# **Bioinformatics Lab**

## Ivan Gesteira Costa & Zhijian Li Institute for Computational Genomics



# **Objectives**

- Hands on introduction to bioinformatics programing
- Review basic biological/computational aspects
  - 1. basics of molecular biology
  - 2. basics of sequencing
  - 3. basics bioinformatics problems
    - short sequences read alignment
    - differential peak calling
    - motif detection
    - footprint detection



# **Objectives**

- Introduction to Bioinformatics Frameworks/Tools
  - 1. biological sequence data formats/handling
    - Biopython, Pysam
  - 2. bioinformatics tools
    - BWA (aligner), MACS/ODIN(Peak caller), RGT/ bedtools (interval algebra), regulatory genomics toolbox (RGT)



# **Grading/Online material**

#### **Evaluation:**

- 20% prototypes
- 60% final project
- 20% presentation

Extra-work (anyone from media informatics?):

research report

#### **References/Courses Online**

http://costalab.org/teaching/bioinformatics-software-lab-2018/



# **Introduction to Molecular Biology**



- How is genetic information inherited?
- How the genetic information influence cellular processes?
- How genes work together to promote particular molecular functions?



# **Genetic Information - DNA**



## DNA (Deoxyribonucleic)

- chain of nucleic acids
- 4 bases: A;C;G;T
- forms DNA duplexes with paring A = T e C = G



# **Central Dogma - Transcription**



## Transcription

• DNA to RNA

## RNA (ribonucleic acid)

- single stranded
- 4 bases: A;C;G;U
- unstable
- transport of information from nucleus to cytoplasm



# **Central Dogma - Transcription**



Figure 1-5 Molecular Biology of the Cell 5/e (© Garland Science 2008)

#### Transcription - copy of DNA information to RNA (T to U)



# **Central Dogma - Translation**



## Translation

- RNA to Protein
- performed by the ribosome
- follows the genetic code

# Proteins

- single stranded chain
- 20 amino acids
- assumes 3D structure
- main functional entities in the cell



## **Genetic Code - Translation**



Figure 6-50 Molecular Biology of the Cell 5/e (© Garland Science 2008)

#### triples of RNA bases encodes a amino acid



# **Central Dogma**



- Dogma: information flux
  DNA -> mRNA -> Proteins
- Gene: DNA segment coding a protein.
- Transcript: RNA segment associated to a gene.
- Genes is associated to one proteins and one function\*

\* Genes might be associated to many proteins



# **Control of Gene Expression**



Figure 6-19 Molecular Biology of the Cell 5/e (© Garland Science 2008)



# **Gene Expression**





# **Gene Regulation / Motif Search**



## **Regulatory Control – Protein-DNA interaction**





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

## **Motif Search – Computational Approach**





# **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.

MuMHC I	AGGAACT
HuMxA	GGGAACA
HuIFN-β	AGAAAGT
$Mu\beta_2m$	AGGAACT
HuGBP	GAGAAGT
Histone H4	AGGAAGC
HuIFN-α	AGGAACC
	MuMHC I HuMxA HuIFN-β Muβ <sub>2</sub> m HuGBP Histone H4 HuIFN-α



# **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)

AGGAACT
GGGAACA
AGAAAGT
AGGAACT
GAGAAGT
AGGAAGC
AGGAACC
5117731
000002
2660040
000004



# **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







PU.1 PWM

Genome TATCTTTGGAAGTGAAACTACTATCCTGAAACTCGAA















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









PU.1 PWM<sup>®</sup>









# **Example: Binding sites in ID2**

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





# **Epigenetics**



# **Cellular Complexity**



Figure 7-1 Molecular Biology of the Cell 5/e (© Garland Science 2008)

#### Two cells of a organism have exactly\* the same DNA

#### How does this differences arise? How is cell fate remembered?

\* with exception of somatic mutations and rearrangements of immunological loci



# **Epigenetics & Histones**





Modification in histone tails - change strength of DNA binding

- recruit transcription factors



## **Chromatin Structure and Regulation**



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



# **Epigenetics**





#### **Protein-DNA interactions with Next Generation Sequencing**



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









#### Read the bases of a particular DNA/RNA sequence

# **Applications:**

- sequence DNA of known and unknown organism
- detect variants on patients
- sequence the RNA of a cell
- detect location of proteins interacting with DNA
  Problem:
- only short DNA sequences (<1.000 bs) can be read</li>
  Solution: break DNA in several small pieces and
  bioinformatics



# **Next Generation Sequencing**

- NGS take advantage of parallelization
  - reads millions/billions of reads for a time
  - shorter reads (50-100 bps)
  - higher error rates (0.1-1%)
- commercial products:
  - 454
  - **SOLiD**
  - Solexa (Illumina)






# **Illumina Flow Cell - NGS Sequencing**

1- fragment sample DNA, insert adapters, attach to flow cell

2- use (bridge) PCR to copy fragments (close to origin)

3- clusters of single stranded DNA (200m clusters with 2k DNA strands



See video http://www.wellcome.ac.uk/Education-resources/Education-and-learning/Resources/Animation/WTX056051.htm



# **Illumina Flow Cell - NGS Sequencing**

- Iterative evaluation process:
  - 1. add RT-bases, polymerases integrate them
  - 2. wash away all not integrated elements
  - 3. take picture of flow cell to determine current base by dye
  - 4. derive reads from pictures







# **Sequencing Results**



 $P = 10^{(-Q/10)}$ 



## **Next Generation Sequencing**

Improvements in the rate of DNA sequencing over the past 30 years



Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. Nature 458, 719-724 (2009).



## **Sequencing Costs**





# **Sequence Alignment**



# **Sequence Alignment**

## NGS

- reads from DNA fragments
- position in genome is unknown
- solution: alignment

# **DNA Sequencing**

- de-novo assembly
  - construct unknown reference sequence from scratch
- resequencing / mapping
  - reference sequence given (applies to human- and mousestudies)
  - build sequence that is similar but not necessarily identical to reference sequence



# **Alignment Problem**

- a large reference sequence is given (genome)
  - up to billions of base pairs
- short reads (<200bps)
- find most probable position of the read in the genome (by inexact string matching)





- (Unknown) divergent of sample and reference genome
- Repeats in the genome (larger than read size)
- Recombinations
- Poor genome reference quality
- Sequencing/read errors



# Alignment/Mapping is a typical inexact string match problem

**Algorithmic Solutions: ?** 



Alignment/Mapping is a typical inexact string match problem

## **Algorithmic Solutions:**

• Smith & Waterman - dynamic programming (quadratic time/memory)



Alignment/Mapping is a typical inexact string match problem

## **Algorithmic Solutions:**

- Smith & Waterman dynamic programming (quadratic time/memory)
- Blast k-mer search for seeding followed by
  dynamic programming
  - large memory requirement
  - local alignment



- reference sequence is large and fixed
- query sequence (reads) are short and many
   Solution: ?



- reference sequence is large and fixed
- query sequence (reads) are short and many
   Solution: ?
- **1. Use a data structure to represent reference** 
  - k-mer hash table (>40GB)
  - suffix trees (> 4GB)



- reference sequence is large and fixed
- query sequence (reads) are short and many
   Solution: ?
- **1. Use a data structure to represent reference** 
  - k-mer hash table (>40GB)
  - suffix trees (> 4GB)
- 2. Find candidate (k-mer) hits on genome (>100)



- reference sequence is large and fixed
- query sequence (reads) are short and many
   Solution: ?
- **1. Use a data structure to represent reference** 
  - k-mer hash table (>40GB)
  - suffix trees (> 4GB)
- 2. Find candidate (k-mer) hits on genome (>100)
- 3. Improve alignment with Smith-Waterman Methods work on linear time (query sequence)



## Hash based algorithm





# **Computational Epigenomics**



## **Computational Epigenomics - Problems**

## 1. Motif Search (already seen)

- search transcription factor binding sites in genomic sequences

## 2. Peak Calling

- find regions with high number of ChIP-Seq signals

## 3. Digital Footprinting

- find genomic regions with depletion of DNase-seq signals



### **DNA - Protein interactions with NGS**



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

Institute for Computational Genomics 01011011010 1010010010





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































#### 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



#### 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





Counts: 2



#### 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





Counts: 2 4



#### 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







#### 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





#### 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

Counts

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



#### 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/





- which window size to use?
- distinct proteins have distinct peak sizes
- proper quantification of read counts require several further steps: fragment size estimation, CG bias correction, mappability, ...



We will see examples in the next See lectures ....

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


# **Peak Calling - Example for Transcription Factors**

Example of analysis of ChIP-Seq for transcription factors (small peaks)

Histone code	
H3K79me2 - Transcribed	
H3K27ac - Active	
H3K27me3 - Repressed	







# **Peak Calling - Example for Transcription Factors**

Example of analysis of ChIP-Seq for transcription factors (small peaks)





# **Peak Calling - Example for Transcription Factors**

Example of analysis of ChIP-Seq for transcription factors (small peaks)





### **Peak Calling - Example for Histones**

Example of analysis of ChIP-Seq for histones (medium to broad peaks)





#### **DNA - Protein interactions with DNase-Seq**



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









DNAse cleavage



























#### **Open chromatin with ATAC-seq**



Li et. al, unpublished



#### **Open chromatin with ATAC-seq**



Li et. al, unpublished



# **Computational Footprint Methods**





# **Computational Footprint Methods**



Site-centric: Classify motif-predicted binding sites as "bound" or "unbound".



\* only applicable if "true positives" TF are known

Institute for Computational Genomics 01011011010 1010010010

# 16.04.2018 - Introduction to Biology & Bioinformatics 30.04.2018 - Example of NGS pipeline / data formats 07.05.2018 - Intro to RGT / Problem definition ... - Project development 16.07.2015 - Final Presentation



# Thank you!



