



**University of  
Zurich** <sup>UZH</sup>

Institute of Molecular Life Sciences

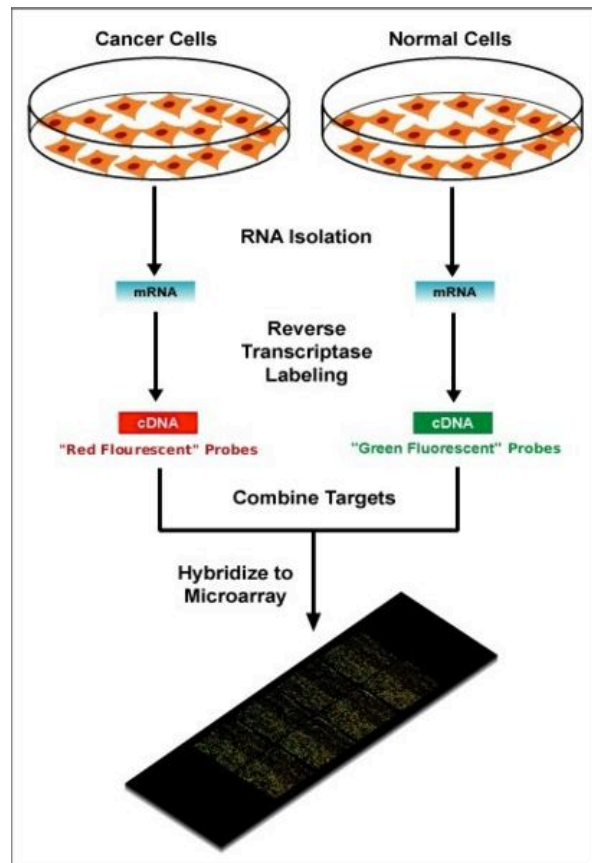
---

**Lecture part: BioC2014**  
**Differential gene- and exon-level  
expression analyses for RNA-seq data using  
edgeR, voom and featureCounts**

Mark D. Robinson, University of Zurich

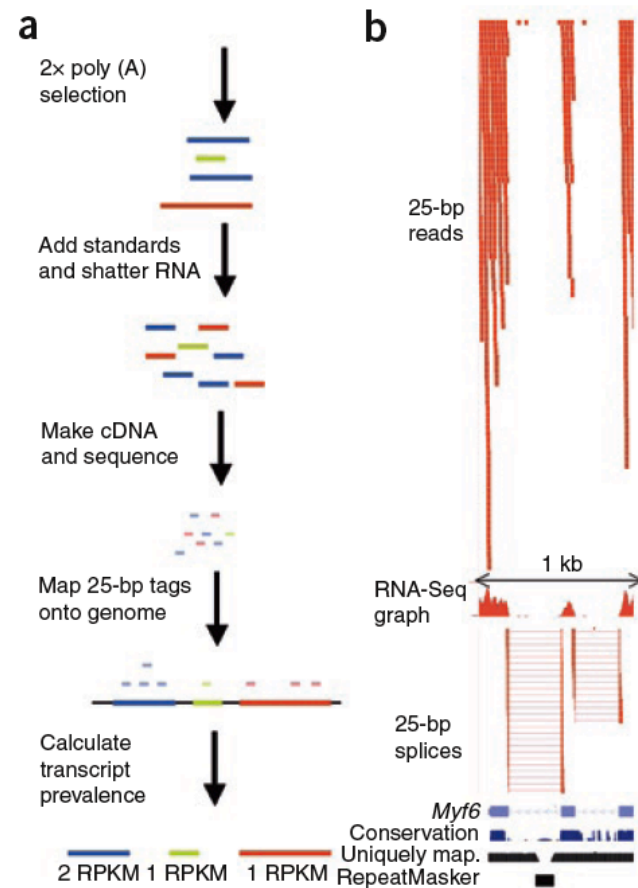


## Abundance by Fluorescence Intensity



[http://en.wikipedia.org/wiki/DNA\\_microarray](http://en.wikipedia.org/wiki/DNA_microarray)

## Abundance by Counting

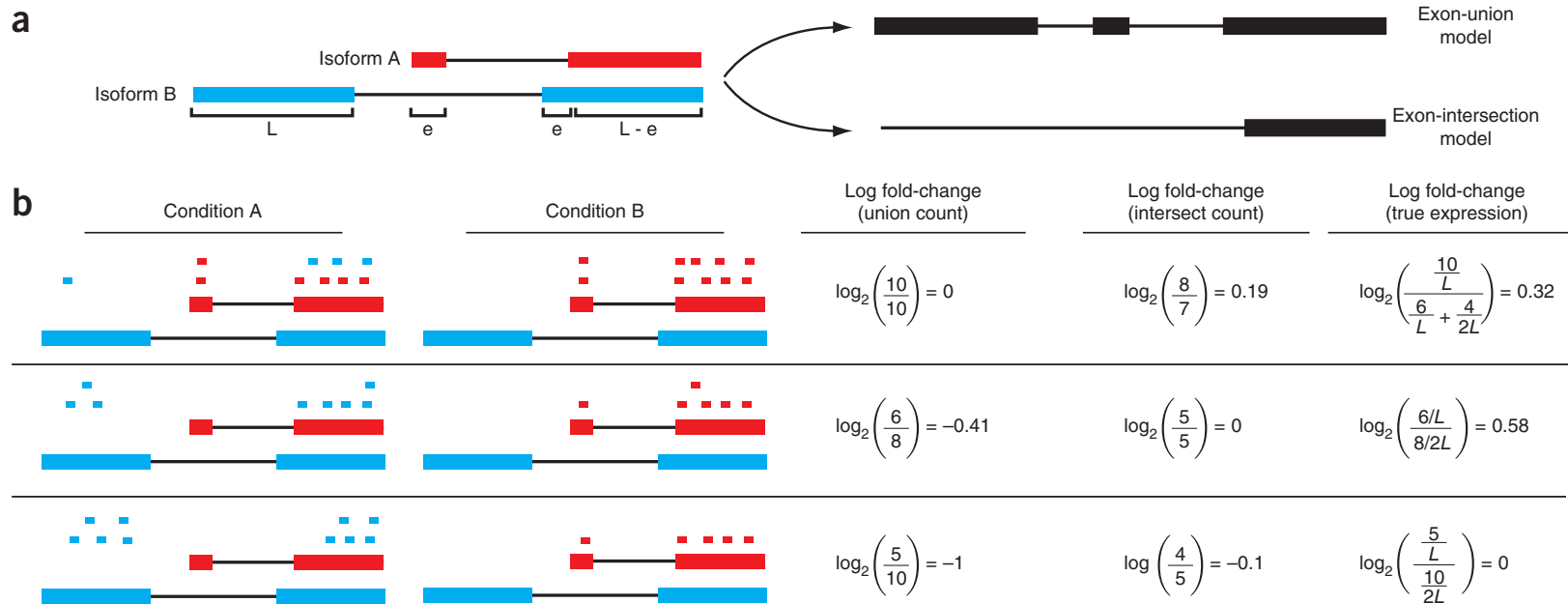


Mortazavi et al., Nature Methods, 2008



# Gene-level counting: issues can be dealt with at second step

Trapnell et al. 2013 Nat Biotech

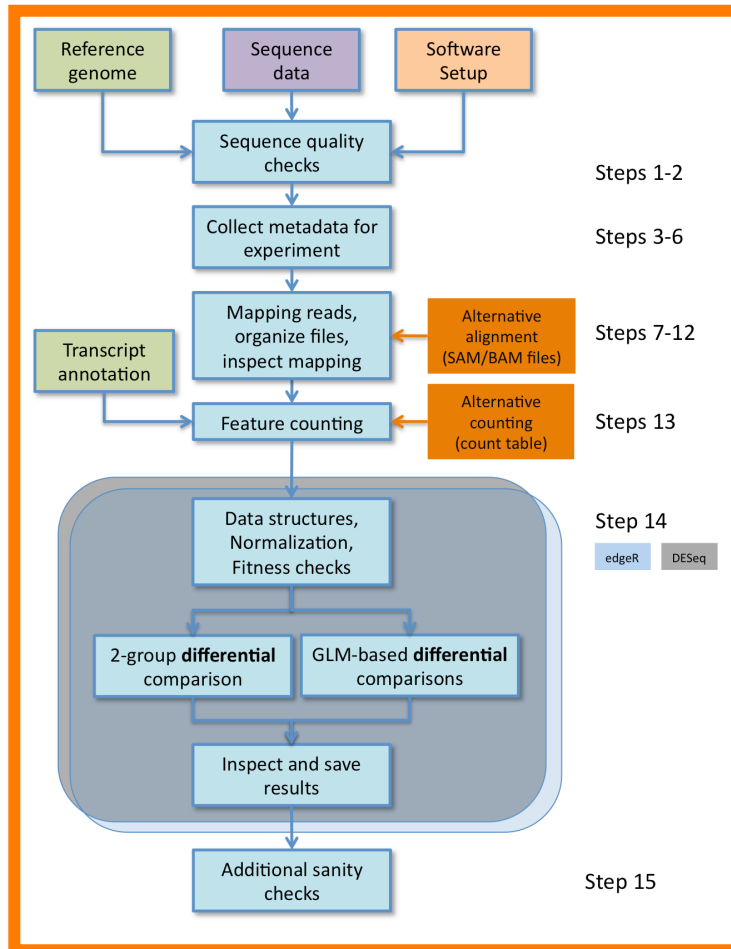


Transcriptome analysis of human tissues and cell lines reveals one dominant transcript per gene

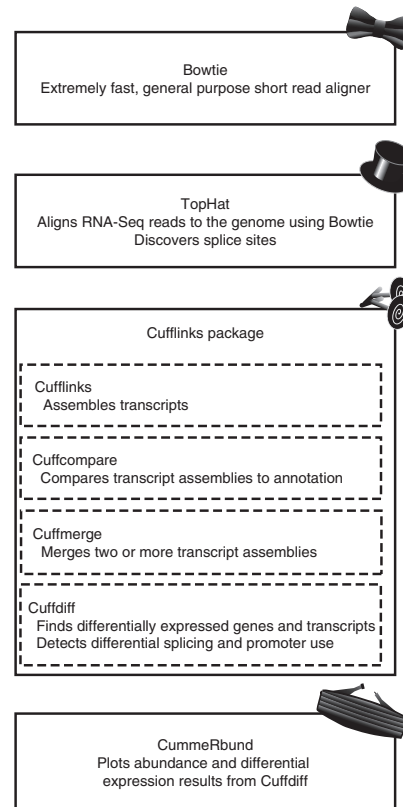
Mar Gonzàlez-Porta<sup>1</sup>, Adam Frankish<sup>2</sup>, Johan Rung<sup>1</sup>, Jennifer Harrow<sup>2</sup> and Alvis Brazma<sup>1\*</sup>



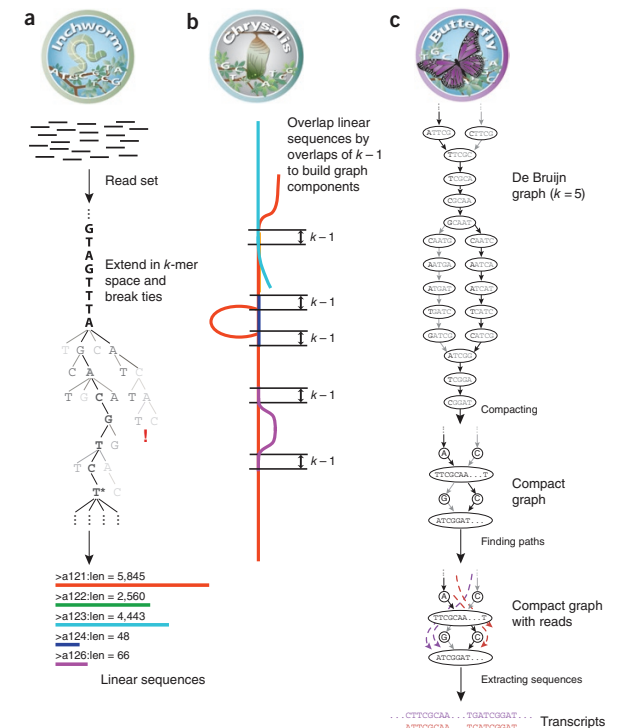
## edgeR, DESeq



## cufflinks, cuffdiff



## Trinity



**Figure 1** Overview of Trinity. (a) Inchworm assembles the read data set (short black lines, top) by greedily searching for paths in a  $k$ -mer graph (middle), resulting in a collection of linear contigs (color lines, bottom), with each  $k$ -mer present only once in the contigs. (b) Chrysalis pools contigs (colored lines) if they share at least one  $k-1$ -mer and if reads span the junction between contigs, and then it builds individual de Bruijn graphs from each pool. (c) Butterfly takes each de Bruijn graph from Chrysalis (top), and trims spurious edges and compacts linear paths (middle). It then reconciles the graph with reads (dashed colored arrows, bottom) and pairs (not shown), and outputs one linear sequence for each splice form and/or paralogous transcript represented in the graph (bottom, colored sequences).



## Counting: a few considerations (gene-level)

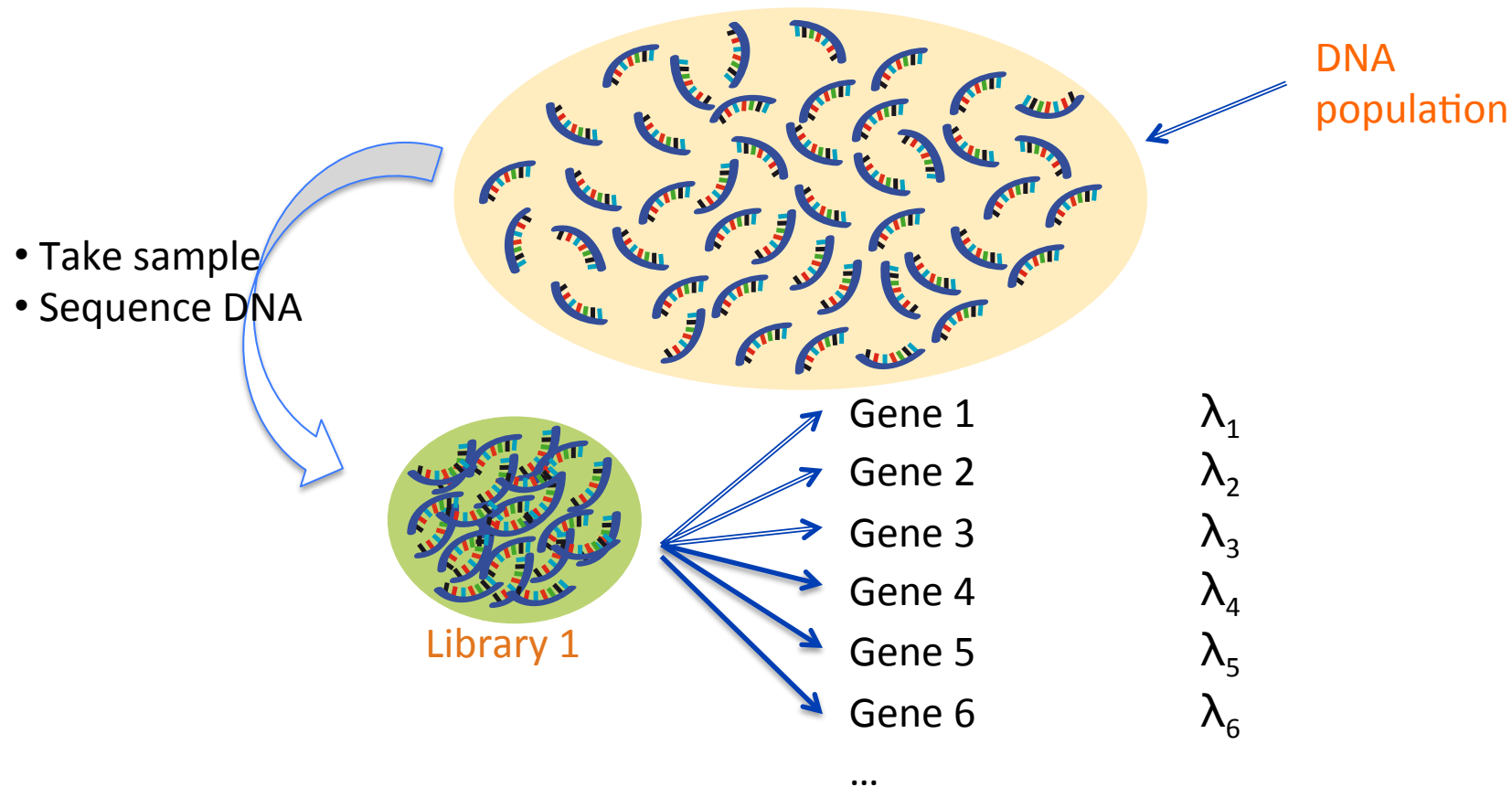
All the downstream statistical methods start with a count table.

- annotation-based? What about novel genes?
- gene-level versus transcript-level? versus exon-level?
- ambiguities

	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous	gene_A	gene_A
	ambiguous	ambiguous	ambiguous



# Sampling reads from population of DNA fragment is multinomial





# For a single gene, it's a coin toss, i.e. Binomial



$$Y_i \sim \text{Binomial}(M, \lambda_i)$$

$Y_i$  - observed number of reads for gene i

$M$  - total number of sequences

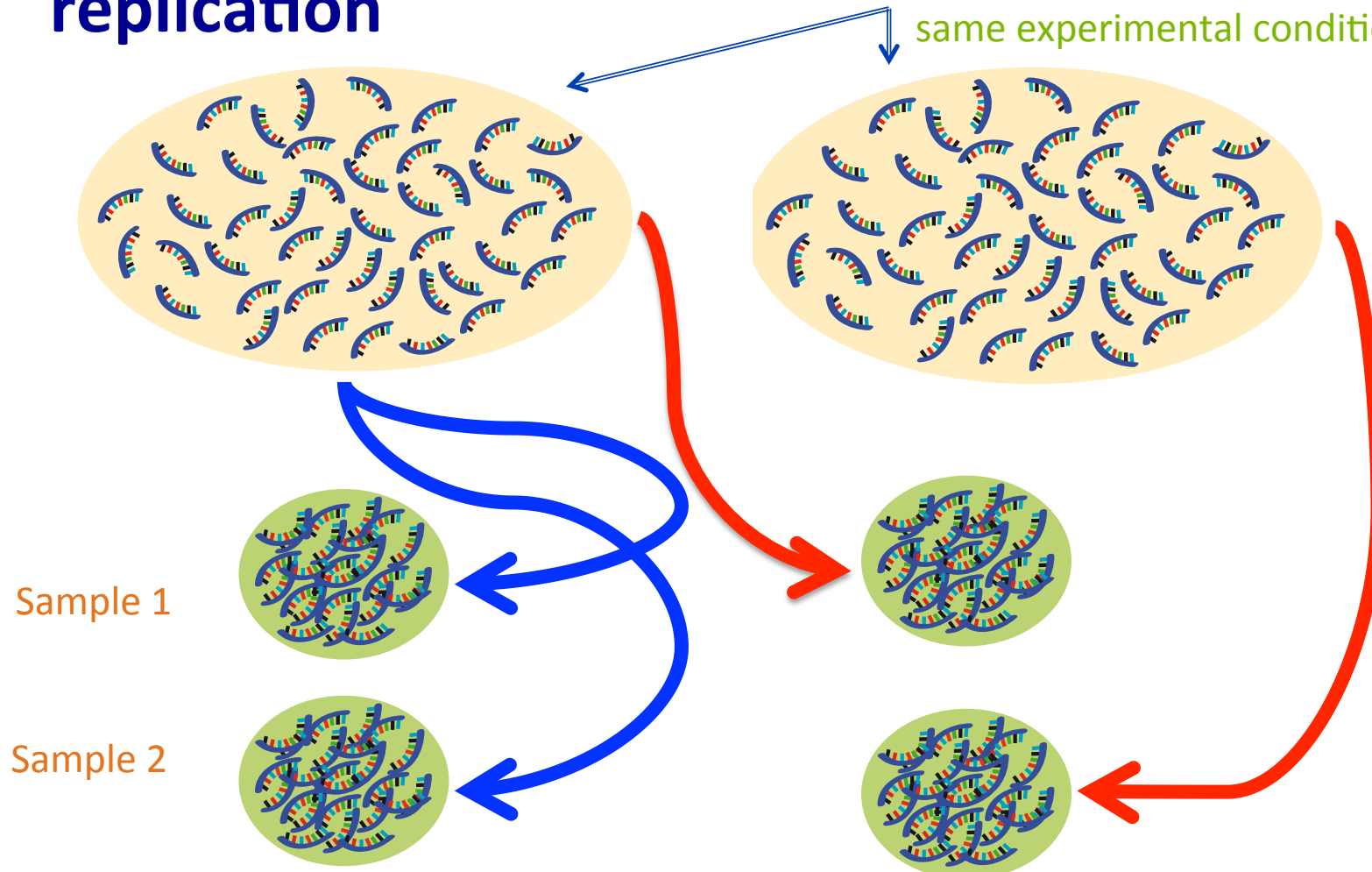
$\lambda_i$  - proportion

Large  $M$ , small  $\lambda_i \rightarrow$  approximated well by Poisson( $\mu_i = M \cdot \lambda_i$ )



# Technical replication versus **biological** replication

Independent DNA populations from same experimental condition

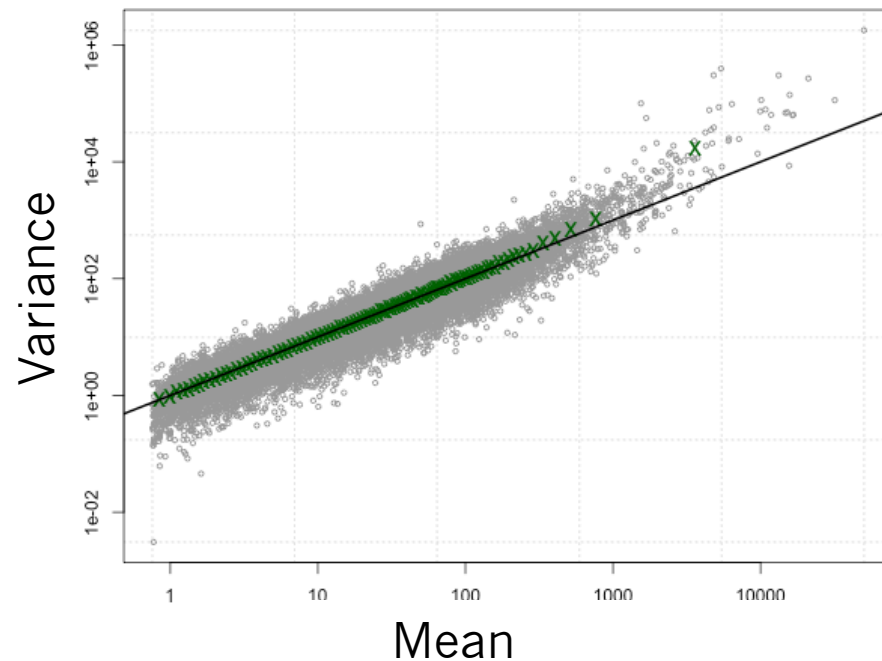






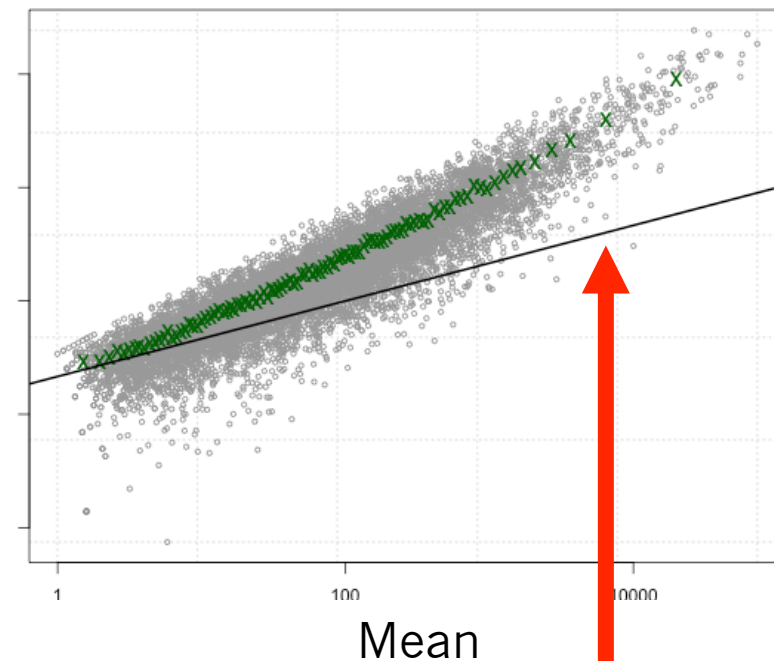
# Mean-Variance plots: What we see in real data

## Technical replicates



Data from Marioni et al. *Genome Research* 2008

## Biological replicates



Data from Parikh et al.  
*Genome Biology* 2010

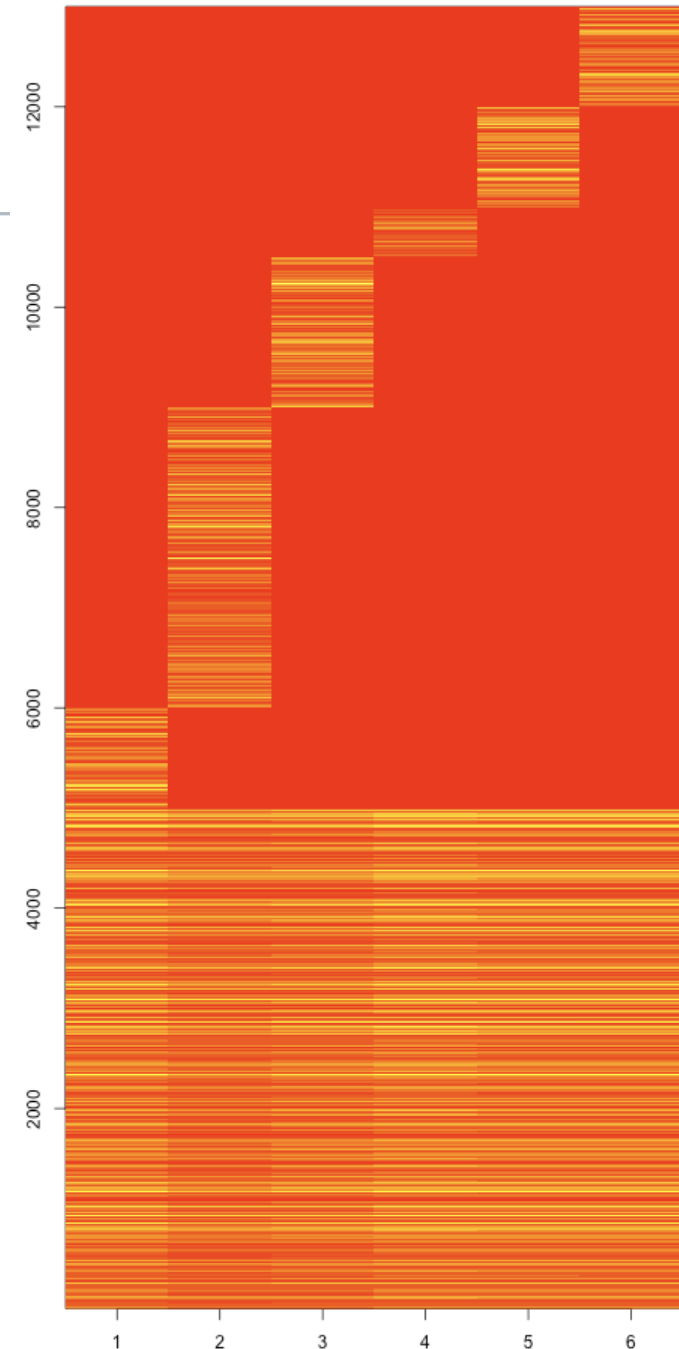
mean=variance  
(Poisson assumption)



## Normalization: “Composition” or “Diversity” can affect read depth

- Hypothetical example: Sequence 6 libraries to the **same** depth, with varying levels of *unique-to-sample* counts
- Read depth is affected not only by expression (and length), but also expression levels of other genes
- Composition can induce (sometimes significant) differences in counts

Red=low, goldenyellow=high





$M_j$  = library size  
 $\lambda_{ij}$  = relative abundance of  
feature  $i$

Poisson describes technical variation:

$$Y_{ij} \sim \text{Pois}(M_j * \lambda_{ij})$$

$$\text{mean}(Y_{ij}) = \text{variance}(Y_{ij}) = M_j * \lambda_{ij}$$

Negative binomial models **biological** variability using the dispersion parameter  $\varphi$ :

$$Y_{ij} \sim \text{NB}(\mu_{ij} = M_j * \lambda_{ij}, \varphi_i)$$

Same mean, variance is quadratic in the mean:

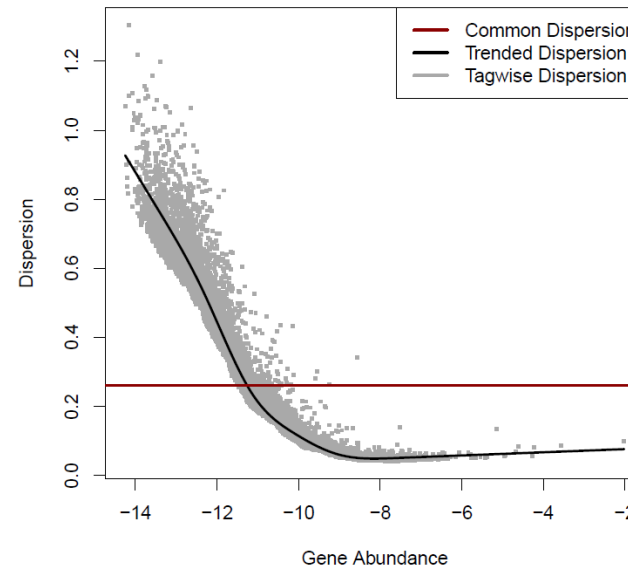
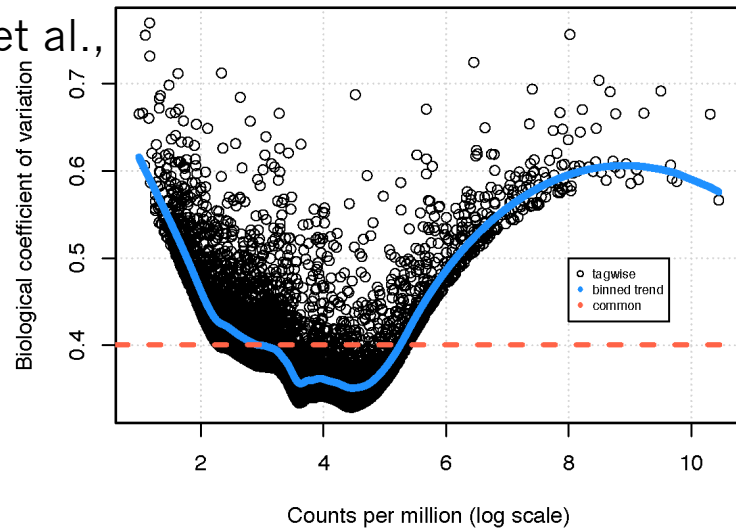
$$\text{variance}(Y_{ij}) = \mu_{ij} (1 + \mu_{ij} \varphi_i)$$

Critical parameter to estimate: dispersion



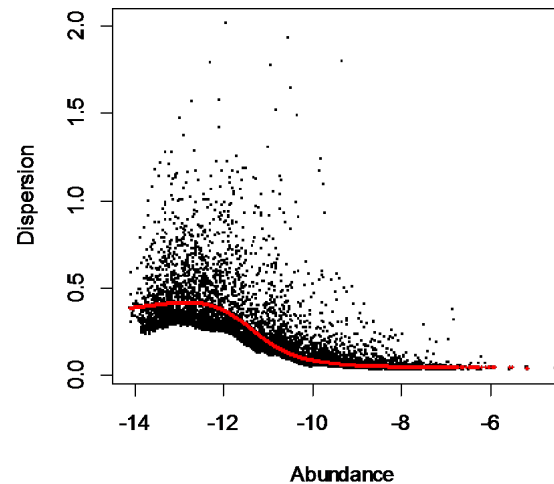
# edgeR dispersion estimation: moderate towards trend

Data:  
Tuch et al.,  
2008



Mouse hemopoietic stem cells, (Samir Taoudi)

Advantage: share information, but genes are allowed to have their own variance.



Mouse lymphomas (Stan Lee)



Response: negative binomial with dispersion fixed (to make it in the exponential family).

Link function (relate mean of response to linear combination of parameters)

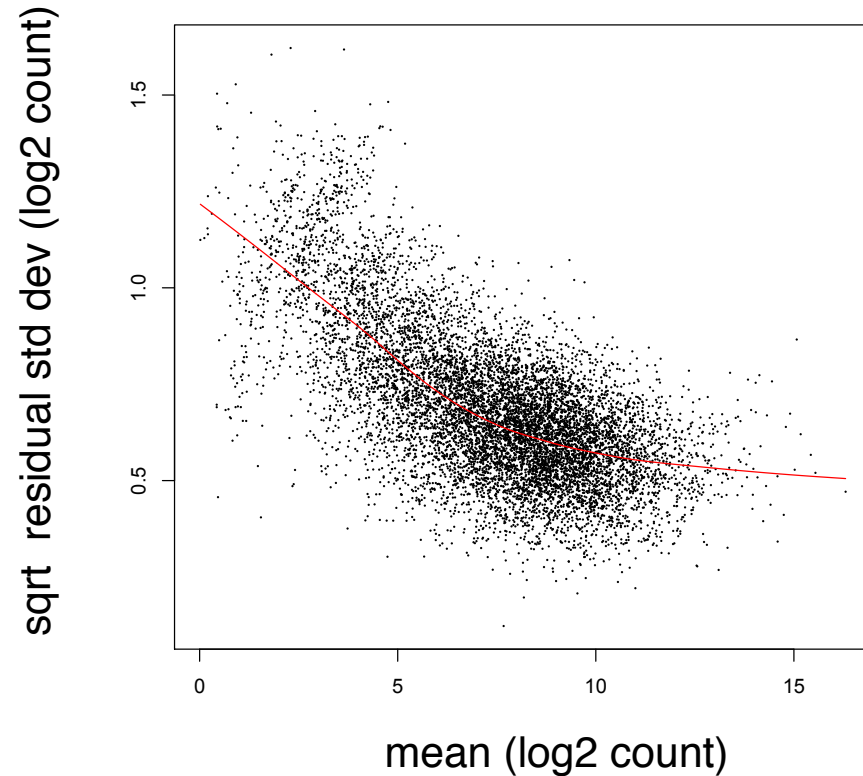
For example:

$$Y_i \sim \text{NB}(\mu_i, \phi) \quad X \quad \text{– design matrix}$$
$$\quad \quad \quad \ln() \quad \text{– link function}$$
$$\mathbf{X}\boldsymbol{\beta} = \ln(\mu) \quad \boldsymbol{\beta} \quad \text{– parameters}$$

Applicability to a wide range of designs



- Converts discrete counts to (log-cpm)
- Removes trend in the variance of counts
- Estimate variances and use inverse as **weight**

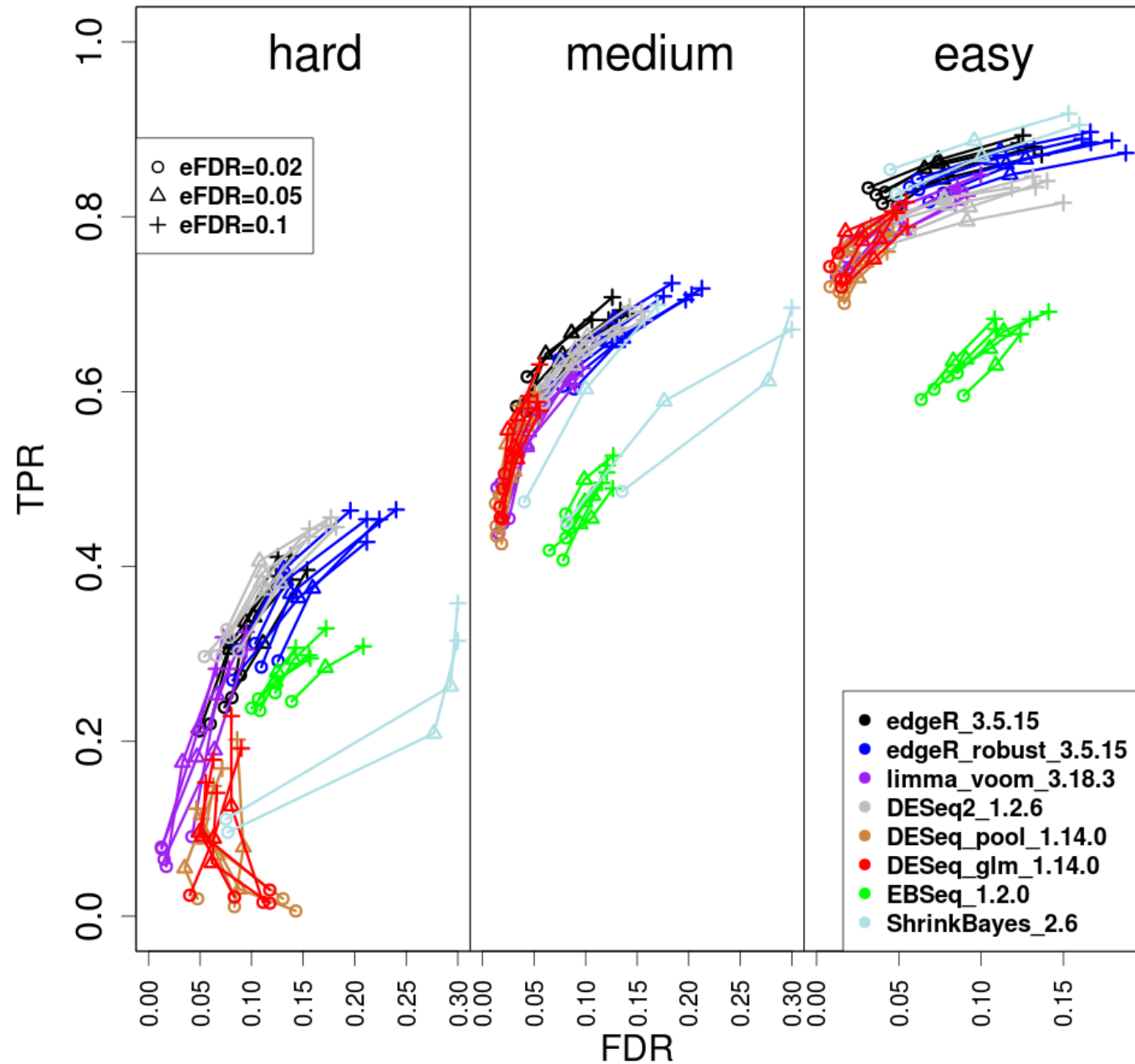


*Voom: precision weights unlock linear model analysis tools for RNA-seq read counts.* Law et al. 2014. *Genome Biology*. 2014, 15:R29.



Simulation-based comparisons: do methods achieve their FDRs?

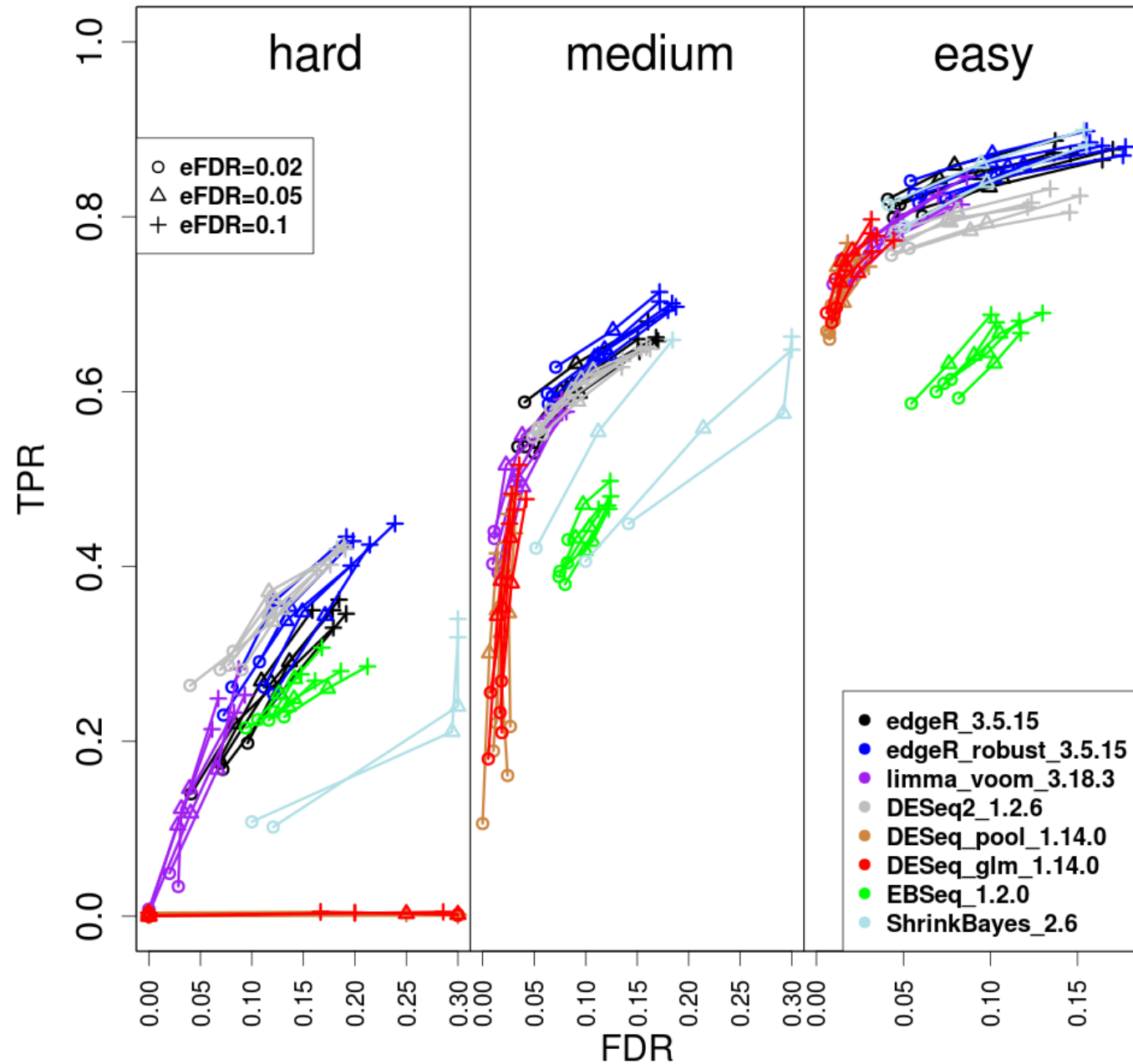
# (a) No outliers/pickrell/5vs5





Simulation-based comparisons: do methods achieve their FDRs?

## (b) 10% outliers/S/pickrell/5vs5





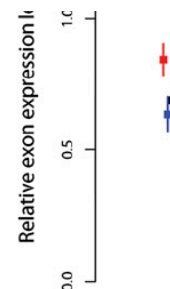


# Beyond differential expression: differential splicing

## Prediction of alternative isoforms from exon expression levels in RNA-Seq experiments

Hugues Richard<sup>1,\*</sup>, Marcel H. Schulz<sup>1,2</sup>, Marc Sultan<sup>3</sup>, Asja Nürnberg<sup>3</sup>, Sabine Schrinner<sup>3</sup>, Daniela Balzereit<sup>3</sup>, Emilie Dagand<sup>3</sup>, Axel Rasche<sup>3</sup>, Hans Lehrach<sup>3</sup>, Martin Vingron<sup>1</sup>, Stefan A. Haas<sup>1</sup> and Marie-Laure Yaspo<sup>3</sup>

<sup>1</sup>Department of Computational Molecular Biology, Max Planck Institute for Molecular Genetics, Ihnestr. 73,  
<sup>2</sup>International Max Planck Research School for Computational Biology and Scientific Computing, and  
<sup>3</sup>Department of Vertebrate Genomics, Max Planck Institute for Molecular Genetics, Ihnestr. 73, 14195 Berlin, Germany

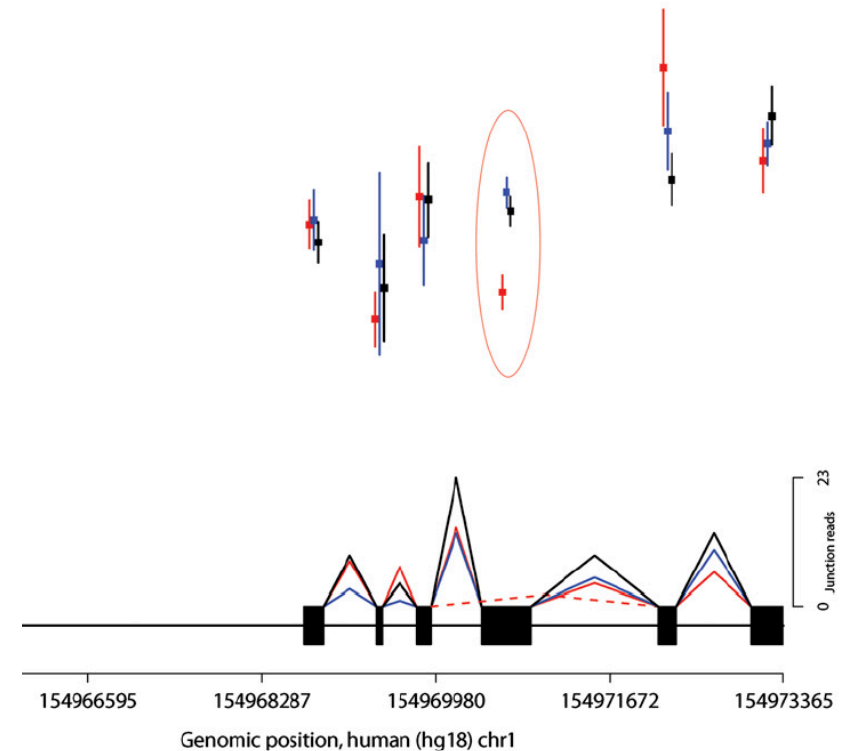


## Sex-specific and lineage-specific alternative splicing in primates

Ran Blekhan, <sup>1,4,5</sup> John C. Marioni, <sup>1,4,5</sup> Paul Zumbo, <sup>2</sup> Matthew Stephens, <sup>1,3,5</sup> and Yoav Gilad <sup>1,5</sup>

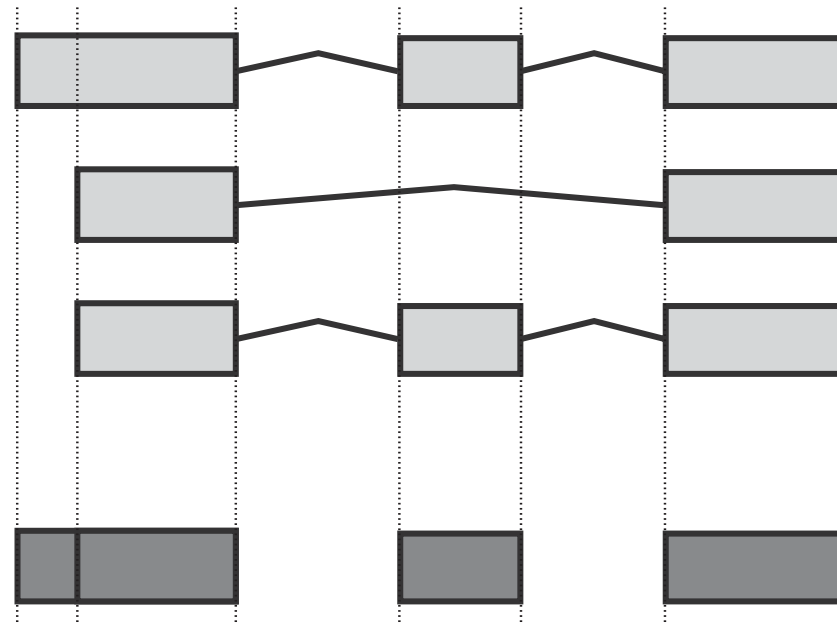
<sup>1</sup>Department of Human Genetics, University of Chicago, Chicago, Illinois 60637, USA; <sup>2</sup>Keck Biotechnology Laboratory, New Haven, Connecticut 06511, USA; <sup>3</sup>Department of Statistics, University of Chicago, Chicago, Illinois 60637, USA

CGI-41





## Counting: a few considerations (exon-level)



**Figure 1.** Flattening of gene models: This (fictional) gene has three annotated transcripts involving three exons (light shading), one of which has alternative boundaries. We form counting bins (dark shaded boxes) from the exons as depicted; the exon of variable length gets split into two bins.



## DEXSeq – general structure

We use generalized linear models (GLMs) (McCullagh and Nelder 1989) to model read counts. Specifically, we assume  $K_{ijl}$  to follow a negative binomial (NB) distribution:

$$K_{ijl} \sim NB\left(\text{mean} = s_j \mu_{ijl}, \text{dispersion} = \alpha_{il}\right), \quad (1)$$

where  $\alpha_{il}$  is the dispersion parameter (a measure of the distribution's spread; see below) for counting bin ( $i$ ,  $l$ ), and the mean is predicted via a log-linear model as

$$\log \mu_{ijl} = \beta_i^G + \beta_{il}^E + \beta_{i\rho_j}^C + \beta_{i\rho_j l}^{EC}. \quad (2)$$

$i$  – gene

$j$  – sample ...  $\rho_j$  is condition (categorical)

$l$  – bin

$\beta^G$  – baseline “expression strength”

$\beta^E$  – “exon” (bin) effect

$\beta^C$  – condition effect

$\beta^{EC}$  – condition x “exon” interaction

---

### Method

## Detecting differential usage of exons from RNA-seq data

Simon Anders,<sup>1,2</sup> Alejandro Reyes,<sup>1</sup> and Wolfgang Huber

*European Molecular Biology Laboratory, 69111 Heidelberg, Germany*