# RETHINKING GENOME BROWSING:
## NAVIGATION BY FUNCTION NOT POSITION

CYDNEY NIELSEN

Postdoctoral Fellow

February 17, 2012

Genome **Medicine**

# The $1,000 genome, the $100,000 analysis?

Elaine R Mardis*

# The challenges of having so much data

Most challenging of all, it will be of critical importance to develop meta-analyses and statistical analysis tools that integrate across disparate data types, …
… thereby enable researchers to collectively interpret these data for all samples in a study and to form testable hypotheses from this discovery phase.

- Elaine Mardis
Anticipating the $1,000 genome

## The challenges of having so much data

Most challenging of all, it will be of critical importance to develop meta-analyses and statistical analysis tools that integrate across disparate data types, …
… thereby enable researchers to collectively interpret these data for all samples in a study and to form testable hypotheses from this discovery phase.

- Elaine Mardis
  Anticipating the $1,000 genome

## Data visualization will be a key player in this domain

# Why visualization?

Data landscape is unknown

# Why visualization?

Data landscape is unknown

- In the discovery phase, not yet clear where the interesting features lie or what they look like

- Features are not sufficiently well defined to be extracted in a purely automated fashion

- Visualization is a powerful approach in such cases:
  Exploit our visual system and knowledge to identify biologically interesting data patterns and subsequently generalize

# Why visualization?

Improves data accessibility to biological community

# Why visualization?

Improves data accessibility to biological community

- More and more data are being produced from large consortia (ENCODE, Epigenome Roadmap Project, etc.)

- Download portals are valuable, but of primary use to computational experts

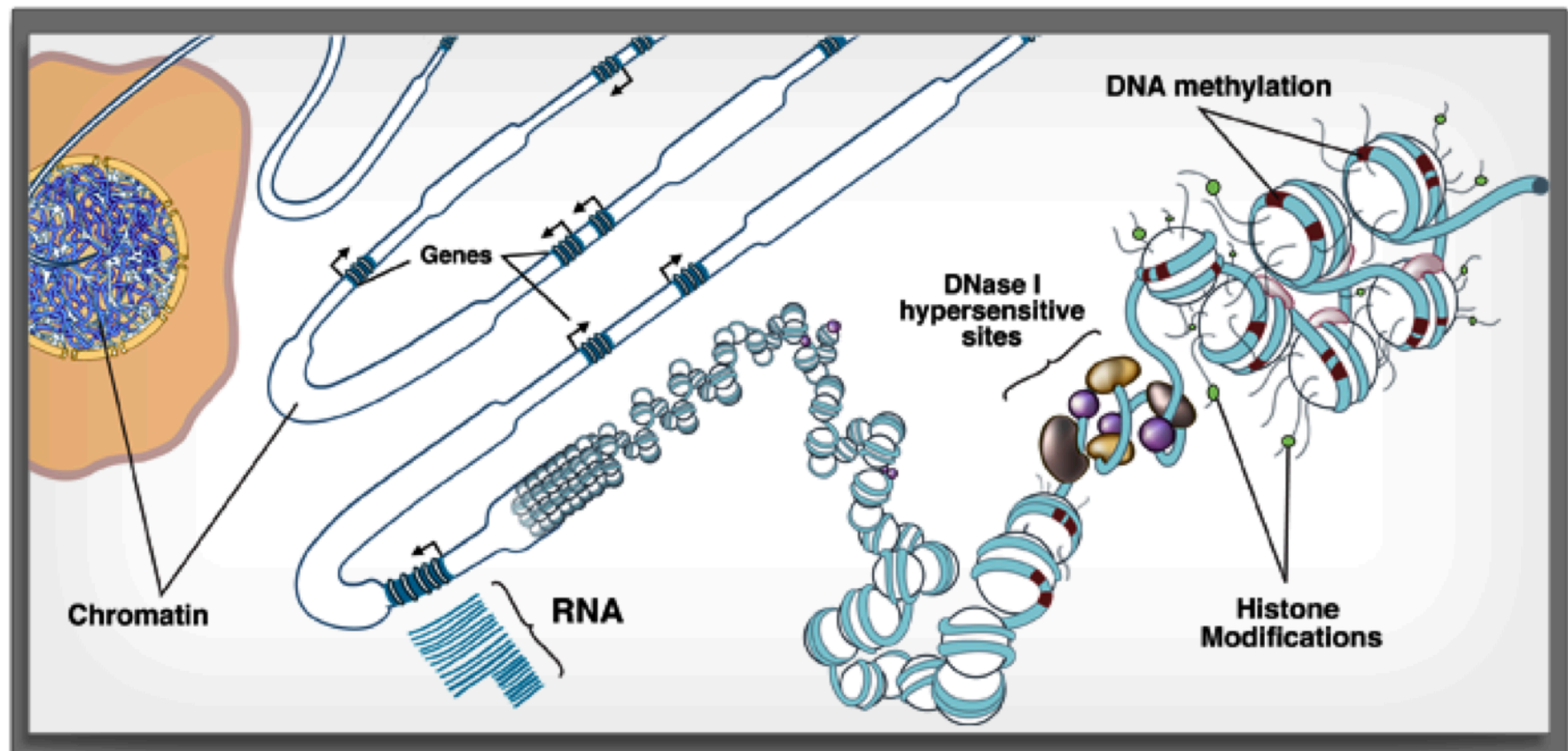- If we are to make the most of these large datasets, need to bridge the computational gap between primary data and biological community – visualization can play a key role

# Outline

- Exploring genome-wide datasets
  - Computational and visual methods to investigate ChIP-Seq data
- Spark
  - A navigational paradigm for interactive genome exploration
- Demo
- Future work

Many cell types with same DNA sequence but different morphologies
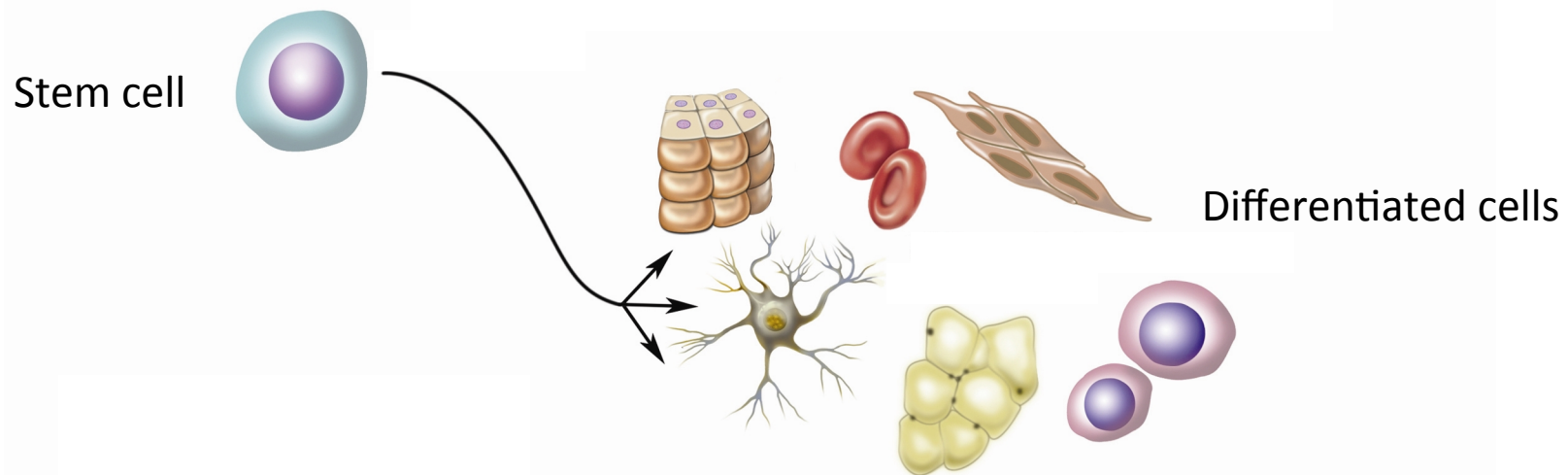
Stem cell

Differentiated cells

www.nationalacademies.org/stemcells

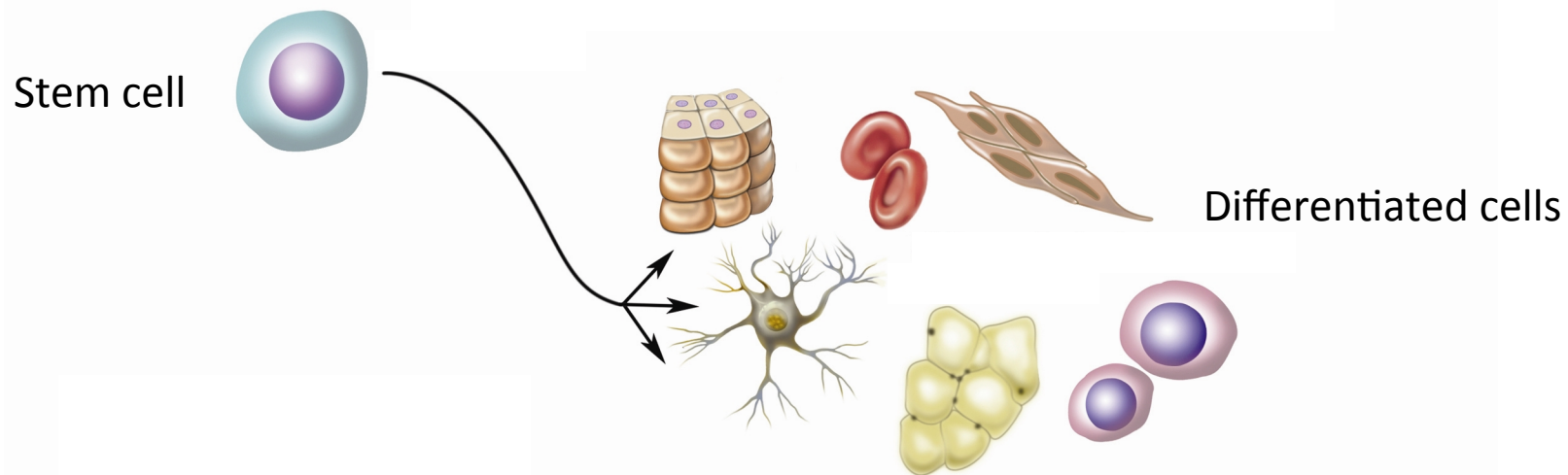Epigenetics - the study of changes in the regulation of gene activity that are not dependent on gene sequence

epi - (above) – genetics

Epigenetics - the study of changes in the regulation of gene activity that are not dependent on gene sequence

epi - (above) – genetics

Epigenetics - the study of changes in the regulation of gene activity that are not dependent on gene sequence

epi - (above) – genetics
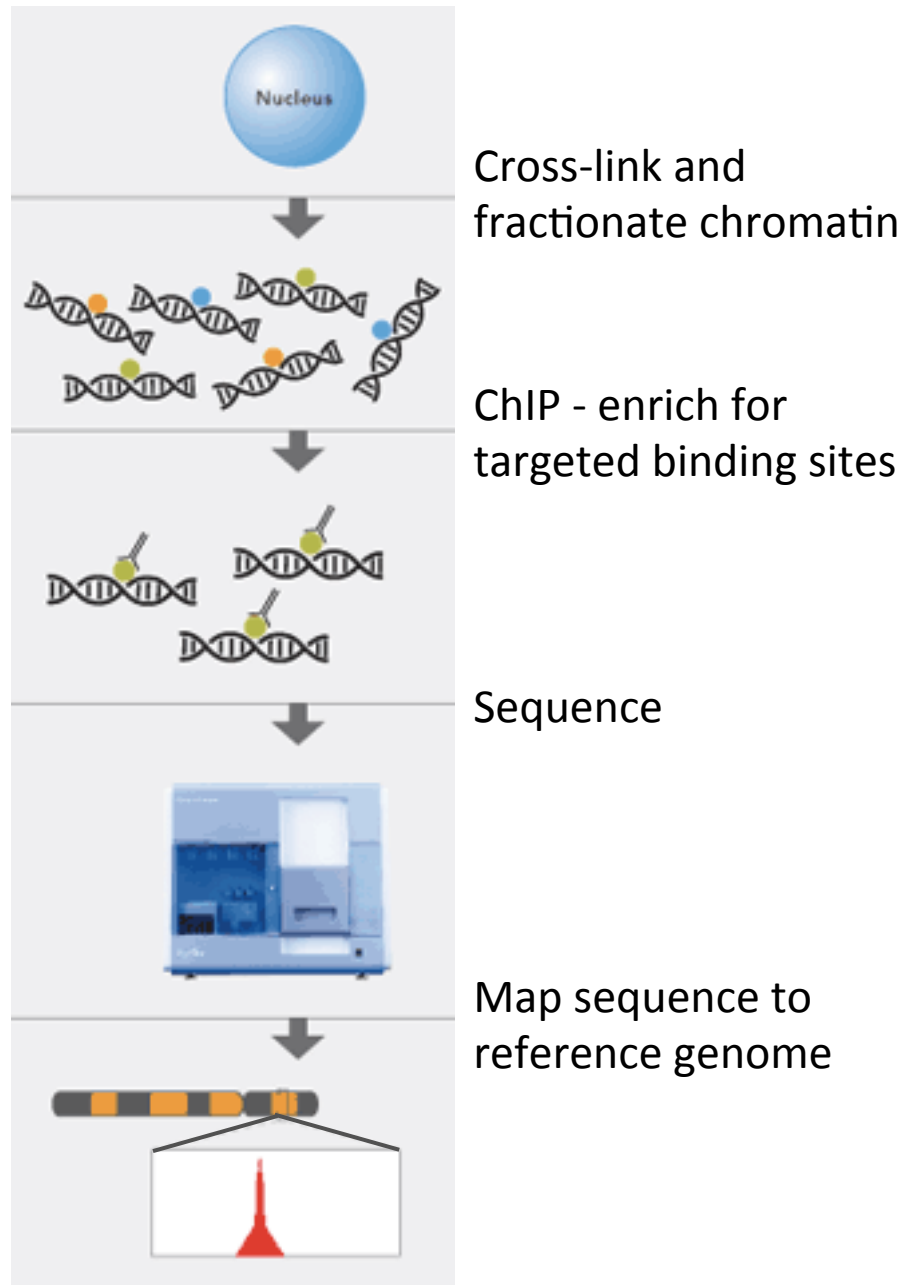
http://bricker.tcnj.edu/Amb/amble9.html

A Roadmap Project goal is to create reference epigenomic maps of many different human cell types (i.e. a map of histone modifications across the genome)

Stem cell

Differentiated cells

A Roadmap Project goal is to create reference epigenomic maps of many different human cell types (i.e. a map of histone modifications across the genome)

Stem cell

Differentiated cells

Developmentally important genes are "poised":
H3K4me3 (active) **AND** H3K27me3 (inactive)

Resolve to either
H3K4me3 **OR** H3K27me3

A Roadmap Project goal is to create reference epigenomic maps of many different human cell types (i.e. a map of histone modifications across the genome)

reversible

Stem cell

Differentiated cells

Developmentally important genes are "poised":
H3K4me3 (active) AND H3K27me3 (inactive)

Resolve to either
H3K4me3 OR H3K27me3

# ChIP-Seq | Chromatin Immunoprecipitation and Sequencing



Cross-link and fractionate chromatin

ChIP - enrich for targeted binding sites

Sequence

Map sequence to reference genome

Peak = putative binding site

illumına®

# Current computational techniques

**Heintzman *et al*. Nature Genetics, 2007**

Motivation – "…the distinguishing chromatin features of promoters and enhancers have yet to be determined, hindering our understanding of a predictive histone code for different classes of regulatory elements."

# Current computational techniques

**Heintzman *et al*. Nature Genetics, 2007**



Clustering data from well annotated regions

# Current computational techniques

**Hon** *et al.*
**PLoS Comput. Biol., 2008**

**ChromaSig**: a probabilistic method that enables discovery of chromatin signatures de novo (no dependence on annotation)

# Current computational techniques

**Ernst and Kellis
Nature Biotechnol.
2010**

Hidden Markov
Model to reveal
'chromatin states'
in human T cells

State display:
abstracted away
from all detail of
primary data

# Current computational techniques

- Require significant computational skill to use (only ChromaSig provides an implementation)

- All produce static overview images and do not support interactive data exploration of individual genes within a signature class

- No integration with downstream processing (e.g. gene ontology enrichments, etc.)

# Current visualization techniques

Kent *et al*. Genome Research, 2002

"The UCSC browser had humble origins. The code originated with a small script in the C programming language, which displayed a splicing diagram for a gene prediction from the nematode *C. elegans* (Kent and Zahler 2000). This web-based splicing display later acquired tracks for mRNA alignments and for homology with the related nematode *Caenorhabditis briggsae*. This was published as the tracks display at http://www.cse.ucsc.edu/~kent/intronerator"
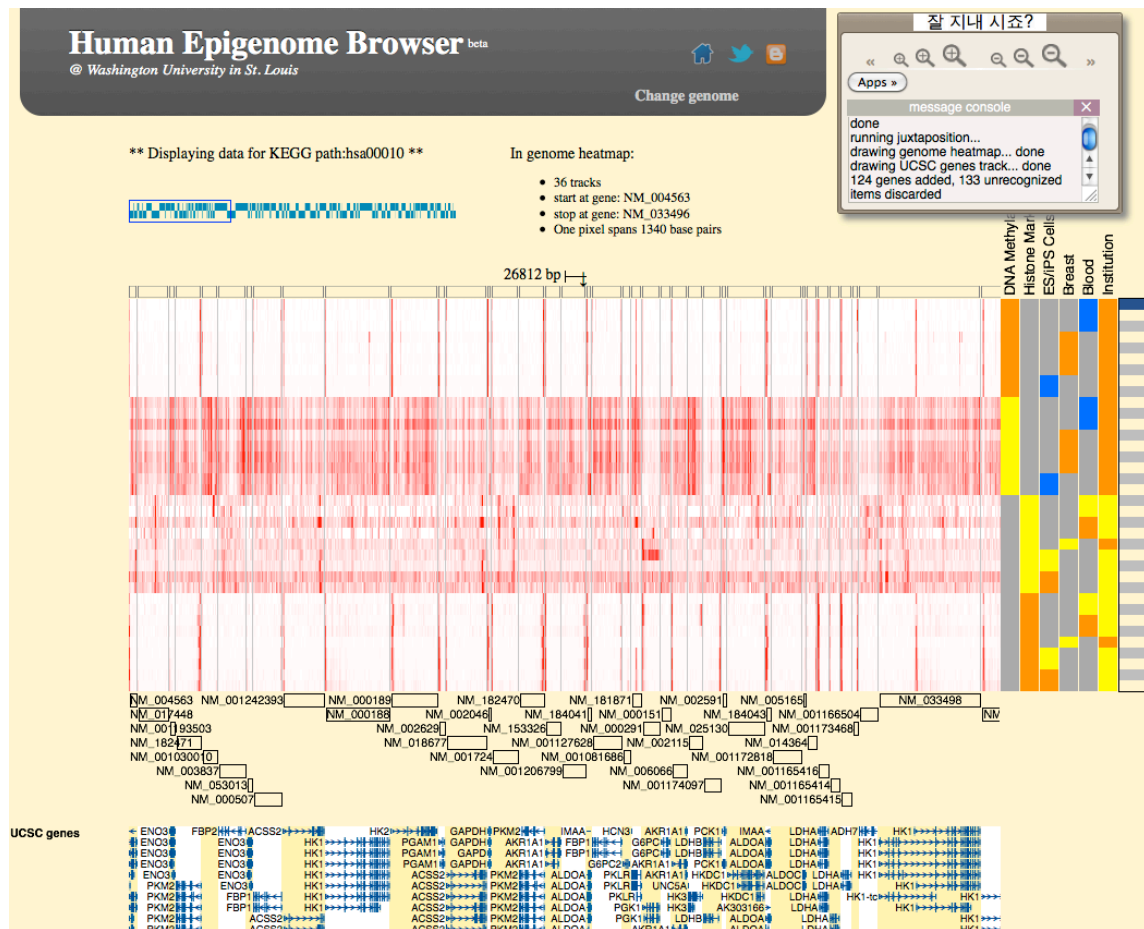
**Kent *et al*. Genome Research, 2002**

# Current visualization techniques

## Genome Browsers

- Intended to provide a local view of a genomic region

- In the absence of other tools, many biologists attempt to informally use them to gather a global overview of data patterns

- In these cases, there is a mismatch between the level of data abstraction at which a biologist reasons about the problem and the level provided by the browser

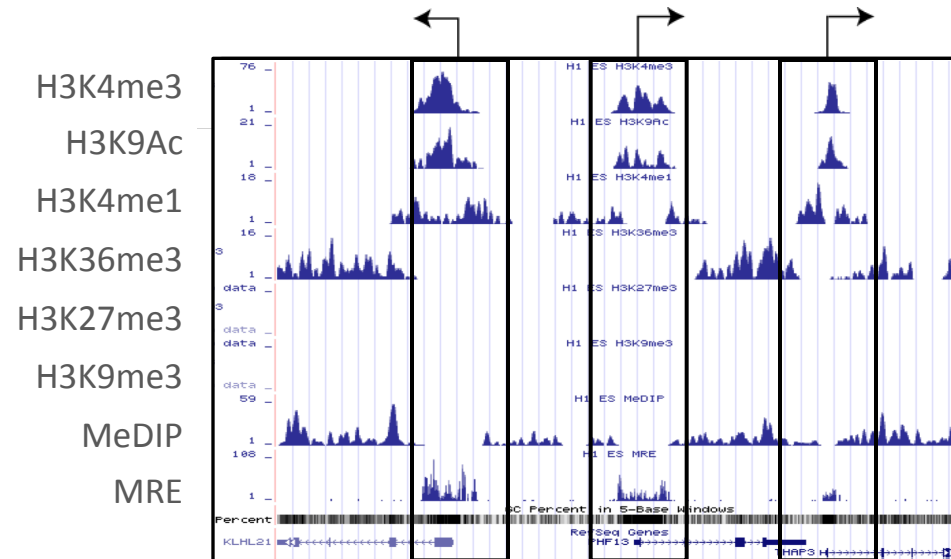- Functional similarity should drive navigation not genomic position

# Current visualization techniques



**Zhou *et al*. Nature Methods, 2011**

Can filter genomic x-axis to just display genes from a pathway of interest (by KEGG ID)
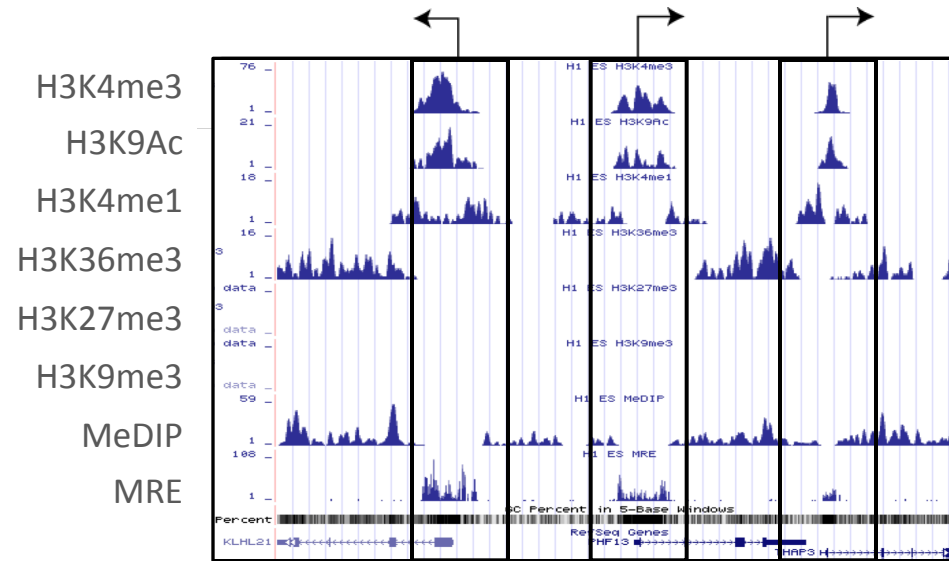
# Spark – A Discovery Tool

# Spark

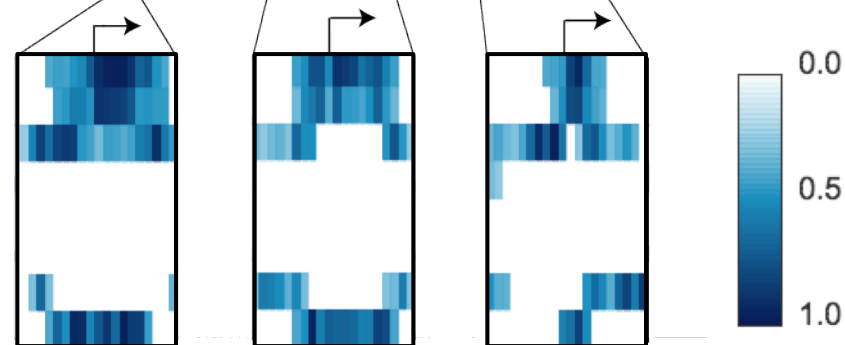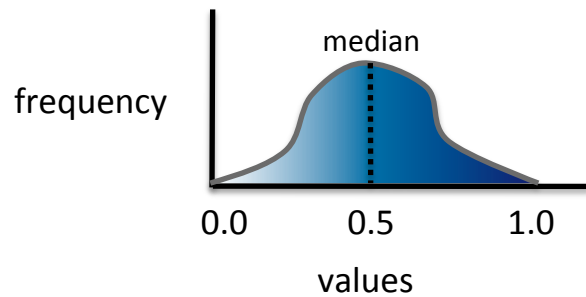1. Focus on regions of interest (e.g. transcriptional start sites (TSSs) +/- 3000 nt)

# Spark

## 1. Focus on regions of interest (e.g. transcriptional start sites (TSSs) +/- 3000 nt)
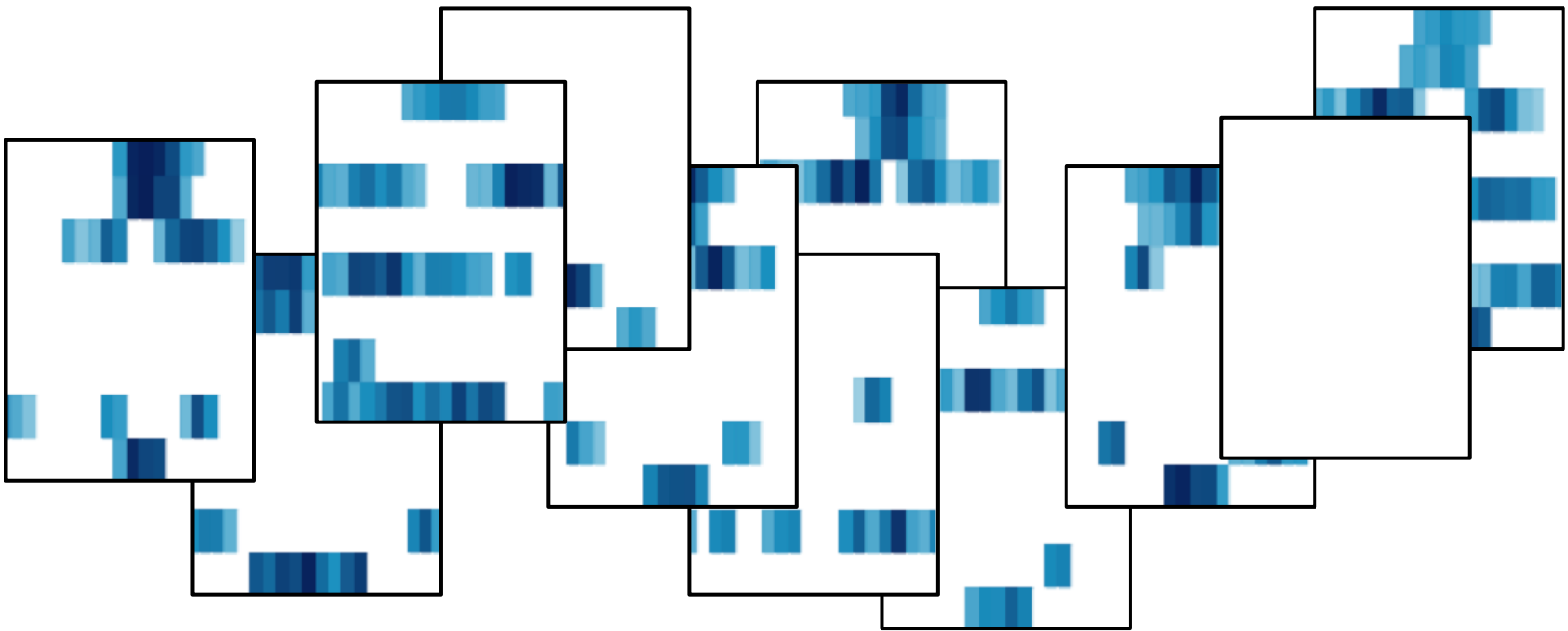


## 2. Extract data matrices

Normalization for bin $i$, sample $h$:



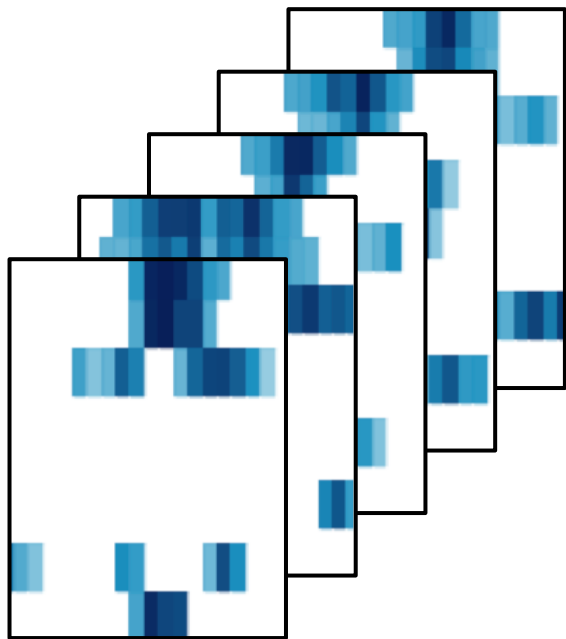$$x'_{h,i} = \frac{1}{1 + e^{-(x_{h,i} - median(x_h))/std(x_h)}}$$

# Spark

3. Cluster matrices (*k*-means)
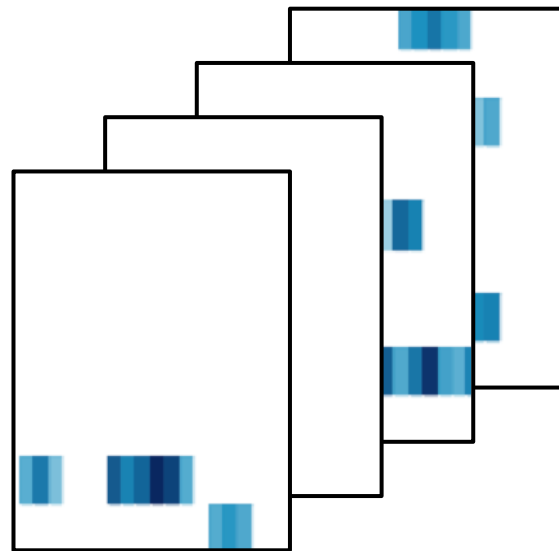
# Spark

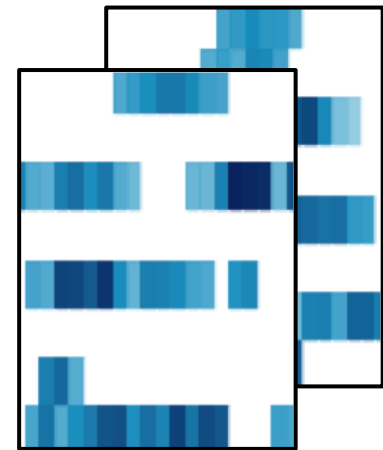3. Cluster matrices (*k*-means)   *k* = 3



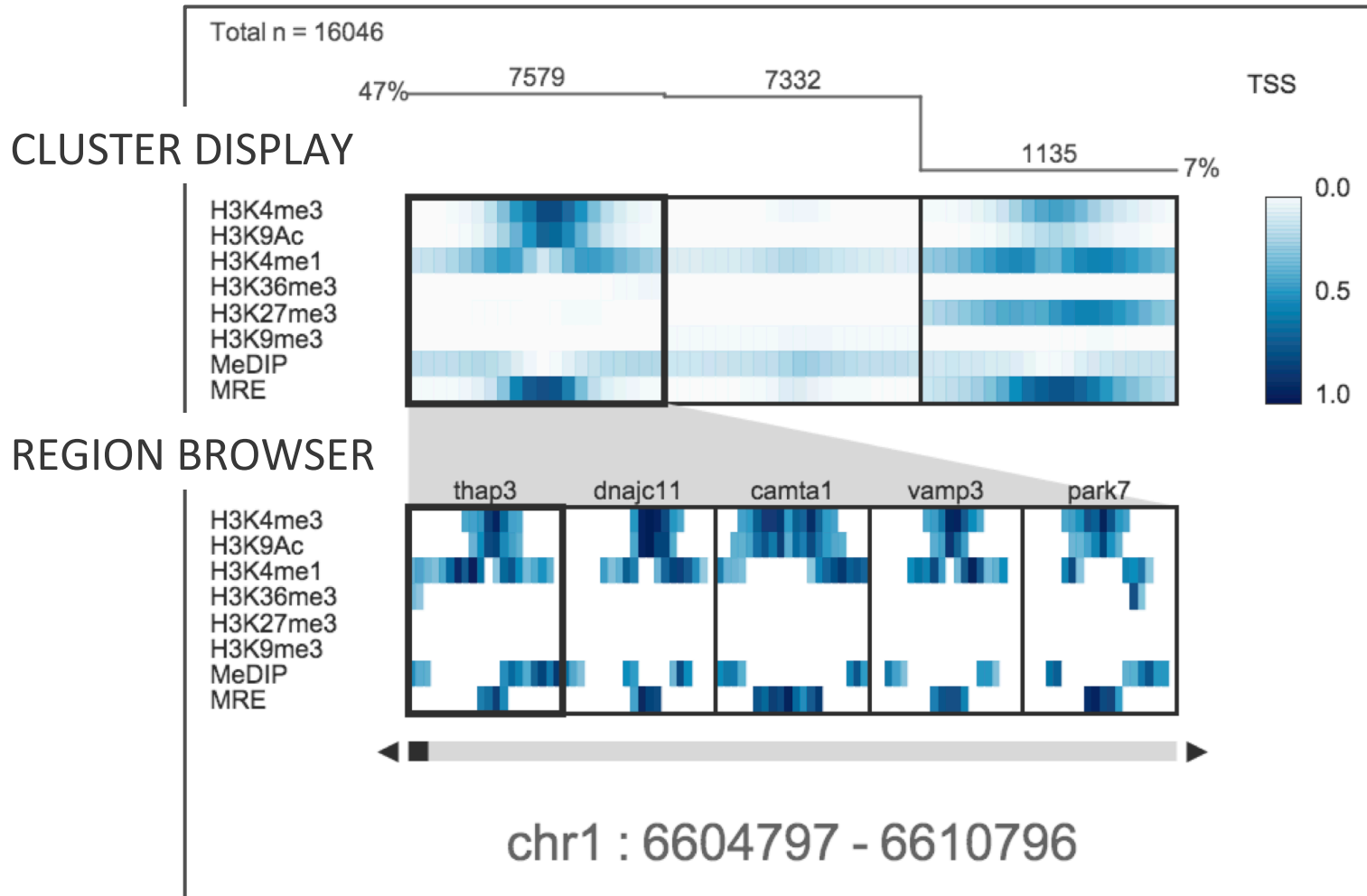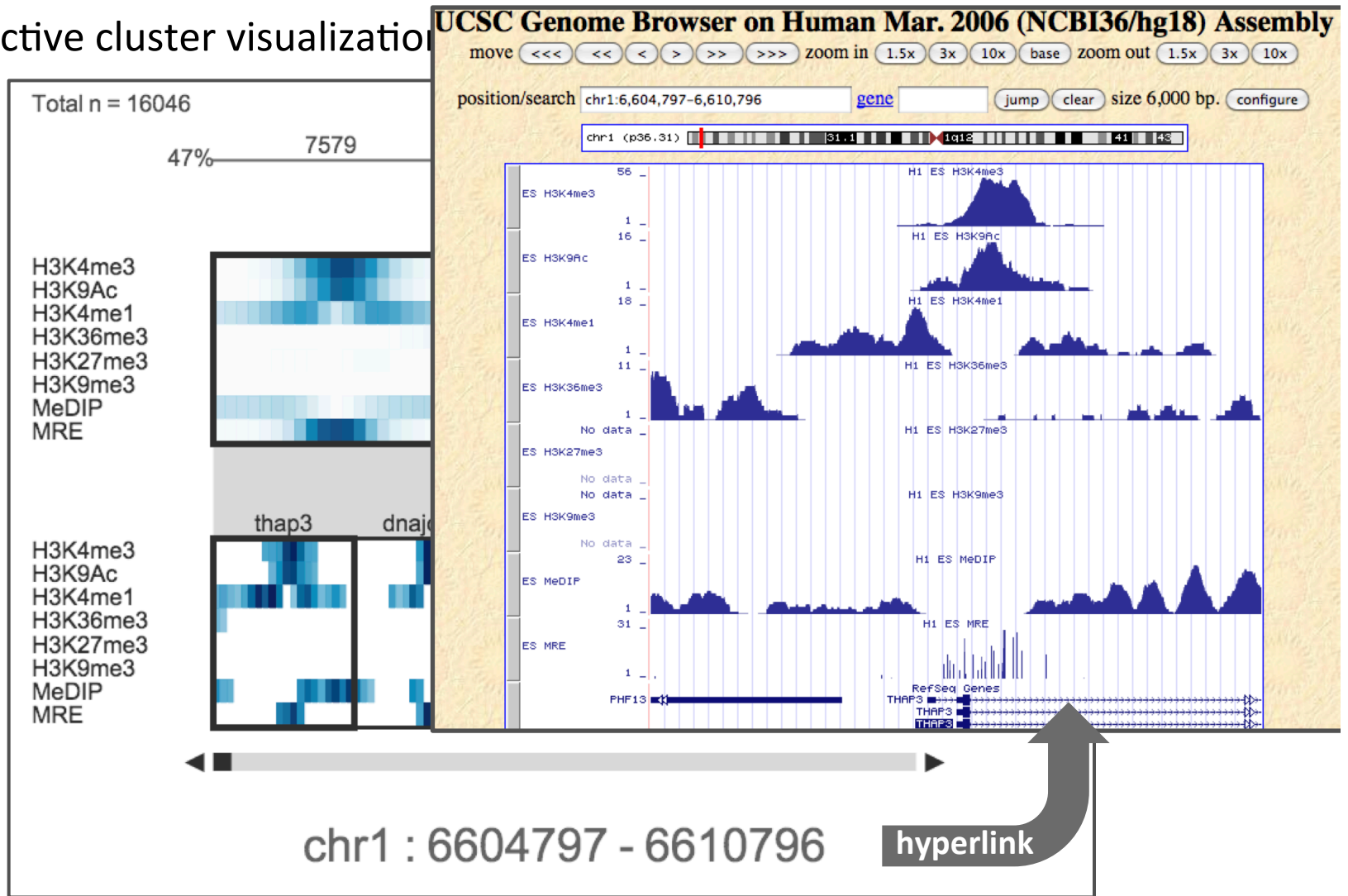Cluster 1                    Cluster 2                    Cluster 3

# Spark

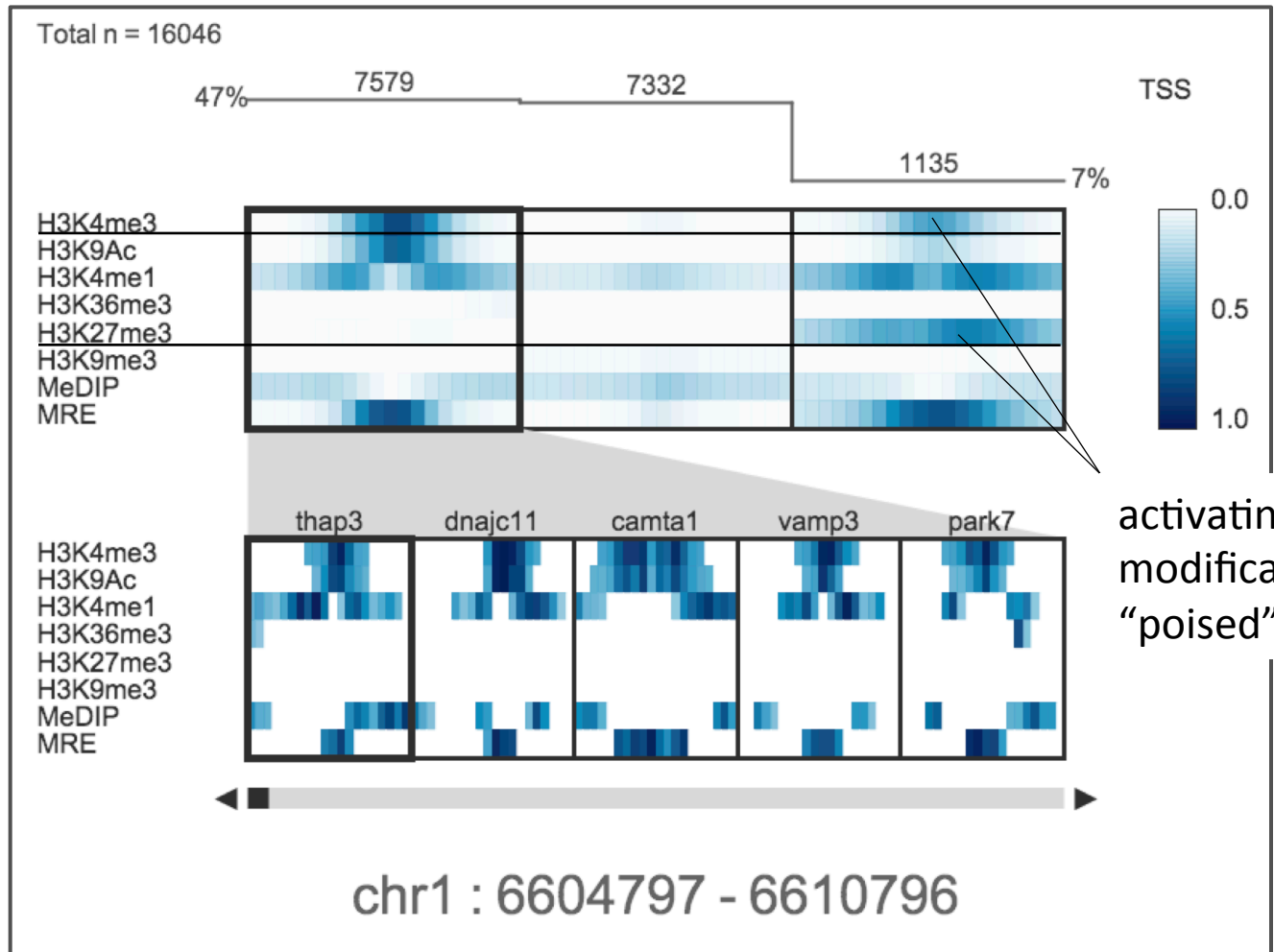4. Interactive cluster visualization - data from human H1 embryonic stem cells



screenshot

# Spark

4. Interactive cluster visualization



Total n = 16046

47%      7579

H3K4me3
H3K9Ac
H3K4me1
H3K36me3
H3K27me3
H3K9me3
MeDIP
MRE

thap3      dnaj

H3K4me3
H3K9Ac
H3K4me1
H3K36me3
H3K27me3
H3K9me3
MeDIP
MRE

chr1 : 6604797 - 6610796

hyperlink

screenshot

UCSC Genome Browser on Human Mar. 2006 (NCBI36/hg18) Assembly

move  <<<  <<  <  >  >>  >>>  zoom in  1.5x  3x  10x  base  zoom out  1.5x  3x  10x

position/search  chr1:6,604,797-6,610,796    gene          jump  clear  size 6,000 bp.  configure

chr1 (p36.31)                31.1        1q12        41    43

56 _                                        H1 ES H3K4me3
ES H3K4me3
1 _
16 _                                        H1 ES H3K9Ac
ES H3K9Ac
1 _
18 _                                        H1 ES H3K4me1
ES H3K4me1
1 _
11 _                                        H1 ES H3K36me3
ES H3K36me3
1 _
No data                                     H1 ES H3K27me3
ES H3K27me3
No data
No data                                     H1 ES H3K9me3
ES H3K9me3
No data
23 _                                        H1 ES MeDIP
ES MeDIP
1 _
31 _                                        H1 ES MRE
ES MRE
1 _
                                            RefSeq Genes
PHF13                                       THAP3
                                            THAP3
                                            THAP3

# Spark

4. Interactive cluster visualization - data from human H1 embryonic stem cells



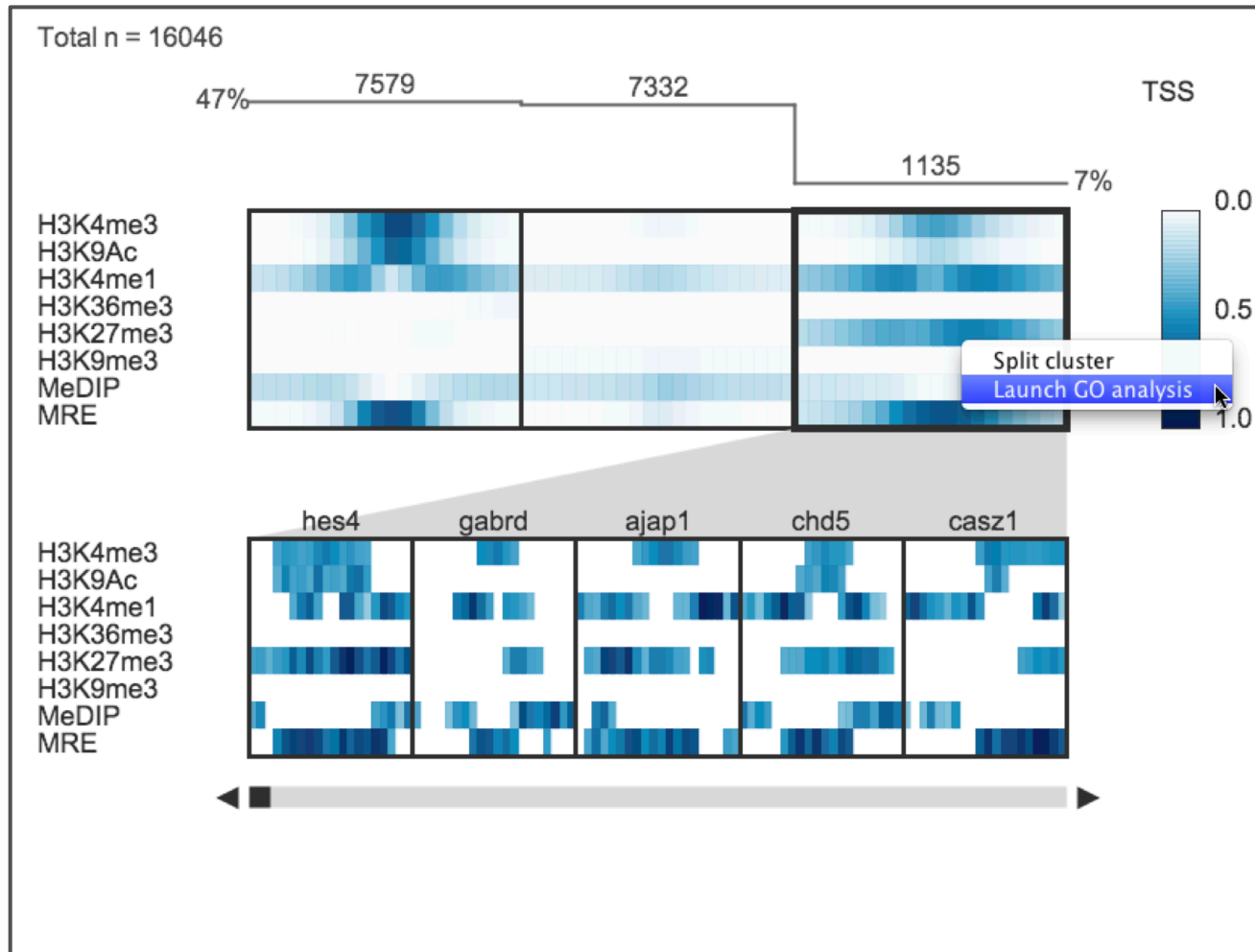activating + silencing
modifications:
"poised"

screenshot

# Spark

Do these clusters have functional meaning?

# Spark

5. Interactive gene ontology analysis



screenshot

# Spark

## 5. Interactive gene ontology analysis

# Spark

## 5. Interactive gene ontology analysis
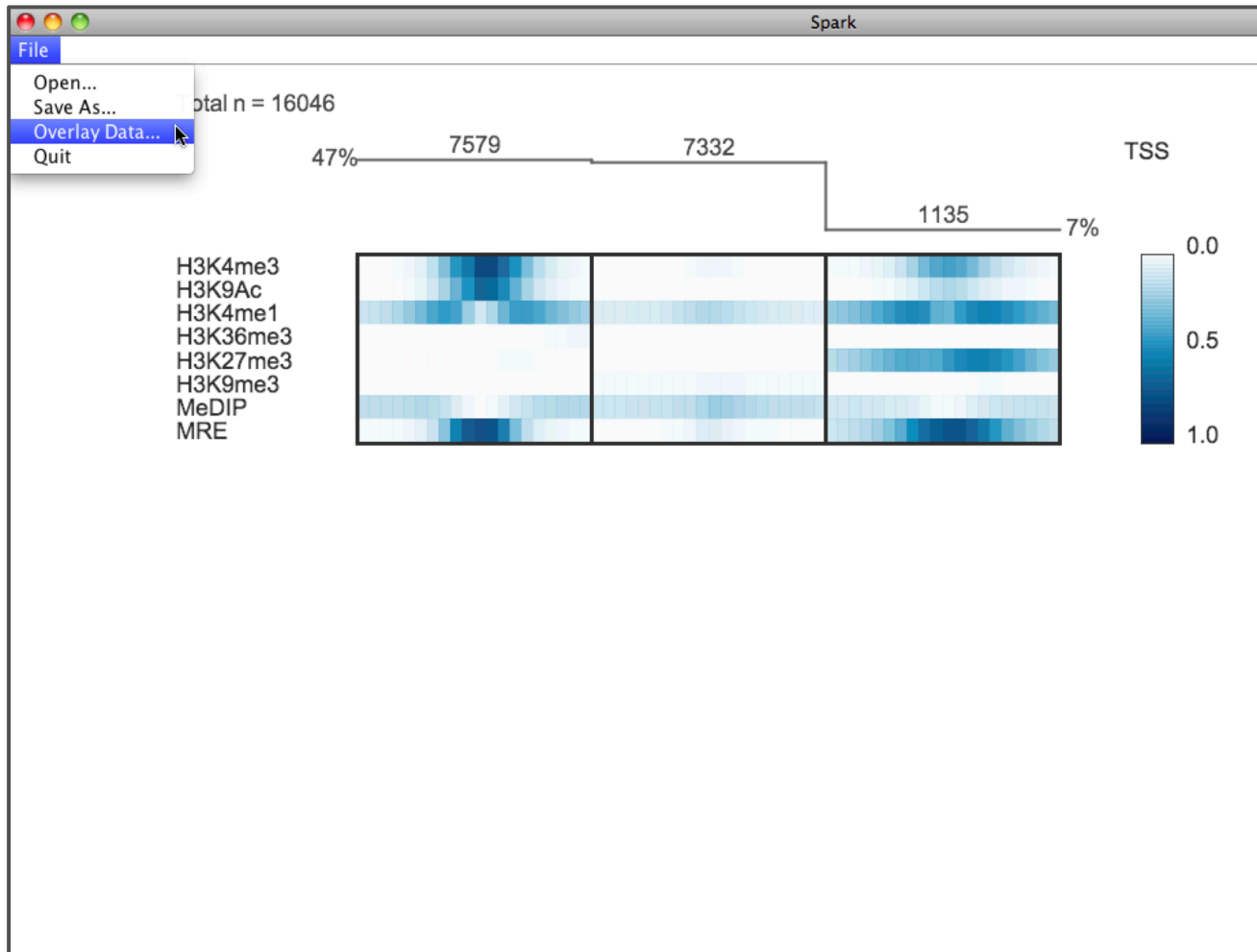
# Spark

Do these clusters have distinct gene expression patterns?
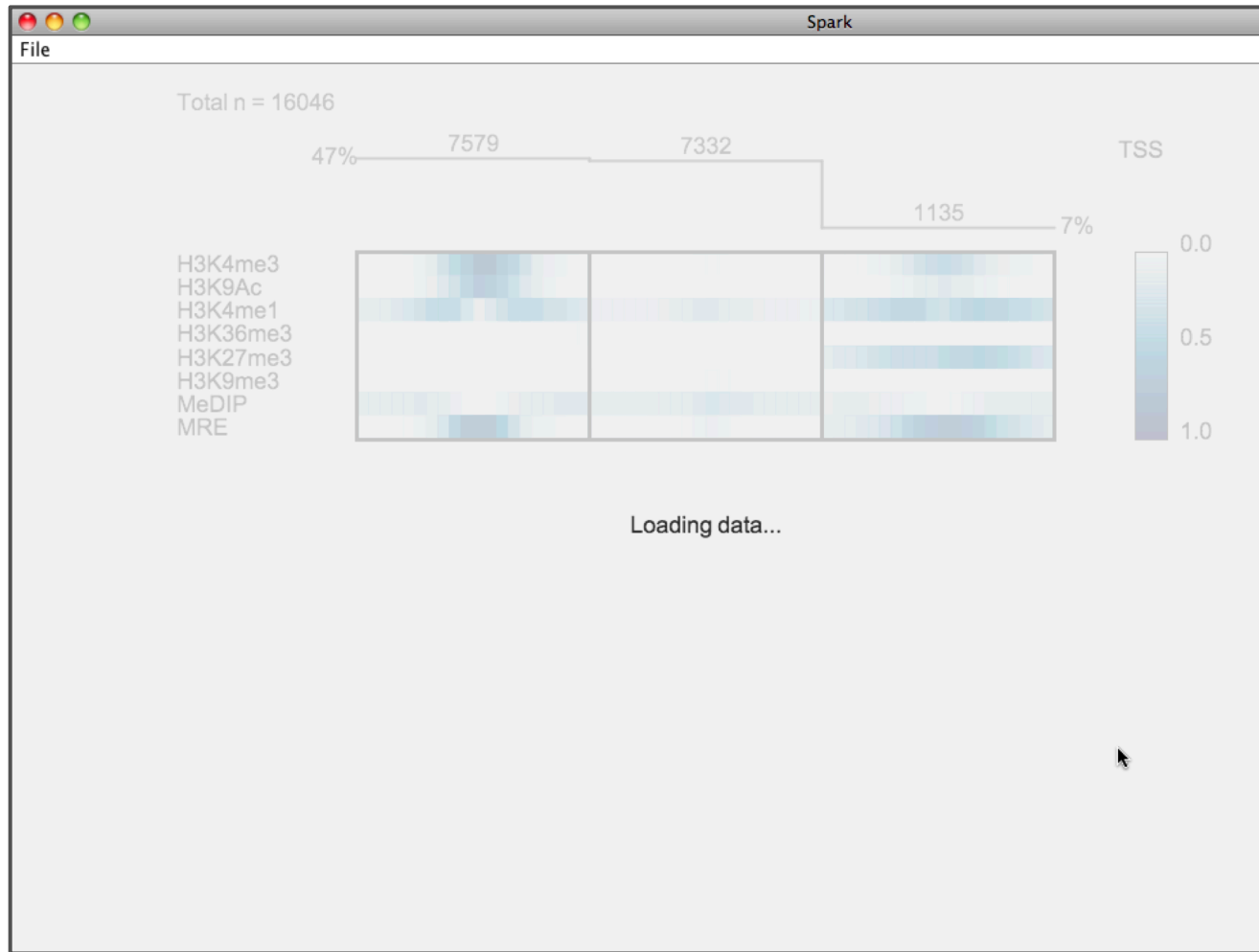
# Spark
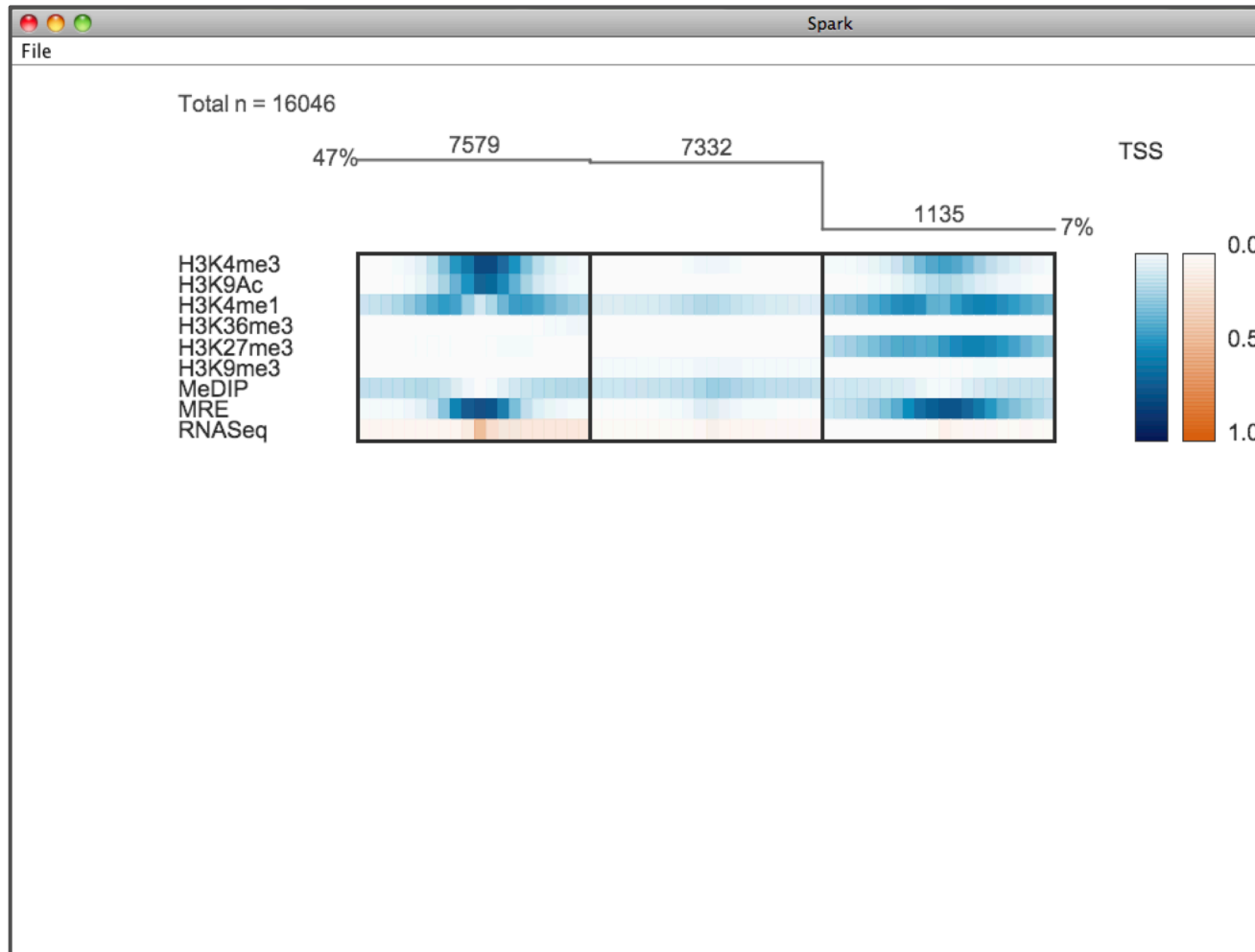
6. Interactive data overlay



screenshot

# Spark

6. Interactive data overlay



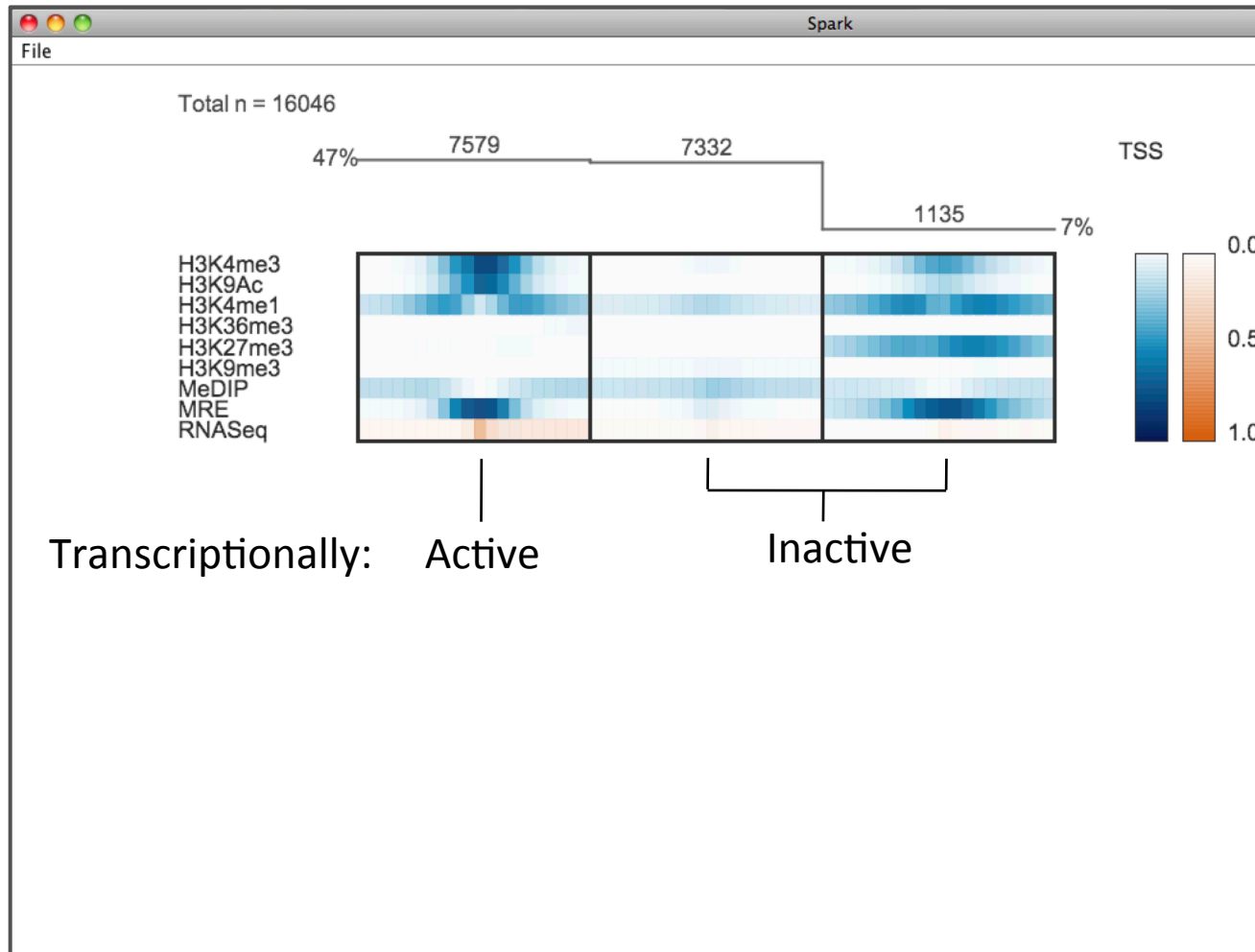screenshot

# Spark

6. Interactive data overlay



screenshot

# Spark

6. Interactive data overlay
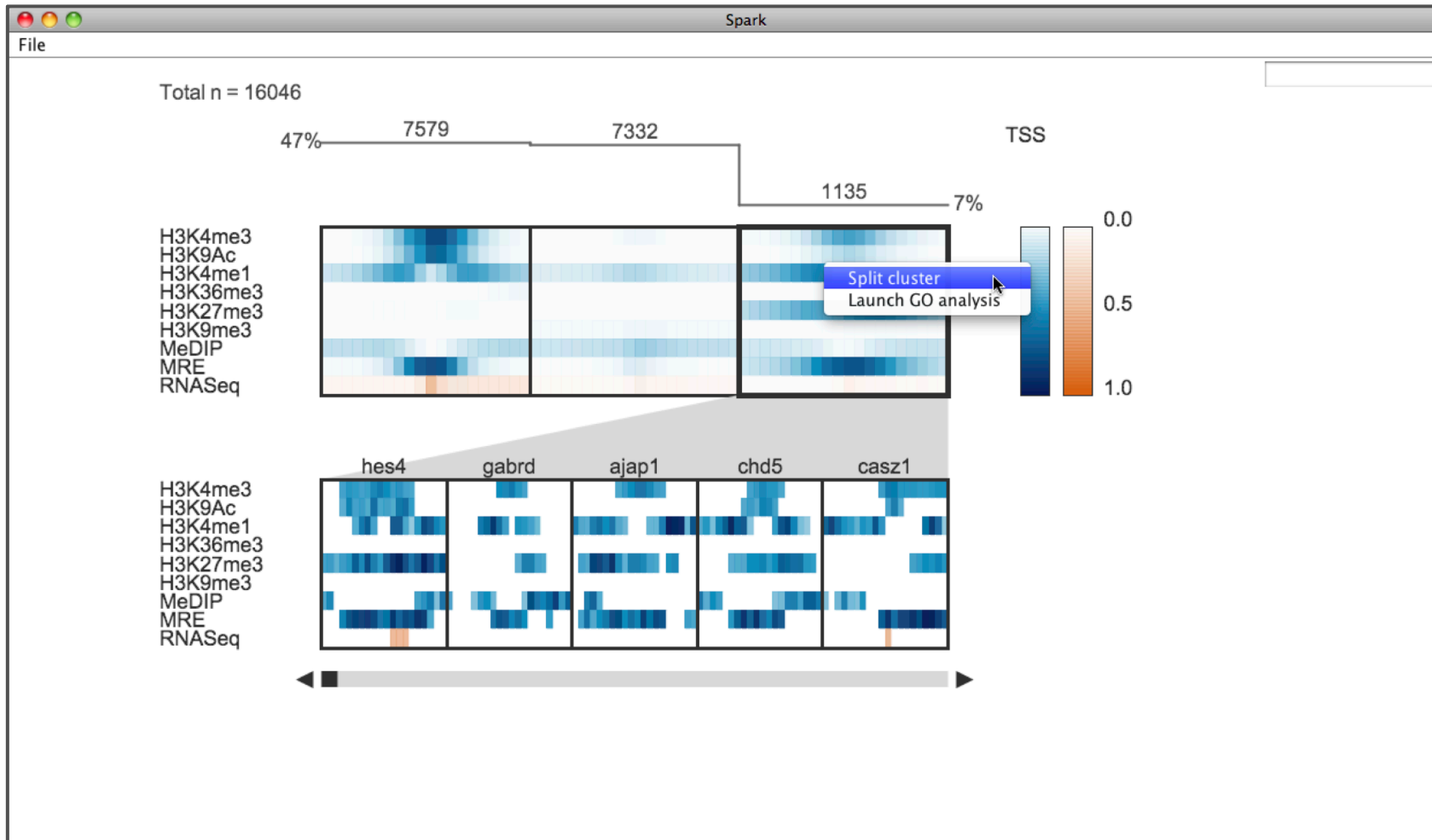


screenshot

# Spark

How many clusters should I generate?

# Spark

## 7. Interactive cluster splitting
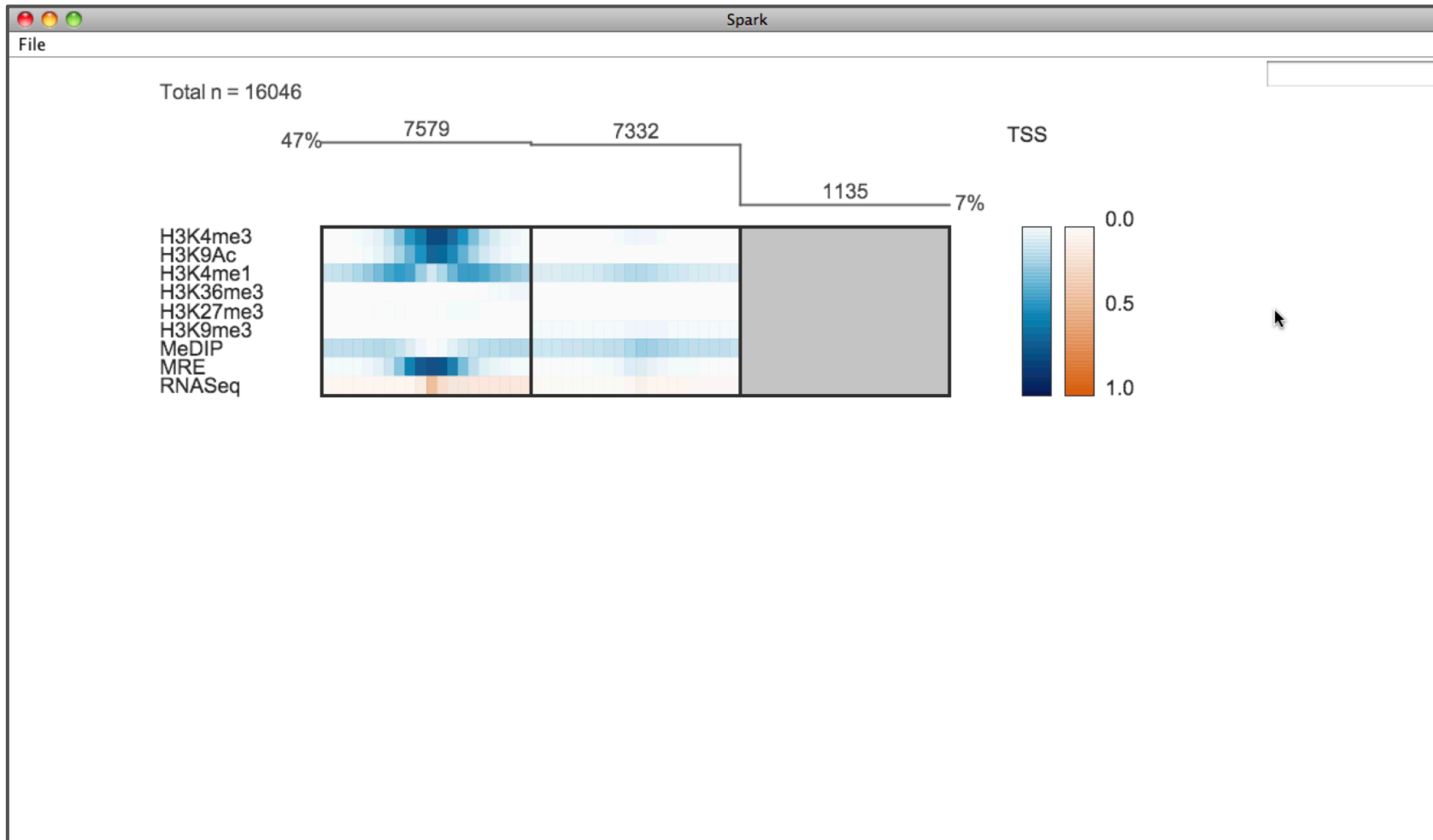


screenshot

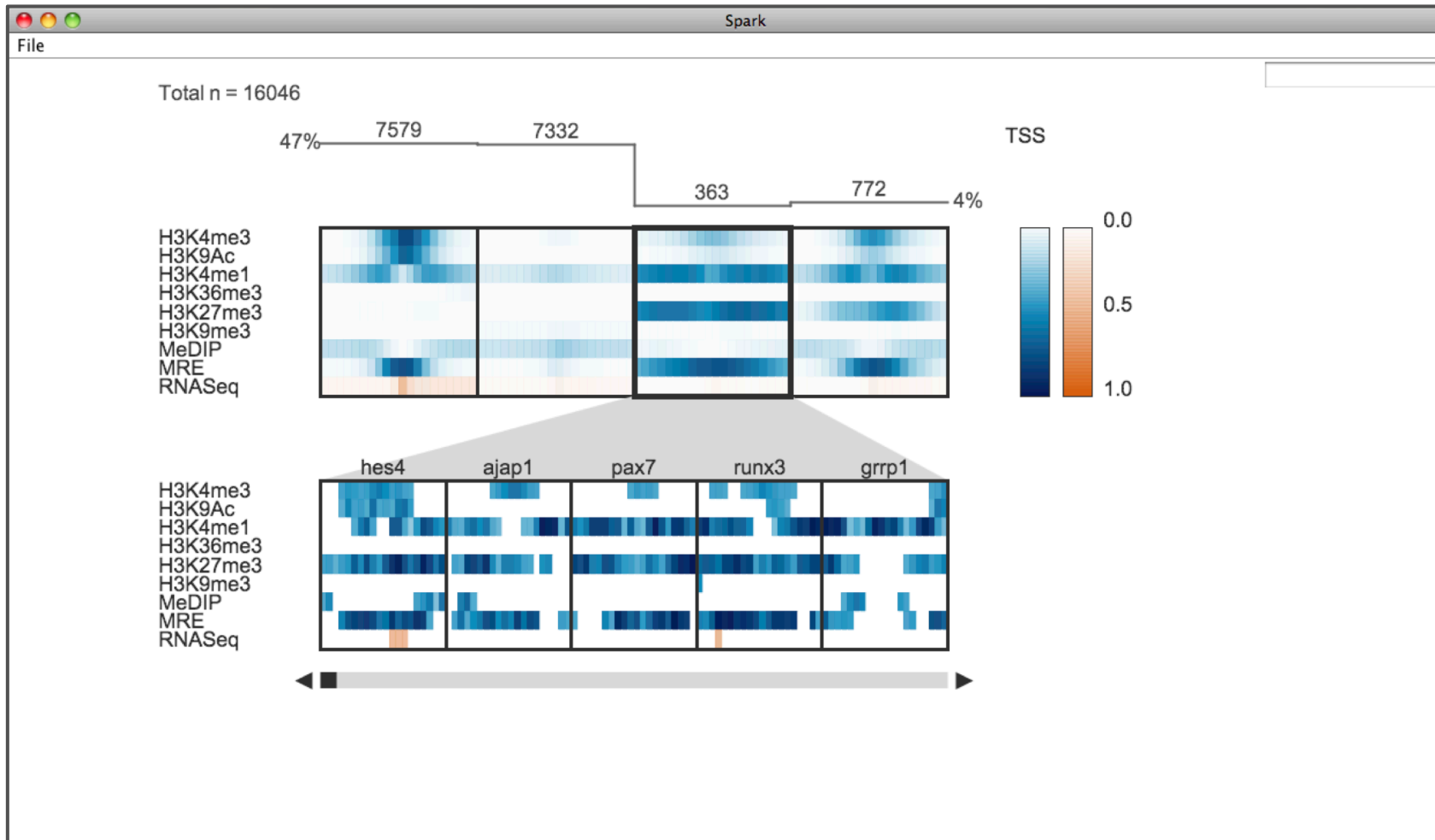# Spark

## 7. Interactive cluster splitting
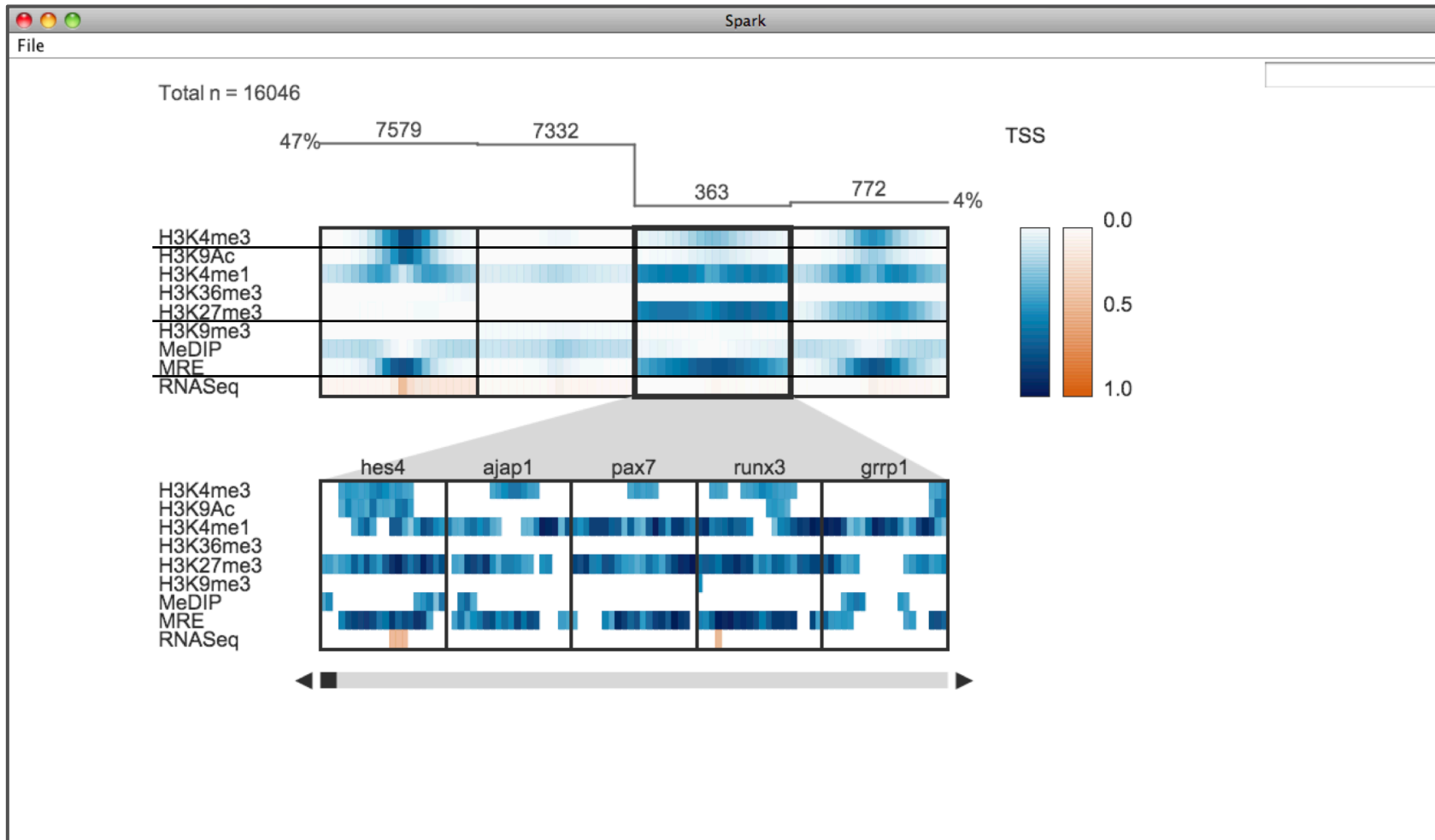


screenshot

# Spark

## 7. Interactive cluster splitting



screenshot

# Spark

7. Interactive cluster splitting
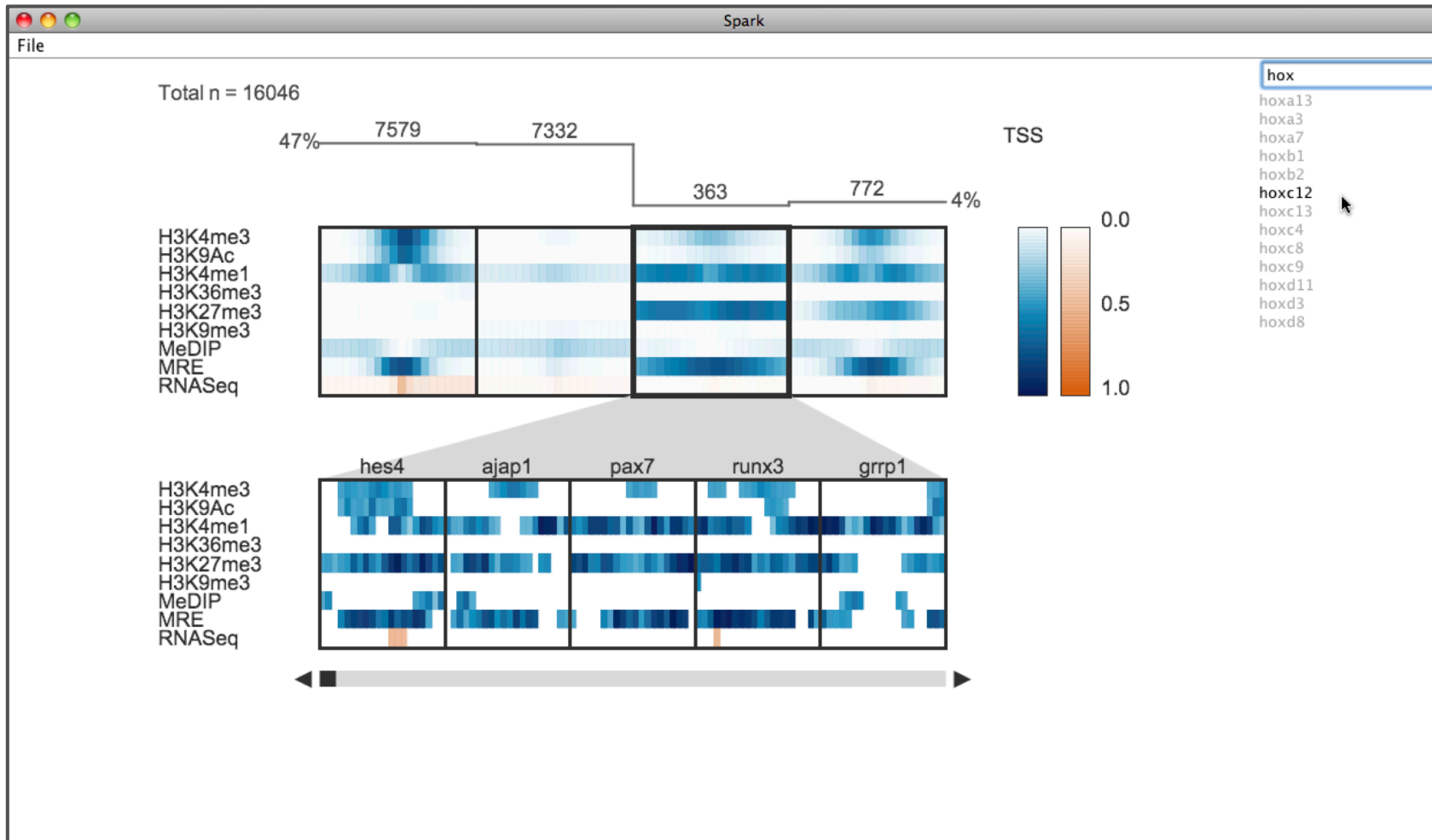


screenshot

# Spark

Where is my favorite gene?

What other genes have similar data patterns?

# Spark

## 8. Finding genes of interest



screenshot

# Spark

## 8. Finding genes of interest



screenshot

# Spark - technical details

# Spark - technical details

# Spark - technical details



Step 1

← Select from the Epigenome Atlas data tracks

← Or browse for local wig files or paste a URL

This Epigenome Atlas Tracks tree above is populated from an XML file containing the directory structure from www.genboree.org/epigenomeatlas and therefore can very easily be updated to reflect new data releases

# Spark - technical details

Your selections can then be added to the input data table below, in which you can edit the sample name and display colour.

# Spark - technical details



Step 2

Select from a provided set of region files

Or browse for local GFF file

# Spark - technical details

**Step 3 of 3: Settings**

Specify the output directory:

/Users/cydneyn/test_atlas

Browse...

Number of bins: 20

Number of clusters: 3

Normalize values: globally

Cancel    Go Back    Finish

**Step 3**

⟵ Specify where to save your analysis

⟵ Adjust clustering parameters or use provided defaults

In cases where the user selects an Epigenome Atlas data track and one of the provided region sets, Spark makes use of pre-computed data files that greatly increase performance.

If you want to use a custom region set with an Epigenome Atlas data track, Spark handles all of the data file downloading and caching to keep things simple.

You're done!
Spark will generate and display your clustering

# Demo

# Future work

- Heatmaps are not the best representation

# Future work

- Heatmaps are not the best representation



Same colour
looks different

Different colour
looks the same

* These rectangles
have the same colour
but look different

Bang Wong, Nature Methods, 2010

# Future work

- Heatmaps are not the best representation
  - Histograms would be better (also capture variation)

# Future work

- Heatmaps are not the best representation
  - Histograms would be better (also capture variation)
- Greater functionality for manipulating clusters
  - Merge clusters
  - Edit current settings

# Future work

- Heatmaps are not the best representation
  - Histograms would be better (also capture variation)
- Greater functionality for manipulating clusters
  - Merge clusters
  - Edit current settings
- Input files likely require filtering
  - could filter data based on an input peak set

# Future work

- Heatmaps are not the best representation
  - Histograms would be better (also capture variation)
- Greater functionality for manipulating clusters
  - Merge clusters
  - Edit current settings
- Input files likely require filtering
  - could filter data based on an input peak set
- More directly support data driven region selection

# Future work

- Heatmaps are not the best representation
  - Histograms would be better (also capture variation)
- Greater functionality for manipulating clusters
  - Merge clusters
  - Edit current settings
- Input files likely require filtering
  - could filter data based on an input peak set
- More directly support data driven region selection
- Spark may not perform well in detected a small subpopulation with a particular pattern

# Spark releases

**Public release (old)**

http://www.bcgsc.ca/platform/bioinfo/software/spark

**Current version (significantly updated; in testing)**

http://www.bcgsc.ca/downloads/spark/v1.1.0/start.jnlp

# Acknowledgements

**BCGSC**
Steve Jones (advisor)
Marco Marra
Martin Hirst
Misha Bilenky
Gordon Robertson
Yongjun Zhao
Martin Krzywinski

Sequencing Team
Bioinformatics Group
LIMS

**UCSF**
Joe Costello
Chibo Hong
Ravi Nagarajan
Thea Tlsty
Philippe Gascard
Mahvash Sigaroudinia
Art Weiss
Terri Kaldecek
Michael McManus
Hunter Richards
Yun Choi
Susan Fisher
Olga Genbacev

**UC Davis**
Peggy Farnham
Henny O'Geen
Lorigail Echipare
Vitalina Komashko
Kimberly Blahnik

**UCSC**
David Haussler
Tracy Ballinger

**Washington Univ**.
Ting Wang