# **Bioinformatics Lab**

Ivan Gesteira Costa, Mingbo Cheng, Zhijian Li, James Nagai, Mina Shaigan Institute for Computational Genomics



# **Computational Epigenomics**



#### **Cell Differentiation**

#### Hematopoiesis





#### **Cell Differentiation**





## **Regulatory Control – Transcription Factor Binding**

![](_page_4_Figure_1.jpeg)

![](_page_4_Picture_2.jpeg)

Source: Alberts, B. et al. (2008) Garland Science, 5th ed.

## **Regulatory Control – Transcription Factor Binding**

![](_page_5_Figure_1.jpeg)

Source: Alberts, B. et al. (2008) Garland Science, 5th ed.

![](_page_5_Picture_3.jpeg)

# **Epigenetics & Histones**

![](_page_6_Figure_1.jpeg)

![](_page_6_Picture_2.jpeg)

Modification in histone tailschange strength of DNA bindingrecruit transcription factors

![](_page_6_Picture_4.jpeg)

### **Chromatin, Regulation and Cellular Memory**

![](_page_7_Figure_1.jpeg)

![](_page_7_Picture_2.jpeg)

Adapted from Lodish, B. et al. (2004) 5th ed.

#### **Chromatin & Histone Code**

![](_page_8_Figure_1.jpeg)

![](_page_8_Picture_3.jpeg)

#### **Chromatin with Next Generation Sequencing**

![](_page_9_Figure_1.jpeg)

Source: Meyer, C.A. and Liu X.S. (2014). Nature Reviews Genetics.

![](_page_9_Picture_4.jpeg)

![](_page_10_Figure_0.jpeg)

Source: Meyer, C.A. and Liu X.S. (2014). Nature Reviews Genetics.

![](_page_10_Picture_3.jpeg)

![](_page_11_Figure_1.jpeg)

![](_page_11_Picture_2.jpeg)

![](_page_11_Picture_3.jpeg)

![](_page_12_Figure_1.jpeg)

![](_page_12_Picture_2.jpeg)

![](_page_12_Picture_3.jpeg)

![](_page_13_Figure_1.jpeg)

![](_page_13_Picture_2.jpeg)

![](_page_14_Figure_1.jpeg)

![](_page_14_Picture_2.jpeg)

![](_page_15_Figure_1.jpeg)

![](_page_15_Picture_3.jpeg)

# **Bioinformatics Pipeline / ATAC-seq**

![](_page_16_Figure_1.jpeg)

Adapted from Rasmussen: http://www.cbs.dtu.dk/courses/27626/programme.php

#### Example of a simple peak caller :

- use a fix window to scan through the genome and obtain a distribution of counts per bin
- define a statistical test to evaluate if the number of reads in higher than expected by change

![](_page_17_Figure_5.jpeg)

Aligned Reads

See for an example of a code for a peak caller http://www.regulatory-genomics.org/rgt/tutorial/implementing-your-own-peak-caller/

![](_page_17_Picture_8.jpeg)

#### Example of a simple peak caller :

- use a fix window to scan through the genome and obtain a distribution of counts per bin
- define a statistical test to evaluate if the number of reads in higher than expected by change

![](_page_18_Figure_5.jpeg)

![](_page_18_Figure_6.jpeg)

![](_page_18_Figure_7.jpeg)

Counts: 2

![](_page_18_Picture_10.jpeg)

#### Example of a simple peak caller :

- use a fix window to scan through the genome and obtain a distribution of counts per bin
- define a statistical test to evaluate if the number of reads in higher than expected by change

![](_page_19_Figure_5.jpeg)

![](_page_19_Figure_6.jpeg)

![](_page_19_Figure_7.jpeg)

Counts: 2 4

![](_page_19_Picture_10.jpeg)

#### Example of a simple peak caller :

- use a fix window to scan through the genome and obtain a distribution of counts per bin
- define a statistical test to evaluate if the number of reads in higher than expected by change

See for an example of a code for a peak caller http://www.regulatory-genomics.org/rgt/tutorial/implementing-your-own-peak-caller/

![](_page_20_Figure_6.jpeg)

![](_page_20_Figure_7.jpeg)

![](_page_20_Picture_9.jpeg)

#### Example of a simple peak caller :

- use a fix window to scan through the genome and obtain a distribution of counts per bin
- define a statistical test to evaluate if the number of reads in higher than expected by change

#### See for an example of a code for a peak caller http://www.regulatory-genomics.org/rgt/tutorial/implementing-your-own-peak-caller/

#### Aligned Reads

![](_page_21_Figure_7.jpeg)

![](_page_21_Picture_9.jpeg)

#### Example of a simple peak caller :

- use a fix window to scan through the genome and obtain a distribution of counts per bin
- define a statistical test to evaluate if the number of reads in higher than expected by change

#### Aligned Reads

![](_page_22_Figure_6.jpeg)

![](_page_22_Figure_7.jpeg)

See for an example of a code for a peak caller

http://www.regulatory-genomics.org/rgt/tutorial/implementing-your-own-peak-caller/

![](_page_22_Picture_11.jpeg)

#### Example of a simple peak caller :

- use a fix window to scan through the genome and obtain a distribution of counts per bin
- define a statistical test to evaluate if the number of reads in higher than expected by change

![](_page_23_Figure_5.jpeg)

![](_page_23_Figure_6.jpeg)

![](_page_23_Figure_7.jpeg)

See for an example of a code for a peak caller

http://www.regulatory-genomics.org/rgt/tutorial/implementing-your-own-peak-caller/

![](_page_23_Picture_11.jpeg)

![](_page_24_Figure_2.jpeg)

See for an example of a code for a peak caller

http://www.regulatory-genomics.org/rgt/tutorial/implementing-your-own-peak-caller/

#### **Peak calling in ATAC-seq**

![](_page_25_Figure_1.jpeg)

- MACS2
  - most frequently used
- HMMRATAC
  - ATAC-seq specific peak caller
  - ignores reads from large fragments / linker cleavage sites

![](_page_25_Picture_7.jpeg)

![](_page_25_Picture_8.jpeg)

# **Bioinformatics Pipeline / ATAC-seq**

![](_page_26_Figure_1.jpeg)

Adapted from Rasmussen: http://www.cbs.dtu.dk/courses/27626/programme.php

![](_page_26_Picture_3.jpeg)

### **Motif Search – Computational Approach**

![](_page_27_Figure_1.jpeg)

![](_page_27_Picture_3.jpeg)

# **Model for DNA-protein binding**

#### PU.1 binding sites

Kanno, Y. et al. (2005) Immune Cell-Specific Amplification of Interferon Signaling by the IRF-4/8-PU.1 Complex.

AGGAACT
GGGAACA
AGAAAGT
AGGAACT
GAGAAGT
AGGAAGC
AGGAACC

![](_page_28_Picture_4.jpeg)

# **Model for DNA-protein binding**

#### PU.1 binding sites

Kanno, Y. et al. (2005) Immune Cell-Specific Amplification of Interferon Signaling by the IRF-4/8-PU.1 Complex.

PU.1 Position

Weight Matrix (PWM)

N	IuMH	IC I	Α	G	G	A	A	C	T
Η	u <b>M</b> xA	4	G	G	G	A	A	С	A
→ H	uIFN	-β	A	G	A	A	A	G	Т
N	luβ <sub>2</sub> m	l	A	G	G	A	A	С	Т
Η	uGBI	P	G	A	G	A	A	G	Т
Η	istone	e H4	A	G	G	A	A	G	С
H	HuIFN-α			G	G	A	A	С	С
		_	↓ ↓	¥	¥	V	V	¥	↓
		A	5	1	1 '	7	7	3	1
		С	0	0	0	0	0	0 2	2
	<b></b>	G	2	6	6	0	0	4	0
		Т	0	0	0	0	0	0	4

![](_page_29_Picture_4.jpeg)

# **Model for DNA-protein binding**

#### PU.1 binding sites

Kanno, Y. et al. (2005) Immune Cell-Specific Amplification of Interferon Signaling by the IRF-4/8-PU.1 Complex.

> PU.1 Position Weight Matrix (PWM)

> > PU.1 Logo

![](_page_30_Figure_4.jpeg)

![](_page_30_Picture_6.jpeg)

PU.1 PWM

Genome TATCTTTGGAAGTGAAACTACTATCCTGAAACTCGAA

![](_page_31_Picture_3.jpeg)

PU.1 PWM Genome TATCTTTGGAAGTGAAACTACTATCCTGAAACTCGAA Score 10.06

![](_page_32_Picture_2.jpeg)

![](_page_33_Figure_1.jpeg)

![](_page_33_Picture_2.jpeg)

![](_page_34_Figure_1.jpeg)

![](_page_34_Picture_2.jpeg)

PU.1 PWM<sup>#</sup>

![](_page_35_Picture_2.jpeg)

![](_page_35_Figure_3.jpeg)

![](_page_35_Picture_4.jpeg)

![](_page_35_Picture_5.jpeg)

PU.1 PWM<sup>#</sup>

![](_page_36_Picture_2.jpeg)

![](_page_36_Figure_3.jpeg)

![](_page_36_Picture_4.jpeg)

**PU.1 PWM**<sup>#</sup>

![](_page_37_Picture_2.jpeg)

![](_page_37_Figure_3.jpeg)

Genome Position (bp)

![](_page_37_Picture_5.jpeg)

## **Example: Binding sites in ID2**

# Motif search for binding sites with 536 PWMs (Jaspar & Uniprobe) and FDR=0,01

![](_page_38_Figure_2.jpeg)

> 3000 predicted binding sites

![](_page_38_Picture_4.jpeg)

![](_page_39_Figure_1.jpeg)

![](_page_39_Picture_2.jpeg)

![](_page_40_Figure_1.jpeg)

![](_page_40_Picture_2.jpeg)

![](_page_40_Picture_3.jpeg)

![](_page_41_Figure_1.jpeg)

![](_page_41_Picture_2.jpeg)

![](_page_42_Figure_1.jpeg)

![](_page_42_Picture_2.jpeg)

![](_page_43_Figure_1.jpeg)

![](_page_43_Picture_2.jpeg)

# **Problem definition**: Find genomic regions (of small size) with depletion in DNase-seq signals

![](_page_44_Figure_2.jpeg)

![](_page_44_Picture_3.jpeg)

# HINT (Hmm-based IdeNtification of Transcription factor footprints)

- generate normalized cleavage signals
- trained with limited supervision
- scan multivariate signals to predict footprints

![](_page_45_Figure_5.jpeg)

![](_page_45_Picture_6.jpeg)

BACK

TOP

FOOT PRINT

UP

Histone

Level

DNase Level

DOWN

# HINT (Hmm-based IdeNtification of Transcription factor footprints)

- generate normalized cleavage signals
- trained with limited supervision
- scan multivariate signals to predict footprints

#### **Prediction Example**

![](_page_46_Figure_6.jpeg)

Gusmao EG *et. al*, (2014), Bioinformatics, 30(22):3143-51. Gusmao EG *et. al*, (2016), Nature Methods, 13, 303–309. Li et al. (2019), Genome Biology, 20:45.

![](_page_46_Picture_8.jpeg)

BACK

TOP

FOOT PRINT

UP

Histone

Level

DNase

Level

DOWN

# HINT (Hmm-based IdeNtification of Transcription factor footprints)

- generate normalized cleavage signals
- trained with limited supervision
- scan multivariate signals to predict footprints

#### **Prediction Example**

![](_page_47_Figure_6.jpeg)

![](_page_47_Picture_8.jpeg)

BACK

TOP

FOOT PRINT

UP

Histone

Level

DNase

Level

DOWN

![](_page_48_Figure_1.jpeg)

![](_page_48_Picture_3.jpeg)

![](_page_49_Figure_1.jpeg)

0101101110100100

Li, ..., Kramann, Costa, Biorvx, https://doi.org/10.1101/865931.

#### **Computational Challenges - Single Cell ATAC**

![](_page_50_Figure_1.jpeg)

![](_page_50_Picture_2.jpeg)

#### **Computational Challenges - Single Cell ATAC**

![](_page_51_Figure_1.jpeg)

![](_page_51_Picture_2.jpeg)

#### **Resume / Single cell clustering**

- Finding groups of single cells require complex pipeline:
  - Cell filtering
  - Normalisation
  - Artefact removal
  - Dimension reduction
  - Integration
  - Clustering
  - Cell annotation / visualisation
- Open points:
  - How to deal with large data sets (millions of cells)?
  - How to detect cells of rare populations?
  - How to deal with sparsity of scATAC seq data?

![](_page_52_Picture_13.jpeg)

### **Clustering of cells / Human Fetal Cell Atlas**

scRNA-seq

scATAC-seq

Single-cell chromatin accessibility profiles 790,957 cells

![](_page_53_Figure_4.jpeg)

![](_page_53_Figure_5.jpeg)

https://descartes.brotmanbaty.org/

- Open points:
  - · How to deal with large data sets (millions of cells)?
  - How to detect cells of rare populations?
  - How to deal with sparsity of scATAC seq data?

![](_page_53_Picture_11.jpeg)

- Review basic biological/computational aspects
  - 1. basics of molecular biology
  - 2. basics of sequencing
  - 3. basics bioinformatics problems
    - short sequences read alignment
    - gene expression quantification
    - single cell approaches
    - computational epigenetic (today)

![](_page_54_Picture_9.jpeg)

Today – Introduction to Bioinformatics, Next Generation Sequencing

- 2.05.2022 Single cell sequencing Practical
- 8.05.2022 Computational Epigenomics / Theory & Practical / Using RWTH HPC/GPU cluster
- 15.05.2022 4.7.2022 Project development
- 11.07.2022 Project Presentation

Communication/discord channel: <u>https://discord.gg/</u> <u>hmGxznNpZH</u>.

![](_page_55_Picture_7.jpeg)

# Thank you!

![](_page_56_Picture_1.jpeg)