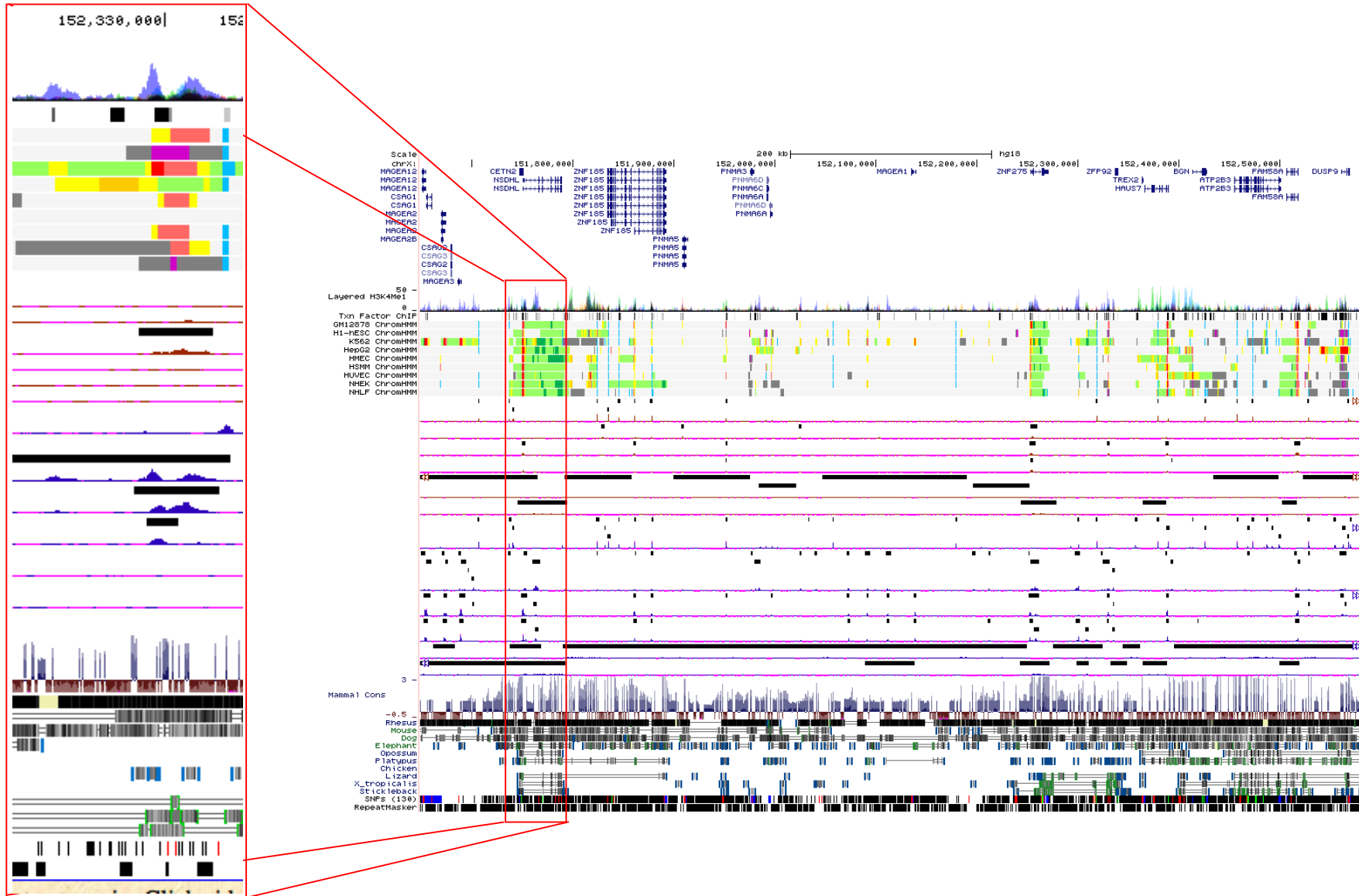  - <u>Characterizing</u> interplay and function of <u>genomic regions</u>
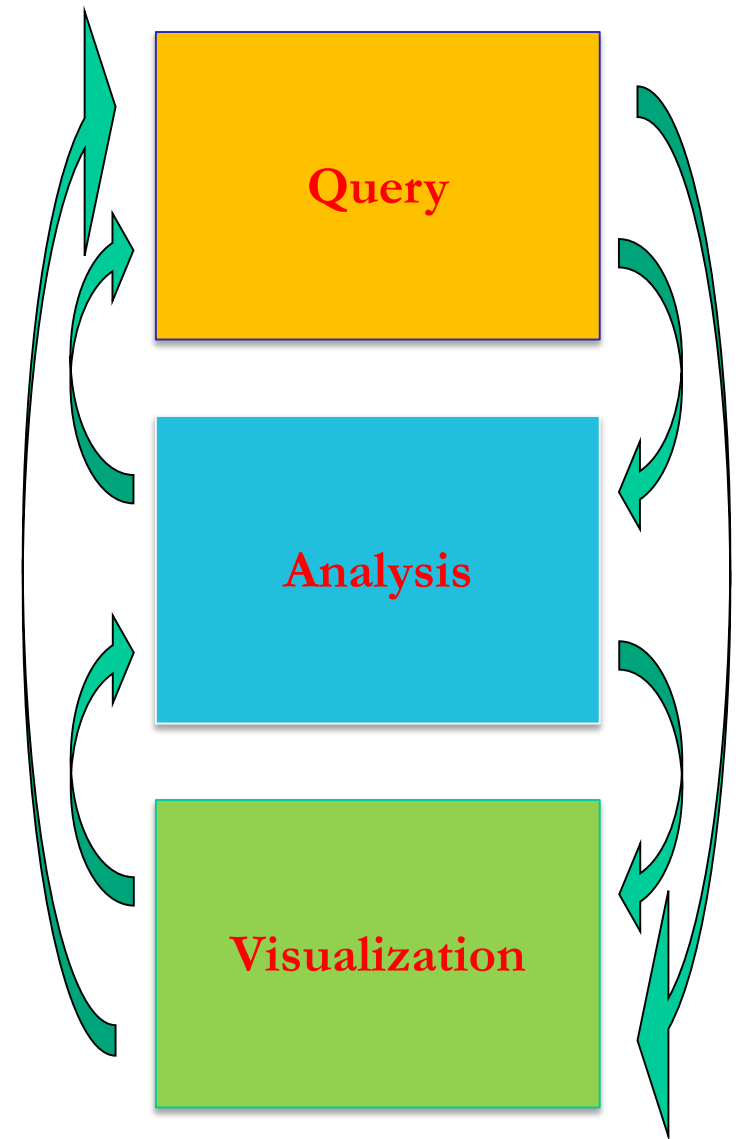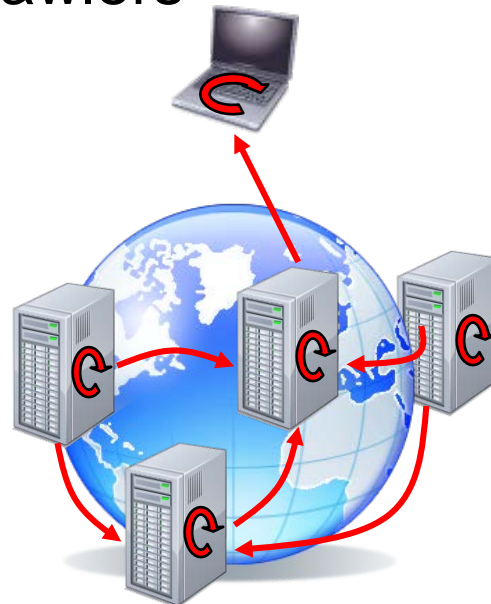
A pattern of genomic features

- **Query**
  - By example
  - Using public DBs & ontologies
  - Search remote data
  - Query / extract remote data
- **Analysis**
  - (un)supervised learning
  - Region finding
  - Motif / pattern finding
- **Visualization**
  - Clustering
  - Long range interactions

- The platform (client & servers) and language should support queries/computations involving different servers
  - Minimizing the information to be transferred among servers and between them and the client
- Each server should expose its own data for access by exploratory search & crawlers

**Overview**: http://www.bioinformatics.deib.polimi.it/genomic_computing/

**GMQL web site**: http://www.bioinformatics.deib.polimi.it/GMQLsystem/
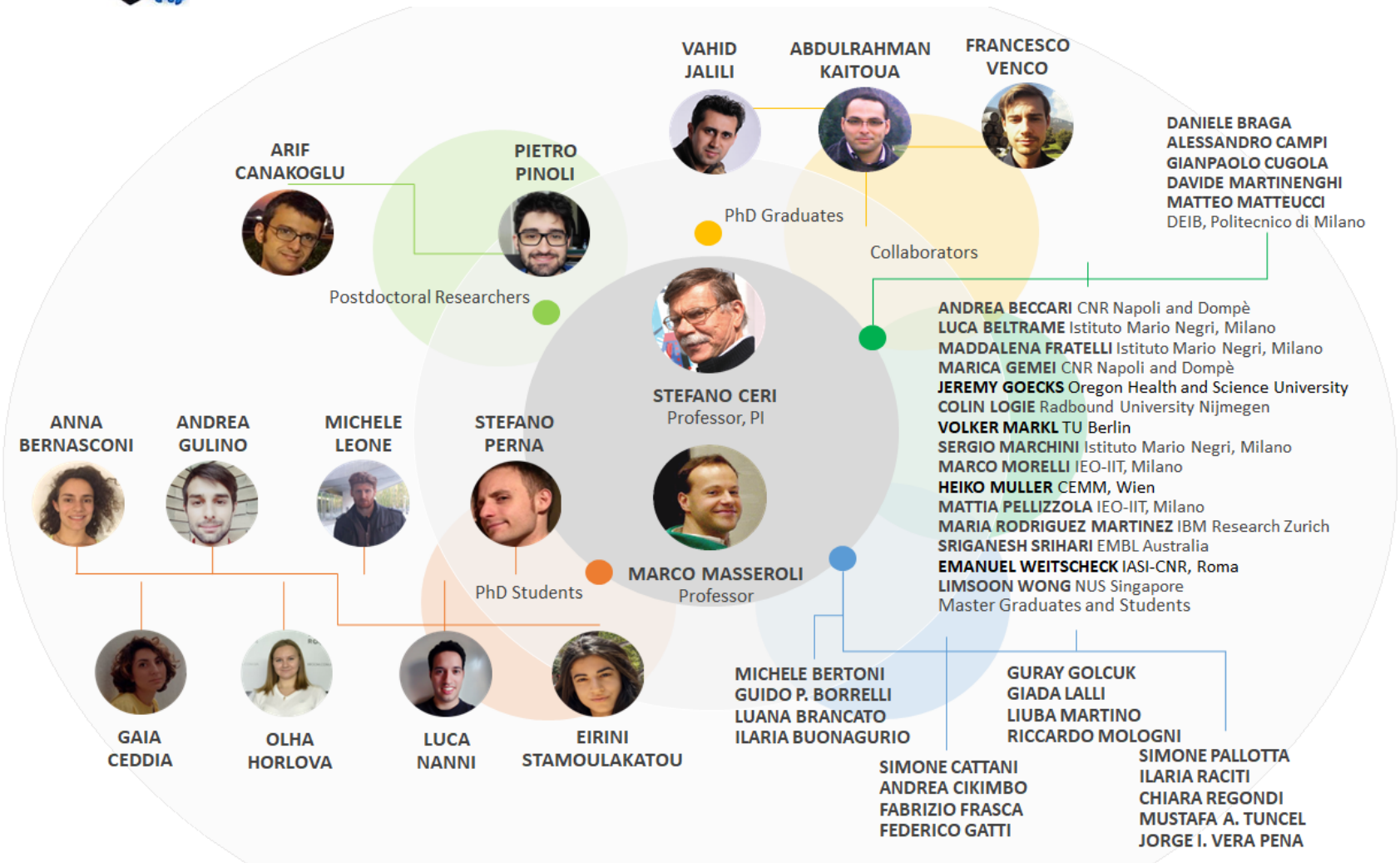Includes:

- Download open source code & documentation
  - GMQL System code & APIs - wiki
  - GMQL Web interface code - wiki & user manual
  - GMQL Package and Quick Start installation
  - GMQL Docker deploy
  - GMQL workspace in the Broad Institute FireCloud platform
  - PyGMQL Python library code & documentation
  - RGMQL R/Bioconductor package code
- Web and REST interfaces: http://www.gmql.eu/
  - User-friendly interface to creating/managing GMQL queries
  - Repository of ENCODE / Roadmap Epigenomics / TCGA datasets

**European Research Council "Data-Driven Genomic Computing"** (**GeCo**) **project**: http://www.bioinformatics.deib.polimi.it/GeCo/

**http://www.bioinformatics.deib.polimi.it/genomic_computing/**

# Thank you for your attention!

## Any question?