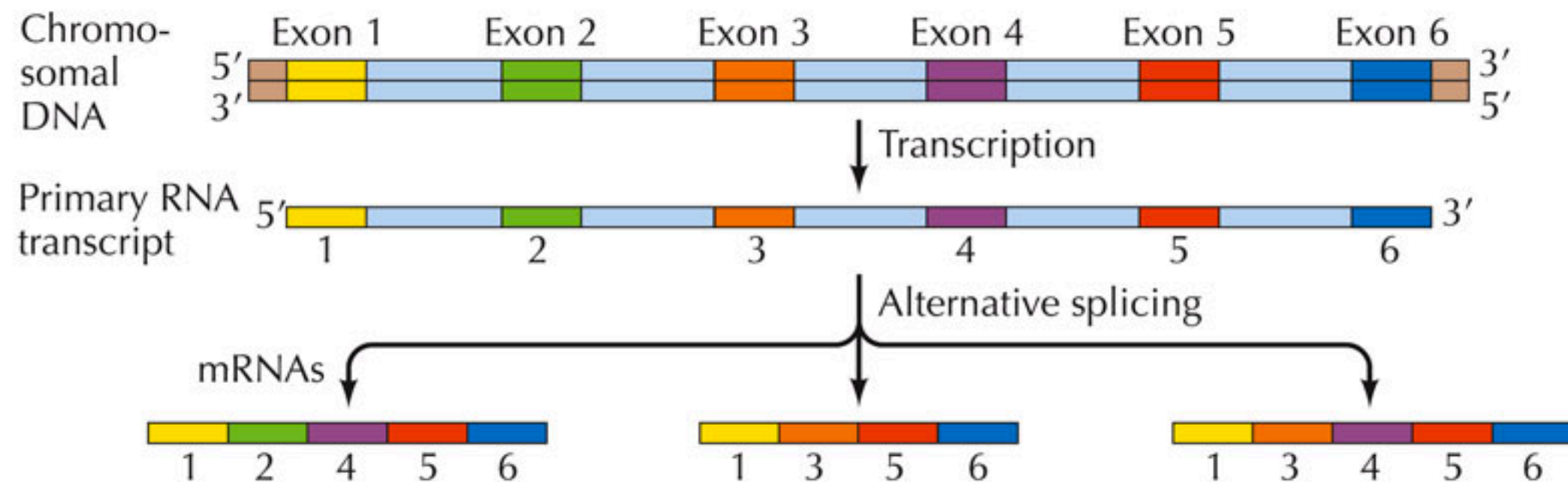


# Differential expression

Alejandro Reyes  
T: @areyesq89

Workshop on Transcriptomics  
September 13th, 2017

# Overview of exons, genes and transcripts



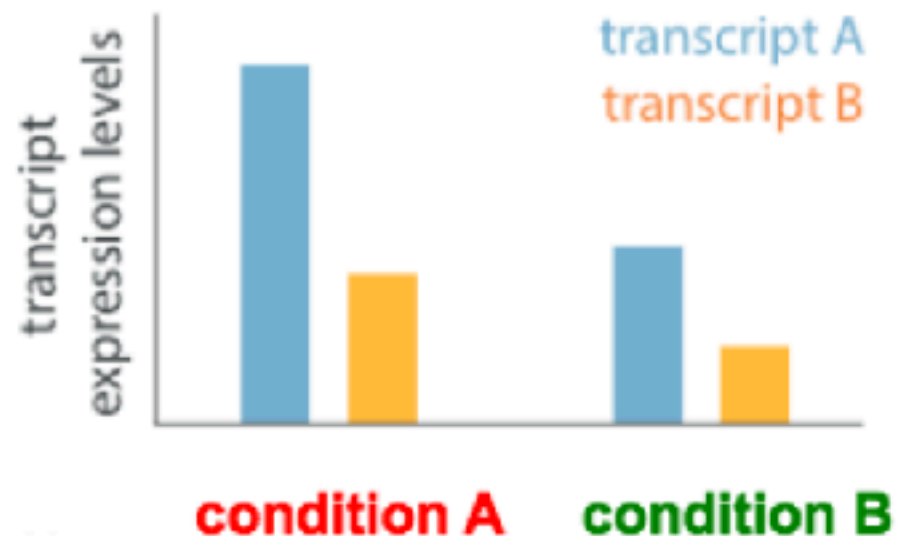
# What is your biological question?

Given a gene, test for:

- Whether transcripts levels change between conditions? (differential exon usage, DGE)
- Whether transcript isoform proportions change between conditions? (differential transcript usage, DTU)
- Whether individuals exons are differentially used? (differential exon usage, DEU)

# Differential problems: DGE

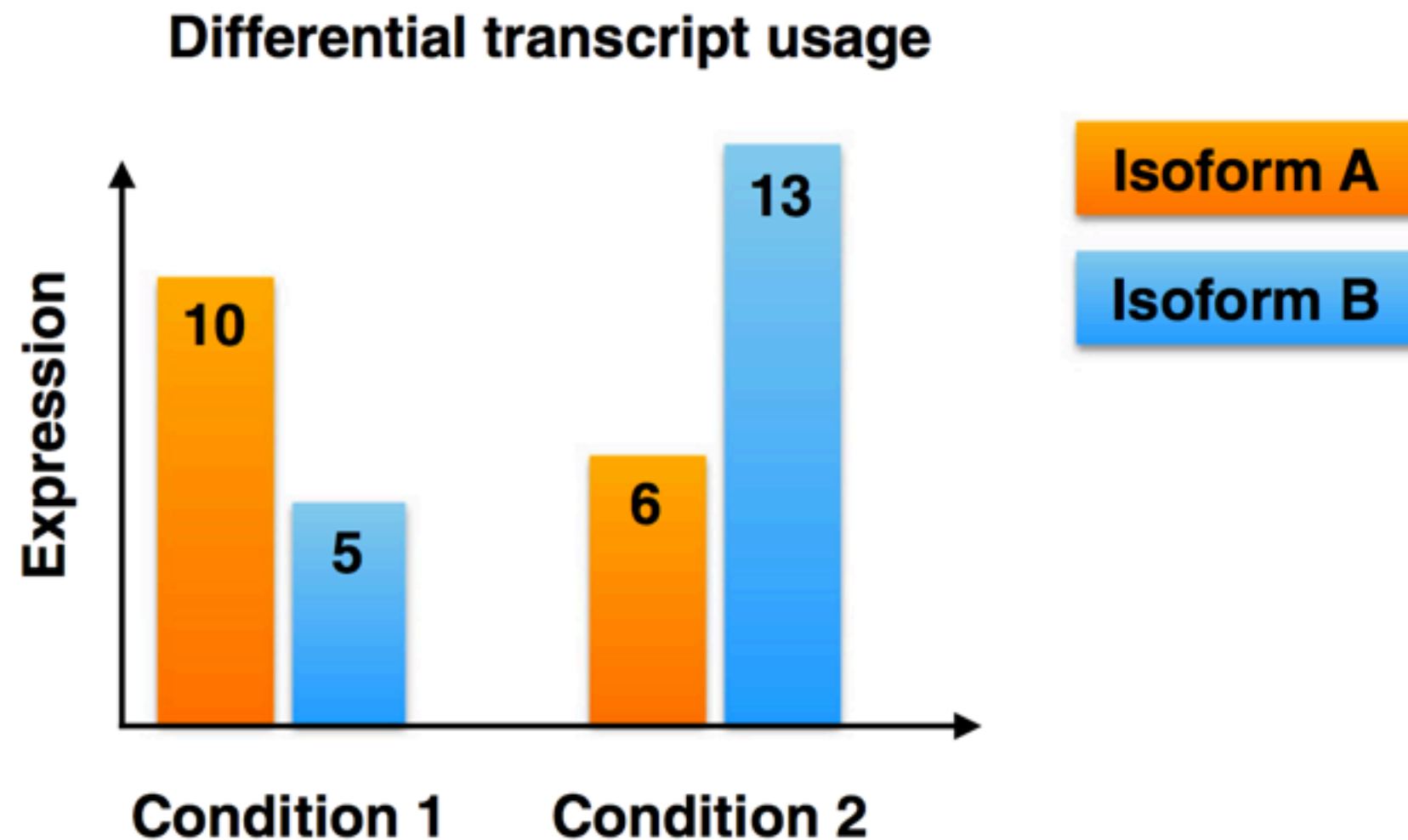
Differential transcript  
expression (DTE)



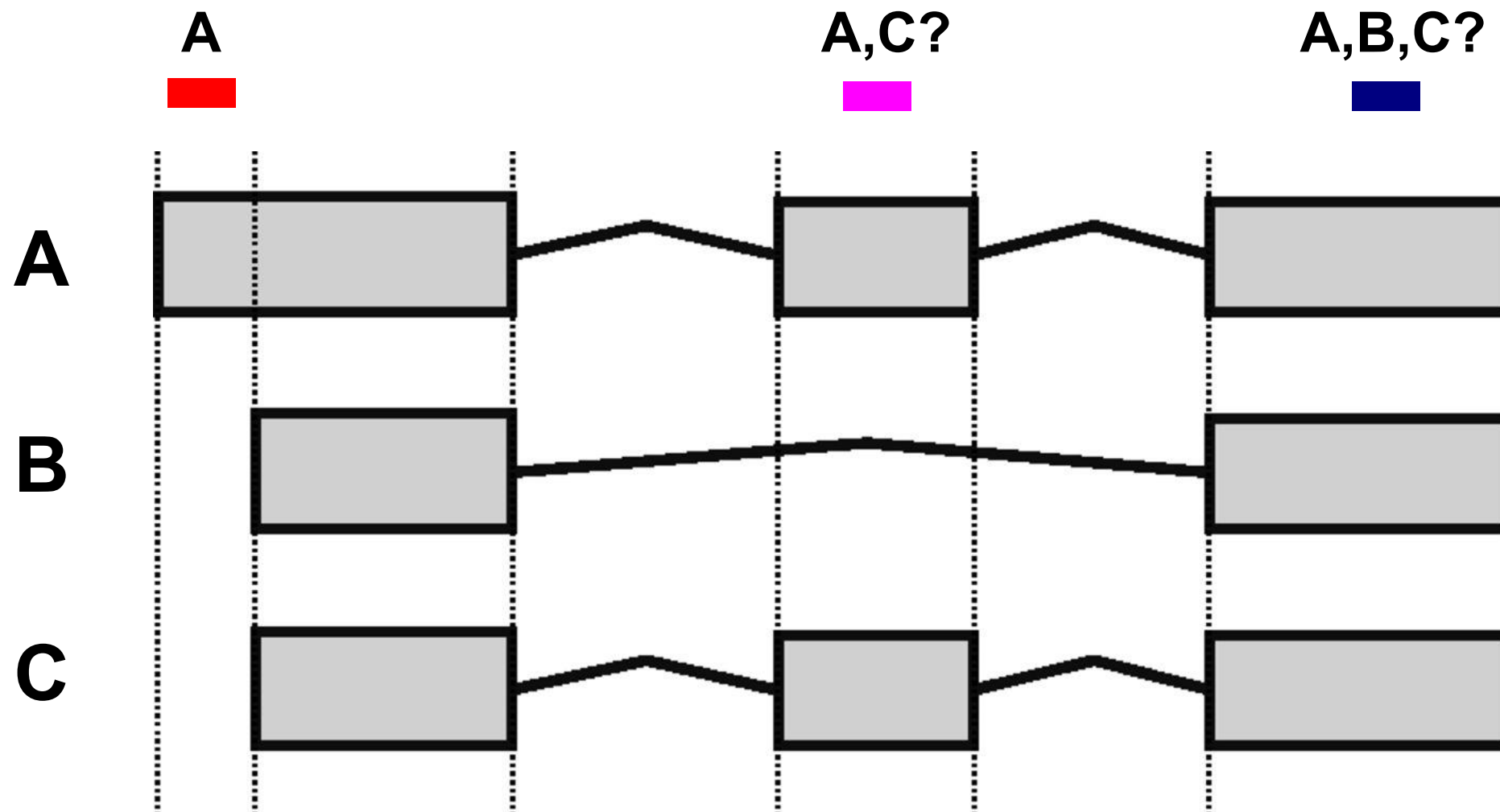
Differential gene  
expression (DGE)



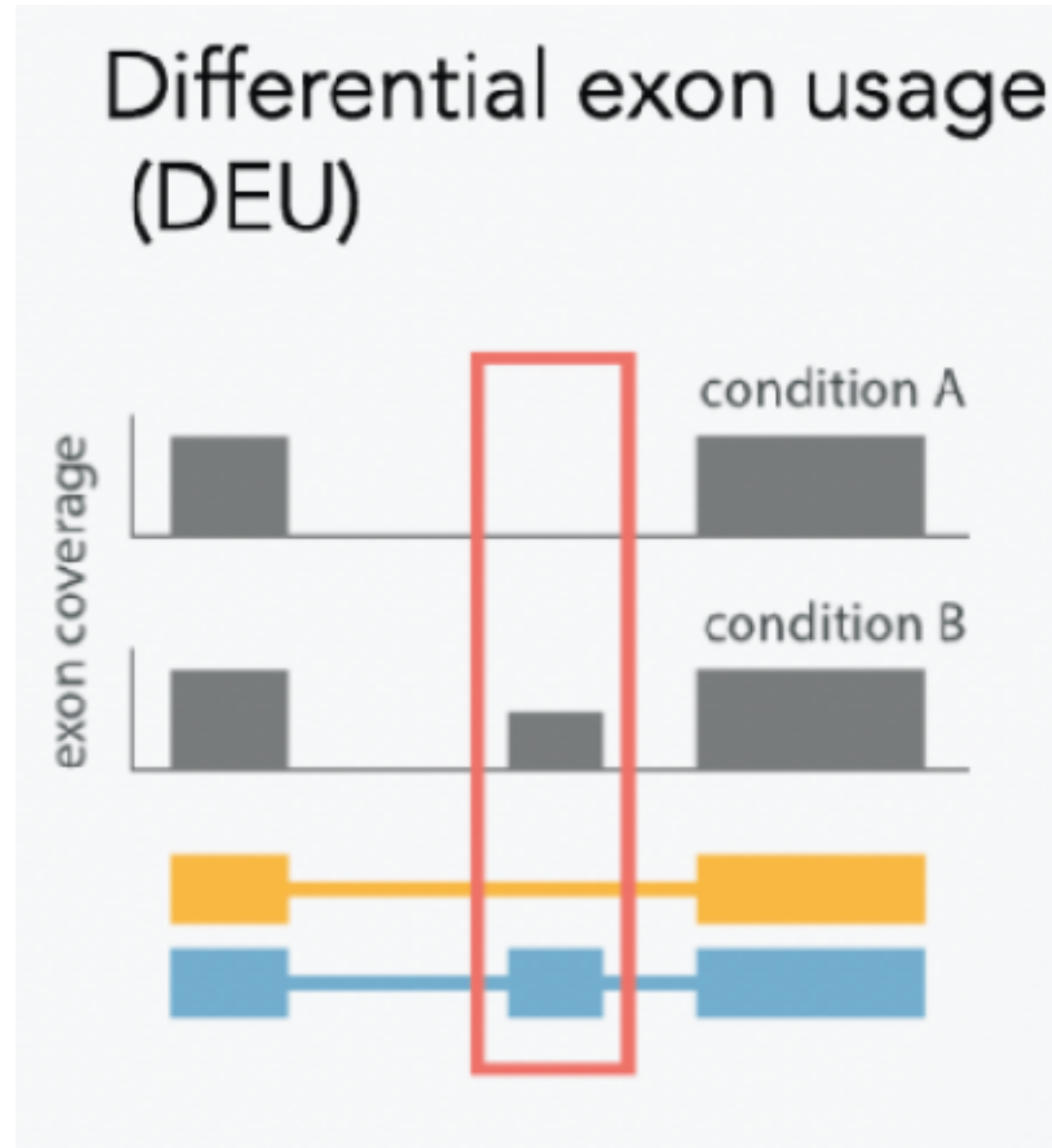
# Differential problems: DTU



# Unambiguous assignment of reads to single transcripts



# Differential problems: DEU



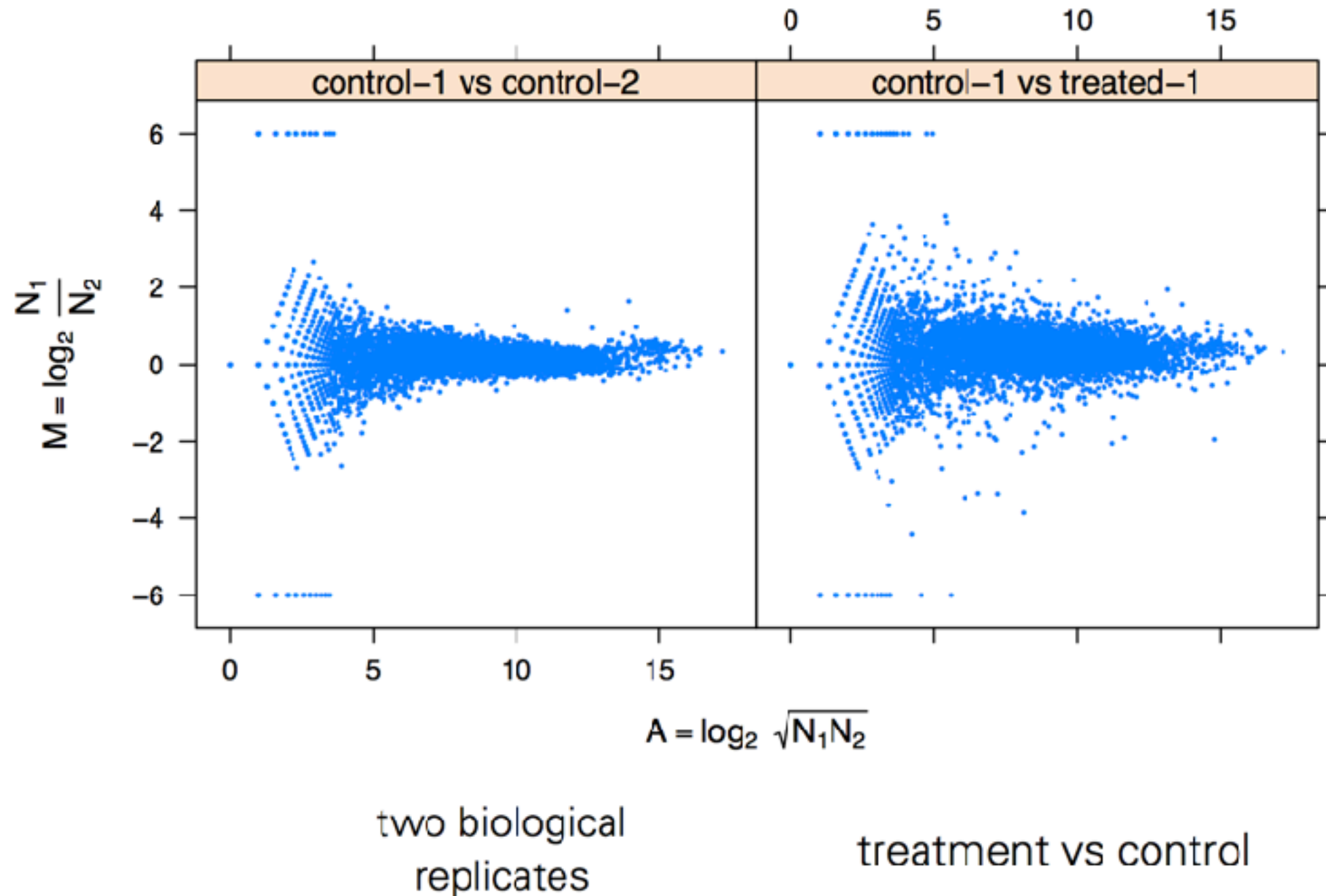
# What is your biological question?

Given a gene, test for:

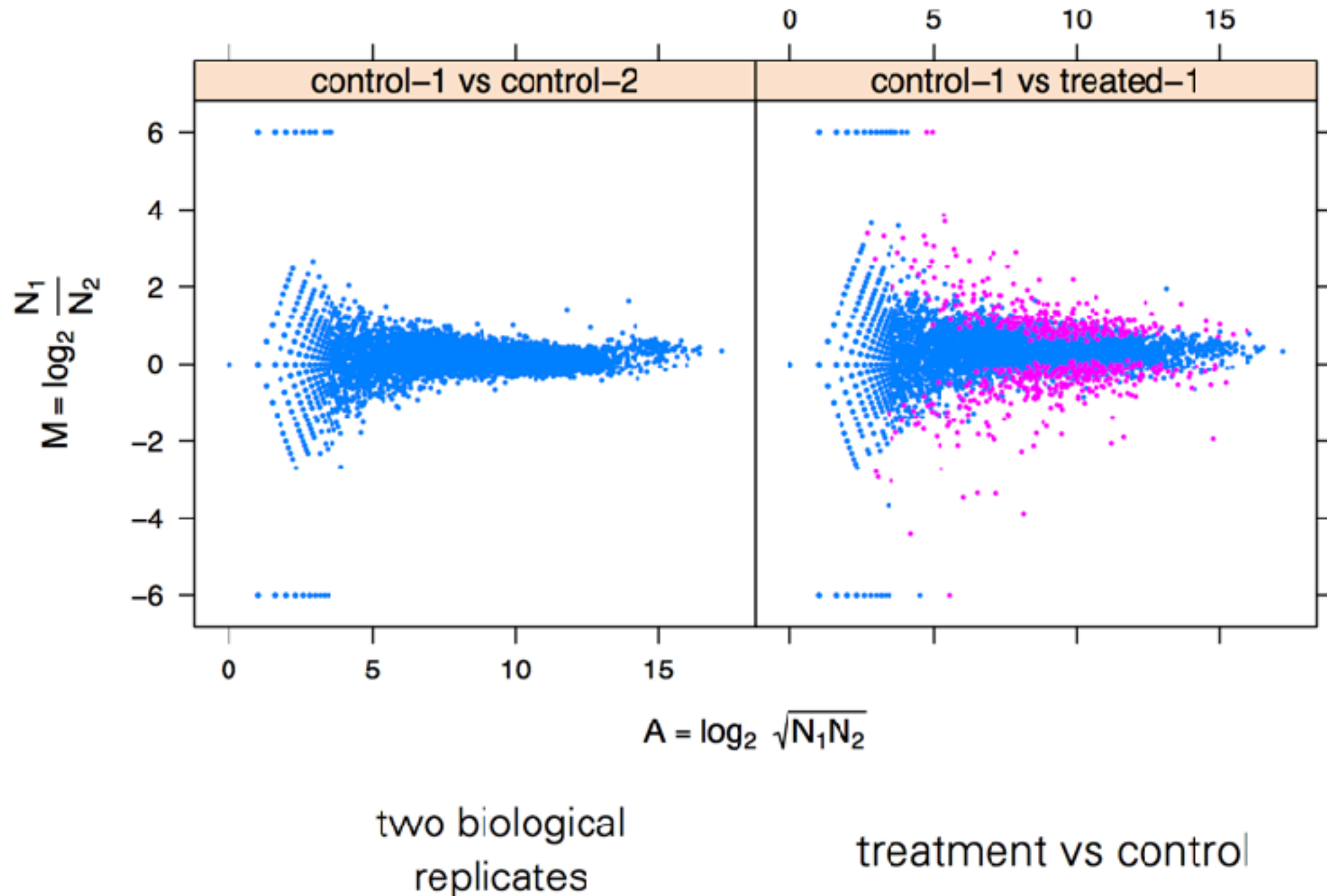
- **Whether transcripts levels change between conditions? (differential gene expression, DGE)**
- Whether transcript isoform proportions change between conditions? (differential transcript usage, DTU)
- Whether individuals exons are differentially used? (differential exon usage, DEU)

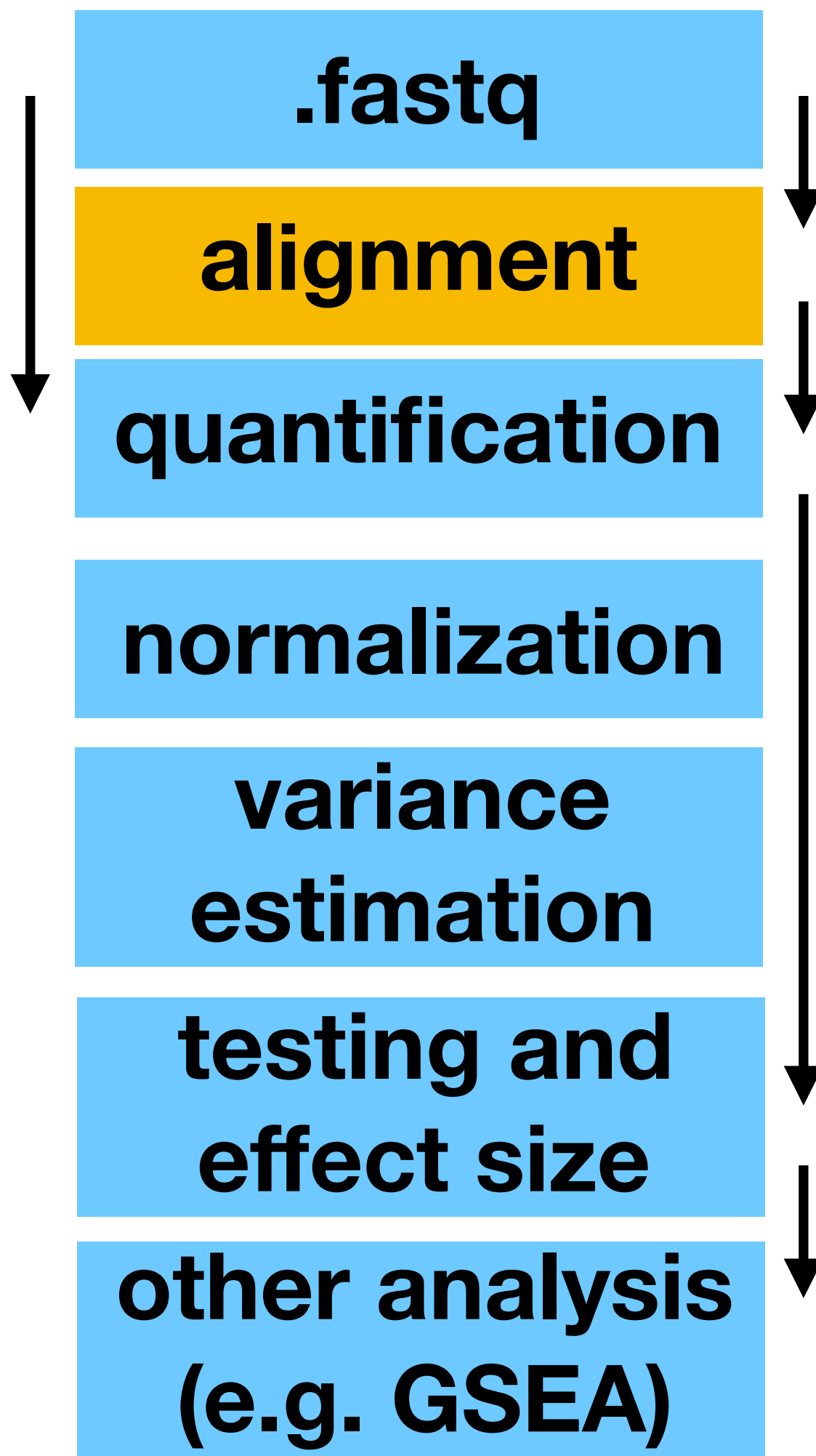


# Differential problems: DGE



# Differential problems: DGE





salmon,  
kallisto,  
...

**.fastq**

**alignment**

**quantification**

**normalization**

**variance  
estimation**

**testing and  
effect size**

**other analysis  
(e.g. GSEA)**



STAR, GSNAP, ...



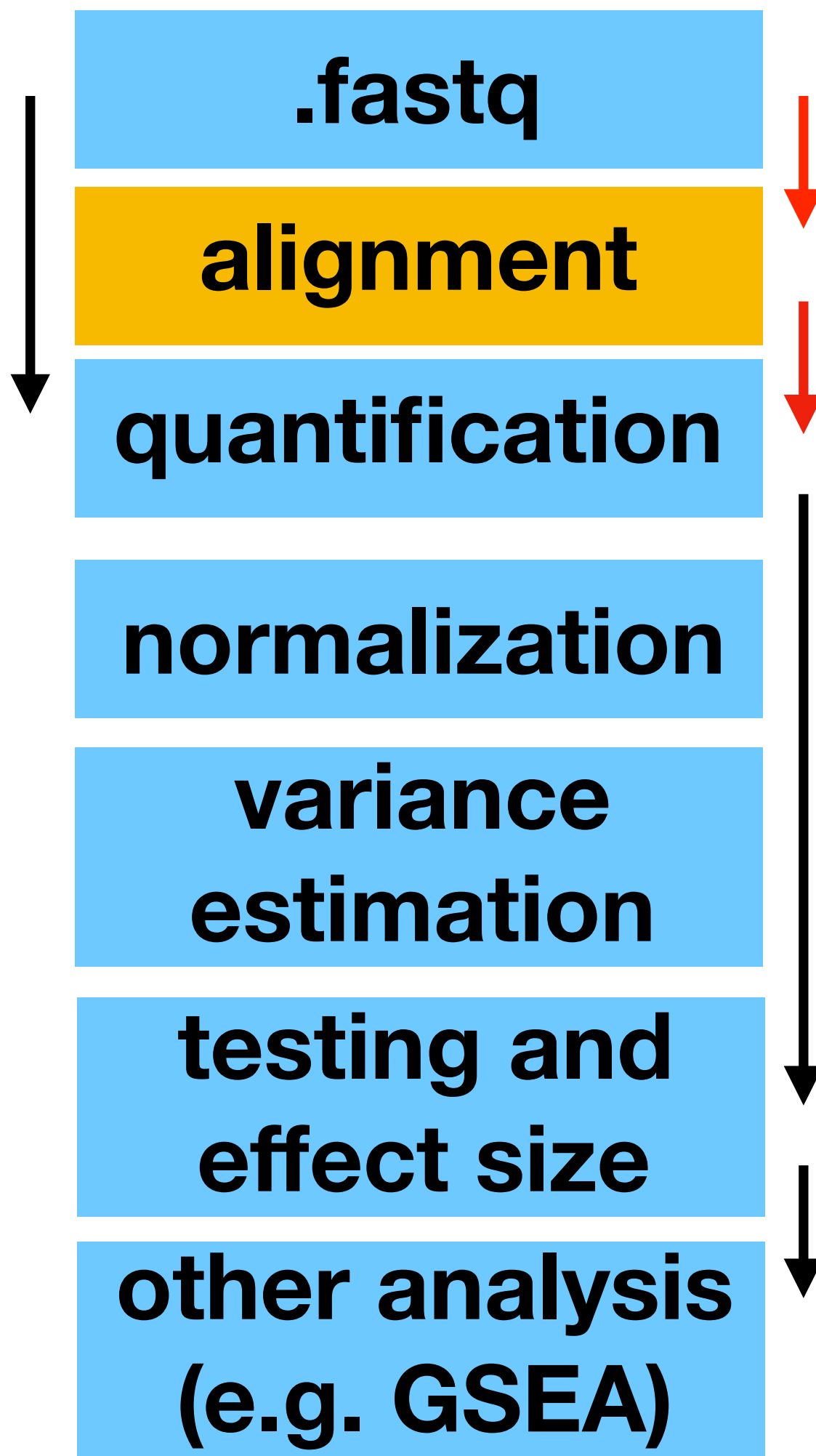
ht-seq, featureCounts,



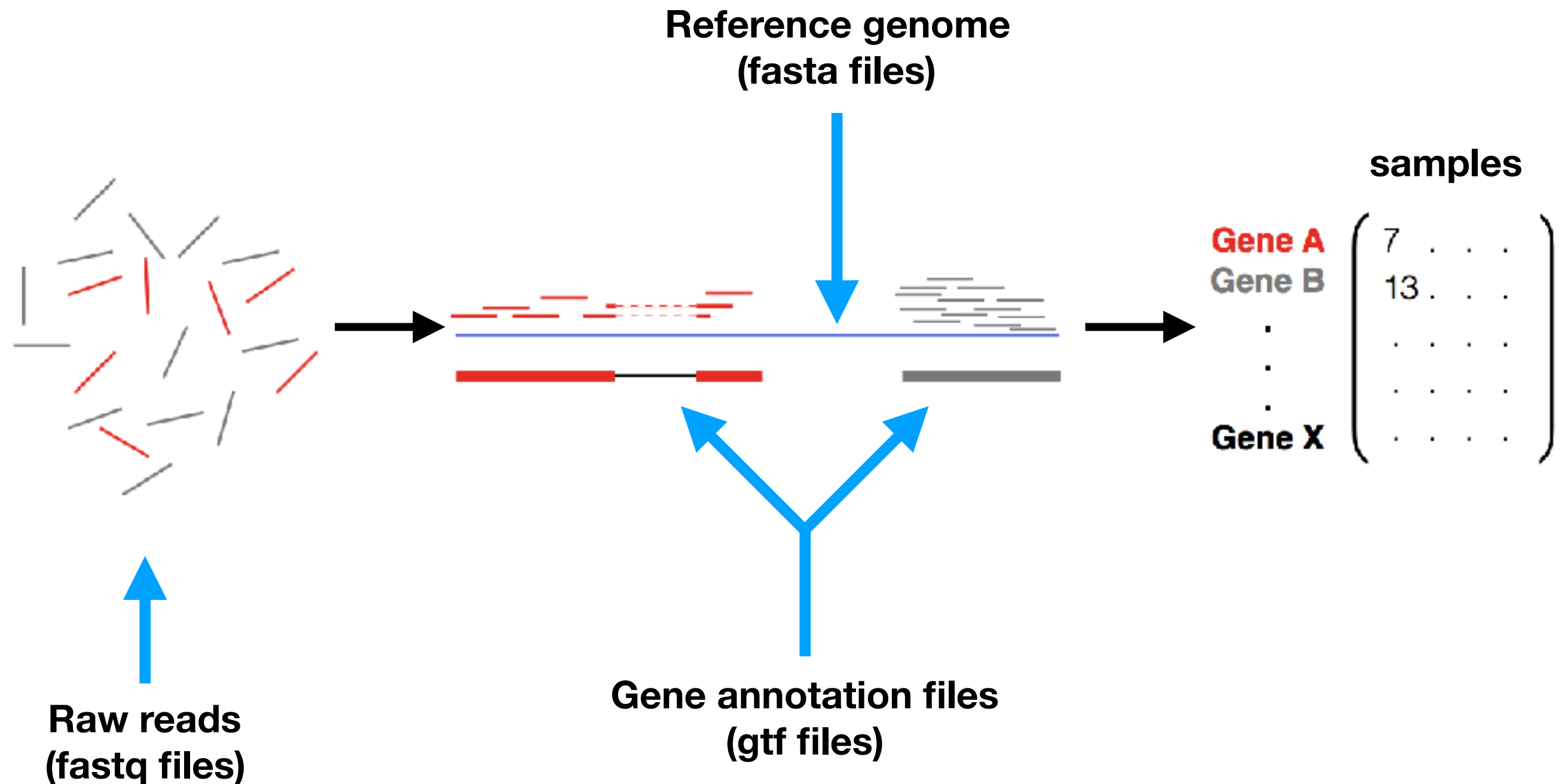
DESeq2,  
edgeR,  
limma-voom  
NOISeq,  
sleuth,...



goseq, roast, ...

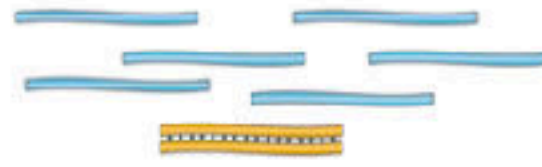


# Alignment-based abundance estimation workflow



**a Data generation**

① mRNA or total RNA



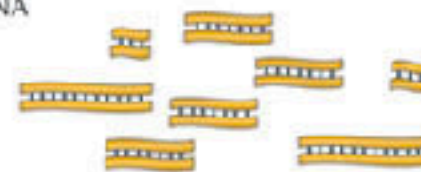
② Remove contaminant DNA



③ Fragment RNA

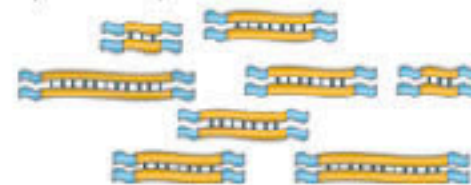


④ Reverse transcribe into cDNA

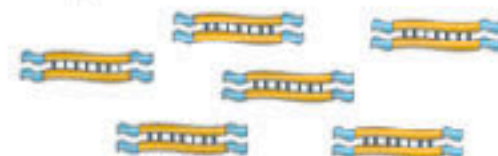


Strand-specific RNA-seq?

⑤ Ligate sequence adaptors



⑥ Select a range of sizes

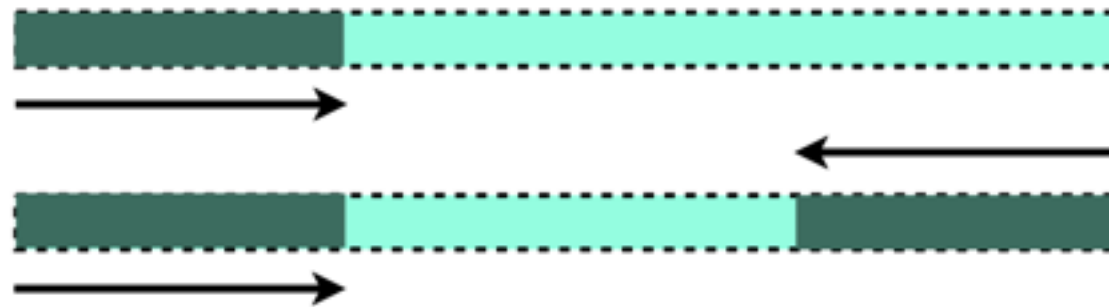


PCR amplification?

⑦ Sequence cDNA ends



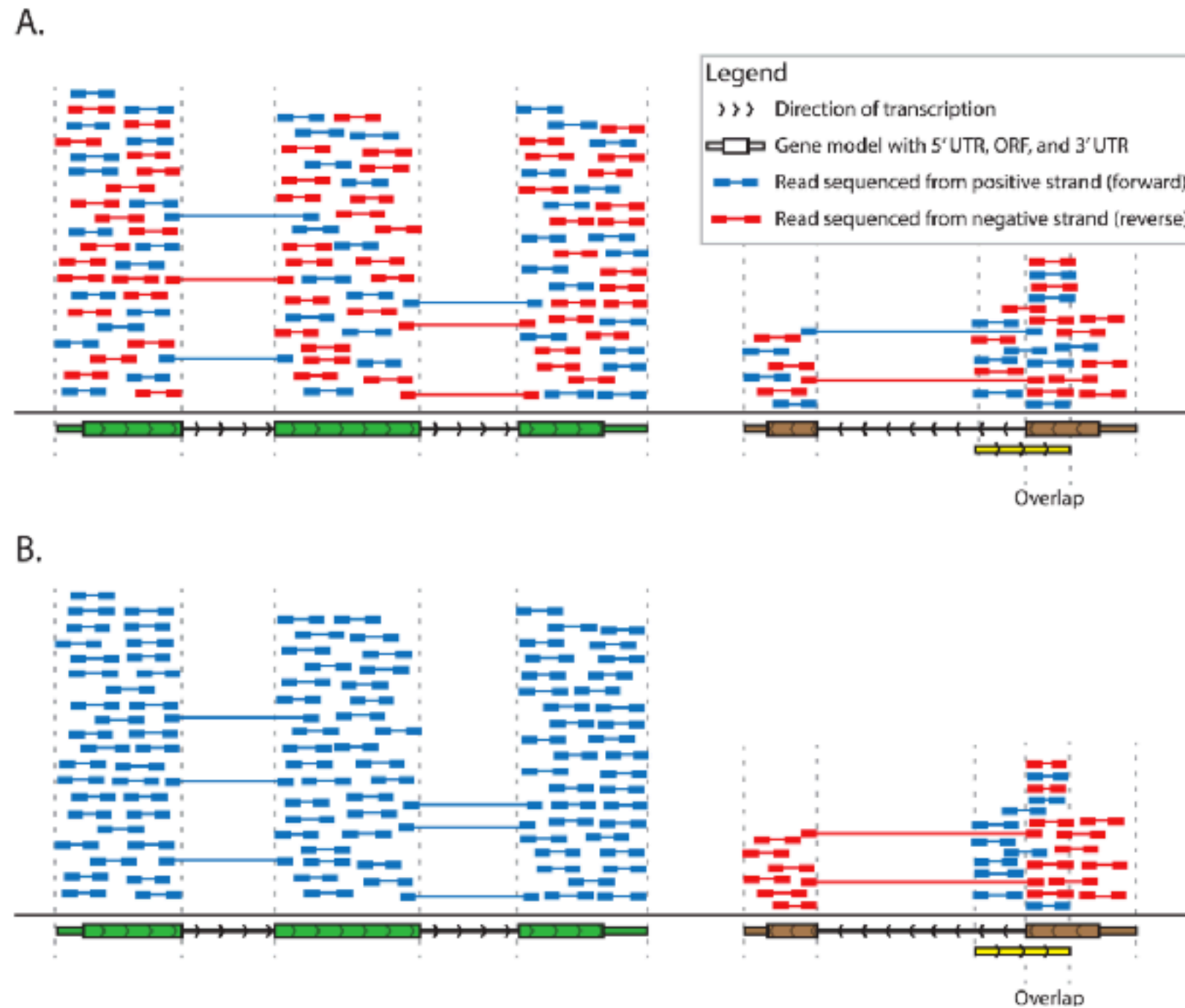
# Are my data single-end or paired-end?



**Tip: Look at the number of files per sample**

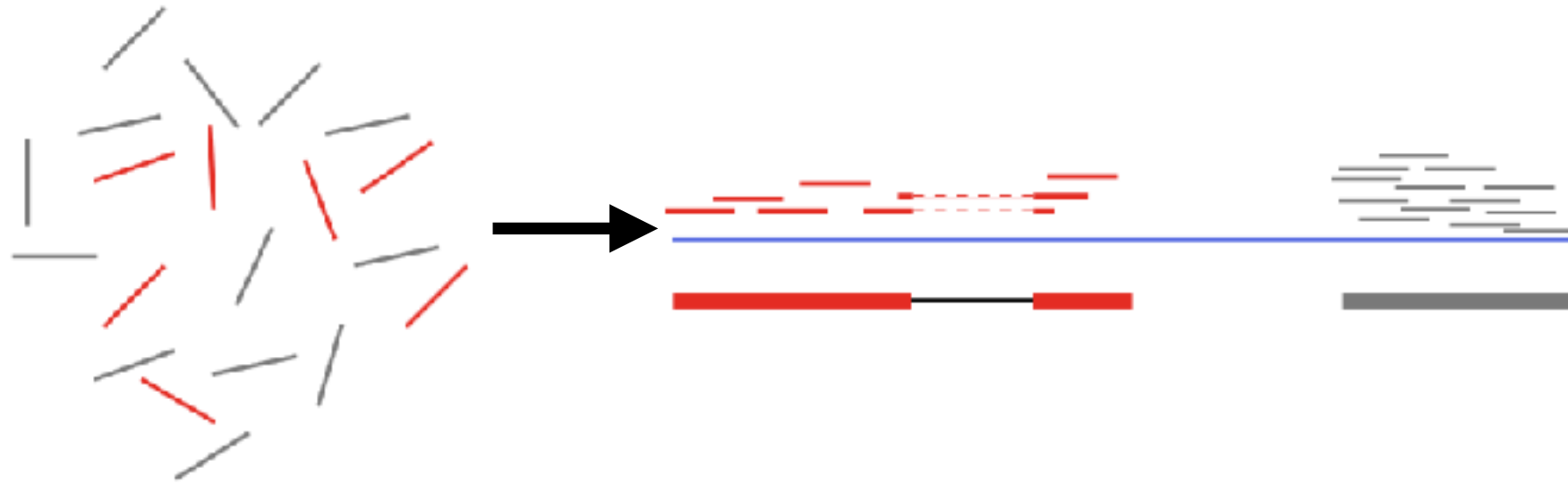


# Are my data strand specific?



**Tip: Visualize mapped reads in a genome browser (e.g. IGV)**

# Alignment



**Need a splice aware aligner  
(e.g. STAR, GSNAP, ...)**

NATURE METHODS | ANALYSIS



## Simulation-based comprehensive benchmarking of RNA-seq aligners

Giacomo Baruzzo, Katharina E Hayer, Eun Ji Kim, Barbara Di Camillo, Garret A FitzGerald & Gregory R Grant

NATURE METHODS | ANALYSIS [OPEN](#)



## Systematic evaluation of spliced alignment programs for RNA-seq data

Pär G Engström, Tamara Steijger, Botond Sipos, Gregory R Grant, André Kahles, The RGASP Consortium, Gunnar Rätsch, Nick Goldman, Tim J Hubbard, Jennifer Harrow, Roderic Guigó & Paul Bertone

# Example using STAR: index generation

```
$ STAR --runThreadN 24 \
--runMode genomeGenerate \
--genomeDir my_genome \
--genomeFastaFiles my_genome.fa \
--sjdbGTFfile my_genes.gtf \
--sjdbOverhang 99
```

number of threads

output folder -  
name according  
to genome

reference  
genome

gene  
annotation  
file

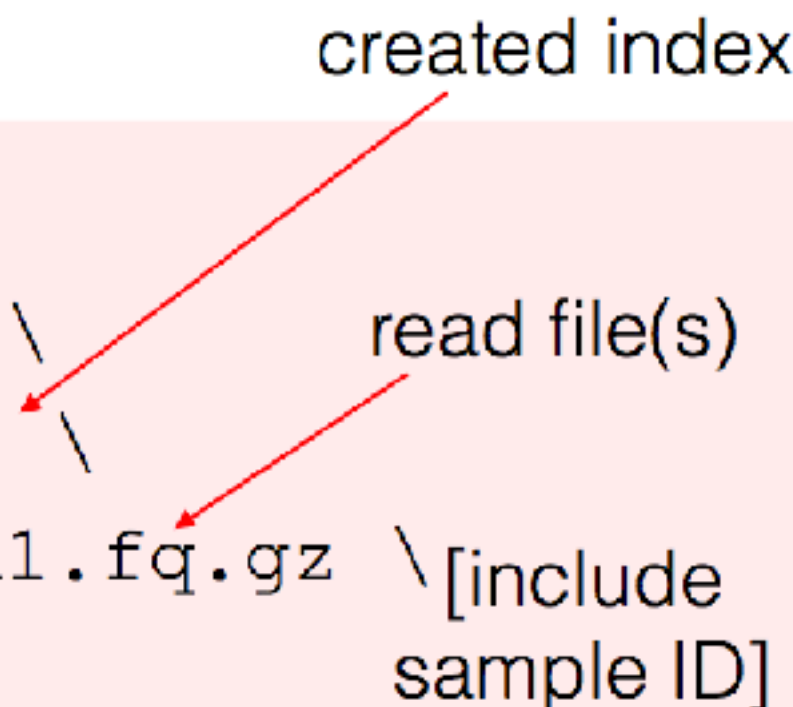
read length - 1

# Example using STAR: alignment


```
$ STAR --runThreadN 24 \  
      --runMode alignReads \  
      --genomeDir my_genome \  
      --readFilesIn S1_read1.fq.gz \[include  
      S1_read2.fq.gz \sample ID]
```

created index

read file(s)



# Aligner output: BAM files




The screenshot shows a file explorer window with a directory containing several files. The files are:

- SRR1039508
- SRR1039509
- SRR1039512
- SRR1039508\_Aligned.sortedByCoord.out.bam
- SRR1039508\_Log.final.out
- SRR1039508\_Log.out

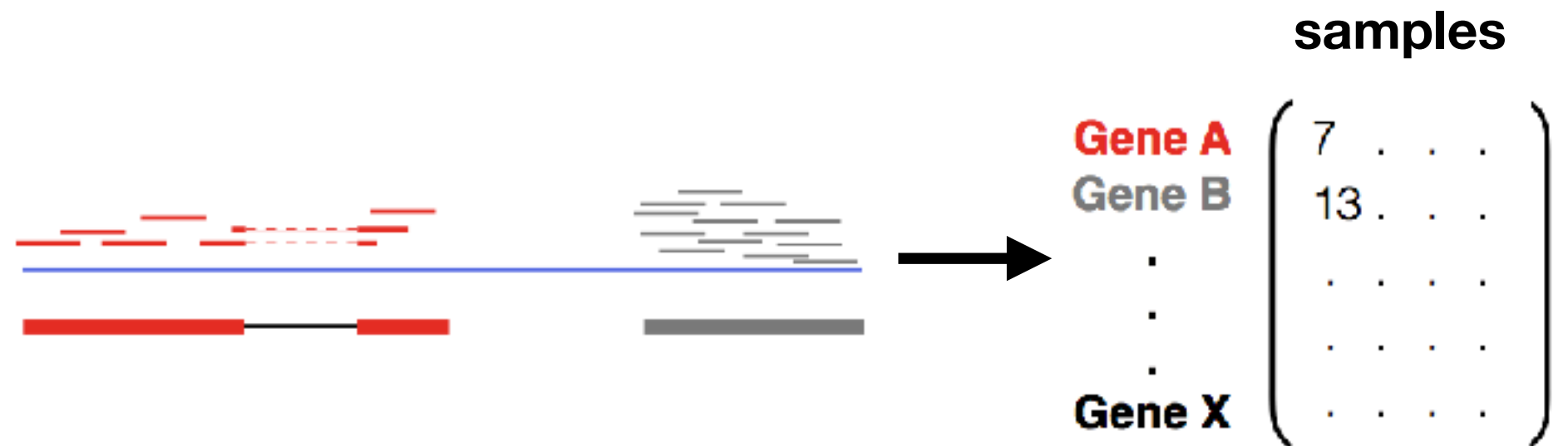
Below the file explorer, a snippet of SAM/BAM alignment data is shown. The data is organized into columns corresponding to the fields in the SAM format:

read name	flag	chr	pos	mapq	CIGAR	sequence	quality
1	1	11307	1	76M	=	11648	417
CTGCAGGGCCCTCTTGCTTACTGTATAGTGGT							
NH:i:4							
HI:i:2							
AS:i:150							
1	1	11579	0	76M	=	11863	359
CAGAATTGTACTGTTCTGTATCCCACCAGCAA							
NH:i:6							
HI:i:3							
AS:i:145							
1	1	11606	0	69M7S	=	11652	460
AGCAATGTCTAGGAATACCTGTTTCTCCACAA							
NH:i:5							
HI:i:4							
AS:i:141							
1	1	11625	0	47M338N29M	=	12083	534
TGTTTCTCCACAAAGTGTTTACTT							
NH:i:6							
HI:i:5							
AS:i:1							
1	1	11648	1	76M	=	11307	-417
TTTGGATTTTGGCCAGTCTAACAGGTGAAGCC							
NH:i:4							
HI:i:2							
AS:i:150							

read name    flag    chr    pos    mapq    CIGAR



# Alignment-based abundance estimation workflow



```
format bam \
--order=pos \
--stranded=no \
--type=exon \
--idattr=gene_id \
--mode=union \
aligned.bam \
my_genes.gtf
```

default=yes

check you  
GTF file!

```
protein_coding exon 5139815 5141712 . - . gene_id "FBgn0020621"; transcript_id  
Btr0112897"; exon_number "10"; gene_name "Pkn"; gene_biotype "protein_coding"; transcript_name "Pkn-RG"  
on_id "FBgn0020621:1";
```

mapped reads

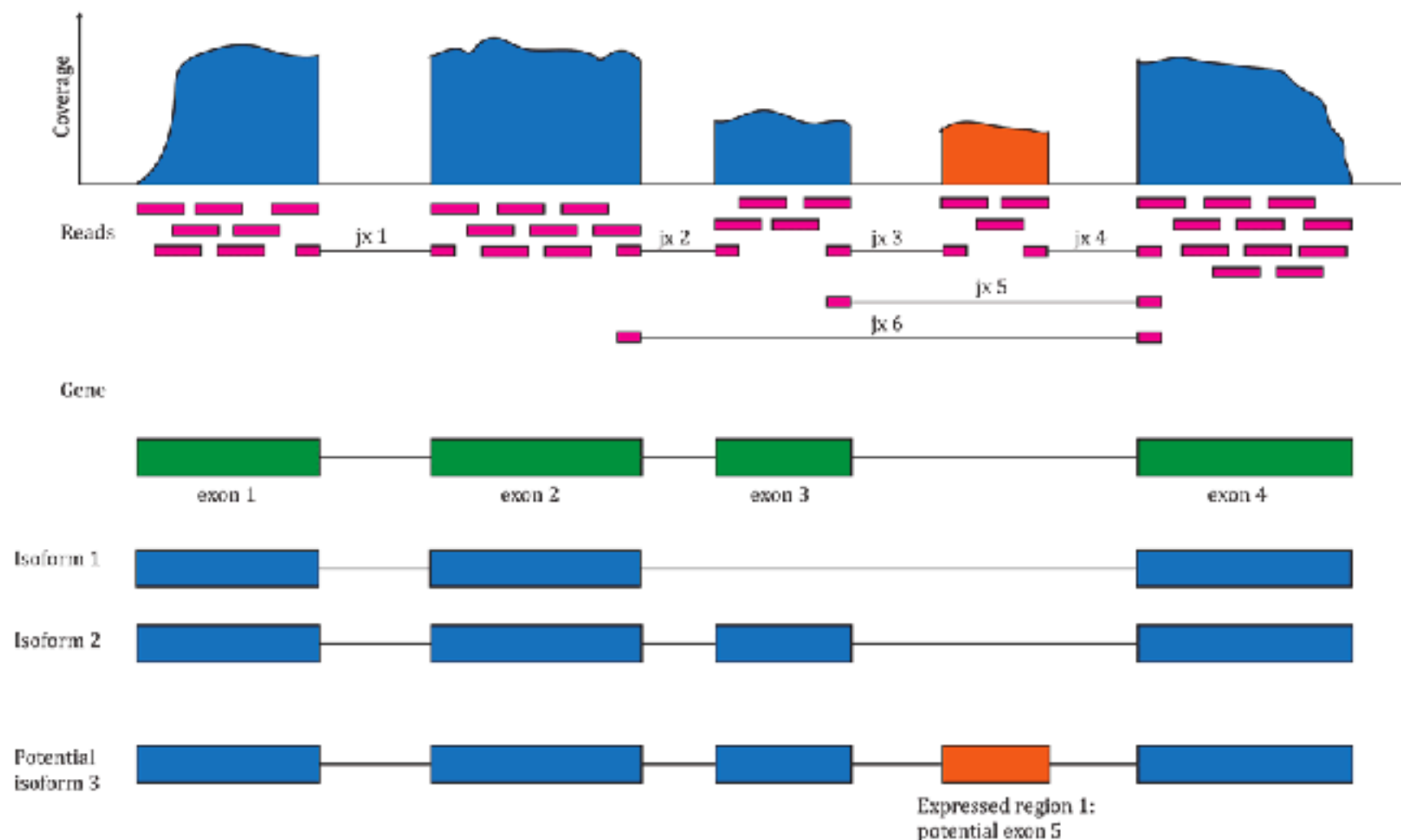
gene

annotation file (GTF)



# recount2

70,000 human RNA-seq samples  
already processed







**.fastq**

**alignment**

**quantification**

**normalization**

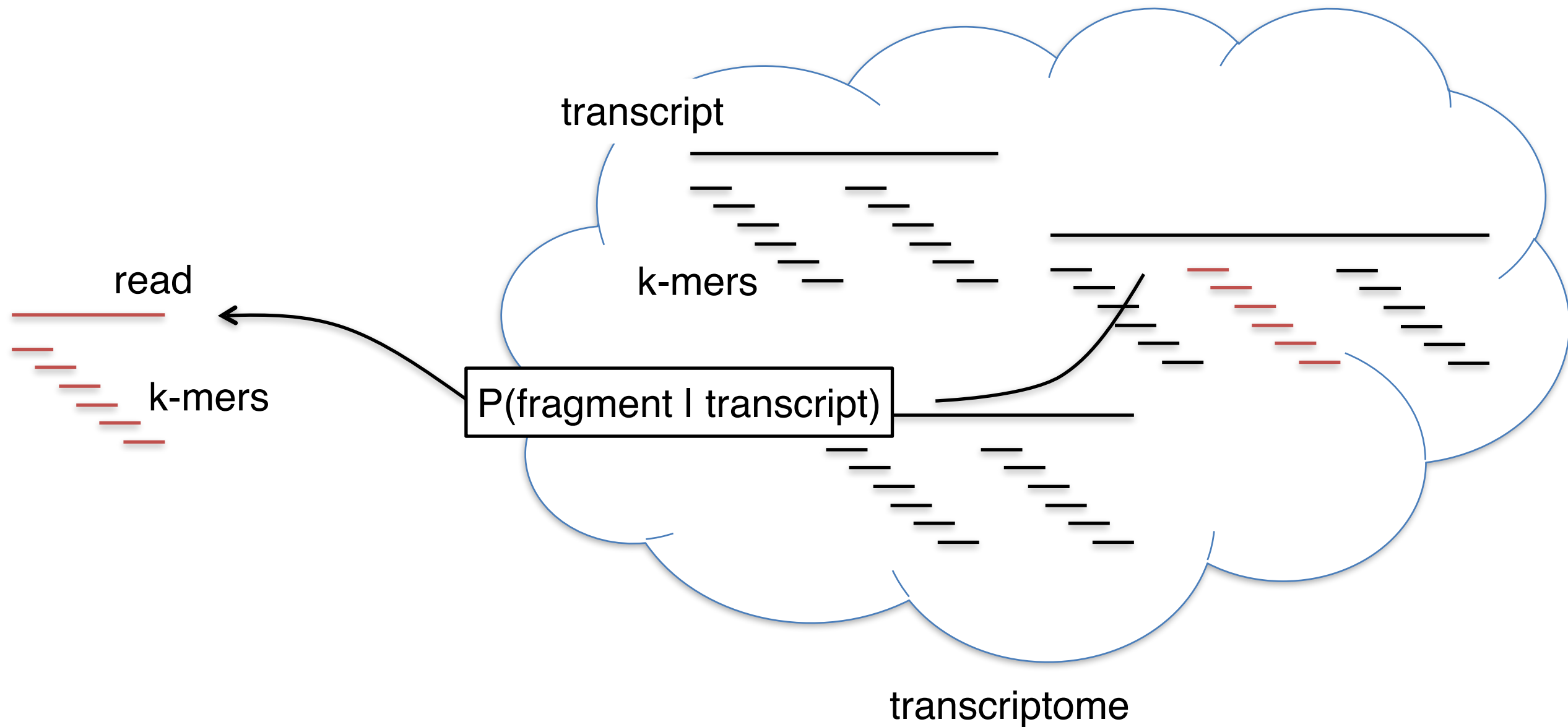
**variance  
estimation**

**testing and  
effect size**

**other analysis  
(e.g. GSEA)**

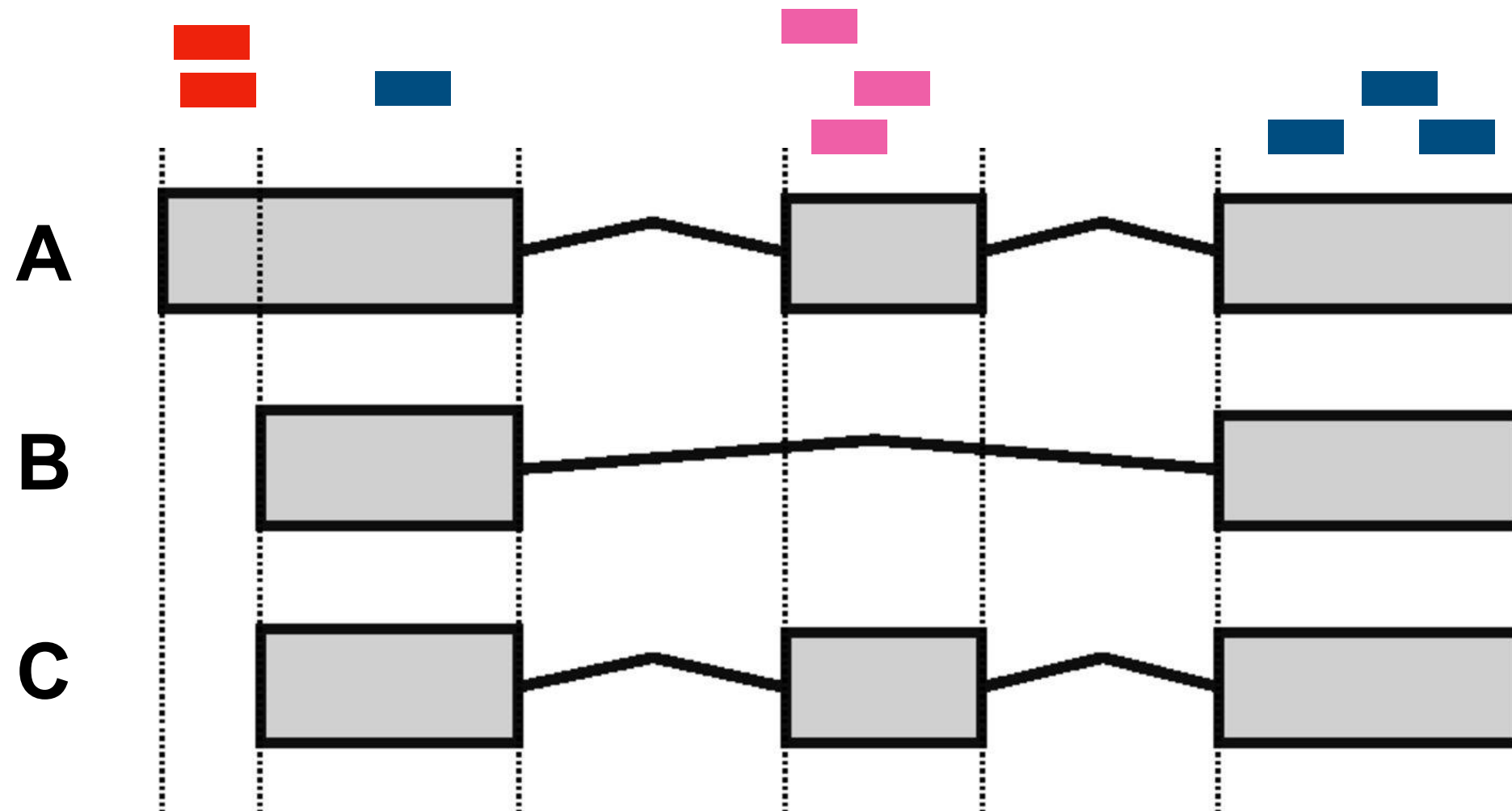


# Alignment-free abundance estimation workflows



**Extremely fast compared to genome alignments!**

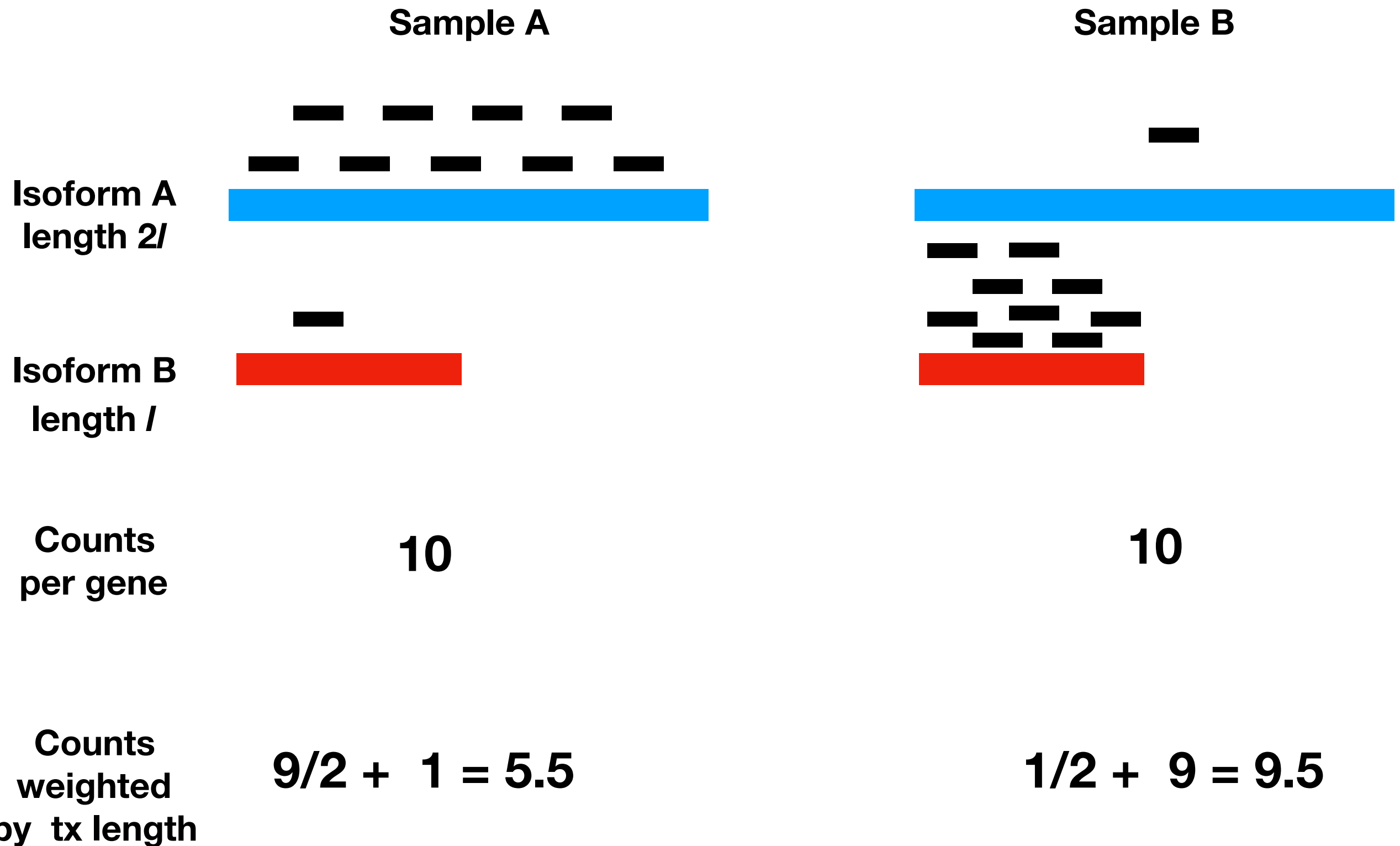
# Alignment-free abundance estimation workflows



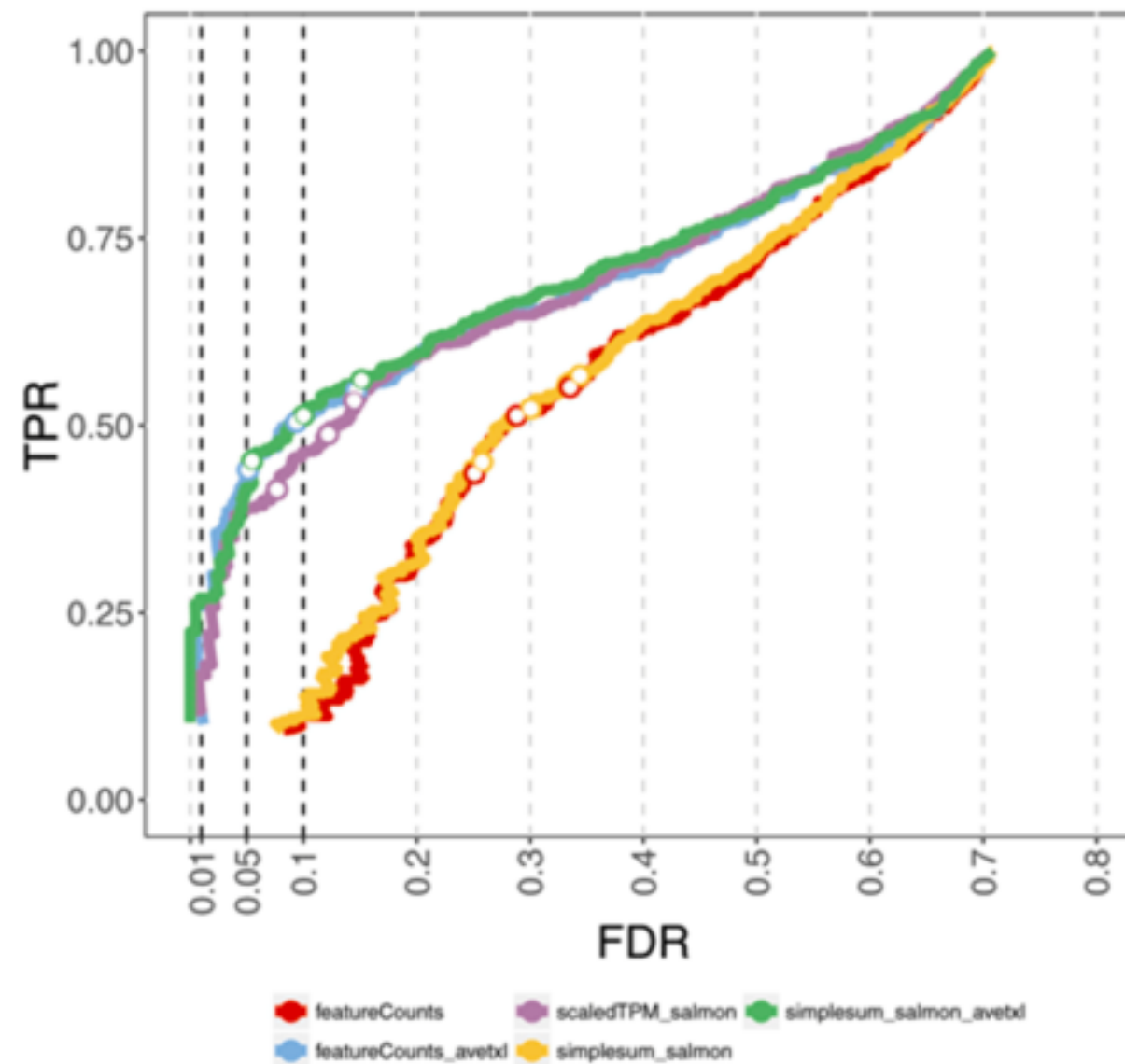
Equivalence classes:  $\{A\} = 2$ ,  $\{A, C\} = 3$ ,  $\{A, B, C\} = 4$

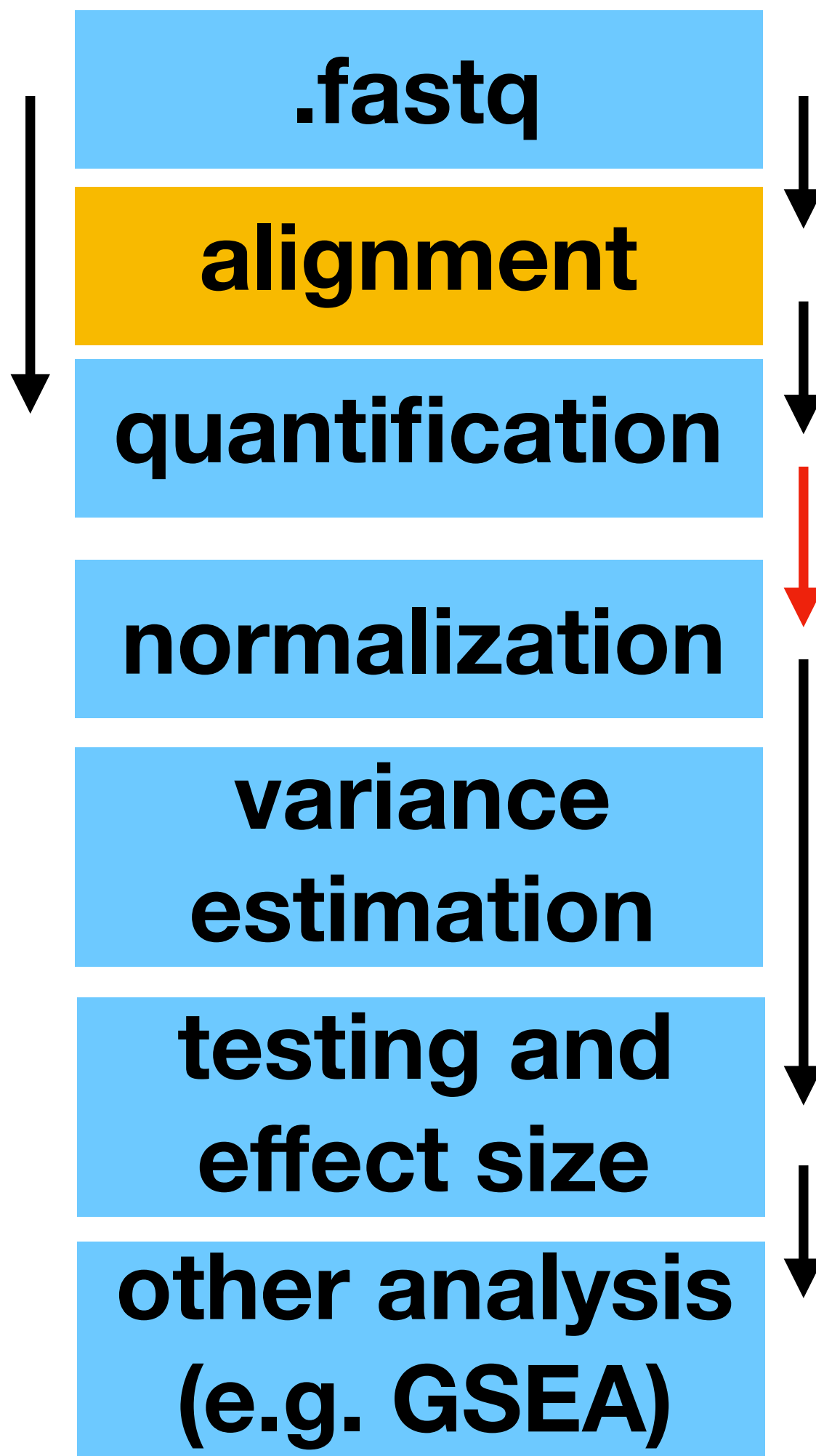
Objective: estimate the transcript abundances that best explain the observed counts for the reference classes.

# Why is it important to consider transcript length information? (hypothetical example)



# Considering average transcript lengths improves estimates of DGE





$$\begin{array}{l}
 \text{Gene A} \\
 \text{Gene B} \\
 \vdots \\
 \text{Gene X}
 \end{array}
 \begin{pmatrix}
 7 & . & . & . \\
 13 & . & . & . \\
 . & . & . & . \\
 . & . & . & . \\
 . & . & . & .
 \end{pmatrix}$$

# Sequencing depth

sample 1

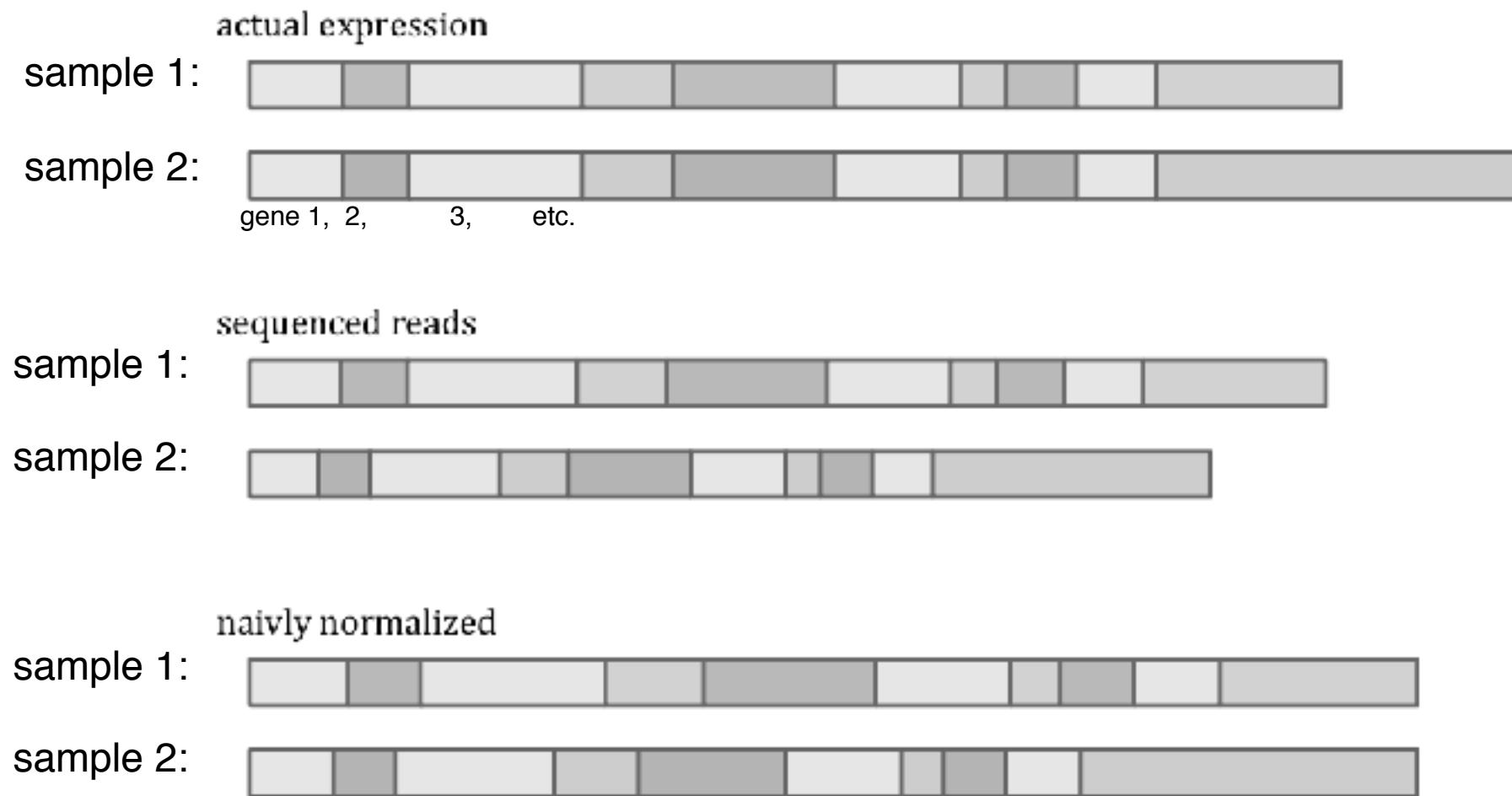


sample 2

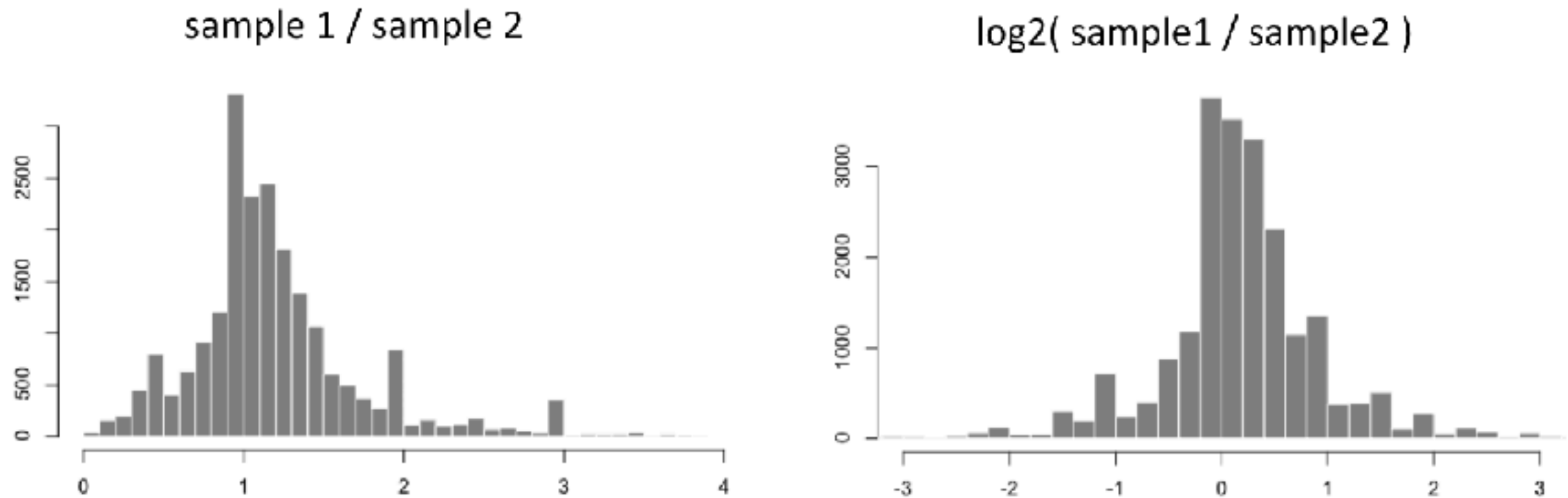




# Normalization for sequencing depth



# DESeq2 uses a median of ratios method



- Create a reference sample by calculating the geometric mean across samples for each gene.
- For each sample, take the ratio with respect to the reference sample.
- Take the median across genes for each sample.

# DESeq2 scaling factors or normalization factors?

$$S_j$$

"size factor"

per sample  $j$

sequencing depth

robustly estimated  
with median ratio

$$S_{ij}$$

"normalization factor"

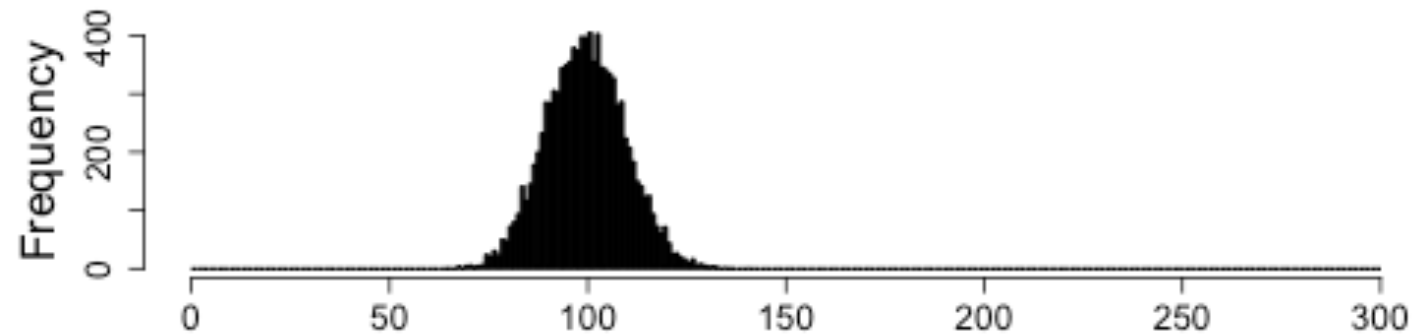
per sample  $j$  & per gene  $i$

sequencing depth and  
other factors differ across samples  
(technical bias: [cqn](#) or [EDASeq](#))  
(gene length: [tximport](#))

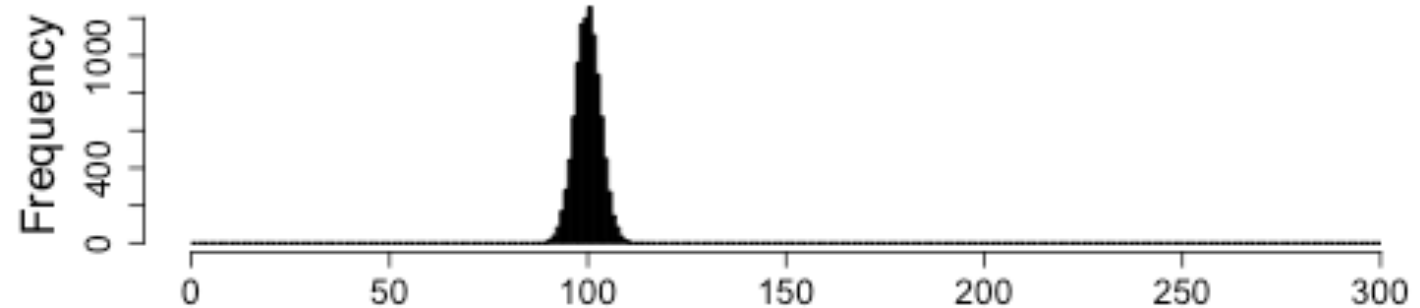
median ratio method for  
sequencing depth can be  
estimated on top

# But, why counts and not other transformed data (e.g. FPKMs)?

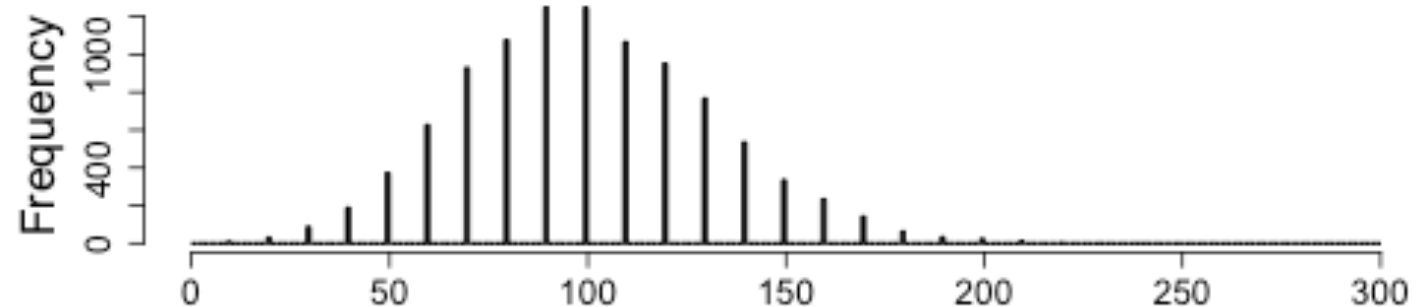
Raw count with mean of 100  
Poisson sampling, so  
SD=10

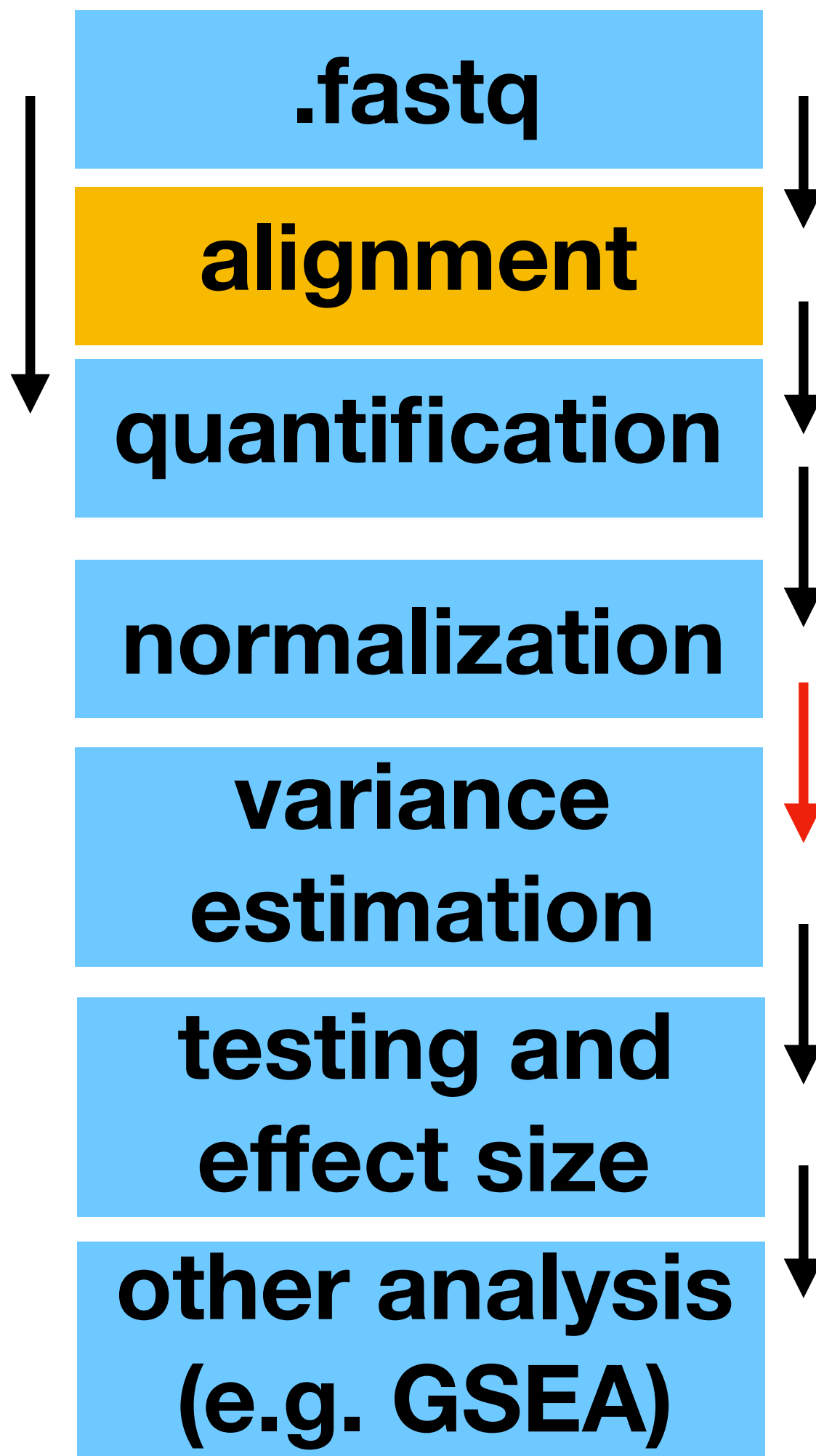


Raw count with mean of 100  
scale by 1/10  
SD = ?

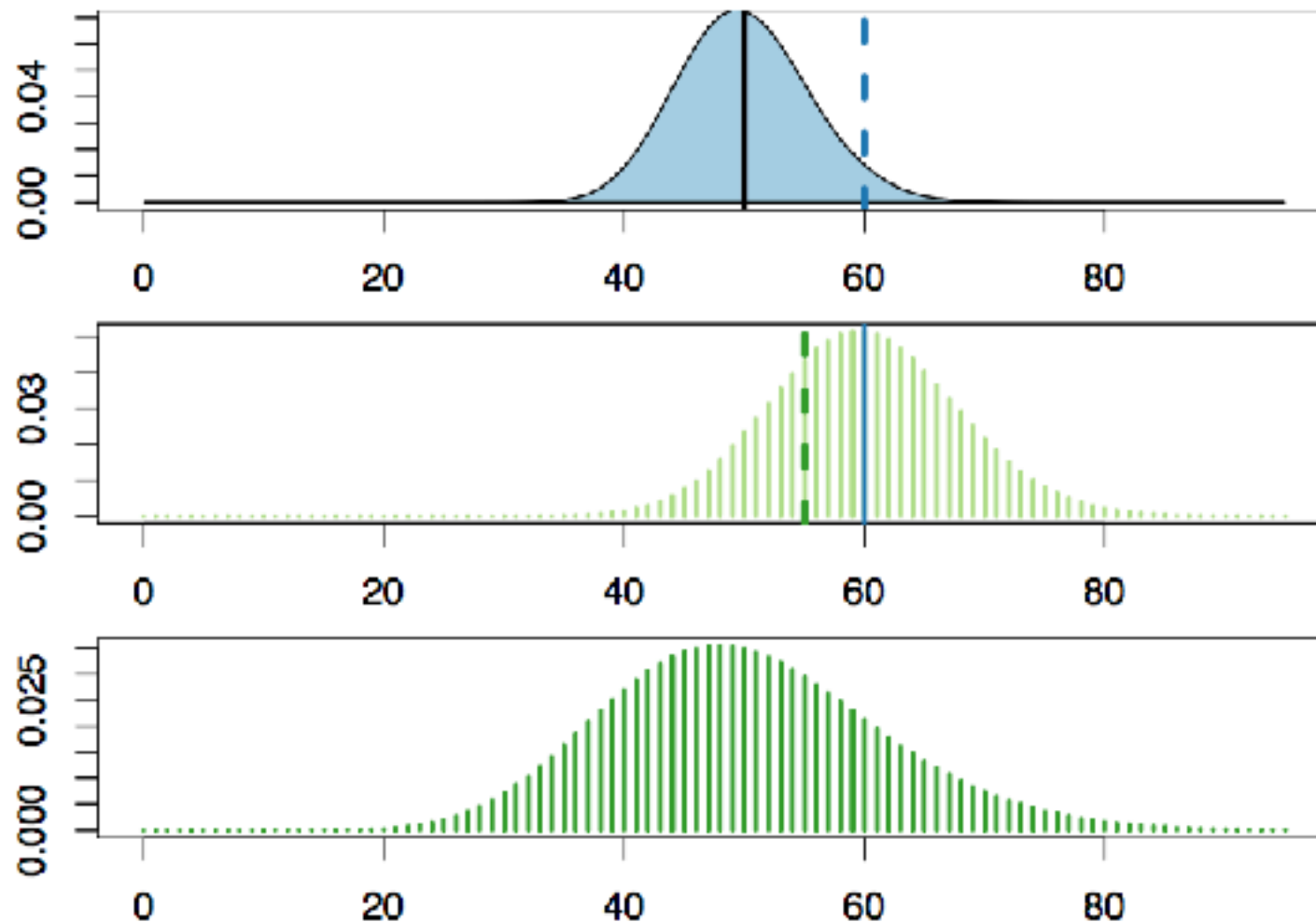


Raw count with mean of 100  
scale by 10  
SD = ?





# Variance of a gene: technical noise + biological noise



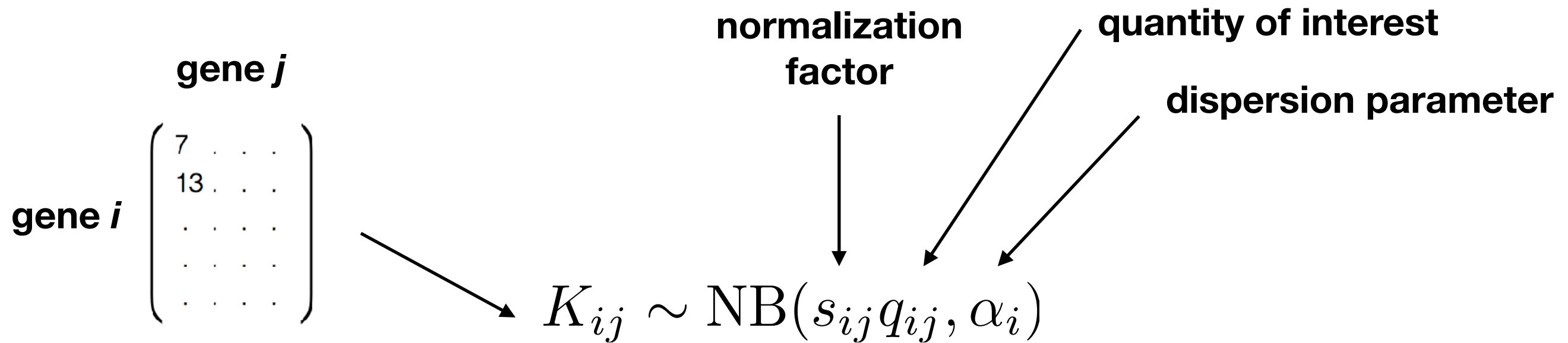
Biological sample with variance  $\Gamma$

Poisson sampling variance =  $\Lambda$

Gamma-poisson process

$$\text{NB}(\mu, \sigma^2 + \mu) = \Lambda(\Gamma(\mu, \sigma^2))$$

Poisson process in which the mean is gamma distributed



$$\text{Var}(K_{ij}) = \underbrace{\mu_{ij}}_{\text{technical noise}} + \alpha_i \underbrace{\mu_{ij}^2}_{\text{extra biological variability}}$$

technical noise

extra biological variability

The equation shows the variance of the count  $K_{ij}$  as the sum of two components. The first component,  $\mu_{ij}$ , is labeled 'technical noise'. The second component,  $\alpha_i \mu_{ij}^2$ , is labeled 'extra biological variability'.

**Challenge: small number of replicates!**

# Dispersion estimation

$$\hat{\alpha}_{\text{MLE}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu}))$$

$$\text{CR}(\alpha) = -\frac{1}{2} \log(\det(X^t W X))$$

$$\hat{\alpha}_{\text{CR}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu}) + \text{CR}(\alpha))$$

$$\text{prior}(\alpha) = f_{\mathcal{N}}(\log(\alpha); \log(\alpha_{\text{fit}}), \sigma_{\alpha\text{-prior}}^2)$$

$$\hat{\alpha}_{\text{CR-MAP}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu}) + \text{CR}(\alpha) + \log(\text{prior}(\alpha)))$$



# Dispersion estimation

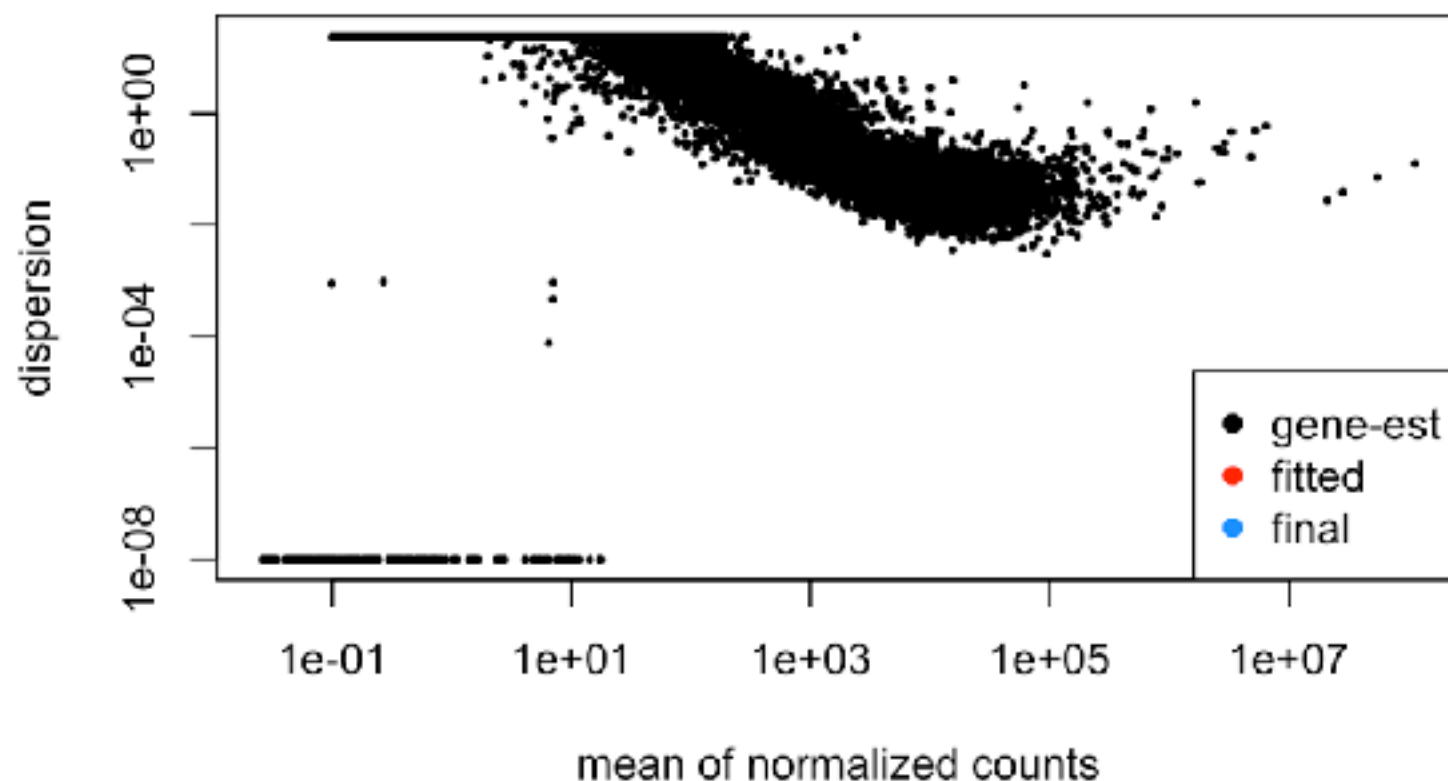
$$\hat{\alpha}_{\text{MLE}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu})) \quad \longleftarrow \quad \text{maximum-likelihood estimates}$$

$$\text{CR}(\alpha) = -\frac{1}{2} \log(\det(X^t W X)) \quad \longleftarrow \quad \text{Cox-Reid bias term}$$

$$\hat{\alpha}_{\text{CR}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu}) + \text{CR}(\alpha)) \quad \longleftarrow \quad \text{Cox-Reid ML estimate}$$

$$\text{prior}(\alpha) = f_{\mathcal{N}}(\log(\alpha); \log(\alpha_{\text{fit}}), \sigma_{\alpha\text{-prior}}^2)$$

$$\hat{\alpha}_{\text{CR-MAP}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu}) + \text{CR}(\alpha) + \log(\text{prior}(\alpha)))$$



# Dispersion estimation

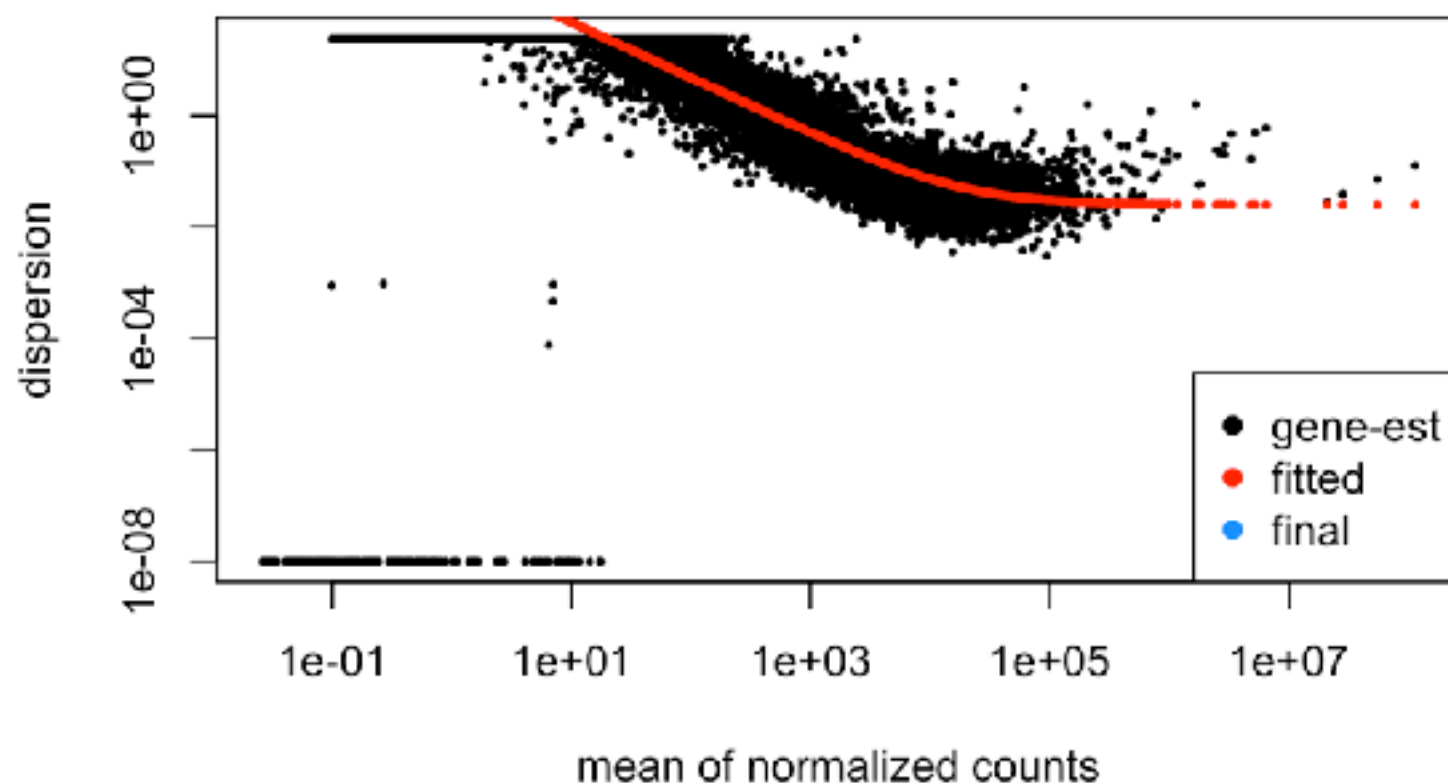
$$\hat{\alpha}_{\text{MLE}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu})) \quad \longleftarrow \quad \text{maximum-likelihood estimates}$$

$$\text{CR}(\alpha) = -\frac{1}{2} \log(\det(X^t W X)) \quad \longleftarrow \quad \text{Cox-Reid bias term}$$

$$\hat{\alpha}_{\text{CR}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu}) + \text{CR}(\alpha)) \quad \longleftarrow \quad \text{Cox-Reid ML estimate}$$

$$\text{prior}(\alpha) = f_{\mathcal{N}}(\log(\alpha); \log(\alpha_{\text{fit}}), \sigma_{\alpha\text{-prior}}^2) \quad \longleftarrow \quad \text{alpha prior by information sharing across genes}$$

$$\hat{\alpha}_{\text{CR-MAP}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu}) + \text{CR}(\alpha) + \log(\text{prior}(\alpha)))$$



# Dispersion estimation

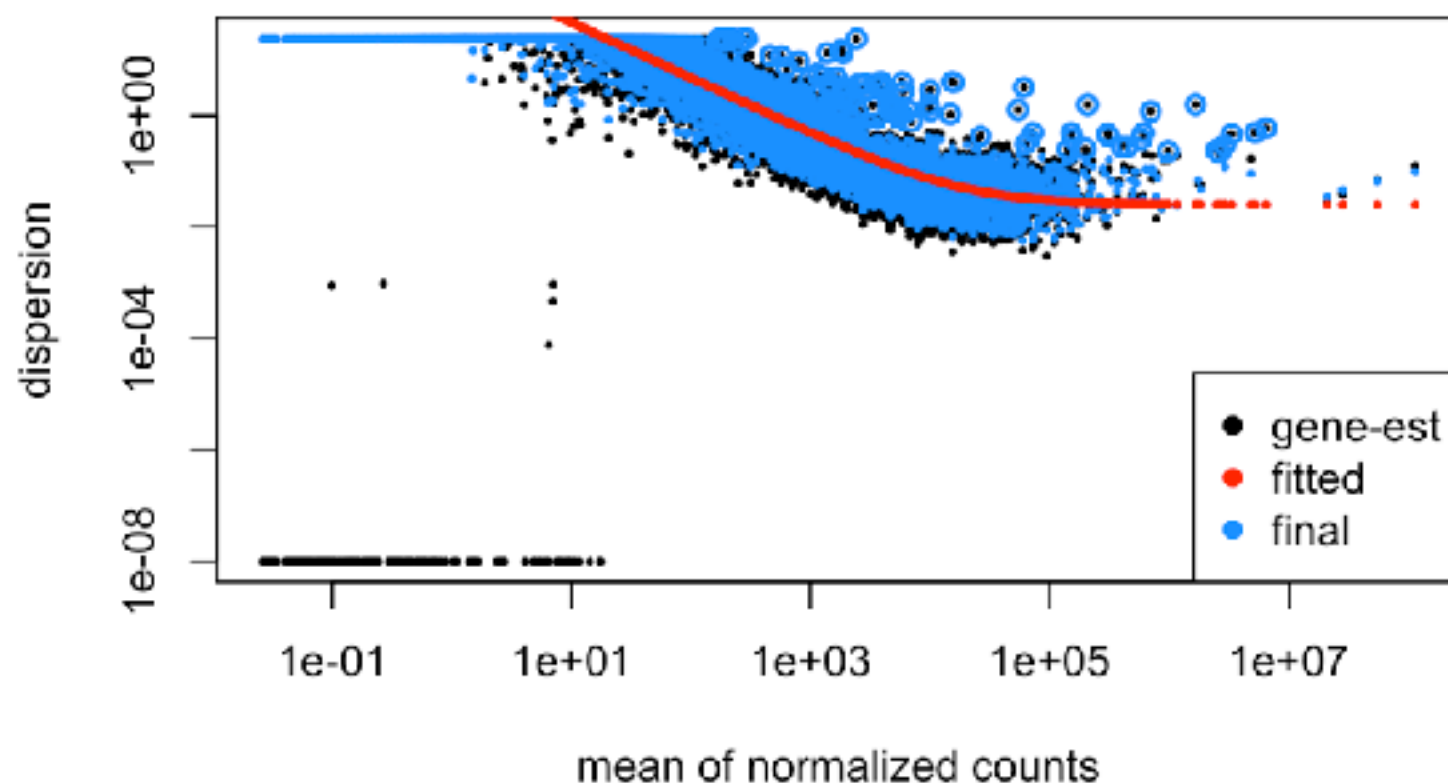
$$\hat{\alpha}_{\text{MLE}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu})) \quad \leftarrow \text{maximum-likelihood estimates}$$

$$\text{CR}(\alpha) = -\frac{1}{2} \log(\det(X^t W X)) \quad \leftarrow \text{Cox-Reid bias term}$$

$$\hat{\alpha}_{\text{CR}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu}) + \text{CR}(\alpha)) \quad \leftarrow \text{Cox-Reid ML estimate}$$

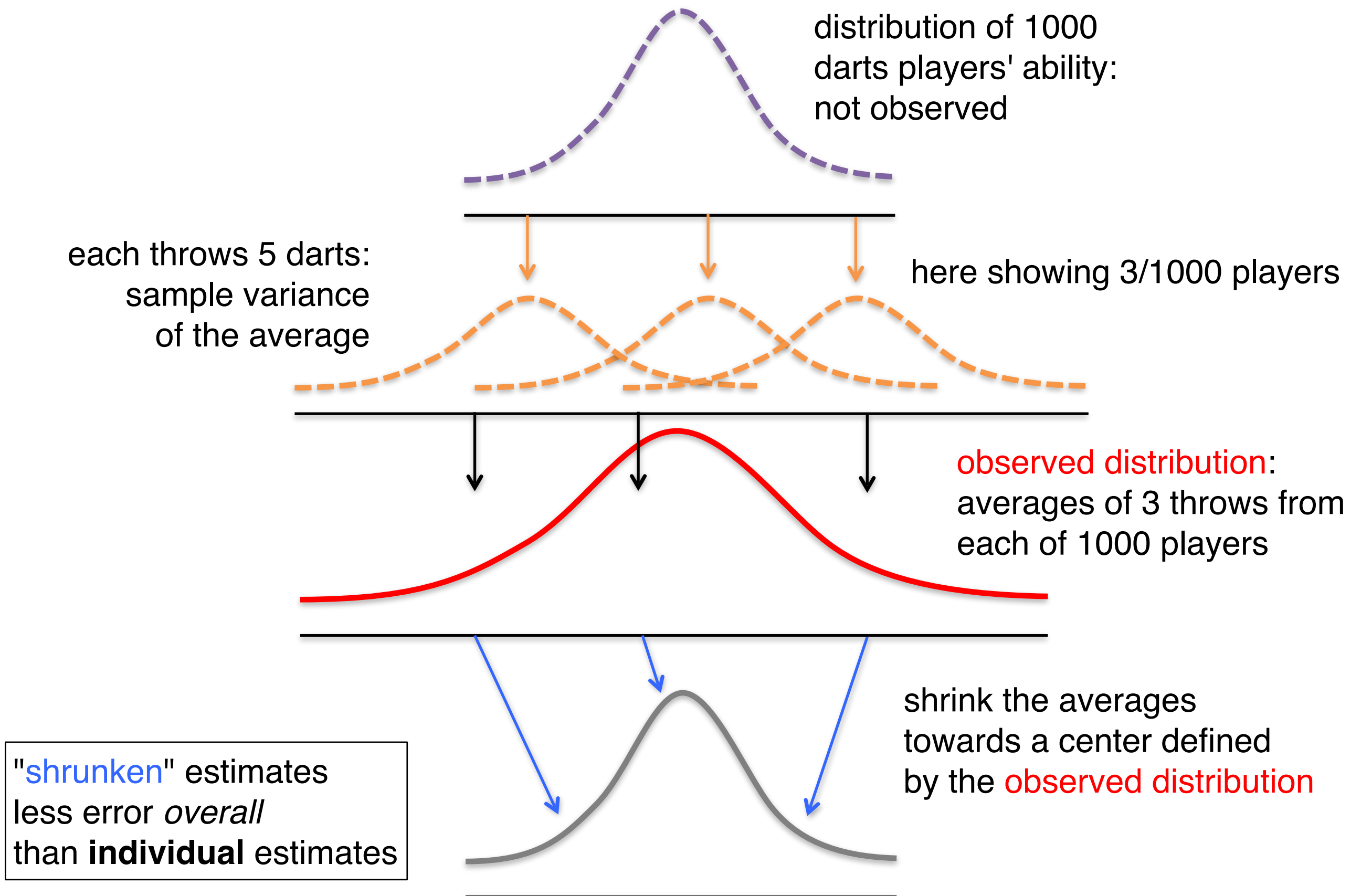
$$\text{prior}(\alpha) = f_{\mathcal{N}}(\log(\alpha); \log(\alpha_{\text{fit}}), \sigma_{\alpha\text{-prior}}^2) \quad \leftarrow \text{alpha prior by information sharing across genes}$$

$$\hat{\alpha}_{\text{CR-MAP}} = \underset{\alpha}{\operatorname{argmax}}(\ell(\alpha|\vec{k}, \hat{\mu}) + \text{CR}(\alpha) + \log(\text{prior}(\alpha)))$$



maximum a posteriori,  
penalized likelihood

# Why shrinkage?



# Shrinkage estimation

population  
distribution

dashed = unobserved

sampling variance  
around true ability

empirical  
distribution

the center defines  
the prior mean

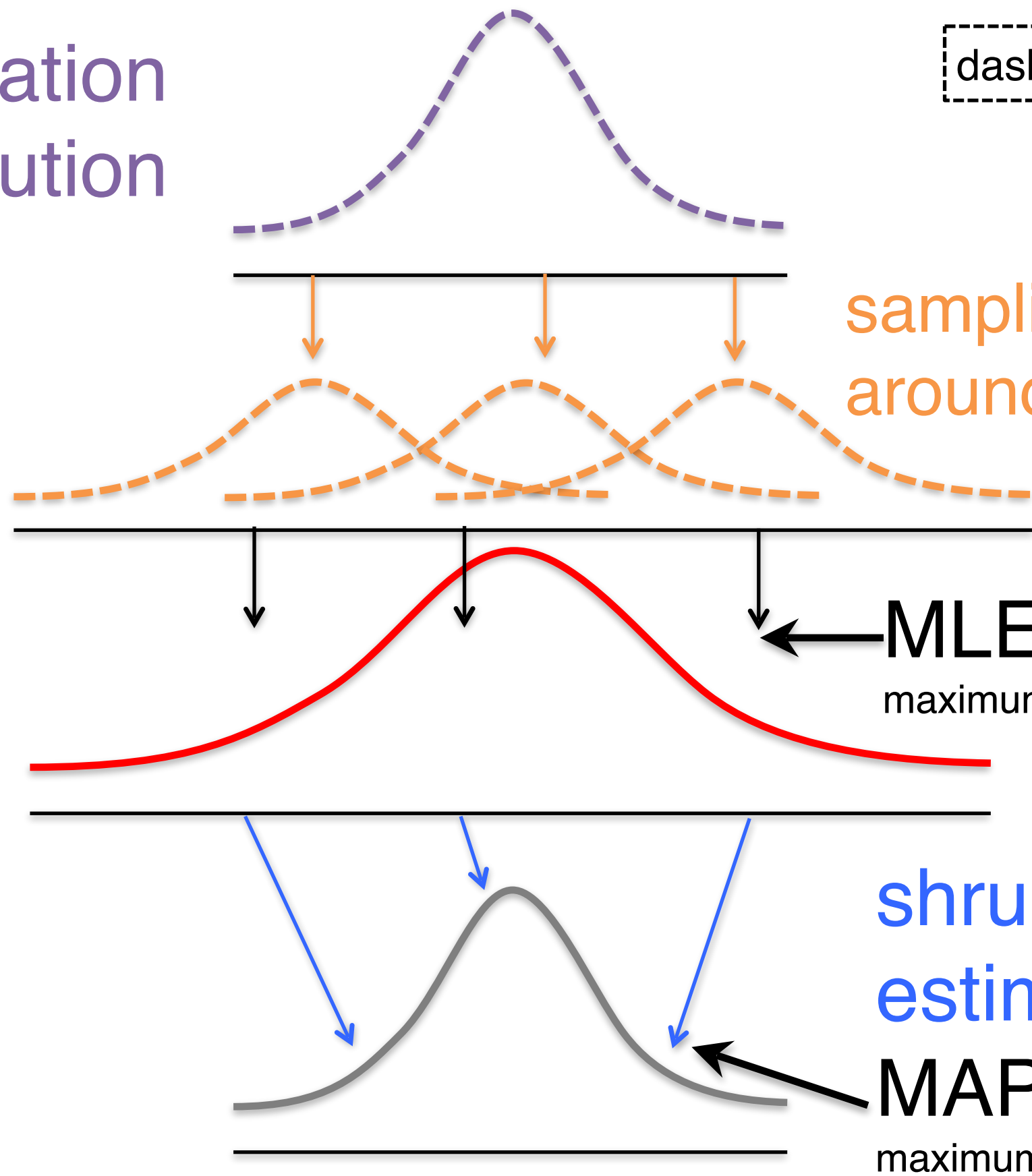
MLE

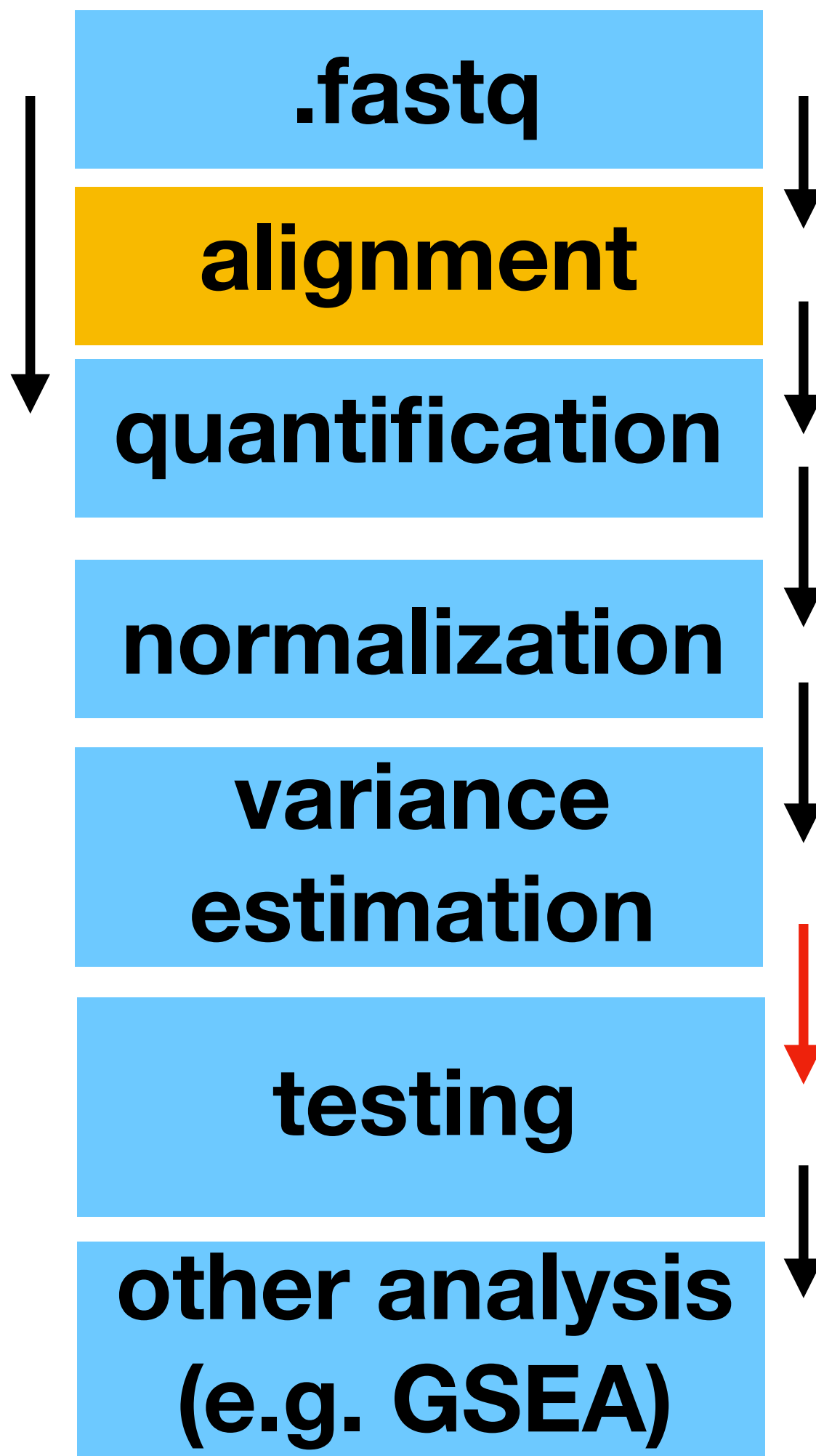
maximum likelihood estimates

shrunk  
estimates or

MAP

maximum a posteriori





# Generalized linear models

$$\log \mu_{ij} = \sum_k \beta_{ik} x_{kj}$$

**response**  
**(expected counts)**



The diagram consists of three arrows pointing upwards towards the equation. The first arrow originates from the text 'response (expected counts)' and points to the term  $\log \mu_{ij}$ . The second arrow originates from the text 'predictors (or differential expression effects)' and points to the coefficient  $\beta_{ik}$ . The third arrow originates from the text 'design matrix' and points to the variable  $x_{kj}$ .

**predictors**  
**(or differential expression effects)**

**design matrix**

# Simplest case: two-group comparison

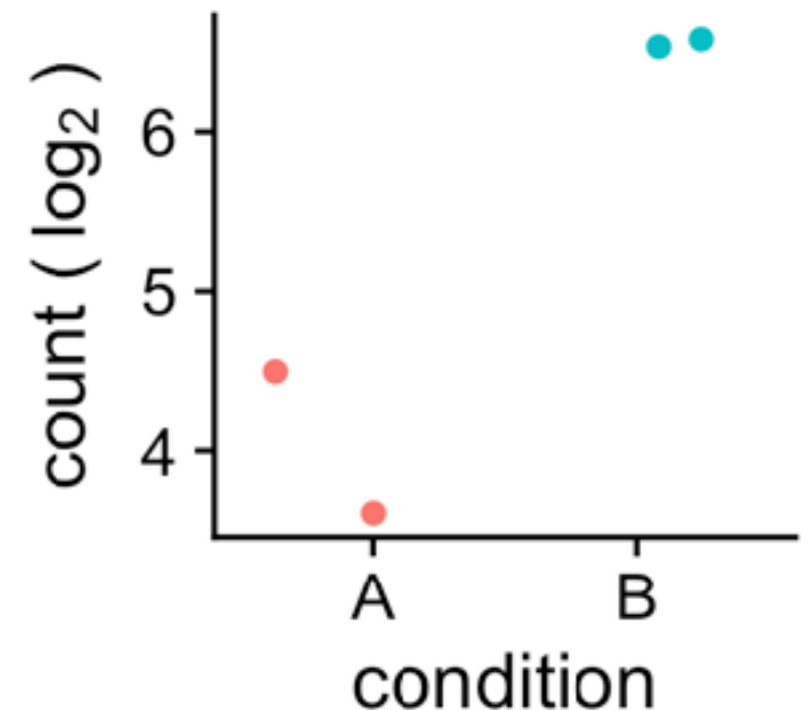
DESeq2 design = ~ condition

$$\begin{array}{cc} \text{condition} & \\ \mathbf{A} & \log_2 u_1 \\ \mathbf{A} & \log_2 u_2 \\ \mathbf{B} & \log_2 u_3 \\ \mathbf{B} & \log_2 u_4 \end{array} = \begin{array}{c} x_1 \\ \left( \begin{array}{cc} 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 1 \end{array} \right) \end{array} \begin{array}{c} \left( \begin{array}{c} \beta_{\text{Intercept}} \\ \beta_{\text{condition\_A\_vs\_B}} \end{array} \right) \end{array}$$

estimated log expression

$$\text{condition A} = \beta_{\text{Intercept}}$$

$$\text{condition B} = \beta_{\text{Intercept}} + \beta_{\text{condition\_A\_vs\_B}}$$





# Multi-level comparisons

DESeq2 design = ~ condition

condition		$x_1$	$x_2$	
<b>A</b>	$\log_2 u_1$	1	0	$\begin{bmatrix} \beta_{\text{Intercept}} \\ \beta_{\text{condition\_A\_vs\_B}} \\ \beta_{\text{condition\_A\_vs\_C}} \end{bmatrix}$
<b>A</b>	$\log_2 u_2$	1	0	
<b>B</b>	$\log_2 u_3$	1	1	
<b>B</b>	$\log_2 u_4$	1	1	
<b>C</b>	$\log_2 u_5$	1	0	
<b>C</b>	$\log_2 u_6$	1	0	

=

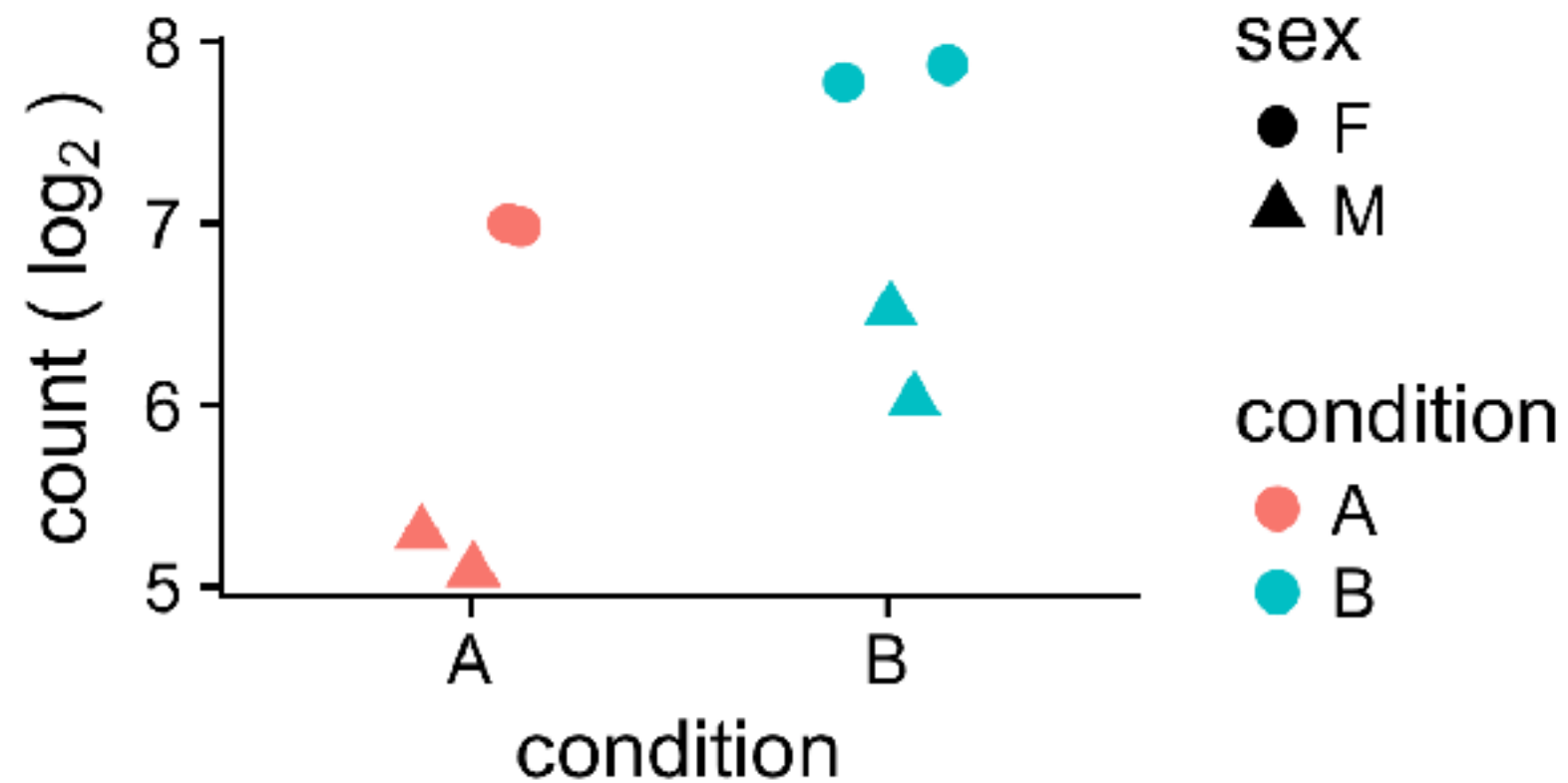
estimated log expression

**condition A** =  $\beta_{\text{Intercept}}$

**condition B** =  $\beta_{\text{Intercept}}$  +  $\beta_{\text{condition\_A\_vs\_B}}$

**condition C** =  $\beta_{\text{Intercept}}$  +  $\beta_{\text{condition\_A\_vs\_C}}$

# Comparisons with blocking factors



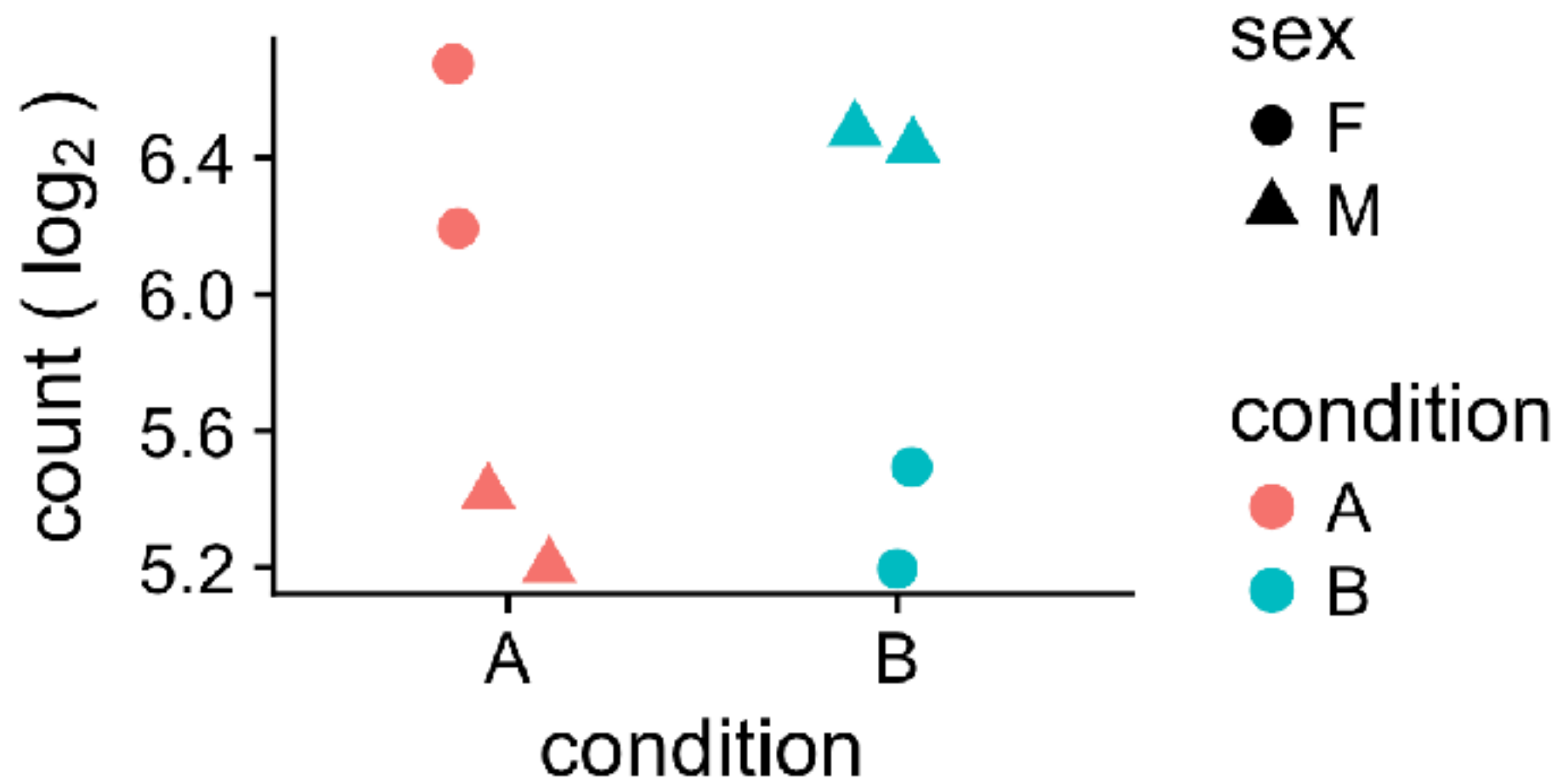
# Comparisons with blocking factors

DESeq2 design = ~ sex + condition

sex	condition				
M	A	$\log_2 u_1$	=	$\begin{pmatrix} & x_1 & x_2 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{pmatrix}$	$\begin{pmatrix} \beta_{\text{Intercept}} \\ \beta_{\text{sex\_female\_vs\_male}} \\ \beta_{\text{condition\_A\_vs\_B}} \end{pmatrix}$
M	A	$\log_2 u_2$			
M	B	$\log_2 u_3$			
M	B	$\log_2 u_4$			
F	A	$\log_2 u_5$			
F	A	$\log_2 u_6$			
F	B	$\log_2 u_7$			
F	B	$\log_2 u_8$			

What would be the predicted log expression for females in condition B?

# Interactions



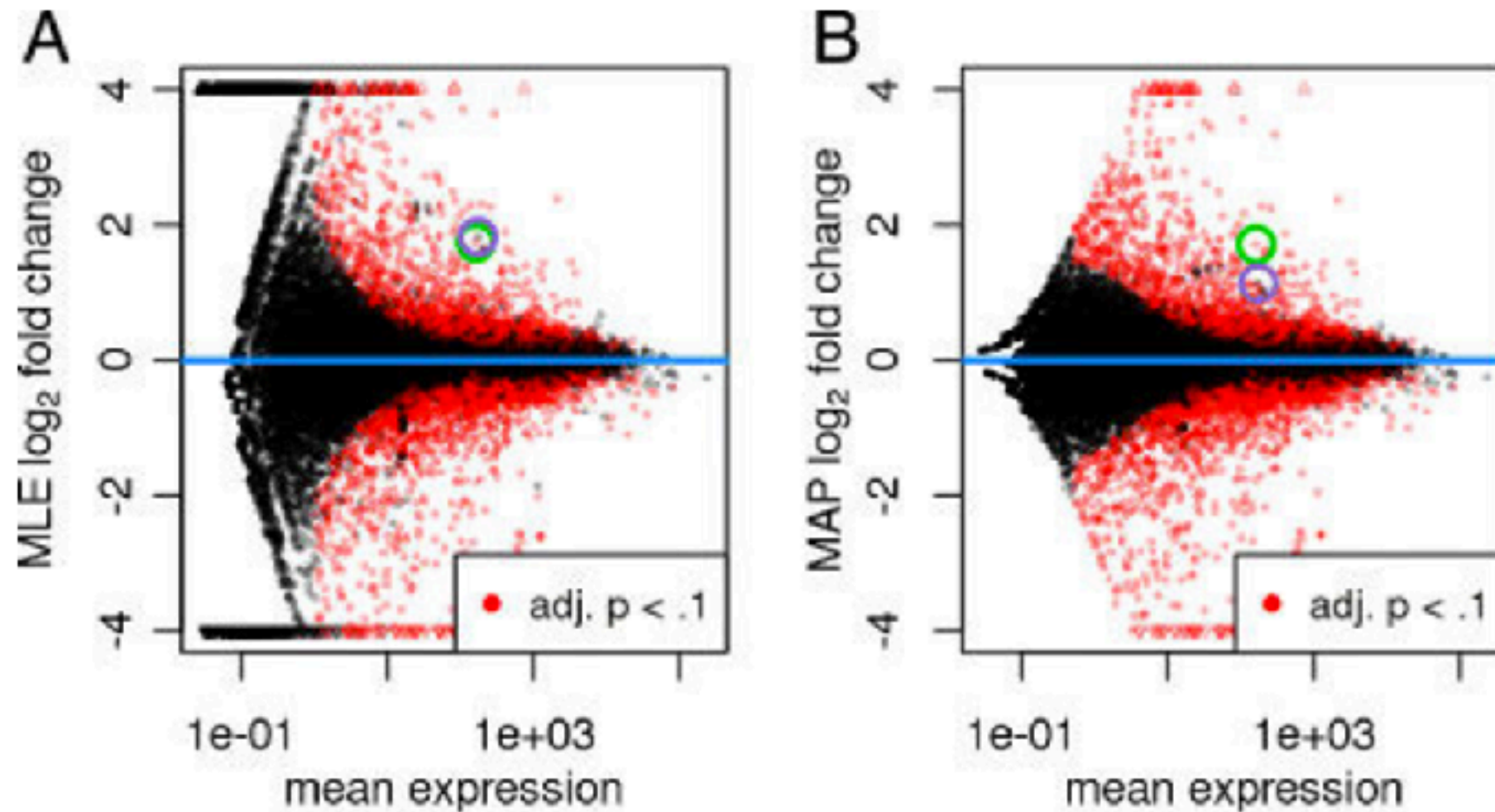
# Comparisons with blocking factors

DESeq2 design = ~ sex + condition + sex:condition

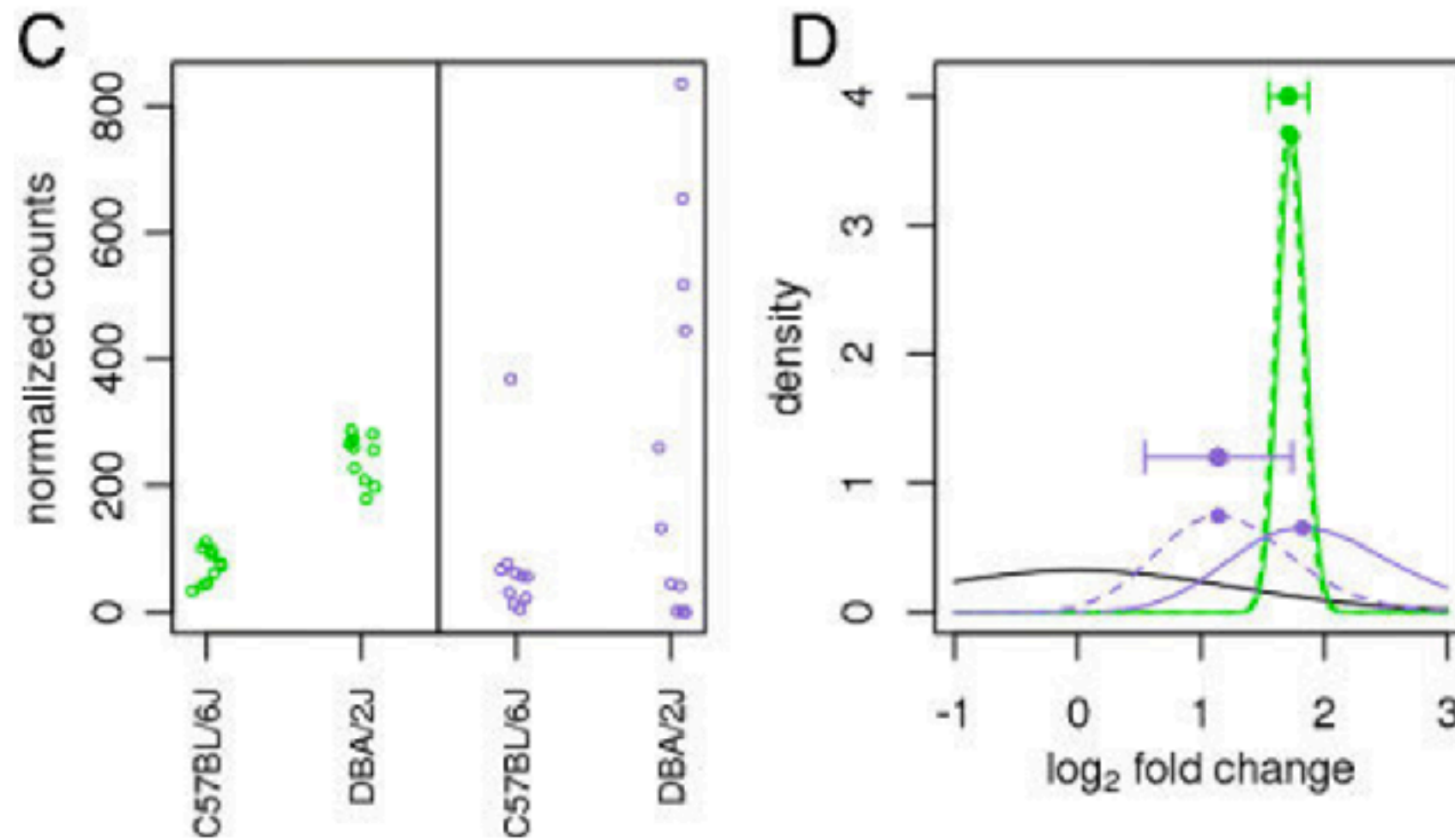
sex	condition	
M	A	$\log_2 u_1$
M	A	$\log_2 u_2$
M	B	$\log_2 u_3$
M	B	$\log_2 u_4$
F	A	$\log_2 u_5$
F	A	$\log_2 u_6$
F	B	$\log_2 u_7$
F	B	$\log_2 u_8$

$$= \begin{pmatrix} & X_1 & X_2 & X_3 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} \beta_{\text{Intercept}} \\ \beta_{\text{sex\_female\_vs\_male}} \\ \beta_{\text{condition\_A\_vs\_B}} \\ \beta_{\text{interaction}} \end{pmatrix}$$

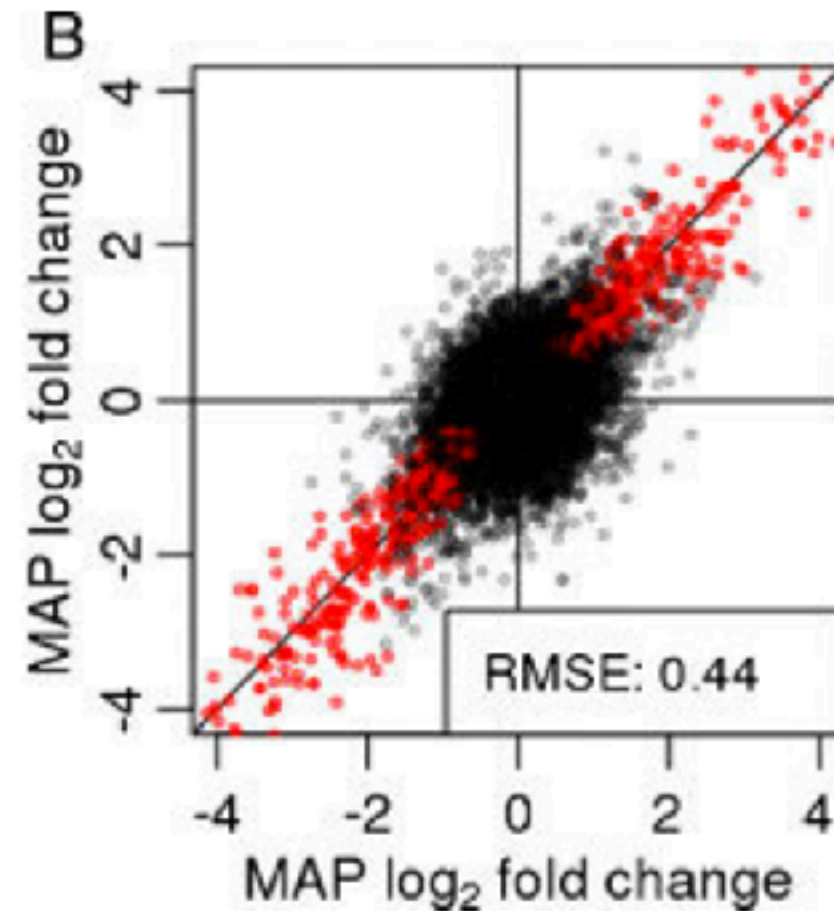
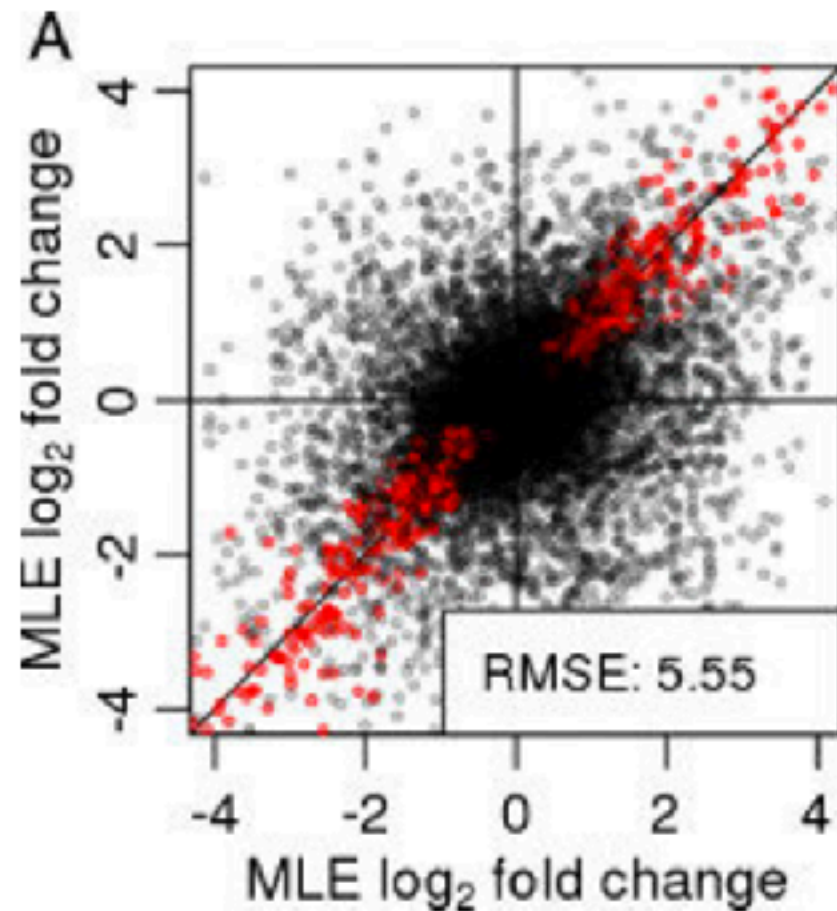
# DESeq2 shrinkage of log fold changes



# DESeq2 shrinkage of log fold changes

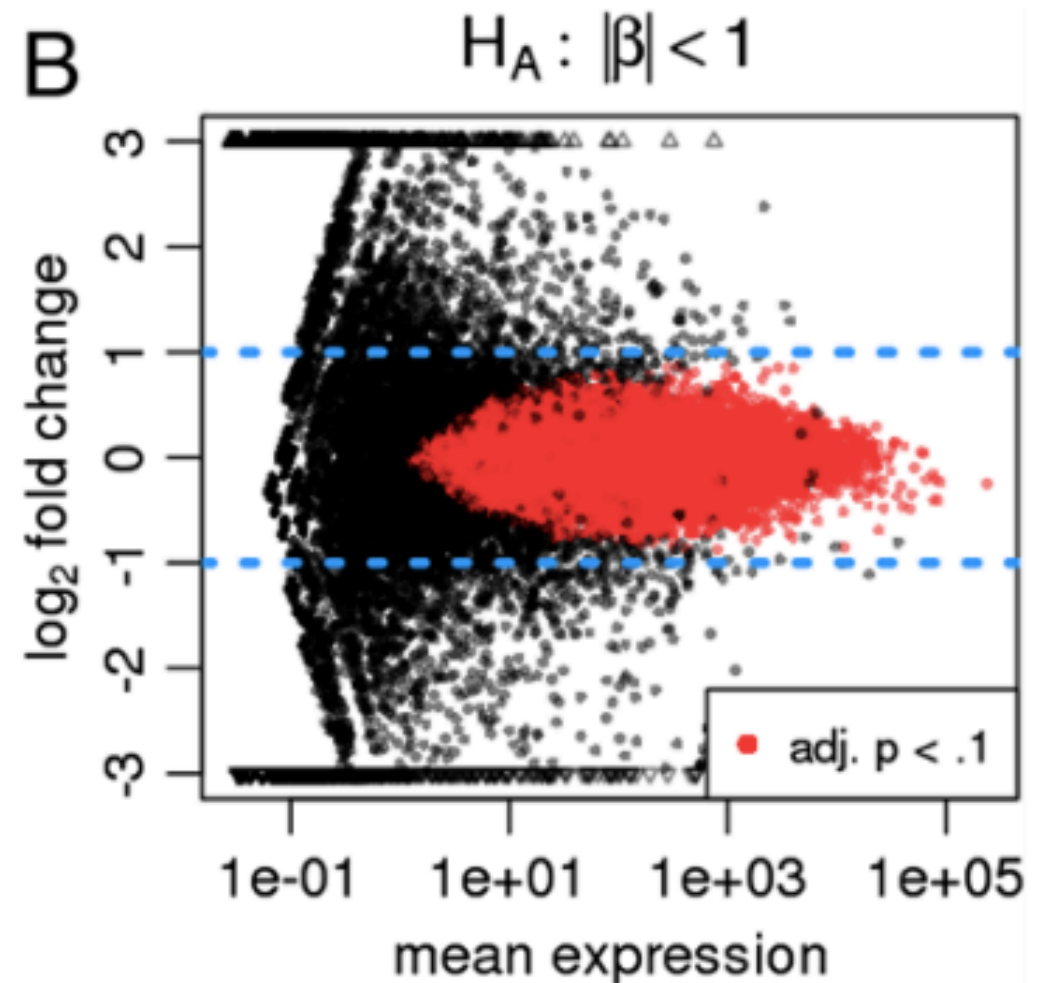
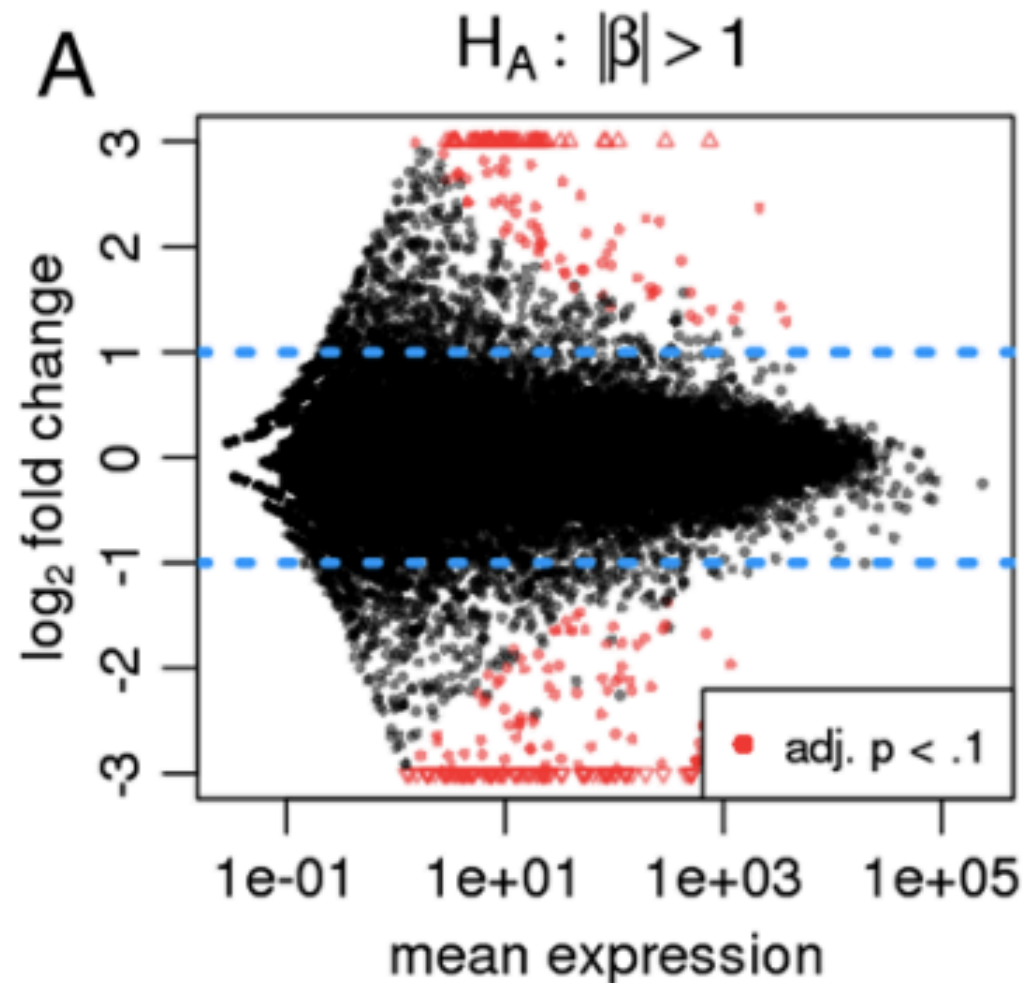


# DESeq2 shrinkage of log fold changes





# Other hypothesis testing options



# Confounders and experimental design

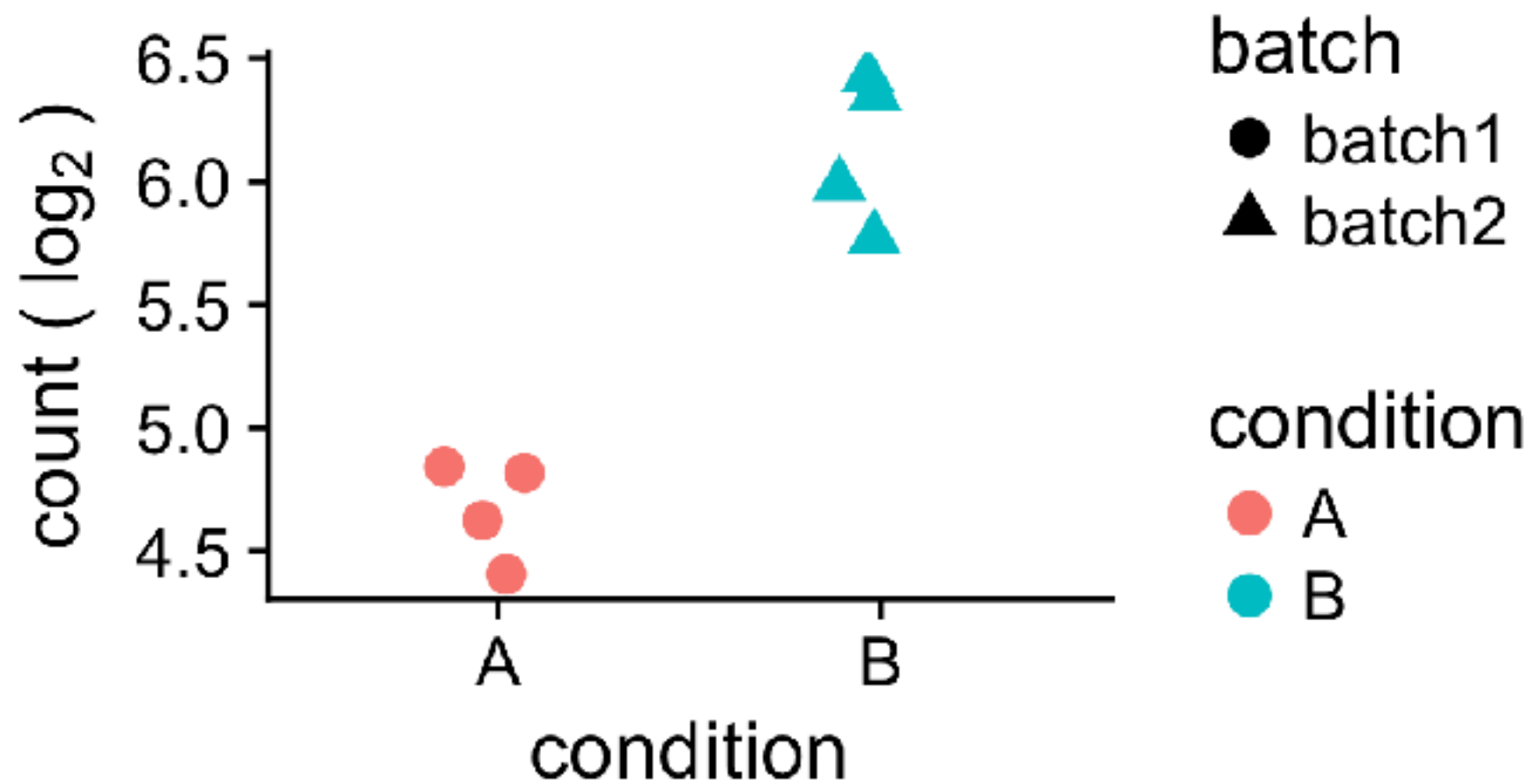
- Spielman et al., Nature Genetics 2007

78% of genes are differentially expressed between human populations

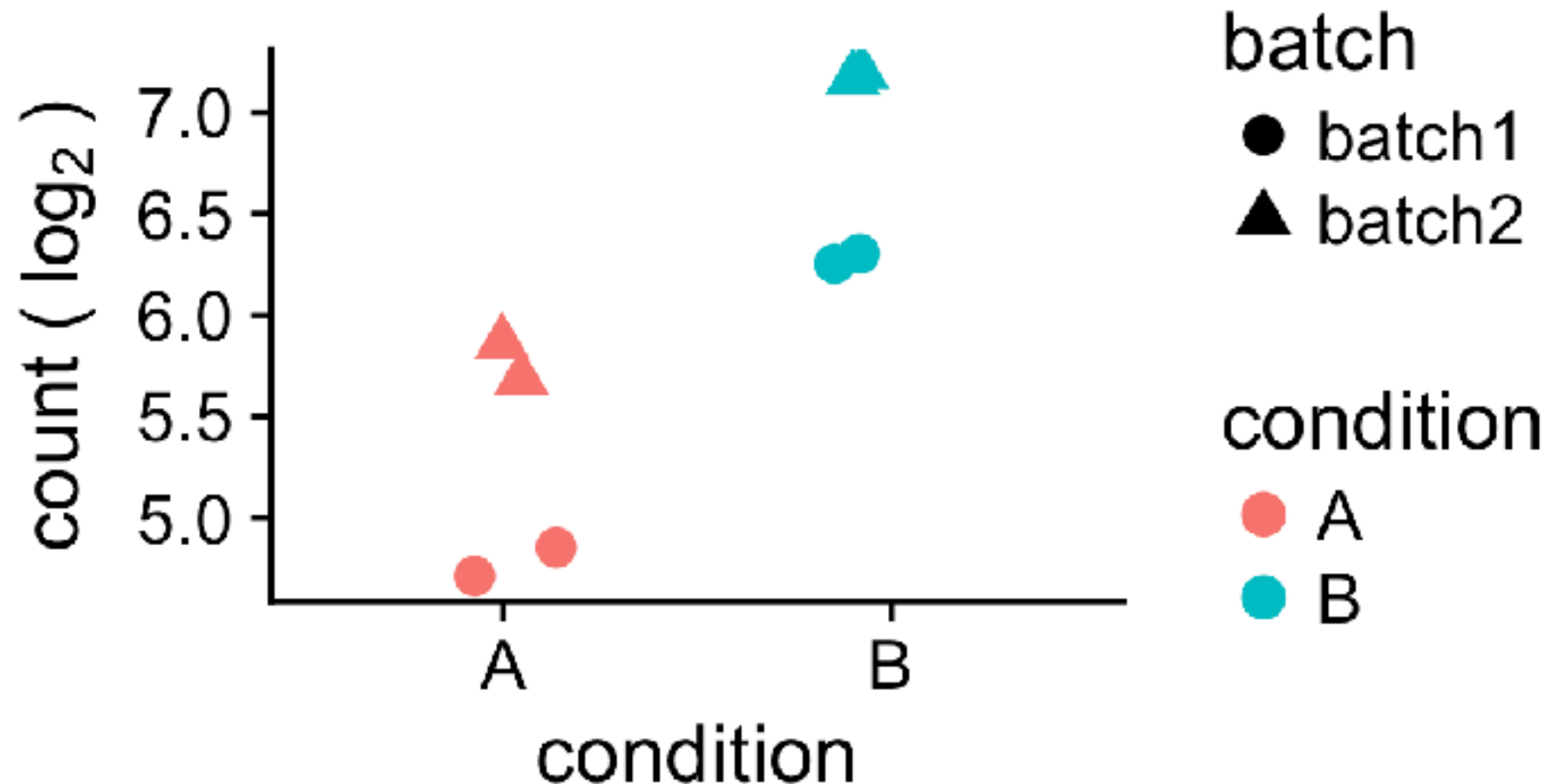
- Lin et al., PNAS, 2014

Differences in gene expression across species are larger than between tissues of a same species.

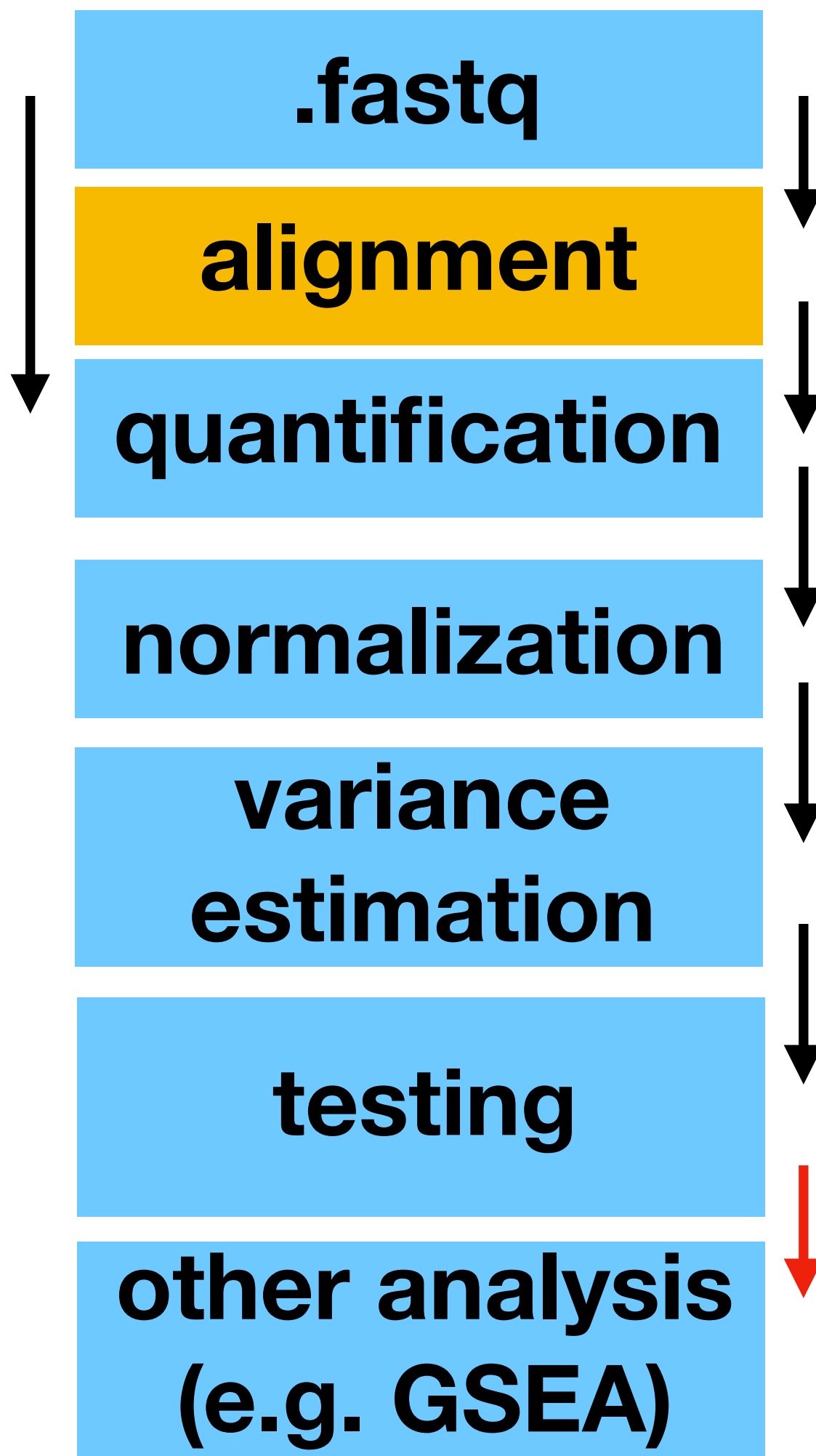
# Confounders experiment



# Confounders and experimental design



- Randomize conditions of interest in batches and include it as a blocking factor.



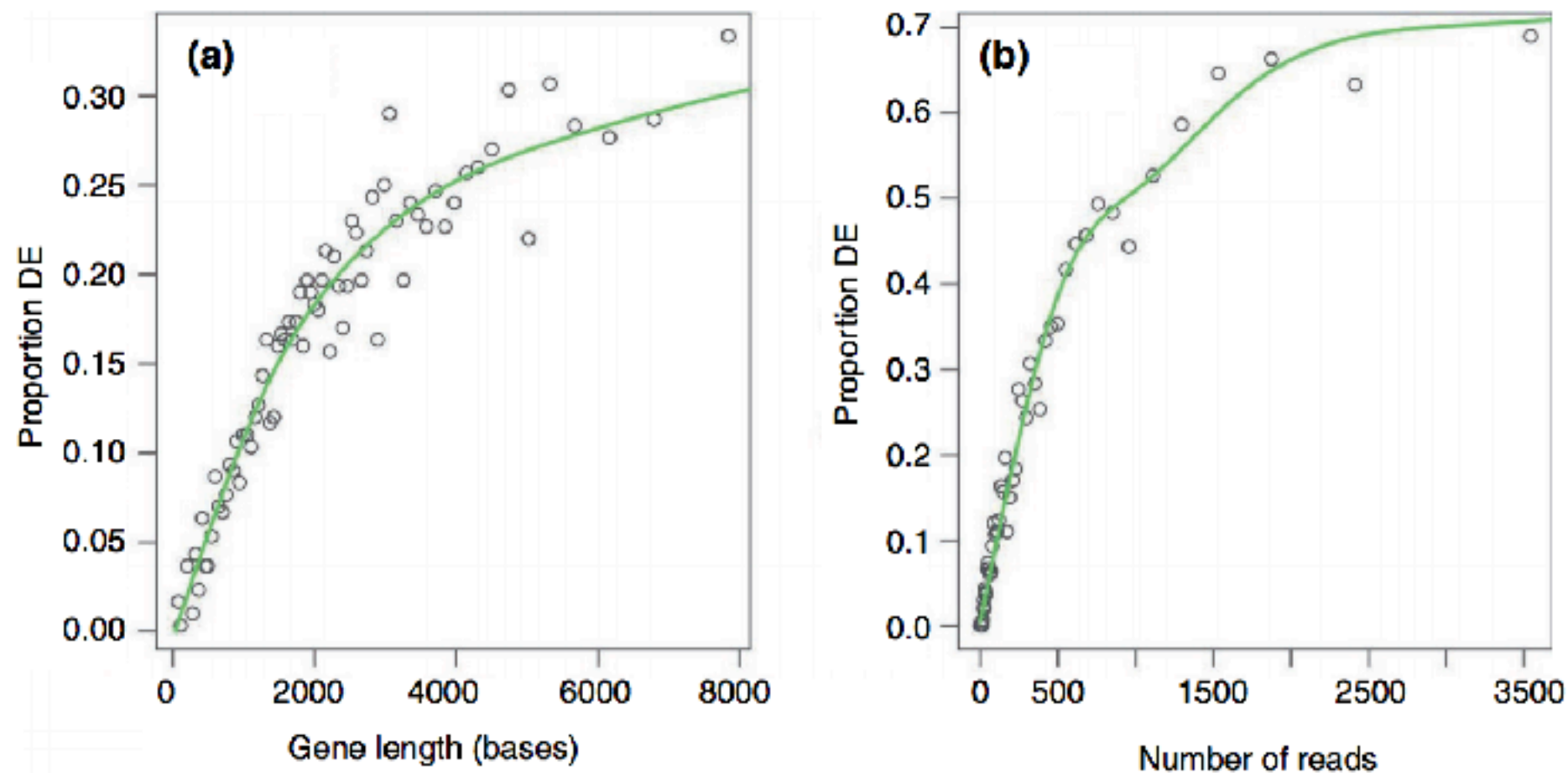
# “Traditional” gene-set enrichment analysis

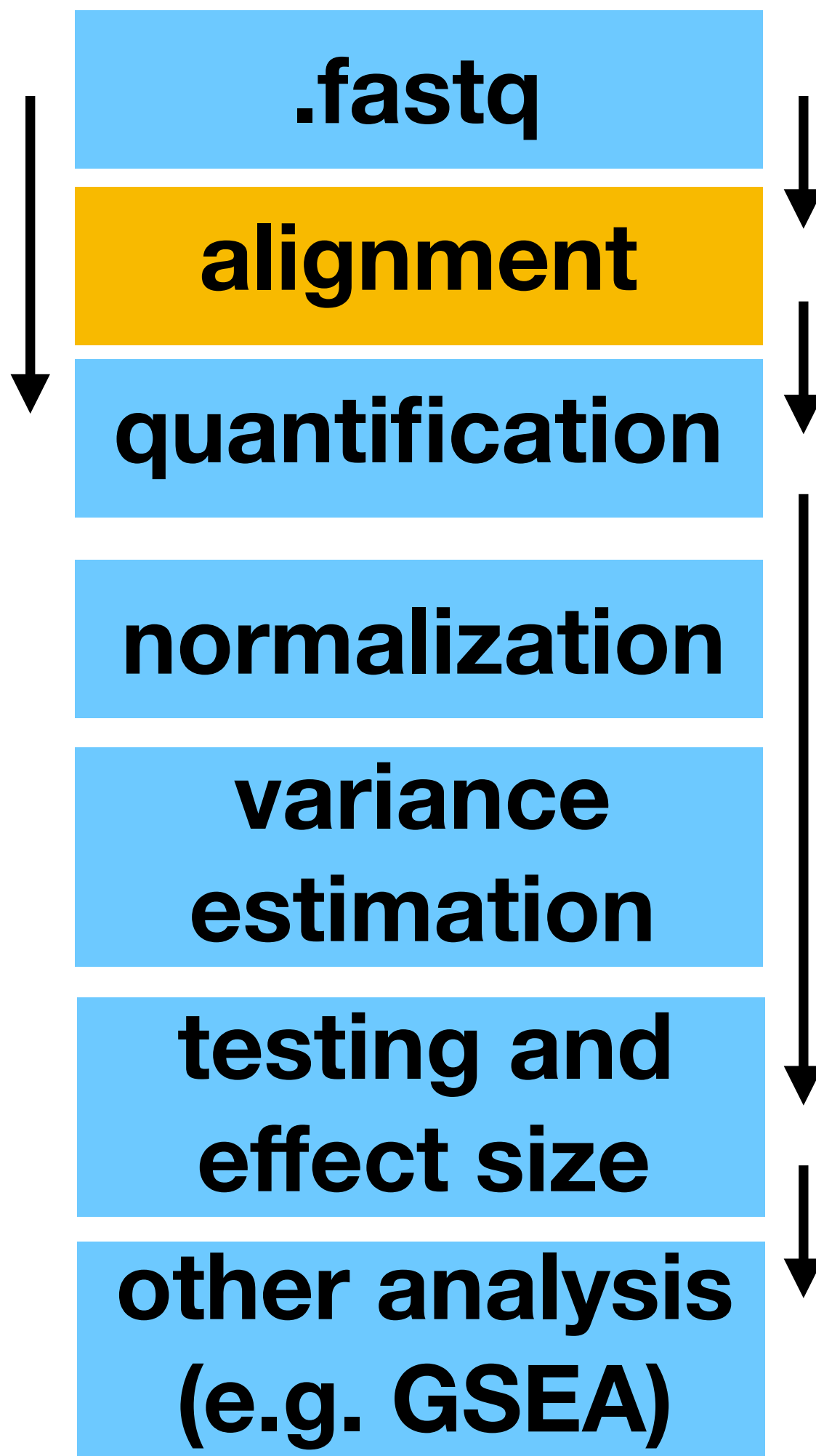
	In gene set	Not in gene set	
DGE	a	b	a + b
Background	c	d	c + d
	<b>a + c</b>	<b>b + d</b>	

$$p = \frac{\binom{a+b}{a} \binom{c+d}{c}}{\binom{n}{a+c}} = \frac{(a+b)! (c+d)! (a+c)! (b+d)!}{a! b! c! d! n!}$$

**Fisher’s test**

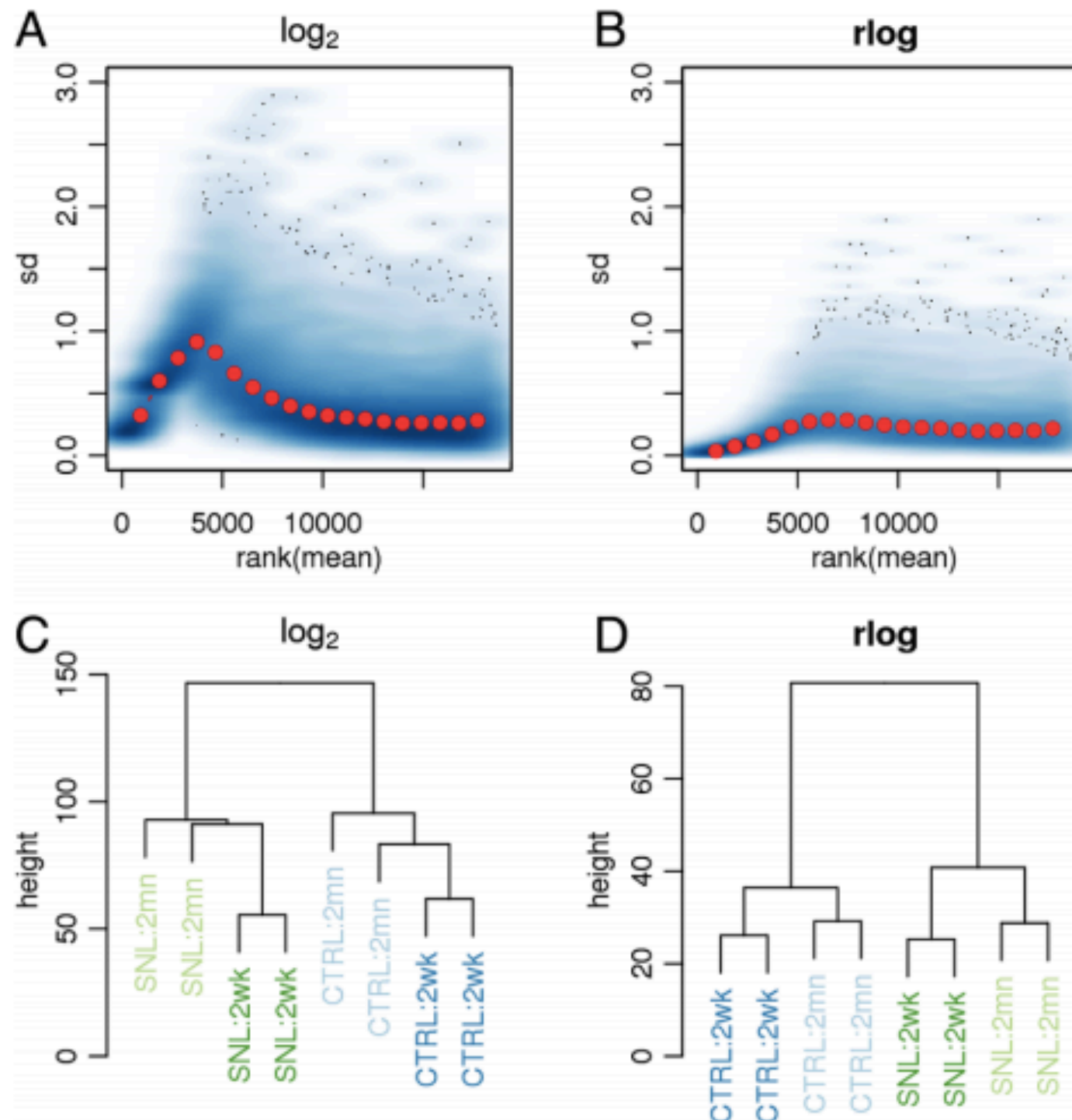
# goseq estimates and corrects for RNA-seq biases in GSEA



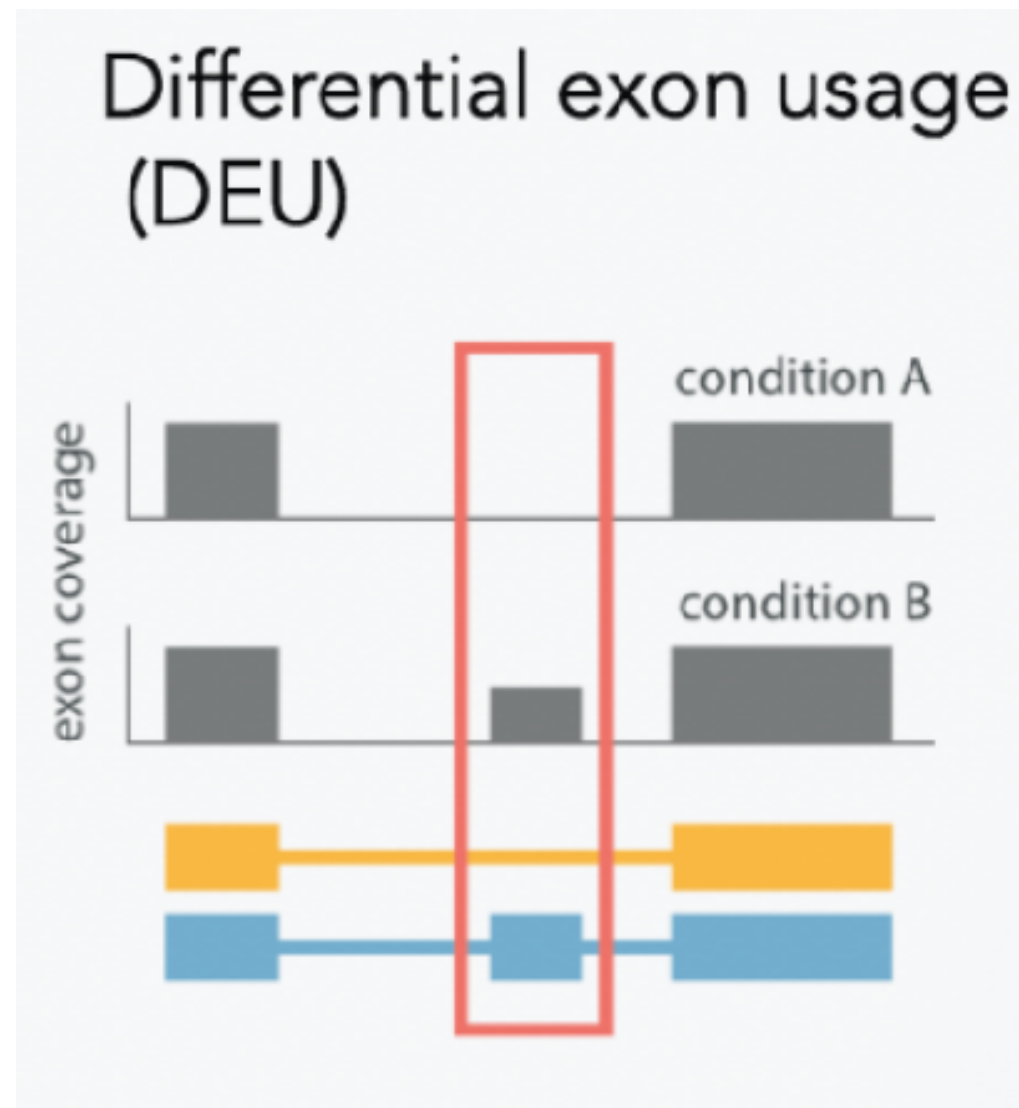




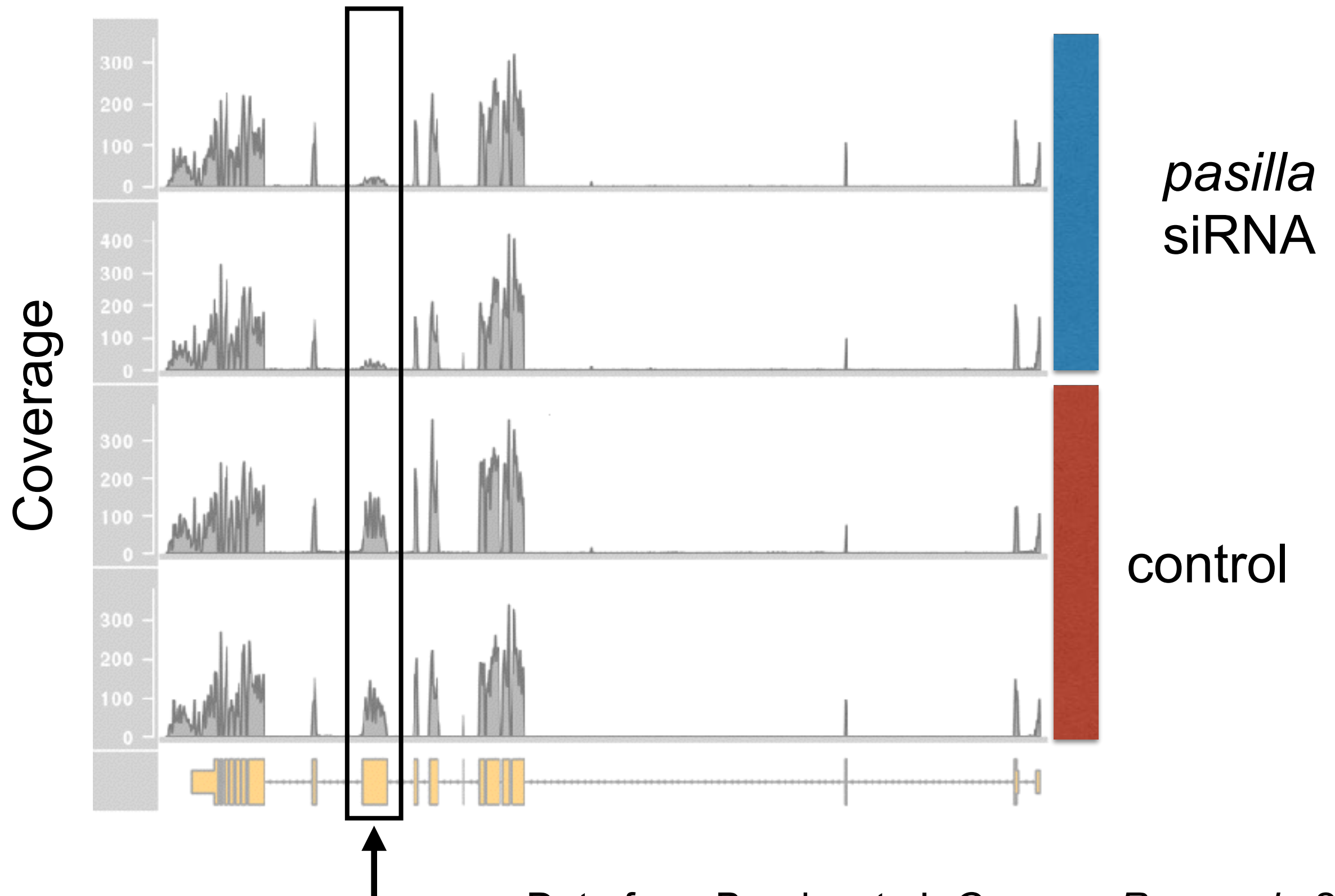
# For exploratory analysis, clustering, PCA: Use rlog or vst data!



# Differential exon usage



# High-throughput RNA sequencing enables an unbiased characterisation of isoform expression



Data from Brooks et al. *Genome Research*, 2011

# DEXSeq: inference of differential exon usage

samples →

exons ↓

	treated1fb	treated2fb	treated3fb	untreated1fb	untreated2fb	untreated3fb	untreated4fb
E001	1997	494	562	1150	2514	570	547
E002	122	112	180	69	203	156	142
E003	276	293	305	190	398	312	259
E004	420	200	182	230	446	183	185
E005	416	217	279	146	170	237	231
E006	486	357	471	190	337	418	364
E007	574	465	536	469	805	480	496
E008	536	417	447	541	832	475	472
E009	191	237	216	217	427	286	222
E010	188	130	96	617	1177	520	508
E011	165	212	210	118	275	294	269
E012	536	437	414	441	792	619	504
E013	72	41	49	40	76	34	38
E014	3	0	33	5	0	2	42

$$\text{exon usage} = \frac{\text{\# of transcripts including an exon}}{\text{\# total transcripts}}$$

# DEXSeq: inference of differential exon usage

exon  $i$ , sample  $j$

$l = 1$  exon under consideration

$l = 0$  sum of gene count

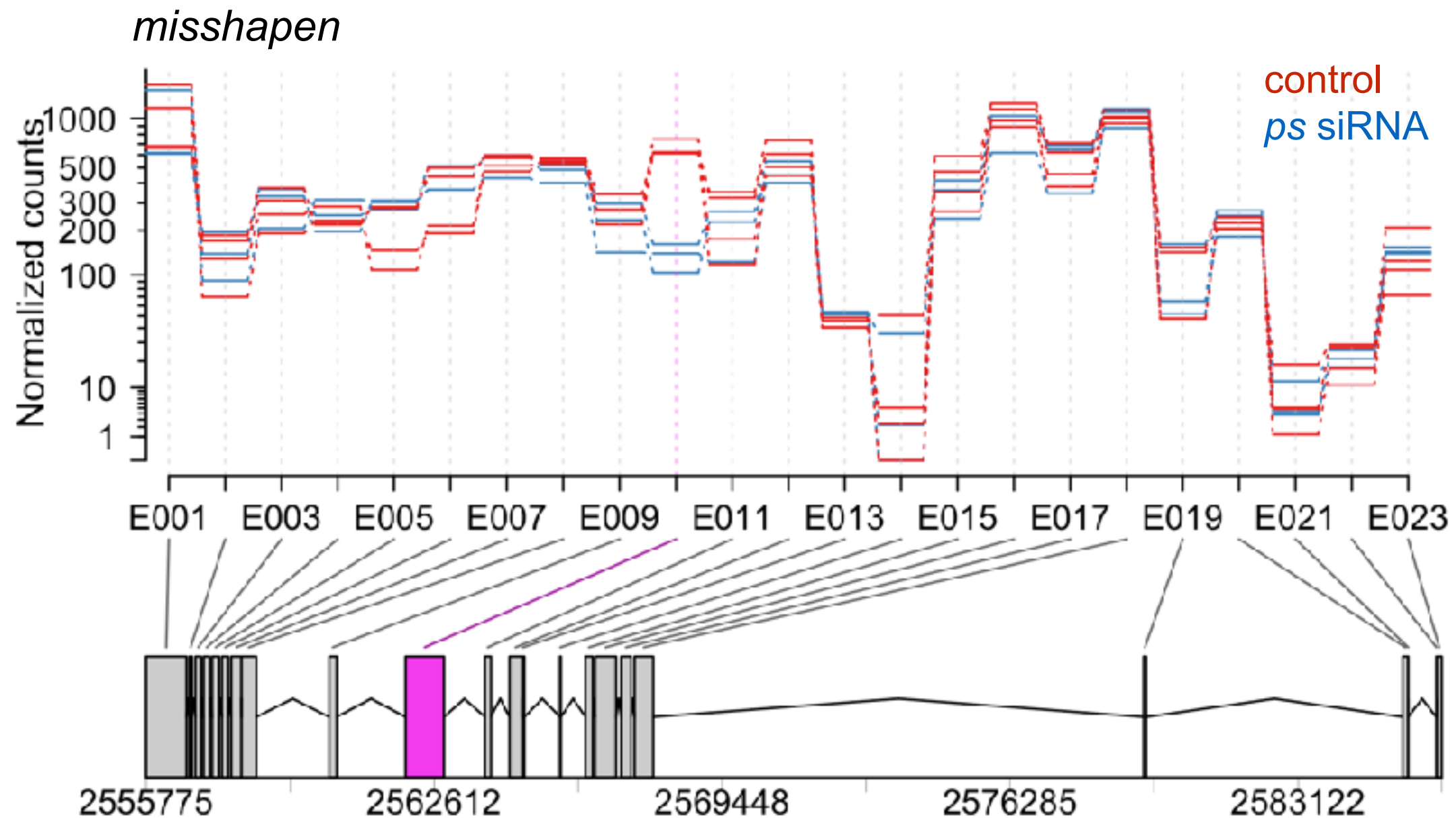
$$\log \mu_{ijl} = \beta_{ij}^S + l\beta_i^E + \beta_{i\rho_j}^{EC}$$

sample specific  
contributions  
(gene expression)

average  
exon usage

changes  
in exon usage  
due to the  
conditions

# DEXSeq: inference of differential exon usage



# Detecting differential splicing vs differential usage of TSS and polyA

