# Gene expression from RNA-Seq
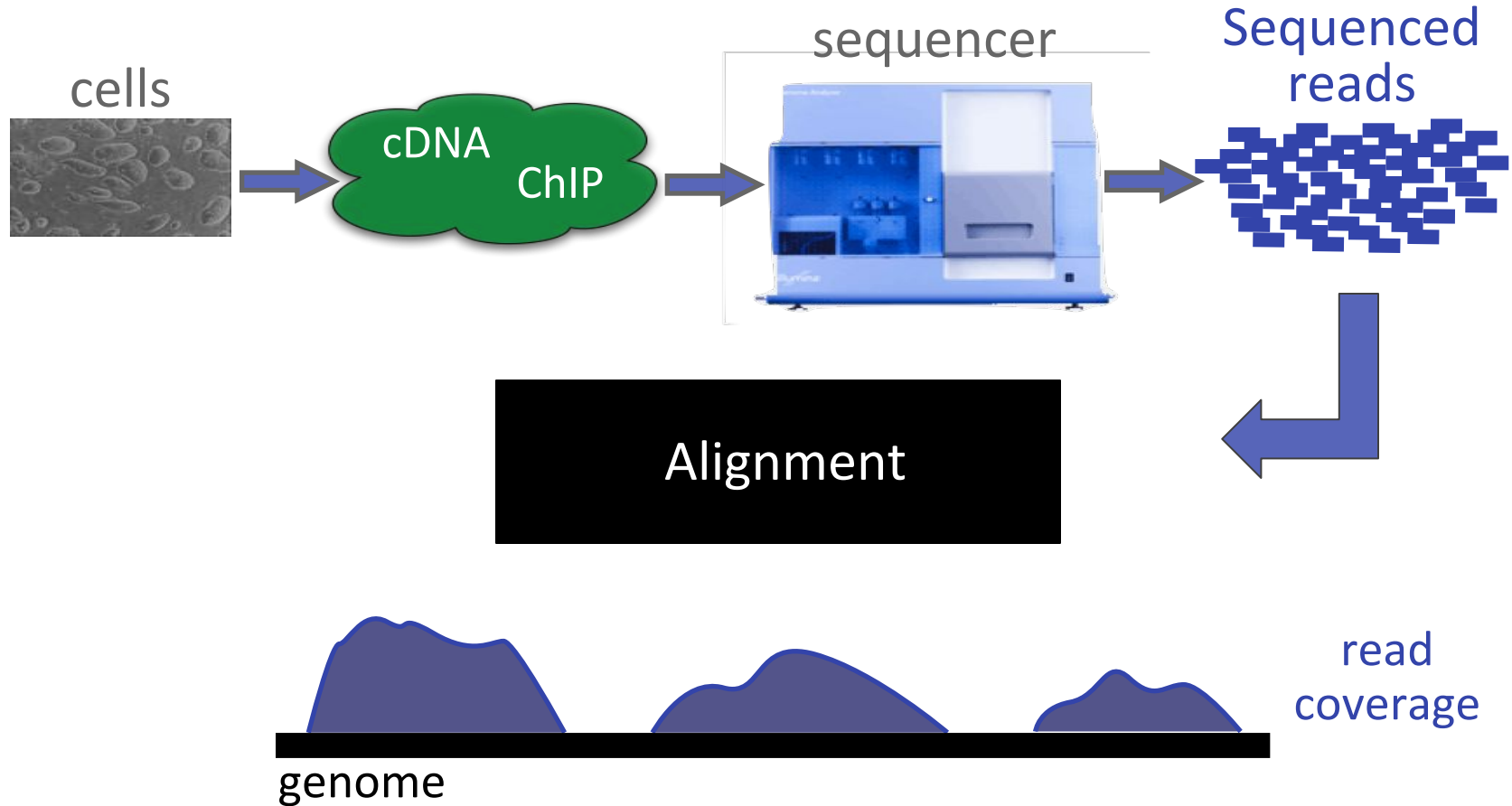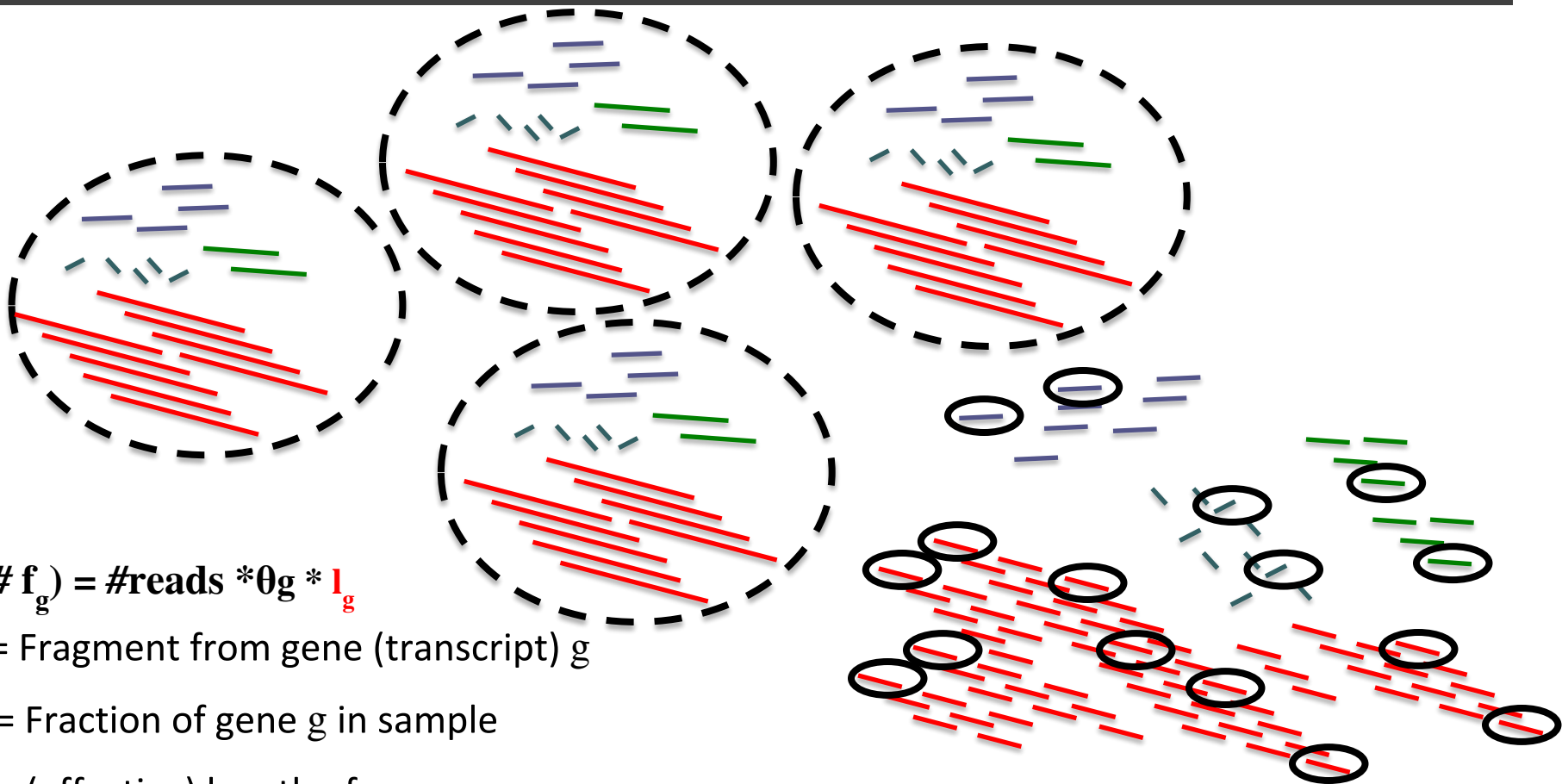
# Once sequenced the problem becomes computational

$$E(\# f_g) = \#reads * \theta g * l_g$$

$f_g$ = Fragment from gene (transcript) g

$\theta_g$ = Fraction of gene g in sample

$l_g$ = (effective) length of gene g

1. **High library complexity**

   - #molecules in library >> #sequenced molecules

2. **Short reads**

   - Read length << sequenced molecule length

Not all applications satisfy this:

- miRNA sequencing

- Small input sequencing (e.g. single cell sequencing)

- Libraries satisfying assumptions 1 & 2 only measure relative abundance

- Key quantity: # fragments sequenced for each transcript.

  **Data:** *Aligned reads*

  **Wanted***: transcript generated the observed read?*


- Isn't this easy?

  - Reads do not uniquely map

  - Genes have different isoforms with overlapping exons

  - Sequencing has a ~ 1% error rate

  - Transcripts are not uniformly sequenced

- Start with a set of previous gene/transcript annotations

- Assume only one isoform per gene

- Assume 1-1 read to transcript correspondence

Let $\Theta = \{\theta_g\}$ the relative abundance of each gene

let $n_g$ the number of reads aligned to gene $g$

$$N = \sum n_g \quad \text{(Sequencing depth)}$$

When a success has probability **θ**, the probability of **n** successes in **N** tries can be calculated by:
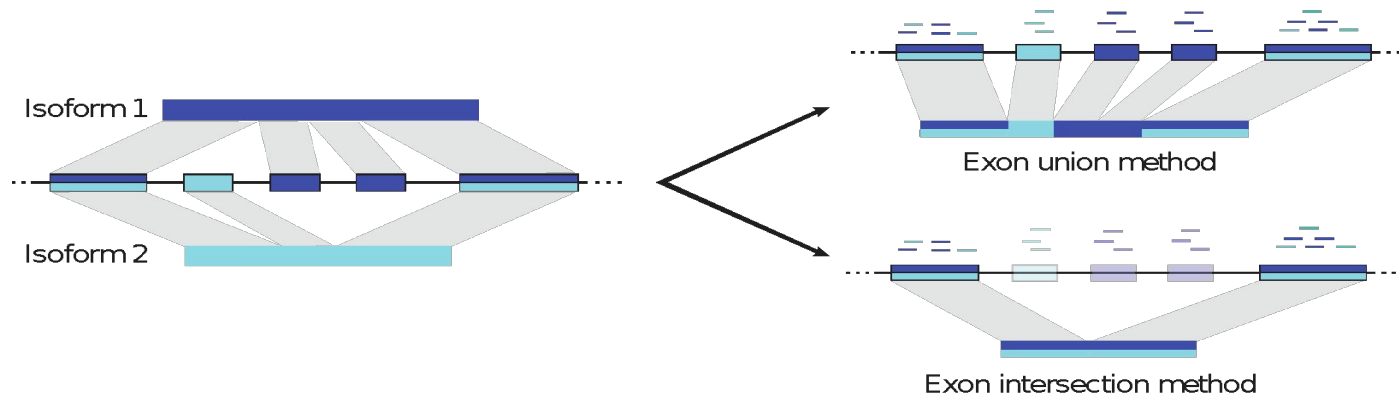
$$P(n_g \mid \theta_g) = \binom{N}{n_g} \theta_g^{n_g} (1 - \theta_g)^{N - n_g}$$

Which, has maximum probability at $\theta_g = \dfrac{n_g}{\sum n_g}$

# The process of RNA-Seq quantification
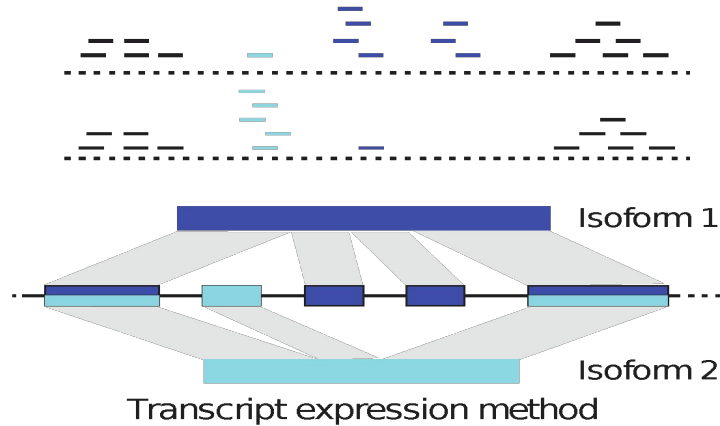
- Sequenced reads are aligned to a reference sequence

  - the species genome or

  - its transcriptome

- Transcript abundance is measured:

  - By counting reads mapped to each transcript (not accurate when multiple isoforms share sequence)

  - By solving a maximum likelihood of the observed mapping given transcript abundance

- To compare samples, the counts need to be normalized

  - Libraries have different sequencing depth

  - Sample composition may be different

  - Most standard normalization: counts → Transcripts per Million (TPM) units

# The gene expression table

Genes are quantified, each gene or isoform has:

- A TPM value

- A (expected) fragment count value

All samples were quantified in the same fashion and arranged into a table of genes (22,000) x samples (24)

- Row i gives the expression of the gene i across all samples

- Column j gives the expression of genes in sample j

| GENE | $L^{\Delta 1,2}$rep1 | $L^{\Delta 1,2}$rep2 | $L^{\Delta 1,2}$rep3 | $L^{\Delta 1}$rep1 | $L^{\Delta 1}$rep2 | $L^{\Delta 1}$rep3 | $L^{\Delta 2}$rep1 | $L^{\Delta 2}$rep2 | $L^{\Delta 2}$rep3 |
|---|---|---|---|---|---|---|---|---|---|
| Mir301 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Cpne2 | 157 | 158.98 | 88.04 | 69 | 111.99 | 114.33 | 93 | 208 | 140 |
| Capn5 | 36 | 65 | 46 | 46 | 69 | 42 | 33 | 58 | 59.01 |
| Lage3 | 313.06 | 241.23 | 276.23 | 218.9 | 285.19 | 359.65 | 269.7 | 359.04 | 417.47 |
| Brd7 | 379 | 358.58 | 390 | 336 | 357.26 | 368.08 | 264 | 564.07 | 476 |
| Dimt1 | 77 | 68 | 58 | 54 | 62 | 60 | 54 | 76 | 97.03 |
| AK017068 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Isoform 1

Isoform 2

Exon union method

Exon intersection method

- Start with a set of previous gene/transcript annotations

- ~~Assume~~ Define only one isoform per gene

- ~~Assume 1-1 read to transcript correspondence.~~

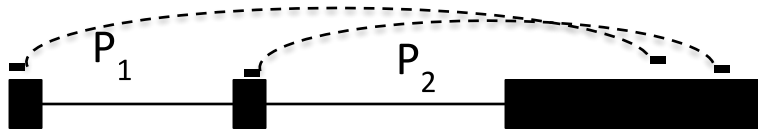- Reads (fragments) are now short, one transcript generates many fragments.

Transcript expression method

Main difference: quantification involves read assignment. Our model must capture read assignment uncertainty.

Objective: Transcript relative abundance

**Unknown!**: Fragment alignment source

Observed variables: N fragment alignments, transcripts, *fragment length distribution*

# We can estimate the insert size distribution



Get all single isoform reconstructions

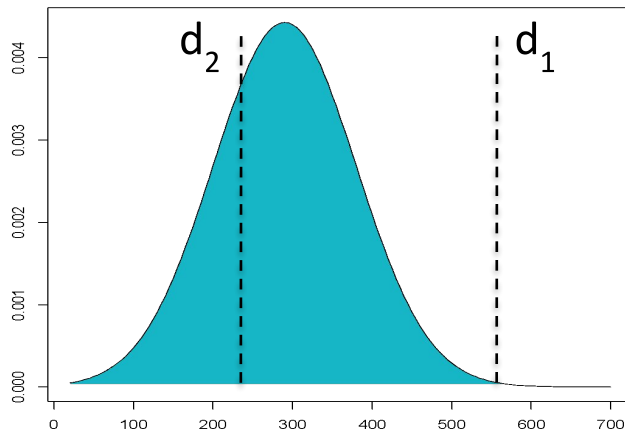Splice and compute insert distance

Estimate insert size empirical distribution

... and use it for probabilistic read assignment

Isoform 1

Isoform 2

Isoform 3

$d_1$

$d_2$

$P(d > d_i)$

$d_2$

$d_1$

For methods such as MISO, Cufflinks and RSEM, it is critical to have paired-end data

# The RNA-Seq quantification problem. Isoform deconvolution



Transcript expression method



Parameters:   Transcript relative abundance

Latent variables: Fragment alignment source

Observed variables: N fragment alignments, transcripts, **fragment length distribution**

$$P(a \in t | D, \theta_t) = \frac{\theta_t \tilde{l}_t}{\sum_{a \in s} \theta_s \tilde{l}_s} P(l(a) | t, D)$$

Probability of the fragment alignment originating from t

$$\mathcal{L}(\Theta \,|\, D, A, G) = \prod_{t \in G} \prod_{a \in t} P(a \in t | D, \theta_t)$$

solvable by expectation maximization

- In its simplest form we assume that reads can be unequivocally mapped
- This allows:

  - Read counts distribute multinomial with rate estimated from the observed counts

- When this assumption breaks, multinomial is no longer appropriate.

- More general models use:

  - Base quality scores

  - Sequence mappability

  - Protocol biases (e.g. 3' bias)

  - Sequence biases (e.g. GC)

- Handling each of these involves a more complex model where reads are assigned probabilistically not only to an isoform but to a *different loci*

Target the start or end of transcripts

Source: End-enriched RNA

- Fragmented then selected
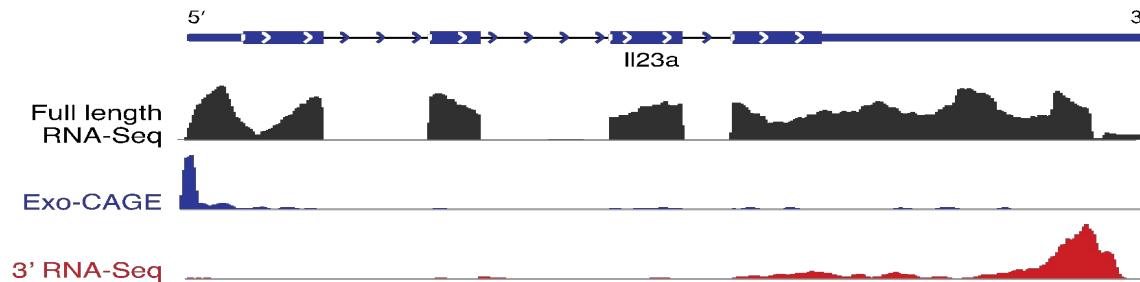
- Fragmented then enzymatically purified

Uses:

- Annotation of transcriptional start sites

- Annotation of 3' UTRs

- Quantification and gene expression

- Depth required 3-8 million reads

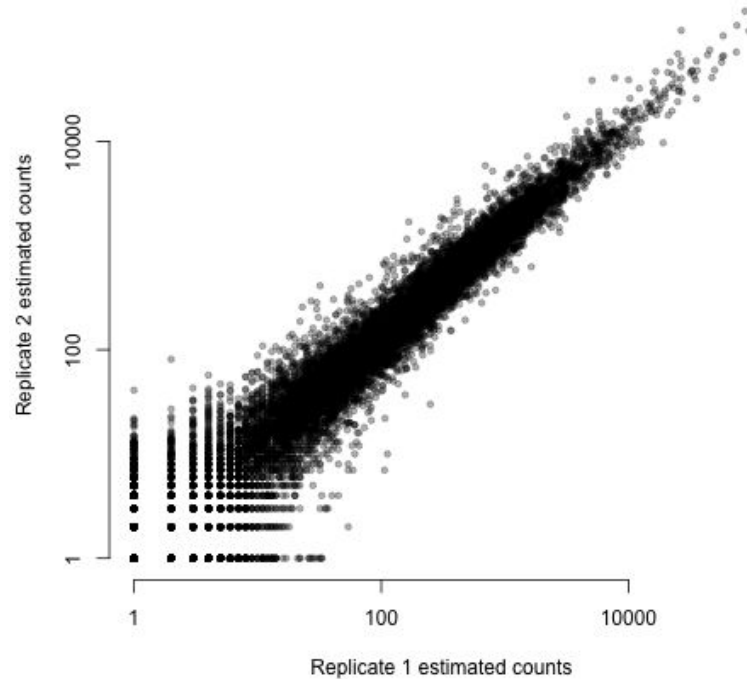- **Low quality RNA samples**

- **Single cell RNA sequencing**

A typical replicate scatter plot

A typical replicate scatter plot

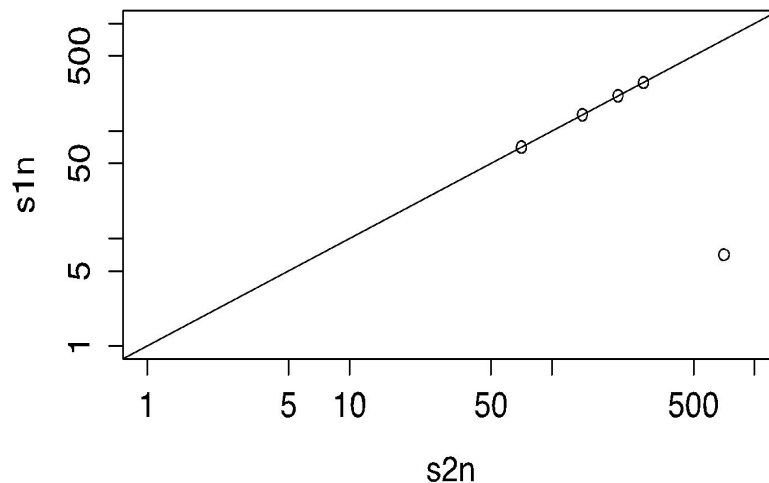$> s1 = c(100, 200, 300, 400, \boxed{10})$
$> s2 = c(50, 100, 150, 200, \boxed{500})$

Similar read number,
one transcript many fold changed

$>$norm=sum$(s2)/$sum$(s1)$
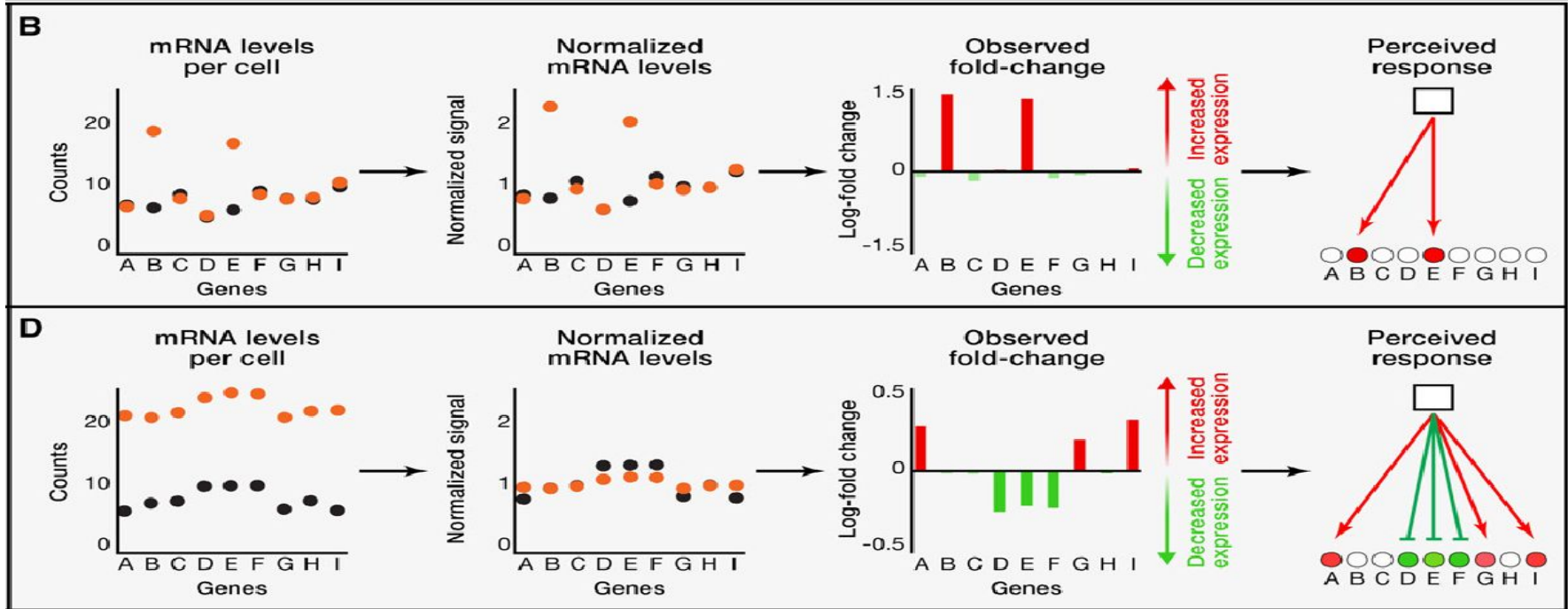$>$plot$(s2, s1*$norm,log="xy")
$>$abline$(a = 0, b = 1)$

$>$g $=$ sqrt$(s1 * s2t)$
$>$s1n $= s1/$median$(s1/g)$;  s2n $= s2/$median$(s2/g)$
$>$plot$(s2n, s1n,$log="xy")
$>$abline$(a = 0, b = 1)$

Library size normalization results
in 2-fold changes in *all* transcripts

# Finding DE genes

- Read mapping (alignment): Placing short reads in the genome

- Quantification:

  - Transcript relative abundance estimation

  - Determining whether a gene is expressed

  - Normalization: Comparing different samples

  - Finding genes/transcripts that are differentially represented between two or more samples.

- Reconstruction: Finding the regions that originated the reads