

# Intro to RNA-Seq

- RNA-Seq replacing microarrays
- Issues with microarrays:
  - designing imaging
- RNA-seq no designing probs just read over the transcriptome
- mRNA  $\longrightarrow$  Short Seq. Reads
- Biologist
- Genome Mapping bioinformaticist

## Advantages of mRNA

- not limited to known genes
- low background signal
- dynamic range of expression

## Computational Issues

- Read mapping
  - ↳ Mapping to the genome
  - ↳ STAR best aligner
- Transcriptome reconstruction
  - ↳ Cufflinks
- Expression quantification
  - ↳ How much expression are we seeing?
  - ↳ Biologists: Cufflinks

## Three Categories for regular RNA

### (a) Genome Mapping

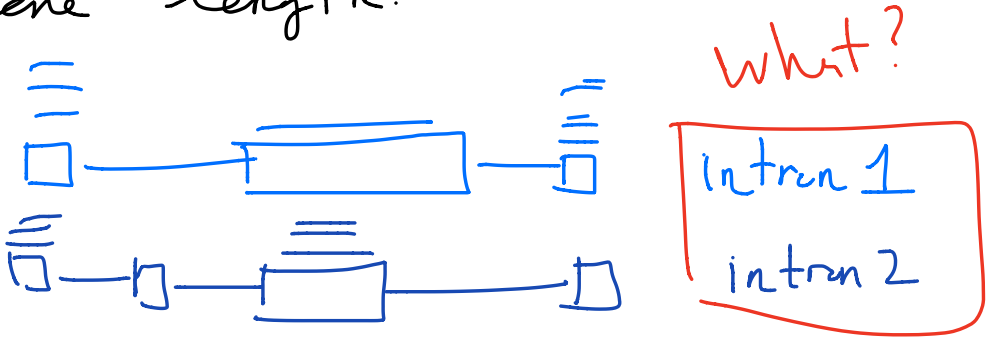
Reads  $\xrightarrow{\text{w/amps}}$  genom. mapping  $\longrightarrow$  trans. id.

### (b) Transcriptome Mapping

Reads  $\longrightarrow$  Transcriptome  $\longrightarrow$  trans. id

... ..

## (1) RefGene - tree assembly

- No mapping exists
  - reconstruct genome from scratch
- Need to normalize counts by gene length.
- 

What?

Intron 1

intron 2
- Usually just combine into transcript union method.
- The counts here are representing the number of the seq. fragments that they find in the mRNA
  - ↳ Describes the gene expression.

- CPM = counts per million

$$= \frac{\# \text{ reads}}{\text{total reads} \times 1,000,000}$$

- Pair versus unpaired read

500bp Unpaired

100bp →

Paired

← 100bp →

Ex:

Gene	Samp1	Samp2	Samp3
A	10	12	30
B			60
C			15
D			1

⇒ # reads RPKM

Tips:

- Ask for BAM
- Samtools used to access these
- (Mapped files in BigWig)
- Map to count files

Modeling

Data:

	Cond A		Cond B	
	Sample 1	...	Sample N	
Gene				

Goal: Compare Group Means.

Issues: Count data, but would  
like to use normal models

Use VST for Pois, NB, binomial.

↳ edge R,  $DES_z 2$

Edge R

$$Y \sim NB(\mu, \phi) \quad E(Y) = \mu$$
$$Var(Y) = \mu(1 + \phi\mu)$$

Let  $Y_i \sim NB(\mu_i = m_i \lambda, \phi)$

library  
size

For equal library sizes

$$Z = \sum_{i=1}^n Y_i \sim NB(nm\lambda, \phi n^{-1})$$

sufficient stat. for  $\lambda$ .

$$l_{Y/Z=Z}(\phi) = (\text{see slides})$$

If the library sizes are not equivalent, then we interpolate

$\tau$   
to a common size, then do  
the same CML.

- $m^* = (\prod m_i)^{1/n}$
- Interpolate counts as if the samp.  
were from  $NB(m^*, \rho)$ 
  - Initialize  $\phi$
  - Estimate  $\lambda$ , given  $\phi$
  - See slides...

This is just one gene. Looking  
to use some Empirical Bayes ideas.

$$WL(\phi_g) = l_g(\phi_g) + \alpha l_c(\phi_a)$$

Weighted likelihood with tuning  $\alpha$ .

Estimate  $\hat{\lambda}_{mom}$  using large

• 1 1

sample to

Approx steps

- Pseudo Counts Interpolation
- Estimating  $\ell$

DES<sub>2</sub>2:

$$Y_{ij} \sim NB(\mu_{ij}, \alpha_i)$$

$$\mu_{ij} = S_i z_{ij}$$

$$\log z_{ij} = \sum_r x_{jr} \beta_{ir}$$

$$S_j \sim \text{Median}_{i: K_i^R \neq 0} \frac{Y_{ij}}{Y_i^R} \quad Y_i^R = \left( \prod Y_{ij} \right)^{1/m}$$

Interested in inferring  $(\alpha, \beta)$

Using a Bayesian model

$$\log \alpha_i \sim N(\log \alpha_{tr}(\bar{\mu}_i), \sigma_\alpha^2)$$

- 1 ...



$$\mu_i = \frac{1}{m} \sum_j \frac{1}{s_j}$$

Parametric  
constraint.

$$\alpha_{tr}(\mu) = \frac{a_1}{\mu} + a_0$$

Using the posterior max gives  $\bar{z}$

$$\beta_{ir} \sim N(0, \sigma_r^2) \quad \sigma_r^2 \leftarrow \text{matches percentiles}$$

$\hat{\beta}_r$

$$\hat{\beta}_i = \operatorname{argmax} \left( \sum_j \log f_{NB}(y_{ij}; \mu_j(\beta), \alpha_i) \right) + \lambda(\beta_i)$$

Issues with DESeq2:

lots of approximations to accommodate NB.

Voom + limma

• LM with  $\log \text{CPM}$

- Really good.

Can we improve by aggregation?

- Still really unclear.

- Read Cufflinks with supp.  
materials