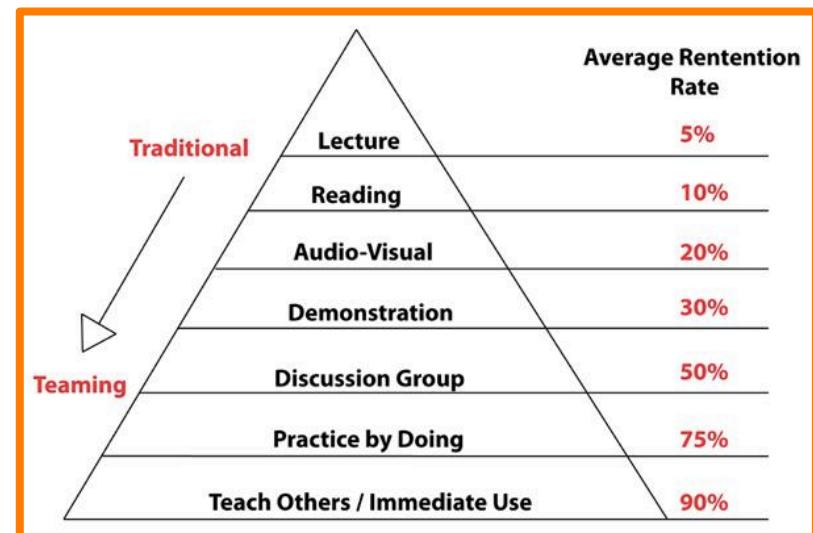




## Computational Genomics

**RNA-seq – 22 Sep 2016**

- 1. WHAT (9.00-10.30)**
- 2. HOW (11.00-12.30)**
- 3. WHY (13.30-15.00)**
- 4. DO (15.30-17.00)**



Mark D. Robinson, Statistical Genomics, IMLS



@markrobinsonca



University of  
Zurich<sup>UZH</sup>

Institute of Molecular Life Sciences

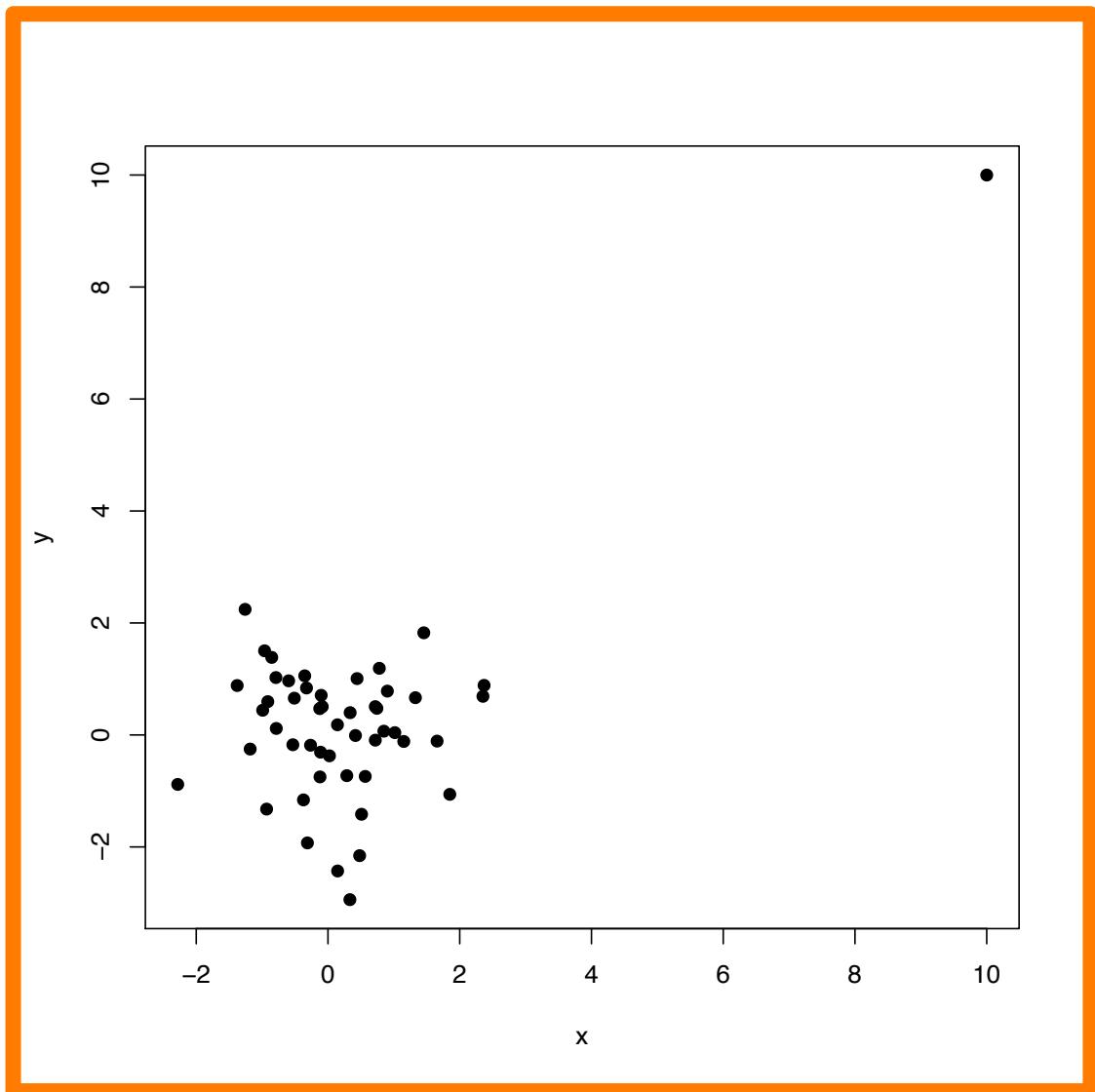
---

**movo.ch**

LU HU DA HE



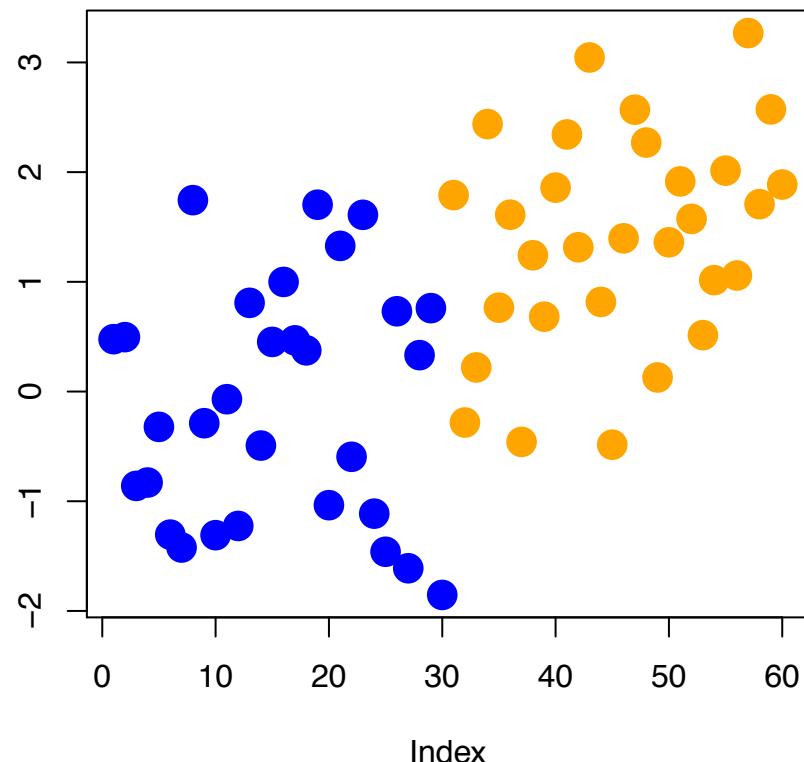
**In your view, what best describes the associations shown in the plot of 'x' and 'y' ?  
(Pearson = linear correlation; Spearman = rank correlation)**



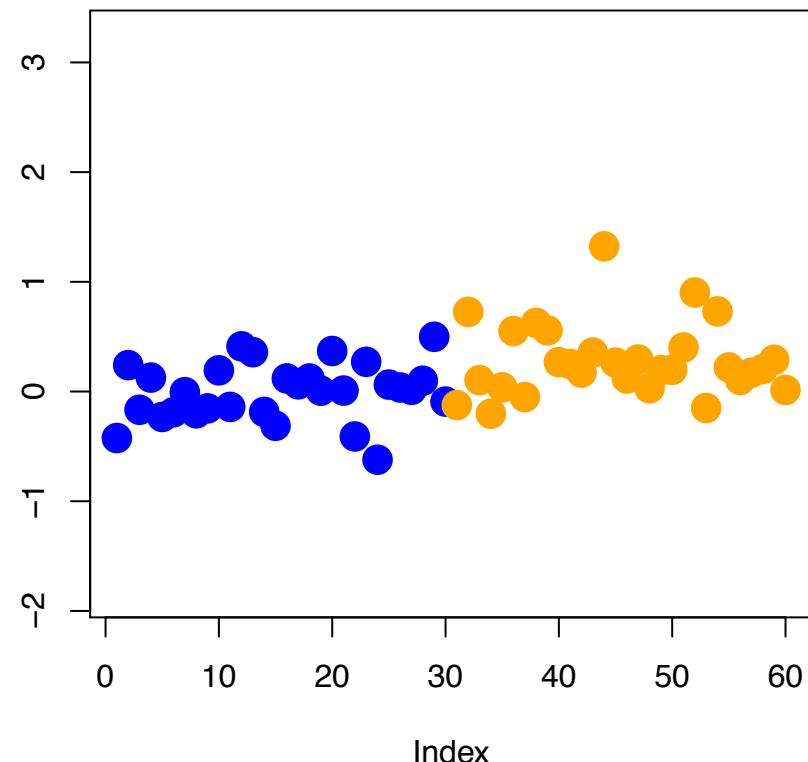


Which plot highlights more statistical evidence for a change in the population means (between orange and blue)?

A



B



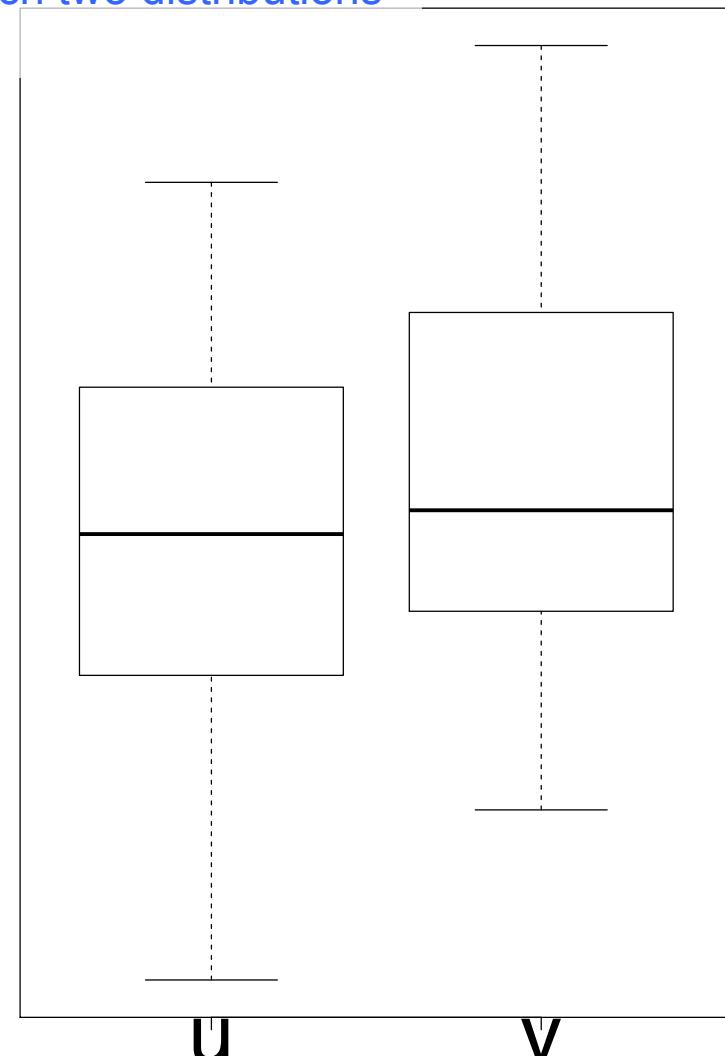
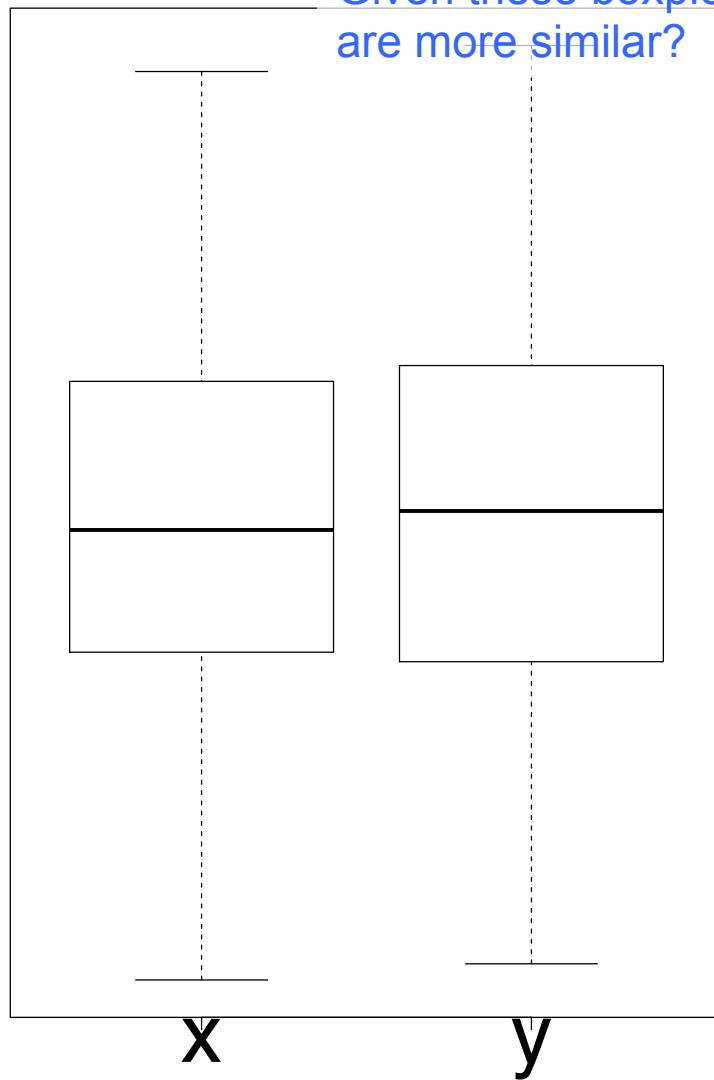


$$X = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$



Given these boxplots, which two distributions are more similar?

75<sup>th</sup> percentile  
median  
25<sup>th</sup> percentile





**1** 
$$\frac{(\bar{x}_1 - \bar{x}_2) - d_0}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$$

**2** 
$$\sum^k \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$$

**3** 
$$\frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\hat{p}(1 - \hat{p})(\frac{1}{n_1} + \frac{1}{n_2})}}$$



University of  
Zurich<sup>UZH</sup>

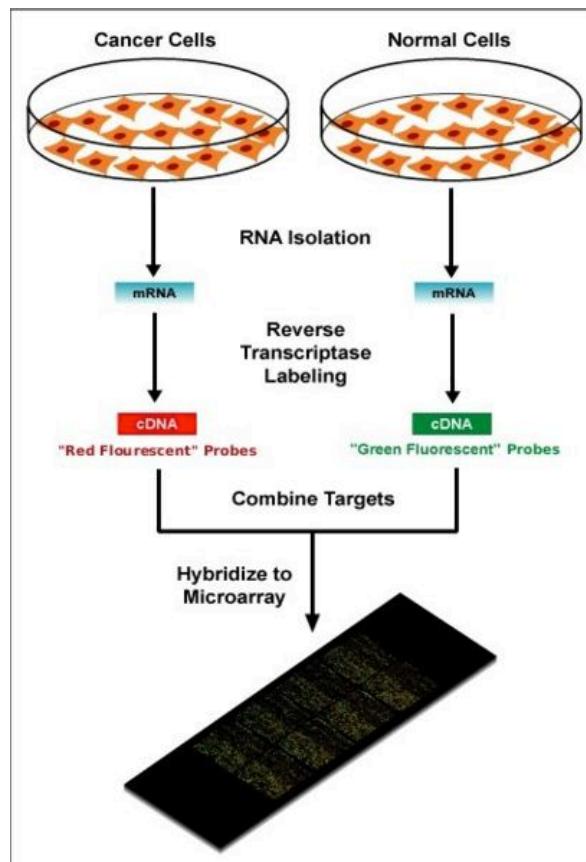
Institute of Molecular Life Sciences

---

# Quick intro

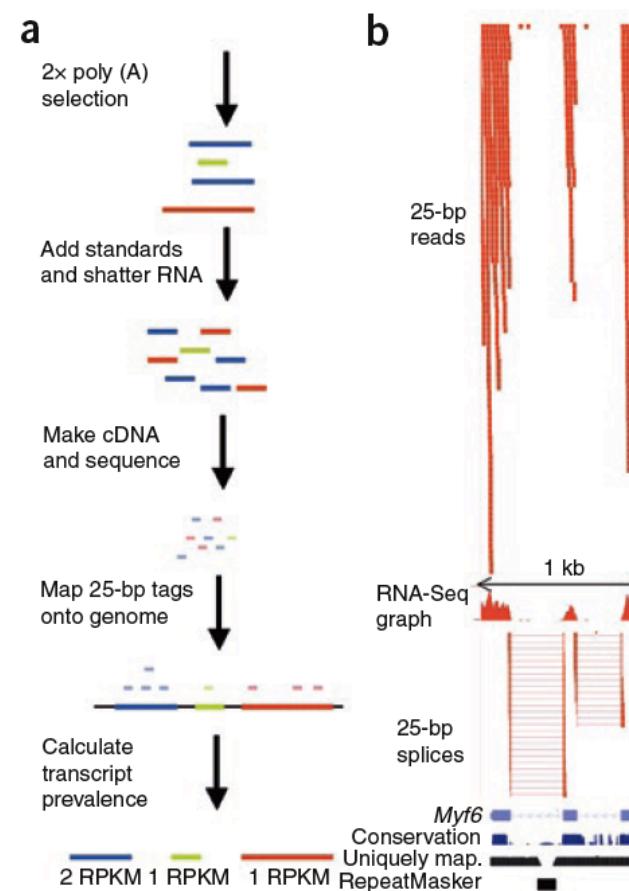


## Abundance by Fluorescence Intensity



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

## Abundance by Counting

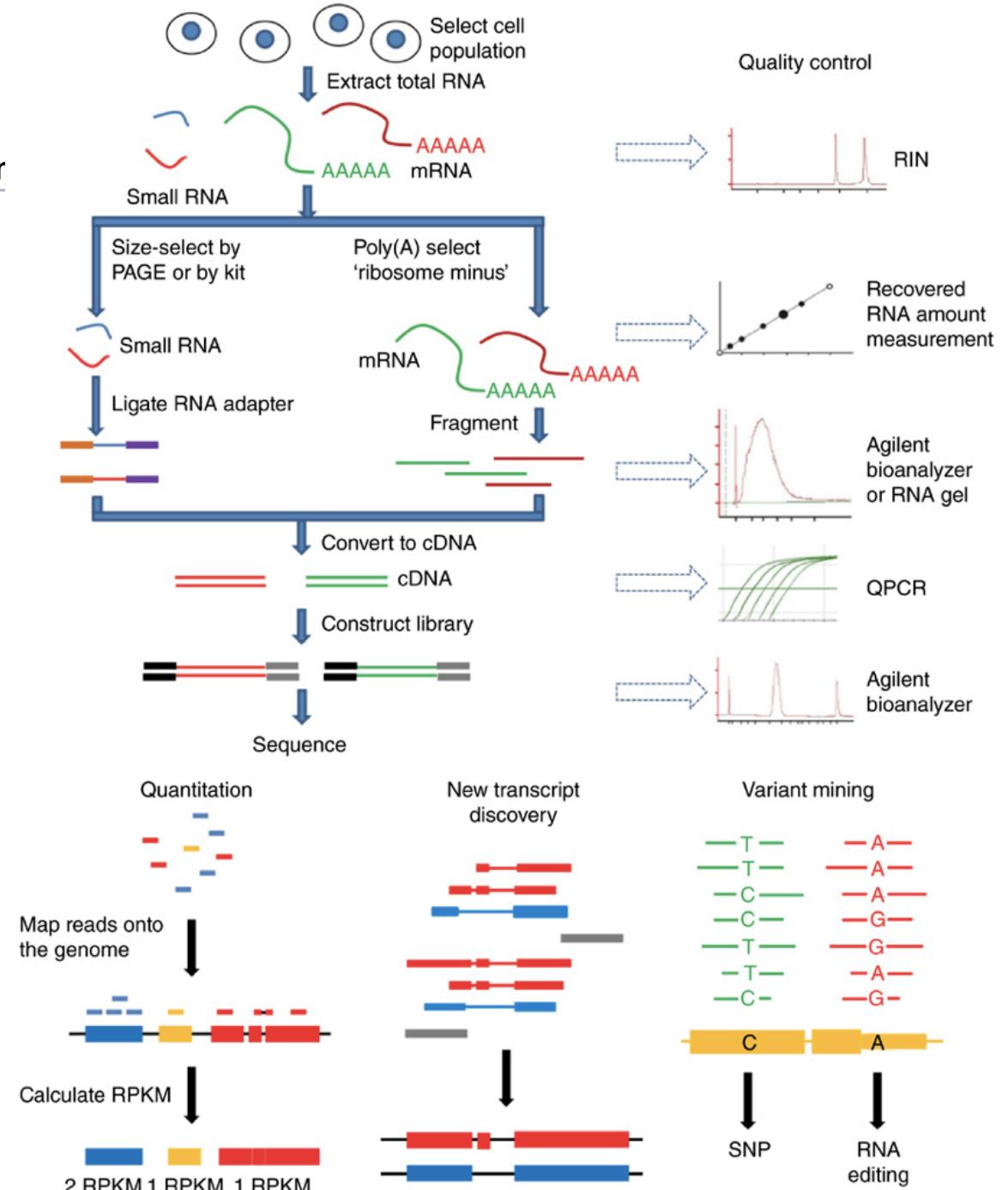


Mortazavi et al., Nature Methods, 2008



University of  
Zurich<sup>UZH</sup>

Institute of Molecular Life Sciences





## Brainstorm session

What all can we do with RNA-seq ?

20-30 minutes: Break into groups. Research/discuss your topic within your group. Find a few slides/figures from the literature/web (and/or some discussion points). A nominated representative from your group can lead discussion. Each group add links/notes to **Etherpad**.

### Groups:

1. Differential expression
2. Differential splicing
3. allele-specific expression
4. *de novo* discovery/assembly
5. expression quantitative trait loci
6. RNA editing



University of  
Zurich<sup>UZH</sup>

Institute of Molecular Life Sciences

---

## Experimental design

Replication

Randomization

Blocking

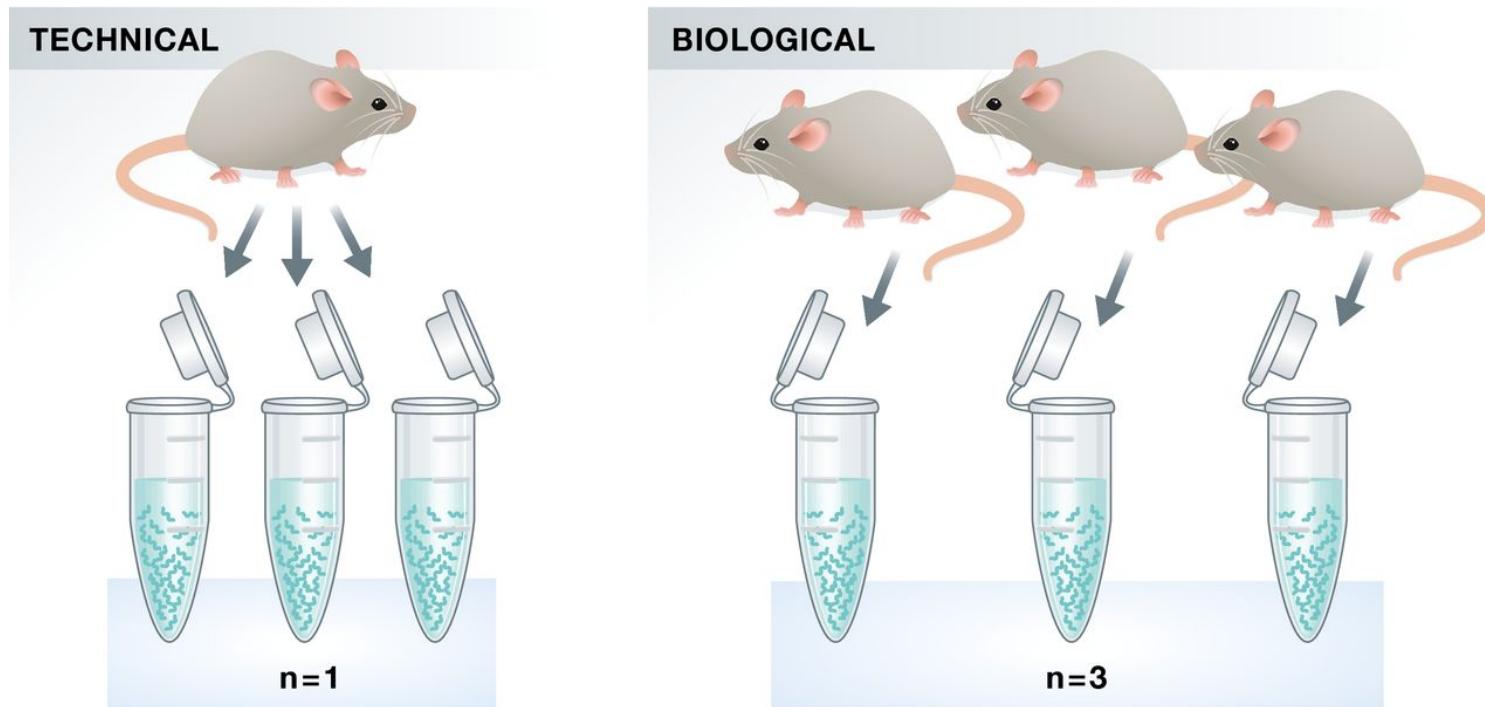


University of  
Zurich<sup>UZH</sup>

From <http://emboj.embopress.org/content/early/2015/09/21/embj.201592958>

Institute of Molecular Life Sciences

## Biological versus technical replication



What do we want? Why?

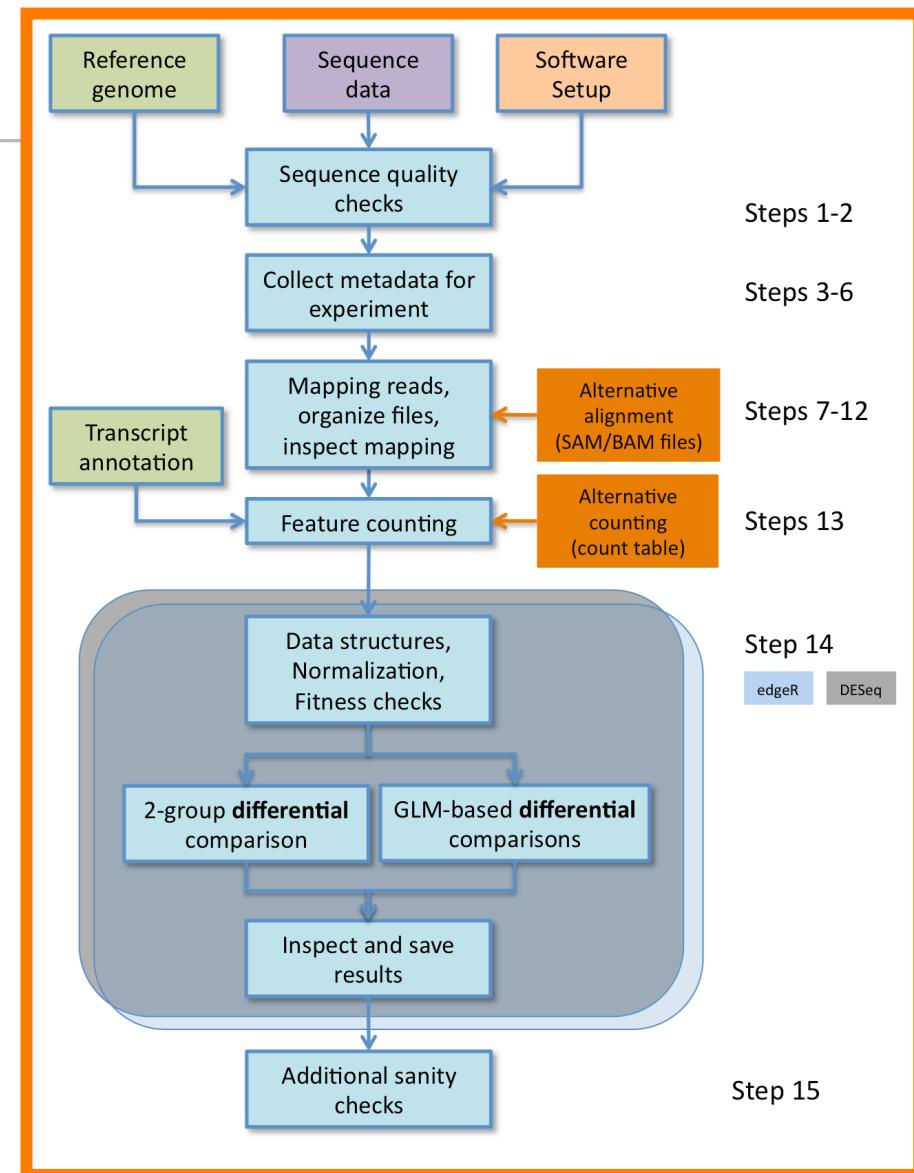


## Standard pipeline

Issues/queries:

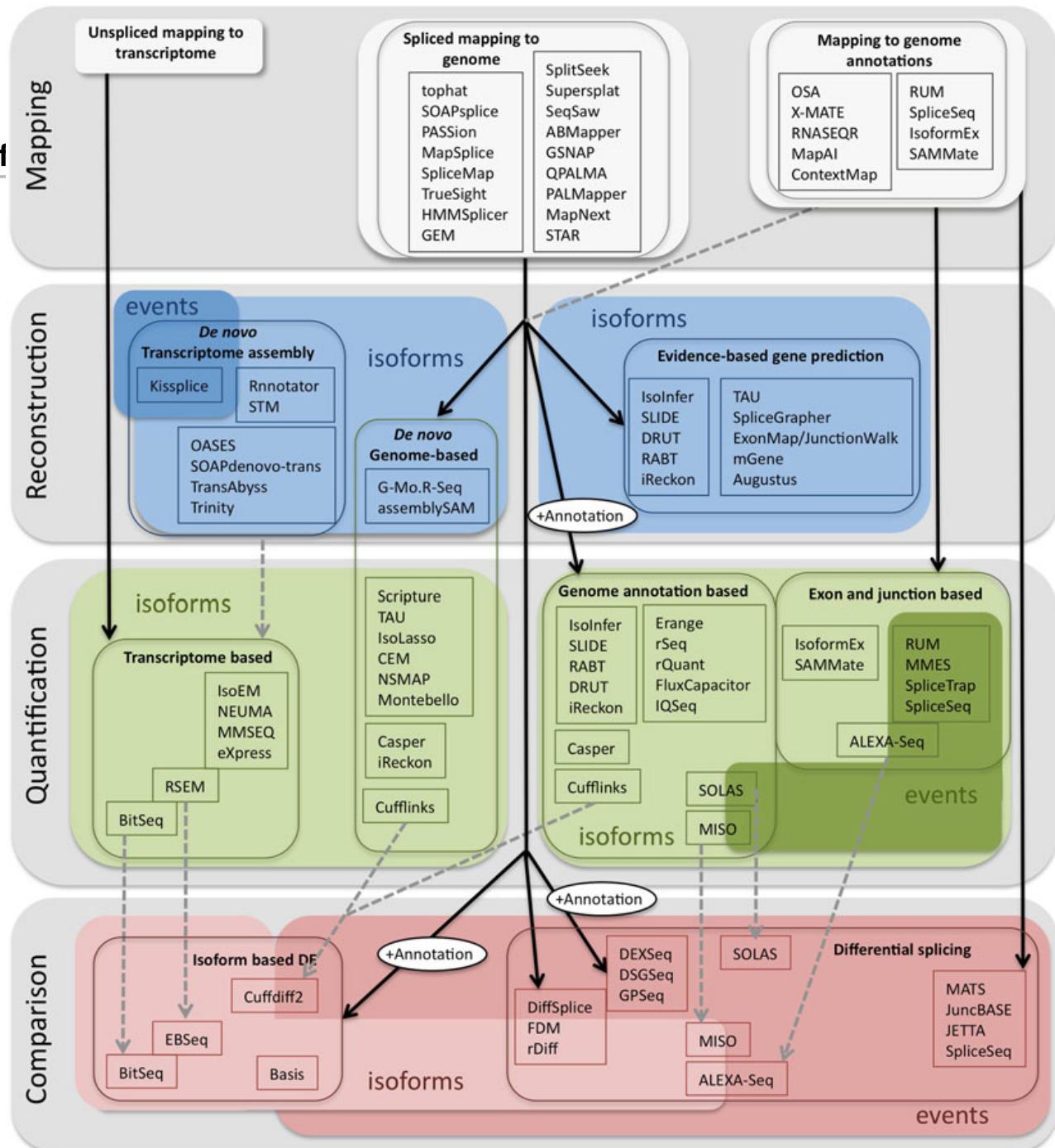
Alignment – to genome? to transcriptome ?  
versus alignment-free or pseudo-alignment ?

“Counting” – what does that really mean ?





## Tools



### Chapter 26

Methods to Study Splicing from High-Throughput RNA Sequencing Data

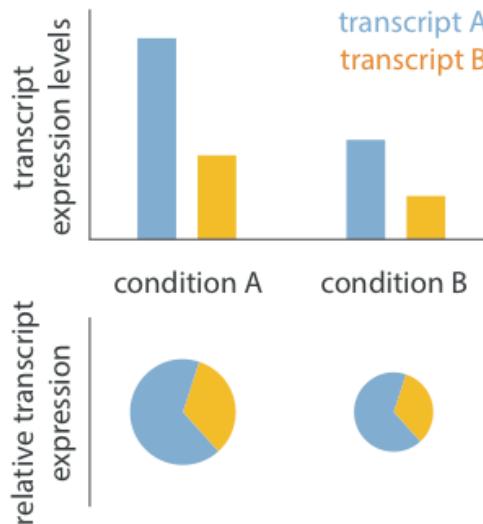
Gael P. Alamancos, Eneritz Agirre, and Eduardo Eyras

New: salmon, kallisto, etc.

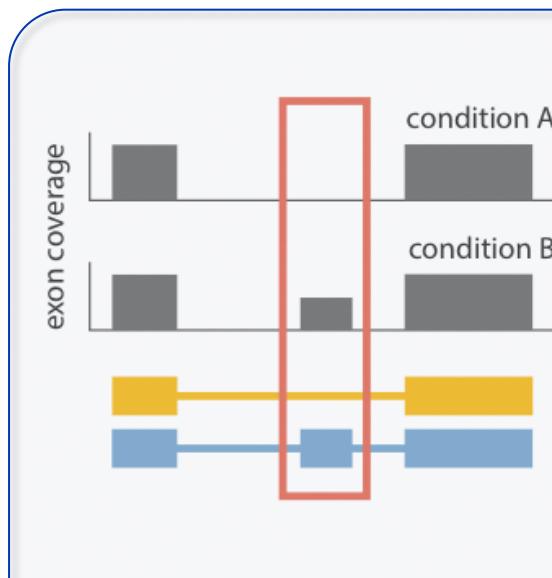


## Some terms: DTE, DEU, DTU .. DGE

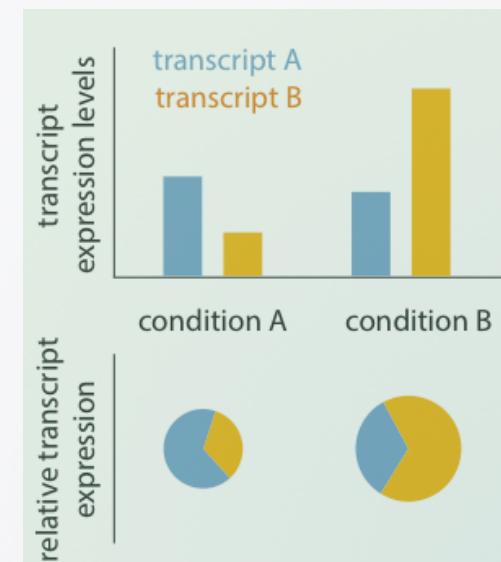
Differential transcript expression (DTE)



Differential exon usage (DEU)



Differential transcript usage (DTU)

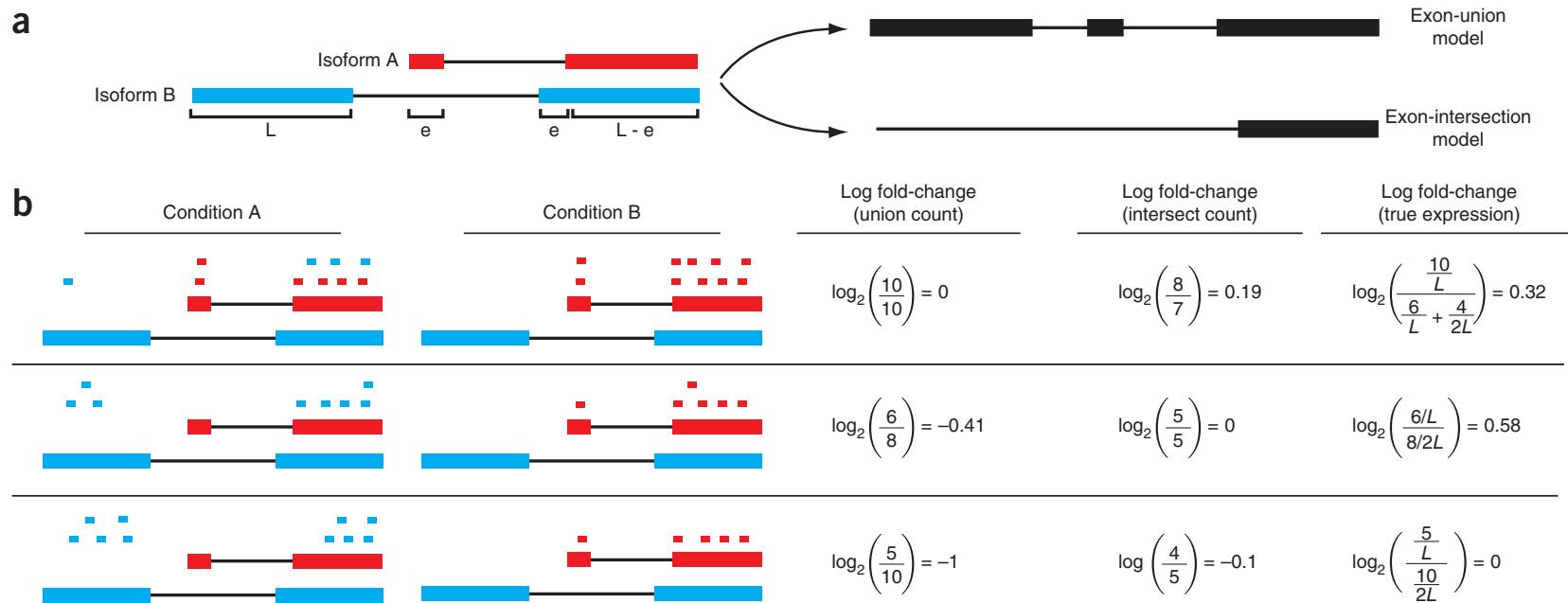


differential splicing



## Caveat: gene-level counting can go wrong, but often not bad

Trapnell et al. 2013 Nat Biotech



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

Mar González-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>



University of  
Zurich<sup>UZH</sup>

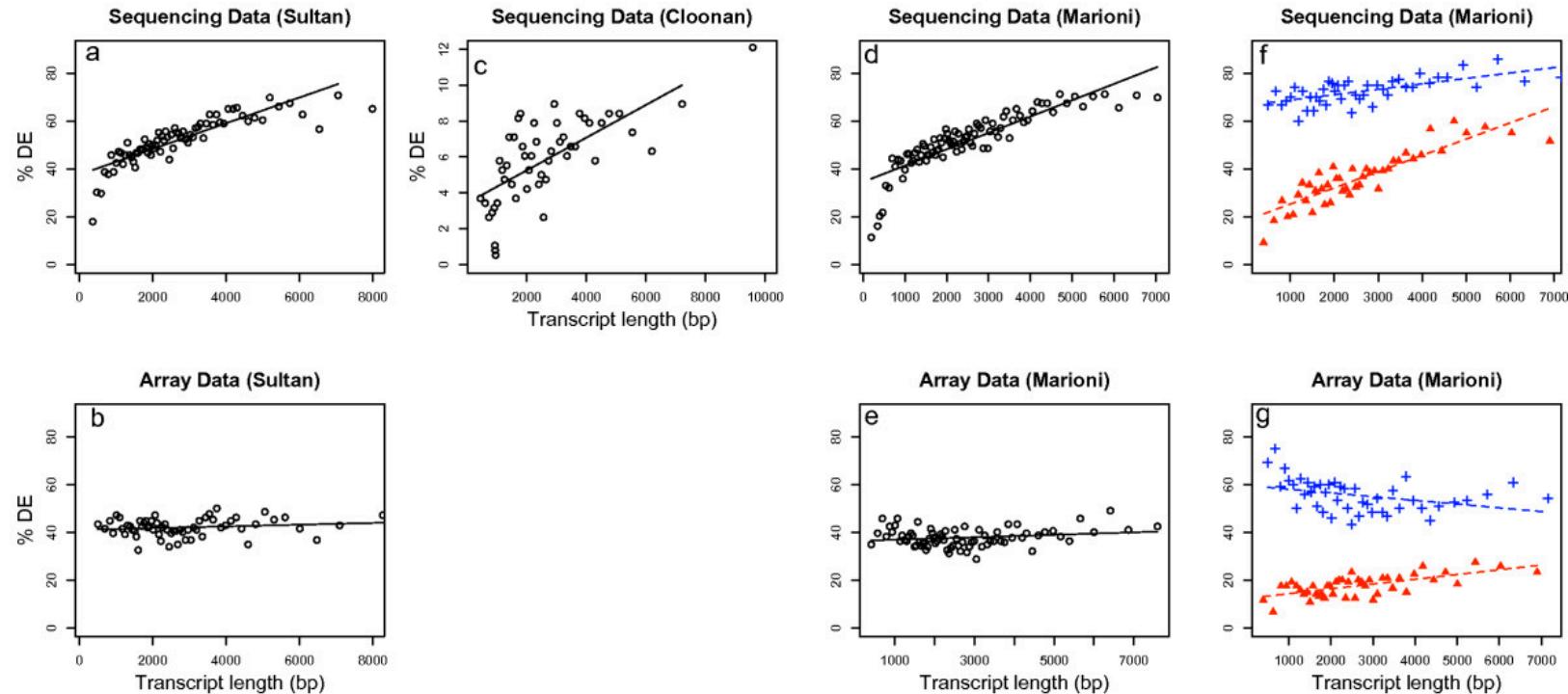
Institute of Molecular Life Sciences

---

# SOME OTHER CAVEATS, LIMITATIONS AND REMARKS



## RNA-seq length bias





# Reconstructing transcripts from current generation RNA-seq?

## Assessment of transcript reconstruction methods for RNA-seq

Tamara Steijger<sup>1</sup>, Josep F Abril<sup>2,11</sup>, Pär G Engström<sup>1,10,11</sup>, Felix Kokocinski<sup>3,11</sup>, The RGASP Consortium<sup>4</sup>, Tim J Hubbard<sup>3</sup>, Roderic Guigó<sup>5,6</sup>, Jennifer Harrow<sup>3</sup> & Paul Bertone<sup>1,7-9</sup>

Nature Methods 2013

“Consequently, the complexity of higher eukaryotic genomes imposes severe limitations on transcript recall and splice product discrimination that are likely to remain limiting factors for the analysis of current-generation RNA-seq data.”



T I Bonner

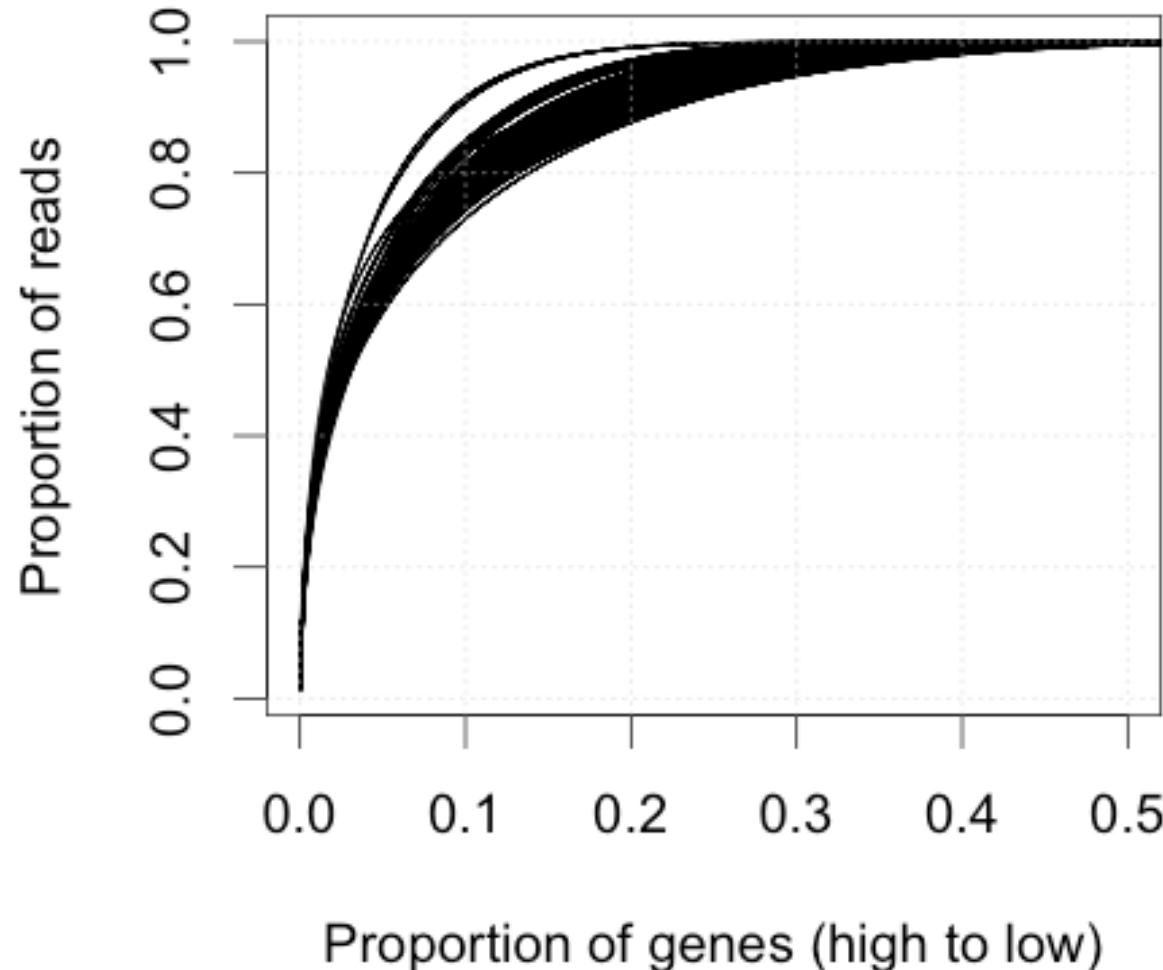
National Institute of Mental Health, Bethesda, MD, USA

## Genome-wide surveys: are we seeing events of interest ?

Alternative splicing of mRNAs occurs in the majority of human genes, and most differential splicing results in different protein isoforms with possibly different functional properties. However, there are many reported splicing variations that may be quite rare, and not all combinatorially possible variants of a given gene are expressed at significant levels. Genes of interest to pharmacologists are frequently expressed at such low levels that they are not adequately represented in genome-wide studies of transcription. In single-gene studies, data are commonly available on the relative abundance and functional significance of individual alternatively spliced exons, but there are rarely data that quantitate the relative abundance of full-length transcripts and define which combinations of exons are significant. A number of criteria for judging the significance of splice variants and suggestions for their nomenclature are discussed.



## How much sequencing goes to highly expressed genes?





University of  
Zurich <sup>UZH</sup>

Institute of Molecular Life Sciences

## Some genes are just difficult

Robert and Watson *Genome Biology* (2015) 16:177  
DOI 10.1186/s13059-015-0734-x



RESEARCH

Open Access



### Errors in RNA-Seq quantification affect genes of relevance to human disease

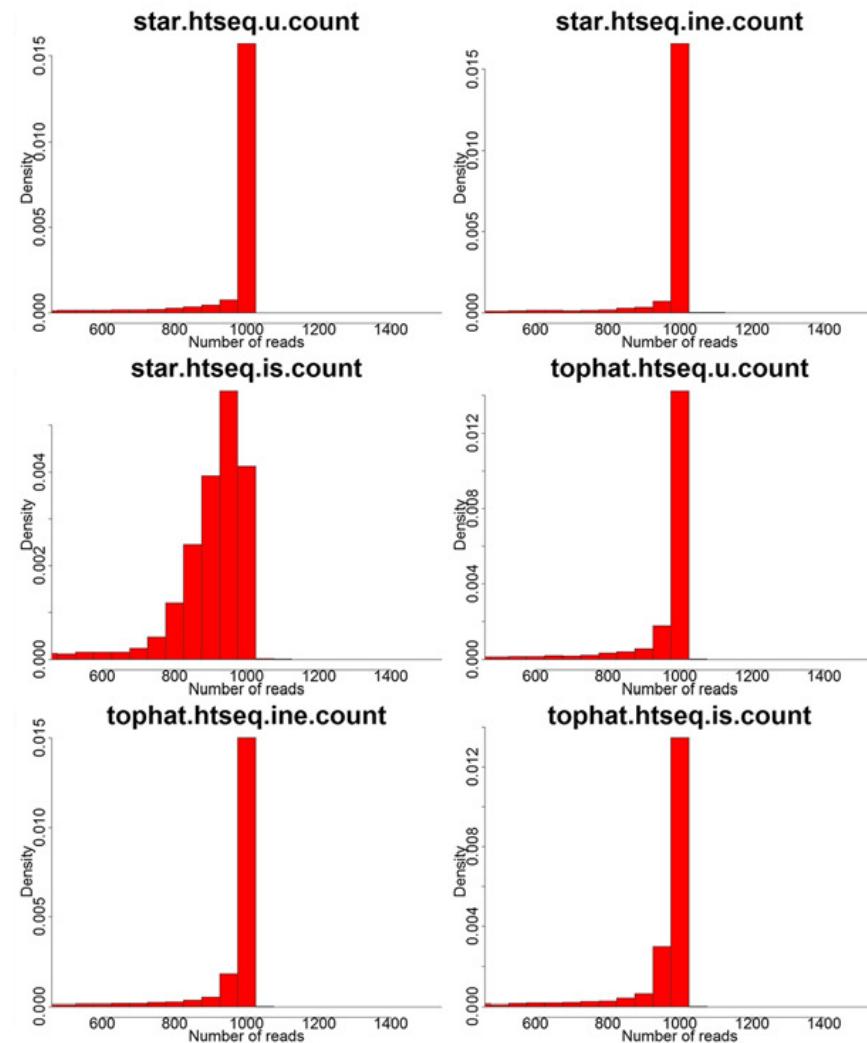
Christelle Robert<sup>1</sup> and Mick Watson<sup>2\*</sup>

#### Abstract

**Background:** RNA-Seq has emerged as the standard for measuring gene expression and is an important technique often used in studies of human disease. Gene expression quantification involves comparison of the sequenced reads to a known genomic or transcriptomic reference. The accuracy of that quantification relies on there being enough unique information in the reads to enable bioinformatics tools to accurately assign the reads to the correct gene.

**Results:** We apply 12 common methods to estimate gene expression from RNA-Seq data and show that there are hundreds of genes whose expression is underestimated by one or more of those methods. Many of these genes have been implicated in human disease, and we describe their roles. We go on to propose a two-stage analysis of RNA-Seq data in which multi-mapped or ambiguous reads can instead be uniquely assigned to groups of genes. We apply this method to a recently published mouse cancer study, and demonstrate that we can extract relevant biological signal from data that would otherwise have been discarded.

**Conclusions:** For hundreds of genes in the human genome, RNA-Seq is unable to measure expression accurately. These genes are enriched for gene families, and many of them have been implicated in human disease. We show that it is possible to use data that may otherwise have been discarded to measure group-level expression, and that such data contains biologically relevant information.

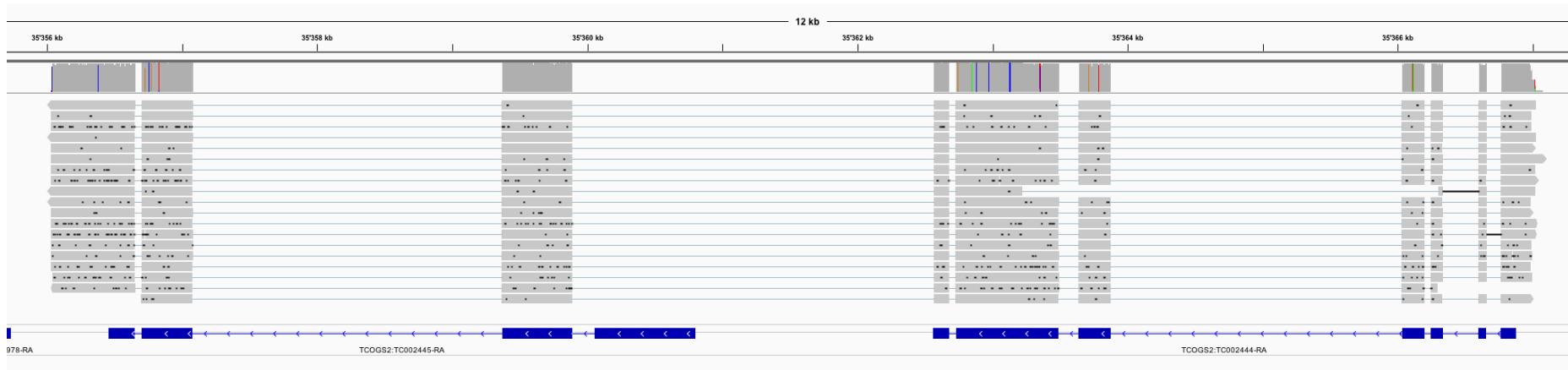




## Soon: counting full-length transcripts?

→ PacBio, Oxford Nanopore, 10x Genomics + Illumina

PacBio IsoSeq data collected  
(with D. Bopp and Functional Genomics Centre Zurich)





## Multiple testing

- In genomics, typically P-values are corrected for multiple testing
- FWER = family-wise error rate (strong)
  - control probability of making 1 FP
- FDR = false discovery rate (weak)
  - control rate of FP amongst the set of positives



University of  
Zurich<sup>UZH</sup>

Institute of Molecular Life Sciences

---

# TRANSCRIPT-LEVEL ESTIMATION

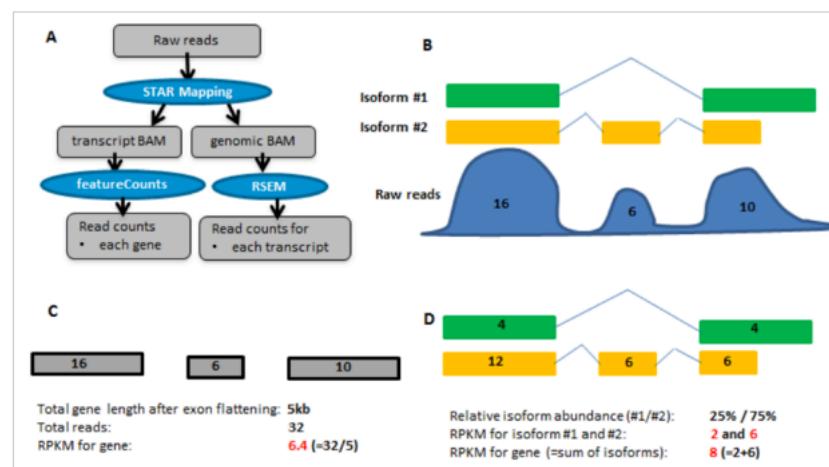


Inst

# You've been doing your RNA-Seq all wrong

Posted by: RNA-Seq Blog In Expression and Quantification November 12, 2015 5,007 Views

In recent years, RNA-seq is emerging as a powerful technology in estimation of gene and/or transcript expression, and RPKM (Reads Per Kilobase per Million reads) is widely used to represent the relative abundance of mRNAs for a gene. In general, the methods for gene quantification can be largely divided into two categories: transcript-based approach and ‘union exon’-based approach. Transcript-based approach is intrinsically more difficult because different isoforms of the gene typically have a high proportion of genomic overlap. On the other hand, ‘union exon’-based approach method is much simpler and thus widely used in RNA-seq gene quantification. Biologically, a gene is expressed in one or more transcript isoforms. Therefore, transcript-based approach is logically more meaningful than ‘union exon’-based approach. Despite the fact that gene quantification is a fundamental task in most RNA-seq studies, however, it remains unclear whether ‘union exon’-based approach for RNA-seq gene quantification is a good practice or not.

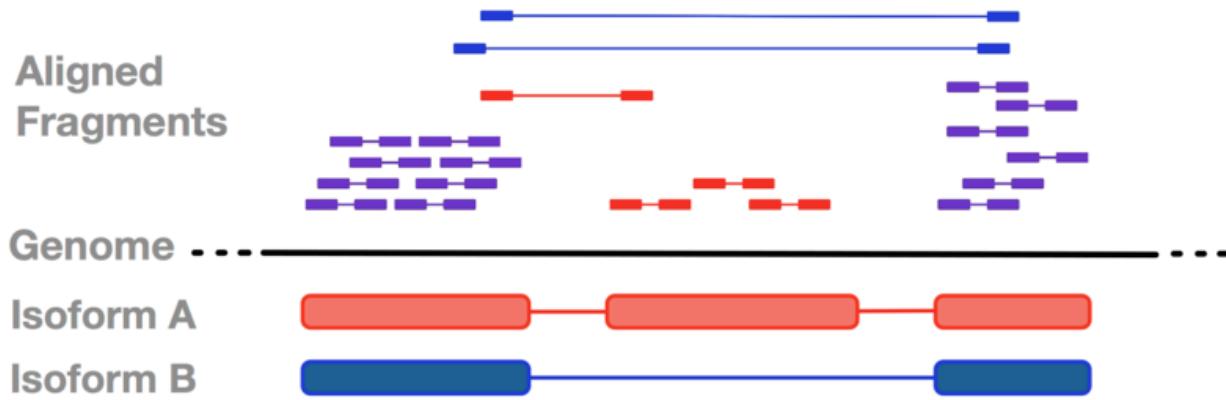


→ Transcript-based versus “union exon”-based quantification



## Transcript-level expression estimation

BitSeq  
CEM  
Cufflinks  
eXpress  
IsoEM  
MMSEQ  
RSEM  
rSEQ  
Sailfish  
Scripture  
TIGAR2  
  
salmon  
kallisto  
  
..



Open question: can you use “estimated counts” into a method that models counts?

Short answer: yes, though improvements that properly take estimation uncertainty into account may be possible.

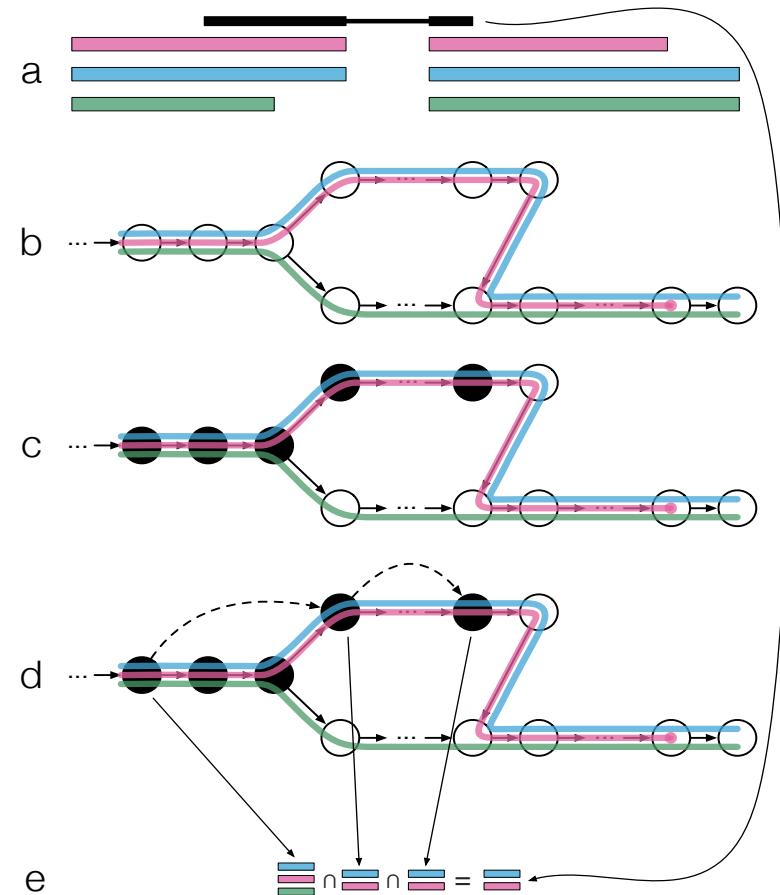
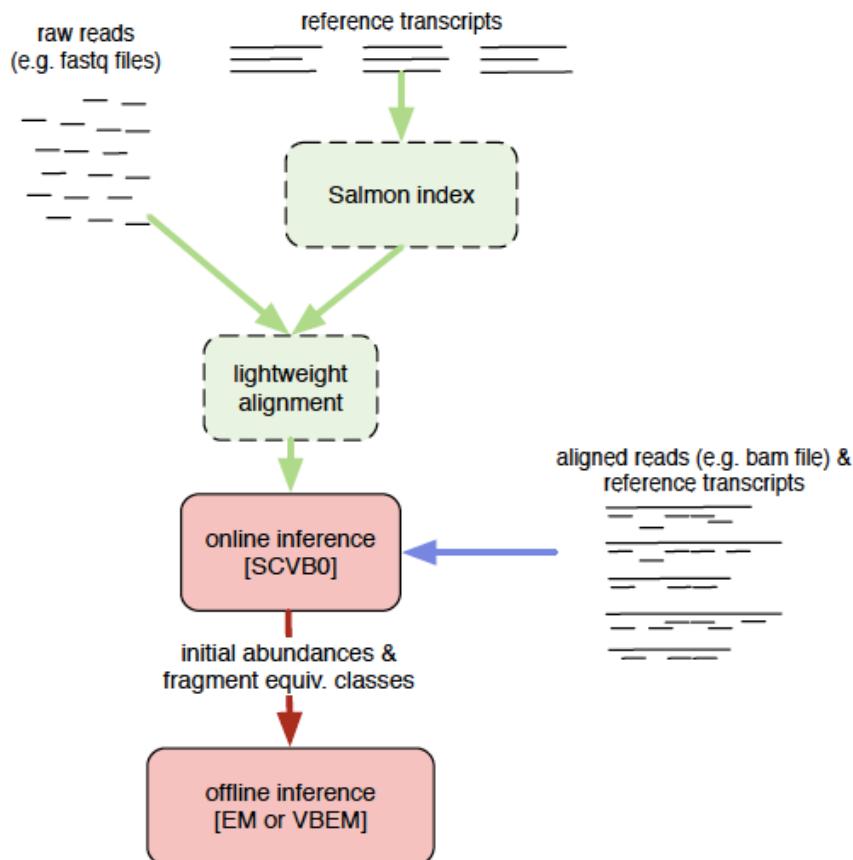


Figure 1: Overview of kallisto. (a) The input consists of a reference transcriptome and reads from an RNA-Seq experiment. (b) An index is constructed by creating the Transcriptome de Bruijn Graph (T-DBG) where nodes are  $k$ -mers, each transcript corresponds to a path and the path cover of the transcriptome induces a  $k$ -compatibility class for each  $k$ -mer. (c) Conceptually, the  $k$ -mers of a read are hashed (black nodes) to find the  $k$ -compatibility class of a read. (d) Skipping uses the information stored in the T-DBG to skip  $k$ -mers that are redundant due to having the same  $k$ -compatibility class. (e) The  $k$ -compatibility class of the read is determined by taking the intersection of the  $k$ -compatibility classes of its constituent  $k$ -mers.



University of  
Zurich<sup>UZH</sup>

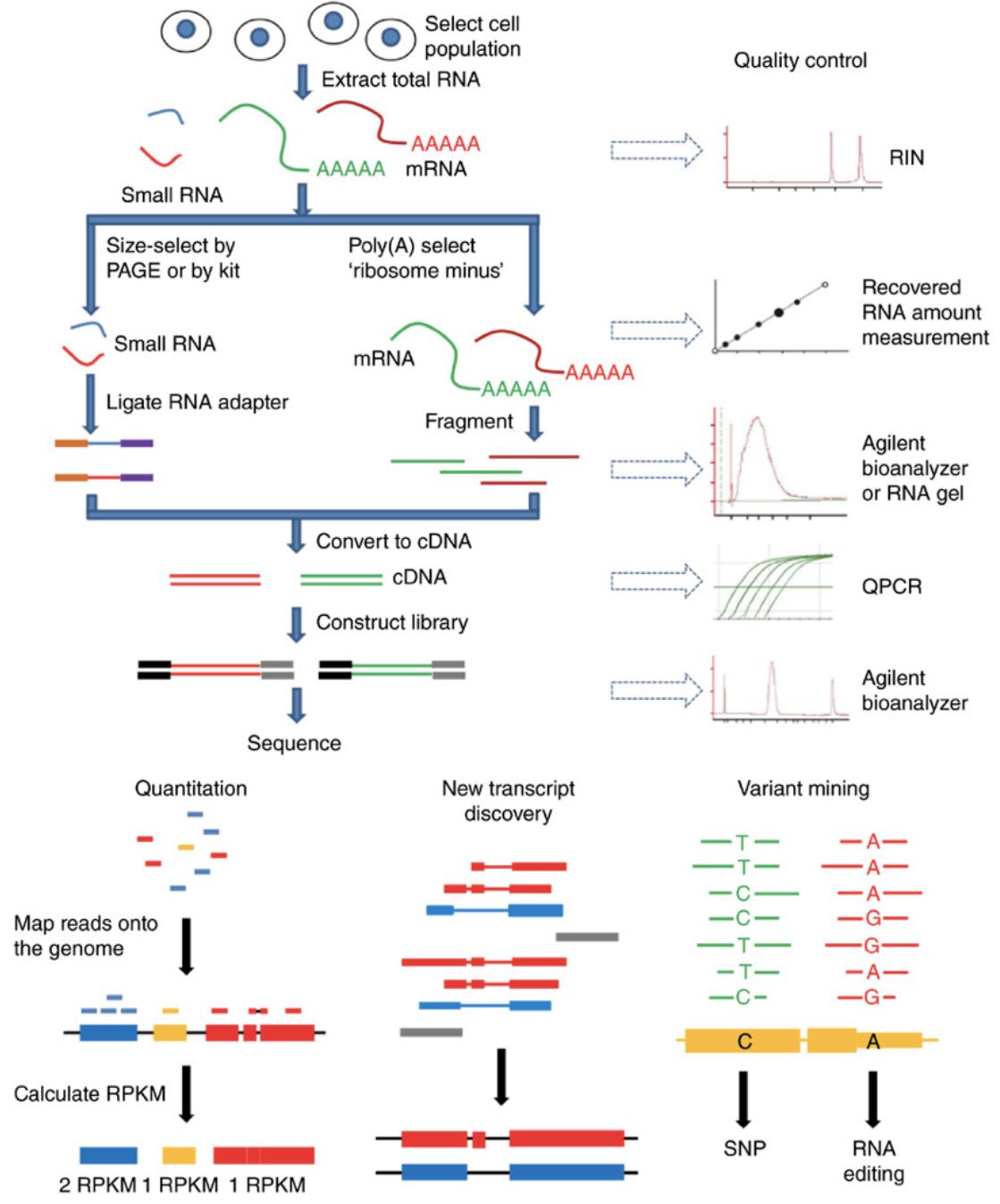
Institute of Molecular Life Sciences

---

# RESEARCH TALK ABOUT DIFFERENTIAL GENE/ TRANSCRIPT ANALYSES



1. Map the reads to reference sequences
2. “Count” reads that map to genes (quantify)
3. Compute DE Statistics





# Not everyone agrees on ..

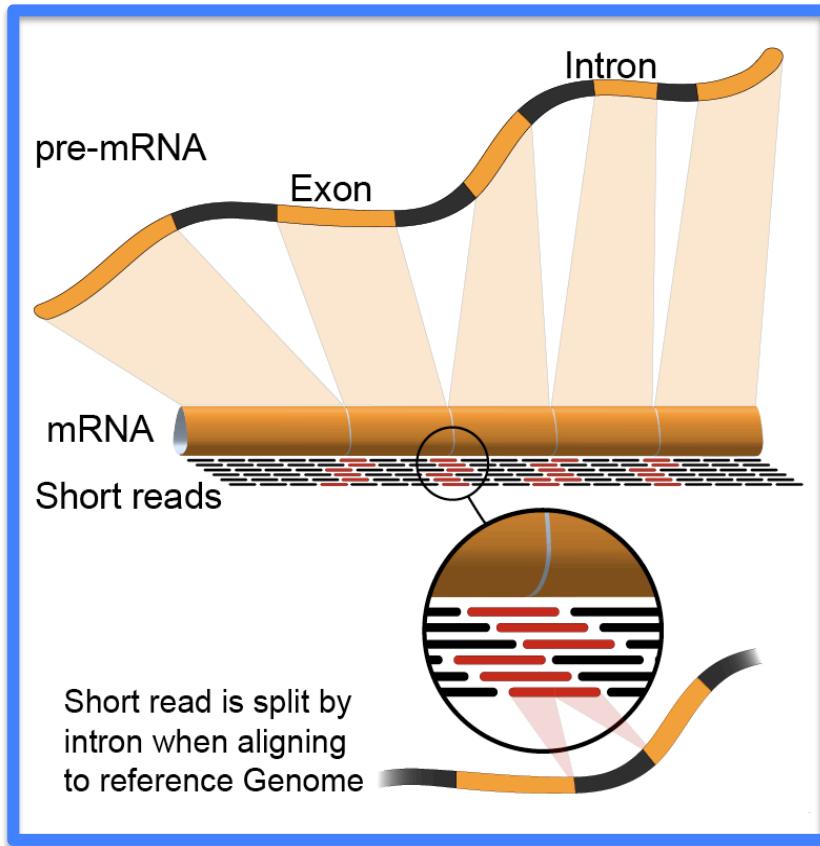
Mapping: full alignments are costly (compute time); “lightweight” alignments are just as accurate but much faster

Counting: 2 distinct schools of thought (exon-union counters, transcript counters)

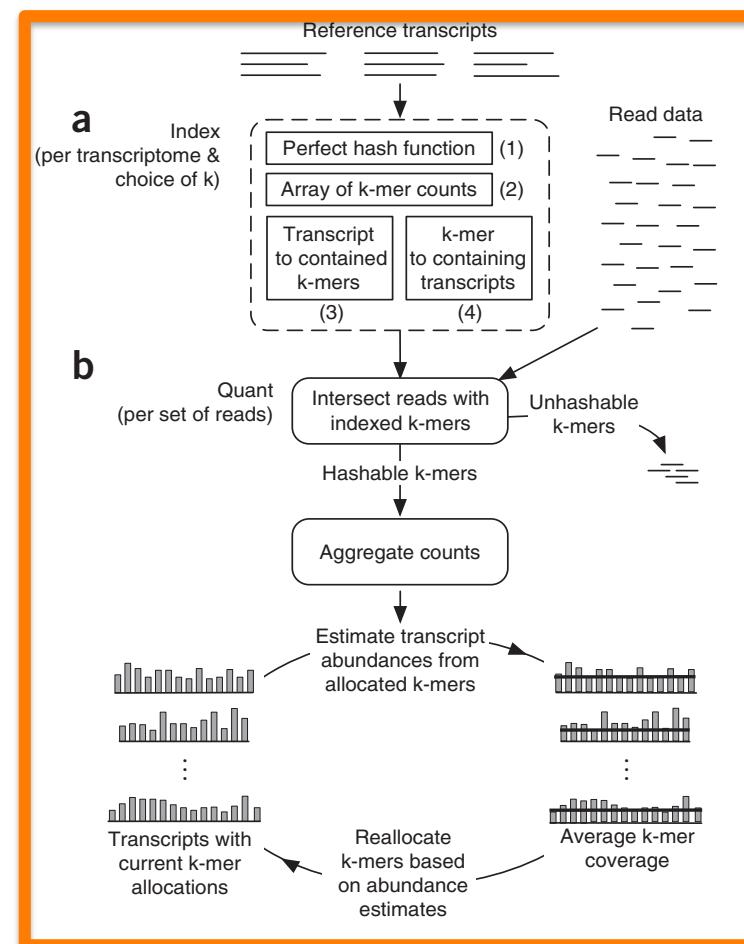
Statistics: data are counts, so models like negative binomial are prominent



# Alignment versus lightweight-alignment



<https://en.wikipedia.org/wiki/RNA-Seq>



sailfish (Patro et al. 2014)



## Counting/Quantification

union counters

→

simple sum of all reads

transcript counters

→

sum of length-normalized reads  
(often unknown which reads  
map to which transcript → portioning)

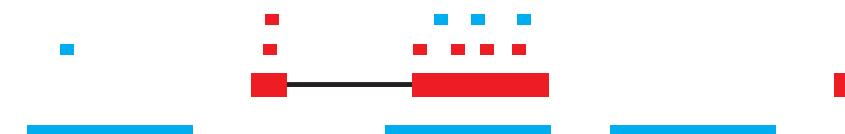
b

Condition A

Condition B

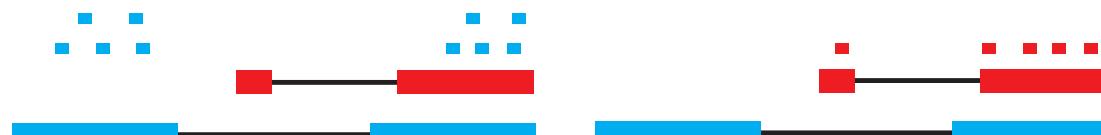
Log fold-change  
(union count)

Log fold-change  
(true expression)



$$\log_2\left(\frac{10}{10}\right) = 0$$

$$\log_2\left(\frac{\frac{10}{L}}{\frac{6}{L} + \frac{4}{2L}}\right) = 0.32$$



$$\log_2\left(\frac{5}{10}\right) = -1$$

$$\log_2\left(\frac{\frac{5}{L}}{\frac{10}{2L}}\right) = 0$$



# Define the differential problem (1)

Differential transcript  
expression (DTE)

Differential gene  
expression (DGE)

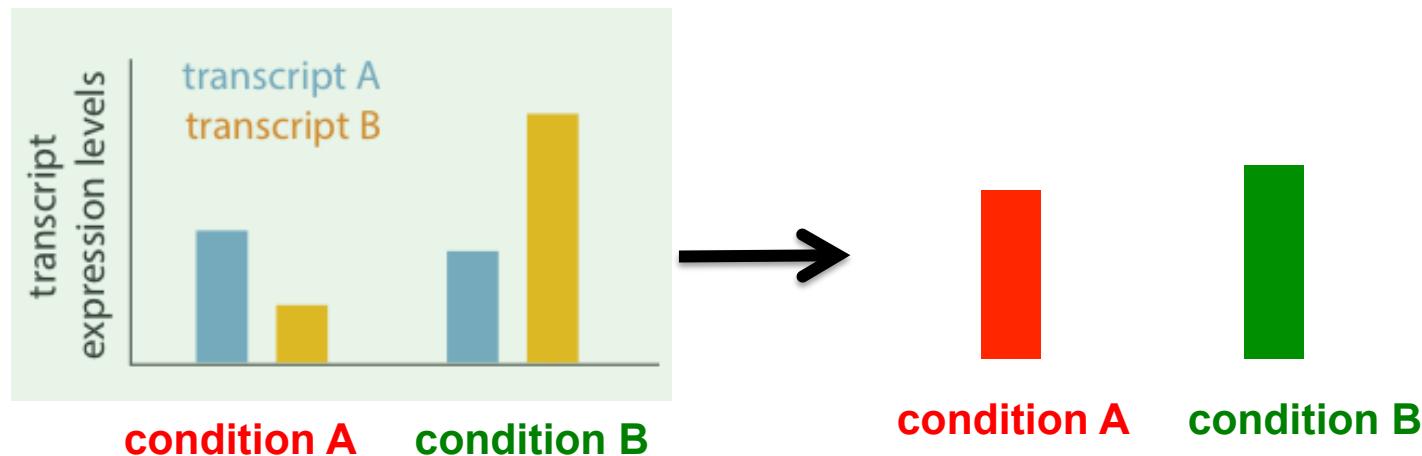




# Define the differential problem (2)

Differential transcript  
usage (DTU)

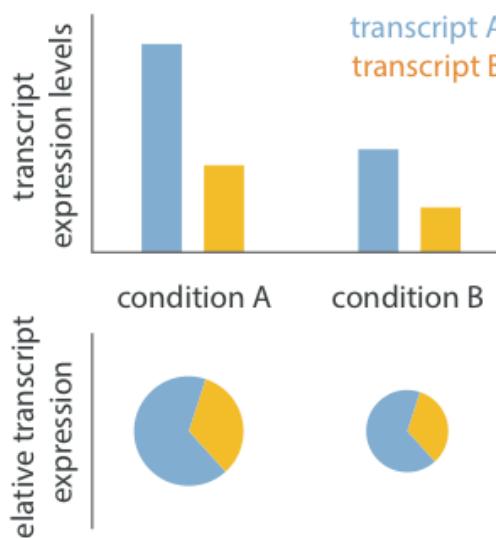
Differential gene  
expression (DGE)



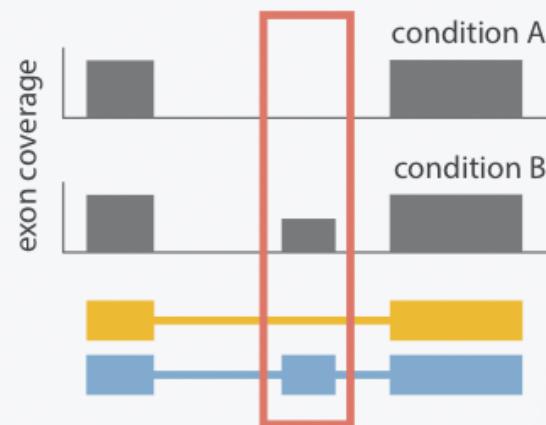


# Define the differential problem (3)

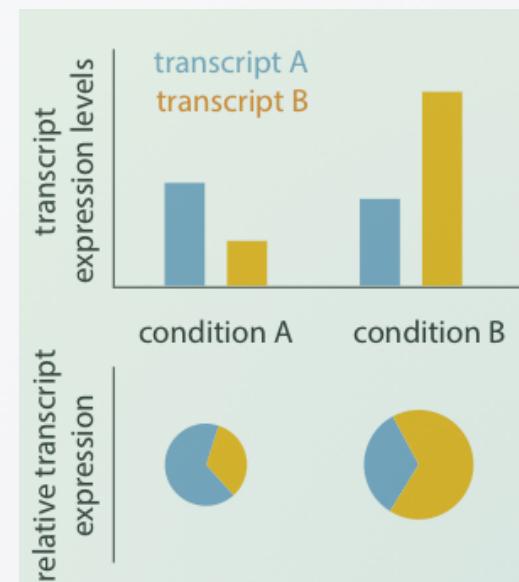
Differential transcript expression (DTE)



Differential exon usage (DEU)



Differential transcript usage (DTU)



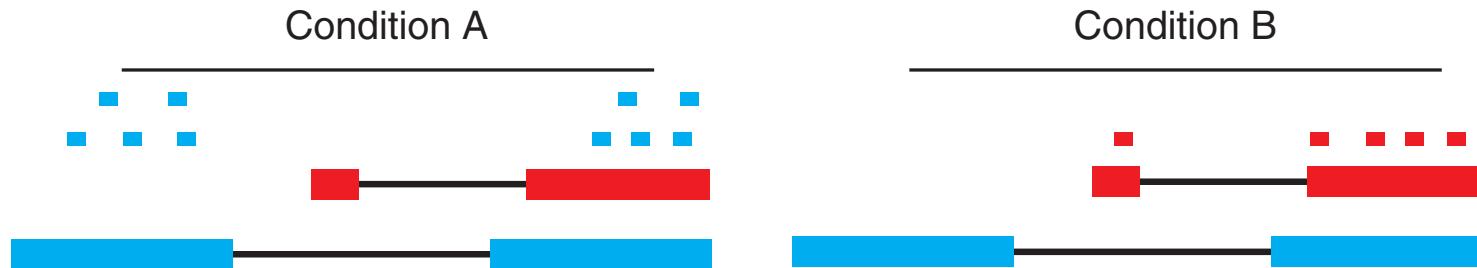
See also Soneson, Matthes et al., 2016,  
Genome Biology (comparison of DTU methods)



# What inference could/should be made?

Do you want to know:

- whether individual transcripts have changed? (DTE)
- whether *any* transcripts in gene have changed? (DTE → G)
- whether the overall output has changed? (DGE)
- whether transcript proportions have changed? (DTU/DEU)



Blue/red transcript changed?  
*Any transcripts changed?*  
Overall expression change?  
Transcript proportions changed?

Yes, Yes  
Yes  
No  
Yes



## Many subtleties here ..

- union counting is easy, even when the transcript catalog is incomplete; transcript counting in complicated genes is difficult
- more transcripts than genes; more statistical tests, higher multiple testing penalty, perhaps lower sensitivity
- do statistical methods (e.g., NB) work well for “raw” and “estimated” counts?
- If a transcript changes, first question often is “did all transcripts change” ?
- Much is hypothetical. Does it matter in practice?

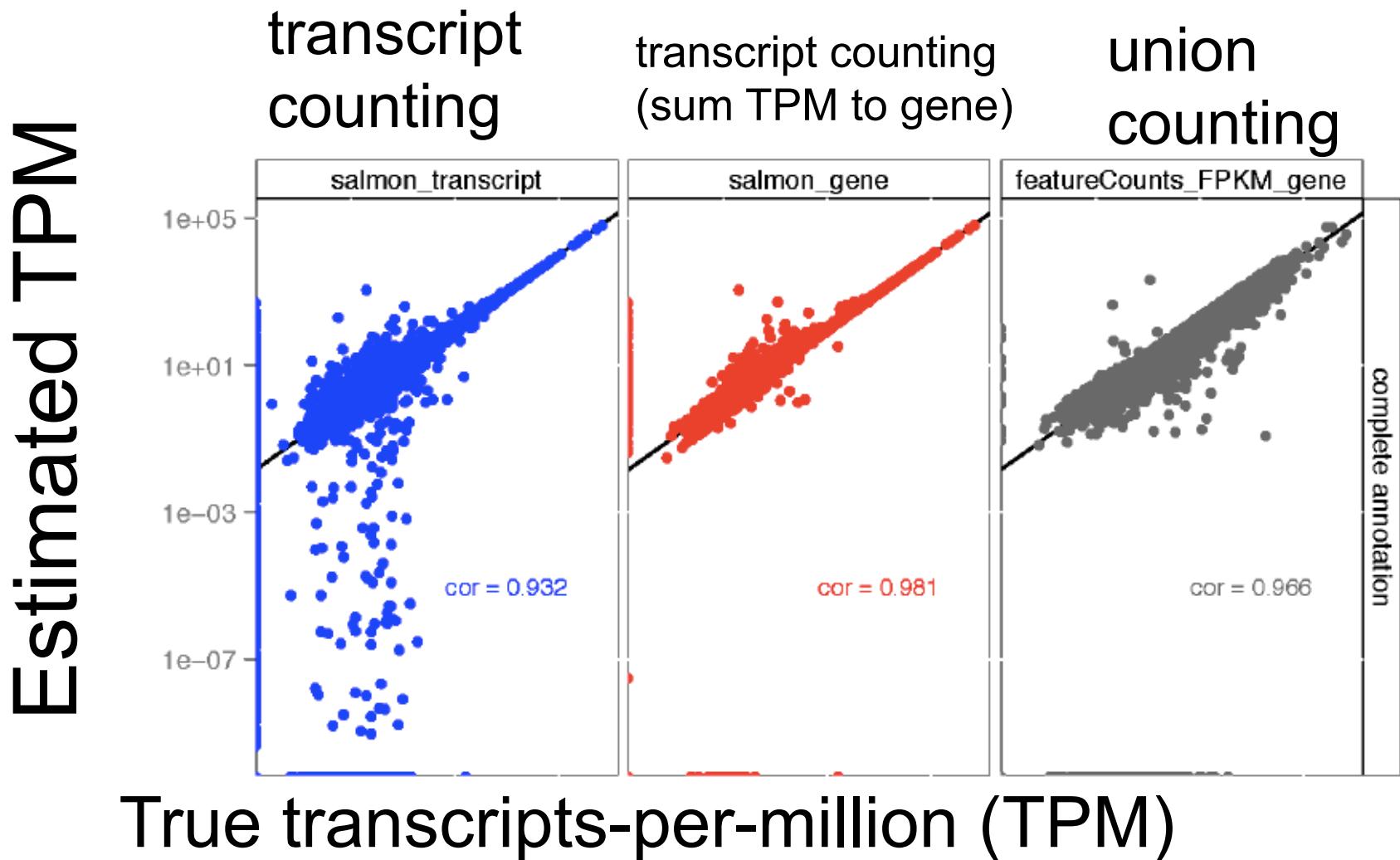


## We make/made some claims

1. expression estimates are considerably more stable at the gene-level
2. statistical engines work well for either gene-level (raw) or transcript-level (estimated) counts
3. more powerful and easier to interpret at the gene-level
4. the difference between union counting and transcript counting is, on average, small (depends on dataset)

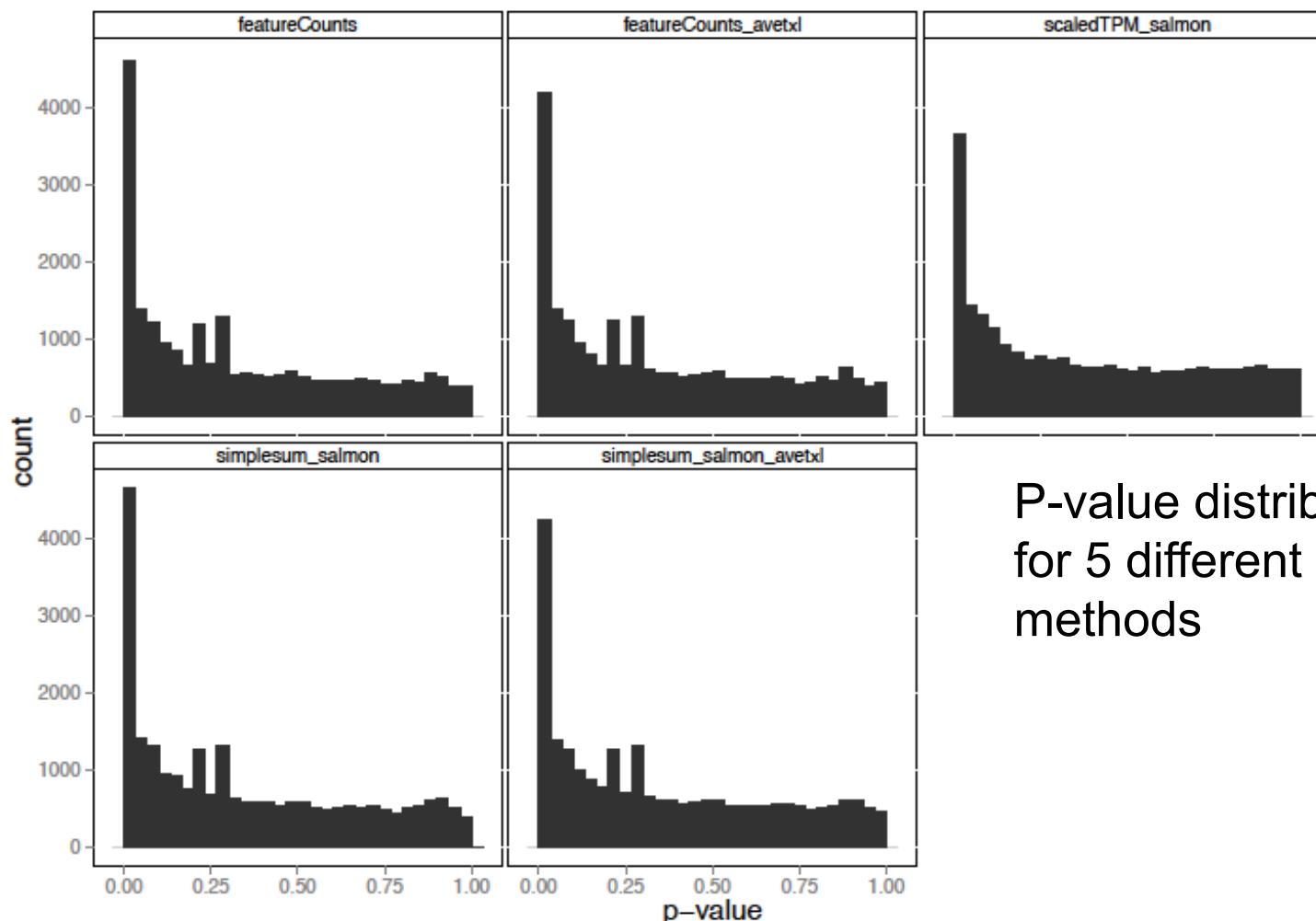


Claim 1: gene-level estimation is easier (simulated RNA-seq data)



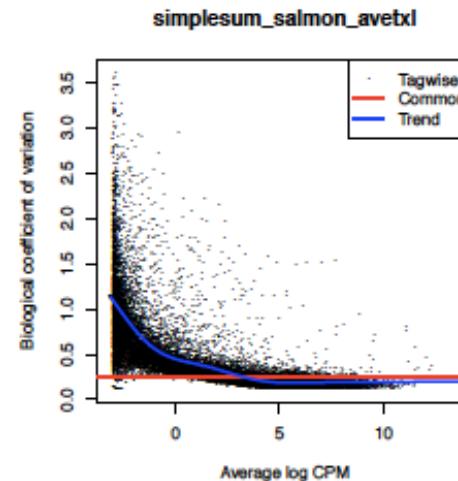
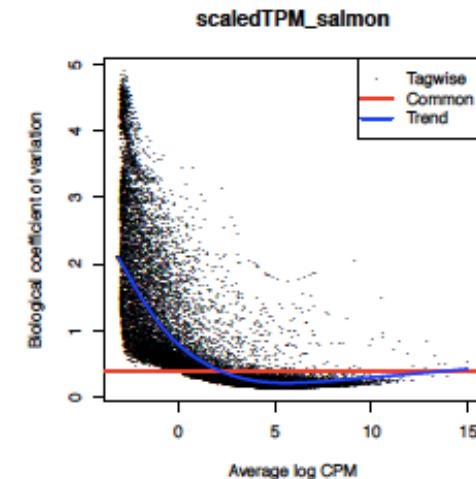
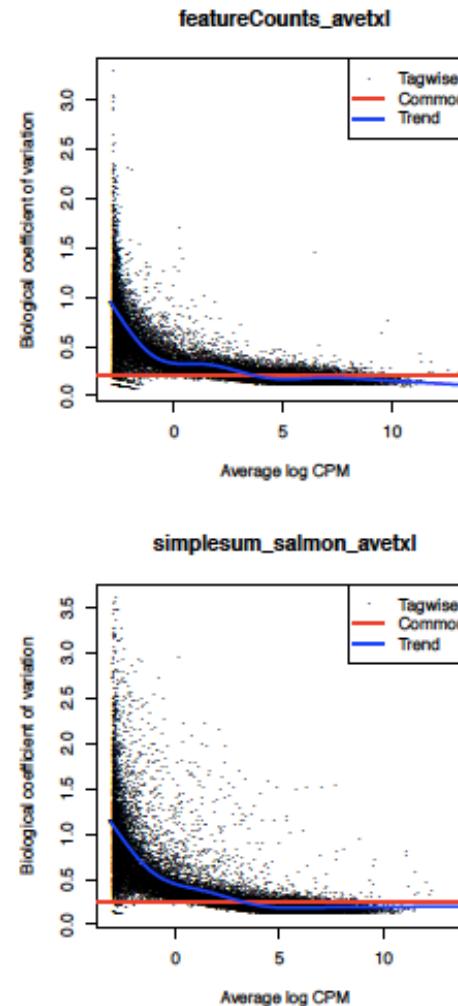
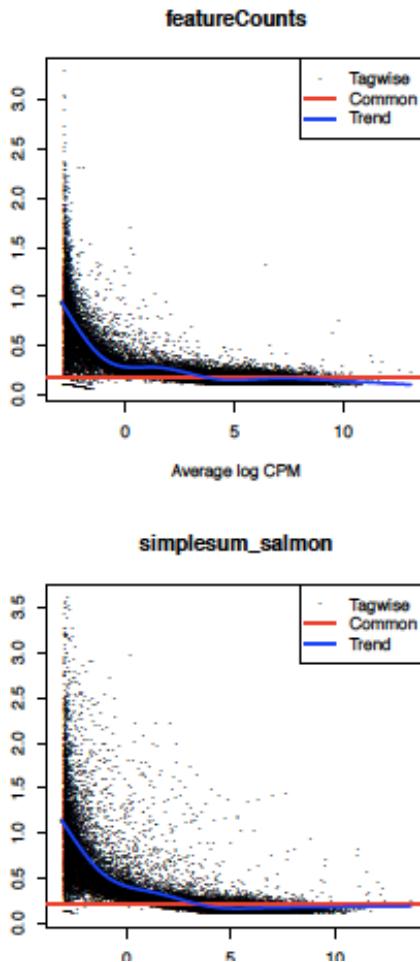


## Claim 2a: statistical methods are equally “healthy”





dispersion



mean



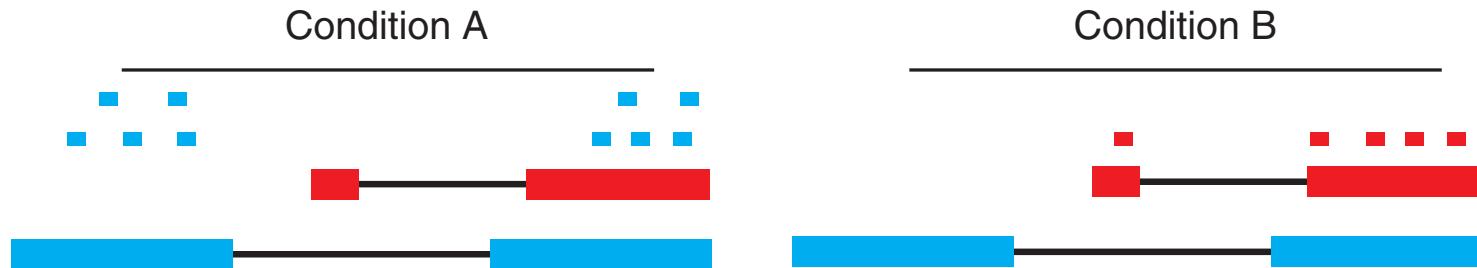
5 different counting  
methods



# What inference could/should be made?

Do you want to know:

- whether individual transcripts have changed? (DTE)
- whether *any* transcripts in gene have changed? (DTE → G)
- whether the overall output has changed? (DGE)
- whether transcript proportions have changed? (DTU/DEU)

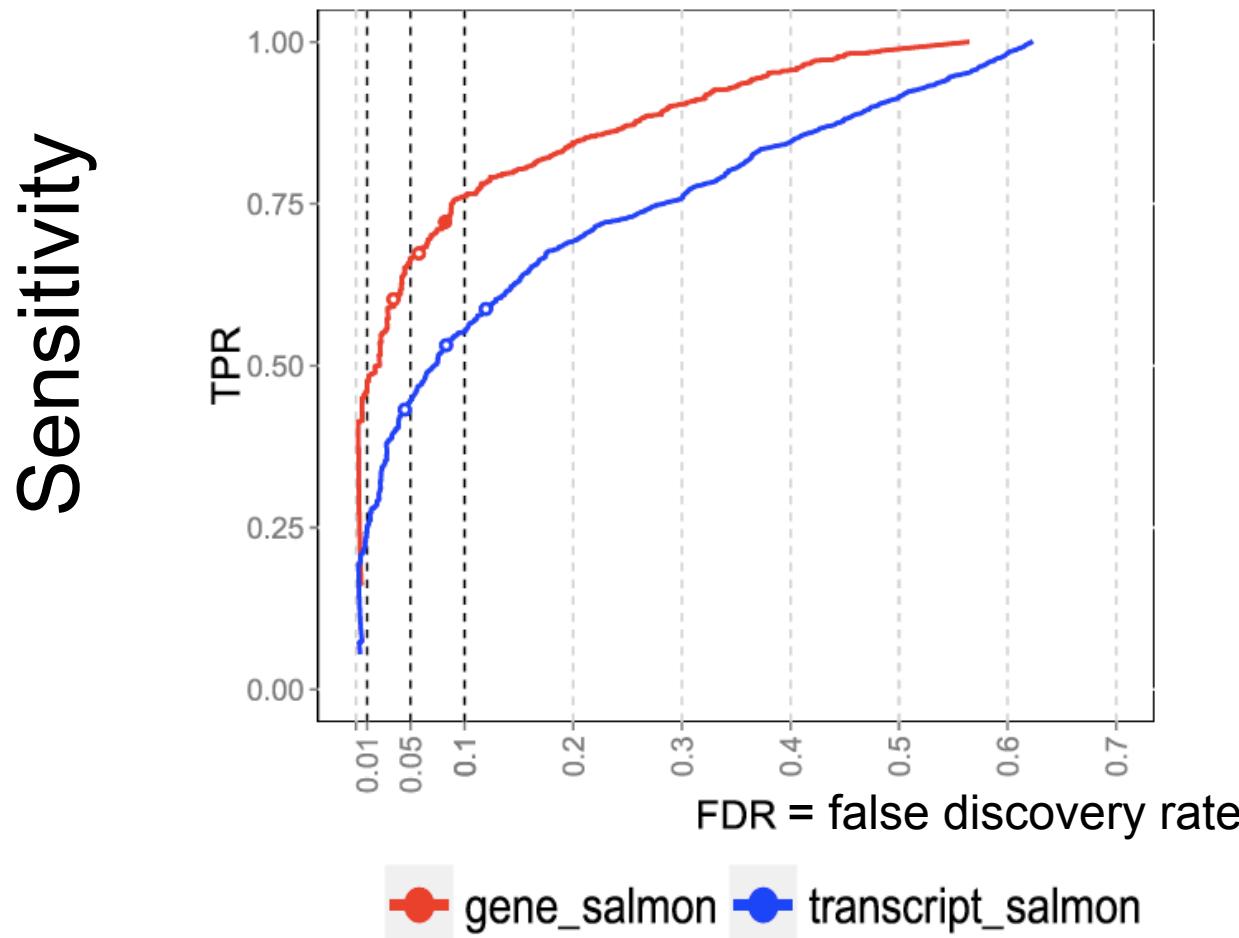


Blue/red transcript changed?  
*Any transcripts changed?*  
Overall expression change?  
Transcript proportions changed?

Yes, Yes  
Yes  
No  
Yes



## Claim 3a: DE analyses more powerful at gene-level (DTE → G)



N.B.: answering  
different questions

DTE – do any  
transcripts change?  
(harder)

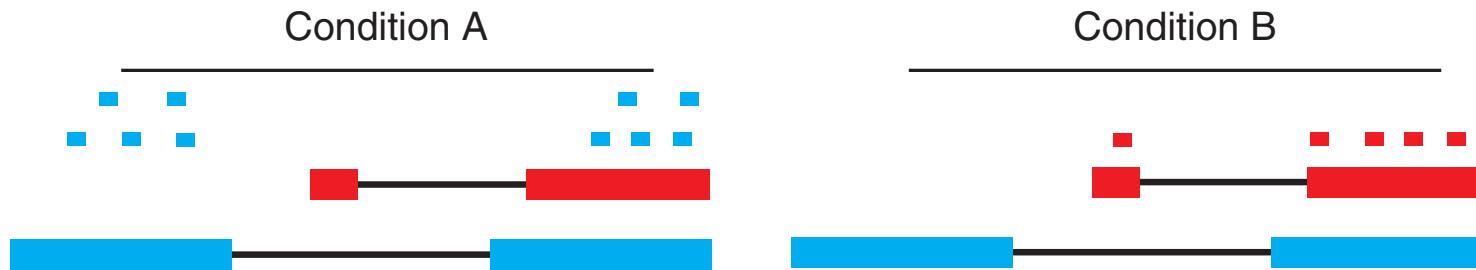
DTE → G – any  
changes within  
transcripts of this  
gene? (easier)



( Claim 3a:  $DTE \rightarrow G > DTE$  )

Claim 3b:  $DTU+DGE > DTE$

- individual transcripts have changed ? (DTE)
- the overall output has changed ? (DGE)
- transcript proportions have changed ? (DTU)



Blue/red transcript changed?  
Any transcripts changed?  
Overall expression change?  
Transcript proportions changed?

Yes, Yes  
Yes  
No  
Yes



# Some situations strictly call for DTE

## BACKGROUND

The androgen-receptor isoform encoded by splice variant 7 lacks the ligand-binding domain, which is the target of enzalutamide and abiraterone, but remains constitutively active as a transcription factor. We hypothesized that detection of androgen-receptor splice variant 7 messenger RNA (AR-V7) in circulating tumor cells from men with advanced prostate cancer would be associated with resistance to enzalutamide and abiraterone.

## METHODS

We used a quantitative reverse-transcriptase–polymerase-chain-reaction assay to evaluate AR-V7 in circulating tumor cells from prospectively enrolled patients with metastatic castration-resistant prostate cancer who were initiating treatment with either enzalutamide or abiraterone. We examined associations between AR-V7 status (positive vs. negative) and prostate-specific antigen (PSA) response rates (the primary end point), freedom from PSA progression (PSA progression–free survival), clinical or radiographic progression–free survival, and overall survival.

## RESULTS

A total of 31 enzalutamide-treated patients and 31 abiraterone-treated patients were enrolled, of whom 39% and 19%, respectively, had detectable AR-V7 in circulating tumor cells. Among men receiving enzalutamide, AR-V7–positive patients had lower PSA response rates than AR-V7–negative patients (0% vs. 53%,  $P=0.004$ ) and shorter PSA progression–free survival (median, 1.4 months vs. 6.0 months;  $P<0.001$ ), clinical or radiographic progression–free survival (median, 2.1 months vs. 6.1 months;  $P<0.001$ ), and overall survival (median, 5.5 months vs. not reached;  $P=0.002$ ). Similarly, among men receiving abiraterone, AR-V7–positive patients had lower PSA response rates than AR-V7–negative patients (0% vs. 68%,  $P=0.004$ ) and shorter PSA progression–free survival (median, 1.3 months vs. not reached;  $P<0.001$ ), clinical or radiographic progression–free survival (median, 2.3 months vs. not reached;  $P<0.001$ ), and overall survival (median, 10.6 months vs. not reached,  $P=0.006$ ). The association between AR-V7 detection and therapeutic resistance was maintained after adjustment for expression of full-length androgen receptor messenger RNA.

## CONCLUSIONS

Detection of AR-V7 in circulating tumor cells from patients with castration-resistant prostate cancer may be associated with resistance to enzalutamide and abiraterone. These findings require large-scale prospective validation. (Funded by the Prostate Cancer Foundation and others.)

## ORIGINAL ARTICLE

### AR-V7 and Resistance to Enzalutamide and Abiraterone in Prostate Cancer

Emmanuel S. Antonarakis, M.D., Changxue Lu, Ph.D., Hao Wang, Ph.D., Brandon Luber, Sc.M., Mary Nakazawa, M.H.S., Jeffrey C. Roeser, B.S., Yan Chen, Ph.D., Tabrez A. Mohammad, Ph.D., Yidong Chen, Ph.D., Helen L. Fedor, B.S., Tamara L. Lotan, M.D., Qizhi Zheng, M.D., Angelo M. De Marzo, M.D., Ph.D., John T. Isaacs, Ph.D., William B. Isaacs, Ph.D., Rosa Nadal, M.D., Channing J. Paller, M.D., Samuel R. Denmeade, M.D., Michael A. Carducci, M.D., Mario A. Eisenberger, M.D., and Jun Luo, Ph.D.

Cited 484 times since Sept 2014



# Simulation

synthetic data set: 3,858 genes and 15,677 transcripts (human chr1;  
ArrayExpress: E-MTAB-4119)

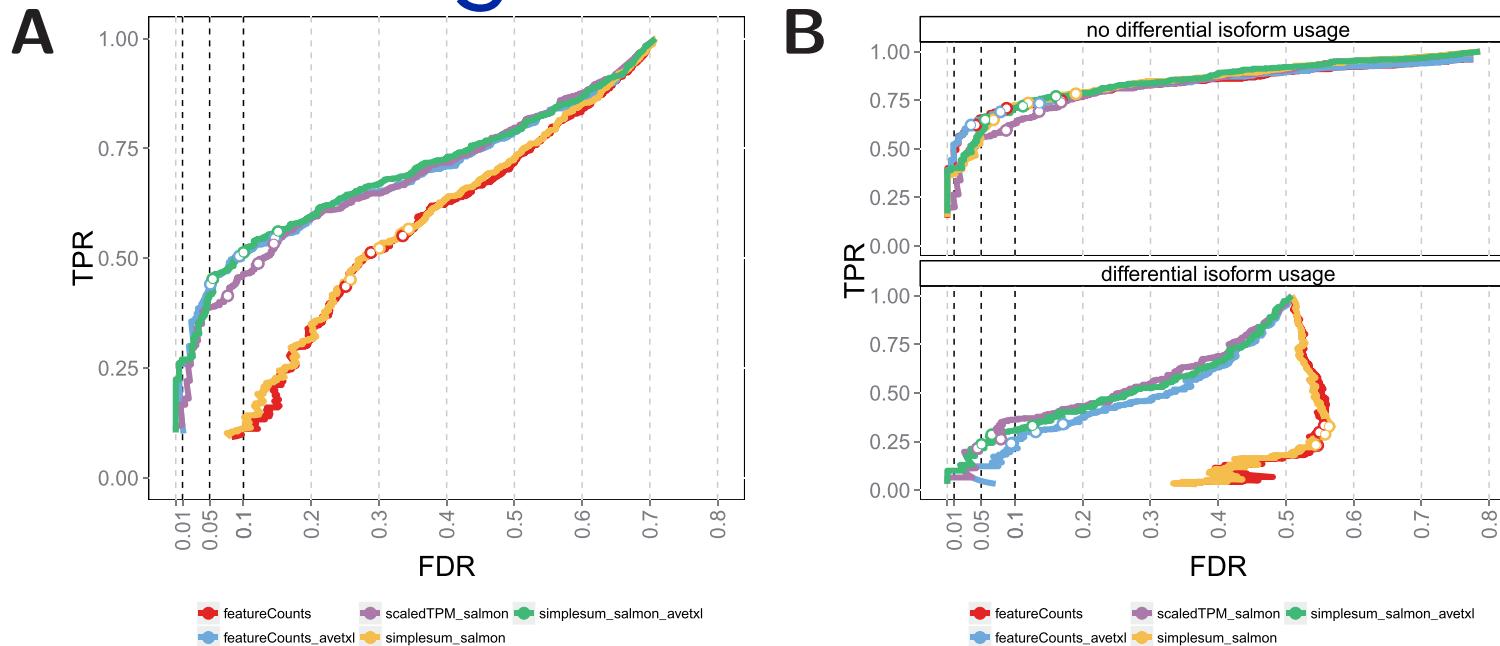
two conditions with three biological replicates (empirical distribution  
from real dataset)

True changes come in three disjoint sets:

- overall DGE → all transcripts of gene w/ same fold change (420 genes)
- DGE=0; DTU → total output constant but transcript proportions changed (420 genes)
- DTE → 10% of the transcripts of each affected gene modified (422 genes, 528 transcripts)



# Claim 4: DTU can derail inference based on union counting



**Figure 3 (sim2). A:** DGE detection performance of *edgeR* applied to three different count matrices (simplesum, scaledTPM, featureCounts), with or without including an offset representing the average transcript length (for simplesum and featureCounts). Including the offset or using the scaledTPM count matrix leads to improved FDR control compared to using simplesum or featureCounts matrices without offset. The curves trace out the observed FDR and TPR for each significance cutoff value. The three circles mark the performance at adjusted p-value cutoffs of 0.01, 0.05 and 0.1. **B:** stratification of the results in **A** by the presence of differential isoform usage. The improvement in FDR control seen in **A** results from an improved treatment of genes with differential isoform usage, while all methods perform similarly for genes without differential isoform usage.

If lots of DTU → union counters underperform



University of  
Zurich<sup>UZH</sup>

Institute of Molecular Life Sciences

---

# THEORY



## Differential expression, small sample inference

- Table of data (e.g., microarray gene expression data OR RNA-seq gene expression counts with replicates of each of condition A, condition B)
  - *rows* = features (e.g., genes), *columns* = experimental units (samples)
- Most common problem in statistical bioinformatics: want to infer whether there is a **change in the response** → a statistical test for each row of the table.

What test might you use? Why is this hard? What issues arise? How much statistical power is there [1] ?

```
> head(y)
      group0      group0      group0      group1      group1      group1
gene1 -0.1874854  0.2584037 -0.05550717 -0.4617966 -0.3563024 -0.03271432
gene2 -3.5418798 -2.4540999  0.11750996 -4.3270442 -5.3462622 -5.54049106
gene3 -0.1226303  0.9354707 -1.10537767 -0.1037990  0.5221678 -1.72360854
gene4 -2.3394536 -0.3495697 -3.47742610 -3.2287093  6.1376670 -2.23871974
gene5 -3.7978820  1.4545702 -7.14796503 -4.0500796  4.7235714 10.00033769
gene6  1.4627078 -0.3096070 -0.26230124 -0.7903434  0.8398769 -0.96822312
```

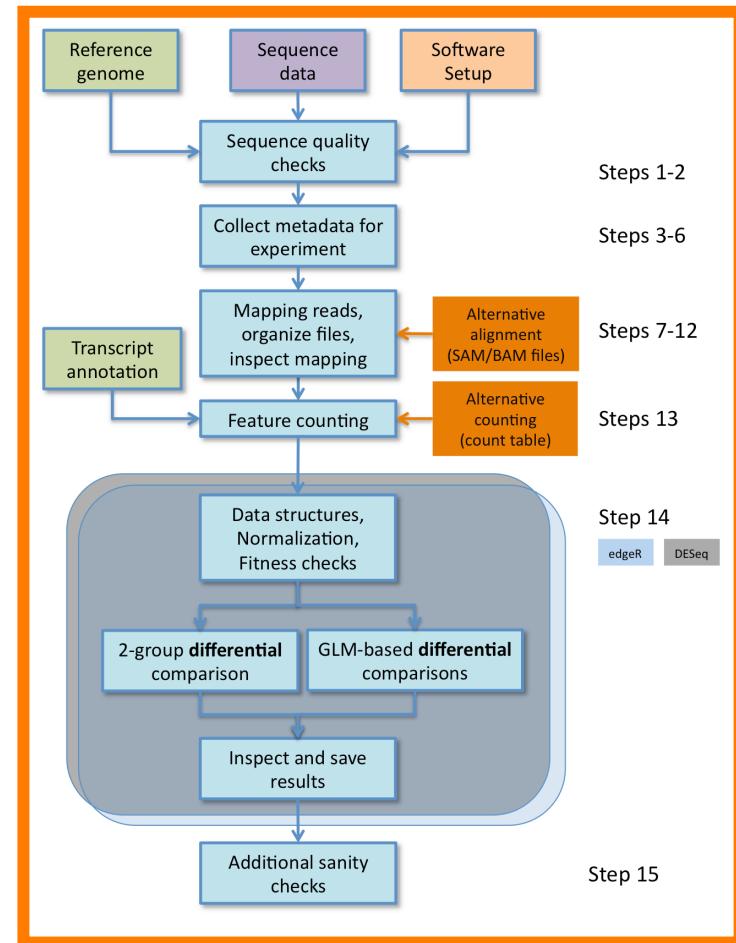
[1] <http://www.stat.ubc.ca/~rollin/stats/ssize/n2.html>



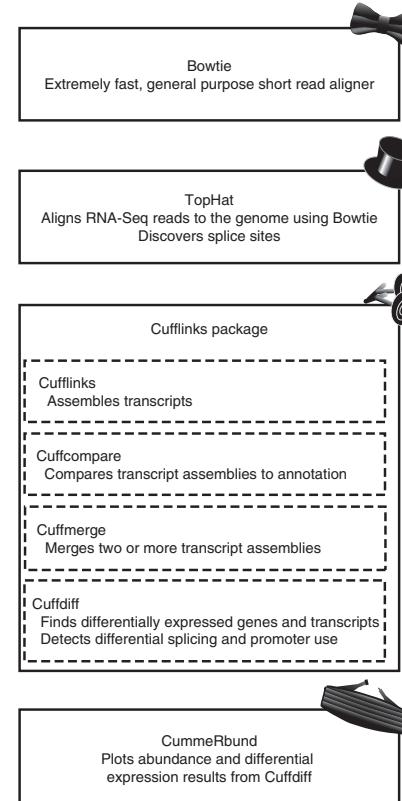
# Data analysis pipelines for RNA-seq differential expression

Institute of Molecular Life Sciences

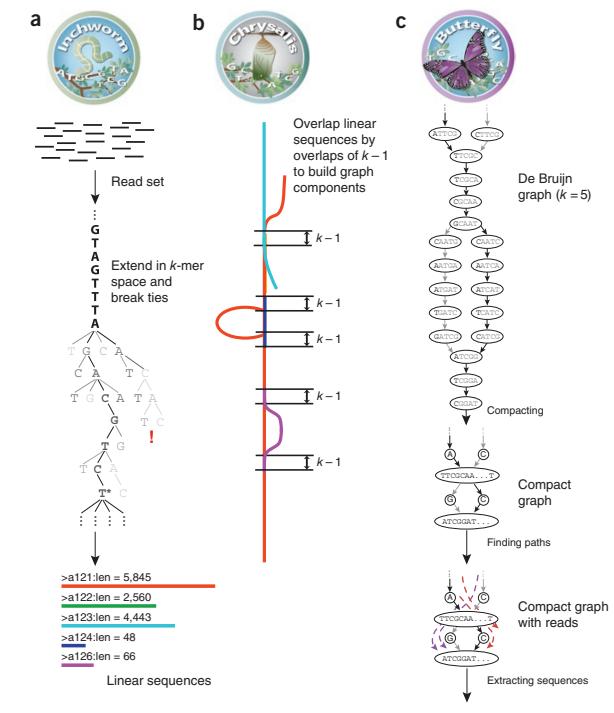
## edgeR, DESeq



## 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).

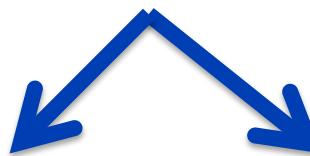


## Differential expression: why not use methods developed for microarrays?

Count data is discrete, not continuous.

Methods designed for microarrays are not directly applicable and suboptimal (**more on this later**)

Two options:

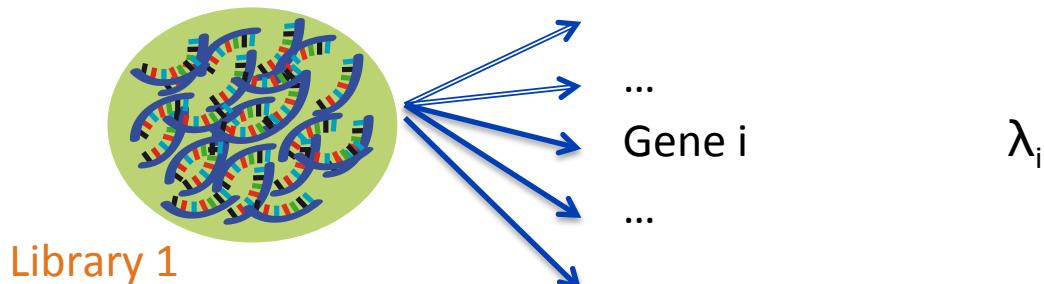


Transform count data  
and apply standard  
methodology

Analyze using  
models for count  
data



## 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  $i$

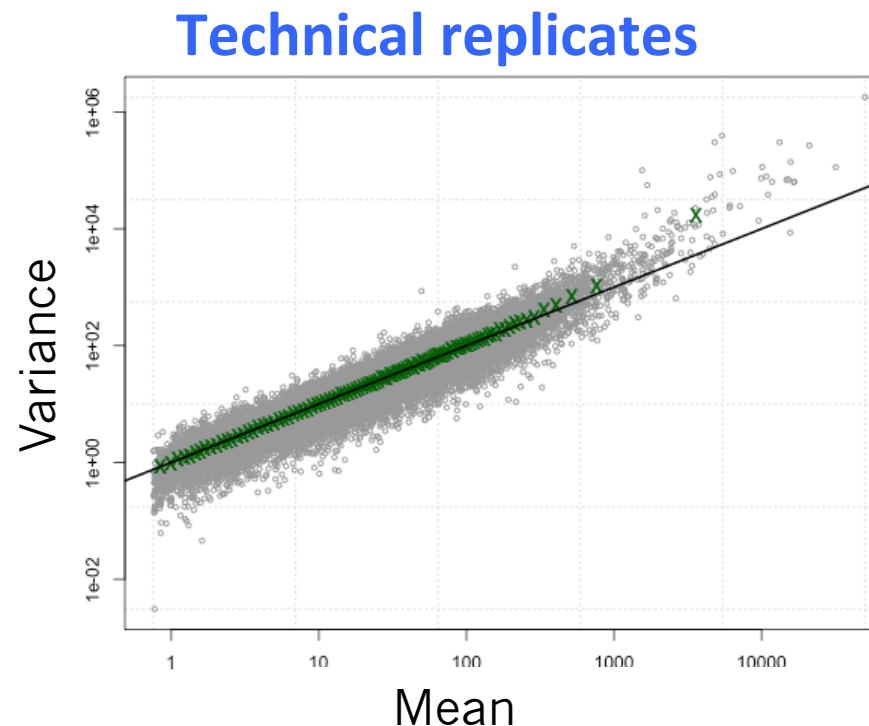
$M$  - total number of sequences

$\lambda_i$  - proportion

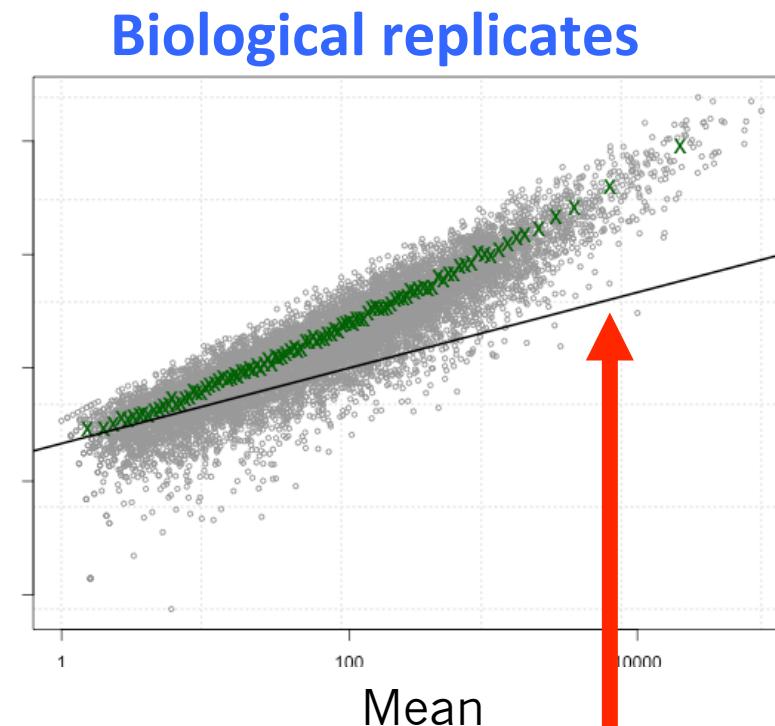
Large  $M$ , small  $\lambda_i \rightarrow$  approximated well by Poisson(  $\mu_i = M \cdot \lambda_i$  )



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



Data from Marioni et al. *Genome Research* 2008



Data from Parikh et al.  
*Genome Biology* 2010

mean=variance  
(Poisson assumption)



## Count data modeling assumptions

Poisson adequately describes technical variation

$$Y_i \sim \text{Pois}(M * \lambda_i)$$

$$\text{mean}(Y_i) = \text{variance}(Y_i) = M * \lambda_i$$

Negative binomial (gamma-Poisson) model is a natural extension that allows biological variability:

$$Y_i \sim \text{NB}(\mu_i = M * \lambda_i, \phi_i)$$

Same mean, variance is quadratic in the mean:

$$\text{variance}(Y_i) = \mu_i (1 + \mu_i \phi_i)$$

$M$  = library size

$\lambda_i$  = relative contribution of gene  $i$



## Analogy to t-tests (microarray setting)

$$t_g = \frac{\bar{y}_{\text{mu}} - \bar{y}_{\text{wt}}}{s_g c}$$

Feature-specific

$$\tilde{t}_g = \frac{\bar{y}_{\text{mu}} - \bar{y}_{\text{wt}}}{\tilde{s}_g u}$$

Moderated

$$t_{g,\text{pooled}} = \frac{\bar{y}_{\text{mu}} - \bar{y}_{\text{wt}}}{s_0 c}$$

Common

Why did moderated-t work well for microarrays ?



## Moderated estimates: let's try the same strategy with counts

At one extreme, assume all genes have same dispersion (too strong)

At other extreme, estimate dispersion separately/independently for each gene (poor estimates)

Shrink individual estimates toward common/trend (how?)

No hierarchical model (e.g. limma) to do this:  
**approximations, weighted likelihood**

No t-distribution theory to formulate statistical tests.



## Second challenge: Moderate dispersion estimate

Weighted likelihood -- individual log-likelihood plus a weighted version of the **common** log-likelihood:

$$WL(\phi_g) = l_g(\phi_g) + \alpha l_C(\phi_g)$$

↑  
 $(1-\alpha)$

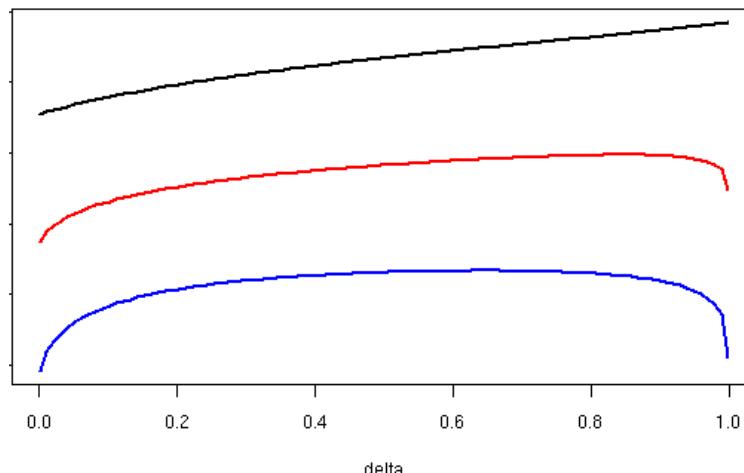
$L_g$  - quantile-adjusted conditional likelihood

**Black:** single tag

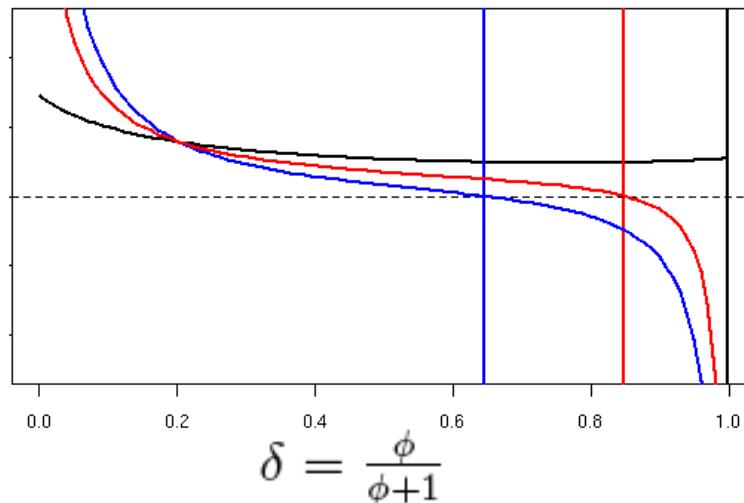
**Blue:** common dispersion

**Red:** Linear combination of the two

Log-Likelihood



Score (1<sup>st</sup> derivative of LL)

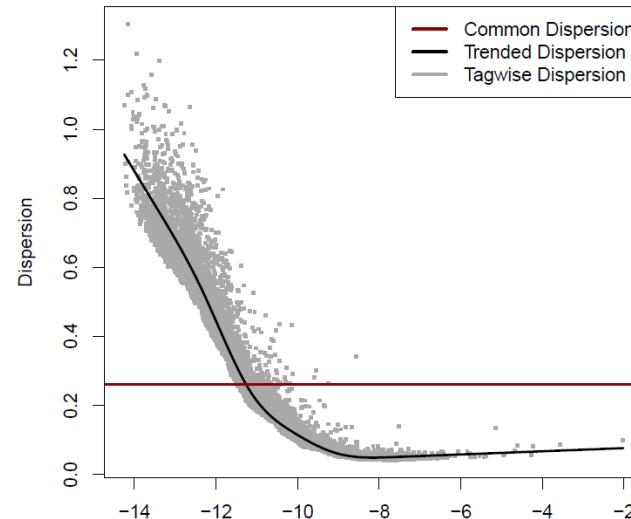
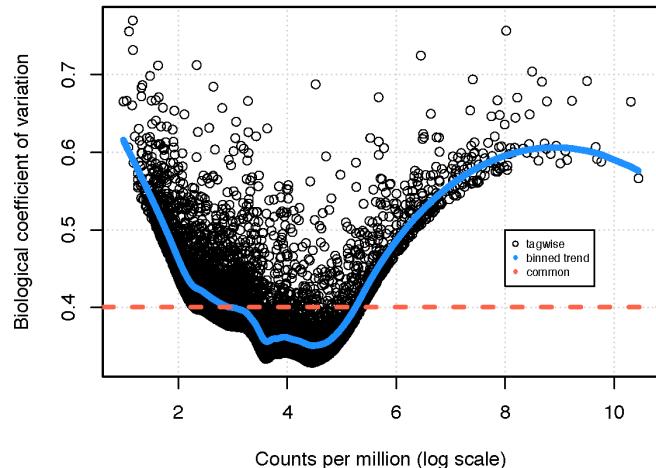




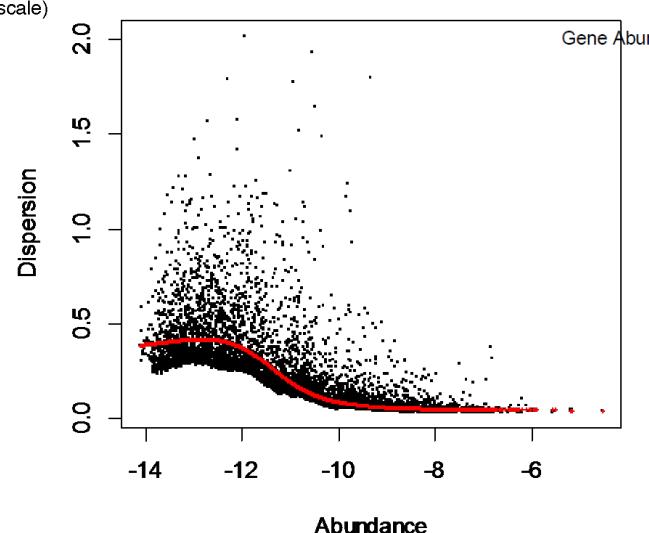
# Dispersion varies with mean: moderate dispersion towards trend

Data:

Tuch et al.,  
2008



Mouse  
hemopoietic  
stem cells



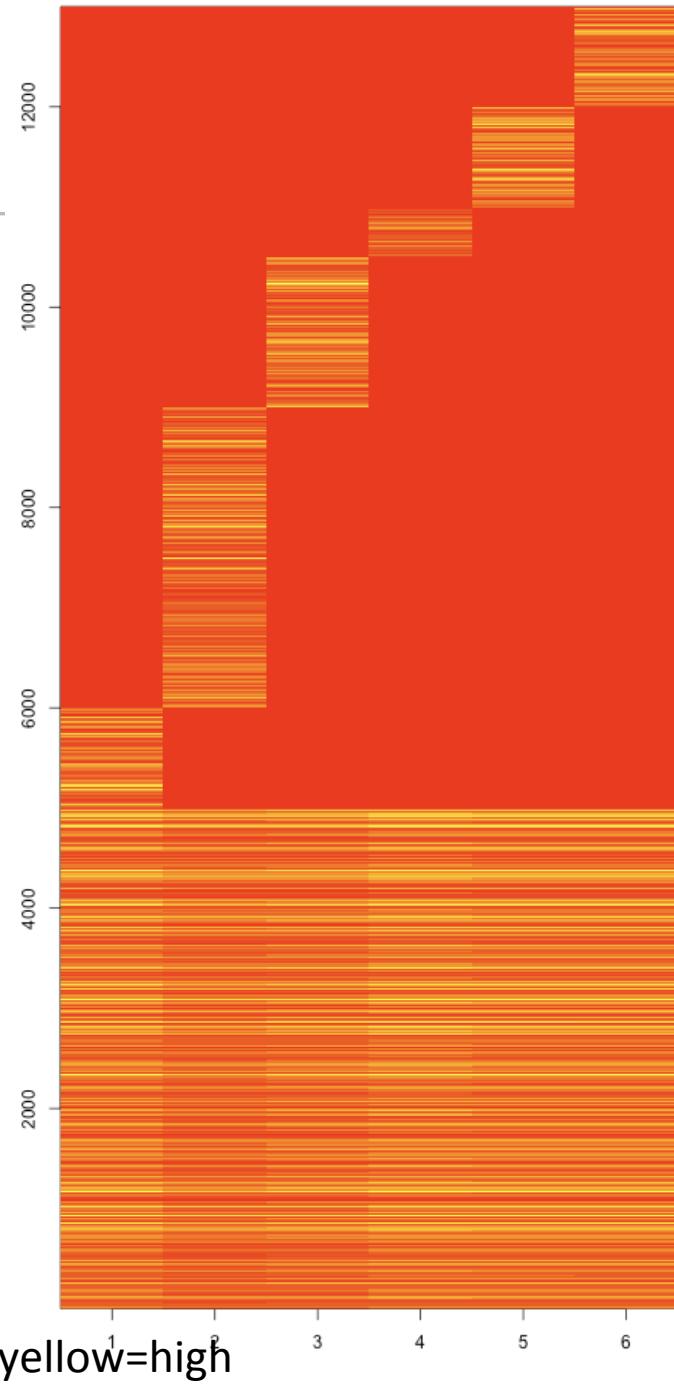
Mouse  
lymphomas

Advantage: genes are  
allowed to have their  
own variance.



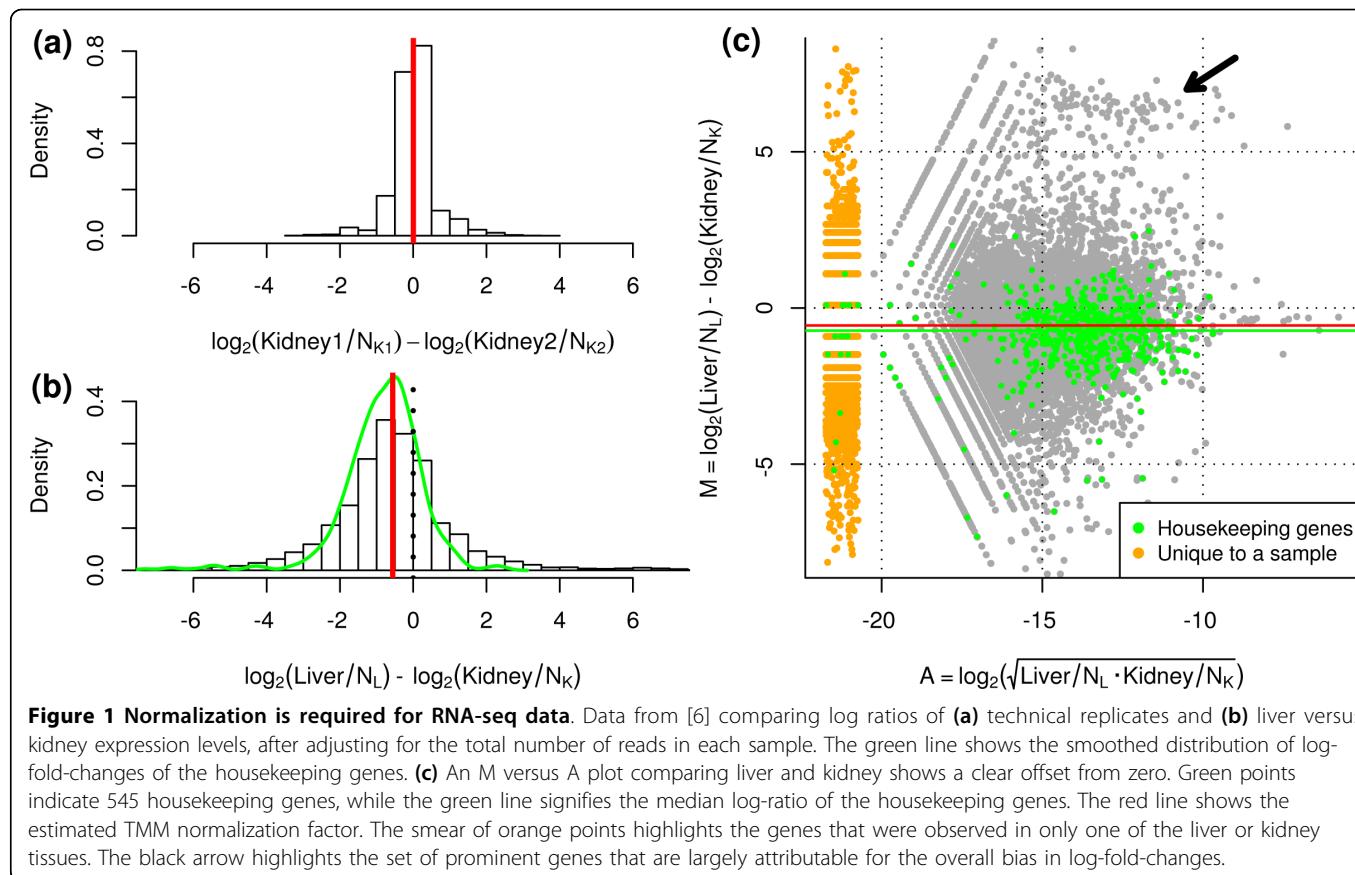
## “Composition” or “Diversity” can affect read depth

- 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





# Kidney and Liver RNA have very different composition





## What does transformation do to M-V relationship?

For Poisson data, square-root should stabilize

Logarithm is too strong – variance decreases to asymptote (dispersion; Neg Bin) or 0 (Poisson)

How to pick? Doesn't matter ... voom

voom: mean-variance modeling at the observational level

voom

package:limma

R Documentation

Transform RNA-Seq Data Ready for Linear Modelling

Description:

Transform count data to log2-counts per million, estimate the mean-variance relationship and use this to compute appropriate observational-level weights. The data are then ready for linear modeling.



## Model log counts per million

log counts per million:

$$z_{gi} = \log_2 \left( 1e6 \frac{\text{count}_{gi} + 0.5}{\text{libsize}_{gi} + 1.0} \right) = \log_2 \left( 1e6 \frac{y_{gi} + 0.5}{M_{gi} + 1.0} \right)$$

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

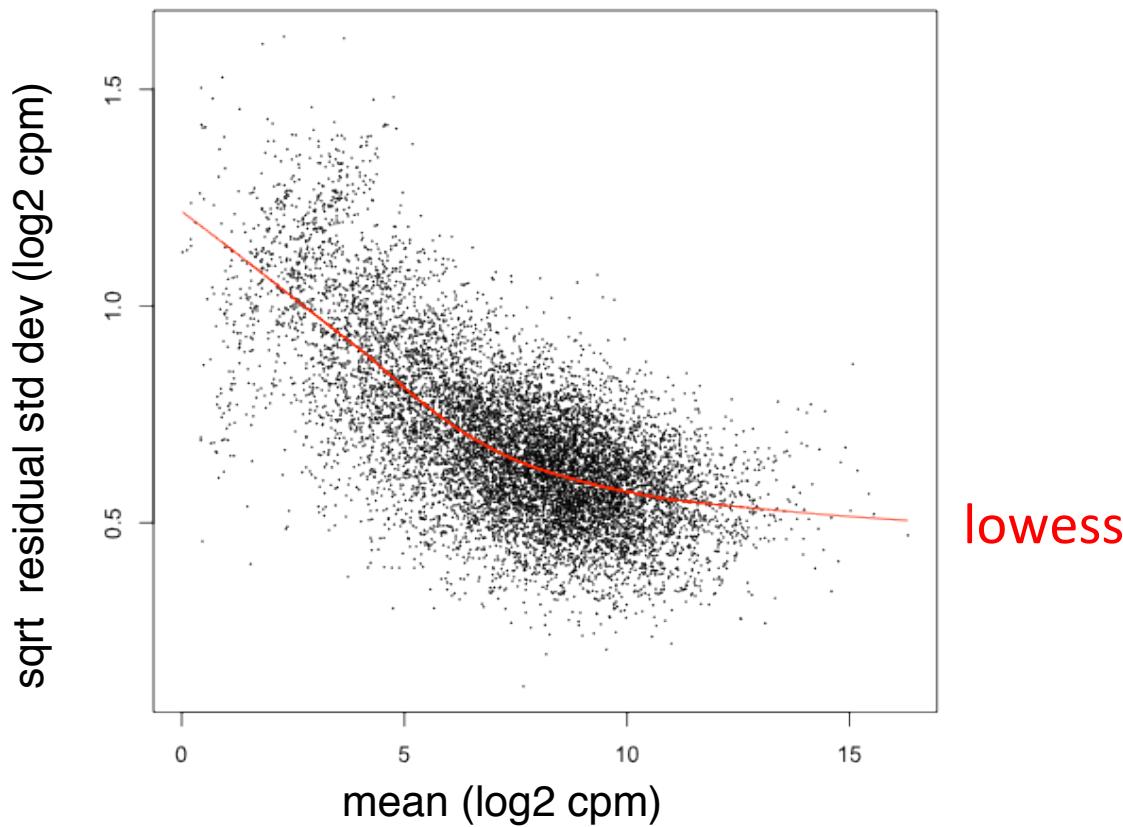
Linear modelling:

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

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



***voom* fits a lowess trend to the mean-variance relationship ...**



→ Use weights ( $1/\text{var}$ ) in limma analysis



## **Other useful resources**

RNA-Seq Methods and Algorithms (Part I – Intro and overview of RNA-Seq) 2015 UC Davis Workshop

<https://www.youtube.com/watch?v=96yBPM8lEt8>

RNA-Seq Methods and Algorithms (Part II – Alignment Algorithms) 2015 UC Davis Workshop

<https://www.youtube.com/watch?v=b4tVokh6Law>

RNA-Seq Methods and Algorithms (Part III – Quantification) 2015 UC Davis Workshop

[https://www.youtube.com/watch?v=ztyjiCCt\\_IM](https://www.youtube.com/watch?v=ztyjiCCt_IM)

RNA-Seq Methods and Algorithms (Part IV – Differential Expression) 2015 UC Davis Workshop

<https://www.youtube.com/watch?v=BRWj6re9iGc>