# Nebula – A web-server for advanced ChIP-seq data analysis

# Tutorial

## by Valentina BOEVA

# Content

# Our web server: http://nebula.curie.fr/

Our web service, Nebula, is based on the Galaxy open source framework.



Tool box      Work field      History

Main Galaxy server: [ http://main.g2.bx.psu.edu/ ] **does not** include **all** our ChIP-seq analysis tools, but you can use it for other occasions.

# Create your account

- Each registered user have a 50Gb quota and unregistered user have a 15Gb quota (which is enough to run the tutorial with examples).

- We would prefer you to register even if you don't use your real email address.

# Download the test dataset to the history





- Select and import all datasets:



- Then go back by clicking "Analyze Data"

# Alternative way to download your dataset to the history



**This way you will use outside of this tutorial**

- To upload files larger than 2GB, the user has to use the URL method through FTP/HTTP protocol. The user must have access to an open web server or ftp server where he should upload his data. If the user does not have access to any web or ftp server, he can install his own web server.

  The following servers are free and can be easily installed:

  *Web servers:*
  MAMP for Mac (http://www.mamp.info/en/index.html)
  WAMP for Windows (http://www.wampserver.com/)

  *Ftp servers:*
  FileZilla Server for Windows (http://filezilla-project.org/download.php?type=server)
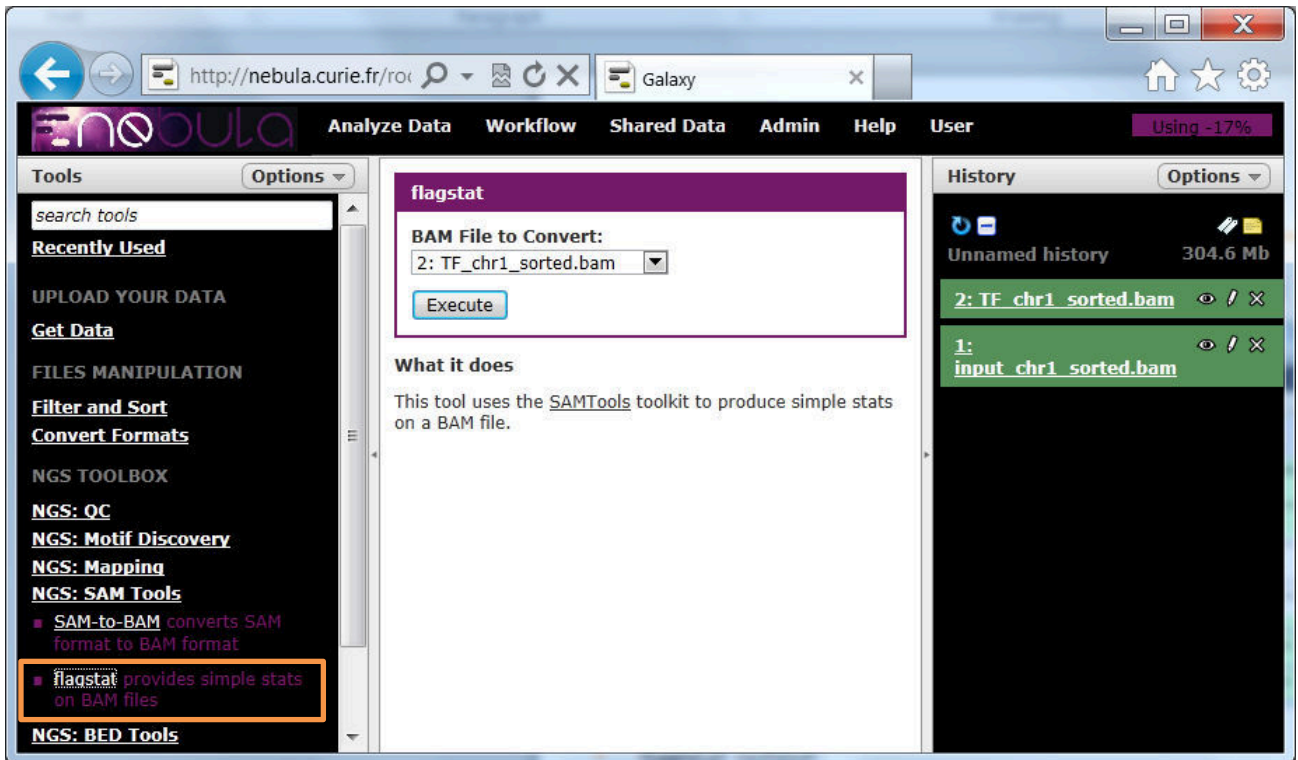  Pure-FTPd for Mac (http://www.pureftpd.org/project/pure-ftpd)

  Once the user has his own server installed, he can put his data on the server, copy the URL to the file (http://publicIP/path/to/file or ftp://user:passw@publicIP/path/to/file) and paste the URL into the URL Text box of the upload tool. After clicking on "execute", the upload will start.

- **A more complete tutorial can be found at the main Galaxy server: https://main.g2.bx.psu.edu/ -> Live Quickies: Uploading Data using FTP, Galactic quickie #17**

# Read statistics

- Run "flagstat" – to see how many reads were mapped



- flagstat output:

# Check read quality before calling peaks

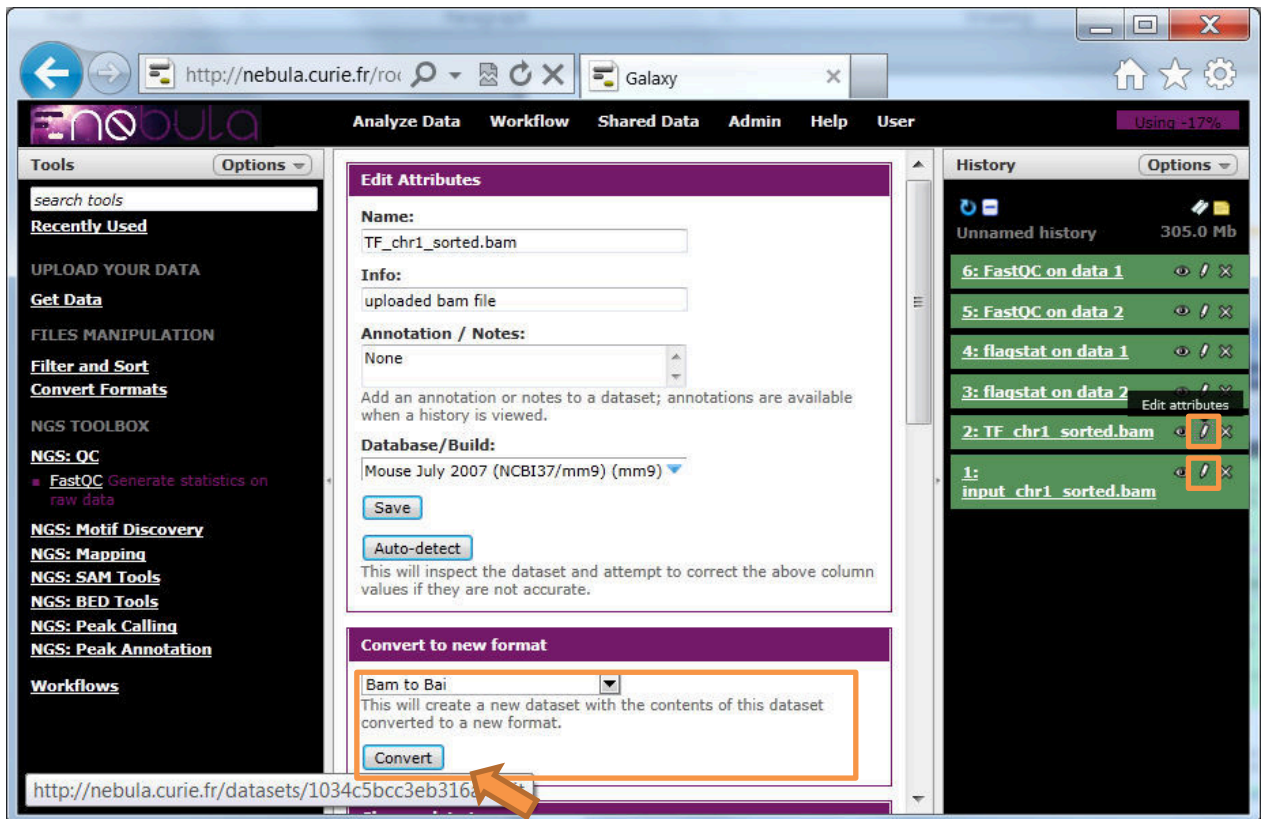- Run FASTQC – to see statistics on read quality



- Check FASTQC output:

# Check read quality before calling peaks

- Check how many reads you have in total by looking at the output of 'flagstat' (p. 7).
- How many reads were promised by the sequencing facilities? ☺
- I would say that 20 million mapped reads should be OK. In our example we have more then 2 million reads on chr1 (0.07 of the total mouse genome), this corresponds to about 30 million reads for the whole genome.
- Check the proportion of duplicate reads ('FASTQC', p. 8). High level of PCR duplicates means that you provided to little material for sequencing.
- Check whether you will have enough reads when you filter out duplicates. In our case we have about 30% of reads which are duplicates. Looking at the graph, we can say that filtering of duplicate reads will remove about 20% of reads. So it is still OK to continue our analysis and do peak calling. (MACS and FindPeaks will remove duplicate reads for you).

# Visualize .SAM/.BAM files in UCSC

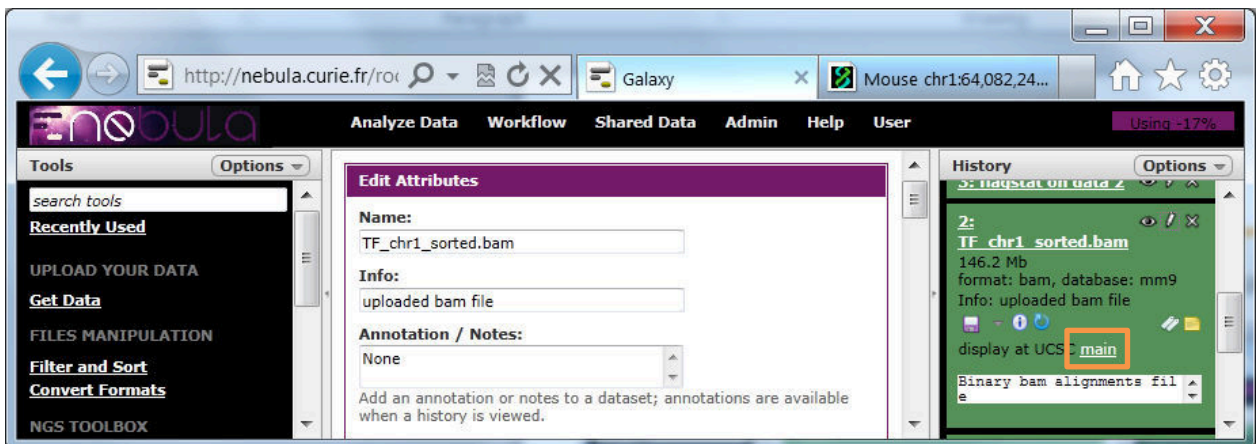- First create an index (.bai) for .BAM files



- Do it both for **TF_chr1_sorted.bam** and **input_chr1_sorted.bam** !

# Visualize .SAM/.BAM files in UCSC

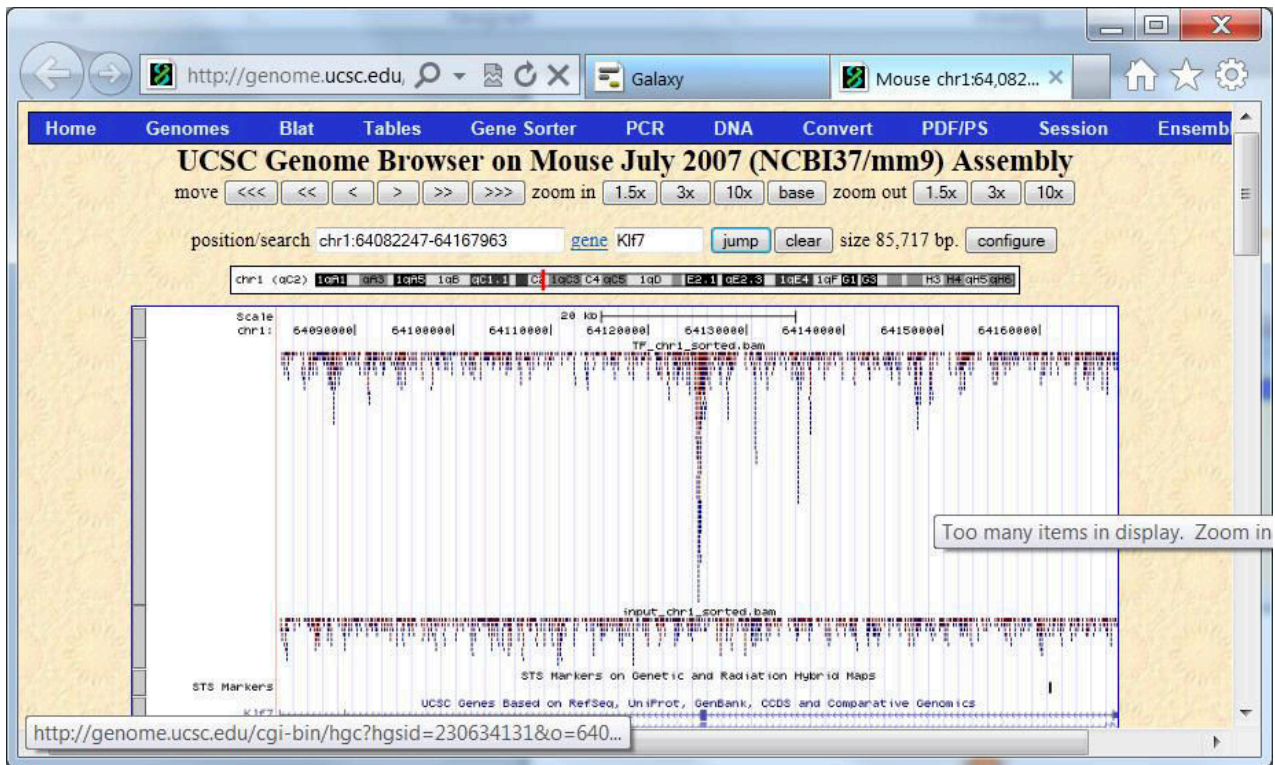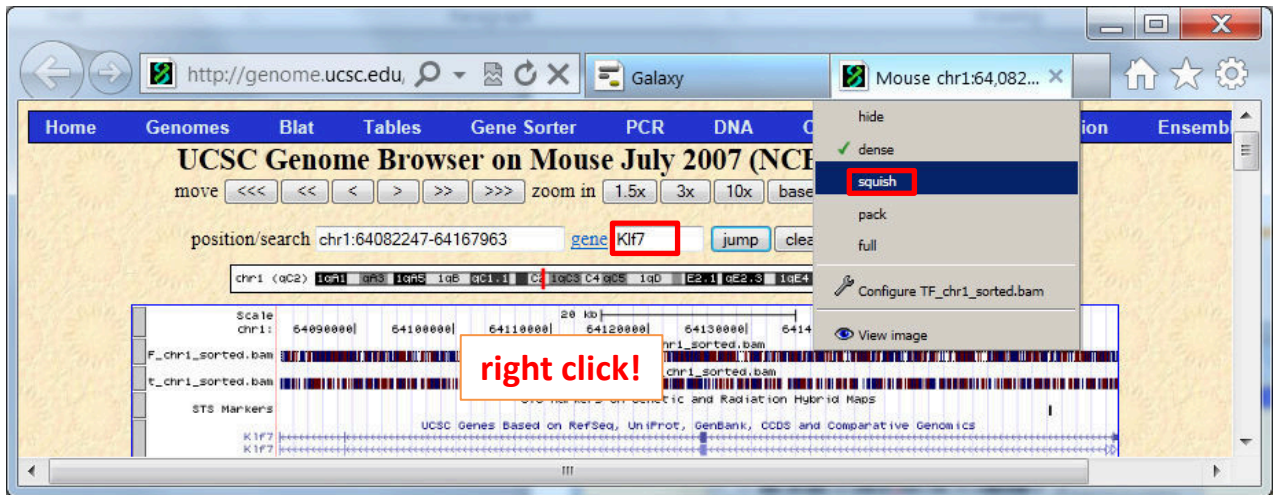- Click on 'main' of the initial .BAM file to visualize it in UCSC



- Do it twice: **TF_chr1_sorted.bam** and **input_chr1_sorted.bam** !

- Uploaded tracks will stay in your UCSC for several days. You can close and open the UCSC browser when you want and you won't lose your tracks.
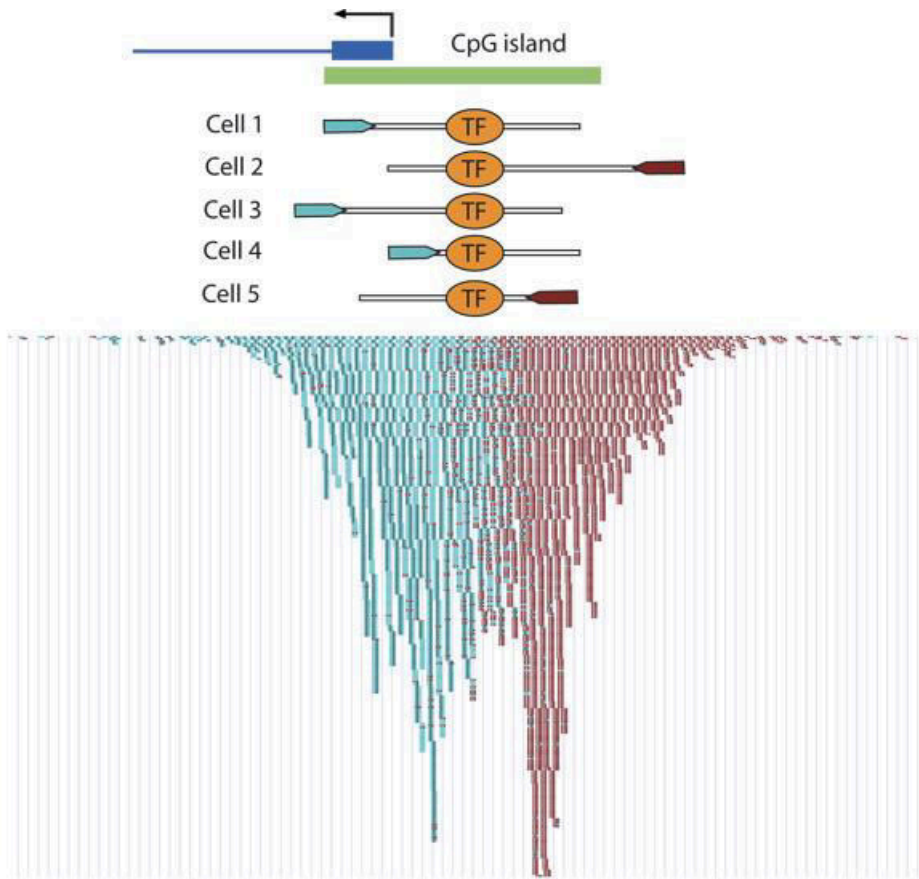
# Visualize .SAM/.BAM files in UCSC

- Go to the **Klf7** gene and change **view of the track**

# How does a good binging site look like?



(from Valouev et al., Nat Methods 2008)

- In our case the separation of forward and revers reads is not as clear. This is because it is SOLiD reads and we performed double sonication (one before and one after immunoprecipitation)

# There exist two main ways to extract the signal (construct peaks)

- ## Tools:

  - FindPeaks
  - QuEST
  - F-Seq
  - SICER

  - CisGenome
  - GLITR
  - PeakSeq
  - Spp

  - Useq
  - MACS
  - ERANGE
  - SiSSRs

- ## Main methods:



Adopted from S. Pepke et al., 2009 Nat Methods

(FindPeaks)                    (MACS)

# Run MACS (if you want to compare its output with the output of FindPeaks. **You can skip this step**.)



- Band width: This value is only used while building the shifting model. Should be ≥ DNA fragment lengths.
- For transcription factors, it is important to check '**Parse xls files into into distinct interval files**' to get the locations of peak summits for peak annotation.

# Run FindPeaks

- Create a subset of the control dataset if there are more reads for the control sample than for the ChIP sample



- This command will
  1. filter out duplicate reads from your ChIP and Control datasets,
  2. randomly select reads from the Control sample so that the total number of reads in both sample were equal.
  3. Transform .BAM into .SAM, because for some unknown reason FindPeaks does not like some .BAM …

- If you have the same number of reads in the ChIP and the control sample, you will be able to compare their outputs later on and filter out peaks detected in both datasets. Imagine, you have 10 times more reads in the control? – Then your real signal in the ChIP can appear weak…

# Run FindPeaks

- Run FindPeaks on the TF and Input sample (twice!)

# Calculate peak height distribution – immunoprecipitation quality control

- You should enter FindPeaks output files **(.peaks)** for the **TF** and **Input**



- You can the select the minimal peak height for further analyses using the calculated evaluation of false discovery rate:

# Calculate peak height distribution – immunoprecipitation quality control

- More high peaks in the ChIP sample – the better the immunoprecipitation was preformed

# Filter FindPeaks' output using peaks from the control dataset

- The actual peak shapes is replaced by triangles (start, end, maximum and height). Then, the height ($x$) of maximal overlap is calculated. The ChIP peak is rejected if its height ($h1$) divided by $x$ is less than or equal to a given threshold.



$h1/x > 2$? ➡ Keep the peak

# Convert FindPeaks output (.peaks) into Bed
## if you did **not** select output in .BED at the previous step

- Convert "filtered" peaks into .BED (.BED is a standard format for genomic intervals):



- Convert the "**control**" peaks too. One should use a low threshold on peak height (we will further use these peaks as "random" control for peak location distribution):

# Visualize .bed in UCSC

- For .BED: visualize directly in UCSC

# Visualize .wig in UCSC

- For .wig.gz (output of FindPeaks): visualize directly in UCSC



- For .Wig (output of MACS): you need to convert .wig to .bw (Big Wig) first and then you will visualyze the BigWig file:

# Get .fasta sequences to find over-represented motifs

- Create .bed with coordinates of **central regions** of peaks (FindPeaks output: use .bed file)



- If you want to extract central regions for MACS use "peaks: interval" file instead of "peaks: bed", since the former contains information about peak summits:

# Get .fasta sequences to find over-represented motifs

- Extract .fasta



- Extract .fasta for MACS central peak regions too if you used MACS peak calling:

# Run motif finding on central regions of peaks



- If you also created .fasta for the MACS peaks (p. 24), you can run motif finding on them too:

# Run motif finding on central regions of peaks

- Motifs found in peaks identified by FindPeaks (200bp central region, use UnZoom to see it better):

## Mask sequences ("filter")

**Looking for several motifs of one TF**

~~cgatcgaga~~**CAGGAATG**~~gct~~**agat**~~a~~
~~cacatgtac~~**CAGGAATC**~~cg~~**agat**~~at~~
~~acg~~**agat**~~cg~~**CAGGAAAG**~~gctacgat~~
~~cacat~~**agat**~~CCGGAATG~~~~cgatgcat~~
~~actgcgctg~~**CAGGAATG**~~agct~~**agat**
cac**agat**GGAAGGAAGGAAatgcat
**agat**cgcGGAAGGAAGGAActagca

| CAGGAATG | GGAAGGAAGGAA |
| CAGGAATC | GGAAGGAAGGAA |
| CAGGAAAG | |
| CCGGAATG | |
| CAGGAATG | |

| Motif 1 | Motif 2 |

## Mask motifs ("mask")

**Looking for motifs of co-factors**

cgatcgaga**CAGGAATG**gct**agat**a
cacatgtac**CAGGAATC**cg**agat**at
acg**agat**cg**CAGGAAAG**gctacgat
cacat**agat**CCGGAATGcgatgcat
actgcgctg**CAGGAATG**agct**agat**
cac**agat**GGAAGGAAGGAAatgcat
**agat**cgcGGAAGGAAGGAActagca

| CAGGAATG | **agat** |
| CAGGAATC | **agat** |
| CAGGAAAG | **agat** |
| CCGGAATG | **agat** |
| CAGGAATG | **agat** |
| | **agat** |
| | **agat** |

| Motif 1 | Motif 2 |

- Here we used the "mask" mode for motif finding (p. 25)

# Which known transcription factors correspond to identified motifs?



- Run TOMTOM: http://meme.sdsc.edu/meme/cgi-bin/tomtom.cgi

- TOMTOM can be found by **google** using "**TOMTOM motif**"

- Select the type of motif "mask"

- Copy-paste your motif from Galaxy (remove motif name and nucleotide names in the beginning of each line)

# Calculate distribution of peak locations around gene TSS

- Select the **.bed** files for FindPeaks filtered peaks and control.
- Use ProbeSets_FC1.5_10022011.txt file with information about activated/repressed genes (this file you have uploaded to your history in the very beginning)



ChIP          Control          ChIP vs Control

# Format of gene expression/modulation file

- Tab-delimited:
1. Some field
2. Gene symbol
3. Some value (expression or fold change)
4. Gene feature: e.g., expressed, regulated, etc.

# Annotate peaks with genomic features

# Annotate peaks with genomic features

# Run motif finding for peaks with selected genomic features



- **Use:** (c7=='promoter' or c7=='immediateDownstream') and c9=='up-regulated'



- Select central regions of peaks:

# Run motif finding for peaks with selected genomic features

- Get .fasta



- Run Motif Finding



- Visualize motifs

# Annotate genes with peak information

# Our workflow

# Analysis of histone data

- Peak calling
- Peak visualization
- Peak statistics
- Gene annotation

Test Data: H3K27me3 (H.Ashoor, 2013) for a bladder cancer cell line. Data provided for chromosome 1 only.

# Upload data to the history

Test Data: H3K27me3 (H.Ashoor, 2013) for a bladder cancer cell line. ONLY chr1.

Download the test dataset:





Check read number and sequencing quality                          pp. 7-9

Visualize .BAM files in UCSC                                      pp. 10-13

# Read statistics

- Use Samtools "flagstat" or "FastQC" to get statistics about reads

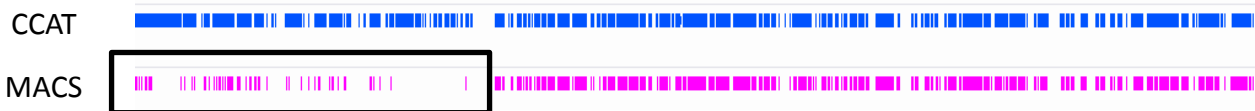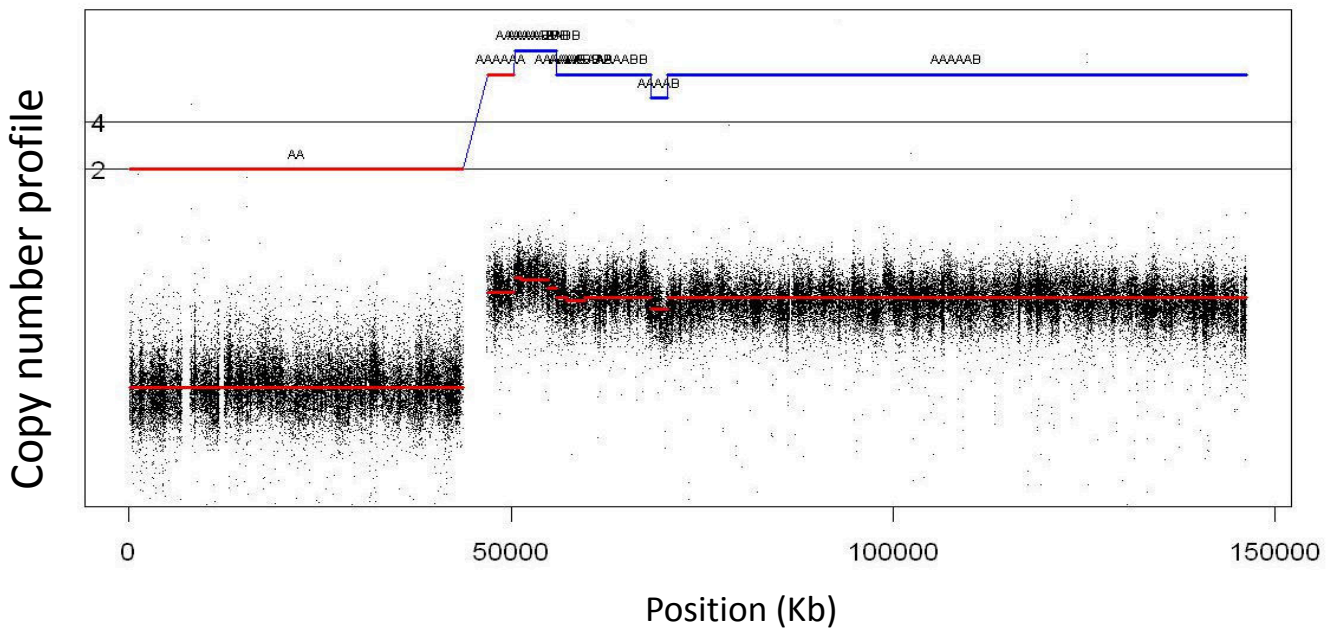Check read number and sequencing quality pp. 7-9

Visualize .BAM files in UCSC pp. 10-13

# Peak calling for histone data

- To call peaks, use MACS or CCAT
- For cancer datasets, use CCAT since it does not show copy number bias:



low density of sites in predictions of MACS in the region of low copy number

- MACS generates .wig files

# Peak calling with MACS for histone data

# Visualize the .WIG

# Visualize the .WIG



Unzoom

Big zoom 1

Big zoom 1

# Visualize MACS peak as well as .WIG profiles

- Visualize the .BED file too:

# Peak calling with CCAT for histone data
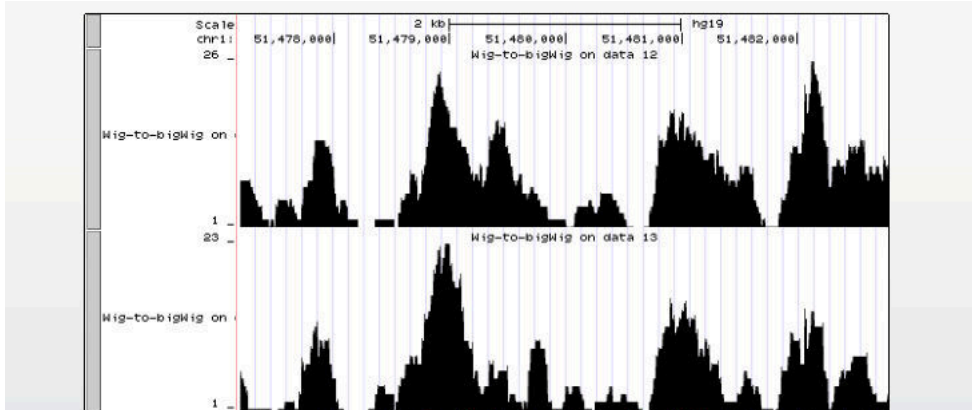
- Transform BAM to BED



**Repeat for ChIP and control .bam!!!**

- Run CCAT



Change!

# Peak calling with CCAT for histone data

- Run CCAT:

# Peak calling with CCAT for histone data

- CCAT provides important information about noise in the ChIP-seq data

# Peak calling with CCAT for histone data

- Transform strange output format of CCAT (chr center start end reads_chip reads_control FC FDR) into .Bed:

# Peak calling with CCAT for histone data

- Visualize .Bed file:

# Annotate predicted histone marks

- ## Narrow marks:
  - ### H3K4me3
  - ### H3K4me1

  > Where are they located respectively to the gene transcription starts?

  > The same kind of analysis as for transcription factors!

- ## Large marks
  - ### HK36me3
  - ### H3K27me3
  - ### H3K9ac

  > Are gene transcription starts covered by these marks? How much of the gene body is covered?

# Annotate predicted histone marks

- ## Narrow marks:



peak center => putative binding

peak in promoter

gene

Promoter region

- ## Large marks



peak in promoter & in gene body

gene

Promoter region

# Format of gene expression/modulation file

- Tab-delimited:
1. Some field
2. Gene symbol
3. Some value (expression or fold change)
4. Gene feature: e.g., expressed, regulated, etc.

# Check peak distribution around gene transcription start sites (TSS)

# Check peak presence in genomic regions

# Check peak presence in genomic regions



Gene TSS region

Gene body

Promoter

Enhancer

Immediate downstream

Intragenic

Gene downstream

# Annexes

**File formats:**

- .BAI – index for a .BAM file (to visualize .BAM in UCSC)
- .BAM – aligned reads, binary .SAM
- .BED – genomic coordinates

```
chr1    23386792    23387348    23387022    8.336    +
chr1    24005015    24005240    24005133    8.680    +
chr1    36187196    36187544    36187322    12.0     +
```

- .BW (BIG WIG) – signal profile (to visualize .WIG in UCSC)
- .CSFASTA – read sequences in color code

```
>921_41_109_R17C7_F3
T2130102221132101221333213002121321220223132222222
```

- .FASTA – DNA sequences

```
>chr1:3525467-3526150
ACTGGGTAAATAGCAGGTAGCAATTTTATGCAGAGGTTGGAGCTCACTTGGAACACACTTCCACCTTTG
```

- .FASTQ – read sequences and qualities

```
@921_29_592_R17C7
T1002011.2202001221203112001111321213112113222200220
+
&2*8411!6%'#,)##)$'#*&-5&&-&4-&%&,$&+%*$$+-,&'&4)#
```

- .SAM – aligned reads

```
119_171_1134_R17C1_    0    chr1    3000539  255    50M    *
286_1719_1498_R17C1_   0    chr1    3000539  255    50M    *
391_1794_580_R17C1_    0    chr1    3000539  255    50M    *
```

- .QUAL – read qualities

```
>921_41_132_R17C7_F3
8 13 3 5 6 4 4 4 10 6 2 13 5 8 6 5 6 11 7 3 5 13 15 5 5 2 5 10 6 15 12
```

- .WIG – signal profile (for visualization)

```
variableStep chrom=chr1 span=15
3000181 1
3000196 1.1
3000211 1.3
3000226 1.8
3000241 1.4
```