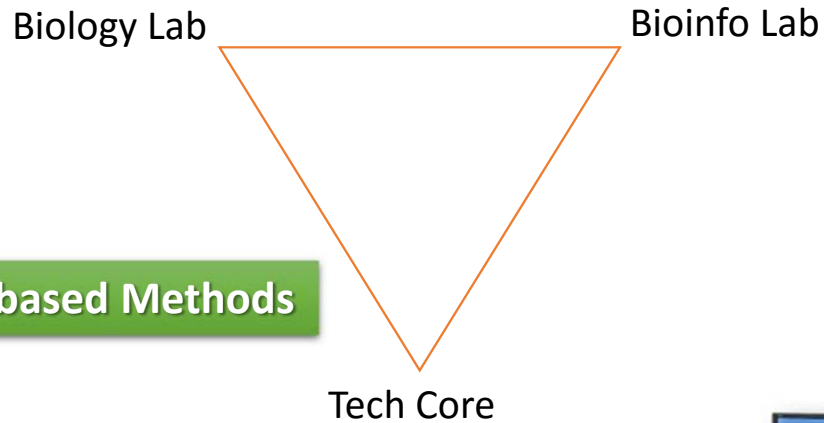# Multi-Omics onLine Analysis System for Gene Expression Profiling and Whole Methylome

Life Science Library Training Course
2016/06/09

Chen, Shu-Hwa
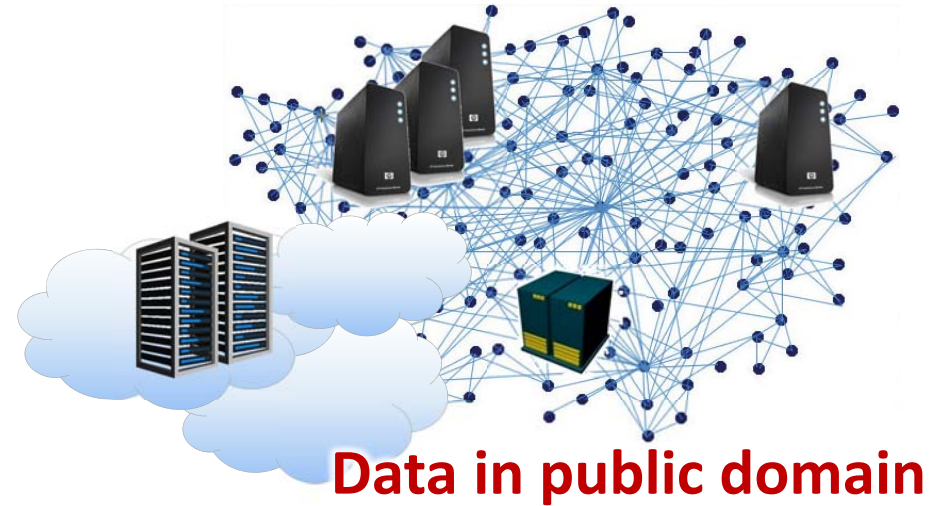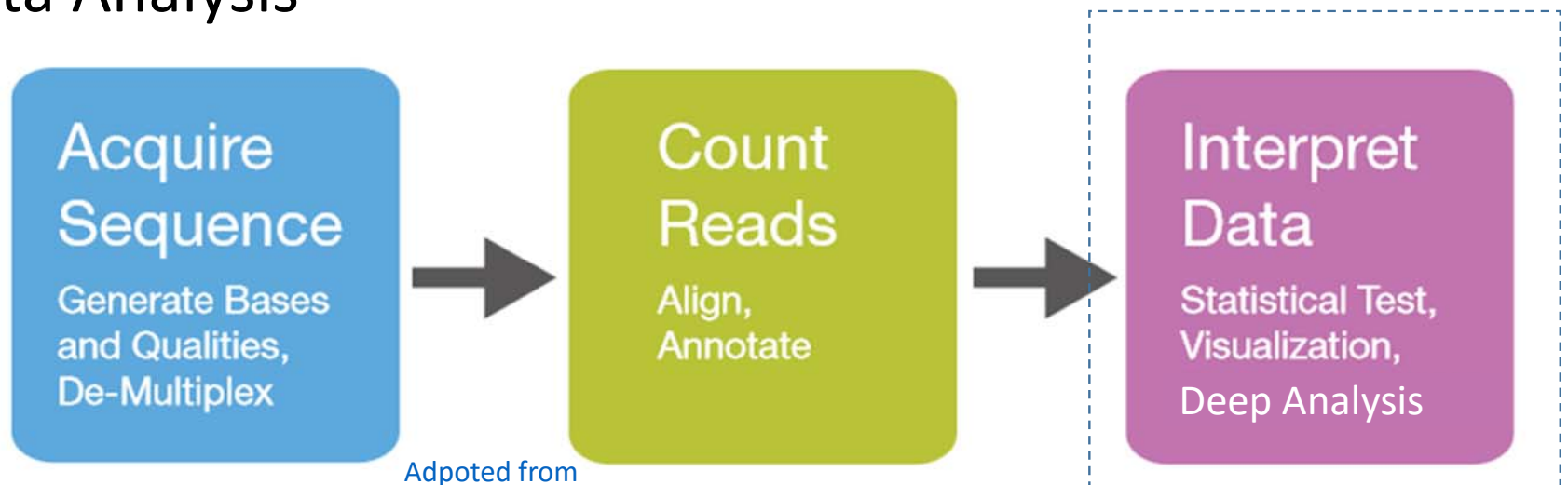
IIS, Academia Sinica

# High-throughput Methods

Biology Lab

Bioinfo Lab

Tech Core

**Hybridization-based Methods**

**Microarray**

**Data in public domain**

**Sequencing-based Methods**

Roche GS-FLX

Life Technologies SOLiD

Illumina HiSeq

Life Technologies Ion Torrent

PacBio

# Microarray Data Analysis

**Acquire Probe Intensity**

Generate labeled cRNA,
hybridization
Scan slide
Report probe intensity

**Interpret Data**

Statistical Test,
Visualization,
Deep Analysis

# RNASeq Data Analysis

**Acquire Sequence**

Generate Bases
and Qualities,
De-Multiplex

**Count Reads**

Align,
Annotate

**Interpret Data**

Statistical Test,
Visualization,
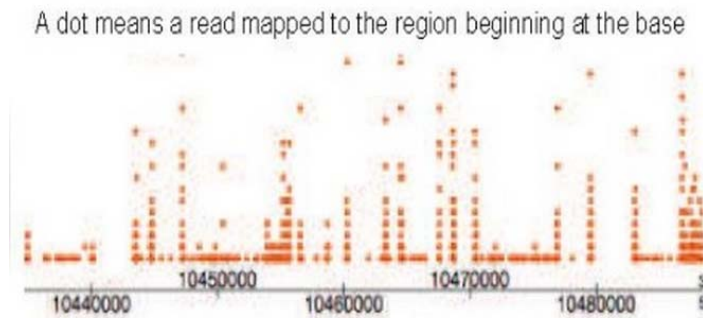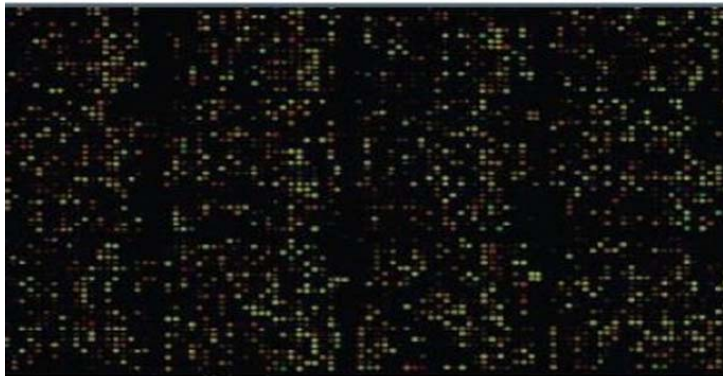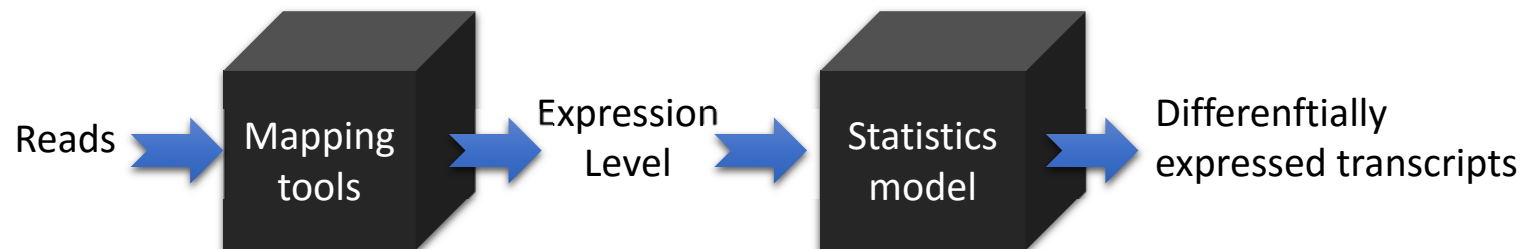Deep Analysis

Adpoted from
http://www.illumina.com/documents/products/datasheets/datasheet_rnaseq_analysis.pdf

# Analog signal vs Digital signal



A dot means a read mapped to the region beginning at the base

10440000   10450000   10460000   10470000   10480000

http://www.slideshare.net/ueb52/uebuat-bioinformatics-course-session-23-vhir-barcelona

Reads → **Mapping tools** → Expression Level → **Statistics model** → Differenftially expressed transcripts

# A typical RNA-Seq experiment



Prepare RNA samples

QC

Prepare sequencing libraries

QC

Millions of shot-gun reads

QC
Data preprocessing

Map to reference

Read mapping
…computing intensive jobs !

Convert to expression level

Intensive analysis to Interpretate Biological Meanings

mRNA

RNA fragments

or

cDNA

EST library with adaptors

ATCACAGTGGGACTCCATAAATTTTTCT
CGAAGGACCAGCAGAAACGAGAGAAAAA
GGACAGAGTCCCCAGCGGGCTGAAGGGG
ATGAAACATTAAAGTCAAACAATATGAA
……

Short sequence reads

ORF
Coding sequence

Exonic reads

Junction reads

poly(A) end reads

Mapped sequence reads

**Base-resolution expression profile**

RNA expression level

Nucleotide position

Nature Reviews | Genetics

http://www.nature.com/nrg/journal/v10/n1/full/nrg2484.html

# FastQ format

- Start with "@"
- Four lines: "+" w/ or w/o seq head, quality scores

| seq head | @EAS139:136:FC706VJ:2:5:1000:12850 1:N:18:ATCACG |
| --- | --- |
| seq letters | ACTTCAGGAGATTGTACATTTAGAGACAAAAAAAA |
| + | + |
| quality score | BBBBCCCC?<A?BC?7@@???????DBBA@@@@A@@ |

FASTQ files from CASAVA-1.8 Should have the following READ-ID format:

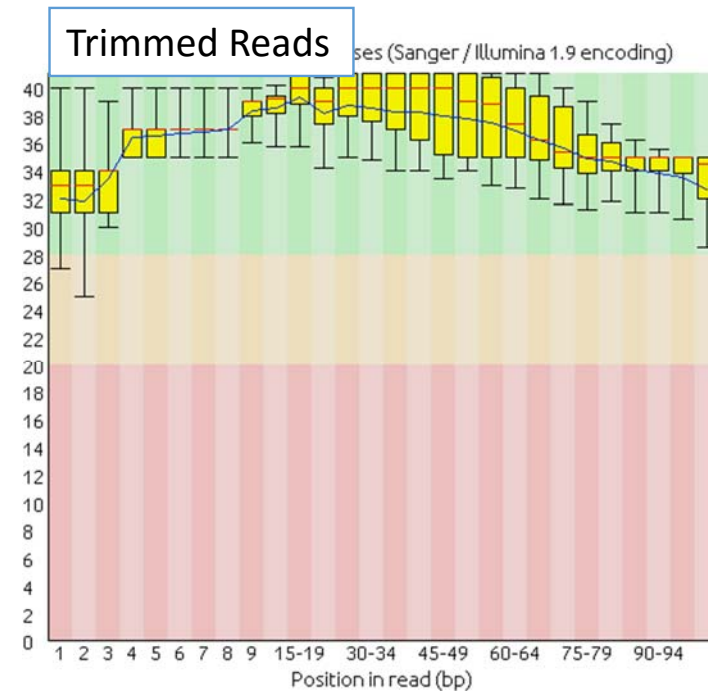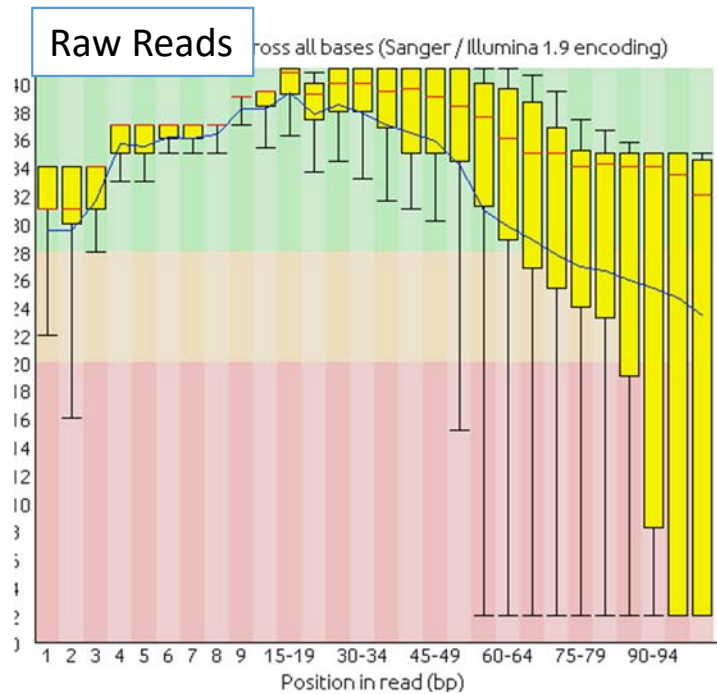@<instrument>:<run number>:<flowcell ID>:<lane>:<tile>:<x-pos>:<y-pos>
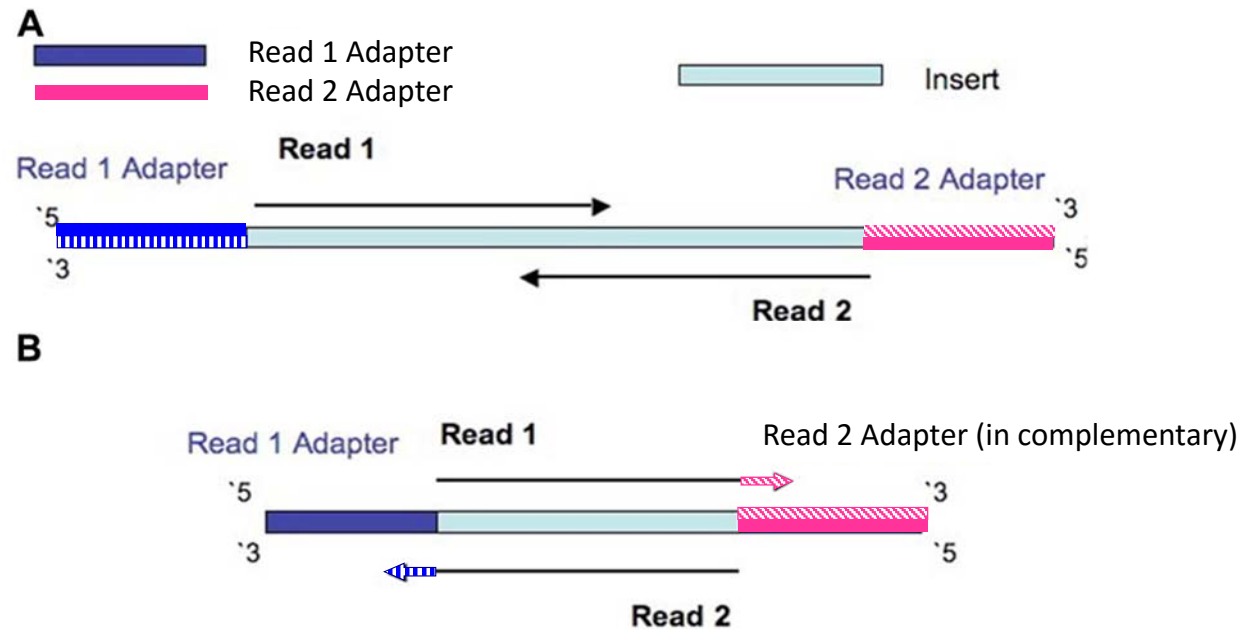
<read>:<is filtered>:<control number>:<index sequence>

# Read preprocessing

| Phred Quality Score | Probability of incorrect base call | Base call accuracy |
|---|---|---|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.9% |
| 40 | 1 in 10,000 | 99.99% |

• Trimming: by base quality



Raw Reads

Trimmed Reads

http://www.bioinformatics.babraham.ac.uk/projects/fastqc/

# Read preprocessing

- Trimming: adapter contamination

# Expression Level
## by Gene or by Transcript?



Read (pairs)

Reference Genome
Seq :fasta file
Annotation : gtf / gff file

mapper

splice-aware mapper

exprs level by Gene

exprs level by Transcripts & Gene

Read alignment
(maps reads to a reference genome)

TopHat    STAR

Exon A | Exon B | Exon C
Processed mRNA

Exon A | Exon B | Exon C
Genomic mapping and splicing alignment

Post-alignment quality control

RNA-SeQC

Additional quality assessment,
e.g., duplication rates, GC bias,
rRNA contamination

ngs.plot

Visualization of reads across
genes—TSSs to TESs—
to assess 3′ bias

Quantification and analyses

HTSeq

Transcript-level summary counts
(e.g., Bland-Altman plot)

Cufflinks

Transcriptome assembly,
isoform expression, splice
variants and differential analysis

MISO

Splicing events quantification
and differential analysis

http://www.nature.com/neuro/journal/v17/n11/fig_tab/nn.3816_F1.html

# Other issues

- Stranded or not?

- PolyA tailed or rRNA depletion?

- Have reference genome? Novel transcripts? Fusion transcripts?

- Special protocols that need extra bioinformatical works?

- Trimmed read length? Low complexity repeats? Other sources of contamination?

# Normalization is a Necessary Evil

- Between samples:

  Initial Input ; Volume of Reads

  

  Library 1     Library 2

- Within sample:

  length effect

  

  seq1

  seq2

- Count the mapped read number, normalized to library size

  cpm: count per million reads

- Count the mapped read number, normalized to BOTH library size and (target seq) length

  ✓TPM: transcripts per million reads

  ✓RSEM: RNA-Seq by Expectation-Maximization

  ✓RPKM: reads Per kilobase of exon per million mapped reads

  ✓FPKM: fragments per kilobase of exon per million fragments mapped

http://www.slideshare.net/mikaelhuss/rnaseq-differential-expression-analysis

# How many replicates does the experiment design includes?

- Theoretically………… **BUT!** in reality ……………

- Borrowing information among genes to get better estimates.
- Count-based model
  - edgeR, DESeq ….. etc.
  - Use "read count" (or estimated count from RSEM) and enforced a normalization model to fit data to the statistic assumption
  - Want to provide an analysis with statistic power
- Programs like SAMSeq (rank-based model, only applicable for large replicates) and limma are fine with continuous values (like FPKM). Limma takes more cares about weak mean-variance relationship (stabilizing variation).

# The Usage

Demo: http://molas.iis.sinica.edu.tw/grch38/

# 多重體學線上分析平台
# Multi-Omics onLine Annotation System (MOLAS)



To view and analyse your RNASeq experiment

# All you need is an expression file

Input file

- A tab-delimited text file generated by other software (e.g. cufflink, EdgeR, RSEM) in ensembl transcript id (grch38 and grcm38)



| #tracking_id | GA120-2_0 | GA120-3_0 |
|---|---|---|
| ENST00000591062 | 0 | 0.159246 |
| ENST00000376259 | 0 | 3.96794 |
| ENST00000235878 | 0.287651 | 0 |
| ENST00000299596 | 0.0300576 | 0.0146675 |
| ENST00000625158 | 6.08204 | 7.03465 |
| ENST00000321949 | 4.24507 | 4.28616 |
| ENST00000258484 | 0 | 6.00768 |
| ENST00000625157 | 0.0134854 | 0.00783917 |
| ENST00000321944 | 6.44635 | 5.25123 |
| ENST00000321945 | 0.907242 | 1.13444 |

Read (cleaned)

Reference genome

fasta gtf

Splice aware mapper

# GTF: the Gene Tranfer Format



1	ensembl_havana	transcript	4344146 4360314 .	-	.	gene_id "ENSMUSG000000259 00"; gene_version "6"; transcript_id "ENSMUST00000027032"; transcript_version "5"; gene_name "Rp1"; gene_source "ensembl_havana"; gene_biotype "protein_coding"; transcript_name "Rp1-001"; transcript_source "ensembl_havana"; transcript_biotype "protein_coding"; tag "CCDS"; ccds_id "CCDS14804";

MOLAS compatible GTF

grch38
grch37

grcm38

converter

GTF

transcript ID

GO

KEGG

HGNC

ensembl

# New Submssion

# New Submission

# Project Profile



This project is a transcriptome study on
grch38 reference genome (transcripts #:197523,library#:2)

## Project Info

**Project Name** grch38 demo (limit to 50 words)

Brief on this Project [?] :

grch38 demo

Upload an website logo (image file in jpg,gif,or png format)

選擇檔案 未選擇任何檔案
[?]

Name of Sub-directory: http://molas.iis.sinica.edu.tw/ grch38 [?]

Contact E-mail as Account: molas.iis@gmail.com [?]

Password: •••• [?]

Open to Public:            ⦿Yes
                           ○No ☐share this project data to my friends with this secret word: [?]

# Deployment Success

Dear User:

You have completed the submission. There are 8 libraries in your submission.
The whole system will be ready few minutes later after data deployment.
Please check the website below to start your journey on data analysis.

http://molas.iis.sinica.edu.tw/grch38    _    **Data Deployment Success!**

Thanks for your using our platform to deep your research.

MOLAS administrator

# Browse project and …….

# Fuzzy Search

# Pairwise Comparis

## Select libraries you v

Total:17764 input gene symbol. hit:5382 used. nohit:12382 excluded. [Heatmap]

Show [10 ▼] entries                                    Search: [          ]  [CSV] [PDF]

| Pathway name | Knumbers frequency | Background frequency | P-value ▲ | Genename associated to the term |
|---|---|---|---|---|
| Protein processing in endoplasmic reticulum | 128 out of 4307 knumbers | 128 out of 4598 knumbers | 0.00021 | ATF6 BCL2 ▶ |
| RNA transport | 120 out of 4307 knumbers | 120 out of 4598 knumbers | 0.00035 | AAAS CYFIP1 ▶ |
| Spliceosome | 111 out of 4307 knumbers | 111 out of 4598 knumbers | 0.00064 | BCAS2 CDC40 ▶ |
| Epstein-Barr virus infection | 146 out of 4307 knumbers | 147 out of 4598 knumbers | 0.00064 | AKAP8L AKT2 ▶ |
| Cell cycle | 105 out of 4307 knumbers | 105 out of 4598 knumbers | 0.00096 | ABL1 ANAPC11 ▶ |
| Parkinson's disease | 101 out of 4307 knumbers | 101 out of 4598 knumbers | 0.00126 | APAF1 ATP5A1 ▶ |
| Viral carcinogenesis | 131 out of 4307 knumbers | 132 out of 4598 knumbers | 0.00160 | ACTN3 ACTN4 ▶ |

## Pairwise Comparison

1. Select the data for comparison
   Present grouping:

| Pool | Dataset |
|---|---|
| pool a: | sample_1,sample |
| pool b: | sample_5 |

2. Apply Data Filter
   ⦿ Summation of expression levels by PoolA and PoolB  [>=
   ○ PoolA expression level [>= ▼] [     ]  PoolB express

3. Set the comparing scheme
   fold change cutoff: [>= ▼] [3]    expression pattern: [

4. Select Analytic Approach:
   ○ Show Gene List
   ○ Calculate GO term enrichment    default p value cutoff [0.1 ▼]
   ⦿ Calculate KEGG pathway enrichment
   ○ Draw heatmap with 2D clustering
   ○ Map on Protein Network (Max 600 transcripts)

# Clustering

If some samples have similar properties, clustering can help group them together and perform gene expression profile analysis.

# Clustering Results

# KEGG Pathway

# Enrichment Analysis

Insert a list of interesting genes to see which pathway they are involved.

# KEGG Global View

KEGG Global View provide an overview picture of KEGG pathway of human (hg19) and mouse (mm10) organisms. You can investigate specific metabolic pathway by exploring each category.

# Demo

Hands on practice on MOLAS

- Build your own project
- Browse project and conduct a study

http://molas.iis.sinica.edu.tw/human_demo_grch38/

http://molas.iis.sinica.edu.tw/mouse_demo_grcm38/

內容設定

◯ 允許所有網站顯示彈出式視窗

◉ 不允許任何網站顯示彈出式視窗 (建議)

管理例外情況...

# What to do if you have no replicates?

## Suggestions from edgeR authors

- Be satisfied with a descriptive analysis, that might include an MDS plot and an analysis of fold changes. Do not attempt a significance analysis. This may be the best advice.

- Simply pick a reasonable dispersion value, based on your experience with similar data, and use that for DE detection
  - In edgeR (empirically):
    - 0.4 human data (genetically "not" identical)
    - 0.1 for "genetically identical" strains of model organisms
    - 0.01 for technical replicates

- estimate dispersion from dataset reducing one (less critical) experiment factor

- estimate dispersion from a sizeable number of control transcripts that should not be DE if there exists

edgeR paper http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2796818/
menu http://www.bioconductor.org/packages/release/bioc/vignettes/edgeR/inst/doc/edgeRUsersGuide.pdf

# Limitations

- Assumption of "Uniformity" of all expressed transcripts may not always true

- Uncertain problems in mapping
  - Transcripts length issue
  - Redundance seq in genome
  - Reference is never a perfect match to the actual biological source of RNA being sequenced

- Reference  & no Reference

- Lag in analytic tools.

- No single robotic analylsis scheme fits all kind of needs

- Cost !!

Shu-Hwa Chen

Institute of Information Science
Academia Sinica, Taiwan
2016/10/27

# Epigenetic Modification

- Epigenetic Modification: Reversible modifications on genome components to affect gene expression without changing the DNA sequence



Histone acetylation, phosphorylation, methylation

DNA methylation

miRNA

piRNA

lncRNA

mRNA stability and translation

**One Genome**

**Many Phenotypes**

Adopted from McEwen BS et al., Nature Neuroscience 18, 1353–1363 (2015)

# Methylated Cytosine: the Fifth Base

The most common and stable epigenetic marks in nucleotide level



C ≡ G

C$^m$ ≡ G

- Involved in
  - Genomic imprinting
  - Cell Fate Determination / Reprogramming
  - Transposon genes silencing

- In vertebrates, 1-6% of genomic cytosine are methylated
- In plants, the proportion of methylated cytosine is even higher
- But……..

# Whole Genome Shotgun Bisulfite Sequencing



Reproduced and modified from Fig 1 in Curr Protoc Nucleic Acid Chem (2008) Chapter 6:Unit 6.10.

# Mapping BS-Seq Reads to Reference Genome

Reference genome

The Bisulfite converted,
PCR amplified library

5' –   T C G C C G G T A C -
       a g c g g c c a t g

indexed database

(C ➜ T/ G➜A)

Indexed in 3 bases/ wild card method

**BS Seq mappers**

T T G C T G G T A T
a a c g a c c a t a

T C G C C G G T A T
a g c g g c c a t a

A T T C C A G G A G C T C G C C G G T A C C T C A C C A
A G G A G C T T G C C G G
G G A G C T T G C T G G T A C C T
G A G C T C G C C G G T A C C T C A C
A G C T C G C C G G T A C C T C A C C A A T A
C T T G C T G G T A C C T

**Reference**

Short
Reads

# Difficulty to Access BS Seq Data/ Methylome

- Complicated Contents

**By Context**

-CG-     -CHG-    -CHH-

H=A, T or C

**By Location**

Gene1    Gene2    Gene3

- Promoter
- Gene Body

- Visualization

Methylated CG island

Chromosome–wide View of DNA Methylation Distribution

# The Workflow

# TEA
# The epigenomic platform for Arabidopsis

**Reference Genome**

A T T C C A G G A G C T C G C C G G T A C C T C A C C A

**Reads**

AGGAGTTTGTCGG
GGAGTTTGCTGGTATT
GAGTTCGTCGGTATTCAT
AGTTCGTCGGTATTTATTAATA
TTTGTTGGTATTT

- Type: **CG**
- Total observation (Read depth): 5
- Methylated C: 2, Unmethylated C: 3
  ➔ score of this C: 2/5 = 0.4

- Type: **CHH**
- Total observation (Read depth): 4
- Methylated C: 0, Unmethylated C: 4
  ➔ score of this C: 0

- Type: **CHG**
- Total observation (Read depth): 5
- Methylated C: 3, Unmethylated C: 2
  ➔ score of this C: 0.6

**Reference Genome**

**Reads**

- Scored gene / promoter: # observed bases >=5

**By Context**  **By Location**

Average DNA methylation level in promoter or gene body $= \dfrac{\sum\limits_{i \in X} c_i}{\sum\limits_{i \in X} 1}$  (1.2)

$X$ = promoter or gene body

**Reference Genome**

ATTCCAGGAGCT C G C C GGTA C CTCACCA

**Reads**

AGGAGTTTGTCGG
GGAGTTTGCTGGTATTT
GAGTTCGTCGGTATTTCAT
AGTTCGTCGGTATTTTATTAATA
TTTGTTGGTATTT

- Observed event for each C: >=4
- Scored gene / promoter: # observed bases >=5
- Supporting Mapper: BS-Seeker2 and Bismark

**mtable**

| gene_id | pmt_CG | gene_CG | pmt_CHG | gene_CHG | pmt_CHH | gene_CHH |
|---------|--------|---------|---------|----------|---------|----------|
| AT1G01010 | 0.011463 | 0.053009 | 0.010000 | 0.011635 | 0.021765 | 0.012631 |
| AT1G01020 | 0.000000 | 0.081519 | 0.006957 | | 0.003614 | 0.007521 |
| AT1G01030 | 0.005385 | 0.012800 | | | 0.003116 | 0.016939 |
| AT1G01040 | 0.011200 | | | 0.015773 | 0.016944 | 0.011699 |
| AT1G01046 | 0.765250 | 0.385000 | 0.022500 | 0.058750 | 0.014325 | 0.047727 |

The Methylation Landscape

**Input/ Mapping Report**

Read depth (C+T) >=4

Assign the methylation level m=NaN

Calculate the methylation level m= (C)/ (C+T)

Classify Cs by the sequence context

GTF

Classify Cs by locations (Promoter/Gene) & Type of Methylation (CG/CHG/CHH)

C with numeric value m >=5

Assign the methylation score to NaN

Calculate the methylation score in average

**mtable**

**Inputs**

■ **BS-Seq mapping report**

➤ **CGmap** from BS- Seeker2 /
➤ **CX_report.txt** from Bismark
➤ Or an equivalent from other BS-Seq mappers

■ **GTF of the reference genome**

| Gene_id | pmt_CG | gene_CG | pmt_CHG | gene_CHG | pmt_CHH | gene_CHH |
|---------|--------|---------|---------|----------|---------|----------|
| AT1G01010 | 0.005 | 0.068448 | 0.00375 | 0.028333 | 0.004739 | 0.024981 |
| AT1G01020 | 0.012353 | 0.092468 | 0.013182 | 0.015667 | 0.013614 | 0.019084 |

# TEA Website

Demo site: http://tea.iis.sinica.edu.tw/

Demo site: http://symbiosis.iis.sinica.edu.tw/tea/molas.html

# Project Summary

## Project Briefs

Datasets from DOMAINS REARRANGED METHYLTRANSFERASE3 controls DNA methylation and regulates RNA polymerase V transcript abundance in Arabidopsis study http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4311829/
Project Name: Demo published Arabidopsis dataset

There are 5 datasets uploaded to build this project. We summarized the mapping conditions in below:

| Sample Label | Uploaded IDs | Mapped IDs | mapped in tair10 geneid |
|---|---|---|---|
| Col_1 | 33602 | 100.0% (33602/33602) | 100.0% (33602/33602) |
| Col_2 | 33602 | 100.0% (33602/33602) | 100.0% (33602/33602) |
| drm2 | 33602 | 100.0% (33602/33602) | 100.0% (33602/33602) |
| drm3 | 33602 | 100.0% (33602/33602) | 100.0% (33602/33602) |
| nrpe1 | 33602 | 100.0% (33602/33602) | 100.0% (33602/33602) |

Poor ID mapping rate?!

Check the gtf version

# Project Summary

We further summarized the number of analyzable genes/promoters for different methylated C sequence contexts each sample :

| Sample Label | CG | | CHG | | CHH | |
|---|---|---|---|---|---|---|
| | promoter | gene | promoter | gene | promoter | gene |
| Col_1 | 28260 | 33387 | 28252 | 33437 | 28290 | 33485 |
| | 84.0% | 99.0% | 84.0% | 99.0% | 84.0% | 99.0% |
| Col_2 | 28233 | 33342 | 28228 | 33390 | 28281 | 33443 |
| | 84.0% | 99.0% | 84.0% | 99.0% | 84.0% | 99.0% |
| drm2 | 28160 | 33207 | 28137 | 33222 | 28207 | 33320 |
| | 83.0% | 98.0% | 83.0% | 98.0% | 83.0% | 99.0% |
| drm3 | 28183 | 33244 | 28160 | 33276 | 28191 | 33321 |
| | 83.0% | 98.0% | 83.0% | 99.0% | 83.0% | 99.0% |
| nrpe1 | 28291 | 33424 | 28288 | 33462 | 28326 | 33508 |
| | 84.0% | 99.0% | 84.0% | 99.0% | 84.0% | 99.0% |

Missing Data ?!

Check the (1) read mapping rate (2) throughput

# Gene Central View

## AT5G27150: NHX1

**Gene: NHX1**

### Gene Central View

| NHX1 Sodium/hydrogen exchanger 1 [Source:UniProtKB/Swiss-Prot;Acc:Q68KI4] | |
|---|---|
| Ensembl ID | Gene_Biotype |
| AT5G27150 | protein_coding |
| Synonym/ prev Symbol | chromosome location |
| | **ch5**: 9,553,438-9,557,513 forward strand. |

### The methylation level of NHX1 in all libraries

**Layout 1: by sequence type**    Layout 2: by location

### Layout 1
Main categories in methylC sequence contexts (CG/CHG/CHH)

| Methylation Level | | | | | | |
|---|---|---|---|---|---|---|
| AT5G27150 | pmt_CG | gene_CG | pmt_CHG | gene_CHG | pmt_CHH | gene_CHH |
| Col_1 | 0.380513 | 0.366392 | 0.303947 | 0.013275 | 0.262383 | 0.017167 |
| Col_2 | 0.357317 | 0.344938 | 0.285128 | 0.012076 | 0.228465 | 0.012457 |
| drm2 | 0.375405 | 0.386733 | 0.186757 | 0.007299 | 0.015115 | 0.009421 |
| drm3 | 0.370256 | 0.362956 | 0.305405 | 0.015357 | 0.19905 | 0.019717 |
| nrpe1 | 0.301026 | 0.32378 | 0.018378 | 0.012773 | 0.021895 | 0.012926 |

The methylation level of NHX1 in all libraries

Layout 1: by sequence type | Layout 2: by location

## Layout 1
Main categories in methylC sequence contexts (CG/CHG/CHH)

| Methylation Level | | | | | | |
|---|---|---|---|---|---|---|
| AT5G27150 | pmt_CG | gene_CG | pmt_CHG | gene_CHG | pmt_CHH | gene_CHH |
| Col_1 | 0.380513 | 0.366392 | 0.303947 | 0.013275 | 0.262383 | 0.017167 |
| Col_2 | 0.357317 | 0.344938 | 0.285128 | 0.012076 | 0.226465 | 0.012457 |
| drm2 | 0.375405 | 0.366733 | 0.186757 | 0.007299 | 0.015115 | 0.009421 |
| drm3 | 0.370256 | 0.362956 | 0.305405 | 0.015357 | 0.19905 | 0.019717 |
| nrpe1 | 0.301026 | 0.32378 | 0.018378 | 0.012773 | 0.021895 | 0.012926 |



Measures of Methylation

Genome Browser

# Data Analysis Modules

http://tea.iis.sinica.edu.tw/tea/access_project.html

# Find Genes by Value

DMGs : Select differentially methylated genes by the interested methylation score



Threshold : Select genes by a cutoff value on the methylation score

# Gene List and Data Visualization

# *Questions?*

# Future Works

- A more sophisticate measures that highlight the pattern of methylation


- Multi-Omic Integration

# Sequencing Platforms



|  | ABI 3730xl Sanger Sequencing | 454 Life Sciences pyrosequencing | SOLiD + Illumina | Pacific Biosciences, Oxford Nanopore etc Single-molecule sequencing |
|---|---|---|---|---|
| **Length/read** | 800 bp | 400 bp | 100 bp | 20 000+ bp |
| **Reads/run** | 96 | 1 million | 2 billion | 5 million |
| **Bases/run** | 60 kbp | 400 Mbp | 500 Gbp | 100 Gbp |
| **Speed** | 10 years/HG | 1 month/HG | 1 day/HG | 10 min/HG |
|  | "old school" | "2nd gen" |  | "3rd gen" |

http://www.slideshare.net/COST-events/rnaseq-analysis-17037153?next_slideshow=1

# FastQ format

- Start with "@"
- Four lines:,, "+" w/ or w/o seq head, quality scores

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS....................................
.......................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX......................
...............................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII..................
.................................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.................
LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL....................................
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
|                        |   |       |                                    |           |
33                       59  64      73                                   104         126

0........................26...31.......40
                         -5....0........9.............................40
                               0........9.............................40
                                  3.....9.........................40
0.2......................26...31........41

S - Sanger        Phred+33,  raw reads typically (0, 40)
X - Solexa        Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64,  raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64,  raw reads typically (3, 40)
    with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
    (Note: See discussion above).
L - Illumina 1.8+ Phred+33,  raw reads typically (0, 41)
```

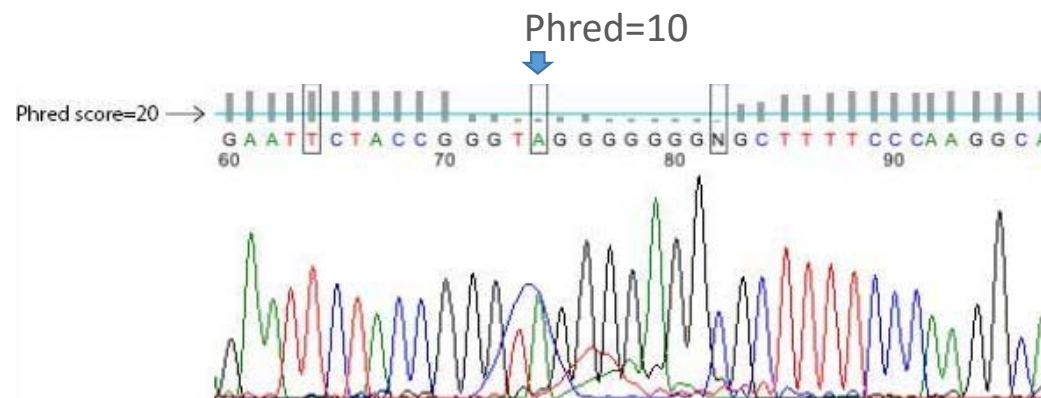## Phred quality scores are logarithmically linked to error probabilities

| Phred Quality Score | Probability of incorrect base call | Base call accuracy |
|---|---|---|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.9% |
| 40 | 1 in 10,000 | 99.99% |
| 50 | 1 in 100,000 | 99.999% |
| 60 | 1 in 1,000,000 | 99.9999% |

Phred=10

Phred score=20 →

G A A T T C T A C C G G G T A G G G G G G G N G C T T T T C C CA A G G C A
60                    70                   80                    90

http://en.wikipedia.org/wiki/Phred_quality_score