

Introduction to NGS analyses

Giorgio L Papadopoulos

Institute of Molecular Biology and Biotechnology

Bioinformatics Support Group

04/12/2015

Overview

- 1 Introduction
 - DNA sequencing
 - Applications

Overview

- 1 Introduction
 - DNA sequencing
 - Applications
- 2 Primary Analysis
 - ChIPseq
 - RNAseq

Overview

- 1 Introduction
 - DNA sequencing
 - Applications
- 2 Primary Analysis
 - ChIPseq
 - RNAseq
- 3 Interpretation
 - Visualization
 - Analysis Software
 - Downstream Analysis

Overview

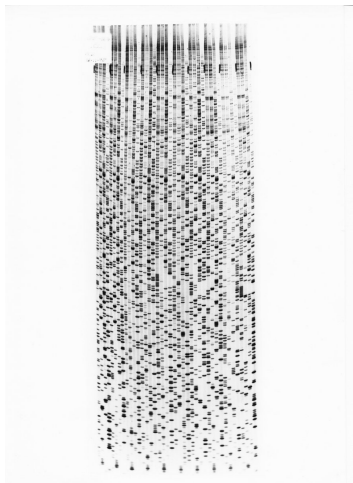
- 1 Introduction
 - DNA sequencing
 - Applications
- 2 Primary Analysis
 - ChIPseq
 - RNAseq
- 3 Interpretation
 - Visualization
 - Analysis Software
 - Downstream Analysis
- 4 Online Resources
 - ENCODE
 - GEO
 - ENA



Fred Sanger
Late 1970s

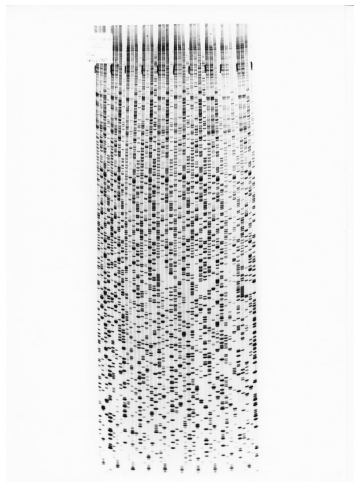


Fred Sanger
Late 1970s





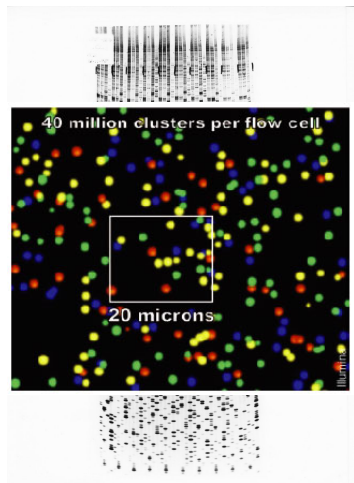
Fred Sanger
Late 1970s



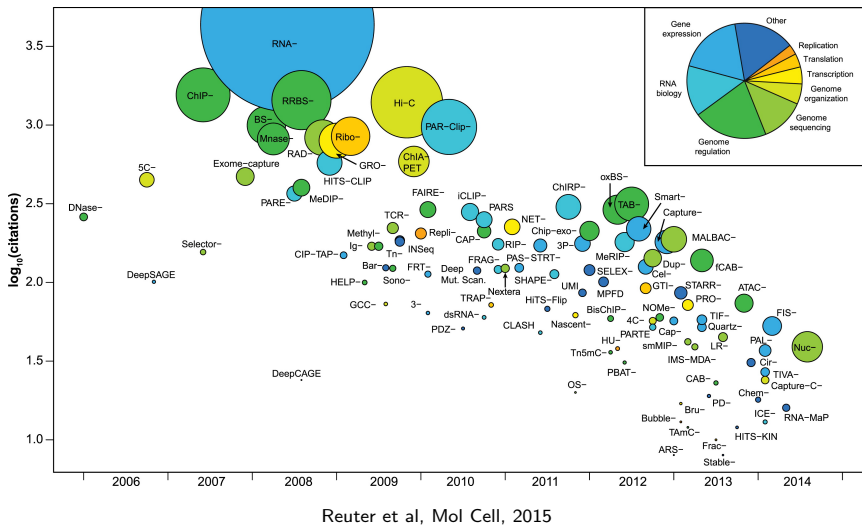
... Several million times ...
Late 2000s

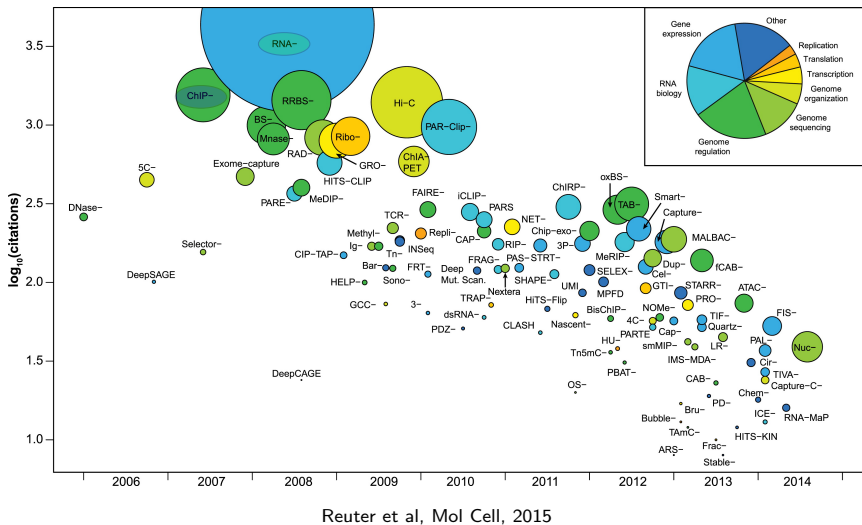


Fred Sanger
Late 1970s

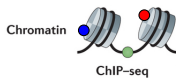


... Several million times ...
Late 2000s

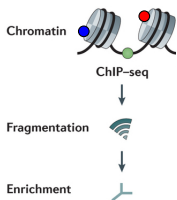




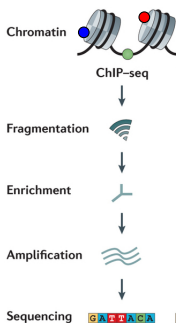
Genome Wide Occupancy Profiling



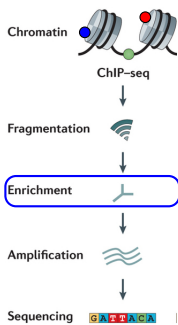
Genome Wide Occupancy Profiling



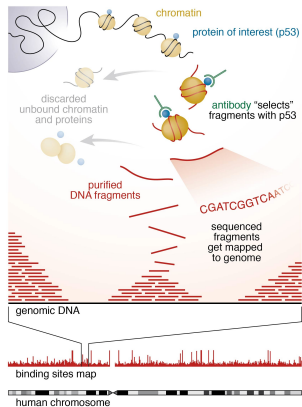
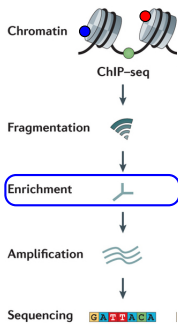
Genome Wide Occupancy Profiling



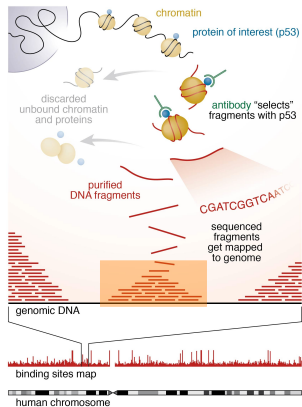
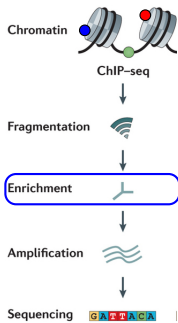
Genome Wide Occupancy Profiling



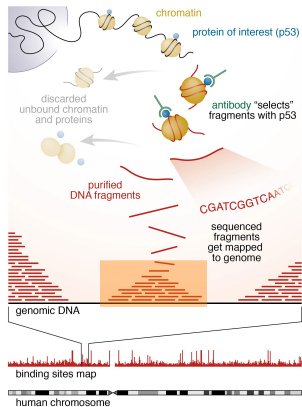
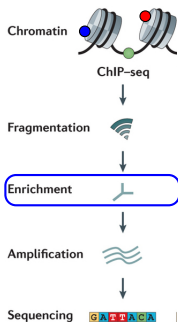
Genome Wide Occupancy Profiling



Genome Wide Occupancy Profiling



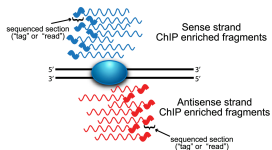
Genome Wide Occupancy Profiling



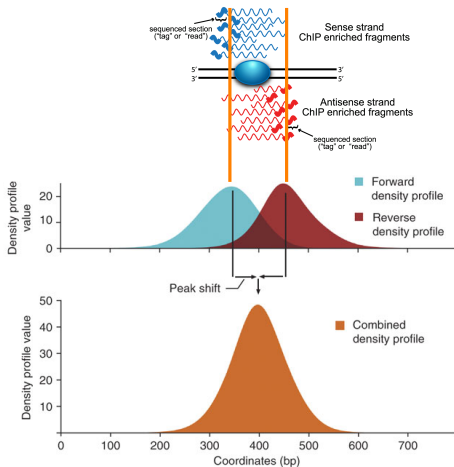
Overlapping reads come from different cells!!!

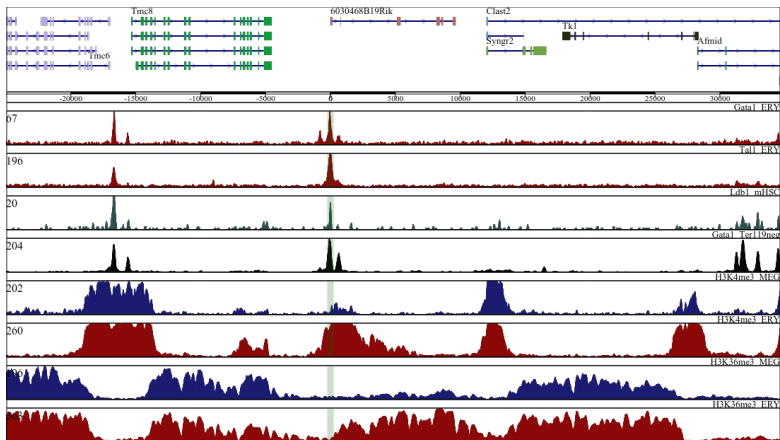
NGS is a cell population readout (50 million cells)

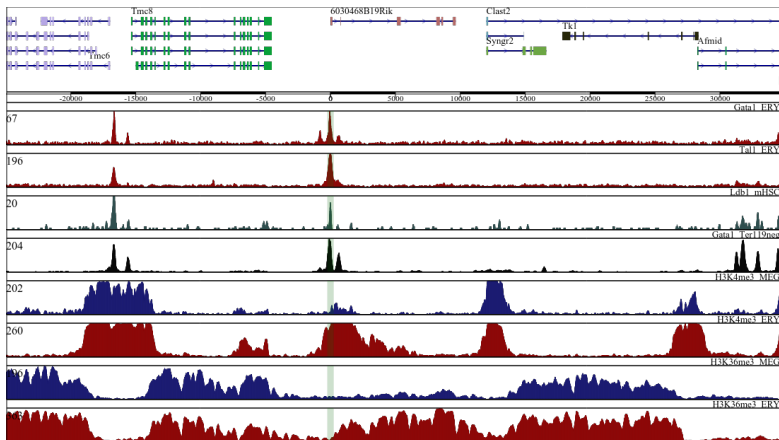
Primary Sequencing Data



Primary Sequencing Data

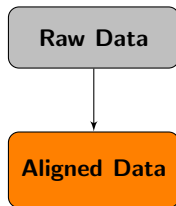


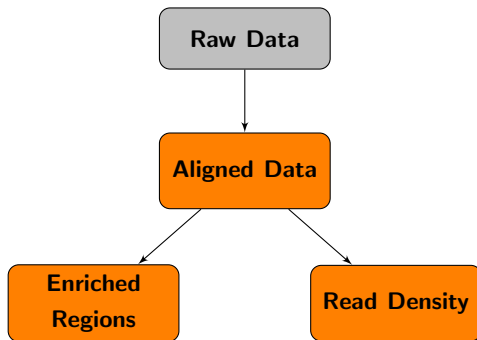


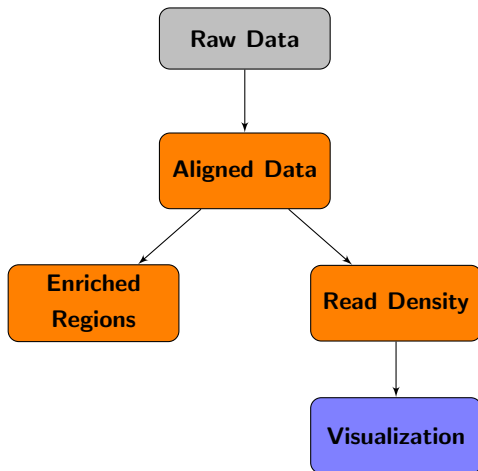


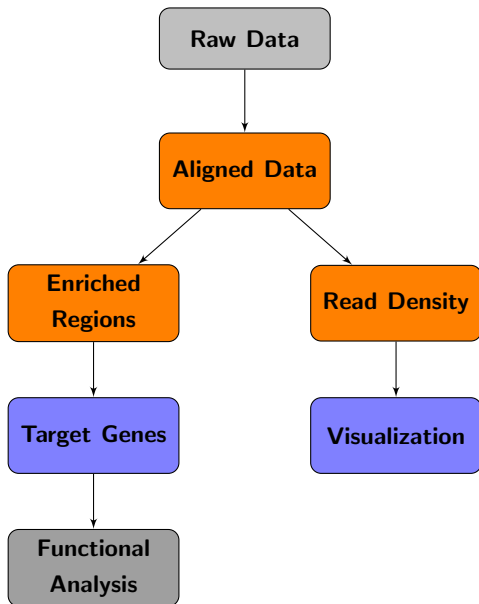
Genome wide, highly accurate representation of genomic states
 Fairly unbiased
 Complementary and overlapping information

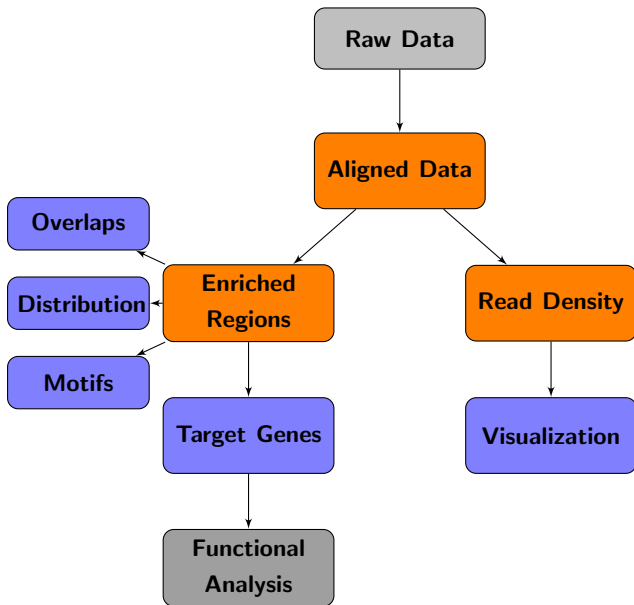
Raw Data











Raw Data

Raw Data

What is FASTQ?

- Text-based format for storing both biological sequences and corresponding quality scores.
- FASTQ = FASTA + QUALITY
- A FASTQ file uses four lines per sequence.

```
1 @SEQ_ID
2 GATTTGGGGTTCAAAGCAGTATCGATCAAA
3 +SEQ_ID(Optional)
4 !'*(((('*'+))%%#+) (%%%) .1**
```

Raw Data

What is FASTQ?

- Text-based format for storing both biological sequences and corresponding quality scores.
- FASTQ = FASTA + QUALITY
- A FASTQ file uses four lines per sequence.

```
1 @SEQ_ID
2 GATTTGGGGTTCAAAGCAGTATCGATCAAA
3 +SEQ_ID(Optional)
4 !'!*((( (***) )%%%+) (%%%) .1**
```

Operations:
DeMultiplexing
Quality Control

Raw Data

What is FASTQ?

- Text-based format for storing both biological sequences and corresponding quality scores.
- FASTQ = FASTA + QUALITY
- A FASTQ file uses four lines per sequence.

```
1 @SEQ_ID
2 GATTTGGGGTTCAAAGCAGTATCGATCAAA
3 +SEQ_ID(Optional)
4 !'!*((( (***) )%%#+) (%%%) .1**
```

FastQC

Operations:
DeMultiplexing
Quality Control

Raw Data

What is FASTQ?

- Text-based format for storing both biological sequences and corresponding quality scores.
- FASTQ = FASTA + QUALITY
- A FASTQ file uses four lines per sequence.

```












1 @SEQ_ID
2 GATTTGGGGTTCAAAGCAGTATCGATCAAA
3 +SEQ_ID(Optional)
4 !'!*((( (***) )%%%+) (%%%) .1**

```

Operations:

DeMultiplexing
Quality Control

FastQC

-  [Basic Statistics](#)
-  [Per base sequence quality](#)
-  [Per sequence quality scores](#)
-  [Per base sequence content](#)
-  [Per base GC content](#)
-  [Per sequence GC content](#)
-  [Per base N content](#)
-  [Sequence Length Distribution](#)
-  [Sequence Duplication Levels](#)
-  [Overrepresented sequences](#)
-  [Kmer Content](#)

Raw Data

What is FASTQ?

- Text-based format for storing both biological sequences and corresponding quality scores.
- FASTQ = FASTA + QUALITY
- A FASTQ file uses four lines per sequence.

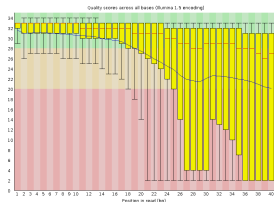
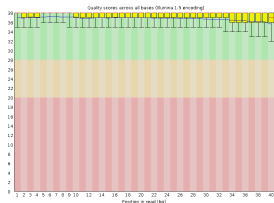
```

1 @SEQ_ID
2 GATTTGGGGTTCAAAGCAGTATCGATCAAA
3 +SEQ_ID(Optional)
4 !'!*((( (**++) )%%%+) (%%%) .1**

```

Operations:
DeMultiplexing
Quality Control

FastQC



Raw Data

What is FASTQ?

- Text-based format for storing both biological sequences and corresponding quality scores.
- FASTQ = FASTA + QUALITY
- A FASTQ file uses four lines per sequence.

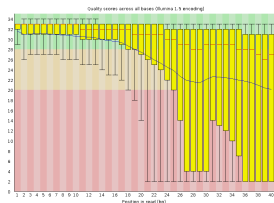
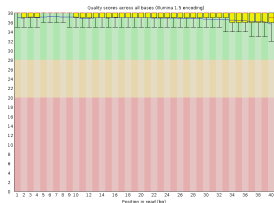
```

1 @SEQ_ID
2 GATTTGGGGTTCAAAGCAGTATCGATCAAA
3 +SEQ_ID(Optional)
4 !'!*((( (**++) )%%%%) .1**

```

Operations:
DeMultiplexing
Quality Control

FastQC



Align Data...

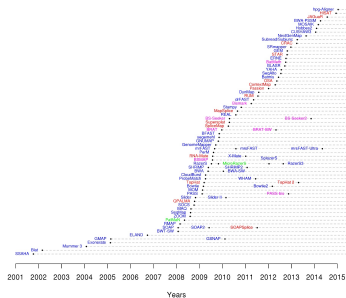
Align Data...

Align Data...

Aligner (Mapper)

Align Data...

Aligner (Mapper)



http://www.ebi.ac.uk/~nf/hts_mappers/

Align Data...

Aligner (Mapper)

Reference Genome

Human (hg18, hg19)

Mouse (mm9, mm10)

...

(iGenomes, UCSC, NCBI, EBI...)

**Different mappers will use
different Index Formats...**

Align Data...

Aligner (Mapper)

bowtie.2

Reference Genome

Human (hg18, hg19)

Mouse (mm9, mm10)

...

(iGenomes, UCSC, NCBI, EBI...)

**Different mappers will use
different Index Formats...**

Align Data...

Aligner (Mapper)

bowtie.2

Reference Genome

Human (hg18, hg19)

Mouse (mm9, mm10)

...

(iGenomes, UCSC, NCBI, EBI...)

Fast...

Versatile...

Documented and supported...

**Different mappers will use
different Index Formats...**

Align Data...

Aligner (Mapper)

bowtie.2

Reference Genome

Human (hg18, hg19)

Mouse (mm9, mm10)

...

(iGenomes, UCSC, NCBI, EBI...)

**Different mappers will use
different Index Formats...**

% of alignment...

Align Data...

Aligner (Mapper)

bowtie.2

Reference Genome

Human (hg18, hg19)

Mouse (mm9, mm10)

...

(iGenomes, UCSC, NCBI, EBI...)

**Different mappers will use
different Index Formats...**

% of alignment...

Low percentage?

Contaminations,
adapter/barcode trimming,
wrong reference...

Align Data...

Aligner (Mapper)

bowtie.2

Reference Genome

Human (hg18, hg19)

Mouse (mm9, mm10)

...

(iGenomes, UCSC, NCBI, EBI...)

**Different mappers will use
different Index Formats...**

% of alignment...

Low percentage?

Contaminations,
adapter/barcode trimming,
wrong reference...

Bad Library...

Aligned Data



Aligned Data

SAM ⇔ **BAM** ⇒ **BED**

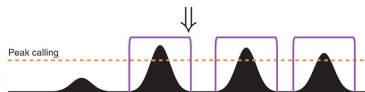
```
graph TD; A[Aligned Data] --- B[SAM ↔ BAM ⇒ BED]; B --> C[Enriched Regions];
```

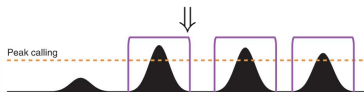
Aligned Data

SAM ↔ BAM ⇒ BED



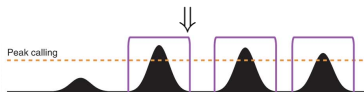
**Enriched
Regions**

Aligned Data**SAM** \Leftrightarrow **BAM** \Rightarrow **BED****Enriched
Regions**

Aligned DataSAM \Leftrightarrow BAM \Rightarrow BED**Enriched
Regions****Chr Start Stop ID Score****SET of
GENOMIC COORDINATES
(PEAKS...)**

Aligned DataSAM \Leftrightarrow BAM \Rightarrow BED**Enriched
Regions**

How Many???



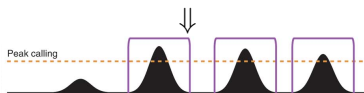
Chr Start Stop ID Score

**SET of
GENOMIC COORDINATES
(PEAKS...)**

Aligned DataSAM \Leftrightarrow BAM \Rightarrow BED**Enriched
Regions**

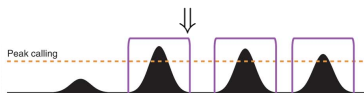
How Many???

Where???



Chr Start Stop ID Score

**SET of
GENOMIC COORDINATES
(PEAKS...)**

Aligned DataSAM \Leftrightarrow BAM \Rightarrow BED**Enriched
Regions**

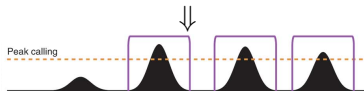
Chr Start Stop ID Score

**SET of
GENOMIC COORDINATES
(PEAKS...)**

How Many???

Where???

With Who???

Aligned DataSAM \Leftrightarrow BAM \Rightarrow BED**Enriched
Regions**

Chr Start Stop ID Score

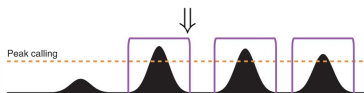
**SET of
GENOMIC COORDINATES
(PEAKS...)**

How Many???

Where???

With Who???

Function???

Aligned DataSAM \Leftrightarrow BAM \Rightarrow BED**Enriched
Regions**

Chr Start Stop ID Score

**SET of
GENOMIC COORDINATES
(PEAKS...)**

How Many???

Where???

With Who???

Function???

HOW???

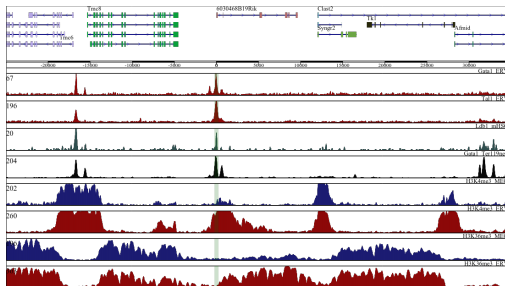
“Peak calling is a computational method used to identify areas in a genome that have been enriched with aligned reads“

“Peak calling is a computational method used to identify areas in a genome that have been enriched with aligned reads“

Challenge: Not one single 'peak shape'...

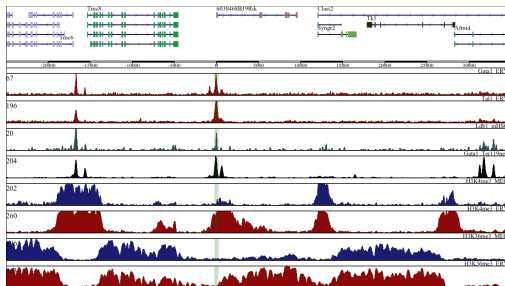
“Peak calling is a computational method used to identify areas in a genome that have been enriched with aligned reads“

Challenge: Not one single 'peak shape'...



“Peak calling is a computational method used to identify areas in a genome that have been enriched with aligned reads”

Challenge: Not one single 'peak shape'...



Number of peaks

Peak distribution

Peak height

Range of measurements

Experimental variation (IP...)

Solutions:

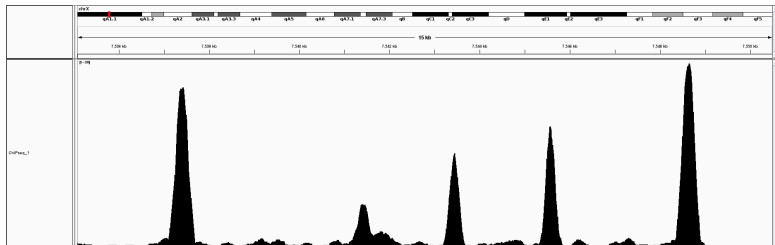
Different Peak Callers

Optimal Parameters

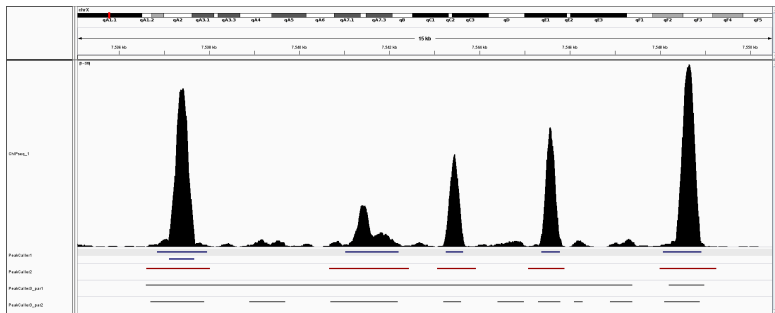
Correct Experimental Design

Good ChIP...

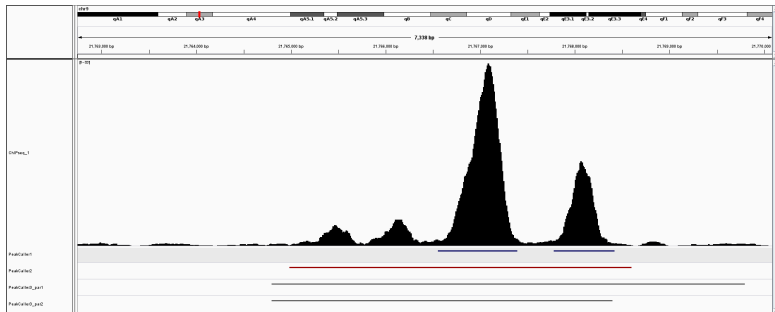
Solutions:
Different Peak Callers
Optimal Parameters
Correct Experimental Design
Good ChIP...



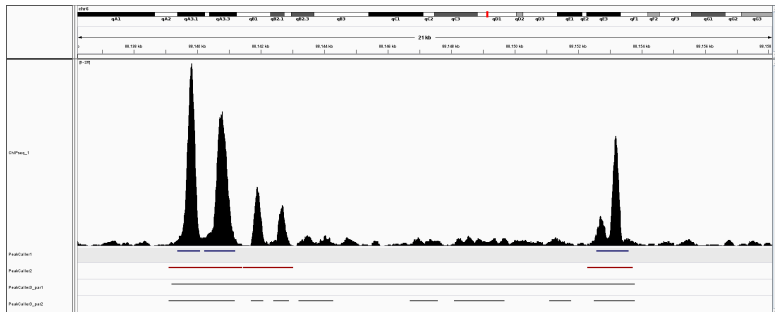
Solutions:
Different Peak Callers
Optimal Parameters
Correct Experimental Design
Good ChIP...



Solutions:
Different Peak Callers
Optimal Parameters
Correct Experimental Design
Good ChIP...



Solutions:
Different Peak Callers
Optimal Parameters
Correct Experimental Design
Good ChIP...



Solutions:

Different Peak Callers

Optimal Parameters

Correct Experimental Design

Good ChIP...

Peaks represent only a 'summary' of a ChIPseq experiment...

Solutions:

Different Peak Callers

Optimal Parameters

Correct Experimental Design

Good ChIP...

Peaks represent only a 'summary' of a ChIPseq experiment...

Useful metrics:**FRiP:** Fraction of Reads in Peaks ($>1\%$)**IDR:** Irreproducible Discovery Rate**Significance:** pValue, ChIPtoINPUT ratio, # of reads...

Solutions:

Different Peak Callers

Optimal Parameters

Correct Experimental Design

Good ChIP...

Peaks represent only a 'summary' of a ChIPseq experiment...

Empirical:**Known Motif:** Enriched in peaks, Central Position**Distribution:** Near TSS, Within GeneBody, Enhancers**Target Genes:** Gene expression changes, Functional Annotation

Solutions:

Different Peak Callers

Optimal Parameters

Correct Experimental Design

Good ChIP...

Peaks represent only a 'summary' of a ChIPseq experiment...

Empirical:

Known Motif: Enriched in peaks, Central Position

Distribution: Near TSS, Within GeneBody, Enhancers

Target Genes: Gene expression changes, Functional Annotation

Browser inspection and previously known sites

TFs: MACS

HMs: SICER

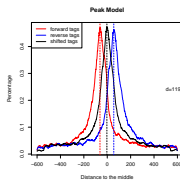
TFs: MACS

HMs: SICER

Running MACS: `macs14 -t ChIP -c Input -g mm -n Name`

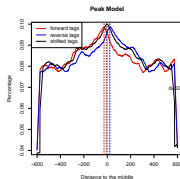
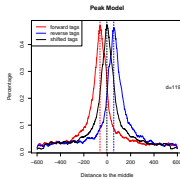
TFs: MACS
HMs: SICER

Running MACS: `macs14 -t ChIP -c Input -g mm -n Name`



TFs: MACS
HMs: SICER

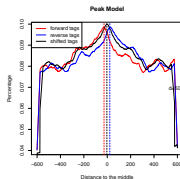
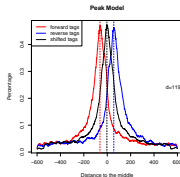
Running MACS: `macs14 -t ChIP -c Input -g mm -n Name`



TFs: MACS

HMs: SICER

Running MACS: `macs14 -t ChIP -c Input -g mm -n Name`



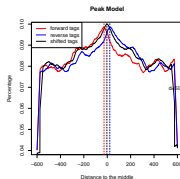
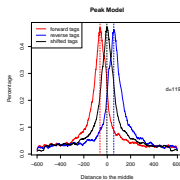
```
# This file is generated by MACS version 1.4.2 20120305
# ARGUMENTS LIST:
# name = Gata1_ERYneg
# format = AUTO
# CHIP-seq file = Gata1_ERYneg.bam
# control file = Inputs/ERYneg_INPUT.bam
# effective genome size = 1.87e+09
# band width = 300
# model fold = 10
# pvalue cutoff = 1.00e-05
# Large dataset will be scaled towards smaller dataset.
# Range for calculating regional lambda is: 1000 bps and 10000 bps

# tag size is determined as 51 bps
# total tags in treatment: 22404993
# tags after filtering in treatment: 19156302
# maximum duplicate tags at the same position in treatment = 1
# Redundant rate in treatment: 0.14
# total tags in control: 10891587
# tags after filtering in control: 10220192
# maximum duplicate tags at the same position in control = 1
# Redundant rate in control: 0.06
# d = 188
```

| chr | start | end | length | summit | tags | $-10 \cdot \log_{10}(\text{pvalue})$ | fold_enrichment | FDR(%) |
|------|---------|---------|--------|--------|------|--------------------------------------|-----------------|--------|
| chr1 | 3042341 | 3043008 | 668 | 472 | 38 | 87.161 | 30.2 | 2.67 |
| chr1 | 3049560 | 3050013 | 454 | 267 | 47 | 173.5 | 41.73 | 0.57 |
| chr1 | 3435094 | 3436610 | 1027 | 551 | 146 | 414.38 | 61.7 | 0.33 |

TFs: MACS
HMs: SICER

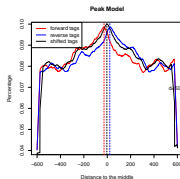
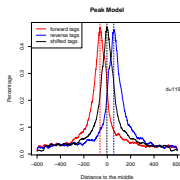
Running MACS: macs14 -t ChIP -c Input -g mm -n Name



| Chromosome | Start | Stop | ID | Score |
|------------|-----------|-----------|-----------------|--------|
| chr1 | 3042340 | 3043008 | MACS_peak_1 | 87.16 |
| chr1 | 3049559 | 3050013 | MACS_peak_2 | 173.5 |
| chr1 | 3435583 | 3436610 | MACS_peak_3 | 414.38 |
| . | . | . | . | . |
| . | . | . | . | . |
| chrX | 165596178 | 165596764 | MACS_peak_26347 | 271.4 |
| chrX | 165658009 | 165658680 | MACS_peak_26348 | 253.88 |

TFs: MACS
HMs: SICER

Running MACS: `macs14 -t ChIP -c Input -g mm -n Name`



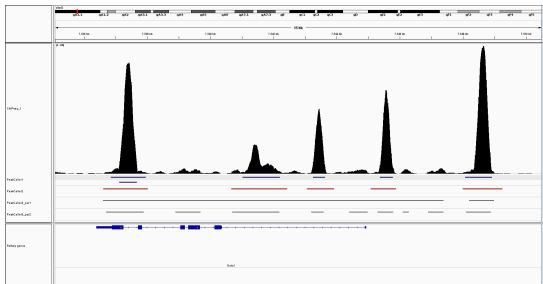
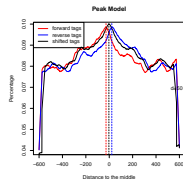
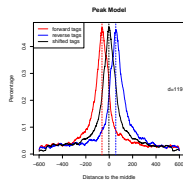
WIG file: Normalized Read Density text format



bigWig file: Binary Read Density format

TFs: MACS
 HMs: SICER

Running MACS: `macs14 -t ChIP -c Input -g mm -n Name`



General Guidelines:

Biological Replicates (2x)

INPUT (noIP sample)

20-25 mln reads per sample

Optimize ChIP conditions before library preparation

Minimize experimental variation

Look at the data

General Guidelines:

Biological Replicates (2x)

INPUT (noIP sample)

20-25 mln reads per sample

Optimize ChIP conditions before library preparation

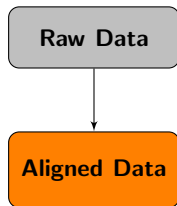
Minimize experimental variation

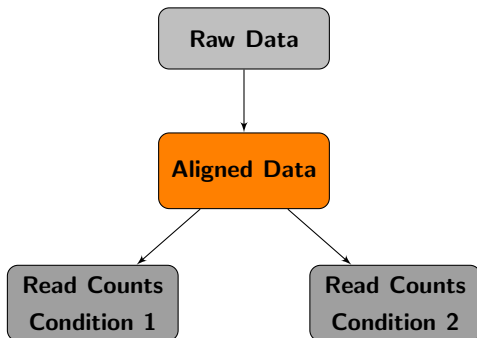
Look at the data

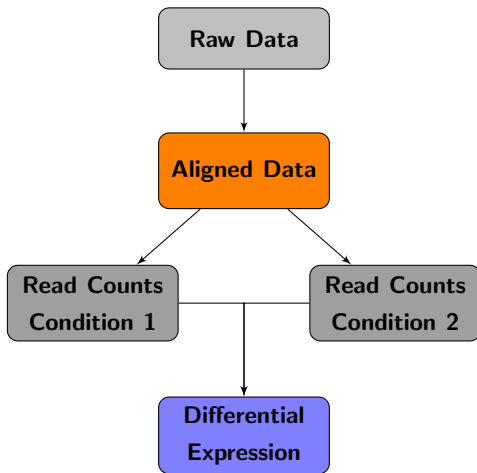
'ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia'
Landt et al, Genome Res, 2012

'Practical Guidelines for the Comprehensive Analysis of ChIP-seq Data'
Bailey et al, PLoS Comp Biol, 2013

Raw Data







Aligner (Mapper)

Aligner (Mapper)

Reference Genome

Human (hg18, hg19)

Mouse (mm9, mm10)

...

Spliced alignment

**Build reference based on
known transcripts**

Aligner (Mapper)

STAR

Reference Genome

Human (hg18, hg19)

Mouse (mm9, mm10)

...

Spliced alignment

**Build reference based on
known transcripts**

Aligner (Mapper)**STAR****Reference Genome**

Human (hg18, hg19)

Mouse (mm9, mm10)

...

Spliced alignment**Build reference based on
known transcripts**

VERY Fast...

(45 million paired reads
per hour per processor)

Aligner (Mapper)**STAR****Reference Genome**

Human (hg18, hg19)

Mouse (mm9, mm10)

...

Spliced alignment**Build reference based on
known transcripts****VERY Fast...**
(45 million paired reads
per hour per processor)**TopHat2**

Aligner (Mapper)**STAR****Reference Genome**

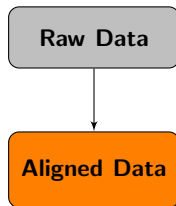
Human (hg18, hg19)

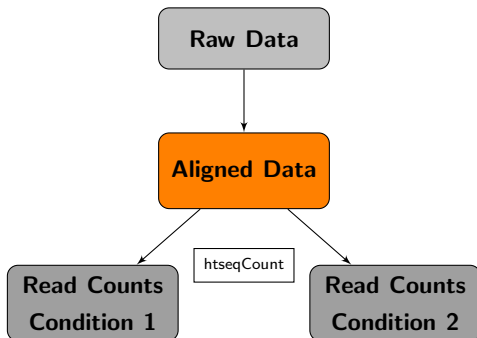
Mouse (mm9, mm10)

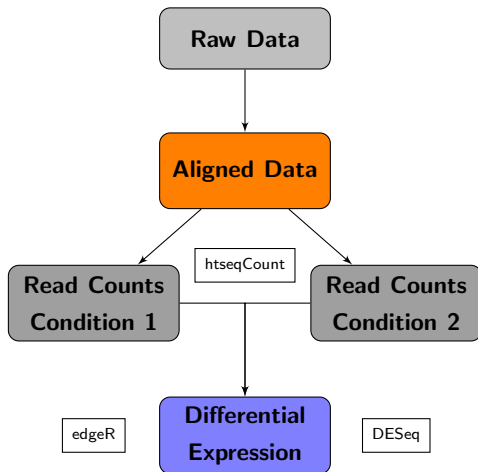
...

Spliced alignment**Build reference based on
known transcripts**VERY Fast...
(45 million paired reads
per hour per processor)**TopHat2**

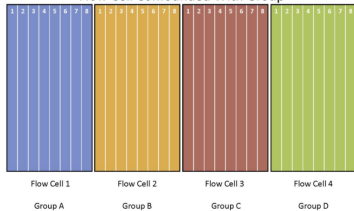
If you want to use CuffLinks...







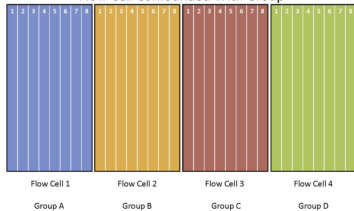
Raw Data

Raw Data**Differential Expression Across Groups**
Flow Cell Confounded With Group

Raw Data

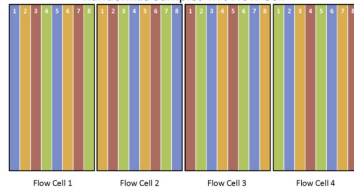
Differential Expression Across Groups

Flow Cell Confounded With Group



Differential Expression Across Groups

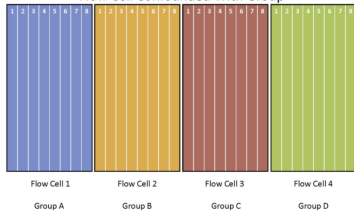
Randomize Samples wrt Flow Cell



Raw Data

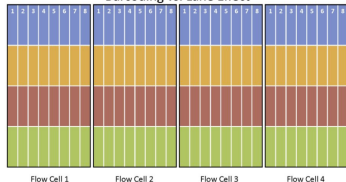
Differential Expression Across Groups

Flow Cell Confounded With Group



Differential Expression Across Groups

Barcoding vs. Lane Effect



Local

Local

Online

Local**Online**

IGV (Integrative Genomics Viewer)

Local

IGV (Integrative Genomics Viewer)

Online

UCSC Genome Browser

Local

IGV (Integrative Genomics Viewer)

Cross Platform (Win, Mac, Linux)

Different file formats

Fast

Online

UCSC Genome Browser

Local

IGV (Integrative Genomics Viewer)

Cross Platform (Win, Mac, Linux)

Different file formats

Fast

Online

UCSC Genome Browser

Web based

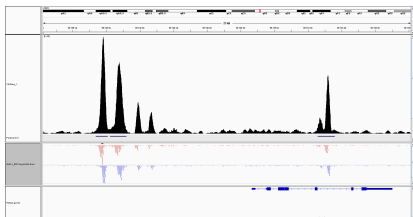
Slow

Comprehensive database

Local

IGV (Integrative Genomics Viewer)

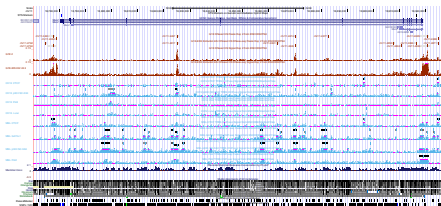
Cross Platform (Win, Mac, Linux)
 Different file formats
 Fast



Online

UCSC Genome Browser

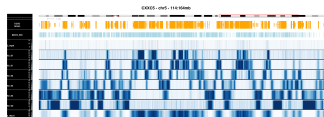
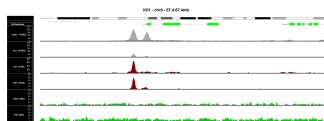
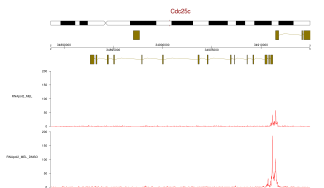
Web based
 Slow
 Comprehensive database



Local

Online

GenomeGraphs (R package)

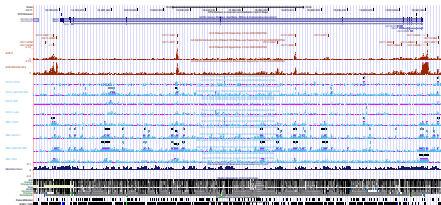


UCSC Genome Browser

Web based

Slow

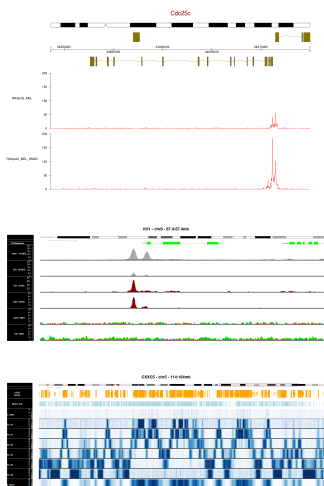
Comprehensive database



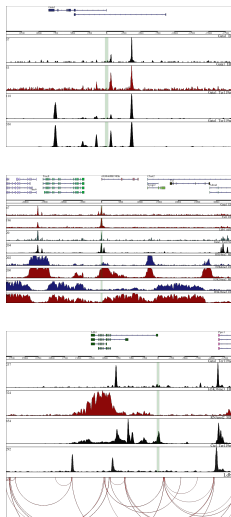
Local

Online

GenomeGraphs (R package)



Ariadne (GeneViewer)



GALAXY server

GALAXY server



Online collection of widely used
bioinformatics tools
(e.g. FastQC, bowtie.2, CuffLinks,
MACS, SICER ...)

GALAXY server



Online collection of widely used
bioinformatics tools
(e.g. FastQC, bowtie.2, CuffLinks,
MACS, SICER ...)

Learn Galaxy...



GALAXY server



Online collection of widely used
bioinformatics tools
(e.g. FastQC, bowtie.2, CuffLinks,
MACS, SICER ...)

Learn Galaxy...



<https://galaxyproject.org/>

GALAXY server



Online collection of widely used
bioinformatics tools
(e.g. FastQC, bowtie.2, CuffLinks,
MACS, SICER ...)

Learn Galaxy...



<https://galaxyproject.org/>

Chipster

GALAXY server



Online collection of widely used
bioinformatics tools
(e.g. FastQC, bowtie.2, CuffLinks,
MACS, SICER ...)

Learn Galaxy...

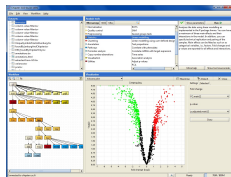
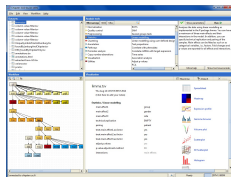


<https://galaxyproject.org/>

Chipster



User-friendly analysis software for
high-throughput data (GUI)



Genomic Regions Enrichment of Annotations Tool (GREAT)

Genomic Regions Enrichment of Annotations Tool (GREAT)



Functional annotation of
cis-regulatory regions

Genomic Regions Enrichment of Annotations Tool (GREAT)



Functional annotation of
cis-regulatory regions

Input: BED file

Genomic Regions Enrichment of Annotations Tool (GREAT)



Functional annotation of
cis-regulatory regions

Input: BED file

Output:

List of potential target genes
Peak distribution plots
Geneset enrichment analysis
(GO, Pathways, Phenotypes ...)

Genomic Regions Enrichment of Annotations Tool (GREAT)



Functional annotation of
cis-regulatory regions

Input: BED file

Output:

List of potential target genes

Peak distribution plots

Geneset enrichment analysis

(GO, Pathways, Phenotypes ...)

[http:](http://bejerano.stanford.edu/great/public/html/index.php)

[//bejerano.stanford.edu/great/public/html/index.php](http://bejerano.stanford.edu/great/public/html/index.php)

Genomic Regions Enrichment of Annotations Tool (GREAT)

The logo for GREAT, featuring the word "GREAT" in a stylized font with a blue triangle above the letter "A".

Functional annotation of
cis-regulatory regions

Input: BED file

Output:

List of potential target genes

Peak distribution plots

Geneset enrichment analysis

(GO, Pathways, Phenotypes ...)

http:

[//bejerano.stanford.edu/great/public/html/index.php](http://bejerano.stanford.edu/great/public/html/index.php)

MEME-ChIP

Genomic Regions Enrichment of Annotations Tool (GREAT)



Functional annotation of
cis-regulatory regions

Input: BED file

Output:

List of potential target genes

Peak distribution plots

Geneset enrichment analysis

(GO, Pathways, Phenotypes ...)

[http:](http://bejerano.stanford.edu/great/public/html/index.php)

[//bejerano.stanford.edu/great/public/html/index.php](http://bejerano.stanford.edu/great/public/html/index.php)

MEME-CHIP



Perform motif discovery, motif
enrichment analysis and clustering

Genomic Regions Enrichment of Annotations Tool (GREAT)



Functional annotation of
cis-regulatory regions

Input: BED file

Output:

List of potential target genes

Peak distribution plots

Geneset enrichment analysis

(GO, Pathways, Phenotypes ...)

[http:](http://bejerano.stanford.edu/great/public/html/index.php)

[//bejerano.stanford.edu/great/public/html/index.php](http://bejerano.stanford.edu/great/public/html/index.php)

MEME-CHIP



Perform motif discovery, motif
enrichment analysis and clustering

Input: FASTA file

Genomic Regions Enrichment of Annotations Tool (GREAT)



Functional annotation of
cis-regulatory regions

Input: BED file

Output:

List of potential target genes
Peak distribution plots
Geneset enrichment analysis
(GO, Pathways, Phenotypes ...)

[http:](http://bejerano.stanford.edu/great/public/html/index.php)

[//bejerano.stanford.edu/great/public/html/index.php](http://bejerano.stanford.edu/great/public/html/index.php)

MEME-CHIP

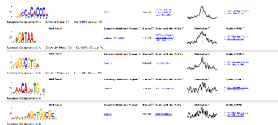


Perform motif discovery, motif
enrichment analysis and clustering

Input: FASTA file

Output:

Enriched motifs
Motif distribution plots



Genomic Regions Enrichment of Annotations Tool (GREAT)



Functional annotation of
cis-regulatory regions

Input: BED file

Output:

List of potential target genes
Peak distribution plots
Geneset enrichment analysis
(GO, Pathways, Phenotypes ...)

[http:](http://bejerano.stanford.edu/great/public/html/index.php)

[//bejerano.stanford.edu/great/public/html/index.php](http://bejerano.stanford.edu/great/public/html/index.php)

MEME-CHIP

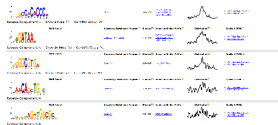


Perform motif discovery, motif
enrichment analysis and clustering

Input: FASTA file

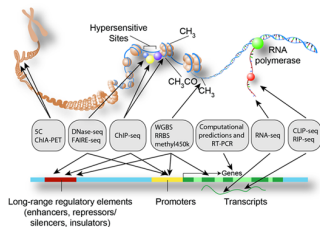
Output:

Enriched motifs
Motif distribution plots



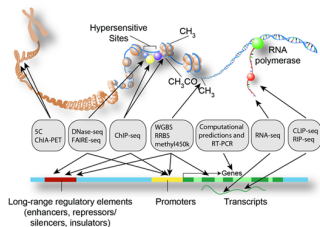
<http://meme-suite.org/tools/meme-chip>

Encyclopedia of DNA Elements



Aim: build a comprehensive parts list of functional elements in the human genome

Encyclopedia of DNA Elements



Aim: build a comprehensive parts list of functional elements in the human genome

New website is extremely easy to use!

Encyclopedia of DNA Elements

ENCODE

Data ▾

Methods ▾

About ▾

Help ▾

Search...



Encyclopedia of DNA Elements

Assay

| | |
|--|------|
| ChIP-seq | 4540 |
| RNA-seq | 1080 |
| DNase-seq | 654 |
| shRNA knockdown followed by RNA-seq | 333 |
| transcription profiling by array assay | 293 |

[+ See more...](#)

Experiment status

| | |
|----------|------|
| released | 8538 |
| revoked | 25 |

Genome assembly (visualization)

| | |
|------|------|
| hg19 | 3222 |
| mm9 | 570 |
| mm10 | 551 |
| dm3 | 108 |

Organism

| | |
|--------------------------------|------|
| <i>Homo sapiens</i> | 6971 |
| <i>Mus musculus</i> | 1282 |
| <i>Drosophila melanogaster</i> | 197 |

Target of assay

| | |
|----------------------|------|
| histone | 2646 |
| histone modification | 2610 |
| transcription factor | 1221 |
| control | 844 |
| RNA binding protein | 567 |

[+ See more...](#)

Biosample type

| | |
|-------------------------------|------|
| immortalized cell line | 3087 |
| tissue | 2413 |
| primary cell | 1761 |
| stem cell | 619 |
| in vitro differentiated cells | 425 |

[+ See more...](#)

ENCODE [Data](#) [Methods](#) [About](#) [Help](#)

Search

Encyclopedia of DNA Elements

Assay

CHIP-seq 4540

RNA-seq 1080

DNase-seq 654

shRNA knockdown followed by RNA-seq 333

transcription profiling by array assay 293

[+ See more...](#)

Experiment status

released 8538

revoked 25

Genome assembly (visualization)

hg19 3222

mm9 570

mm10 551

dm3 108

Organism

Homo sapiens 6971

Mus musculus 1282

Drosophila melanogaster 197

Target of assay

histone 2646

histone modification 2610

transcription factor 1221

control 844

RNA binding protein 567

[+ See more...](#)

Biosample type

immortalized cell line 3087

tissue 2413

primary cell 1761

stem cell 619

in vitro differentiated cells 425

[+ See more...](#)

ENCODE [Data](#) [Methods](#) [About](#) [Help](#)

Search...

Assay

CHIP-seq 6

Experiment status

released 6

Genome assembly (visualization)

hg19 65

mm9 6

Organism

Mus musculus 6

Target of assay

histone 13

histone modification 13

transcription factor 6

control 4

Biosample type

immortalized cell line 111

tissue 33

in vitro differentiated cells 16

stem cell 11

primary cell 6

Organ

bone element 2

Life stage

adult 4

embryonic 2

Available data

bam 6

bed broadPeak 6

bigBed broadPeak 6

bigWig 6

fastq 6

Showing 6 of 6 experiments [Viewstate \(F\)](#) [Download](#)

ChIP-seq of bone marrow macrophage (*Mus musculus*, adult 8 week) Experiment ENCSR000CFJ released

Target: CTCF
Lab: Bing Ren, UCSD
Project: ENCODE

ChIP-seq of bone marrow macrophage (*Mus musculus*, adult 8 week) Experiment ENCSR000CFK released

Target: POLR2A
Lab: Bing Ren, UCSD
Project: ENCODE

ChIP-seq of embryonic fibroblast (*Mus musculus*, adult) Experiment ENCSR000CBX released

Target: POLR2A
Lab: Bing Ren, UCSD
Project: ENCODE

ChIP-seq of embryonic fibroblast (*Mus musculus*, adult) Experiment ENCSR000CBW released

Target: CTCF
Lab: Bing Ren, UCSD
Project: ENCODE

ChIP-seq of erythroblast (*Mus musculus*, embryonic 14.5 day) Experiment ENCSR000DB released

Target: TAL1
Lab: Ross Hardison, PennState
Project: ENCODE

ChIP-seq of erythroblast (*Mus musculus*, embryonic 14.5 day) Experiment ENCSR000DL released

Target: GATA1
Lab: Ross Hardison, PennState
Project: ENCODE

Encyclopedia of DNA Elements

Assay

| | |
|--|------|
| ChIP-seq | 4540 |
| RNA-seq | 1080 |
| DNase-seq | 654 |
| shRNA knockdown followed by RNA-seq | 333 |
| transcription profiling by array assay | 293 |

[+ See more...](#)

Experiment status

| | |
|----------|------|
| released | 8538 |
| revoked | 25 |

Genome assembly (visualization)

| | |
|------|------|
| hg19 | 3222 |
| mm9 | 570 |
| mm10 | 551 |
| dm3 | 108 |

Organism

| | |
|-------------------------|------|
| Homo sapiens | 6971 |
| Mus musculus | 1282 |
| Drosophila melanogaster | 197 |

Target of assay

| | |
|----------------------|------|
| histone | 2646 |
| histone modification | 2610 |
| transcription factor | 1221 |
| control | 844 |
| RNA binding protein | 567 |

[+ See more...](#)

Biosample type

| | |
|-------------------------------|------|
| immortalized cell line | 3087 |
| tissue | 2413 |
| primary cell | 1761 |
| stem cell | 619 |
| in vitro differentiated cells | 425 |

[+ See more...](#)

| ENCODE | | | | | | | | | | |
|--|-----------|----------------------|---------------------|-------------|--------------|------------|------------------|--------------------------|------------|--|
| Data - Methods - About - Help - | | | | | | | | | | |
| Search... | | | | | | | | | | |
| Raw data | | | | | | | | | | |
| Accession | File type | Biological replicate | Technical replicate | Read length | Run type | Paired end | Mapping assembly | Lab | Date added | |
| ENCFF001MMN Download 1.35 GB | testq | 1 | 1 | 36 nt | single-ended | | | Ross Hardison, PennState | 2011-11-05 | |
| ENCFF001MMR Download 3.97 GB | testq | 2 | 1 | 36 nt | single-ended | | | Ross Hardison, PennState | 2011-11-05 | |

Encyclopedia of DNA Elements

ENCODE
Data ▾ Methods ▾ About ▾ Help ▾

Search...

Processed data

| Accession | File type | Output type | Biological replicate(s) | Technical replicate | Mapping assembly | Genome annotation | Lab | Date added |
|---|---------------------|-------------|-------------------------|---------------------|------------------|-------------------|--------------------------|------------|
| ENCF001MAM Download 4.05 GB | bigWig | signal | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAG Download 1.68 GB | bam | alignments | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAJ Download 4.55 GB | bam | alignments | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAK Download 284 kB | bigBed broadPeak | peaks | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAL Download 244 kB | bigBed broadPeak | peaks | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAM Download 266 kB | bigBed broadPeak | peaks | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAO Download 4.89 GB | bigWig | signal | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MMP Download 5.62 GB | bigWig | signal | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEN Download 126 kB | bed broadPeak | peaks | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEO Download 133 kB | bed broadPeak | peaks | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEP Download 99.7 kB | bed broadPeak | peaks | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |

[Contact](#)
[Terms of Use](#)
[Submitter sign-in](#)

©2015 Stanford University

Assay

| | |
|--|------|
| ChIP-seq | 4540 |
| RNA-seq | 1060 |
| DNase-seq | 654 |
| shRNA knockdown followed by RNA-seq | 333 |
| transcription profiling by array assay | 293 |

[+ See more...](#)

Experiment status

| | |
|----------|------|
| released | 8538 |
| revoked | 25 |

Genome assembly (visualization)

| | |
|------|------|
| hg19 | 3222 |
| mm9 | 570 |
| mm10 | 551 |
| dm3 | 108 |

Organism

| | |
|-------------------------|------|
| Homo sapiens | 6971 |
| Mus musculus | 1282 |
| Drosophila melanogaster | 197 |

Target of assay

| | |
|----------------------|------|
| histone | 2646 |
| histone modification | 2610 |
| transcription factor | 1221 |
| control | 844 |
| RNA binding protein | 567 |

[+ See more...](#)

Biosample type

| | |
|-------------------------------|------|
| immortalized cell line | 3087 |
| tissue | 2413 |
| primary cell | 1761 |
| stem cell | 619 |
| in vitro differentiated cells | 425 |

[+ See more...](#)

Encyclopedia of DNA Elements

ENCODE
Data ▾ Methods ▾ About ▾ Help ▾

Search... Q

Processed data

| Accession | File type | Output type | Biological replicate(s) | Technical replicate | Mapping assembly | Genome annotation | Lab | Date added |
|--|---------------------|-------------|-------------------------|---------------------|------------------|-------------------|--------------------------|------------|
| ENCF001MM1M Download 4.05 GB | bigWig | signal | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MM1G Download 1.68 GB | bam | alignments | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MM1J Download 4.55 GB | bam | alignments | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MM1K Download 284 kB | bigBed broadPeak | peaks | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MM1L Download 244 kB | bigBed broadPeak | peaks | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MM1E Download 266 kB | bigBed broadPeak | peaks | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MM1O Download 4.89 GB | bigWig | signal | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MM1P Download 5.62 GB | bigWig | signal | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEN Download 126 kB | bed broadPeak | peaks | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEO Download 133 kB | bed broadPeak | peaks | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEP Download 99.7 kB | bed broadPeak | peaks | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |

[Contact](#)
[Terms of Use](#)
[Submitter sign-in](#)

©2015 Stanford University

Assay

| | |
|--|------|
| ChIP-seq | 4540 |
| RNA-seq | 1060 |
| DNase-seq | 654 |
| shRNA knockdown followed by RNA-seq | 333 |
| transcription profiling by array assay | 293 |

[+ See more...](#)

Experiment status

| | |
|----------|------|
| released | 8538 |
| revoked | 25 |

Genome assembly (visualization)

| | |
|------|------|
| hg19 | 3222 |
| mm9 | 570 |
| mm10 | 551 |
| dm3 | 108 |

Organism

| | |
|-------------------------|------|
| Homo sapiens | 6971 |
| Mus musculus | 1282 |
| Drosophila melanogaster | 197 |

Target of assay

| | |
|----------------------|------|
| histone | 2646 |
| histone modification | 2610 |
| transcription factor | 1221 |
| control | 844 |
| RNA binding protein | 567 |

[+ See more...](#)

Biosample type

| | |
|-------------------------------|------|
| immortalized cell line | 3087 |
| tissue | 2413 |
| primary cell | 1761 |
| stem cell | 619 |
| in vitro differentiated cells | 425 |

[+ See more...](#)

Encyclopedia of DNA Elements

ENCODE
Data ▾ Methods ▾ About ▾ Help ▾

Search... Q

Processed data

| Accession | File type | Output type | Biological replicate(s) | Technical replicate | Mapping assembly | Genome annotation | Lab | Date added |
|---|---------------------|-------------|-------------------------|---------------------|------------------|-------------------|--------------------------|------------|
| ENCF001MAM Download 4.05 GB | bigWig | signal | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAG Download 1.68 GB | bam | alignments | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAJ Download 4.55 GB | bam | alignments | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAK Download 284 kB | bigBed broadPeak | peaks | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAL Download 244 kB | bigBed broadPeak | peaks | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAM Download 266 kB | bigBed broadPeak | peaks | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAO Download 4.89 GB | bigWig | signal | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MMP Download 5.62 GB | bigWig | signal | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEN Download 126 kB | bed broadPeak | peaks | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEO Download 133 kB | bed broadPeak | peaks | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEP Download 99.7 kB | bed broadPeak | peaks | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |

[Contact](#)
[Terms of Use](#)
[Submitter sign-in](#)

©2015 Stanford University

Assay

| | |
|--|------|
| ChIP-seq | 4540 |
| RNA-seq | 1060 |
| DNase-seq | 654 |
| shRNA knockdown followed by RNA-seq | 333 |
| transcription profiling by array assay | 293 |

[+ See more...](#)

Experiment status

| | |
|----------|------|
| released | 8538 |
| revoked | 25 |

Genome assembly (visualization)

| | |
|------|------|
| hg19 | 3222 |
| mm9 | 570 |
| mm10 | 551 |
| dm3 | 108 |

Organism

| | |
|-------------------------|------|
| Homo sapiens | 6971 |
| Mus musculus | 1282 |
| Drosophila melanogaster | 197 |

Target of assay

| | |
|------------------------------|------|
| histone | 2646 |
| histone modification | 2610 |
| transcription factor control | 1221 |
| control | 844 |
| RNA binding protein | 567 |

[+ See more...](#)

Biosample type

| | |
|-------------------------------|------|
| immortalized cell line | 3087 |
| tissue | 2413 |
| primary cell | 1761 |
| stem cell | 619 |
| in vitro differentiated cells | 425 |

[+ See more...](#)

Encyclopedia of DNA Elements

ENCODE
Data ▾ Methods ▾ About ▾ Help ▾

Search... Q

Processed data

| Accession | File type | Output type | Biological replicate(s) | Technical replicate | Mapping assembly | Genome annotation | Lab | Date added |
|---|---------------------|-------------|-------------------------|---------------------|------------------|-------------------|--------------------------|------------|
| ENCF001MAM Download 4.05 GB | bigWig | signal | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAG Download 1.68 GB | bam | alignments | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAJ Download 4.55 GB | bam | alignments | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAK Download 284 kB | bigBed broadPeak | peaks | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAL Download 244 kB | bigBed broadPeak | peaks | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAM Download 266 kB | bigBed broadPeak | peaks | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MAO Download 4.89 GB | bigWig | signal | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001MMP Download 5.62 GB | bigWig | signal | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEN Download 126 kB | bed broadPeak | peaks | | | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEO Download 133 kB | bed broadPeak | peaks | 1 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |
| ENCF001YEP Download 99.7 kB | bed broadPeak | peaks | 2 | 1 | mm9 | | Ross Hardison, PennState | 2011-11-05 |

[Contact](#)
[Terms of Use](#)
[Submitter sign-in](#)

©2015 Stanford University

Assay

| | |
|--|------|
| ChIP-seq | 4540 |
| RNA-seq | 1060 |
| DNase-seq | 654 |
| shRNA knockdown followed by RNA-seq | 333 |
| transcription profiling by array assay | 293 |

[+ See more...](#)

Experiment status

| | |
|----------|------|
| released | 8538 |
| revoked | 25 |

Genome assembly (visualization)

| | |
|------|------|
| hg19 | 3222 |
| mm9 | 570 |
| mm10 | 551 |
| dm3 | 108 |

Organism

| | |
|-------------------------|------|
| Homo sapiens | 6971 |
| Mus musculus | 1282 |
| Drosophila melanogaster | 197 |

Target of assay

| | |
|----------------------|------|
| histone | 2646 |
| histone modification | 2610 |
| transcription factor | 1221 |
| control | 844 |
| RNA binding protein | 567 |

[+ See more...](#)

Biosample type

| | |
|-------------------------------|------|
| immortalized cell line | 3087 |
| tissue | 2413 |
| primary cell | 1761 |
| stem cell | 619 |
| in vitro differentiated cells | 425 |

[+ See more...](#)

Gene Expression Omnibus



Gene Expression Omnibus



GEO is a public functional genomics data repository

Gene Expression Omnibus



GEO is a public functional genomics data repository

All NGS datasets have to be deposited to GEO prior to publication

Gene Expression Omnibus



GEO is a public functional genomics data repository

All NGS datasets have to be deposited to GEO prior to publication

All datasets included in GEO are publicly available for academic use

Gene Expression Omnibus



GEO is a public functional genomics data repository

All NGS datasets have to be deposited to GEO prior to publication

All datasets included in GEO are publicly available for academic use

The search function is pretty close to awful

Gene Expression Omnibus



GEO is a public functional genomics data repository

All NGS datasets have to be deposited to GEO prior to publication

All datasets included in GEO are publicly available for academic use

The search function is pretty close to awful

Best way to query the GEO database is by Accession Number
(declared in the manuscript, GSE30142, GSM746581)

Gene Expression Omnibus



GEO is a public functional genomics data repository

All NGS datasets have to be deposited to GEO prior to publication

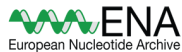
All datasets included in GEO are publicly available for academic use

The search function is pretty close to awful

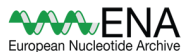
Best way to query the GEO database is by Accession Number
(declared in the manuscript, GSE30142, GSM746581)

FTP access to SRA experiment (Sequence Read Archive)

European Nucleotide Archive



European Nucleotide Archive



ENA provides a comprehensive record of the world's nucleotide sequencing information

European Nucleotide Archive



ENA provides a comprehensive record of the world's nucleotide sequencing information

It covers raw sequencing data, sequence assembly information and functional annotation

European Nucleotide Archive



ENA provides a comprehensive record of the world's nucleotide sequencing information

It covers raw sequencing data, sequence assembly information and functional annotation

The search function is pretty close to awful

European Nucleotide Archive



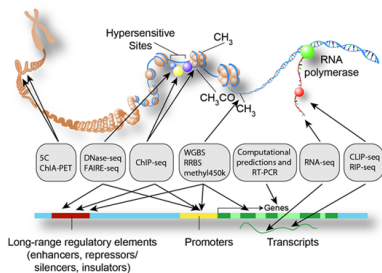
ENA provides a comprehensive record of the world's nucleotide sequencing information

It covers raw sequencing data, sequence assembly information and functional annotation

The search function is pretty close to awful

Large overlap with GEO database

Comprehensive database of genomic features



Systemic understanding
of biological processes

Thank you...