# Multi-Omics onLine Analysis System for Profiling Gene Expression
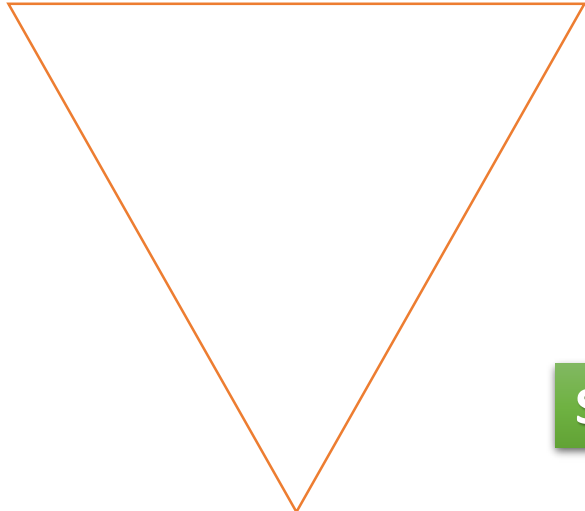


Life Science Library Training Course
2018/12

Chen, Shu-Hwa

IIS, Academia Sinica
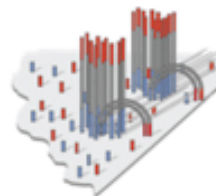
# High-Throughput Methods

**Data in the public domain**
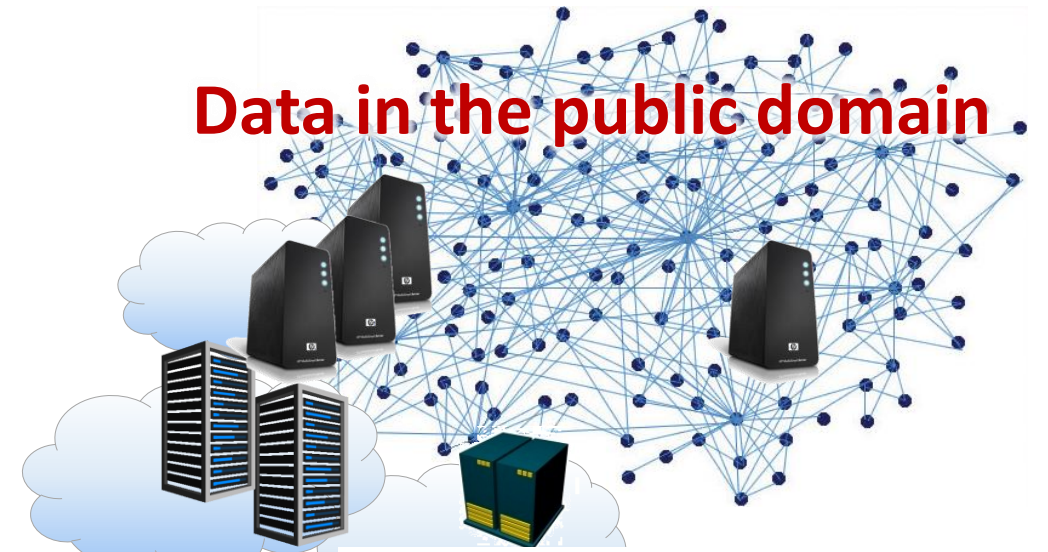
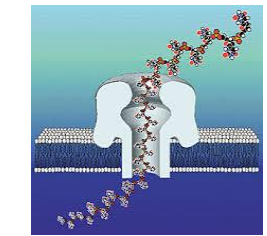Biology Lab       Bioinfo Lab

Tech Core

NCBI    Site map  |  All databases  |  Search

Sequence Read Archive

NIH   THE CANCER GENOME ATLAS
National Cancer Institute
National Human Genome Research Institute

https://www.ncbi.nlm.nih.gov/sra     https://cancergenome.nih.gov/

**Sequencing-based Methods**

**Microarray**

**Hybridization-based Methods**

PacBio

# Read in <u>FastQ</u> format

https://en.wikipedia.org/wiki/FASTQ_format

DON'T TRY TO OPEN a fastq file on your desktop PC

- 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@@ |

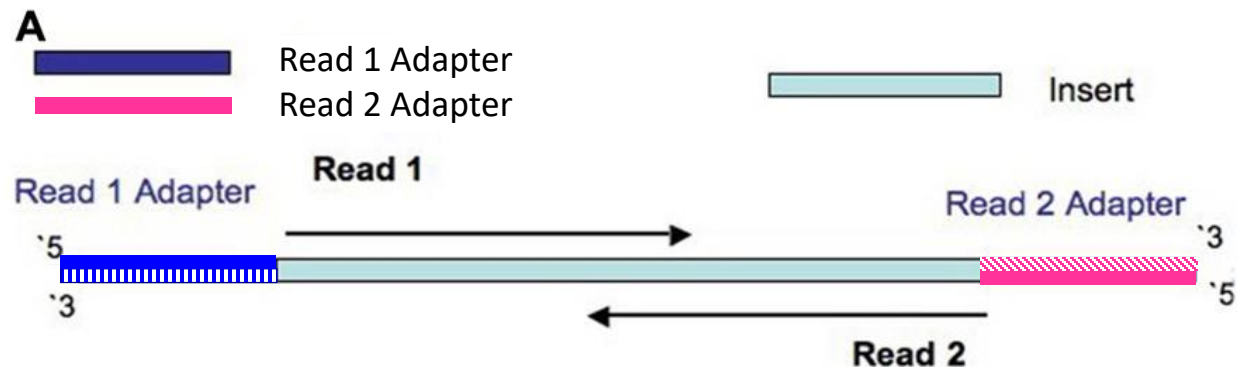Illumina reads

fastq files from a sequencer 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>

R1, R2 can  in separate fastq files
or sometimes in an interlanced fastq file:

@xxxxx:xxxx:xxxx:….. 1:N:0
……………
+
……………
@xxxxx:xxxx:xxxx:….. 2:N:0
……………
+
……………
@xxxxx:xxxx:xxxx:….. 1:N:0
……………
+
……………
@xxxxx:xxxx:xxxx:….. 2:N:0
……………
+
……………

# Illumina Technology



Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

Denaturation leaves single-stranded templates anchored to the substrate.

Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

A

Read 1 Adapter
Read 2 Adapter

Insert

Read 1 Adapter    **Read 1**

Read 2 Adapter

`5                                                                    `3

`3                                                                    `5

**Read 2**
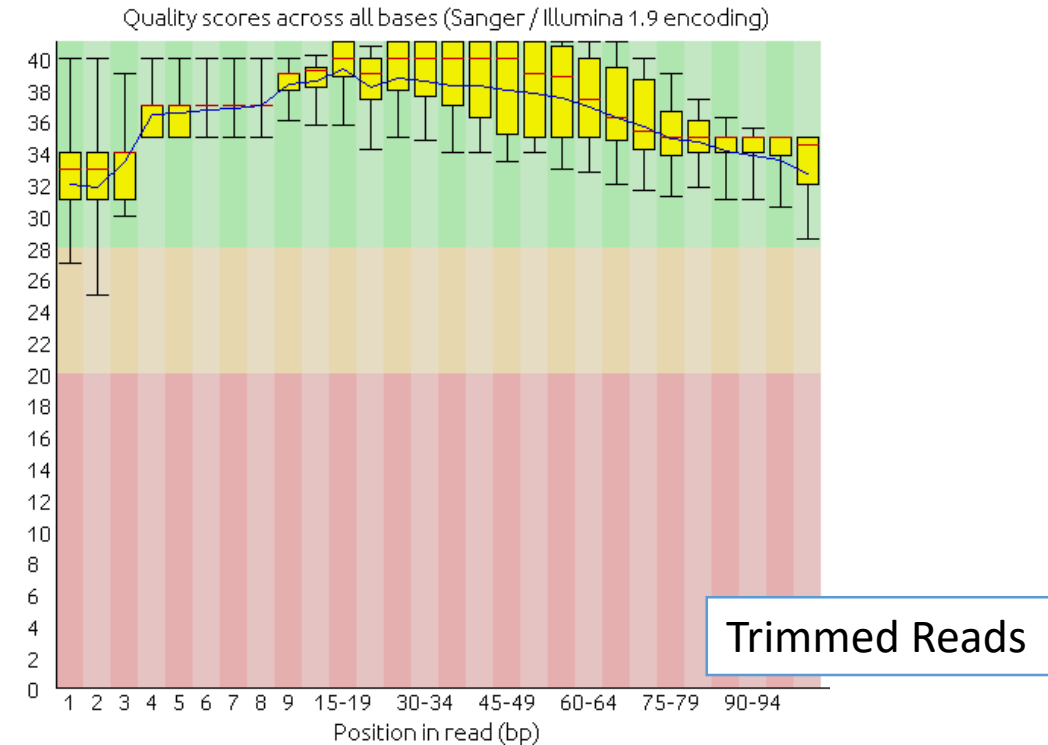
# Check the Quality of Reads

Trimming for base quality

| 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% |

Raw Reads

Trimming
by base quality

Trimmed Reads

# Read Preprocessing (optional):

Trimming for adapter contamination

# Mapping: the Options

## Mapping to a Known Genome

- Working on the target species:
  to profile the gene repertoire on
  a well-defined reference
  genome (fully sequenced and
  annotated) .

- Using a known genome close to
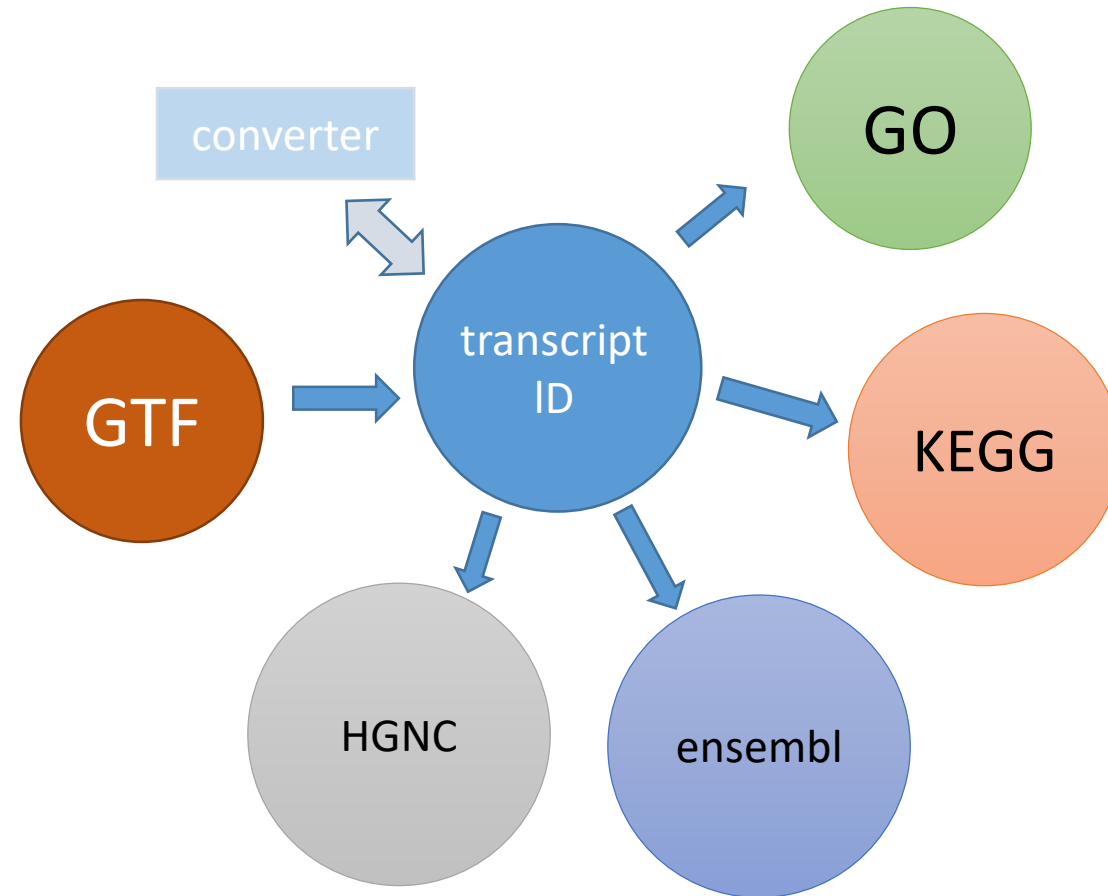  your sample.

## de novo Assembling + Mapping

- Create the reference
  (transcriptome or genome) by
  assembling.

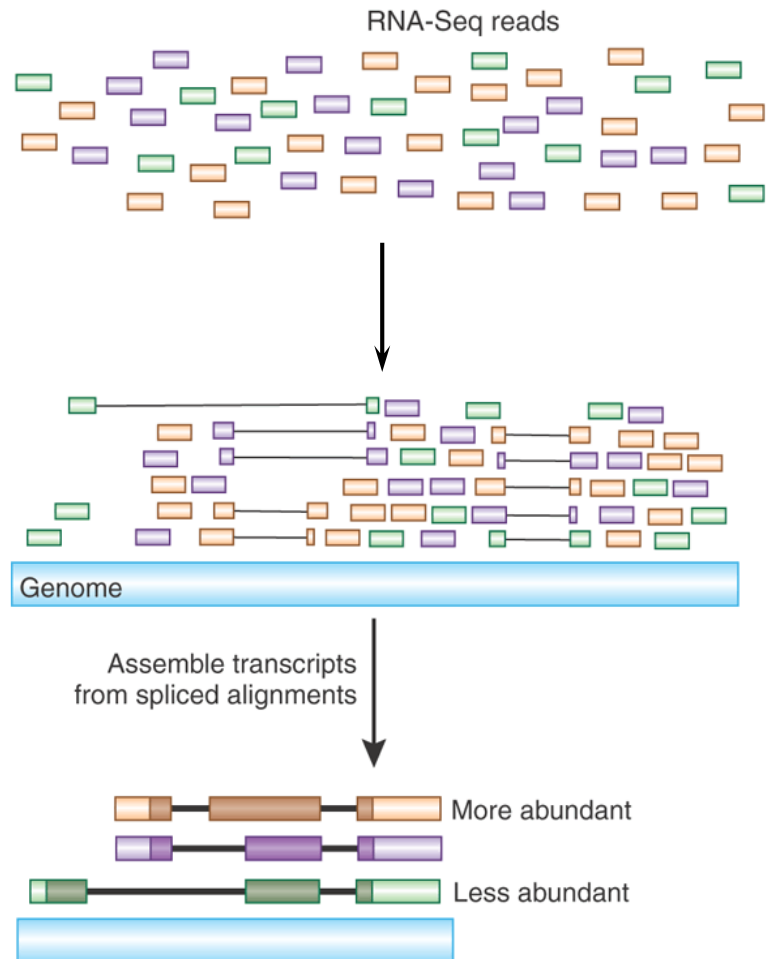- Annotate the new assembled
  reference.

- Map reads to the new
  assembled reference.

# 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 "ens
embl_havana"; transcript_biotype "protein_coding"; tag "CCDS"; ccds_id "CCDS14804";
```
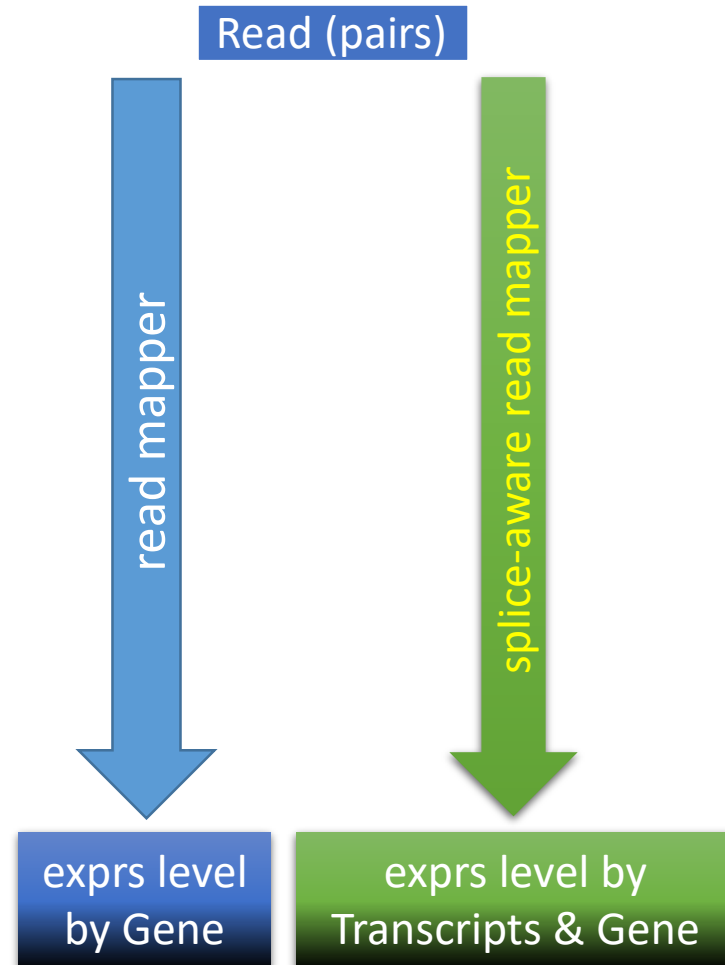
MOLAS compatible GTF



grch38
grch37

grcm38

converter

GTF

transcript ID

GO

KEGG

HGNC

ensembl

# Read Mapping

RNA-Seq reads



Genome

Assemble transcripts
from spliced alignments

More abundant

Less abundant

Reference Genome
- Seq: fasta file / prebuilt index
- Annotation : gtf / gff file

## Expression Level
by Gene or by Transcript?

Read (pairs)

read mapper

splice-aware read mapper

exprs level
by Gene

exprs level by
Transcripts & Gene

# Normalization is a Necessary Evil

- Between samples:

  Initial Input ; Volume of Reads

  

  Library 1: 12M reads    Library 2: 20M reads

- Within sample:

  transcript 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

# The Usage

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

# 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)

Read (cleaned)

Reference
genome

fasta
gtf

Splice aware mapper

| #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 |

# New Submission



MOLAS
Multi-Omics onLine Analysis System

Help     About Us

🏠 Home     📁 Browse Projects     ⊙ New Submission     🔑 Registered User Login     ➕ help

Human, grch37          Human, grch38          Mouse, grcm38

Upload expressed profiling in ○TPM / ⦿FPKM in tab file:

**Example dataset for download:**
For transcript grch37, grch38, grcm38
For genesymbol grch37, grch38, grcm38
For geneid grch37, grch38, grcm38

瀏覽...  未選擇檔案。

*Important: Please read this before submission

By transcript: ⦿EnsemblTranscriptID(ENSMUST00000000001 or ENST00000000233)

By gene:     ○genesymbol(ARF5)  ○EnsemblGeneID(ENSMUSG00000005320 or ENSG00000004059)

For human study ➡  ☐combined with DEMO DB data library.

Submit ✔     Clear All ⟳

# New Submission



MOLAS    About MOLAS    Browse Projects    New Submission    Check Submitted jobs

There are 208244 transcripts annotated in human genome,ensembl grch38.78. In MOLAS, 197912 transcripts are in the database ( transcripts of "small non-coding genes" are excluded. Link to Details)
197523 data entries are found in the uploaded file,in which 14 ensembl transcriptid (0.01%, 14/197523) can not mapped to MOLAS database.
197509 MOLAS database transcript id are mapped (99.8%, calculated by mapped id / molas id: 197509/197912)

FPKM file top 5 lines :

| | | | | operation |
| | | | | ◯ Modify FPKM Sample Name |
| #tracking _id | Sample_ 1 | Sample_ 2 | Sample_ 3 | Sampl e_4 |
|---|---|---|---|---|
| ENST00000380075 | 0 | 0 | 0.909464 | 1.0386 |
| ENST00000380071 | 320.788 | 208.653 | 269.647 | 421.71 |
| ENST00000380079 | 160.909 | 71.0702 | 63.7214 | 0 |
| ENST00000563164 | 11.2517 | 15.5313 | 7.45358 | 14.1989 |
| ENST00000563166 | 0 | 0 | 0 | 1.99288 |

Select library:

Present Selected:

| Dataset | operation |
|---|---|
| Sample_1, Sample_2, Sample_3, Sample_4 | ◯ modify ◯ delete |

Selecting Dataset:

☑Sample_1        ☑Sample_2        ☑Sample_3        ☑Sample_4        Update
                                                                    Reset

# 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



| About MOLAS | Browse Projects | **New Submission** | Check Submitted jobs |

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 .......

Functional Enrichment

Clustering

BLAST Search

Pathway View

Pairwise Comparison

Contig Information

Home | Full-text search on Annotation tables | Pairwise Comparison | Import Genelist | Clustering | KEGG GlobalView | Gene List Analysis

grch38 demo

# Fuzzy Search

# Pairwise Comparison

Select libra

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 ➤ |

Home | Full-text sear

**Dynamic com**

1. Select library:
Present grouping:

| Pool |
|---|
| pool a: |
| pool b: |

2. Select group:
◉ PoolA ∩ PoolB (#
   PoolA FPKM: [>= ▾]

◯ [PoolA Expressed only (F]
   PoolA FPKM: [>= ▾]

3. Select Analytic App
◯ Show Gene List
◉ Functional enrichm
◯ GO

[send] [reset]

# 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 canonical pathway-type overview of genes involved in a particular KEGG pathway.

# Demo

Hands on practice on MOLAS

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

http://molas.iis.sinica.edu.tw/human_grch38_demo/

# 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 claim a significance statistical analysis.
  - In edgeR (empirically): Simply pick a reasonable dispersion value, based on your experience with similar data, and use that for detecting differentially expressed transcripts.
    - 0.4 human data (genetically "not" identical)
    - 0.1 for "genetically identical" strains of model organisms
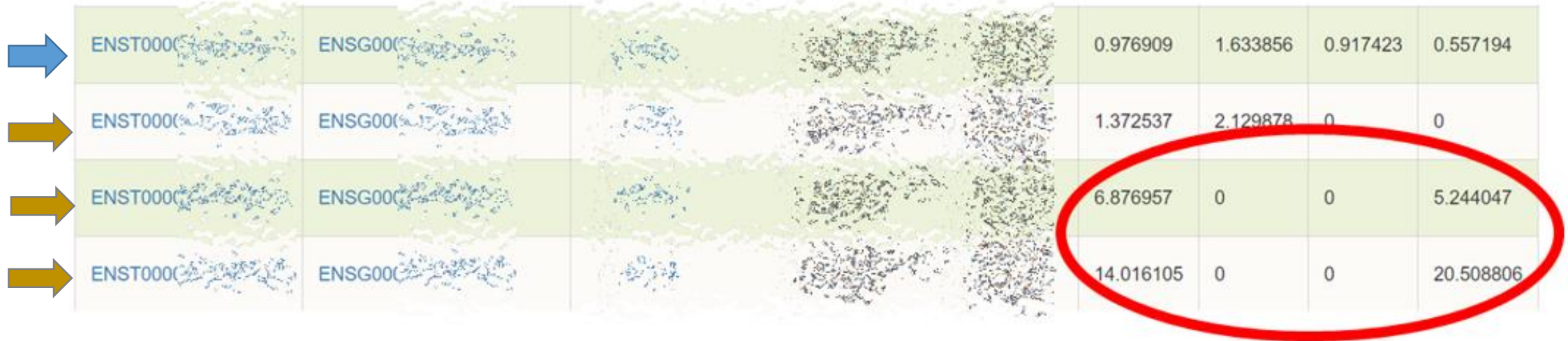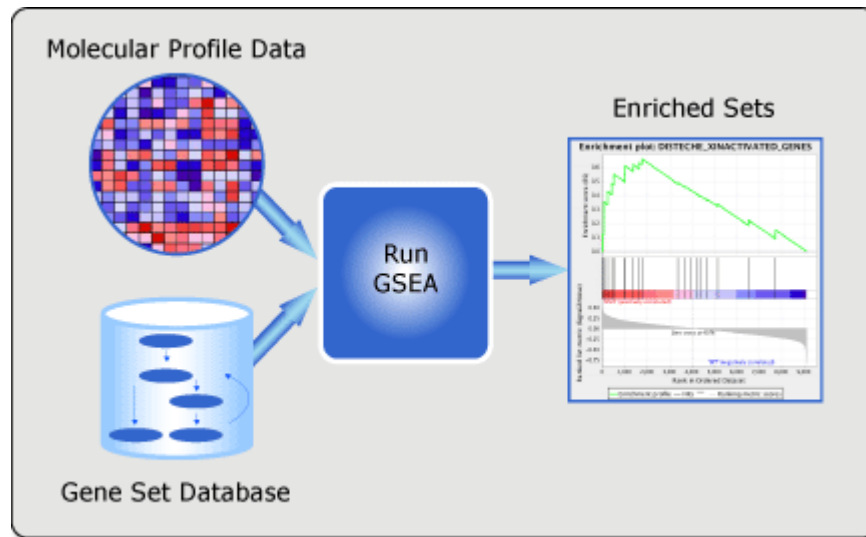    - 0.01 for technical replicates
  - Simulation data: NOISeq

https://f1000research.com/articles/5-1438/v2

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

# Genes with different transcripts…..

**Gene Set Enrichment Analysis**: a computational method that determines whether an a priori defined set of genes shows statistically significant, concordant differences between two biological states (e.g. phenotypes).
Ref: Subramanian, Tamayo, et al. (2005, PNAS 102, 15545-15550)

**the Molecular Signatures Database (MSigDB)**, a collection of annotated gene sets for use with GSEA software.
Ref: Liberzon, et al. (2011, Bionformatics), Liberzon, et al. (2015, Cell Systems)

**H** — **hallmark gene sets** are coherently expressed signatures derived by aggregating many MSigDB gene sets to represent well-defined biological states or processes.

**C1** — **positional gene sets** for each human chromosome and cytogenetic band.

**C2** — **curated gene sets** from online pathway databases, publications in PubMed, and knowledge of domain experts.

**C3** — **motif gene sets** based on conserved cis-regulatory motifs from a comparative analysis of the human, mouse, rat, and dog genomes.

**C4** — **computational gene sets** defined by mining large collections of cancer-oriented microarray data.

**C5** — **GO gene sets** consist of genes annotated by the same GO terms.

**C6** — **oncogenic gene sets** defined directly from microarray gene expression data from cancer gene perturbations.

**C7** — **immunologic gene sets** defined directly from microarray gene expression data from immunologic studies.
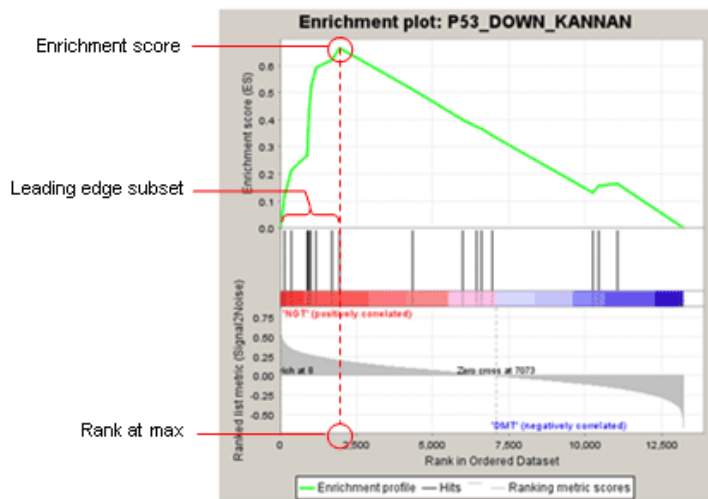
Fig 1: Enrichment plot: P53_DOWN_KANNAN
Profile of the Running ES Score & Positions of GeneSet Members on the Rank Ordered List

http://software.broadinstitute.org/gsea/doc/GSEAUserGuideFrame.html

........ the degree to which a gene set is overrepresented at the top or bottom of a ranked list of genes.

....... Gene sets with a distinct peak at the beginning (such as the one shown here) or end of the ranked list are generally the most interesting.
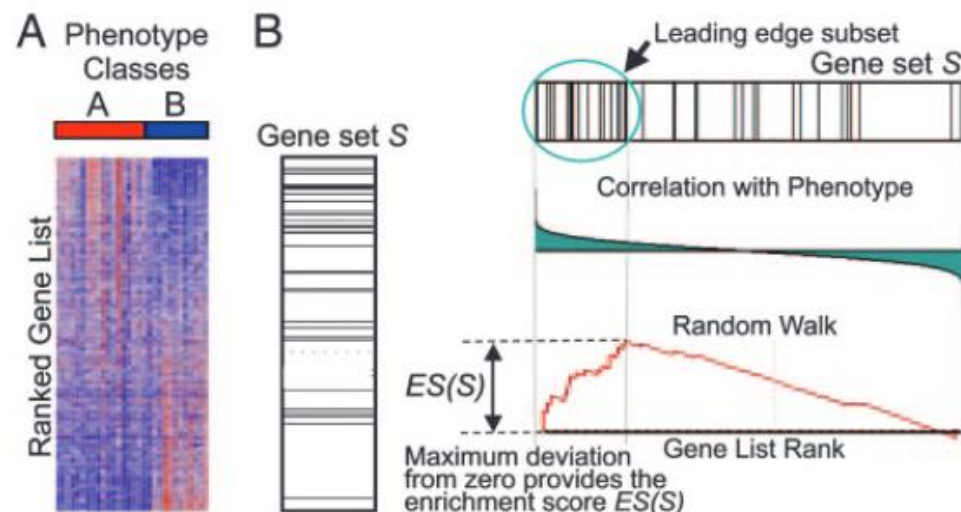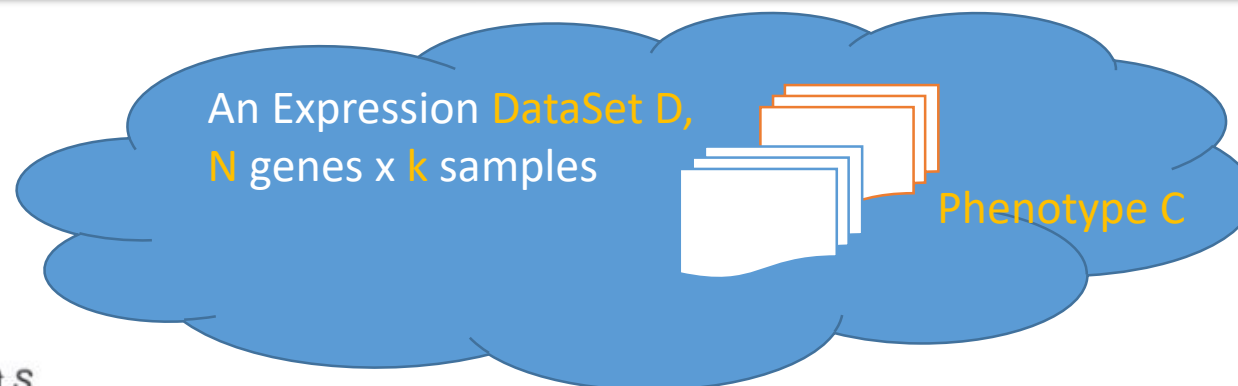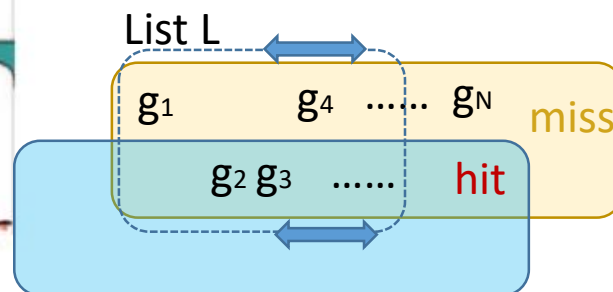
An Expression DataSet D, N genes x k samples

Phenotype C

Gene List L➔ in ranked order L={$g_1$,$g_2$,$g_3$, .......$g_N$} and the correlation to phenotype C, i.e., $r(g_1)=r_1$

List L

$g_1$   $g_4$ ...... $g_N$   miss

$g_2$ $g_3$ ......   hit

Gene Set S
(at least 15 members observed in the DataSet D)

$$P_{\text{hit}}(S, i) = \sum_{\substack{g_j \in S \\ j \leq i}} \frac{|r_j|^P}{N_R}, \quad \text{where } N_R = \sum_{g_j \in S} |r_j|^P$$

$$P_{\text{miss}}(S, i) = \sum_{\substack{g_j \notin S \\ j \leq i}} \frac{1}{(N - N_H)}.$$

ES: max ($P_{\text{hit}}$ - $P_{\text{miss}}$) from zero

a running-sum statistic

Fig. 1. A GSEA overview illustrating the method. (A) An expression data set sorted by correlation with phenotype, the corresponding heat map, and the "gene tags," i.e., location of genes from a set S within the sorted list. (B) Plot of the running sum for S in the data set, including the location of the maximum enrichment score (ES) and the leading-edge subset.

Tamayo, et al. (2005, PNAS 102, 15545-15550)

- **H** (hallmark gene sets, 50 gene sets) ❓
- **C1** (positional gene sets, 326 gene sets) ❓
  - by chromosome: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
- **C2** (curated gene sets, 4762 gene sets) ❓
  - **CGP** (chemical and genetic perturbations, 3433 gene sets) ❓
  - **CP** (Canonical pathways, 1329 gene sets) ❓
  - **CP:BIOCARTA** (BioCarta gene sets, 217 gene sets) ❓
  - **CP:KEGG** (KEGG gene sets, 186 gene sets) ❓
  - **CP:REACTOME** (Reactome gene sets, 674 gene sets) ❓
- **C3** (motif gene sets, 836 gene sets) ❓
  - **MIR** (microRNA targets, 221 gene sets) ❓
  - **TFT** (transcription factor targets, 615 gene sets) ❓
- **C4** (computational gene sets, 858 gene sets) ❓
  - **CGN** (cancer gene neighborhoods, 427 gene sets) ❓
  - **CM** (cancer modules, 431 gene sets) ❓
- **C5** (GO gene sets, 5917 gene sets) ❓
  - **BP** (GO biological process, 4436 gene sets) ❓
  - **CC** (GO cellular component, 580 gene sets) ❓
  - **MF** (GO molecular function, 901 gene sets) ❓
- **C6** (oncogenic signatures, 189 gene sets) ❓
- **C7** (immunologic signatures, 4872 gene sets) ❓

http://software.broadinstitute.org/gsea/msigdb/collections.jsp#C2

## Gene Set: GGCGGCA_MIR371

| | |
|---|---|
| **Standard name** | GGCGGCA_MIR371 |
| **Systematic name** | M15158 |
| **Brief description** | Genes having at least one occurence of the motif GGCGGCA in their 3' untranslated region. The motif represents putative target (that is, seed match) of human mature miRNA hsa-miR-371 (v7.1 miRBase). |
| **Full description or abstract** | |
| **Collection** | C3: motif gene sets<br>MIR: microRNA targets |
| **Source publication** | |
| **Exact source** | |
| **Related gene sets** | |
| **External links** | |
| **Organism** | Homo sapiens |
| **Contributed by** | Xiaohui Xie (Broad Institute) |
| **Source platform** | HUMAN_GENE_SYMBOL |
| **Dataset references** | |
| **Download gene set** | format: grp | text | gmt | gmx | xml |
| **Compute overlaps** ❓ | (show collections to investigate for overlap with this gene set) |
| **Compendia expression profiles** ❓ | Human tissue compendium (Novartis)<br>NCI-60 cell lines (National Cancer Institute) |
| **Advanced query** | Further investigate these 5 genes |
| **Gene families** ❓ | Categorize these 5 genes by gene family |
| **Show members** | (show 5 members mapped to 5 genes) |
| **Version history** | 6.0: Renamed from GGCGGCA,MIR-371 |

See MSigDB license terms here. Please note that certain gene sets have special access terms.

http://software.broadinstitute.org/gsea/msigdb/cards/GGCGGCA_MIR371.html

Shu-Hwa Chen

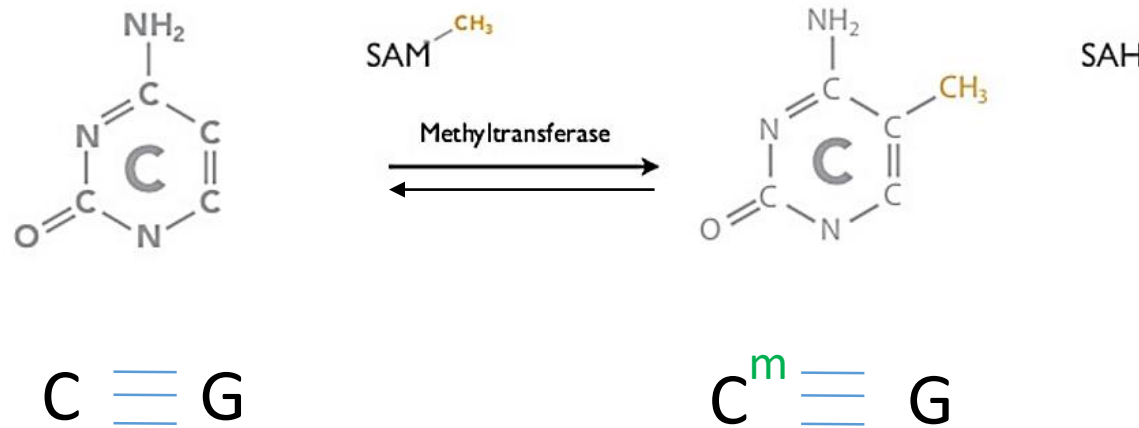Institute of Information Science
Academia Sinica, Taiwan
2018/12

# 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

# Methylated Cytosine: the Fifth Base

The most common and stable epigenetic marks in nucleotide level



- 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

# 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

# The Workflow

# TEA
# The epigenomic platform for Arabidopsis

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

**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
GGAGTTTGCTGGGTATT
GAGTTCGTCGGTATTCAT
AGTTCGTCGGTATTTATTAATA
TTTGTTGGGTATTT

- 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**

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**

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

By Context    By Location

$$\text{Average DNA methylation level in promoter or gene body} = \frac{\sum_{i \in X} c_i}{\sum_{i \in X} 1} \qquad (1.2)$$

X = promoter or gene body

**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**

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

| 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

# TEA Website

# 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 | 0.265128 | 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 |



Measures of Methylation

Arabidopsis thaliana TAIR10  5:9,534,459..9,585,489

Genome Browser

# Data Analysis Modules

# 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

線上基因概況分析平台

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

Raw Reads in a dozen of GB

DOCEXPRESS

FPKM/TPM

Upload to MOLAS

Estimate Expression Profiling in DOCEXPRESS

Unveil the biological secrets hidden behind the big biological data online
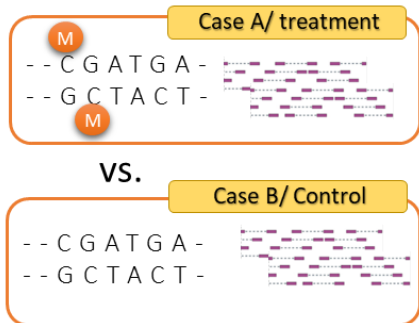
Analytic Platform

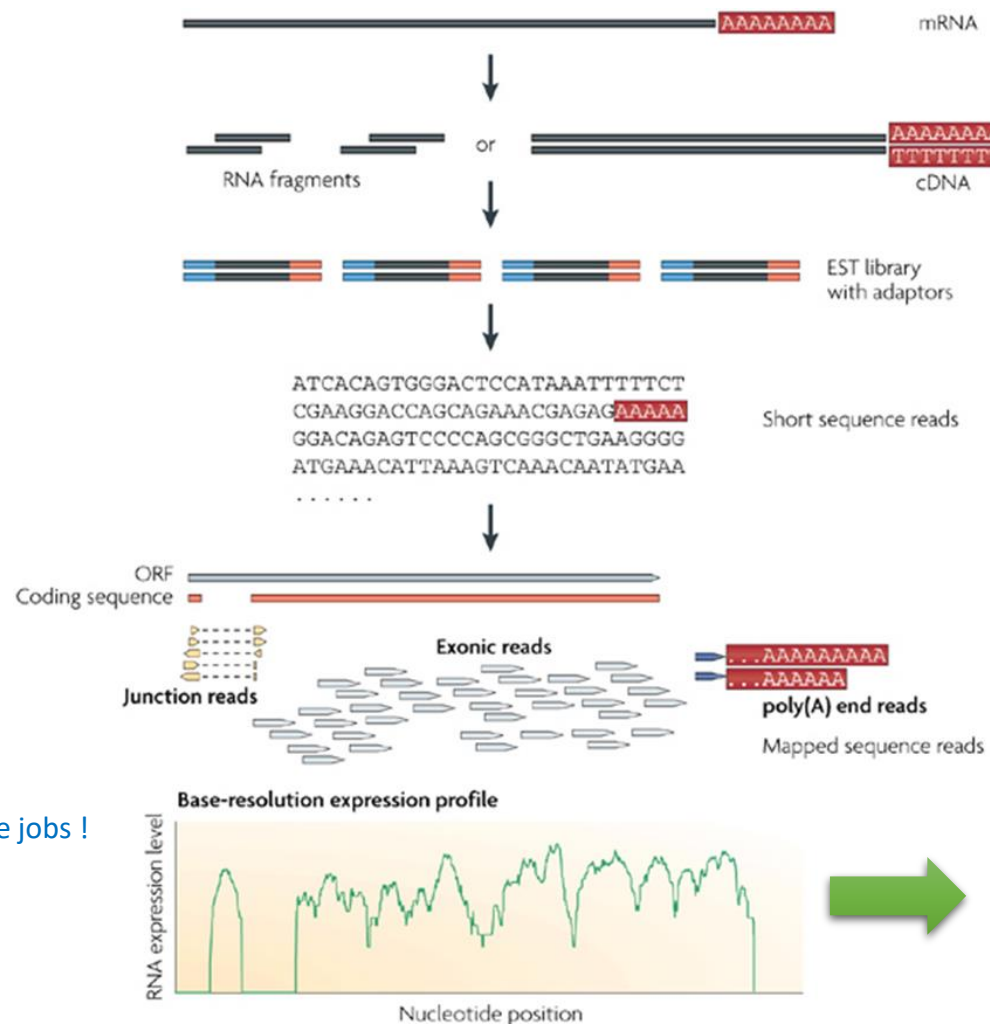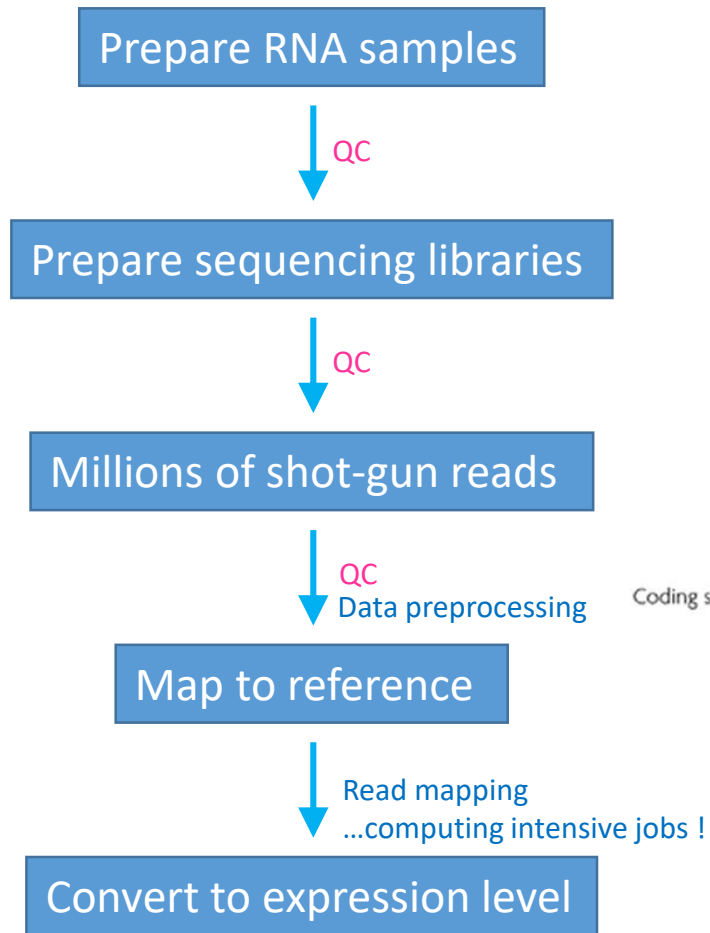Biological Meaning

Biological Big Data

# Questions?

# Other Issues ?

- Experiment Design? Biological Replicates >>> Technical Replicates
- Library Protocols:
  - 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?

# A Typical RNA-Seq Experiment



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