

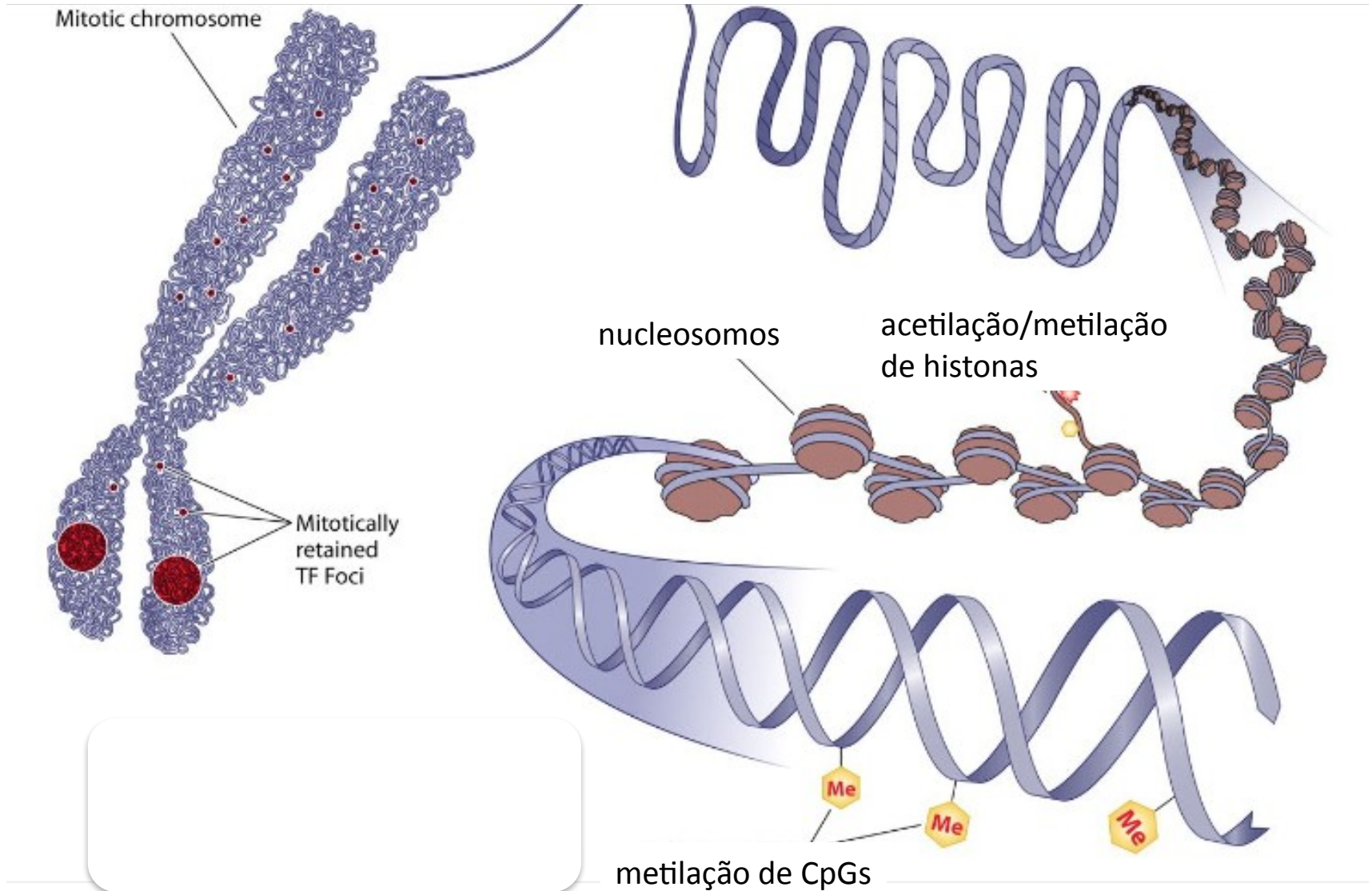
Biologia Molecular Computacional
IBI5035/QBQ2507 - 2023

Análise global de elementos regulatórios da expressão gênica

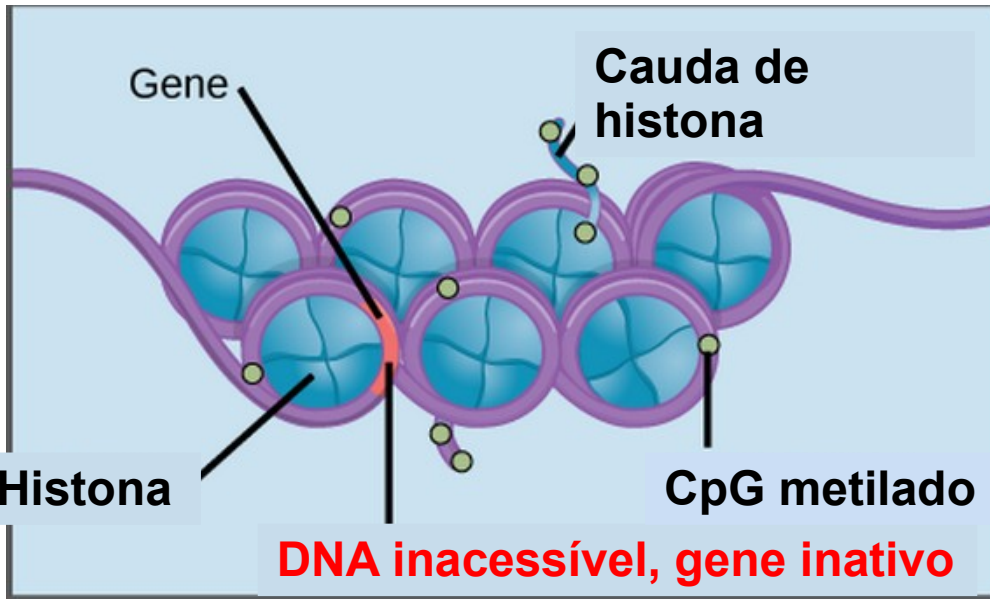
análise de dados de CHIP-Seq

Eduardo Moraes Rego Reis
Instituto de Química - USP

Modificações de histonas e metilação do DNA são mecanismos importantes para a **regulação** da expressão gênica



Modificações covalentes reversíveis em proteínas histonas definem o estados de ativação ou repressão gênica



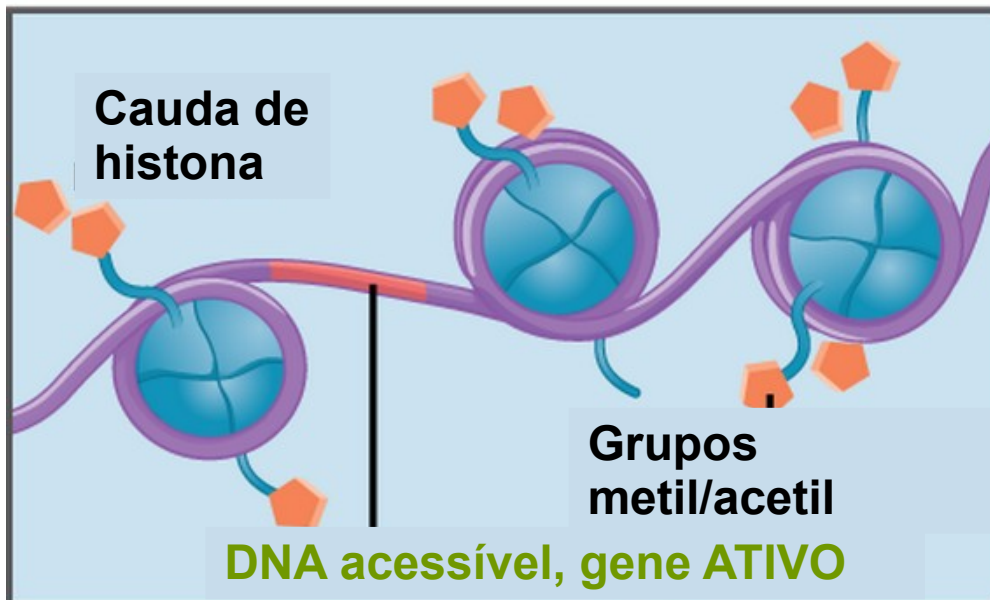
**Histonas deacetiladas/demetiladas,
CpG metilados**



Nucleossomos altamente empacotados.



A maquinaria de transcrição não consegue acessar o DNA.



**Histonas acetiladas/metiladas,
CpGs demetilados**

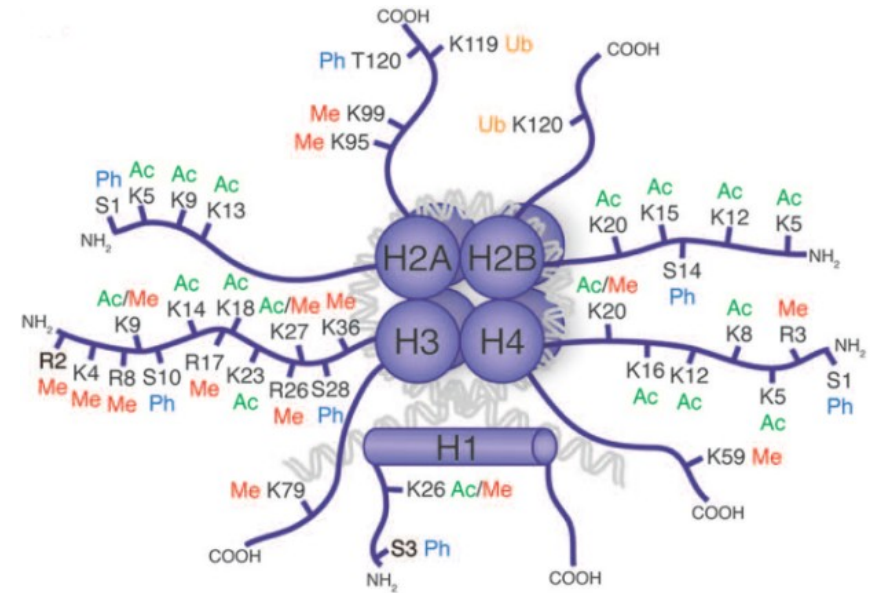


Diminuição do empacotamento dos nucleossomos.



DNA se torna mais acessível a ligação de fatores de transcrição ativadores

Modificações de histonas controlam a compactação do DNA e a expressão gênica



		Papel na transcrição gênica	Resíduos modificados na histona
Acetilação	COCH ₃	ativação	H3 (K9, K14, K18, K56) H4 (K5, K8, K12, K16) H2A/H2B (K6, K7, K16, K17)
Fosforilação	PO ₃	ativação	H3 (S10)
Metilação	CH ₃	ativação repressão	H3 (K4, K36, K79) H4 (K20)
Ubiquitinação	156 aa	ativação repressão	H2B (K123) H2A (K119)
Sumoilação	101 aa	repressão	H3 (?) H4 (K5, K8, K12, K16) H2A (K126) H2B (K6, K7, K16, K17)

- Combinações de diferentes marcas tem efeito combinatorial na cromatina
- Código das histonas ainda não é totalmente conhecido

Promotor



H3K4me3
H3K27ac
H3K9ac

H3K27me3

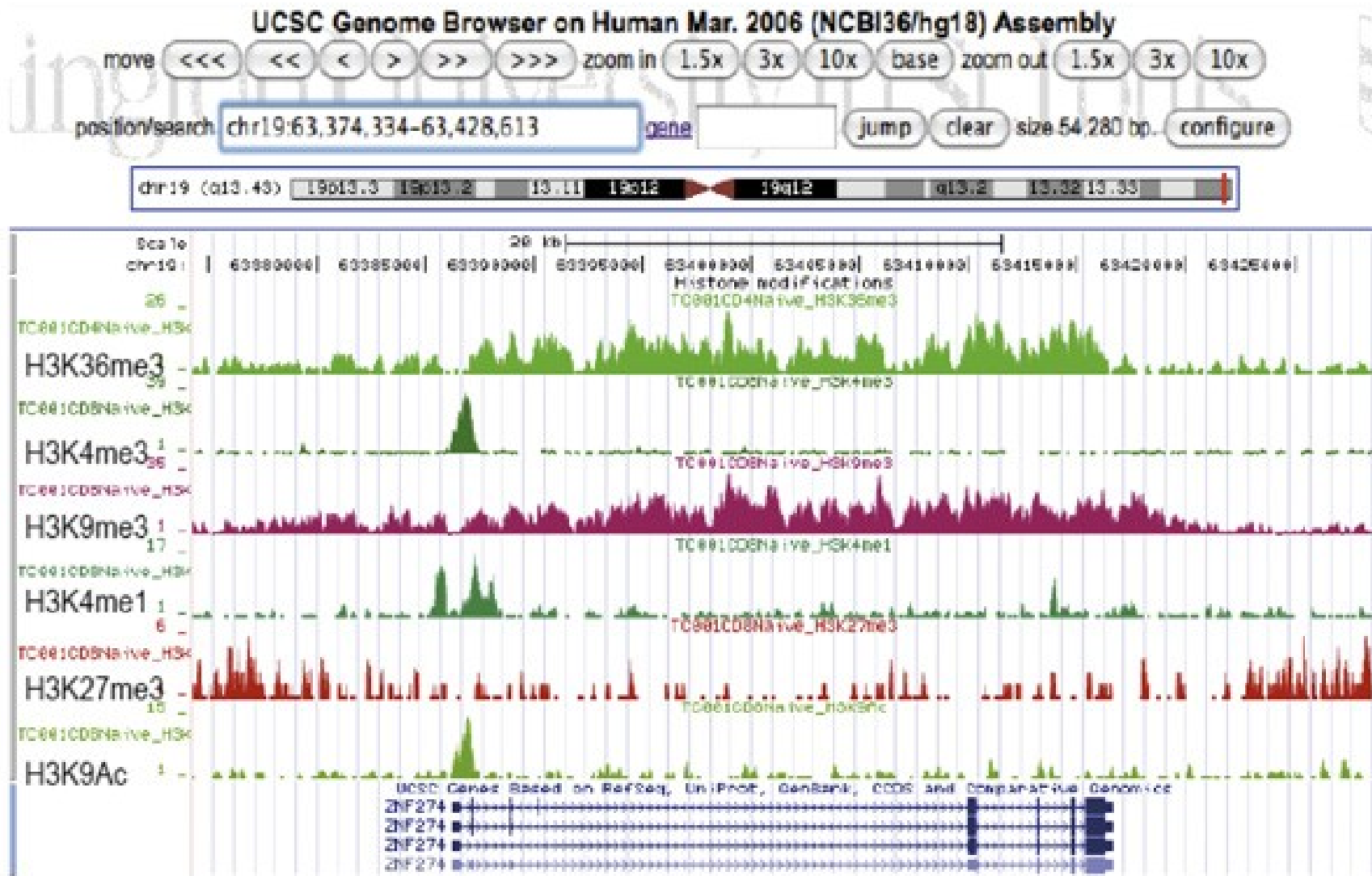
enhancers proximaux e distaux



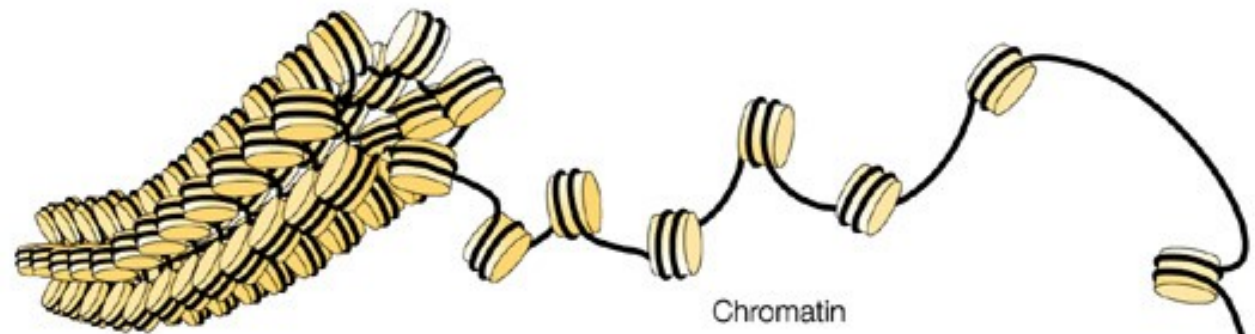
H3K27ac

Mapas de estados da cromatina baseados na localização de marcas de histonas informa sobre a atividade dos genes

visualização de dados de ChIP-seq em browsers genômicos

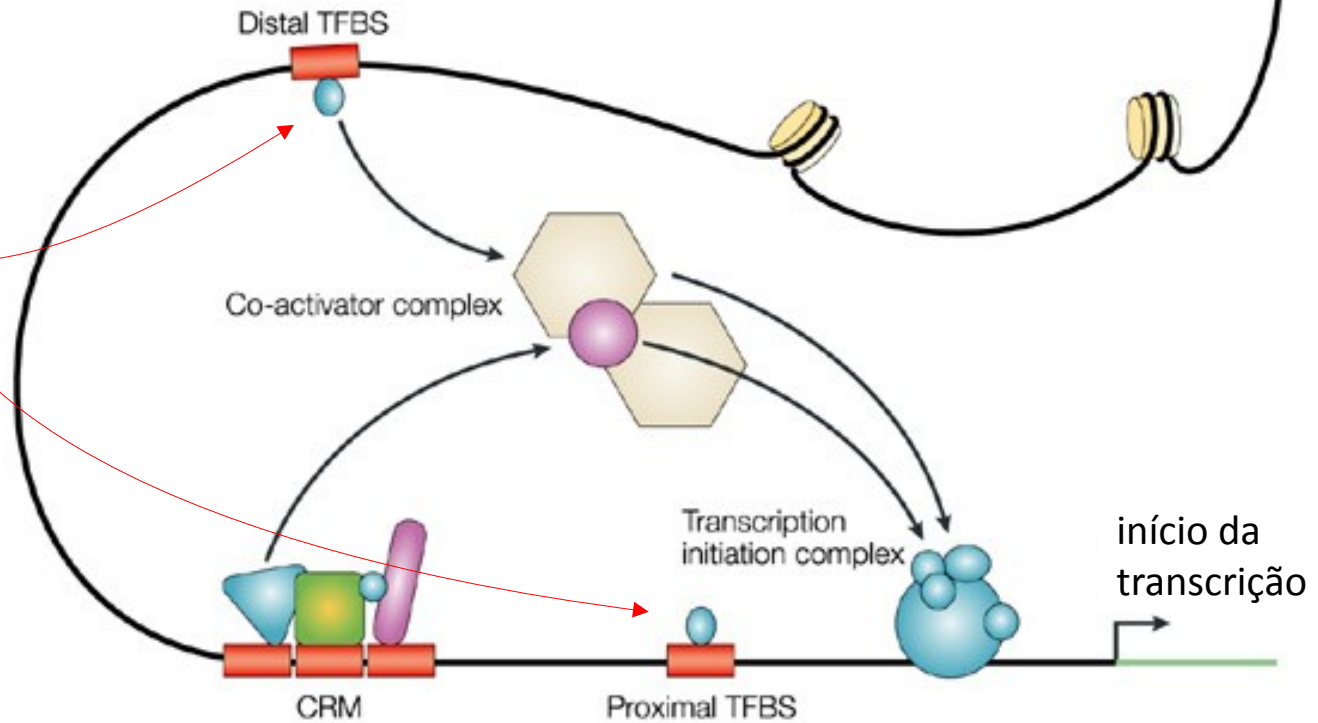


Fatores de transcrição se ligam ao DNA e ativam ou reprimem a transcrição



TFBS:
sítios de ligação
de fatores de
transcrição

Ativadores
da expressão
gênica



ChIP-seq:

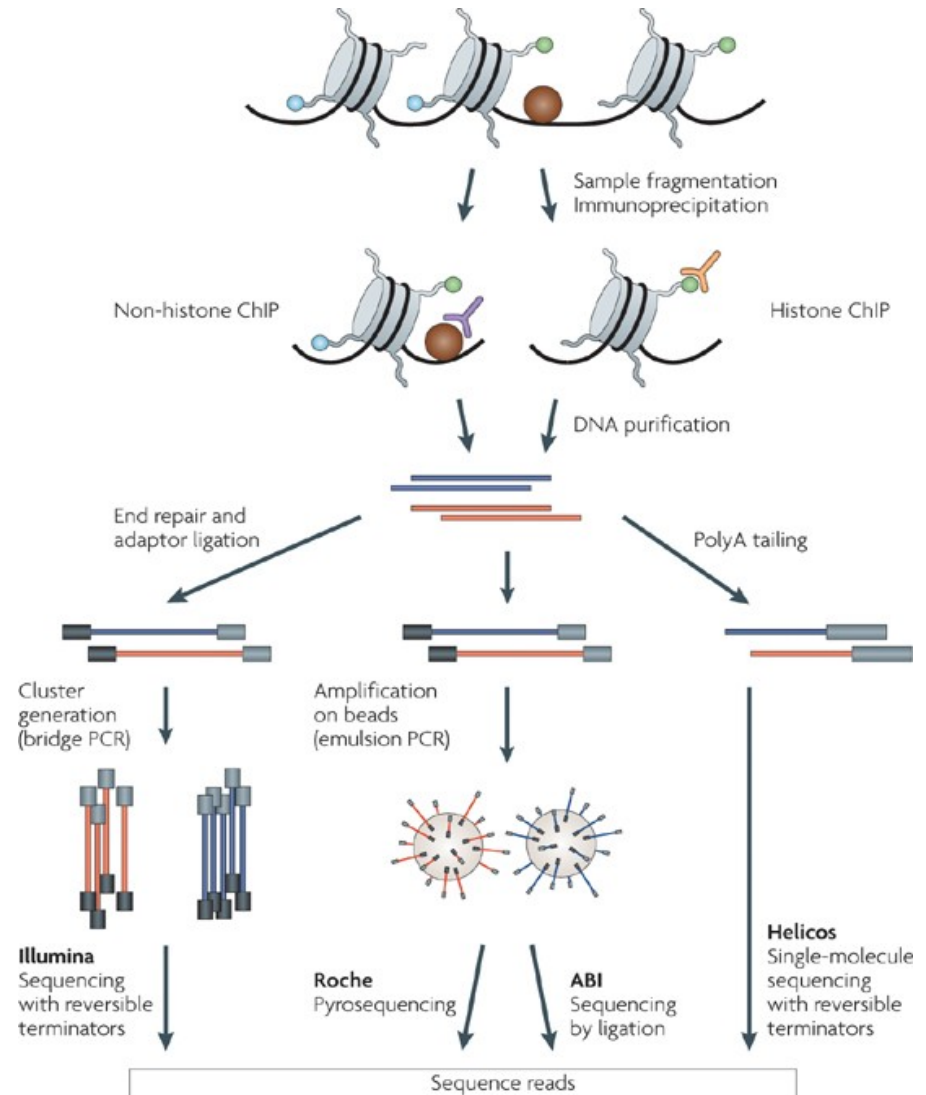
Imunoprecipitação de proteínas associadas a regiões específicas do DNA seguido de sequenciamento

São utilizados anticorpos que reconhecem especificamente a proteína de interesse associada ao DNA.

A cromatina é fragmentada e o DNA associado a proteína-alvo é recuperado por imunoprecipitação e analisado por NGS.

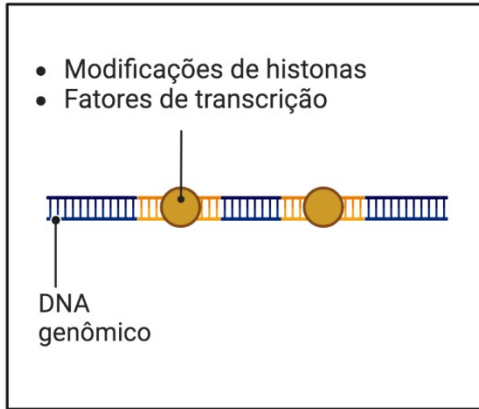
Permite identificar globalmente os sítios de ligação no DNA de :

- acetilação/metilação de histonas
- Fatores de transcrição

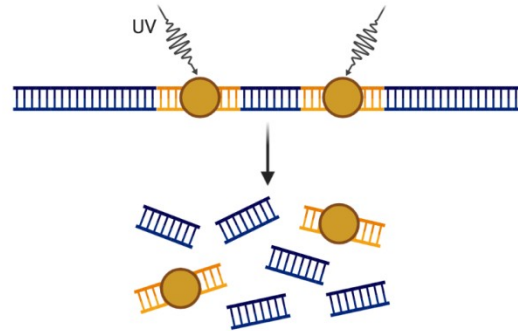


Etapas de um experimento de ChIP-seq

Identificar locais no genoma contendo marcas epigenéticas ou ocupados por proteínas regulatórias da cromatina

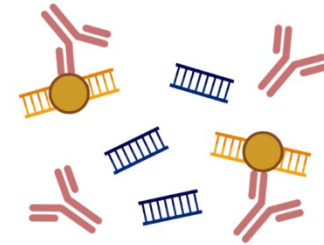


1 Crosslinking e fragmentação do DNA

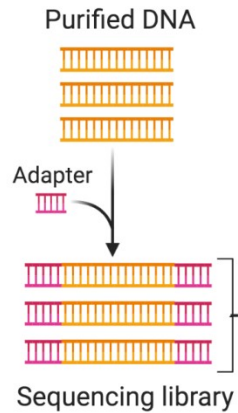


2 Immunoprecipitação da proteína e do DNA ligado

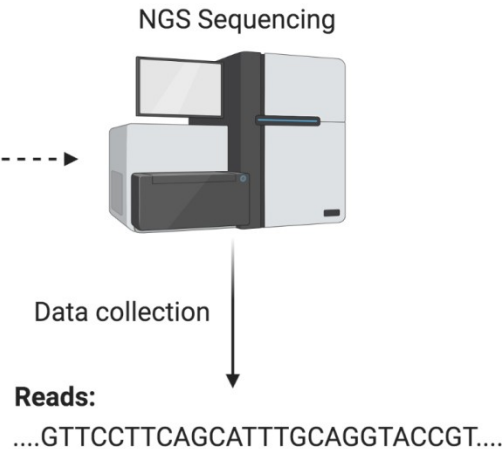
anticorpo específico contra a proteína-alvo



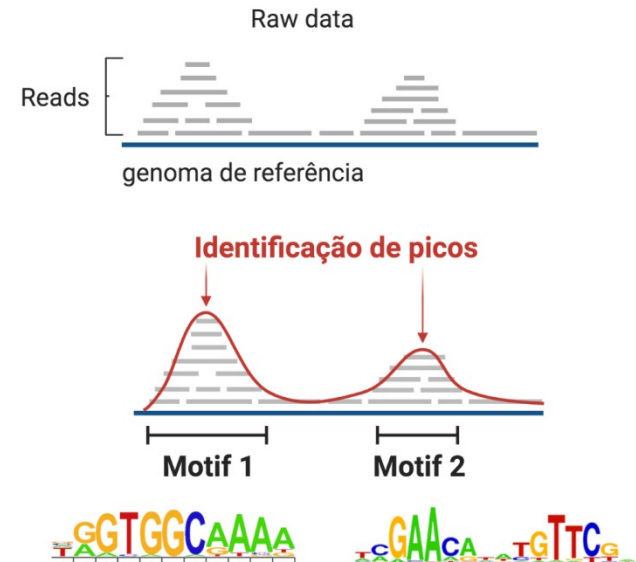
3 purificação do DNA e ligação de adaptadores para NGS



4 Sequenciamento NGS

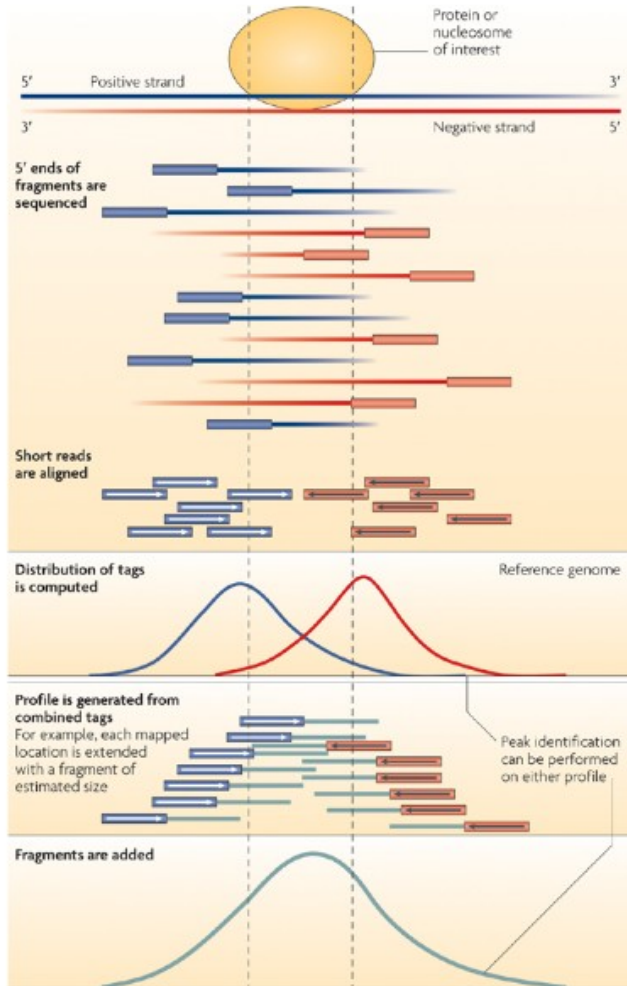


5 Identificação de sítios de ligação no genoma e de motivos no DNA



Exemplo de pipeline computacional para a análise de dados de ChIP-seq

- DNA associado a proteína-alvo co-imunoprecipitado com anticorpo (ChIP)
- DNA total (input)



Nature Reviews | Genetics

Park, P. J., *ChIP-seq: advantages and challenges a maturing technology*, *Nat Rev Genet.* Oct;10(10):669-80 (2009)

Passo 1: verificar a qualidade dos reads (FastQC)

Passo 2: Alinhar os reads no genoma (BWA, Bowtie)

Passo 3: Identificação de picos (MACS2)

Passo 4: anotação de picos

- Identificação de genes sobrepostos
- Identificação de picos comuns/exclusivos entre duas ou mais condições

- Uso do DNA input como referência aumenta a acurácia pois reduz o número de falsos positivos (regiões de cromatina aberta tendem a ser mais detectadas)

Tutorial – ChIP Seq (Galaxy)

No Galaxy Europe, criar um “history” e renomear como “tutorial ChIP-seq”

The screenshot displays the Galaxy Europe web interface. At the top, a dark blue navigation bar contains the Galaxy logo, the text "Galaxy Europe", and a series of icons for navigation: Home, Workflow, Visualize, Shared Data, Help, User, a graduation cap, a bell, a grid, and another bell. A green badge in the top right corner indicates "Using 7%".

On the left side, a "Tools" panel is visible, featuring a search bar with the text "search tools", an "Upload Data" button, and a list of tool categories. The categories include "Get Data", "Send Data", "Collection Operations", "GENERAL TEXT TOOLS", "GENOMIC FILE MANIPULATION", and "VCF/BCF".

On the right side, a "History" panel is shown. It has a search bar with the text "search datasets" and a "Tutorial ChIP-seq" entry. Below the entry, there are icons for a folder, a location pin, and a refresh symbol. A light blue informational message box states: "This history is empty. You can load your own data or get data from an external source."

No modo “tutorial” do Galaxy, navegar até o tutorial “Identification of the binding sites of the T-cell acute lymphocytic leukemia protein 1 (TAL1)”

The screenshot shows the Galaxy Training! interface. The top navigation bar includes 'Galaxy Europe', 'Workflow', 'Visualize', 'Shared Data', 'Help', 'User', and a dropdown menu. The main header contains 'Galaxy Training!', 'Contributors', 'Learning Pathways', 'Help', 'Settings', and a search bar. The main content area is titled 'Welcome to Galaxy Training!' and 'Collection of tutorials developed and maintained by the worldwide Galaxy community'. It is divided into two columns: 'Galaxy for Scientists' and 'Welcome to the GTNI!'. The 'Galaxy for Scientists' column lists various topics, with 'Epigenetics' circled in red. A red arrow points to this circled item. The 'Welcome to the GTNI!' column contains sections for 'Epigenetics', 'Requirements', and 'Material'. The 'Material' section is a table with columns for 'Lesson', 'Slides', 'Hands-on', 'Recordings', 'Input dataset', and 'Workflow'. The row for 'Identification of the binding sites of the T-cell acute lymphocytic leukemia protein 1 (TAL1)' is circled in red. A red arrow points to the 'Epigenetics' link in the left sidebar, and another red arrow points to the 'Epigenetics' link in the main content area.

Galaxy Training!

Contributors Learning Pathways Help Settings Search Tutorials

Welcome to Galaxy Training!

Collection of tutorials developed and maintained by the worldwide Galaxy community

Galaxy for Scientists

Topic

- [Introduction to Galaxy Analyses](#)
- [Assembly](#)
- [Climate](#)
- [Computational chemistry](#)
- [SARS-CoV-2](#)
- [Epigenetics](#)
- [Evolution](#)
- [Genome Annotation](#)
- [Imaging](#)

Welcome to the GTNI!

Epigenetics

DNA methylation is an epigenetic mechanism used by higher eukaryotes and involved in e.g. gene expression, X-Chromosome inactivating, imprinting, and gene silencing of germline specific gene and repetitive elements.

You can view the tutorial materials in different languages by clicking the dropdown icon next to the slides (📄) and tutorial (👤) buttons below.

Requirements

Before diving into this topic, we recommend you to have a look at:

- [Introduction to Galaxy Analyses](#)
- [Sequence analysis](#)
 - Quality Control: 📄 slides - 👤 hands-on
 - Mapping: 📄 slides - 👤 hands-on

Material

Lesson	Slides	Hands-on	Recordings	Input dataset	Workflow
ATAC-Seq data analysis	📄	👤	📺	📄	🔗
DNA Methylation data analysis		👤		📄	🔗
Formation of the Super-Structures on the Inactive X	📄	👤		📄	🔗
Hi-C analysis of Drosophila melanogaster cells using HICEplorer		👤		📄	🔗
Identification of the binding sites of the T-cell acute lymphocytic leukemia protein 1 (TAL1)		👤		📄	🔗
Infinium Human Methylation BeadChip	📄	👤		📄	🔗

Identification of the binding sites of the T-cell acute lymphocytic leukemia protein 1 (TAL1)

Authors:   Mallory Freeberg   Mo Heydarian  Vivek Bhardwaj  Joachim Wolff   Anika Erxleben

Overview

🔍 Questions:

- How is raw ChIP-seq data processed and analyzed?
- What are the binding sites of TAL1?
- Which genes are regulated by TAL1?

🎯 Objectives:

- Inspect read quality with FastQC
- Perform read trimming with Trimmomatic
- Align trimmed reads with BWA
- Assess quality and reproducibility of experiments
- Identify TAL1 binding sites with MACS2
- Determine unique/common TAL1 binding sites from G1E and Megakaryocytes
- Identify unique/common TAL1 peaks occupying gene promoters
- Visually inspect TAL1 peaks with Trackster

✅ Requirements:

- [Introduction to Galaxy Analyses](#)
- [Sequence analysis](#)
 - Quality Control:  [slides](#) -  [hands-on](#)
 - Mapping:  [slides](#) -  [hands-on](#)
- [Trackster](#)

🕒 **Time estimation:** 3 hours

📄 Supporting Materials:

 [Datasets](#)  [Workflows](#)  [FAQs](#)  [Available on these Galaxies](#) ▾



Agenda

In this tutorial, we will deal with:

1. Quality control
2. Trimming and clipping reads
3. Aligning reads to a reference genome
4. Assessing correlation between samples
5. Assessing IP strength
6. Determining TAL1 binding sites
7. Inspection of peaks and aligned data
 1. Inspection of peaks and aligned data with Trackster
 2. Inspection of peaks and aligned data with IGV
8. Identifying unique and common TAL1 peaks between stages
9. Generating Input normalized coverage files
10. Plot the signal on the peaks between samples
11. Additional optional analyses
 1. Assessing GC bias
12. Conclusion

Table 1: Metadata for ChIP-seq experiments in this tutorial. SE: single-end.

Cellular state	Datatype	ChIP Ab	Replicate	SRA Accession	Library type	Read length	Stranded?	Data size (MB)
G1E	ChIP-seq	input	1	SRR507859	SE	36	No	35.8
G1E	ChIP-seq	input	2	SRR507860	SE	55	No	427.1
G1E	ChIP-seq	TAL1	1	SRR492444	SE	36	No	32.3
G1E	ChIP-seq	TAL1	2	SRR492445	SE	41	No	62.7
Megakaryocyte	ChIP-seq	input	1	SRR492453	SE	41	No	57.2
Megakaryocyte	ChIP-seq	input	2	SRR492454	SE	55	No	403.8
Megakaryocyte	ChIP-seq	TAL1	1	SRR549006	SE	55	No	340.3
Megakaryocyte	ChIP-seq	TAL1	2	SRR549007	SE	48	No	356.9

Opção rápida para obter os dados para análise:

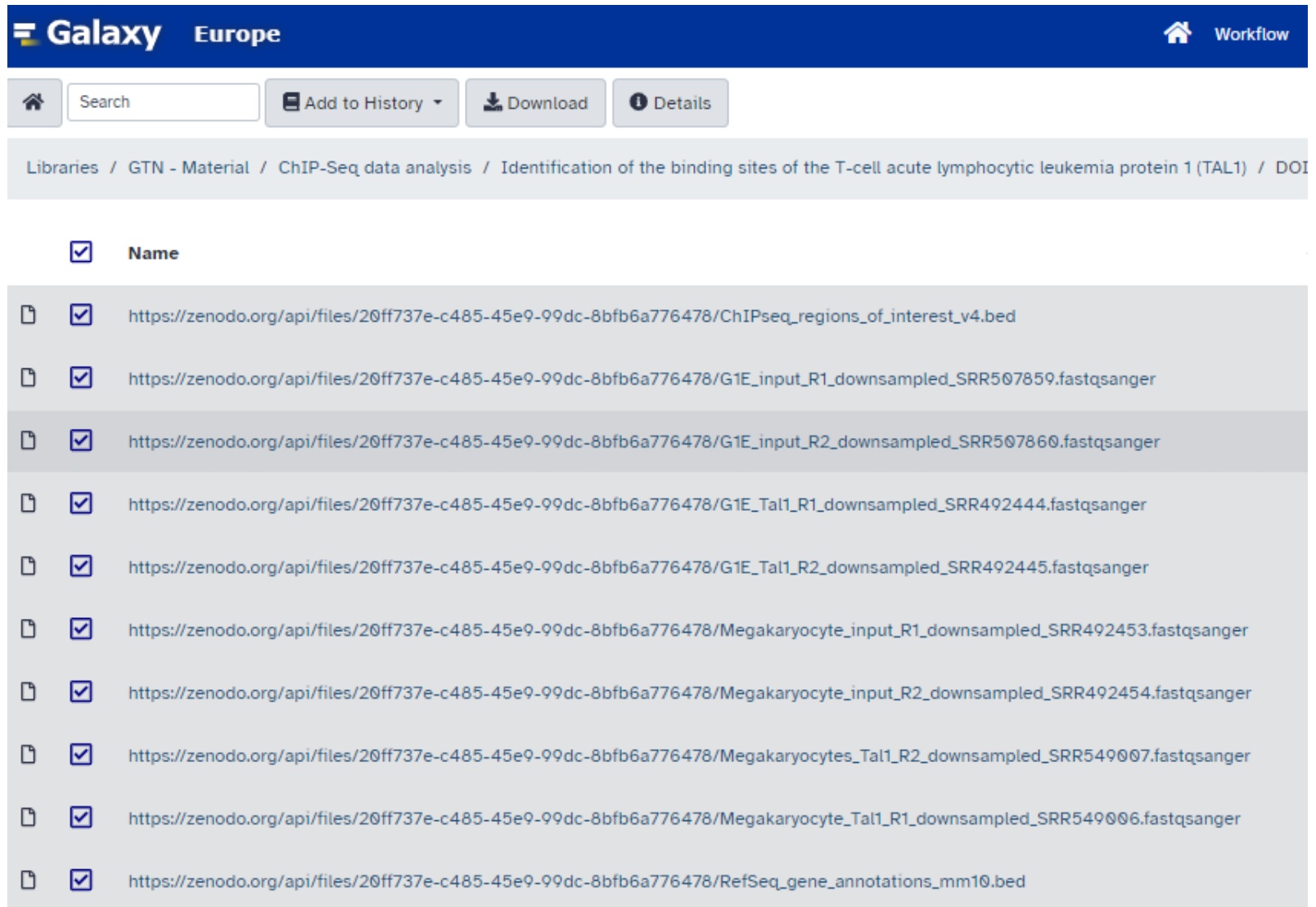
The image shows a navigation menu in a software application. The menu is titled "Shared Data" and contains the following items:

- Data Libraries
- Histories
- Workflows
- Visualizations
- Pages

The "Data Libraries" item is highlighted with a red arrow. Below it, a list of categories is shown, with "ChIP-Seq data analysis" highlighted by a red arrow. A second red arrow points from "ChIP-Seq data analysis" to a specific data item: "DOI: 10.5281/zenodo.197100".

Name
covid-19
Critical Assessment of Metagenome Inter
Earth System Community Modeling
Galaxy courses
Genomes + annot
GTN - Material
IAEA workshop
Mira_test
PGP-UK Open Ac
Street Science Co
Assembly
ChIP-Seq data analysis
diverse
Ecology
Epigenetics
Genome Annotation
Imaging
Introduction to Galat
Metabolomics
Metagenomics
Formation of the Super-Structures on the Inactive X
Identification of the binding sites of the Estrogen receptor
Identification of the binding sites of the T-cell acute lymphocytic leuke
DOI: 10.5281/zenodo.197100

Selecionar todos os arquivos e adicionar ao “history” ChIP-seq



The screenshot shows the Galaxy Europe web interface. At the top, there is a blue header with the Galaxy logo and the text "Europe". On the right side of the header, there is a home icon and the text "Workflow". Below the header, there is a navigation bar with a home icon, a search input field, and three buttons: "Add to History", "Download", and "Details".

Below the navigation bar, there is a breadcrumb trail: "Libraries / GTN - Material / ChIP-Seq data analysis / Identification of the binding sites of the T-cell acute lymphocytic leukemia protein 1 (TAL1) / DOI".

The main content area displays a list of files. Each row starts with a document icon, followed by a checked checkbox, and then the file name. All checkboxes are checked, indicating that all files are selected. The file names are URLs pointing to Zenodo files:




- Name
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/ChIPseq_regions_of_interest_v4.bed
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/G1E_input_R1_downsampled_SRR507859.fastqsanger
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/G1E_input_R2_downsampled_SRR507860.fastqsanger
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/G1E_Tal1_R1_downsampled_SRR492444.fastqsanger
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/G1E_Tal1_R2_downsampled_SRR492445.fastqsanger
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/Megakaryocyte_input_R1_downsampled_SRR492453.fastqsanger
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/Megakaryocyte_input_R2_downsampled_SRR492454.fastqsanger
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/Megakaryocytes_Tal1_R2_downsampled_SRR549007.fastqsanger
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/Megakaryocyte_Tal1_R1_downsampled_SRR549006.fastqsanger
- https://zenodo.org/api/files/20ff737e-c485-45e9-99dc-8bfb6a776478/RefSeq_gene_annotations_mm10.bed

<https://usegalaxy.eu/training-material/topics/epigenetics/tutorials/tal1-binding-site-identification/workflows/>

Aligning reads to a reference genome

To determine where DNA fragments originated from in the genome, the sequenced reads must be aligned to a reference genome. This is equivalent to solving a jigsaw puzzle, but unfortunately, not all pieces are unique. In principle, you could do a BLAST analysis to figure out where the sequenced pieces fit best in the known genome. Aligning millions of short sequences this way, however, can take a couple of weeks. Nowadays, there are many read alignment programs for sequenced DNA, BWA being one of them. You can read more about the BWA algorithm and tool [here](#).

Hands-on: Aligning reads to a reference genome

1.  **BWA** ( Galaxy version 0.7.17.4) : Run BWA to map the trimmed/clipped reads to the mouse genome.
 - "Will you select a reference genome...": **Use a built-in genome index**
 - "Using reference genome": **Mouse (mus musculus) mm10**
 - "Select input type": **Single fastq**
 -  "Select fastq dataset": Select all of the trimmed FASTQ files
2. Rename files to reflect the origin and contents

Determining TAL1 binding sites





Now that BWA has aligned the reads to the genome, we will use the tool MACS2 to identify regions of TAL1 occupancy, which are called "peaks". Peaks are determined from pileups of sequenced reads across the genome that correspond to where TAL1 binds.

MACS2 will perform two tasks:

1. Identify regions of TAL1 occupancy (peaks).
2. Generate bedGraph files for visual inspection of the data on a genome browser.

More information about MACS2 can be found in [Zhang et al. 2008](#).

Hands-on: Determining TAL1 binding sites


1.  **MACS2 callpeak** ( Galaxy version 2.1.1.20160309.6) : Run MACS2 callpeak with the aligned read files from the previous step as Treatment (TAL1) and Control (input).
 - "Are you pooling Treatment Files?": **Yes**
 -  "ChIP-Seq Treatment File": Select all of the replicate ChIP-Seq treatment aligned BAM files for one cell type
 - "Do you have a Control File?": **Yes**
 - "Are you pooling Control Files?": **Yes**
 -  "ChIP-Seq Control File": Select replicate ChIP-Seq control aligned BAM files for the same cell type
 - "Format of Input Files": **Single-end BAM**
 - "Effective genome size": **M. musculus**
 - "Additional Outputs": Select **Peaks as tabular file (compatible with MultiQC)**, **Peak summits**, **Scores in bedGraph files (--bdg)**
2. Rename files to reflect the origin and contents.
3. Repeat for the other cell type.

[Link to here](#) |  [FAQs](#) | [Gitter Chat](#) | [Help Forum](#)

Identifying unique and common TAL1 peaks between stages

We have processed ChIP-seq data from two stages of hematopoiesis and have lists of TAL1-occupied sites (peaks) in both cellular states. The next analysis step is to identify TAL1 peaks that are *shared* between the two cellular states and peaks that are *specific* to either cellular state.

Hands-on: Identifying unique and common TAL1 peaks between states

1.  **bedtools Intersect intervals** (🔗 Galaxy version 2.29.0) : Run bedtools Intersect intervals to find peaks that exist both in G1E and megakaryocytes.
 - "File A to intersect with B": Select the TAL1 G1E narrow peaks BED file
 - "File B to intersect with A": Select the TAL1 Megakaryocytes narrow peaks BED file
 - Running this tool with the default settings will return overlapping peaks of both files.
2.  **bedtools Intersect intervals** (🔗 Galaxy version 2.29.0) : Run bedtools Intersect intervals to find peaks that exist only in G1E.
 - "File A to intersect with B": Select the TAL1 G1E narrow peaks BED file
 - "File B to intersect with A": Select the TAL1 Megakaryocytes narrow peaks BED file
 - "Report only those alignments that ****do not**** overlap the BED file": **Yes**
3.  **bedtools Intersect intervals** (🔗 Galaxy version 2.29.0) : Run bedtools Intersect intervals to find peaks that exist only in megakaryocytes.
 - "File A to intersect with B": Select the TAL1 Megakaryocytes narrow peaks BED file
 - "File B to intersect with A": Select the TAL1 G1E narrow peaks BED file
 - "Report only those alignments that ****do not**** overlap the BED file": **Yes**
4. Rename files to reflect the origin and contents.

Instruções para o relatório

- 1) Descreva sucintamente os passos usados na análise de CHIP-seq
- 2) Responda as seguintes perguntas, explicando como chegou a essas conclusões:

Quantos picos TAL1 são comuns às células G1E indiferenciadas e aos megacariócitos?

Quantos picos são exclusivos das células G1E?

Quantos picos são exclusivos dos megacariócitos?