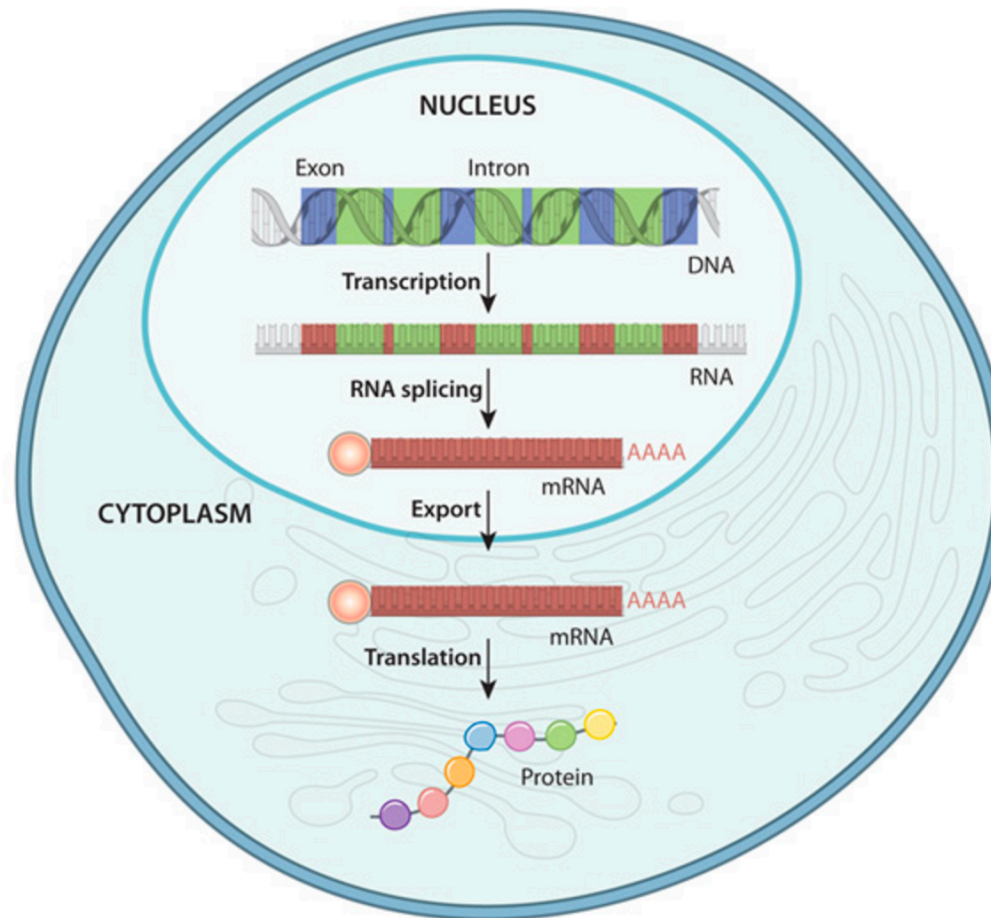# Gene Expression Analysis

P. Tang (鄧致剛); YM Yeh (葉元鳴)

Bioinformatics Center, Chang Gung University.
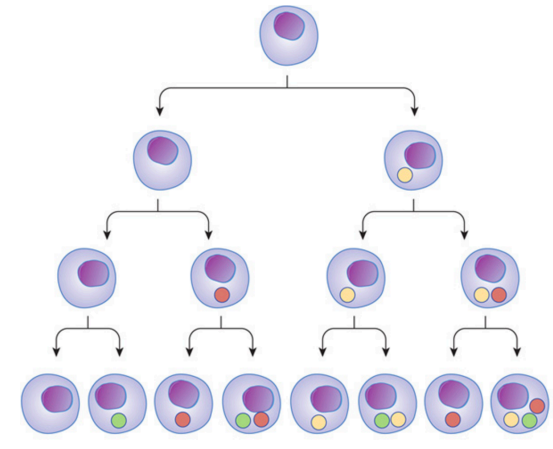
# Gene expression



**Figure 1: An overview of the flow of information from DNA to protein in a eukaryote**
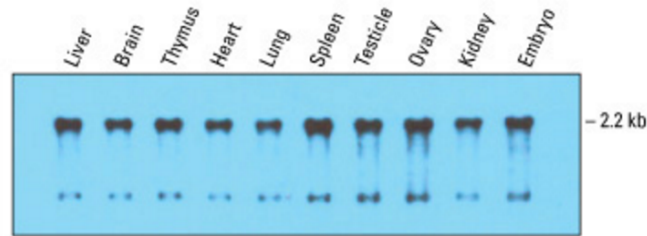First, both coding and noncoding regions of DNA are transcribed into mRNA. Some regions are removed (introns) during initial mRNA processing. The remaining exons are then spliced together, and the spliced mRNA molecule (red) is prepared for export out of the nucleus through addition of an endcap (sphere) and a polyA tail. Once in the cytoplasm, the mRNA can be used to construct a protein.

# (vs. DNA) Why RNA ?

- Differentially expressed – Functional studies
    - Different cell types (muscle cells, fibroblasts)
    - Environmental conditions (heat shock, nutrient deprivation)
    - Developmental phases (embryonic day 12)
    - Cell-cycle stages (S phase)
    - Disease states (tumor cells, virus-infected cells)
- Transcription level – Molecular features
    - Alternative isoforms
    - Fusion transcripts
    - RNA editing
- Prioritizing protein coding somatic mutations (often heterozygous)

# Evolution of transcriptomics technologies



Northern Blot



Real time RT-PCR



Microarray
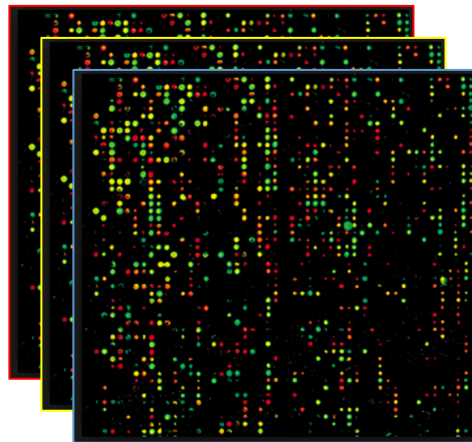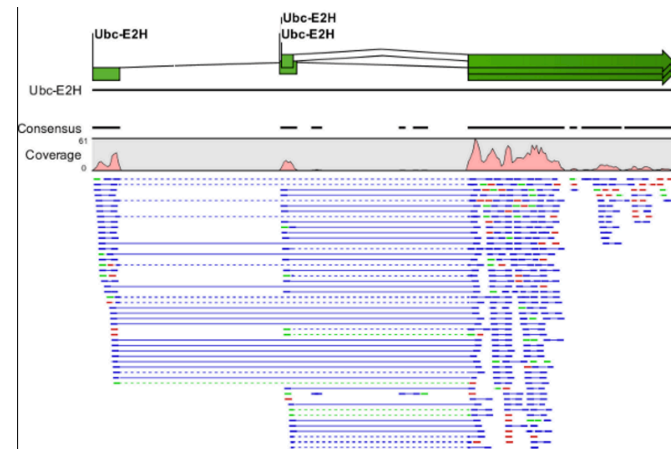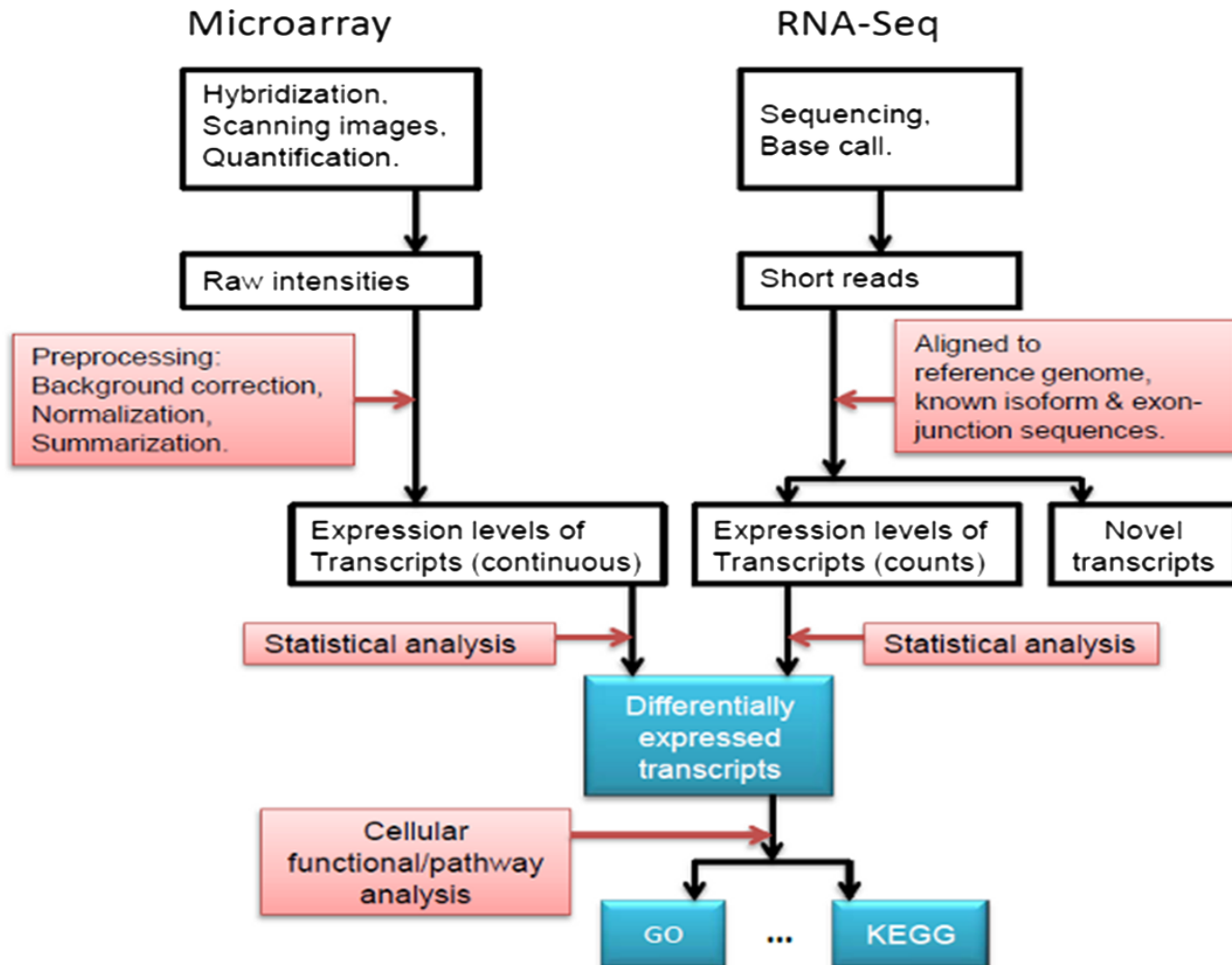


RNA seq

# Overview of analysis workflow for microarray and RNA-seq transcriptional profiling



Fang Z et.al. Cell Biosci. 2012;2:26

# RNA-seq vs. Microarray yield correlated results



Zhing LS et. Al. *Front. Plant Sci.* **5**:802.

# Microarray -> High Throughput Sequencing (HTS)

| Technology | Tiling microarray | EST sequencing | RNA-Seq |
|---|---|---|---|
| **Technology specifications** | | | |
| Principle | Hybridization | Sanger sequencing | High-throughput sequencing |
| Resolution | From several to 100 bp | Single base | Single base |
| Throughput | High | Low | High |
| Reliance on genomic sequence | Yes | No | In some cases |
| Background noise | High | Low | Low |
| **Practical issues** | | | |
| Required amount of RNA | High | High | Low |
| Cost for mapping transcriptomes of large genomes | High | High | Relatively low |
| **Application** | | | |
| Dynamic range to quantify gene expression level | Up to a few hundred-fold | Not practical | >,8000-fold |
| Simultaneously map transcribed regions and gene expression | Yes | Limited for gene expression | Yes |
| Ability to distinguish different isoforms | Limited | Yes | Yes |
| Ability to distinguish allelic expression | Limited | Yes | Yes |

*Wang et al. Nature Rev Genet, 10:57, 2009*

# RNA sequencing

**Samples of interest**

Condition 1 (e.g. tumor)

Condition 2 (e.g. normal)

**Isolate RNAs**

Poly(A) tail

**Generate cDNA, Fragment, size select, add linkers**

**Sequence ends**

100s of millions fo paired reads
10s of billions bases of sequence

**Map to genome, transcriptome, and predicted exon junctions**

Intron    pre-mRNA

Exon

Transcript

Short reads

Short reads split by intron
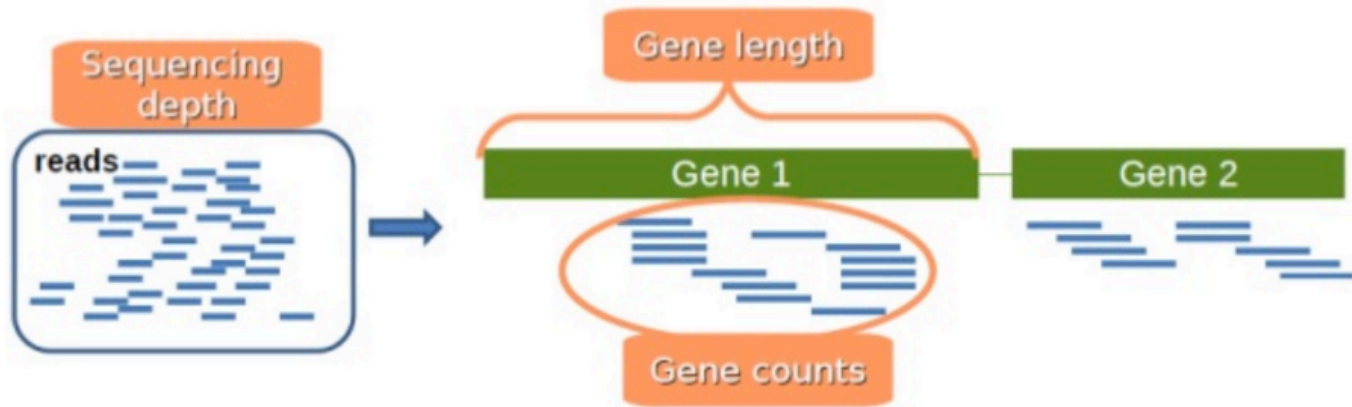
Unsequenced RNA    RNA reads

Short insert

**Downstream analysis**

# Important concepts

- Sequencing depth (X)
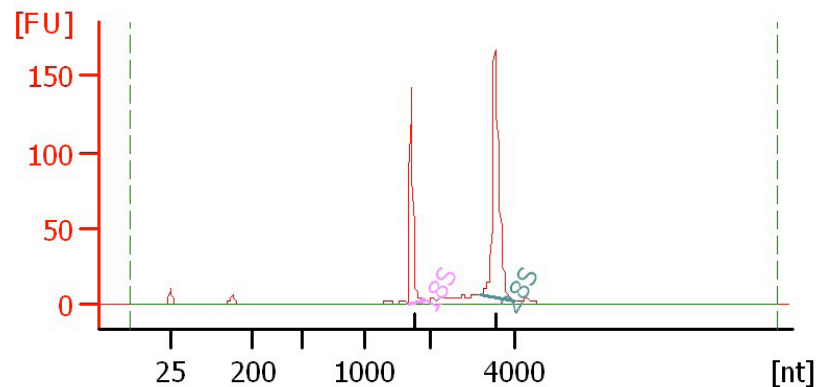- Gene length
- Gene counts

# Challenges

- Sample
  - Purity?, quantity?, quality?
- RNAs consist of small exons that may be separated by large introns
  - Mapping reads to genome is challenging
- The relative abundance of RNAs vary wildly
  - $10^5 - 10^7$ orders of magnitude
  - Since RNA sequencing works by random sampling, a small fraction of highly expressed genes may consume the majority of reads
  - Ribosomal and mitochondrial genes
- RNAs come in a wide range of sizes
  - Small RNAs must be captured separately
  - PolyA selection of large RNAs may result in 3' end bias
- RNA is fragile compared to DNA (easily degraded)

# Agilent example / interpretation

-

- 'RIN' = RNA integrity number
  - 0 (bad) to 10 (good)

# RNA-seq library construction strategies

- Total RNA versus polyA+ RNA?

- Ribo-reduction?

- Size selection (before and/or after cDNA synthesis)
  - Small RNAs (microRNAs) vs. large RNAs?
  - A narrow fragment size distribution vs. a broad one?

- Linear amplification?

- Stranded vs. un-stranded libraries

- Exome captured vs. un-captured

- Library normalization?

- These details can affect analysis strategy
  - Especially comparisons between libraries

# Fragmentation and size selection

# RNA sequence selection/depletion

# Stranded vs. un-Stranded

## A. Depiction of cDNA fragments from an unstranded library

**Legend**

└───→ Transcription start site and direction

⊦◄─── PolyA site (transcription end)

▬▬▬ Read sequenced from positive strand (forward)

▬▬▬ Read sequenced from negative strand (reverse)

## B. Depiction of cDNA fragments from an stranded library

## C. Viewing strand of aligned reads in IGV
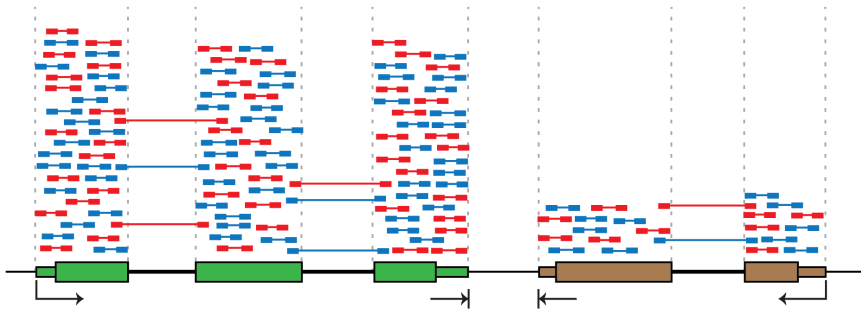
# Replicates

- Technical Replicate
  - Multiple instances of sequence generation
    - Flow Cells, Lanes, Indexes
- Biological Replicate
  - Multiple isolations of cells showing the same phenotype, stage or other experimental condition
  - Some example concerns/challenges:
    - Environmental Factors,
    - Growth Conditions,
    - Time
  - Correlation Coefficient 0.92-0.98



Correlation between expression values from libraries A and B

Correlation = 0.9802 (Spearman)
Loess fit

Library B (BS_U_1a) expression (log2[expression+1])

Library A (BS_U_1b) expression (log2[expression+1])

# Common analysis goals of RNA-Seq

- Gene expression and differential expression
- Alternative expression analysis
- Transcript discovery and annotation
- Allele specific expression
  - Relating to SNPs or mutations
- Mutation discovery
- Fusion detection
- RNA editing

# General themes of RNA-seq workflows

- Each type of RNA-seq analysis has distinct requirements and challenges but also a common theme:

1. **Obtain raw data** (convert format)

2. **Align/assemble reads**

3. **Process alignment** with a tool specific to the goal
    - e.g. 'cufflinks' for expression analysis, 'defuse' for fusion detection, etc.

4. **Post process**
    - Import into downstream software (R, Matlab, Cytoscape, Ingenuity, etc.)

5. **Summarize and visualize**
    - Create gene lists, prioritize candidates for validation, etc.

# How much library depth is needed for RNA-seq?

- Depends on a number of factors:
  - Question being asked of the data.  **Gene expression**? **Alternative expression**?  **Mutation calling**?
  - Tissue type, RNA preparation, quality of input RNA, library construction method, etc.
  - Sequencing type: read length, paired vs. unpaired, etc.
  - Computational approach and resources
- Identify publications with similar goals
- Pilot experiment
- Good news:  1-2 lanes of recent Illumina HiSeq data should be enough for most purposes

# What mapping strategy should I use for RNA-seq?

- Depends on read length

- < 50 bp reads
  - Use aligner like BWA and a genome + junction database
  - Junction database needs to be tailored to read length
    - Or you can use a standard junction database for all read lengths and an aligner that allows substring alignments for the junctions only (e.g. BLAST … slow).
  - Assembly strategy may also work (e.g. Trans-ABySS)

- > 50 bp reads
  - Spliced aligner such as Bowtie/TopHat, STAR, HISAT, etc.

# Reference Sequence Alignment (Mapping)



DNA mappers are plotted in blue, RNA mappers in red, miRNA mappers in green and bisulphite mappers in purple.

# Experiment design
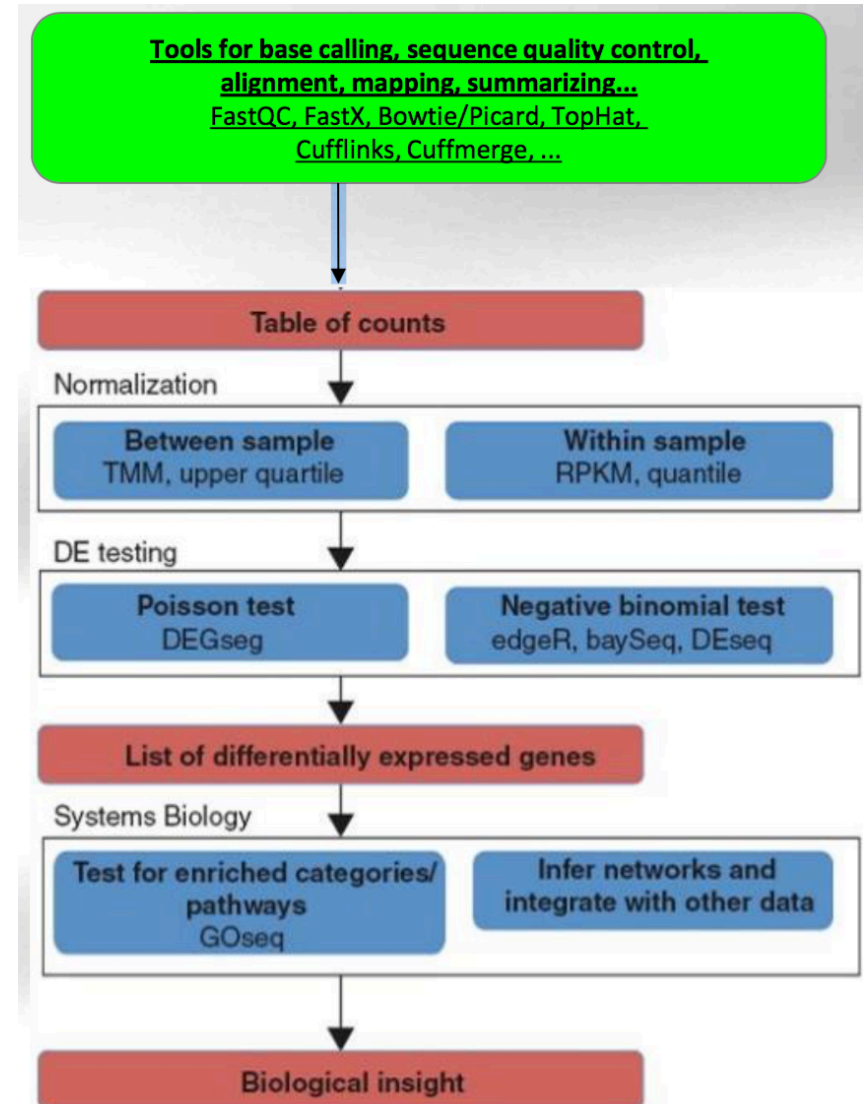
Well begun is
half done.

- Aristotle -

- A clearly defined biological question

- Well control of potential sources of variation

- HTS experimental replicates

- Compliance with the standard of ~~microarray~~ (HTS based) information collection (MINSEQE)
  - ➢ http://fged.org/projects/minseqe/

# RNA seq analysis workflow

- Reads are mapped to the reference genome or transcriptome

- Mapped reads are assembled into expression summaries (tables of counts, showing how may reads are in coding region, exon, gene or junction)

- Data is normalized

- Statistical testing of differential expression (DE) is performed, producing a list of genes with *p-values* and fold changes.

- Similar downstream analysis than microarray results (Functional Annotations, Gene Enrichment Analysis; Integration with other data...)



**Tools for base calling, sequence quality control, alignment, mapping, summarizing...**
FastQC, FastX, Bowtie/Picard, TopHat, Cufflinks, Cuffmerge, ...

Table of counts

Normalization

| Between sample TMM, upper quartile | Within sample RPKM, quantile |

DE testing

| Poisson test DEGseg | Negative binomial test edgeR, baySeq, DEseq |

List of differentially expressed genes

Systems Biology

| Test for enriched categories/ pathways GOseq | Infer networks and integrate with other data |

Biological insight

# Normalization/scaling/transformation: different goals

- **R/FPKM:** (Mortazavi et al. 2008)
    - Correct for: differences in sequencing depth and transcript length
    - Aiming to: compare a gene across samples and diff genes within sample

$$RPKM = \frac{\frac{\text{number of reads in region}}{\text{region length} \times 10^3}}{\text{total reads} \times 10^6}$$

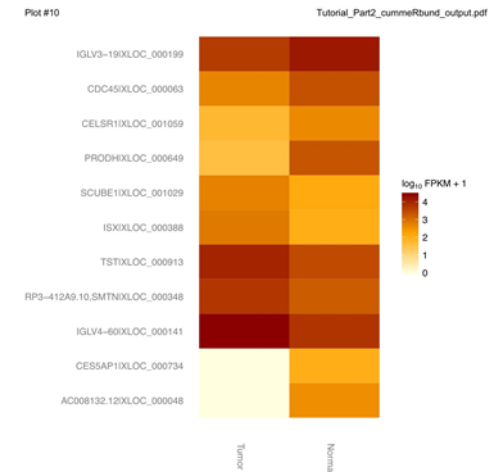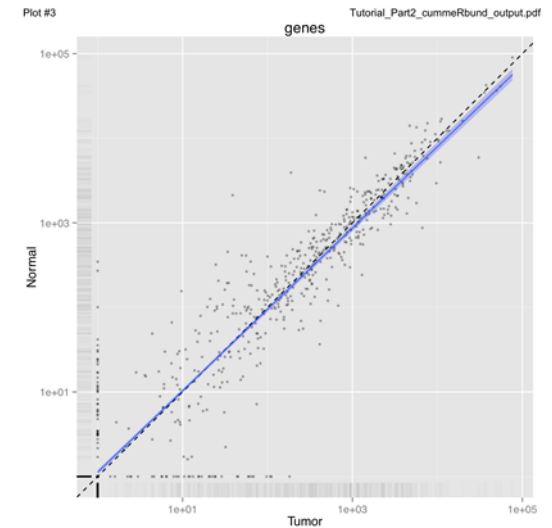- **TMM**: (Robinson and Oshlack 2010)
    - Correct for: differences in transcript pool composition; extreme outliers
    - Aiming to: provide better across-sample comparability

- **TPM**: (Li et al 2010, Wagner et al 2012)
    - Correct for: transcript length distribution in RNA pool
    - Aiming to: provide better across-sample comparability

- **Limma voom (logCPM)**: (Lawet al 2013)
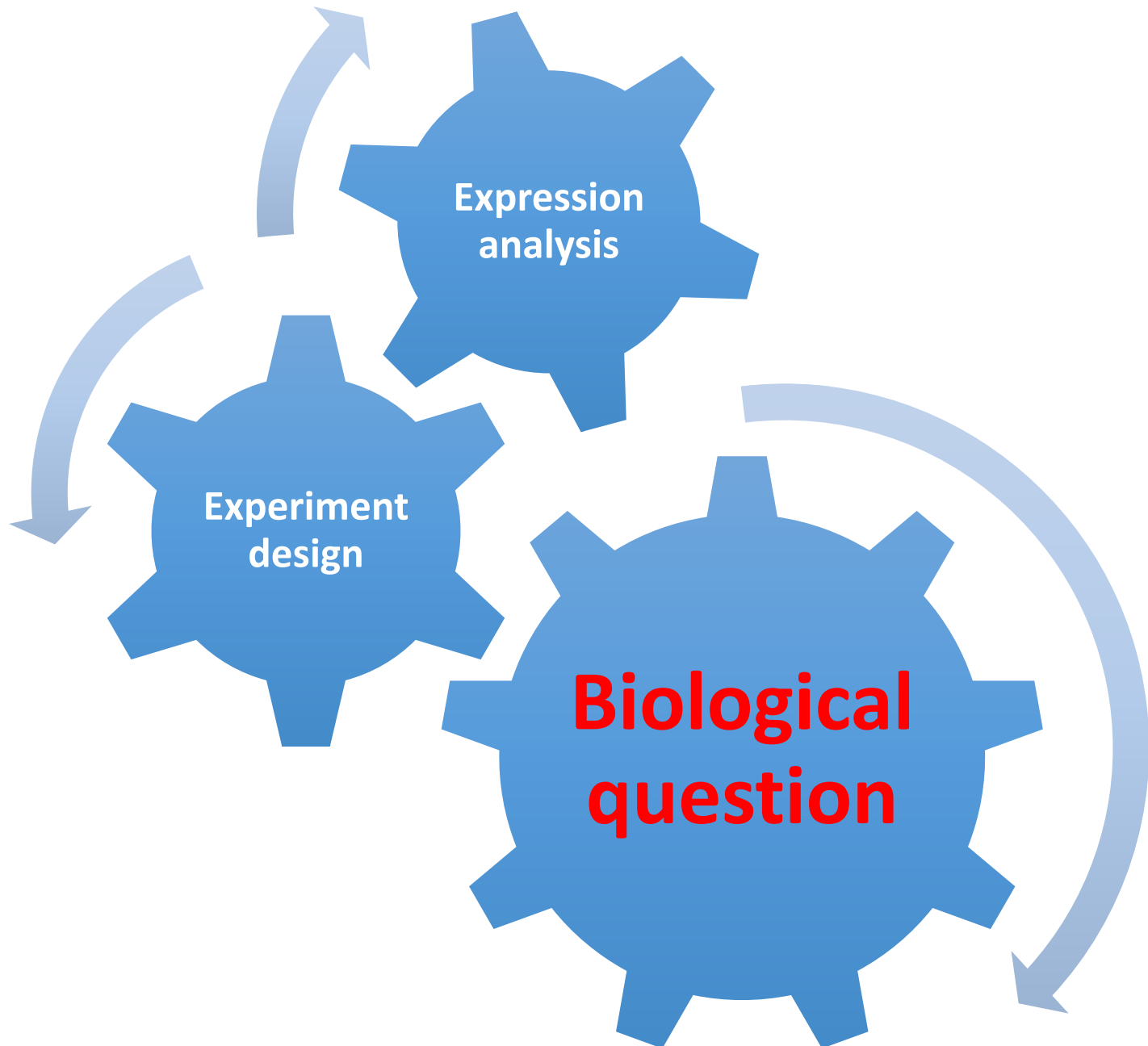    - Aiming to: stabilize variance; remove dependence of variance on the mean

# Differential expression analysis

TABLE 8.1    List of (some) Software Tools for Differential Expression Analysis

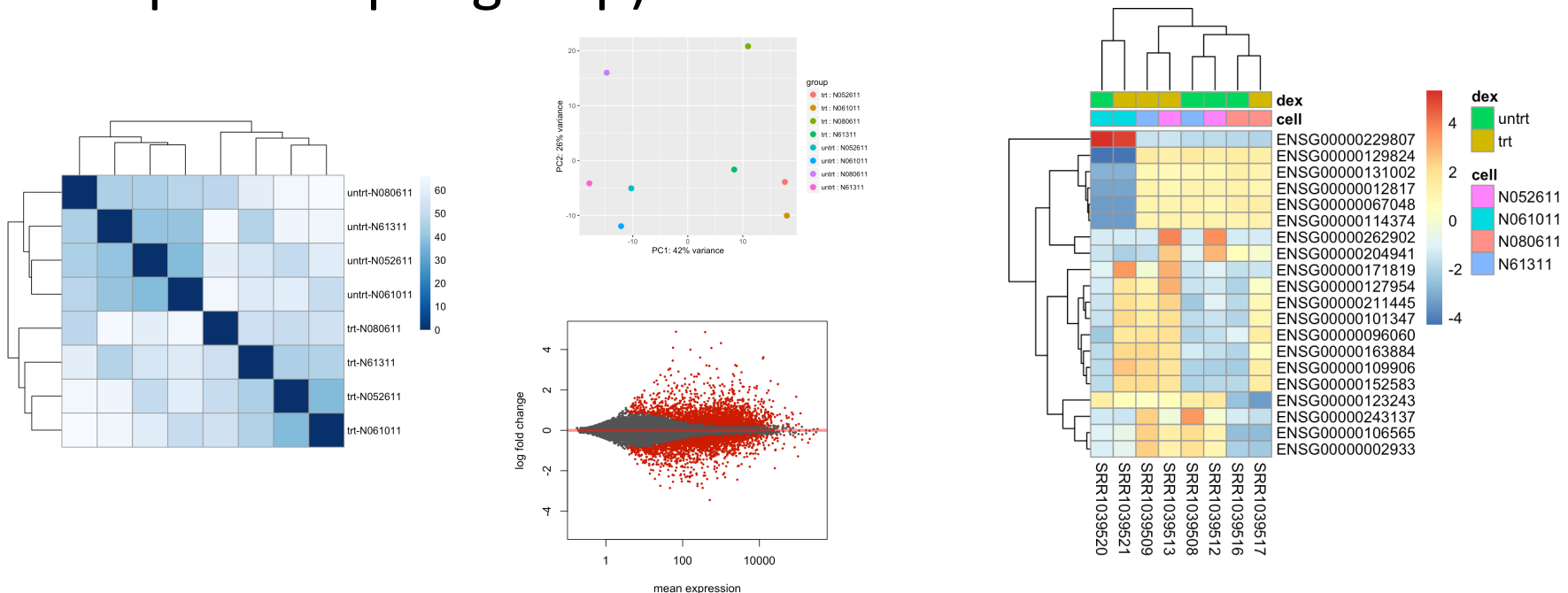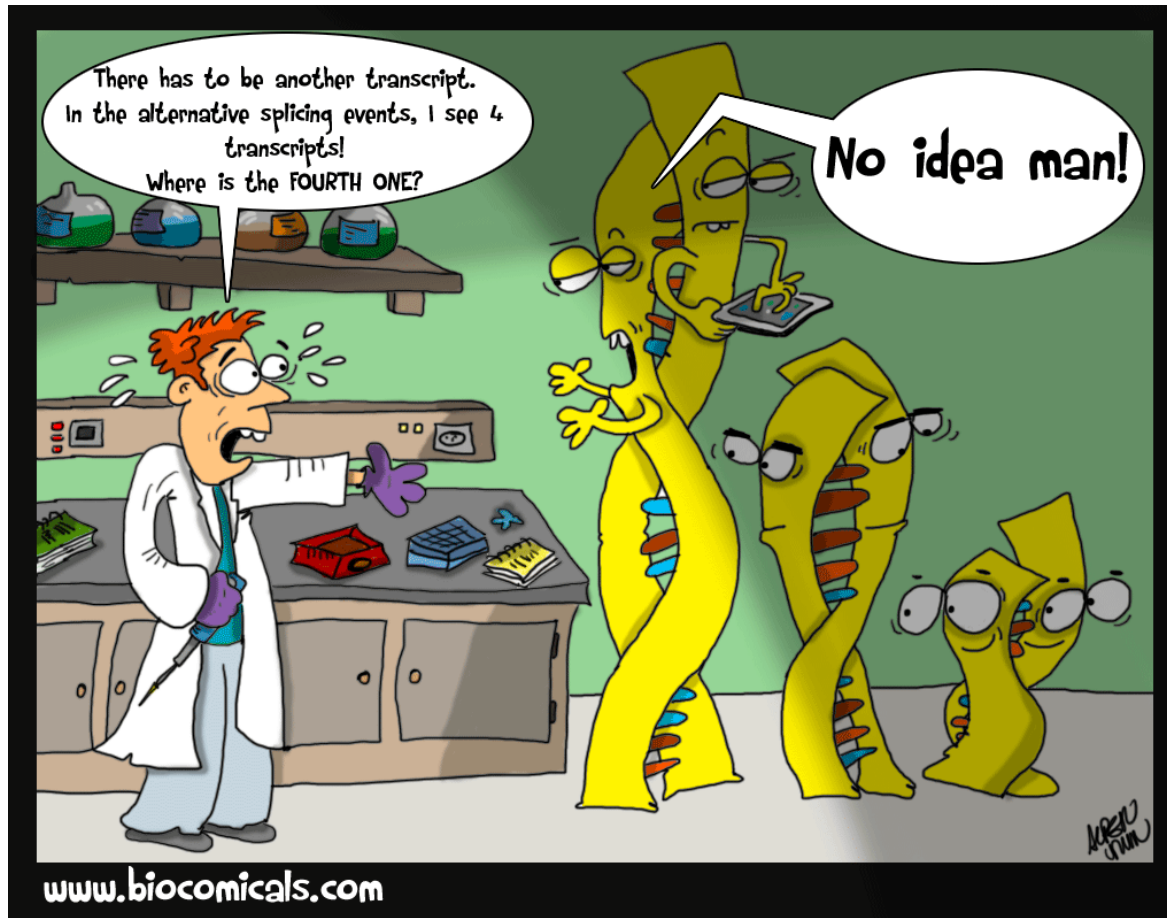| Software Tool | Type of Software | Analysis Approach | Comment |
|---|---|---|---|
| DESeq | R/Bioconductor package | Count-based (negative binomial) | Considered conservative (low false-positive rate) |
| edgeR | R/Bioconductor package | Count-based (negative binomial) | Similar to DESeq in philosophy |
| tweeDESeq | R/Bioconductor package | Count-based (Tweedie distribution family) | More general than DESeq/edgeR, but new and not widely tested |
| Limma | R/Bioconductor package | Linear models on continuous data | Originally developed for microarray analysis, very thoroughly tested. Need to preprocess counts to continuous values |
| SAMSeq (samr) | R package | Nonparametric test | Adapted from the SAM microarray DE analysis approach. Works better with more replicates |
| NOISeq | R/Bioconductor package | Nonparametric test | |
| CuffDiff | Linux command line tool | Isoform deconvolution + count-based tests | Can give differentially expressed isoforms as well as genes (also differential usage of TSS, splice sites) |
| BitSeq | Linux command line tool and R package | Isoform deconvolution in a Bayesian framework | Can give differentially expressed isoforms. Also calculates (gene and isoform) expression estimates |
| ebSeq | R/BioConductor package | Isoform deconvolution in a Bayesian framework | Can give differentially expressed isoforms. Can be used in a pipeline preceded by RSEM expression estimation |

# Bioinformatics task 1
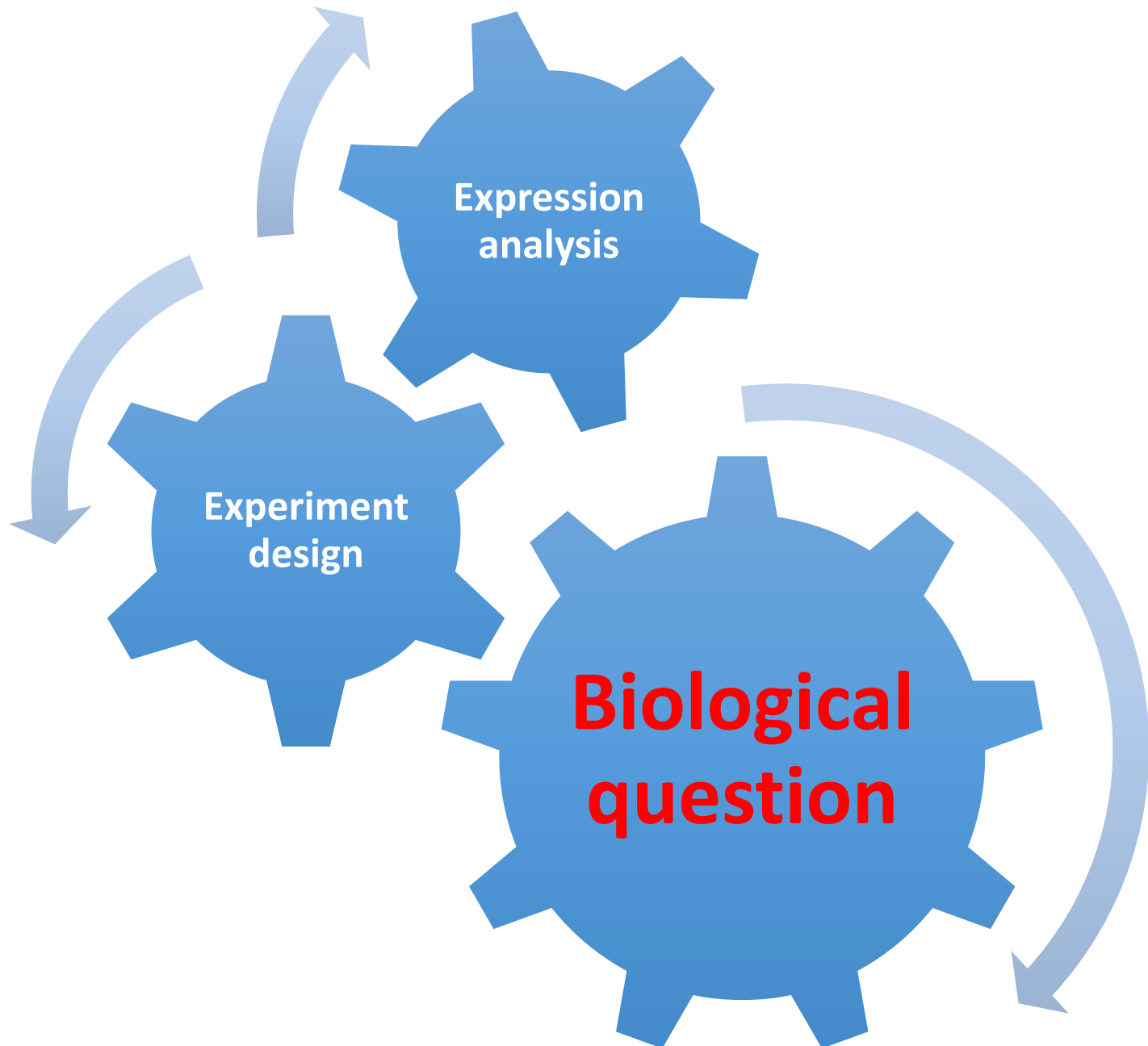
# DESeq analysis as exercise

- Differentially expressed genes
- Complex design (more than one varying factor)
- Simple comparison of groups (less than 5 biological replicates per group)

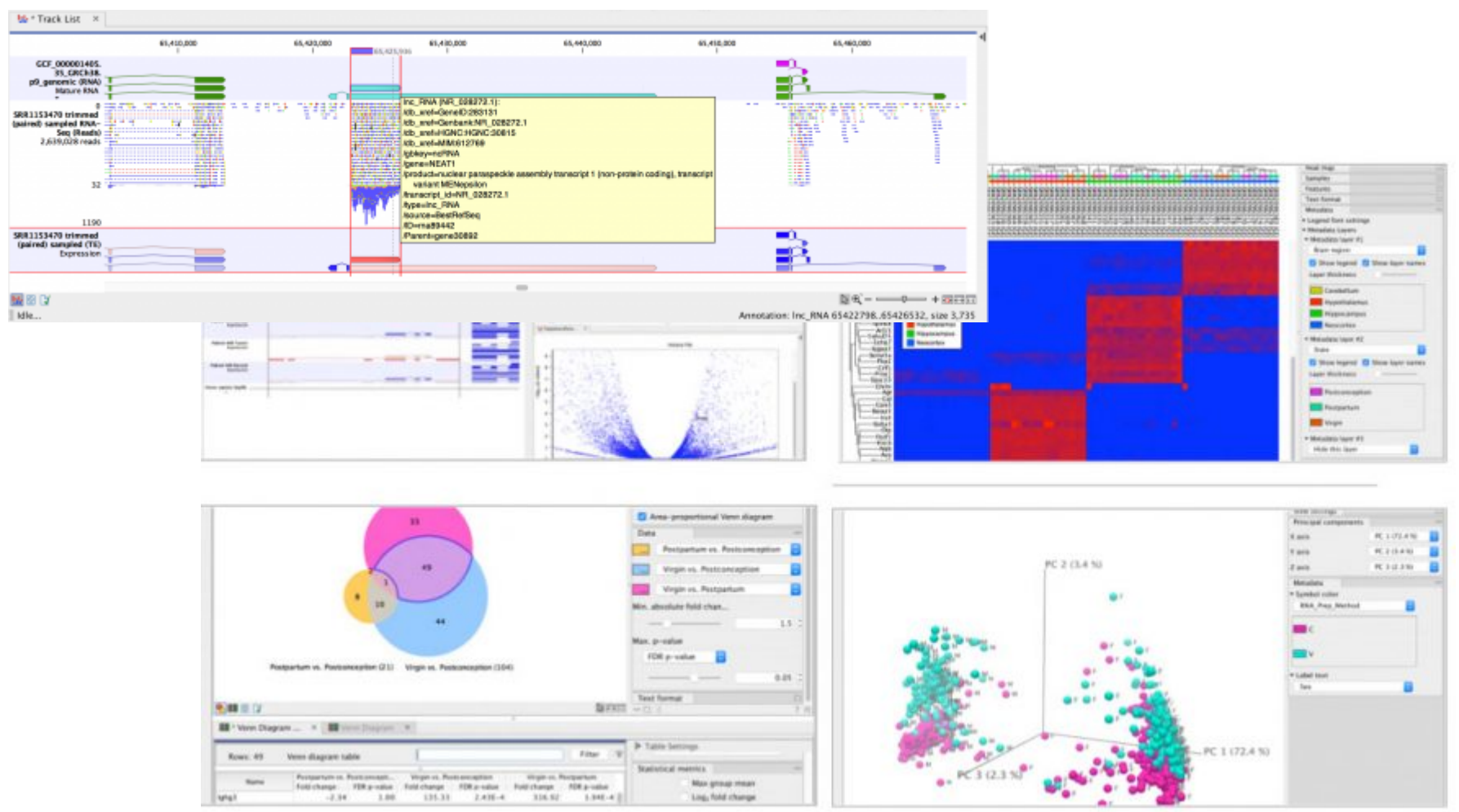**How to spot biological functions embedded in a gene list?**

# Bioinformatics task 2

# Expression Analysis using RNA-Seq

- 10 individuals out of the 726 present in the original dataset [Lin et al., 2016],and have reduced each sample to 250,000 reads mapping to chromosome 2R.
  - The reads from 10 Drosophila samples.
  - An Excel spreadsheet that contains the metadata associated with each individual.
  - SRR_ID is an SRA identifier unique for each individual.
  - DGRP_Number describes the strain of the fly.
  - Sex stated as M for males and F for females.
  - Environment stated as 2 and 3 for different calendar times for collecting the flies.
  - RNA_Prep_Method using QIAGEN RNeasy kit in all cases but following either the
  - centrifuge or the vacuum based protocol.
  - Lane of the sequencer on which the sample was loaded.
  - A workflow to rapidly and efficiently proce.

Comparison of normalization and differential expression analyses using rna-seq data from 726 individual drosophila melanogaster. BMC genomics,17(1):1.

# CLC Genomics Workbench 10

thanks for your attention