

## Analysis of multi-factor RNA / ChIP-Seq experiments with respect to biological variation

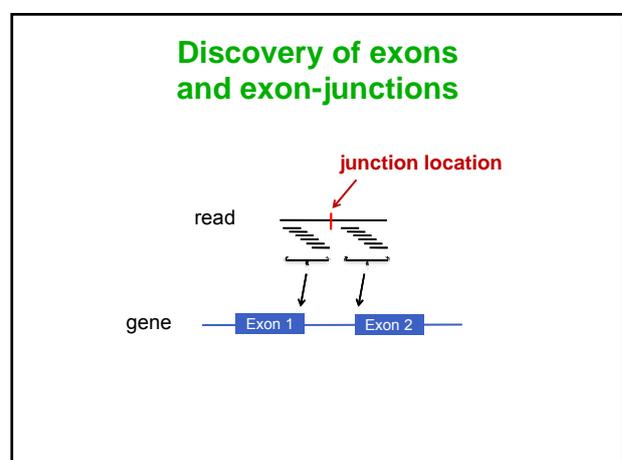
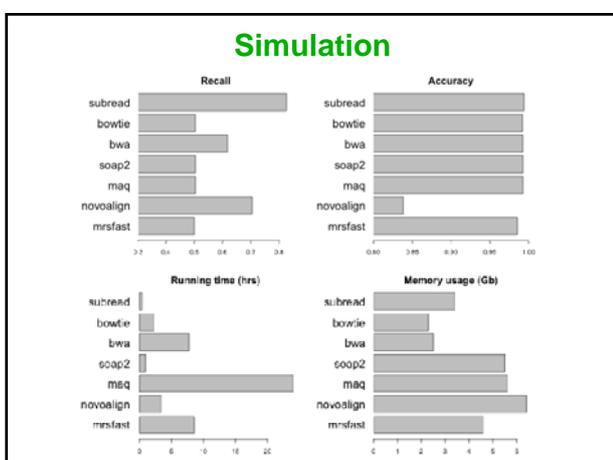
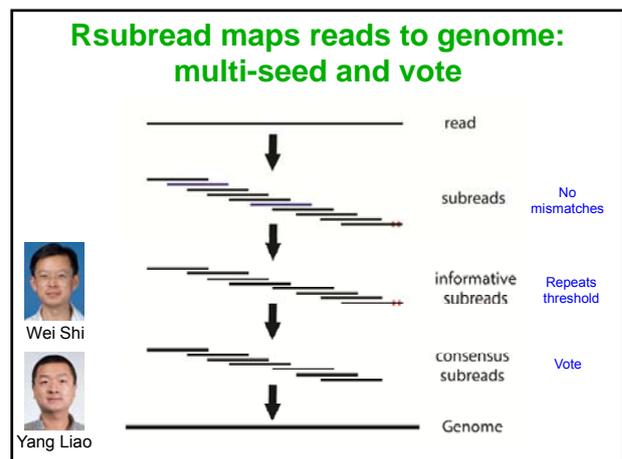
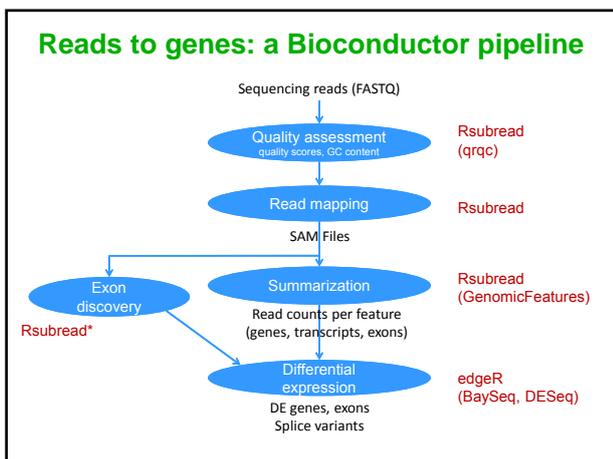
Gordon Smyth  
Bioconductor  
28 July 2011

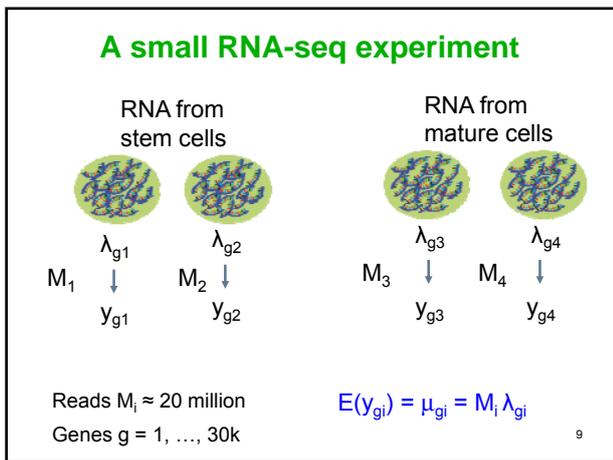
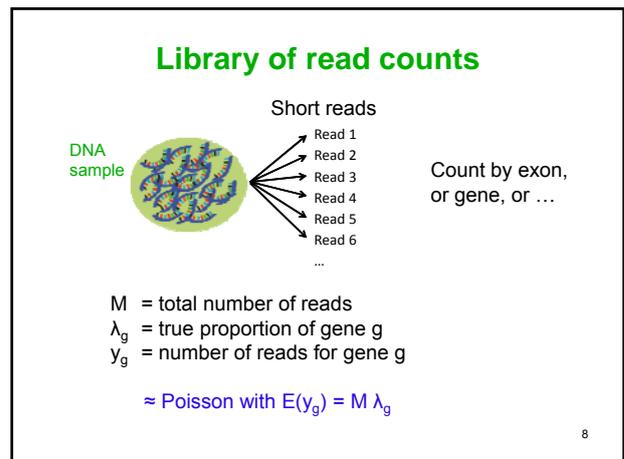
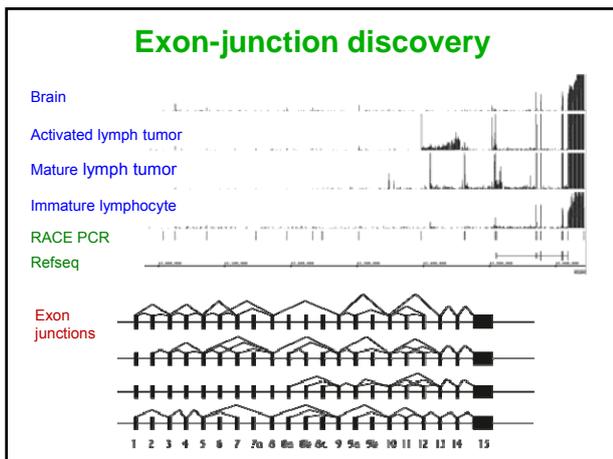


1

## Packages

- **Rsubread**
  - Read alignment
  - Summarization by genomic features
  - Exon discovery
- **limma**
  - Normal-based DE analysis
  - Gene set analysis
- **edgeR**
  - Negative binomial-based DE analysis
  - Detection of splice-variants\*
- **goseq**
  - Gene ontology analysis adjusted for gene length <sup>2</sup>





### log-linear models

$$\log \mu_{gi} = \log \lambda_{ig} + \log M_{gi}$$

$$= \mathbf{x}_i^T \boldsymbol{\beta}_g + \log M_{gi}$$

row of design matrix  $\mathbf{x}_i$

vector of log fold changes  $\boldsymbol{\beta}_g$

(normalized) library size  $M_{gi}$

(unknown)

10

- ### Normalization
- Scale normalization
    - "Effective" library size
  - Nonlinear normalization
    - Quantile normalization
    - Gene length
    - GC content (of reads, of fragments)
- Robinson and Oshlack, Genome Biol 2010
- cqn  
EDASeq
- 11

### Counts show a quadratic mean-variance relationship

$$\text{var}(y_{gi}) = \mu_{gi} + \phi_g \mu_{gi}^2$$

Poisson variation  $\mu_{gi}$

$\text{CV}^2$  of the "true" expression levels  $\lambda_{gi}$  across replicates  $\phi_g \mu_{gi}^2$

$\text{CV} = \text{coefficient of variation} = \text{sd}/\text{mean}$

12

### Biological coefficient of variation

$$\text{Total CV}^2 = \text{Technical CV}^2 + \text{Biological CV}^2$$

From sequencing technology

↓ zero for large counts

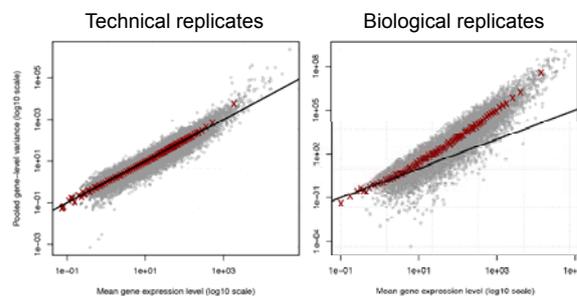
CV of "true" expression level

≈ constant

$$\text{BCV} = \sqrt{\phi_g}$$

13

### Real data show quadratic variances



14

### Statistical properties of read counts

#### Properties

- Integer values (discrete)
- Mean-variance relationship
- Distinguish technical from biological variation

#### Approaches

- log-counts as normal (**limma**)
- counts as negative binomial (**edgeR**)

15

### Limma approach

log-counts:

$$z_{gi} = \log_2 \left( \frac{\text{count}_{gi} + 0.5}{\text{libsiz}_{gi} + 0.5} \right) = \log_2 \left( \frac{y_{gi} + 0.5}{M_{gi} + 0.5} \right)$$

normalize libsize in advance or normalize  $z_{gi}$  as for microarrays.

Linear modelling:

$$E(z_{gi}) = \mu_{gi} = x_i^T \beta_g$$

$$\text{var}(z_{gi}) = s(\mu_{gi}) \sigma_g^2$$

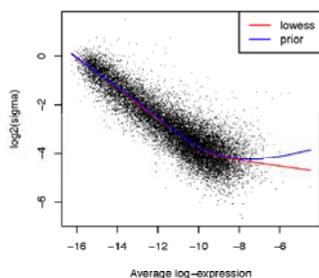
Smooth function of mean



Charity Law

16

### Empirical Bayes with abundance-dependent prior



eBayes(fit, trend=TRUE); plotSA(fit)

17

### Negative binomial approach

If  $\lambda_{gi}$  are gamma distributed, then

$$y_{gi} \sim \text{NegBin}(\mu_{gi}, \phi_g)$$

Once the dispersions are estimated, the log-linear models are **generalized linear models**



Mark Robinson



Davis McCarthy  
Yunshun Chen

18

### Ensuring glm convergence

- Iterative fitting of glms is computationally demanding, and standard glm code can **diverge**
- Pseudo** Newton-Raphson strategy to reduce need for matrix decompositions
- Line searches** to prevent divergence
- Highly **vectorized** code
- Fit genewise glms in a few seconds

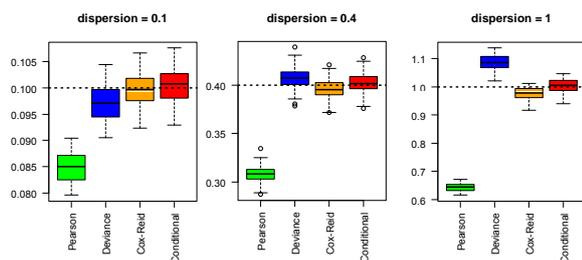
19

### Conditional inference for the dispersions

- Need to adjust for estimation of  $\beta_g$  when forming likelihood for  $\phi_g$
- For two-group comparison, can compute **conditional distributions** given row totals and conduct exact inference
- For more general designs, use **Cox-Reid** adjusted profile likelihood to condition on estimator of  $\beta_g$

20

### Performance of conditional estimators of dispersion



21

### Complexity of dispersion: sharing information between genes

- Separate gene-wise estimation of  $\phi_g$  is impractical
- Common** dispersion (Robinson & Smyth 2008)
- Trended** dispersion (Anders & Huber 2010)
- Gene-wise** by empirical Bayes shrinkage (Robinson & Smyth, 2007)

22

### Common dispersion likelihood

Assume **same dispersion** for all genes  
 $\phi_g = \phi$

Genewise conditional log-likelihood  
 $\ell_g(\phi; y_g)$

Common-dispersion log-likelihood  
 $\ell_c(\phi) = (1/G) \sum_g \ell_g(\phi; y_g)$

Maximized at  $\phi_c$

23

### Empirical Bayes shrinkage for the dispersion

Estimate  $\phi_g$  by empirical posterior mode:

$$\text{Posterior} = \ell_g(\phi_g) + \alpha \ell_c(\phi_c)$$

Genewise likelihood

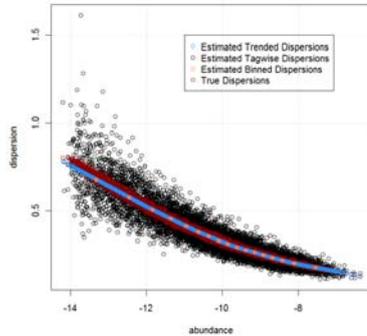
Precision of prior

Empirical prior distribution

Local weighting produces abundance dependent prior

24

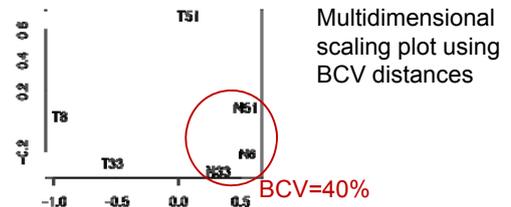
### Estimated dispersions (simulation)



### Oral squamous cancer

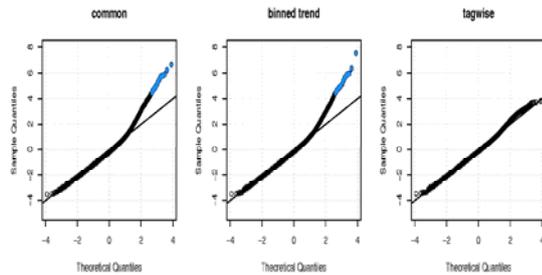
	Normal	Tumour
Patient 8	N8	T8
Patient 33	N33	T33
Patient 51	N51	T51

Tuch et al,  
*PLoS ONE* 2010



26

### Genewise goodness of fit tests



Tagwise dispersion gives the best fit

### Differential expression

Fit models of **increasing complexity**:

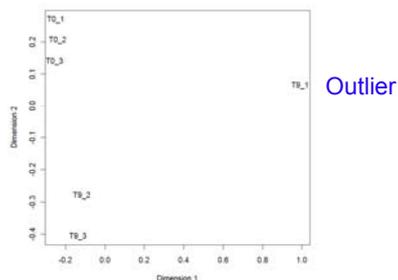
*Patient*  
| LRTs            1271 generally DE genes

*Patient + Tissue Source*  
| LRTs            184 genes specific to individual tumours

*Patient \* Tissue Source*

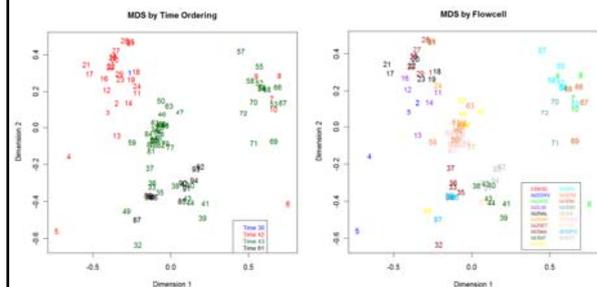
FDR < 0.05    28

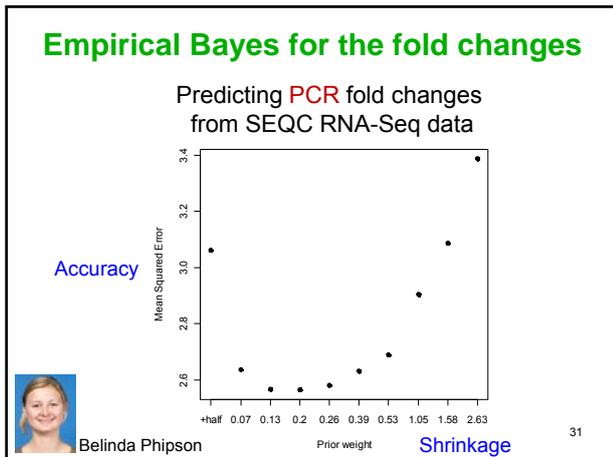
### Multidimensional scaling plots with BCV as distance



29

### Finding technical effects



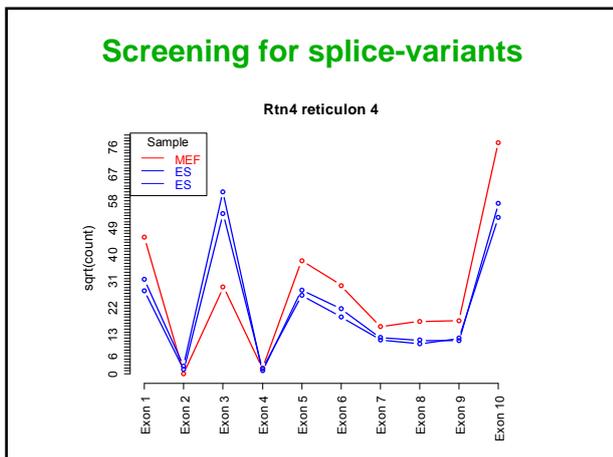


### Screening for splice-variants

- Exon level summaries
- Estimate exon-wise dispersions
- Test **exon x group interaction** for each gene

Compare to:  
Richard et al, NAR 2010  
DEXSeq package

Davis McCarthy



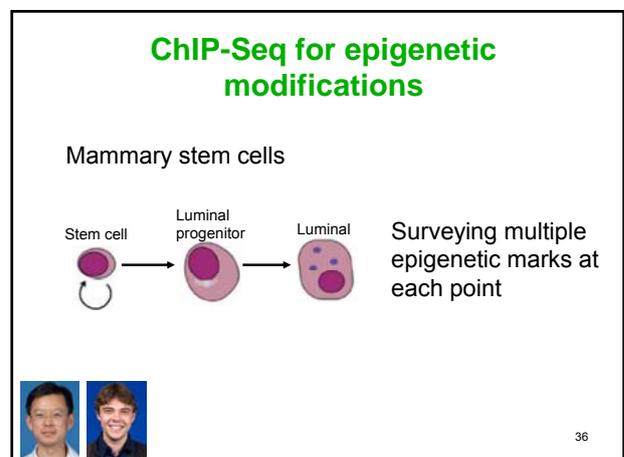
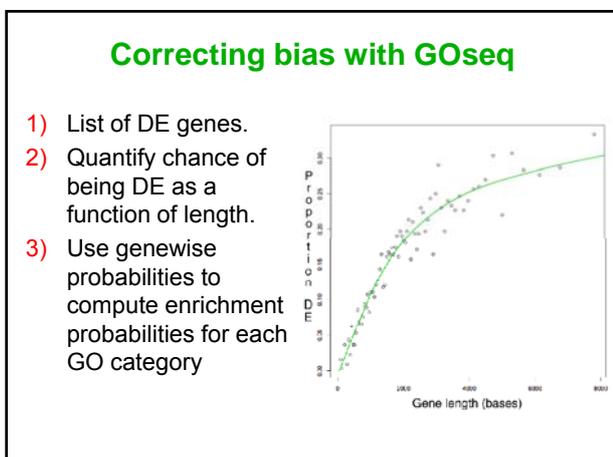
### GOseq

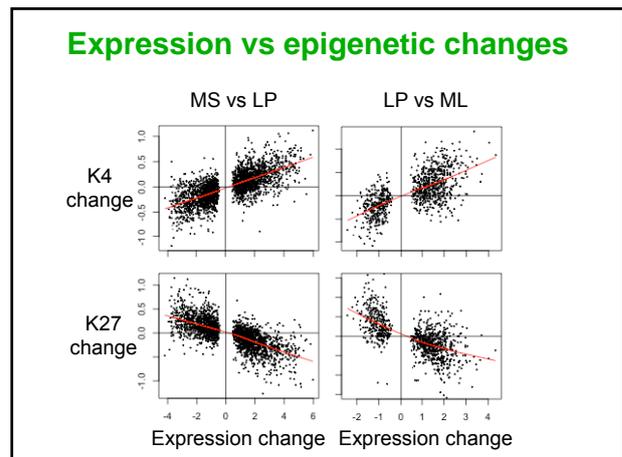
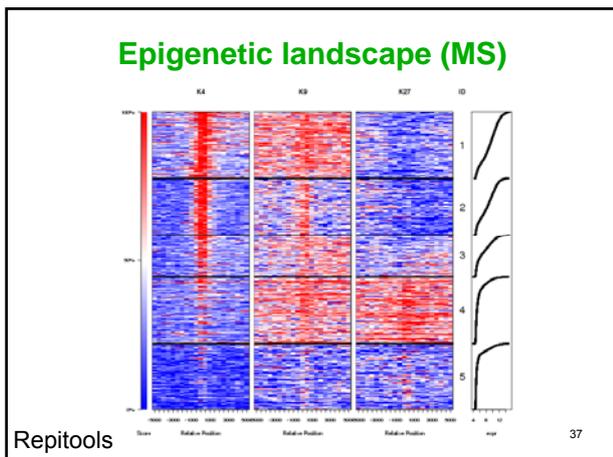
Genes vary in length ...

6X number of fragments

More power to detect DE at a given threshold

Alicia Oshlack Matt Young





- ### Conclusions
- Self-contained pipeline for RNA-Seq close at hand
  - Methods of differential expression analysis of RNA-seq (etc) data based on mean-variance modelling of counts and conditional inference
  - Shared-parameter likelihood priors provide a generally applicable paradigm for parameter shrinkage
- 39

- ### Lab Members
- Matthew Ritchie
  - Wei Shi
  - Belinda Phipson
  - Charity Law
  - Yunshun Chen
  - Joshy George
  - Keith Satterley
  - Yifang Hu
  - Davis McCarthy
  - Cynthia Liu
  - Yang Liao
  - Jenny Dai
- Recent past:
- Alicia Oshlack
  - Di Wu
  - Matthew Young
  - Luke Zappia
  - Carolyn de Graaf
  - Mark Robinson
- 40

- ### Collaborators
- Lynn Corcoran
  - Tim Thomas
  - Anne Voss
  - Bilal Sheikh
  - Samir Taoudi
  - Peter 't Hoen
  - Jane Visvader
  - Geoff Lindeman
  - Bhupinder Pal
  - Marie-Liesse Asselin-Labat
- 41