# Epinomics
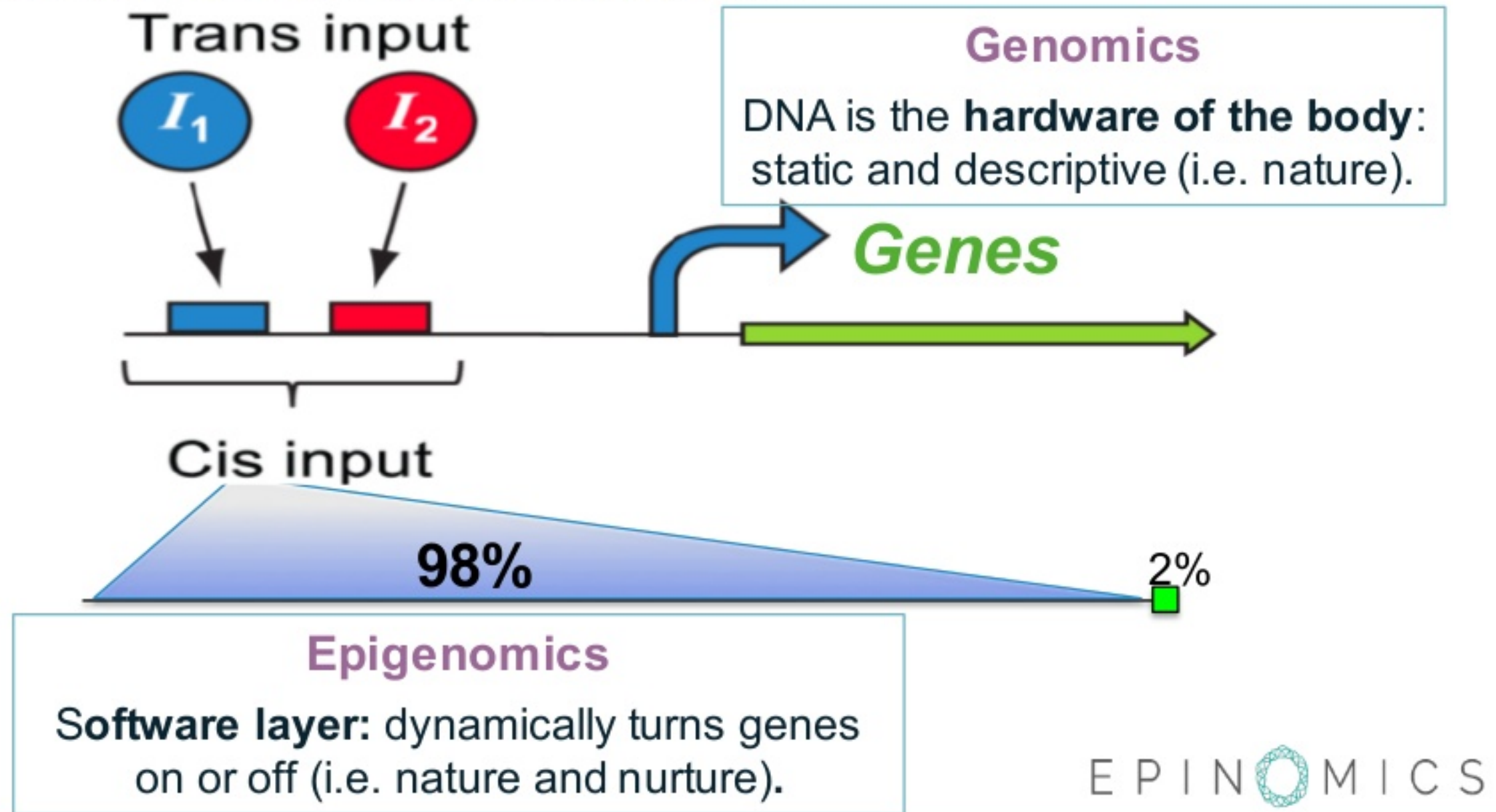
- S
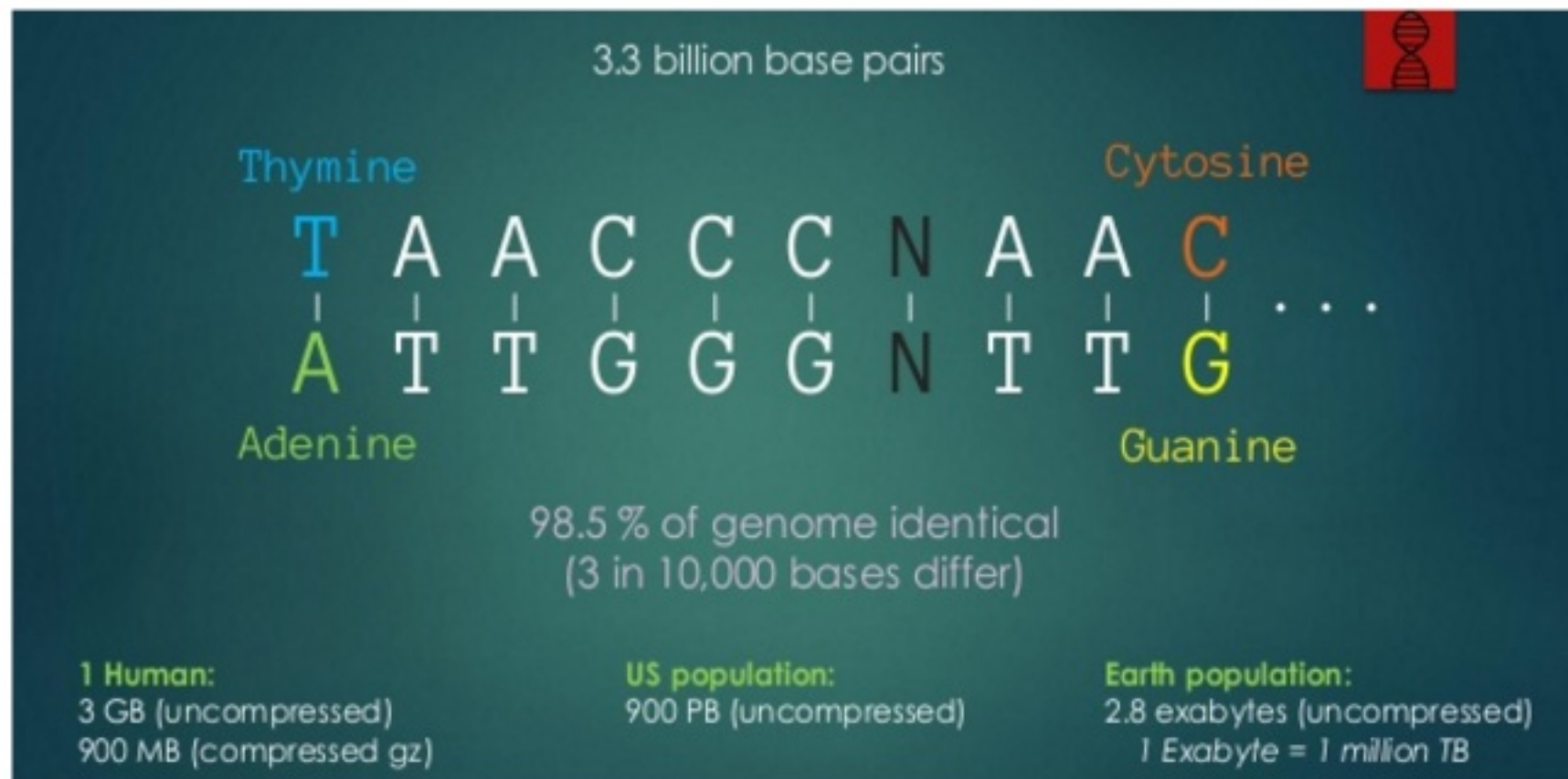


A platform that drives **personalized medicine** by leveraging big data analytics and proprietary **epigenomic** technology.

EPINOMICS

# Typical Genomic data



- Typical genomic sequencing data contains the protein letters **ATCG** .

- Most research work focuses on **variation** from standard genome sequences.
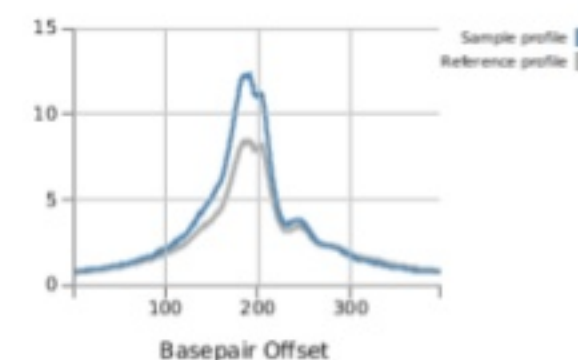
# Epigenomic Data



Peaks of Accessibility

## Fragment Data

Single fragment where DNA was accessible during the experiment.
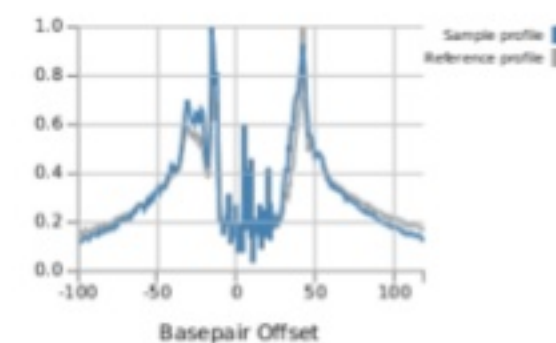
chr1  713701     714600     +
chr1  804976     805650     +
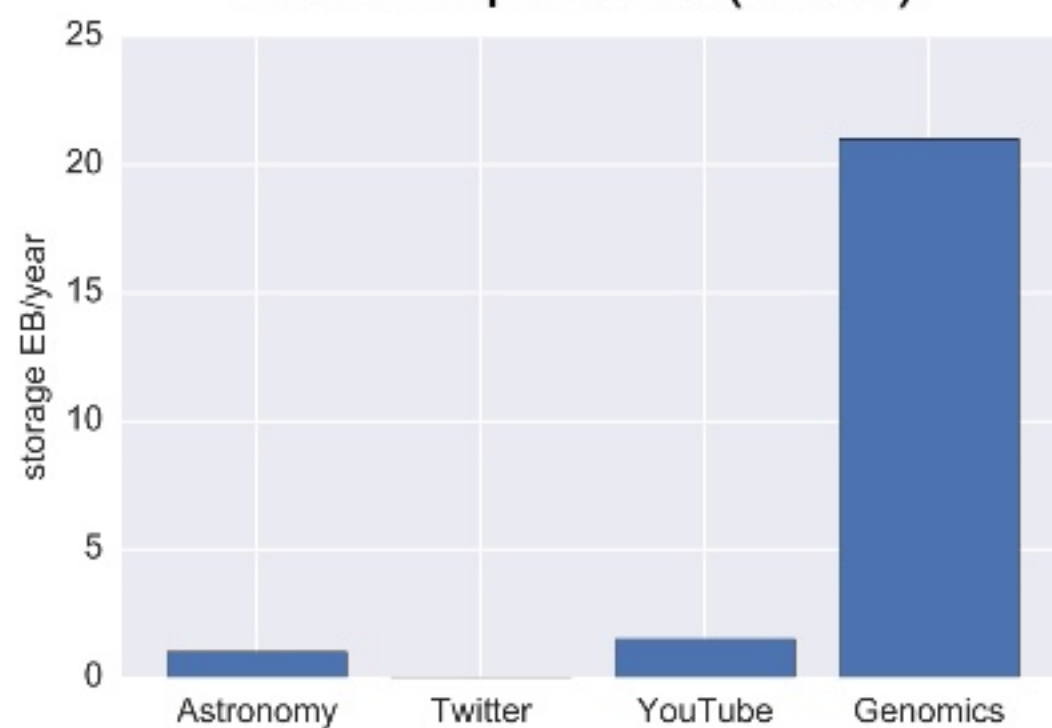


## Peaks Data

Aggregated regions of the genome where DNA was accessible during the experiment.

chr1  713701     714600     peak.1     899  +
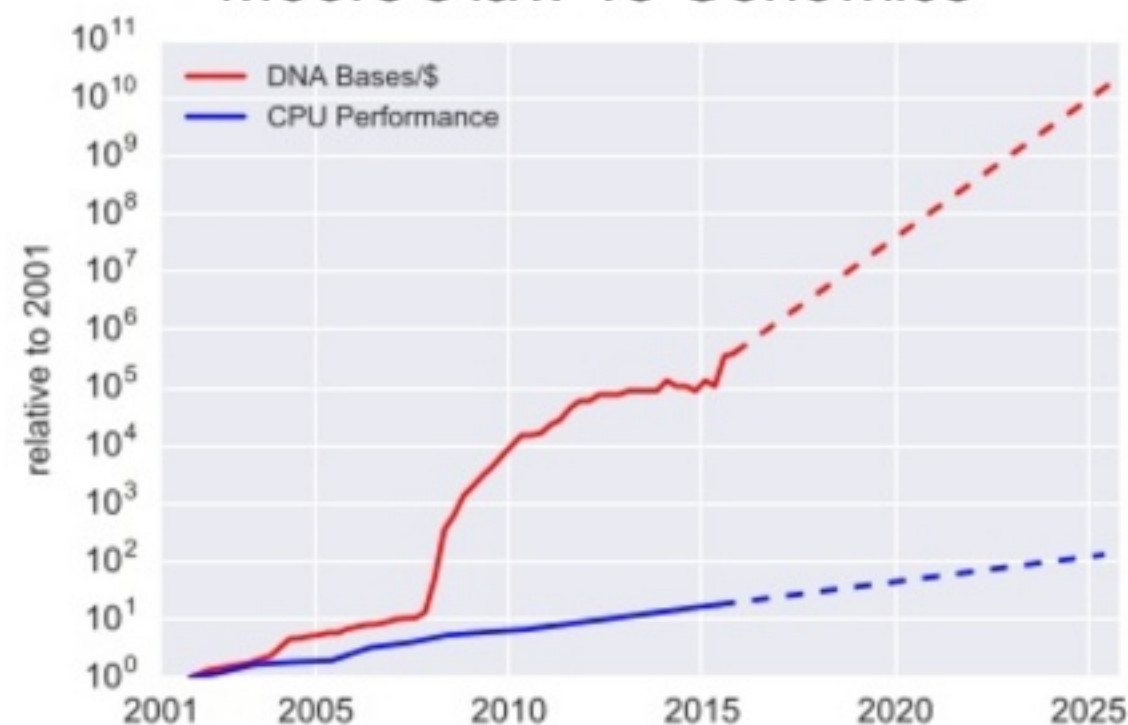chr1  804976     805650     peak.2     674  +



EPINOMICS

# Genomic Data Growth

## Data Acquisition (2015)
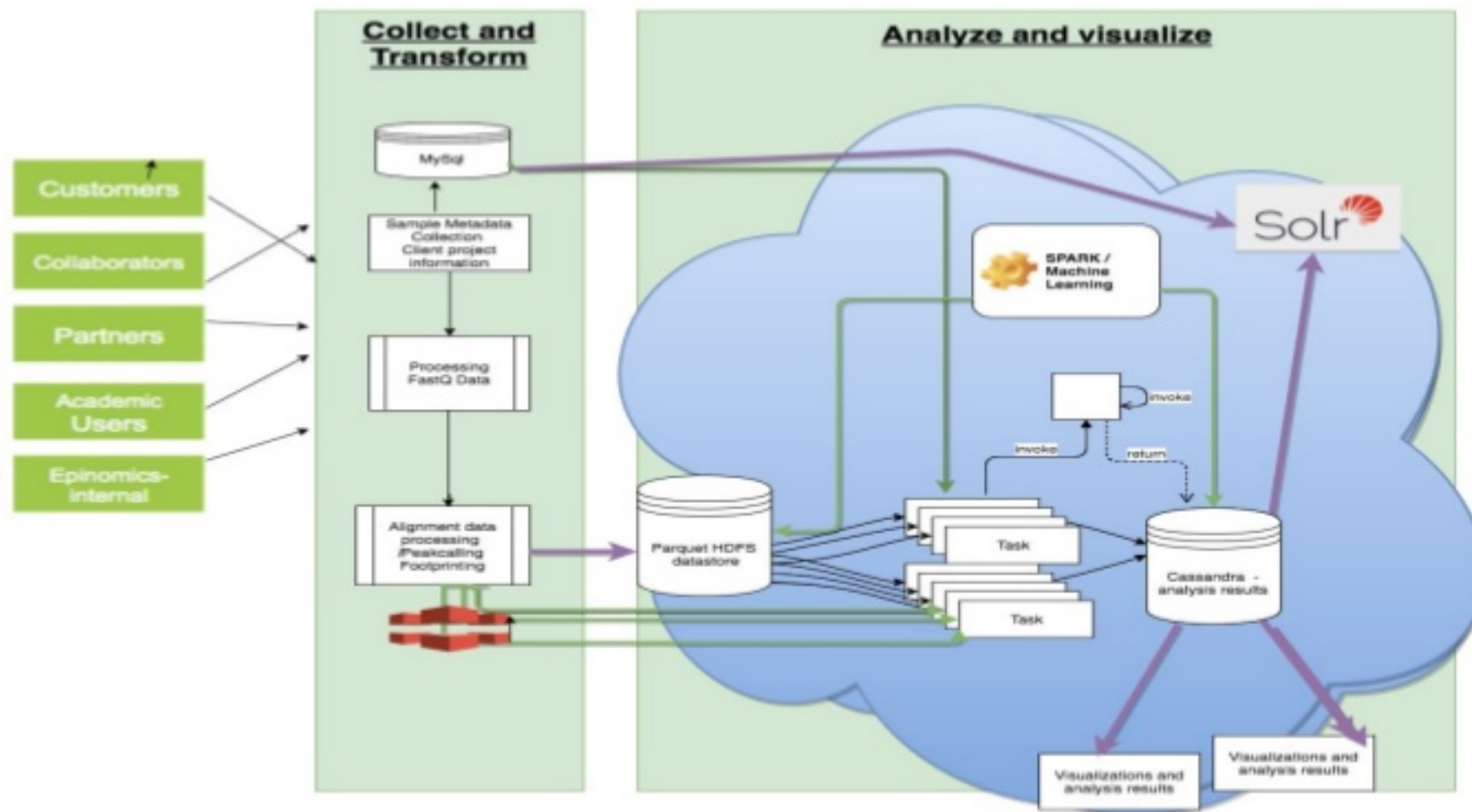


## Moore's law vs Genomics



Stephens, et al., Big Data: Astronomical or Genomical? (2015)

# Data @ Epinomics

# Goal: A Map of Human Health
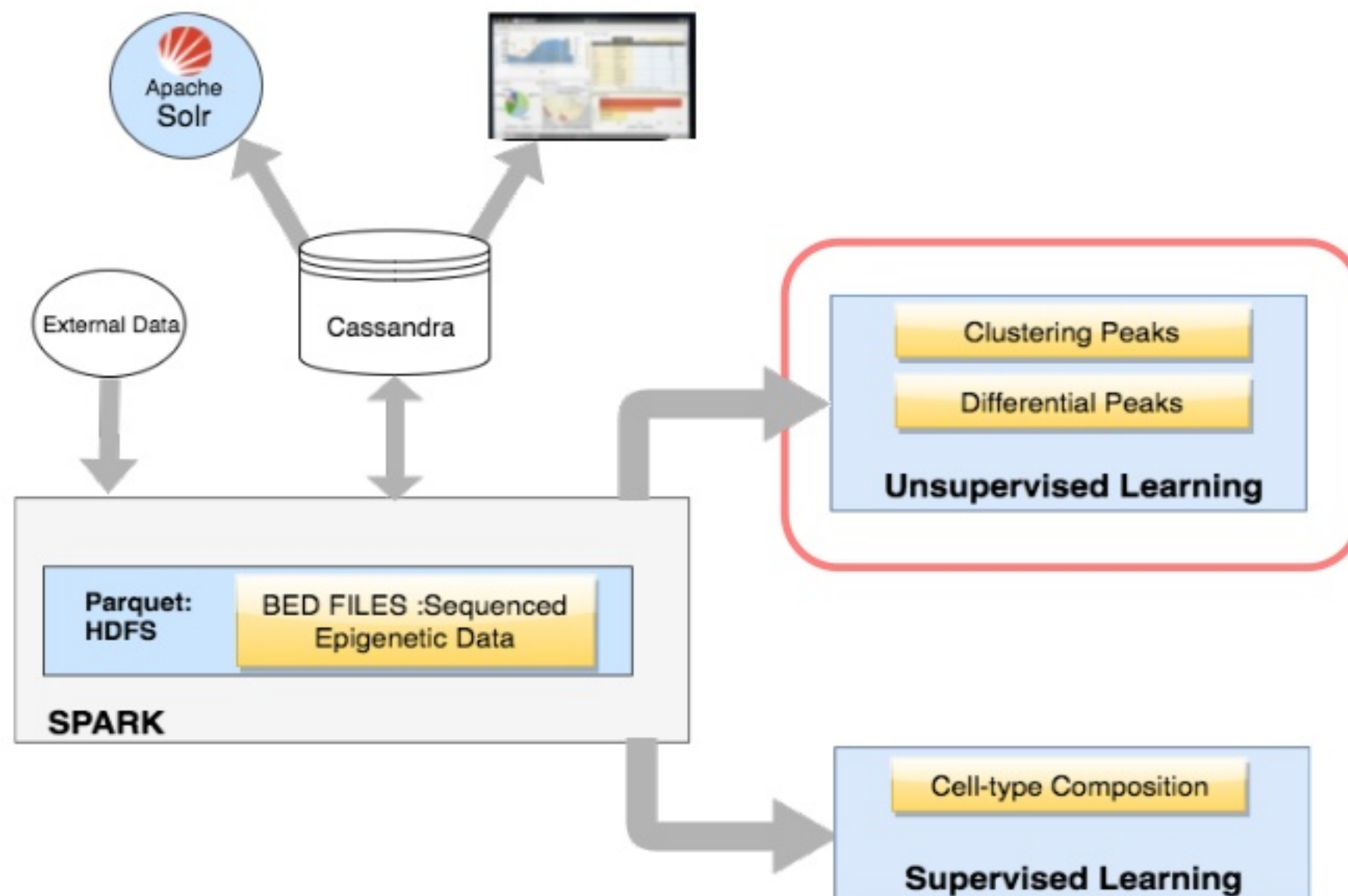
Assessing data quality

Finding patterns in the data

- **Clusters of similar data**

- **Significant differences between groups**

- **Finding unique fingerprints**

Actionable Insight

- Diagnostics, new drugs, dosage, safety

EPINOMICS

# Unsupervised Patterns of Accessibility

Process and
Consolidate Peaks → Store Peaks/Sample → Clustering
Samples based
on Peaks

Find Differences
between Sample
Groups


GraphX
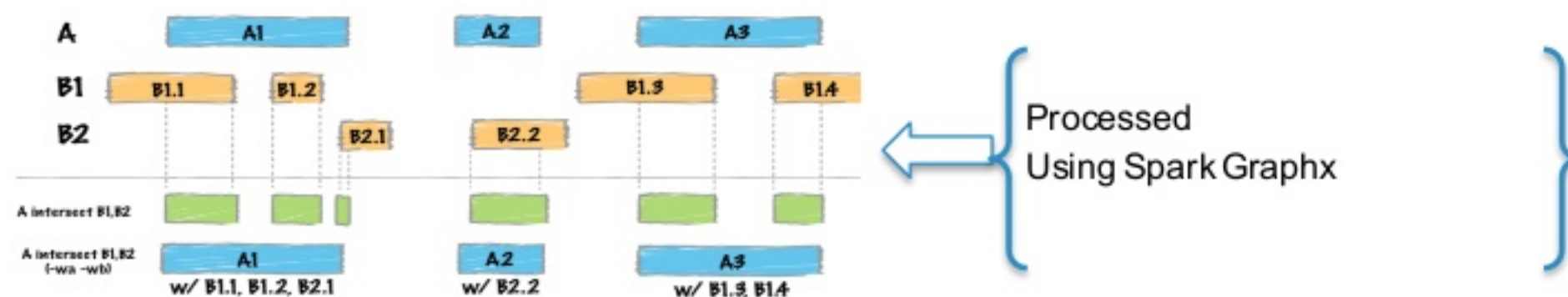

cassandra


Spark
MLlib

EPINOMICS

# Peaks Processing

Each sample will have between 150K to 200K peaks

A typical biological experiment can have between 10 to 200 samples.
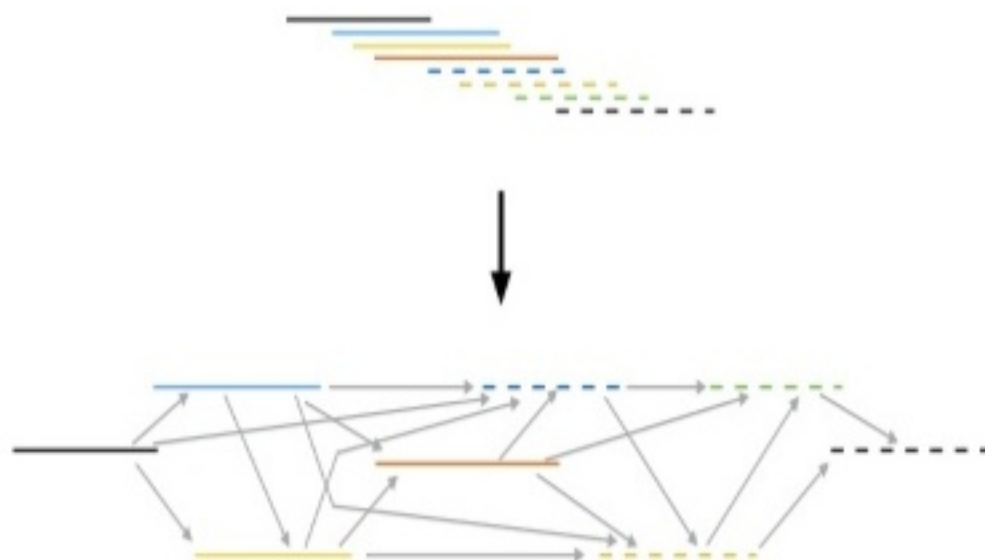
Consolidate and process overlapping peaks



Source -: http://bedtools.readthedocs.io/

A typical experiment will have between 300K to 600k overlapping peaks. (depending on dataset and sequencing depth)

# Peaks Processing

**Merges** overlapping peaks of two <u>genomic ranges</u> <u>vectors using GraphX library</u>



Nodes are peaks and edges are overlaps

1. Map all genomic ranges to Tuples where key is seq name and strand and genomic range
2. All genomic ranges are grouped by key from step above which gives us in next step all sequences with seqname and strand filtered (String, Iterable<GRanges>)
3. Include sorting Iterable<GRanges> by start position in order to implement algorithm, which will help merging ranges
4. Merge CT peaks in the way:
    a. If overlap ratio is >75 then join them into new gene:
    Overlap ratio is calculated for two genomic gRange with same seqname and strand like:
        ratio = overlap_width/width
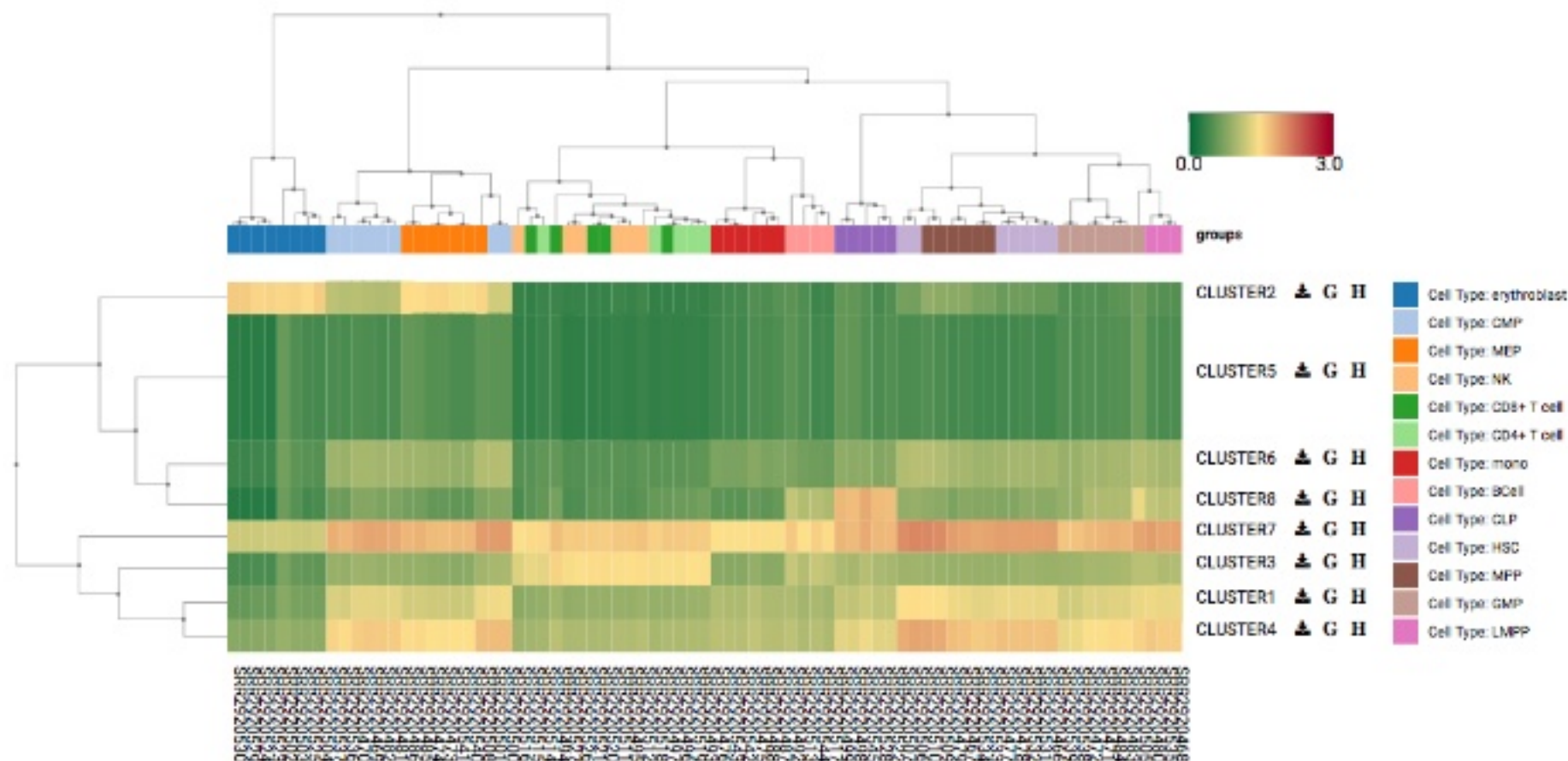    if overlap more than 75% then new gRange is created with range:

```
new IRanges(math.min(_coordinates._start, grange._coordinates._start),
math.max(_coordinates._end, grange._coordinates._end)),
```

5. Empty ranges are removed (where end=start-1)

```
val graph = Graph(jointPeaks,
edges).partitionBy(PartitionStrategy.EdgePartition1D)
    val peakRatio = sc.broadcast(cutoffRatio)
    val subgraphs = graph.connectedComponents().vertices
jointPeaks.join(subgraphs)
        .map(item => item._2.swap)
        .redyceByKey()_
        .map(item => { __.
```
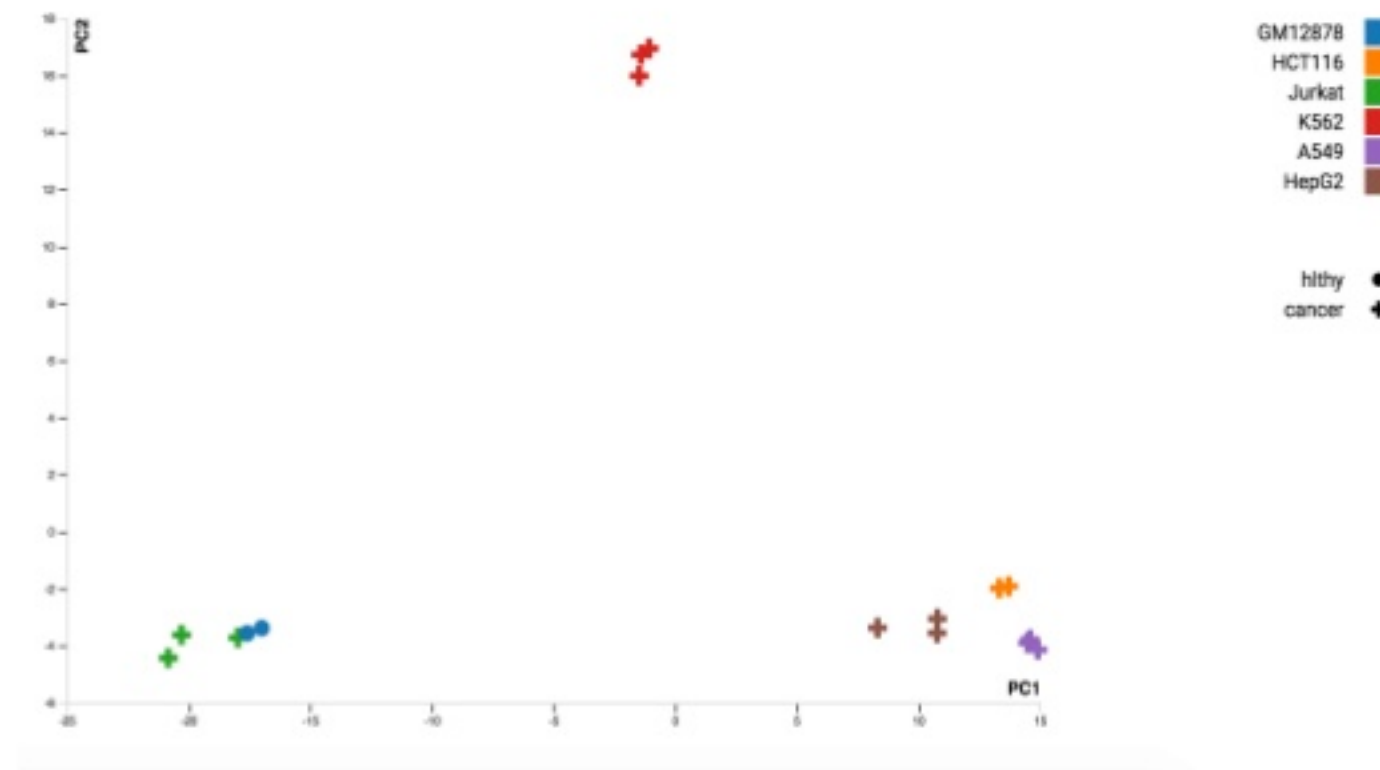
EPINOMICS
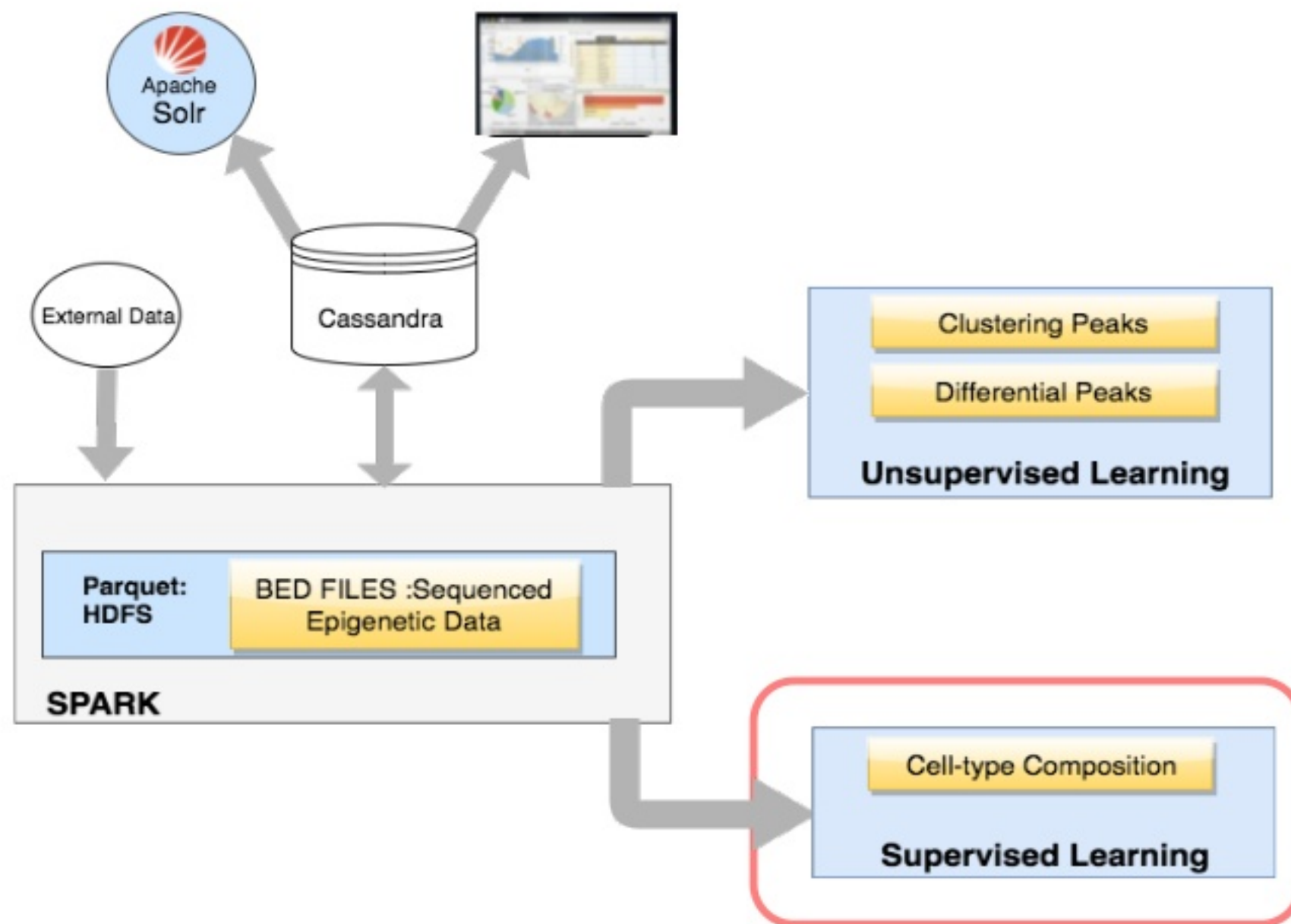
# Unsupervised Learning
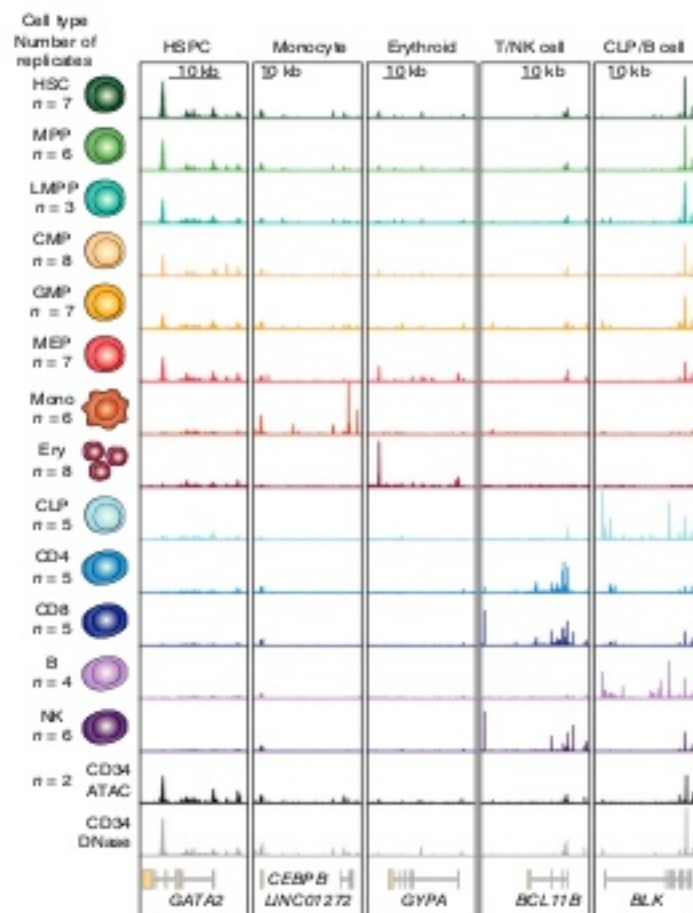


K-means and hierarchical clustering

# Unsupervised Learning



Clustering similar datasets with PCA

# Supervised Learning – Cell composition



Epigenome of each cell-type is
unique fingerprint



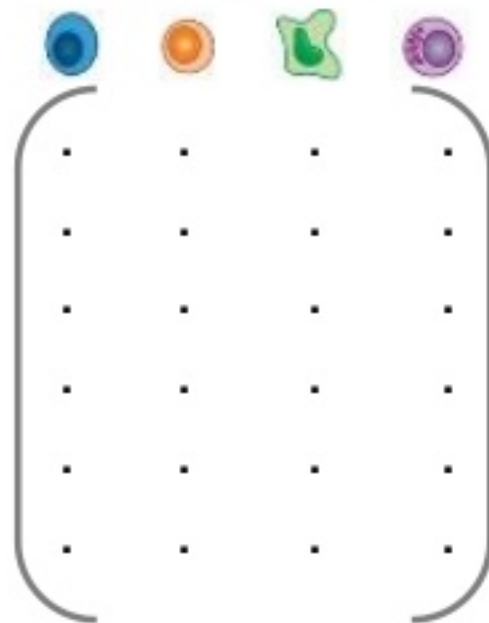**Mixed sample's signature can be
deconvolved into pure cell type signals**

Corces et al. Lineage-specific and single-cell chromatin
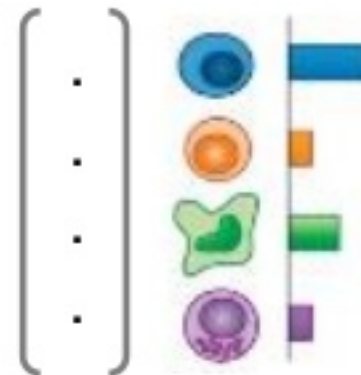accessibility charts human hematopoiesis and leukemia
evolution

# Supervised Learning – Cell composition

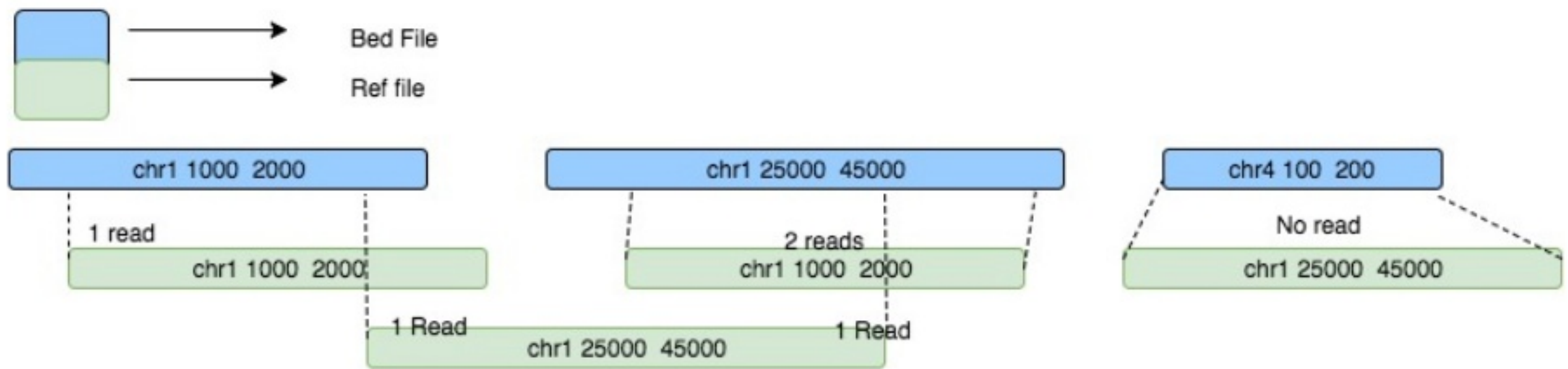**Cell type signature**
Number of reads at
specific sites
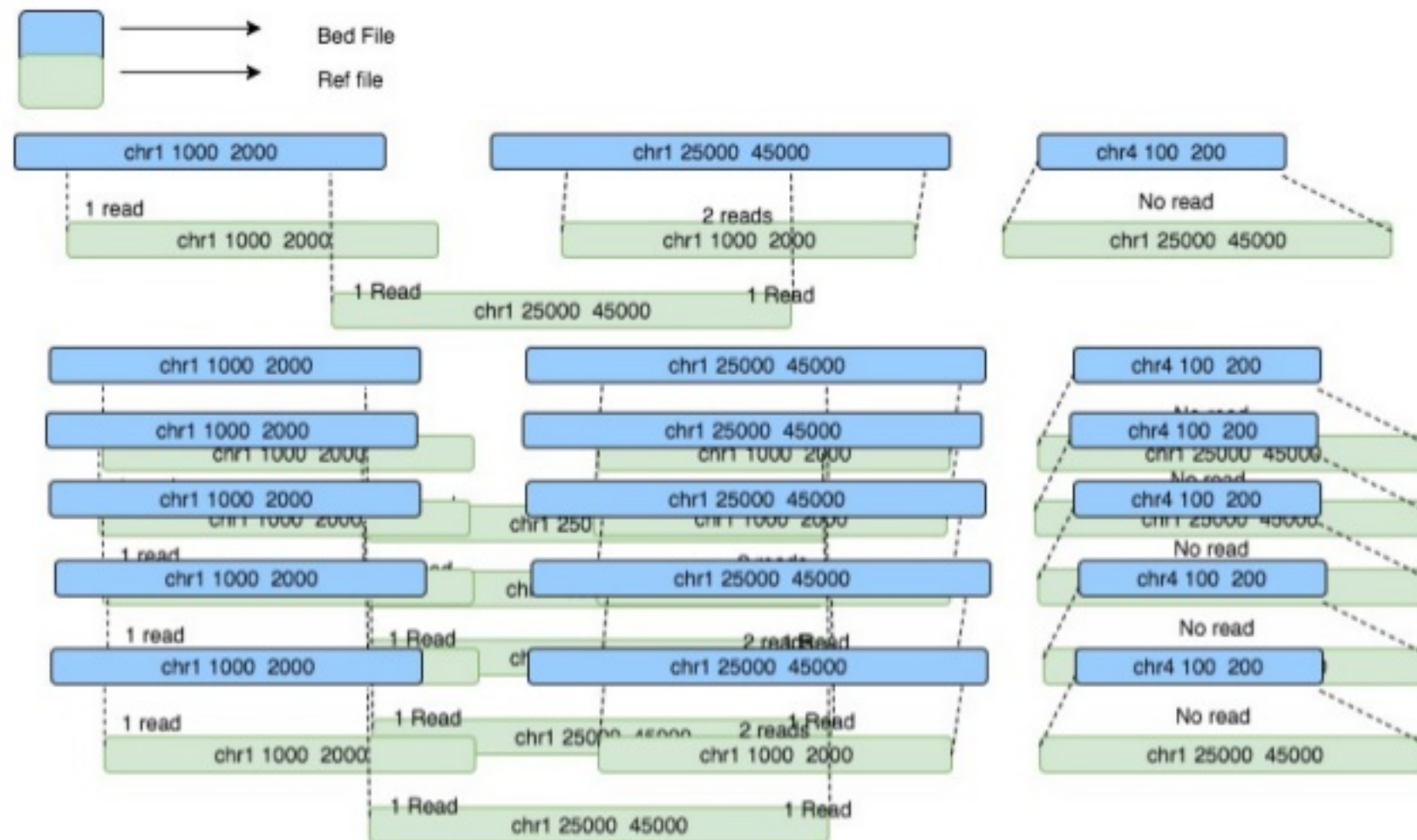
**Sample signature**
Number of reads
at specific sites

# Supervised Learning – Cell composition

Cell-type specific Regions → Count Fragments in these regions per Sample → Deconvolve to describe Cell-type composition

Reference Regions

Clinical sample with mixed cells

Composition of Cells in Sample
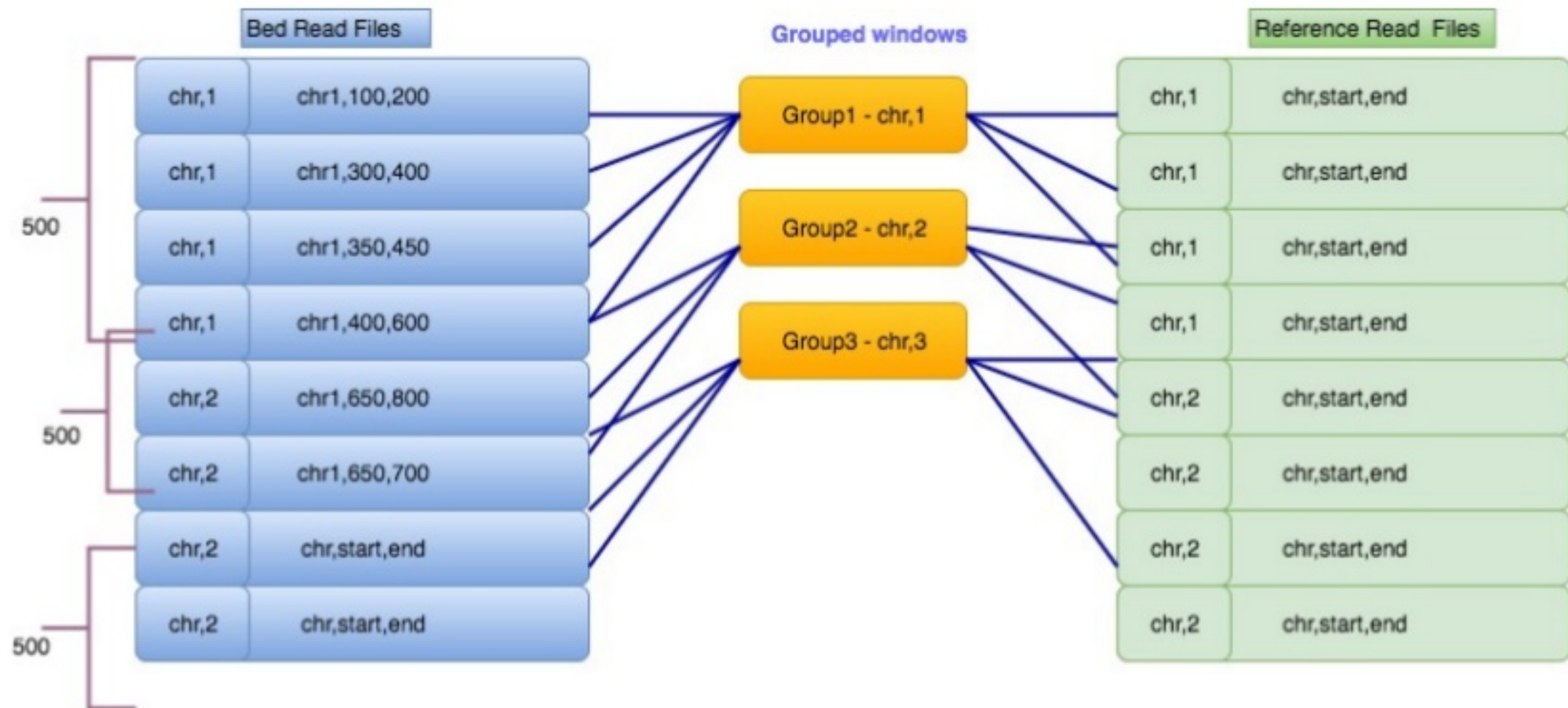
# Counting Reads within Windows
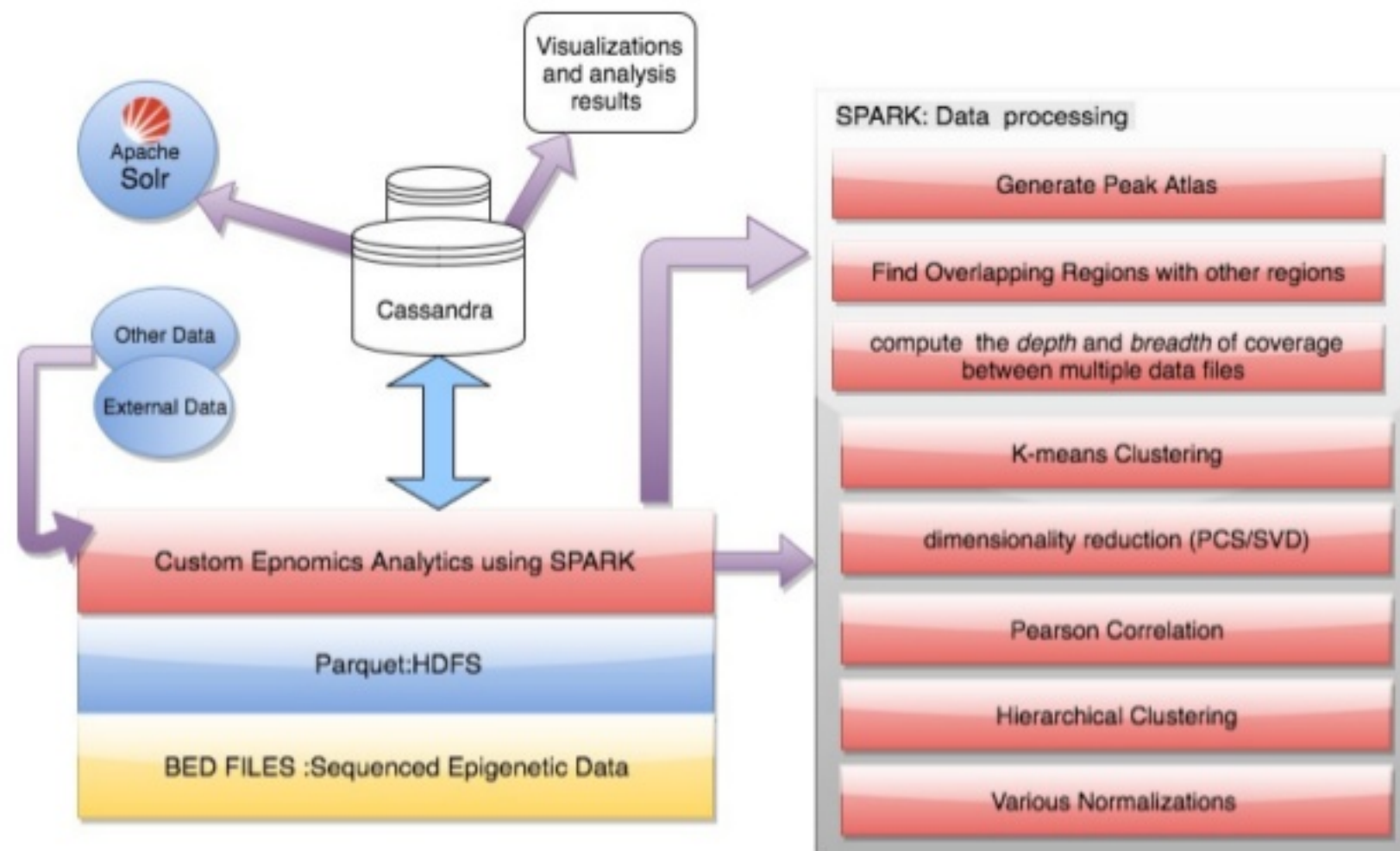
# Counting Reads within Windows

# Counting Reads – Range joins

# Building a Personalized Medicine Workflow

# Conclusion

Epinomics is building **a map of human health** through epigenomics.

ML pipelines combine Spark processing with traditional computing and algorithms.

Spark helps to process **tens of TB of genomic data** for personalized medicine applications.

# Thank You.

Anupama Joshi – anupama.joshi@gmail.com

Matei Negulescu – mnegules@uwaterloo.ca