



# Transcriptomic Data Analysis

Jędrzej Jakub Szymański

IPK Gatersleben & FZJ Jülich



# Transcriptomic Data Analysis

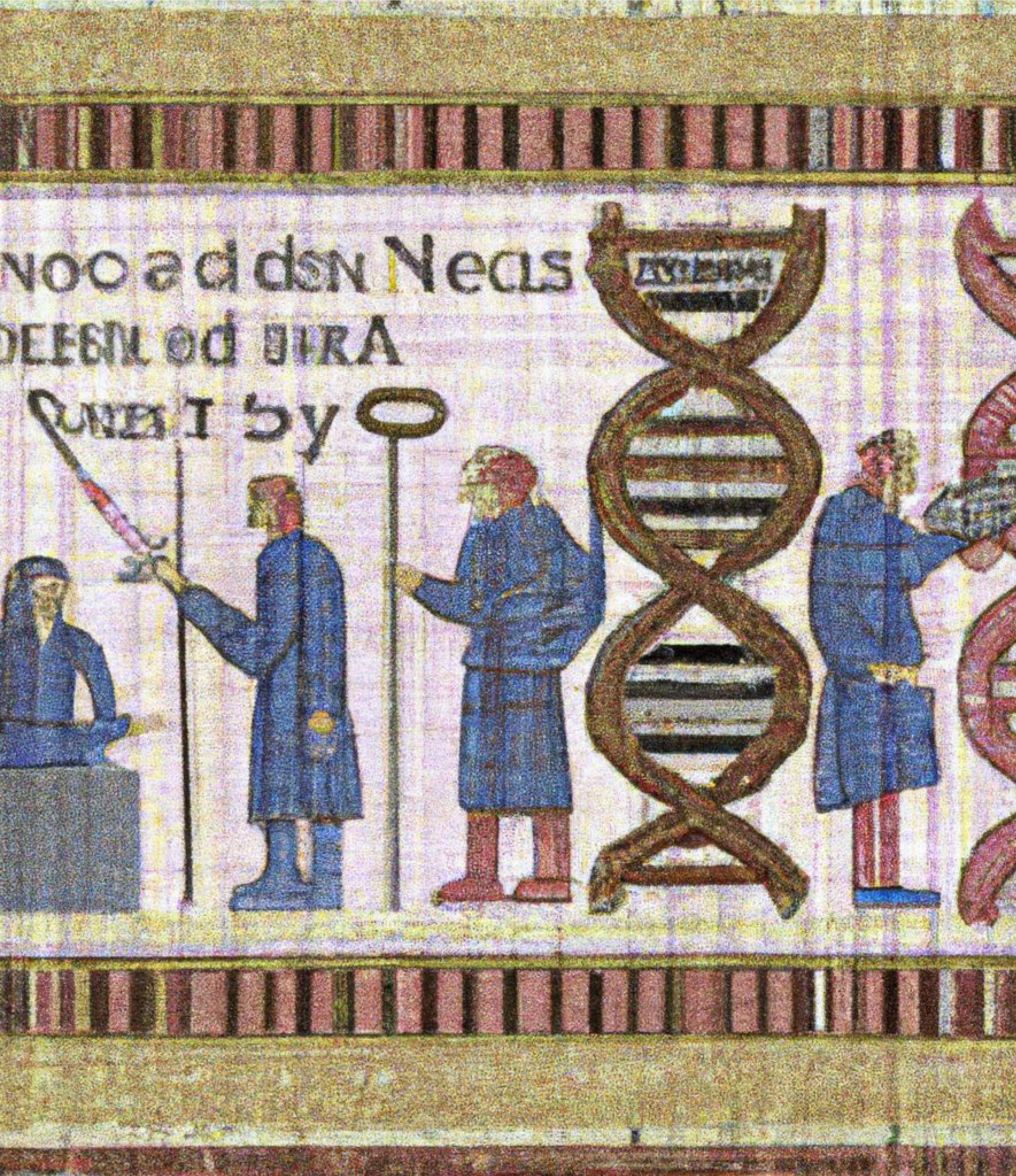
Jędrzej Jakub Szymański

IPK Gatersleben & FZJ Jülich

“Transcriptomic data analysis in the middle ages” DALL-E2 generative AI model

# Lectures in the series

1. Introduction to Transcriptomic Data
2. Study Design, Annotation, and Data Import
3. Filtering and Normalization of Transcriptomic Data
4. Data Analysis - Multivariate Statistics and Visualization
5. Accessing Public Data
6. Differential Gene Expression
7. Functional Enrichment Analysis



# Introduction to Transcriptomic Data

What we **can** learn from transcriptomic data?

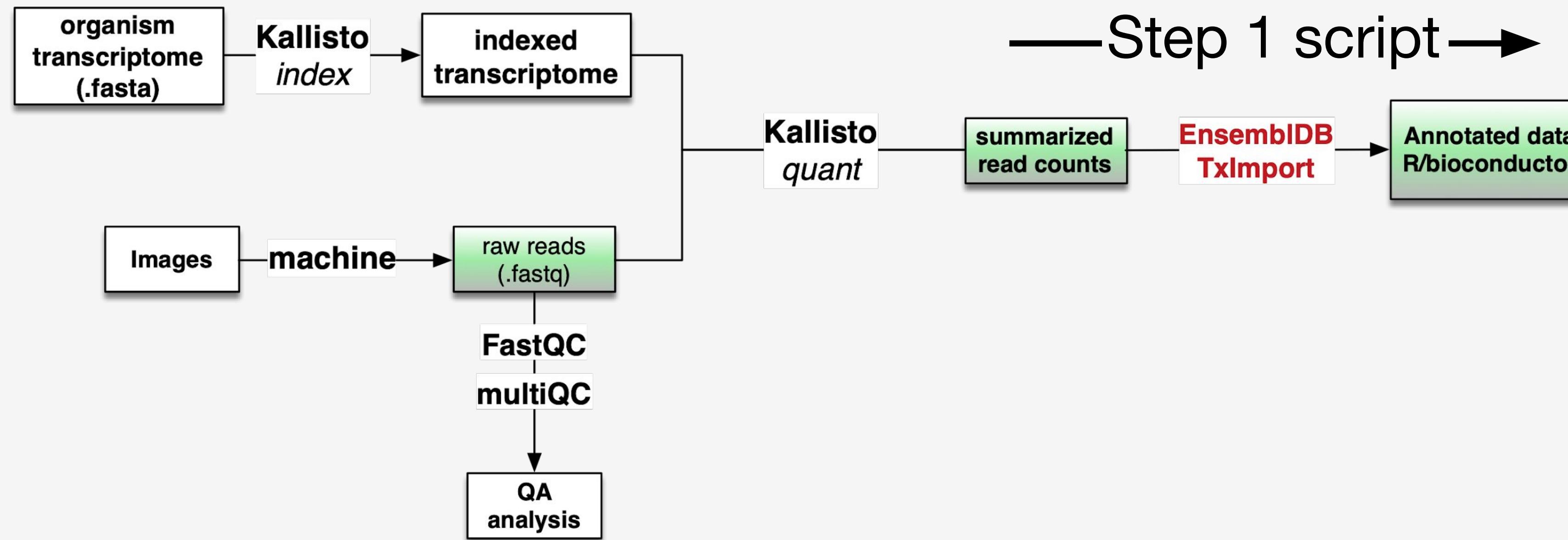
Introduction to data formats

Basic definitions - read, read count, transcript length, RPKM, TPM

Start an RStudio Project directory that we'll use for the rest of the course.

Open and discuss our first script, including installation of packages

# Tracking our workflow in this course



—Step 1 script→

*Don't worry if you've struggled with QA or read mapping!*

# Explore other Kallisto functions on your own

```
(base) daniels-mbp-3:course_dataset danielbeiting$ kallisto
kallisto 0.46.2

Usage: kallisto <CMD> [arguments] ..

Where <CMD> can be one of:

index          Builds a kallisto index
quant         Runs the quantification algorithm
bus            Generate BUS files for single-cell data
pseudo         Runs the pseudoalignment step
merge          Merges several batch runs
h5dump         Converts HDF5-formatted results to plaintext
inspect        Inspects and gives information about an index
version        Prints version information
cite           Prints citation information
```

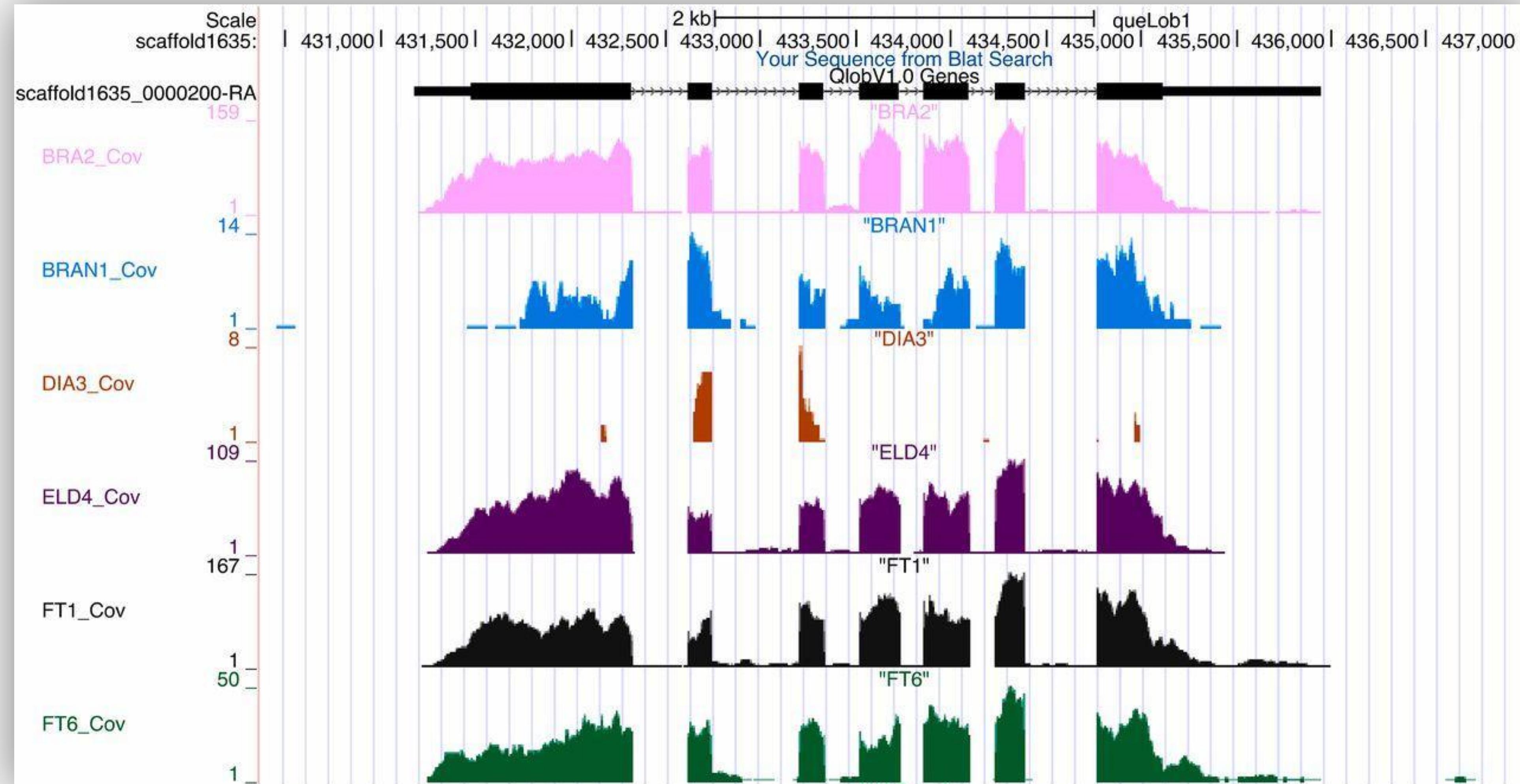
```
kallisto 0.46.0
Computes equivalence classes for reads and quantifies abundances

Usage: kallisto quant [arguments] FASTQ-files

Required arguments:
-i, --index=STRING      Filename for the kallisto index to be used for quantification
-o, --output-dir=STRING Directory to write output to

Optional arguments:
--bias                  Perform sequence based bias correction
-b, --bootstrap-samples=INT Number of bootstrap samples (default: 0)
--seed=INT               Seed for the bootstrap sampling (default: 42)
--plaintext             Output plaintext instead of HDF5
--fusion                Search for fusions for Pizzly
--single                Quantify single-end reads
--single-overhang       Include reads where unobserved rest of fragment is predicted to lie outside a transcript
--fr-stranded           Strand specific reads, first read forward
--rf-stranded           Strand specific reads, first read reverse
-l, --fragment-length=DOUBLE Estimated average fragment length
-s, --sd=DOUBLE          Estimated standard deviation of fragment length (default: -l, -s values are estimated from paired end data, but are required when using --single)
-t, --threads=INT        Number of threads to use (default: 1)
--pseudobam              Save pseudoalignments to transcriptome to BAM file
--genomebam              Project pseudoalignments to genome sorted BAM file
-g, --gtf                 GTF file for transcriptome information (required for --genomebam)
-c, --chromosomes        Tab separated file with chromosome names and lengths (optional for --genomebam, but recommended)
```

# Using Kallisto (pseudoalignment) is a flexible strategy for handling many types of seq experiments



# Pseudoalignment also performs well for lowly expressed transcripts (e.g. lncRNAs)

## Benchmark of long non-coding RNA quantification for RNA sequencing of cancer samples

Hong Zheng  <sup>1</sup>, Kevin Brennan <sup>1</sup>, Mikel Hernaez  <sup>2</sup> and Olivier Gevaert  <sup>1,3,\*</sup>

<sup>1</sup>Stanford Center for Biomedical Informatics Research, Department of Medicine, Stanford University, 1265 University of Illinois at Biomedical Data Science,

“In this benchmarking study, we compared the performance of pseudoalignment methods Kallisto and Salmon, and alignment based methods HTSeq, featureCounts, and RSEM, in lncRNA quantification, by applying them to a simulated RNA-Seq dataset and a pan-cancer RNA-Seq dataset from TCGA.”

“In summary, pseudoalignment methods Kallisto or Salmon in combination with the full transcriptome annotation is our recommended strategy for RNA-Seq analysis for lncRNAs.”

**protocols.hostmicrobe.org/compute/kallisto**

# Kallisto ‘bus tools’ for scRNAseq

```
(base) daniels-mbp-3:course_dataset danielbeiting$ kallisto  
kallisto 0.46.2
```

Usage: kallisto <CMD> [arguments] ..

Where <CMD> can be one of:

index	Builds a kallisto index
quant	Runs the quantification algorithm
bus	Generate BUS files for single-cell data
pseudo	Runs the pseudoalignment step
merge	Merges several batch runs
h5dump	Converts HDF5-formatted results to plain text
inspect	Inspects and gives information about data
version	Prints version information
cite	Prints citation information

The image shows the front page of a Nature Biotechnology journal issue. The header is red with the text 'nature biotechnology' and 'LETTERS'. Below the header, the URL 'https://doi.org/10.1038/s41587-021-00870-2' is visible. On the right, there is a 'Check for updates' button. The main title of the article is 'Modular, efficient and constant-memory single-cell RNA-seq preprocessing'. The authors listed are Páll Melsted<sup>1,8</sup>, A. Sina Booeshaghi<sup>2,8</sup>, Lauren Liu<sup>3</sup>, Fan Gao<sup>4,5</sup>, Lambda Lu<sup>4</sup>, Kyung Hoi (Joseph) Min<sup>6</sup>, Eduardo da Veiga Beltrame<sup>4</sup>, Kristján Eldjárn Hjörleifsson<sup>3</sup>, Jase Gehring<sup>7</sup> and Lior Pachter<sup>3,4</sup>. The abstract begins with: 'We describe a workflow for preprocessing of single-cell RNA-sequencing data that balances efficiency and accuracy. Our workflow is based on the kallisto and bustools programs, and is near optimal in speed with a constant memory requirement providing scalability for arbitrarily large datasets. The workflow is modular, and we demonstrate its flexibility by showing how it can be used for RNA velocity analyses.' The text continues: 'that is untenable given the pace of improvement in technology and the corresponding increase in data volume. In recent work, we introduced a format for scRNA-seq data that makes possible the development of efficient workflows by virtue of decoupling the computationally demanding step of associating reads to transcripts and genes (alignment) from the other steps required for scRNA-seq preprocessing<sup>10</sup>. This format, called BUS'

# Kallisto ‘bus tools’ for scRNAseq

```
(base) daniels-mbp-3:course_dataset danielbeiting$ kallisto  
kallisto 0.46.2
```

Usage: kallisto <CMD> [arguments] ..

Where <CMD> can be one of:

- index
- quant
- bus**
- pseudo
- merge
- h5dump
- inspect
- version
- cite

- Builds a kallisto index
- Runs the quantification algorithm
- Generate BUS files for single-cell
- Runs the pseudoalignment step
- Merges several batch runs
- Converts HDF5 formatted results to



## Modular and efficient pre-processing of single-cell RNA-seq

Posted July 26, 2019.

- [Download PDF](#)
- [Email](#)
- [Share](#)
- [Supplementary Material](#)
- [Data/Code](#)
- [XML](#)
- [Citation Tools](#)
- [Revision Summary](#)

“Our scRNA-seq workflow is up to 51 times faster than CellRanger...[and] requires a small fixed amount of constant memory that is independent of the number of reads being pre-processed. It is therefore suitable for low-cost and environmentally conscious cloud computing”

Analysis of single-cell RNA-seq data begins with pre-processing of sequencing reads to generate count matrices. We investigate algorithm choices for the challenges of pre-

Processed ~400M scRNAseq reads in ~16 min, compared to ~16 hours with Cell Ranger.

is modular, and we demonstrate its flexibility by showing how it can be used for RNA velocity analyses. Documentation and tutorials for using the kallisto I bus workflow are available at <https://www.kallistobus.tools/>.

# Open Kallisto output

*run\_info.json*

Name	Date modified	Type	Size
abundance.h5	11/03/2023 13:12	H5 File	2.834 KB
abundance	11/03/2023 13:12	TSV File	7.896 KB
run_info.json	11/03/2023 13:12	JSON File	1 KB

```
1 ▼ {  
2     "n_targets": 205131,  
3     "n_bootstraps": 0,  
4     "n_processed": 133387934,  
5     "n_pseudoaligned": 78655820,  
6     "n_unique": 18273031,  
7     "p_pseudoaligned": 59.0,  
8     "p_unique": 13.7,  
9     "kallisto_version": "0.48.0",  
10    "index_version": 10,  
11    "start_time": "Mon Jan 23 08:04:39 2023",  
12    "call": "kallisto quant -i Homo_sapiens.GRCh38.cdna.all.index -o SRR8668774 -t 24 --single -l 250  
-s 30 SRR8668774.fastq.gz"  
13 ▲ }  
14
```

**Open Kallisto output**  
*abundance.tsv*

Name	Date modified	Type	Size
abundance.h5	11/03/2023 13:12	H5 File	2.834 KB
abundance	11/03/2023 13:12	TSV File	7.896 KB
run_info.json	11/03/2023 13:12	JSON File	1 KB

abundance - Notepad					
File Edit Format View Help					
target_id	length	eff_length	est_counts	tpm	
ENST00000631435.1	12	3.83025	0	0	
ENST00000415118.1	8	3.3445	0	0	
ENST00000448914.1	13	3.91173	0	0	
ENST00000434970.2	9	3.4951	0	0	
ENST00000632684.1	12	3.83025	0	0	
ENST00000633010.1	16	4.09465	0	0	
ENST00000633009.1	20	4.24848	0	0	
ENST00000632524.1	11	3.73525	0	0	
ENST00000633353.1	31	4.49022	0	0	
ENST00000633765.1	31	4.49022	0	0	
ENST00000633159.1	21	4.27752	0	0	
ENST00000631884.1	17	4.14041	0	0	
ENST00000634070.1	18	4.18068	0	0	
ENST00000633504.1	31	4.49022	0	0	
ENST00000632542.1	18	4.18068	0	0	
ENST00000632619.1	28	4.43409	0	0	
ENST00000631895.1	23	4.32879	0	0	
ENST00000633030.1	19	4.21642	0	0	
ENST00000632911.1	31	4.49022	0	0	
ENST00000632968.1	17	4.14041	0	0	
ENST00000632963.1	20	4.24848	0	0	
ENST00000632473.1	31	4.49022	0	0	
ENST00000633968.1	20	4.24848	0	0	
ENST00000634085.1	16	4.09465	0	0	

abundance - Notepad					
File Edit Format View Help					
ENST00000371754.8	5235	4986	1208.55	4.16607	
ENST00000469991.1	640	391	14.5506	0.639615	
ENST00000371752.5	7212	6963	6866.66	16.9498	
ENST00000396105.6	7209	6960	12876.7	31.7988	
ENST00000371744.5	2964	2715	1972.76	12.4887	
ENST00000455070.1	2122	1873	34.6566	0.318026	
ENST00000509231.1	2224	1975	0	0	
ENST00000646235.1	3341	3092	0	0	
ENST00000644105.2	2821	2572	0	0	
ENST00000646118.1	3411	3162	3	0.016307	
ENST00000647304.1	3707	3458	1075.39	5.3451	
ENST00000395409.7	2192	1943	0	0	
ENST00000540171.2	1307	1058	2	0.0324907	
ENST00000675743.1	648	399	0	0	
ENST00000674723.1	931	682	0	0	
ENST00000674869.1	837	588	0	0	
ENST00000675248.1	577	328	0	0	
ENST00000676427.1	883	634	32.2398	0.874013	
ENST00000330634.11	7578	7329	289.566	0.679076	
ENST00000392634.9	7623	7374	0.956435	0.00222929	
ENST00000675482.1	1020	771	20.6712	0.460816	
ENST00000398337.8	1689	1440	224.596	2.68074	
ENST00000674966.1	859	610	8.63282	0.243242	
ENST00000675207.1	7732	7483	58.9235	0.13534	
ENST00000675616.1	431	182	0	0	

Name	Date modified	Type	Size
abundance.h5	11/03/2023 13:12	H5 File	2.834 KB
abundance	11/03/2023 13:12	TSV File	7.896 KB
run_info.json	11/03/2023 13:12	JSON File	1 KB

abundance - Notepad

target_id	length	eff_length	est_counts	tpm
ENST00000631435.1	12	3.83025	0	0
ENST00000415118.1	8	3.3445	0	0
ENST00000448914.1	13	3.91173	0	0
ENST00000434970.2	9	3.4951	0	0
ENST00000632684.1	12	3.83025	0	0
ENST00000633010.1	16	4.09465	0	0
ENST00000633009.1	20	4.24848	0	0
ENST00000632524.1	11	3.73525	0	0
ENST00000633353.1	31	4.49022	0	0
ENST00000633765.1	31	4.49022	0	0
ENST00000633159.1	21	4.27752	0	0
ENST00000631884.1	17	4.14041	0	0
ENST00000634070.1	18	4.18068	0	0
ENST00000633504.1	31	4.49022	0	0
ENST00000632542.1	18	4.18068	0	0
ENST00000632619.1	28	4.43409	0	0
ENST00000631895.1	23	4.32879	0	0
ENST00000633030.1	19	4.21642	0	0
ENST00000632911.1	31	4.49022	0	0
ENST00000632968.1	17	4.14041	0	0
ENST00000632963.1	20	4.24848	0	0
ENST00000632473.1	31	4.49022	0	0
ENST00000633968.1	20	4.24848	0	0
ENST00000634085.1	16	4.09465	0	0

# The ‘effective length’ of a transcript

The length of a transcript after adjusting for the total number of possible positions a fragment of size X could originate from

$$L_{\text{effective}} = L_{\text{actual}} - L_{\text{fragment}} + 1$$

<i>transcript</i>	ATGCGTAACATG	$L_{\text{actual}} = 12$	
<i>fragment</i>	NNN	$L_{\text{fragment}} = 3$	$L_{\text{effective}} = 10$



# RNAseq gives **relative** quantification of gene expression

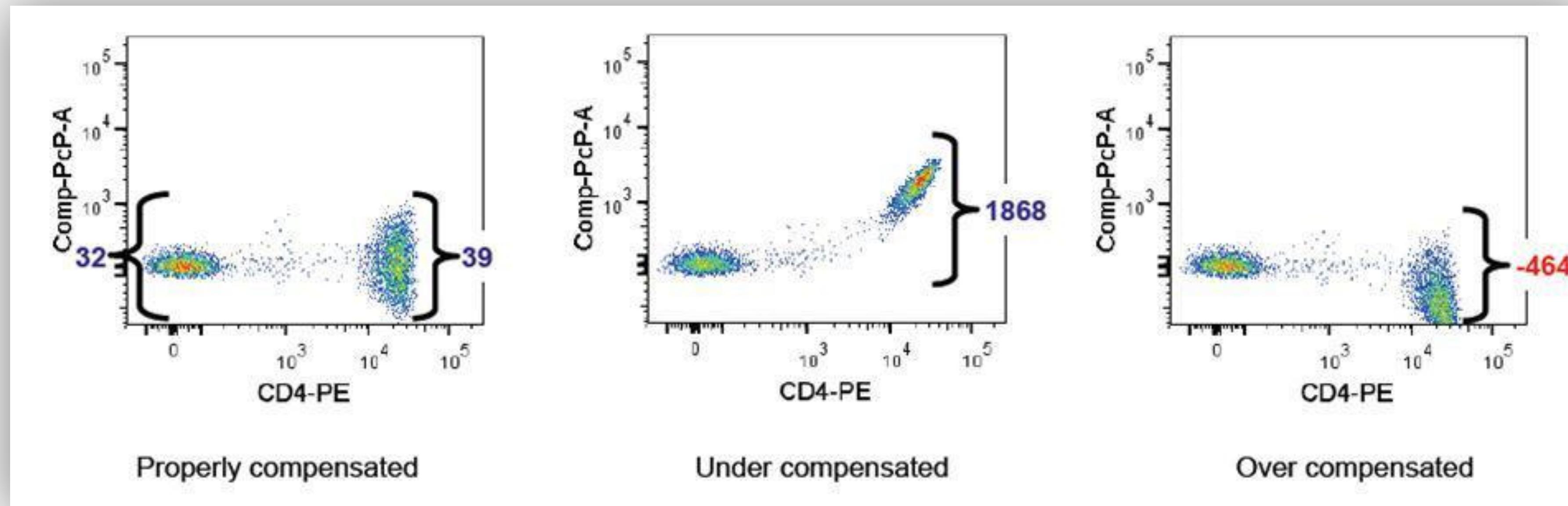
“All commonly used techniques to measure mRNA abundance, including qPCR, microarray signals, as well as reads per kilobase per million reads (RPKM) for RNAseq data, aim at estimating a statistic that is as closely proportional to the **relative molar concentration** as possible.”

- Wagner, *Theories in Biosci.*, 2012

**Understanding units of  
measurement for RNAseq is  
critical to understanding how  
we determine differential  
expression**

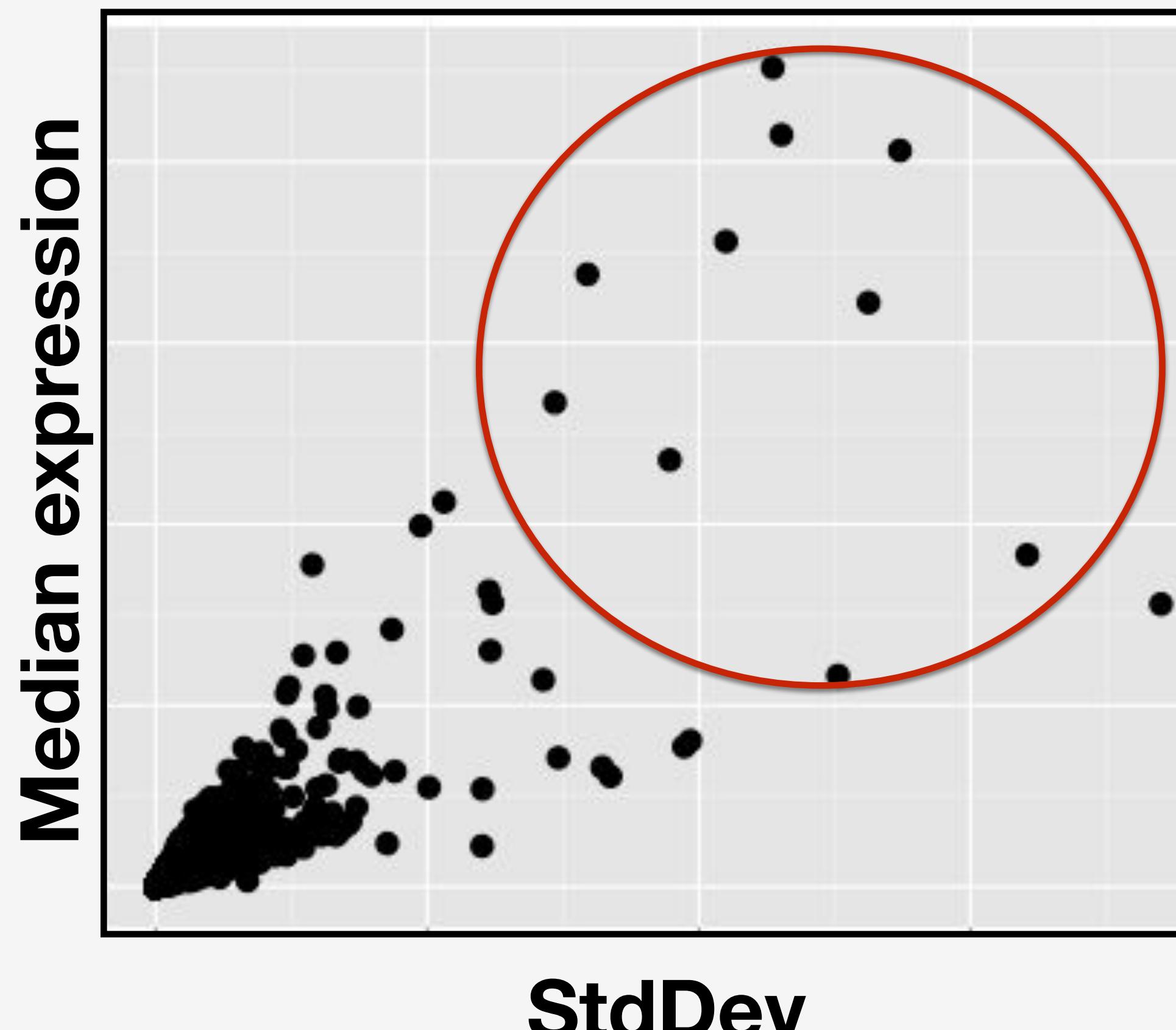
**Normalization**  
*generic term referring to any  
number of ways that a dataset  
is globally altered to improve  
our ability to detect DEGs*

# RNAseq normalization ≈ Flow cytometry compensation



## raw counts

	mouse 1	mouse 2	mouse 3	mouse 4	mouse 5	StDev
gene A	5	10	10	15	10	<b>3.5</b>
gene B	115	110	100	115	118	<b>7.1</b>
gene C	1000	1100	1050	1045	1030	<b>36.4</b>
gene D	8000	9000	10000	6000	7030	<b>1576.4</b>



**genes with higher  
expression have higher  
stdDev**

**Heteroscedasticity**

## Log 2

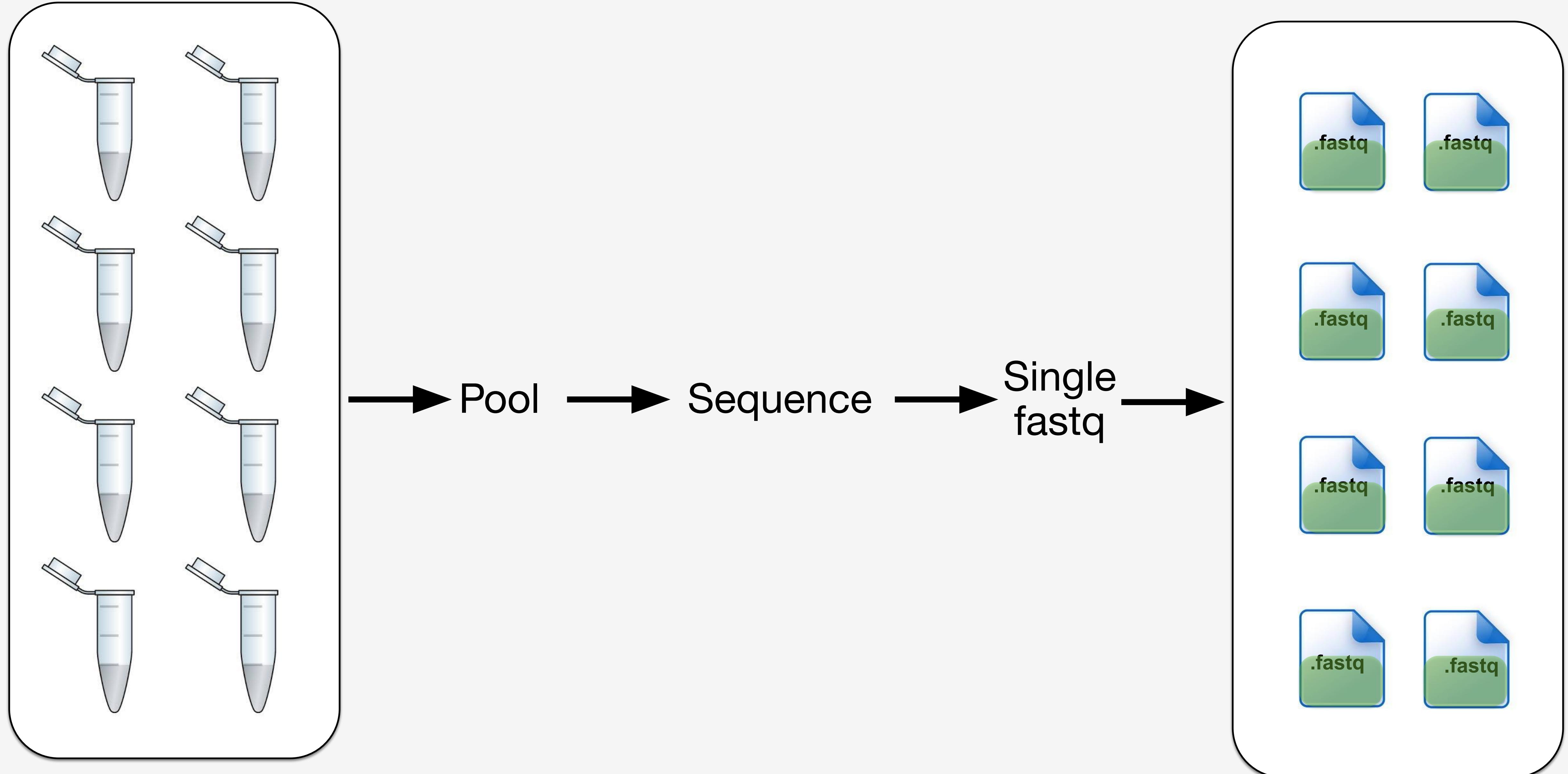
	mouse 1	mouse 2	mouse 3	mouse 4	mouse 5	StDev
gene A	2.3	3.3	3.3	3.9	3.3	<b>0.57</b>
gene B	6.8	6.8	6.6	6.8	6.9	<b>0.09</b>
gene C	10.0	10.1	10.0	10.0	10.0	<b>0.05</b>
gene D	13.0	13.1	13.3	12.6	12.8	<b>0.29</b>

With large datasets, fixing one problem often creates another

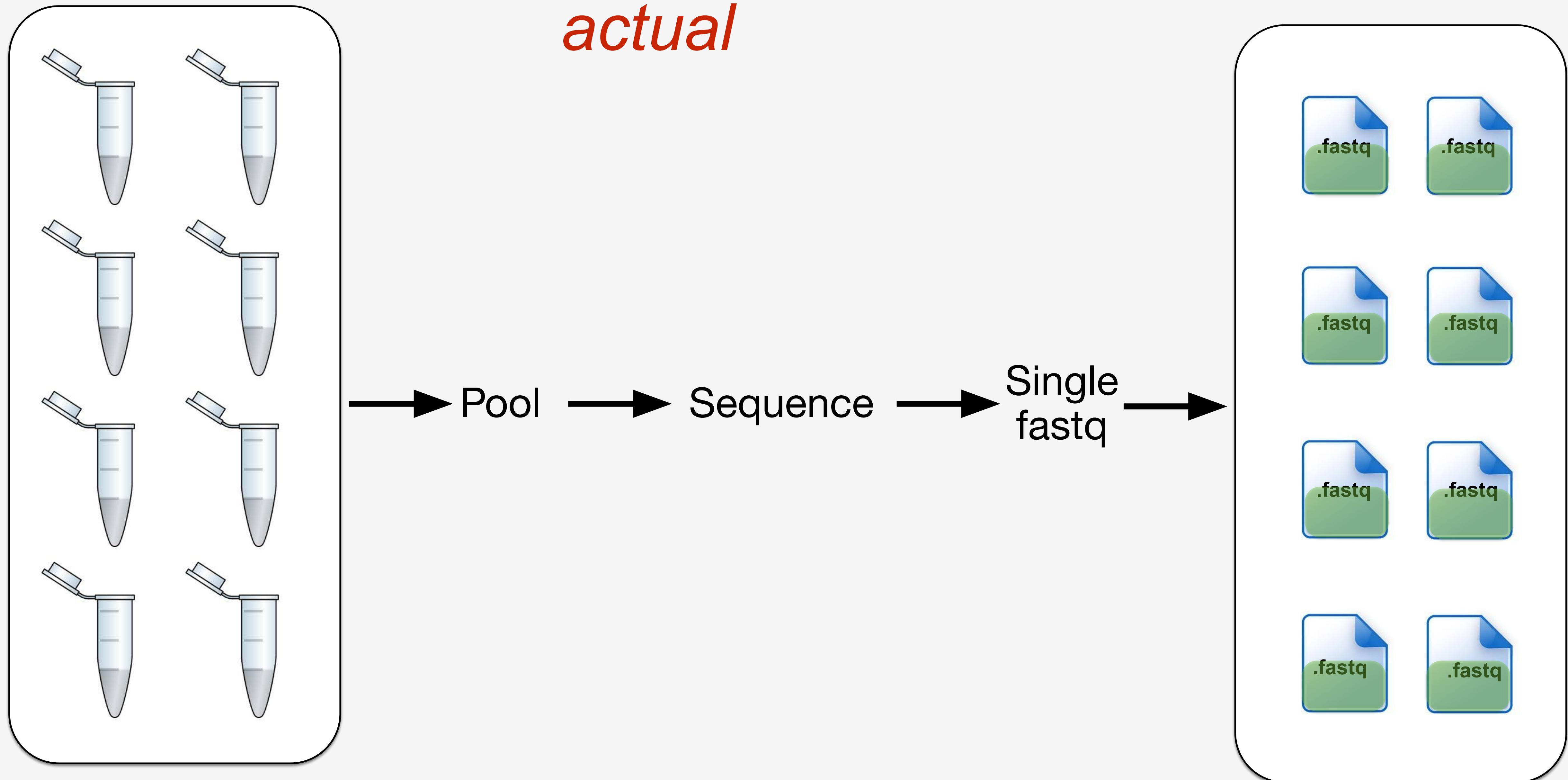


**genes with lower expression have highest StdDev**

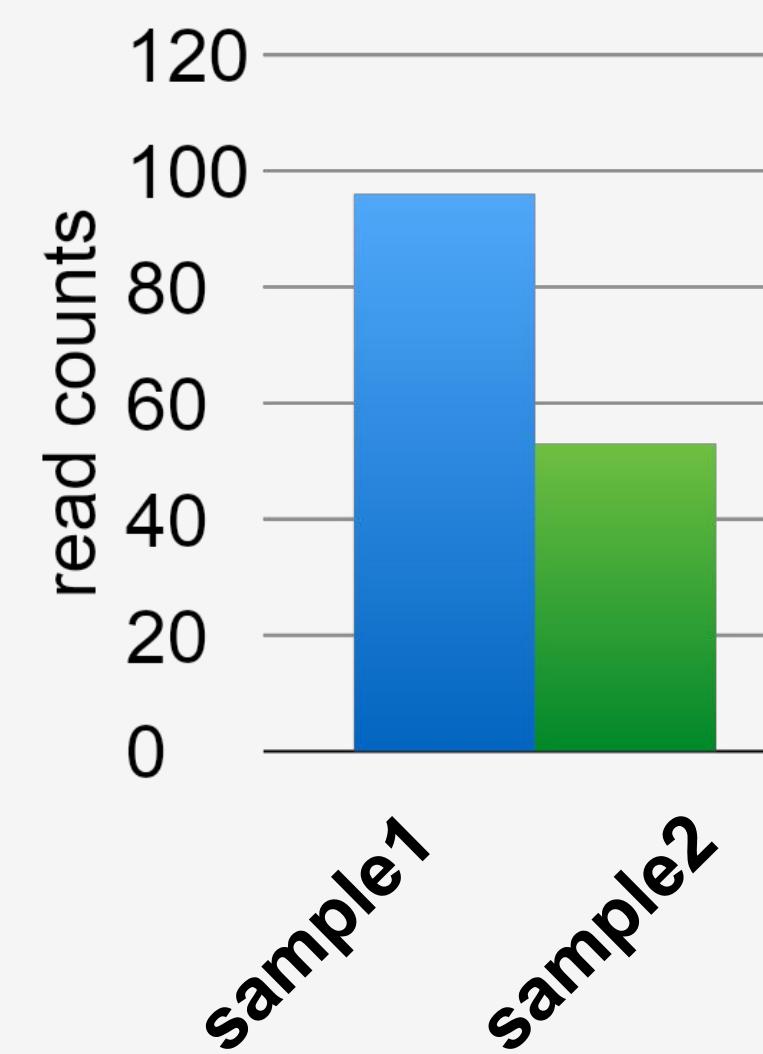
# *expected* reads per sample



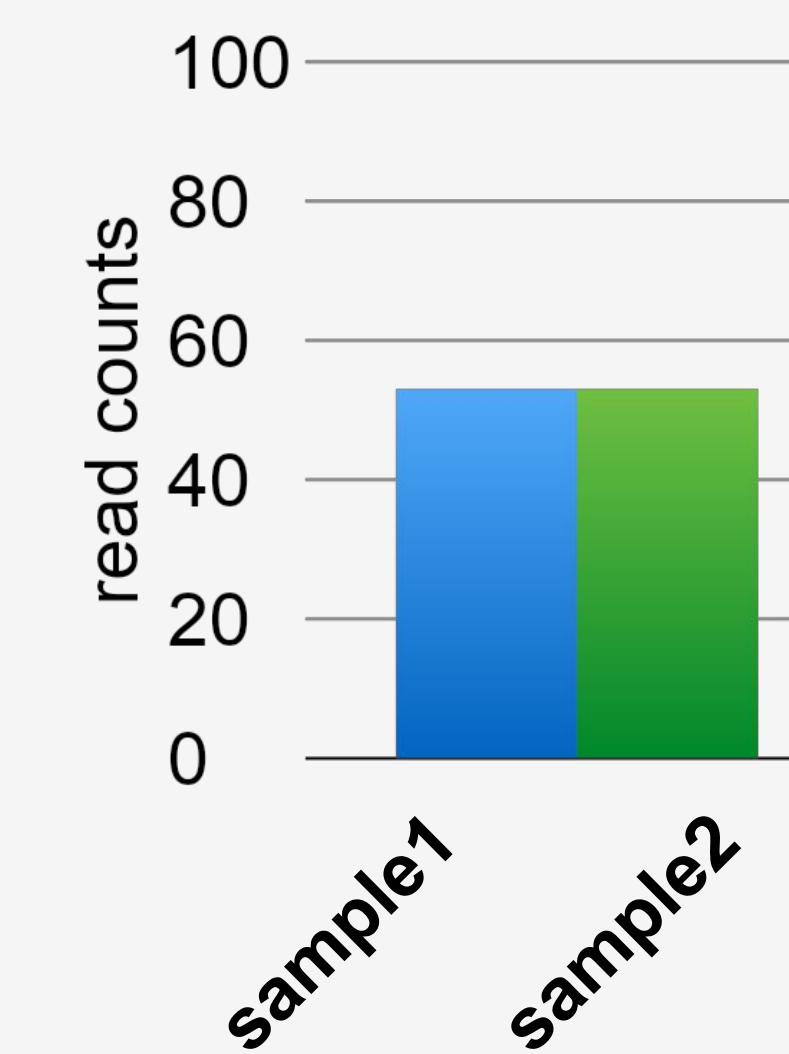
# ~~expected~~<sup>actual</sup> reads per sample



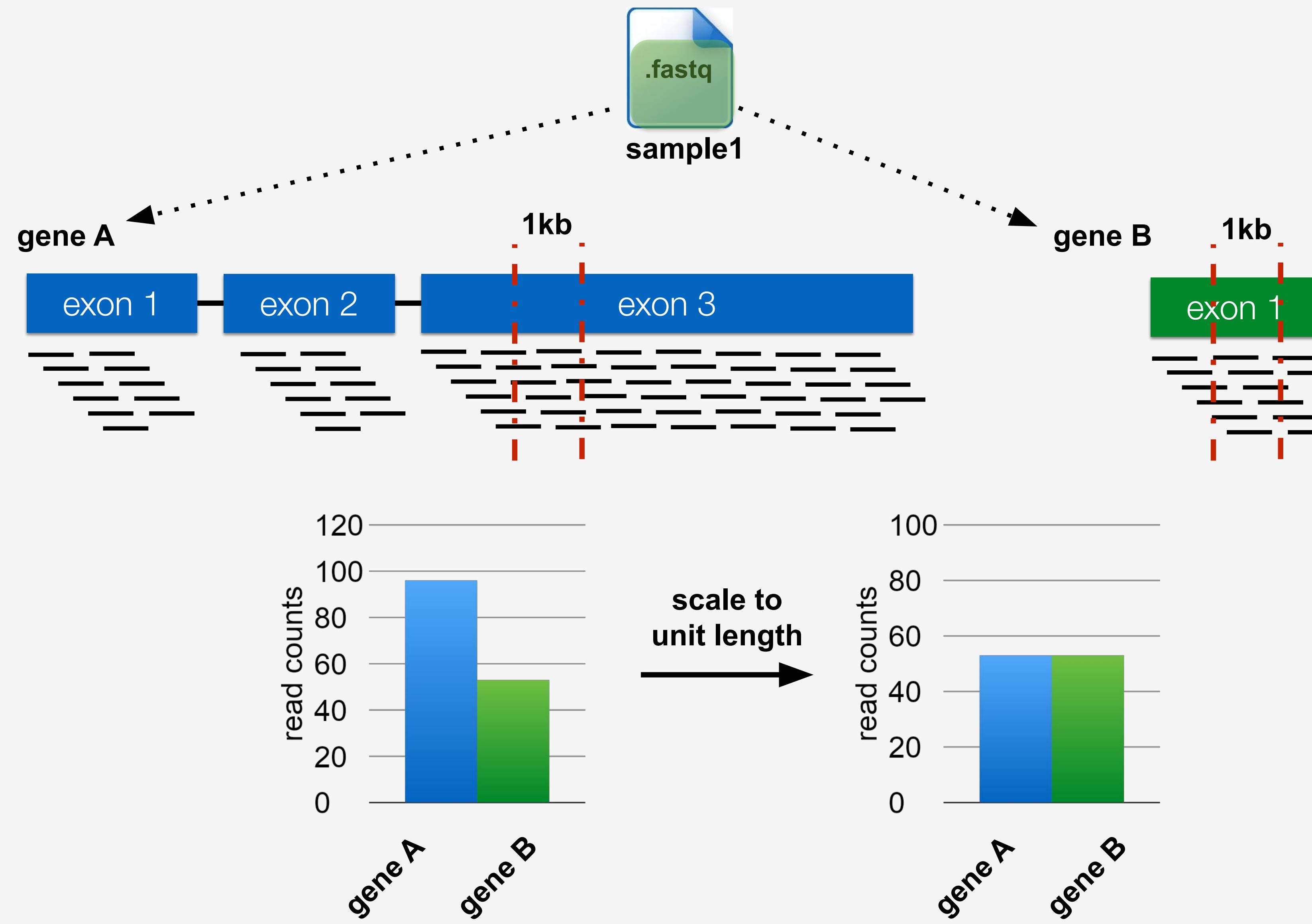
# scaling units for between libraries



scale to  
unit depth

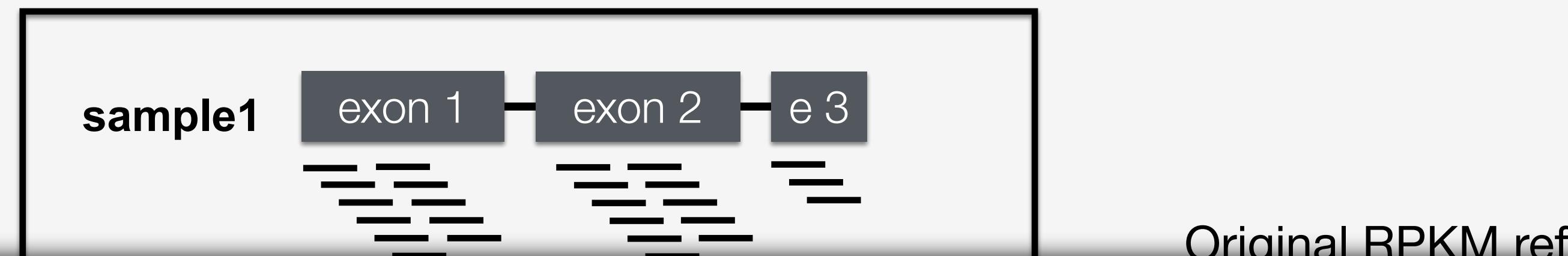


# scaling units for within-sample comparisons



# reads per kilobase, per million reads sequenced (RPKM)

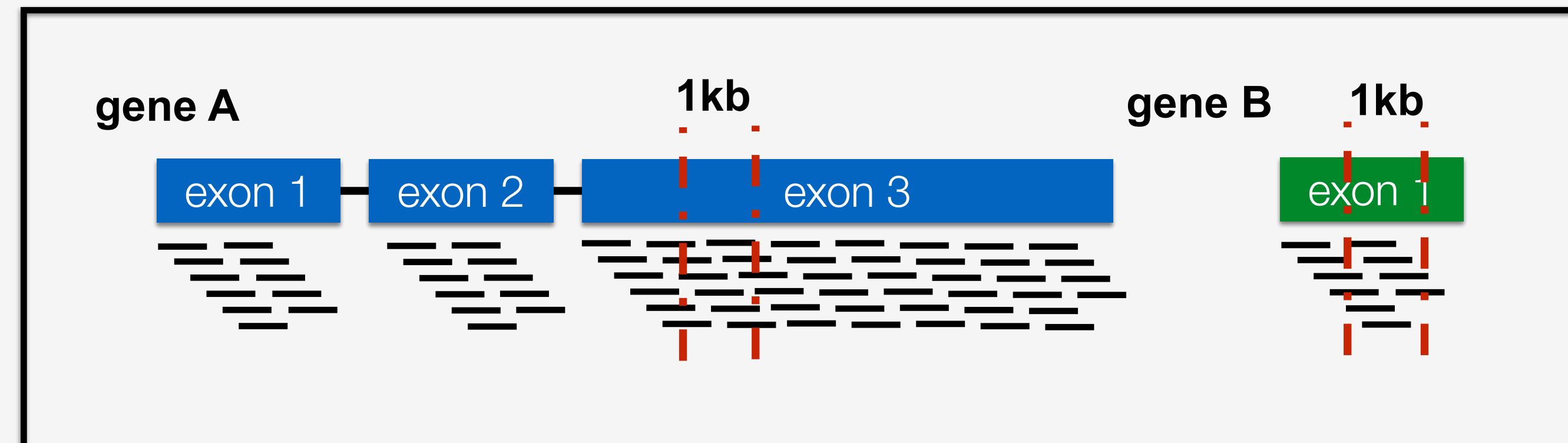
between



RPKM is really just a unit of expression for  
RNAseq data

+

*within*  
sample



# RPKM (FPKM)

$$\text{RPKM} = \frac{\text{\# reads mapped to genomic region}}{(\text{region length in kb})(\text{total \# reads})}$$

# RPKM (FPKM)

$$\text{RPKM} = \frac{\text{\# reads mapped to genomic region}}{(\text{region length in kb})(\text{total \# reads})}$$

$$\text{RPKM} = \frac{1000}{(5)(20000000)} = 0.00001$$

# RPKM (FPKM)

$$\text{RPKM} = \frac{\text{\# reads mapped to genomic region}}{(\text{region length in kb})(\text{total \# reads})} \times 10^6$$

$$\text{RPKM} = \frac{1000}{(5)(20000000)} = 0.00001$$

# RPKM (FPKM)

$$\text{RPKM} = \frac{\text{\# reads mapped to genomic region}}{(\text{region length in kb})(\text{total \# reads})} \times 10^6$$

$$\text{RPKM} = \frac{1000}{(5)(20000000)} = 10$$

# RPKM alone, is not sufficient for normalization

**Table 3:** Summary of comparison results for the seven normalization methods under consideration

Method	Distribution	Intra-Variance	Housekeeping	Clustering	False-positive rate
TC	-	+	+	-	-
UQ	++	++	+	++	-
Med	++	++	-	++	-
<b>DESeq</b>	<b>DESeq</b>	++	++	++	++
<b>TMM</b>	<b>EdgeR</b>	++	++	++	++
Q	++	-	+	++	-
RPKM	-	+	+	-	-

A '-' indicates that the method provided unsatisfactory results for the given criterion, while a '+' and '++' indicate satisfactory and very satisfactory results for the given criterion.

# The problem with RPKM

***read counts from each gene***

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total reads
Sample 1	80	10	6	3	1	100
Sample 2	20	20	10	50	400	500

$$\text{RPKM} = \frac{\# \text{ reads mapped to genomic region}}{(\text{region length in kb})(\text{total } \# \text{ reads})} \times 10^6$$

***RPKM each gene***

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPKM
Sample 1	8000	2000	2400	6000	10000	28400
Sample 2	400	800	800	20000	800000	822000

# RNAseq gives **relative** quantification of gene expression

“All commonly used techniques to measure mRNA abundance, including qPCR, microarray signals, as well as reads per kilobase per million reads (RPKM) for RNAseq data, aim at estimating a statistic that is as closely proportional to the **relative molar concentration** as possible.”

- Wagner, *Theories in Biosci.*, 2012

“The average relative molar concentration for each and every sample of RNA-seq data mapped to the same genome is the same constant value.

- Wagner, *Theories in Biosci.*, 2012

# The problem with RPKM

***read counts from each gene***

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total reads
Sample 1	80	10	6	3	1	100
Sample 2	20	20	10	50	400	500

$$\text{RPKM} = \frac{\# \text{ reads mapped to genomic region}}{(\text{region length in kb})(\text{total } \# \text{ reads})} \times 10^6$$

***RPKM each gene***

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPKM	
Sample 1	8000	2000	2400	6000	10000	28400	$\mu = 5680$
Sample 2	400	800	800	20000	800000	822000	$\mu = 164400$

# Fixing RPKM is easy

***read count***

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total reads
Sample 1	80	10	6	3	1	<b>100</b>
Sample 2	20	20	10	50	400	<b>500</b>

***scale by gene length  
first***



***this becomes a  
normalization factor***

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPK
Sample 1	0.8	0.2	0.24	0.6	1	<b>2.84</b>
Sample 2	0.2	0.4	0.4	10	400	<b>411</b>

# Fixing RPKM is easy

$$\text{TPM} = \frac{\text{reads per Kb}}{\text{total RPK in sample}}$$

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPK
Sample 1	0.8	0.2	0.24	0.6	1	<b>2.84</b>
Sample 2	0.2	0.4	0.4	10	400	<b>411</b>

# Fixing RPKM is easy

$$\text{TPM} = \frac{\text{reads per Kb}}{\text{total RPK in sample}}$$

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPK
Sample 1	0.8/2.84	0.2/2.84	0.24/2.84	0.6/2.84	1/2.84	2.84
Sample 2	0.2/411	0.4/411	0.4/411	10/411	400/411	411

# Fixing RPKM is easy

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total TPM
Sample 1	0.281690141	0.070422535	0.084507042	0.211267606	0.352112676	1
Sample 2	0.000486618	0.000973236	0.000973236	0.0243309	0.97323601	1

$$\text{TPM} = \frac{\text{reads per Kb}}{\text{total RPK in sample}}$$

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPK
Sample 1	0.8/2.84	0.2/2.84	0.24/2.84	0.6/2.84	1/2.84	2.84
Sample 2	0.2/411	0.4/411	0.4/411	10/411	400/411	411

# Fixing RPKM is easy

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total TPM	
Sample 1	281690	70423	84507	211268	352113	1000000	$\mu = 200000$
Sample 2	487	973	973	24331	973236	1000000	$\mu = 200000$

$$\text{TPM} = \frac{\text{reads per Kb}}{\text{total RPK in sample}} \times 10^6$$

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPK	
Sample 1	0.8/2.84	0.2/2.84	0.24/2.84	0.6/2.84	1/2.84	2.84	
Sample 2	0.2/411	0.4/411	0.4/411	10/411	400/411	411	

# Further normalization steps need to be taken before one can compare between samples

“library size scaling is too simple for many biological applications. The number of fragments expected to map to a gene is not only dependent on the expression level and length of the gene, **but also the composition of the RNA population that is being sampled**. Thus, if a large number of genes are unique to, or highly expressed in, one experimental condition, the sequencing 'real estate' available for the remaining genes in that sample is decreased. If not adjusted for, this sampling artifact can force the DE analysis to be skewed towards one experimental condition.”

# Further normalization steps need to be taken before one can compare between samples

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total reads
Sample 1	80	10	6	3	1	100
Sample 2	20	20	10	50	400	500

Calculate RPK

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPK
Sample 1	0.8	0.2	0.24	0.6	1	2.84
Sample 2	0.2	0.4	0.4	10	400	411

Calculate TPM

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPK
Sample 1	281690	70423	84507	211268	352113	1000000
Sample 2	487	973	973	24331	973236	1000000

Differential expression is all messed up! Why?!

# Further normalization steps need to be taken before one can compare between samples

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total reads
Sample 1	80	10	6	3	1	100
Sample 2	20	20	10	50	50	100

Calculate RPK

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPK
Sample 1	0.8	0.2	0.24	0.6	1	2.84
Sample 2	0.2	0.4	0.4	10	10	11

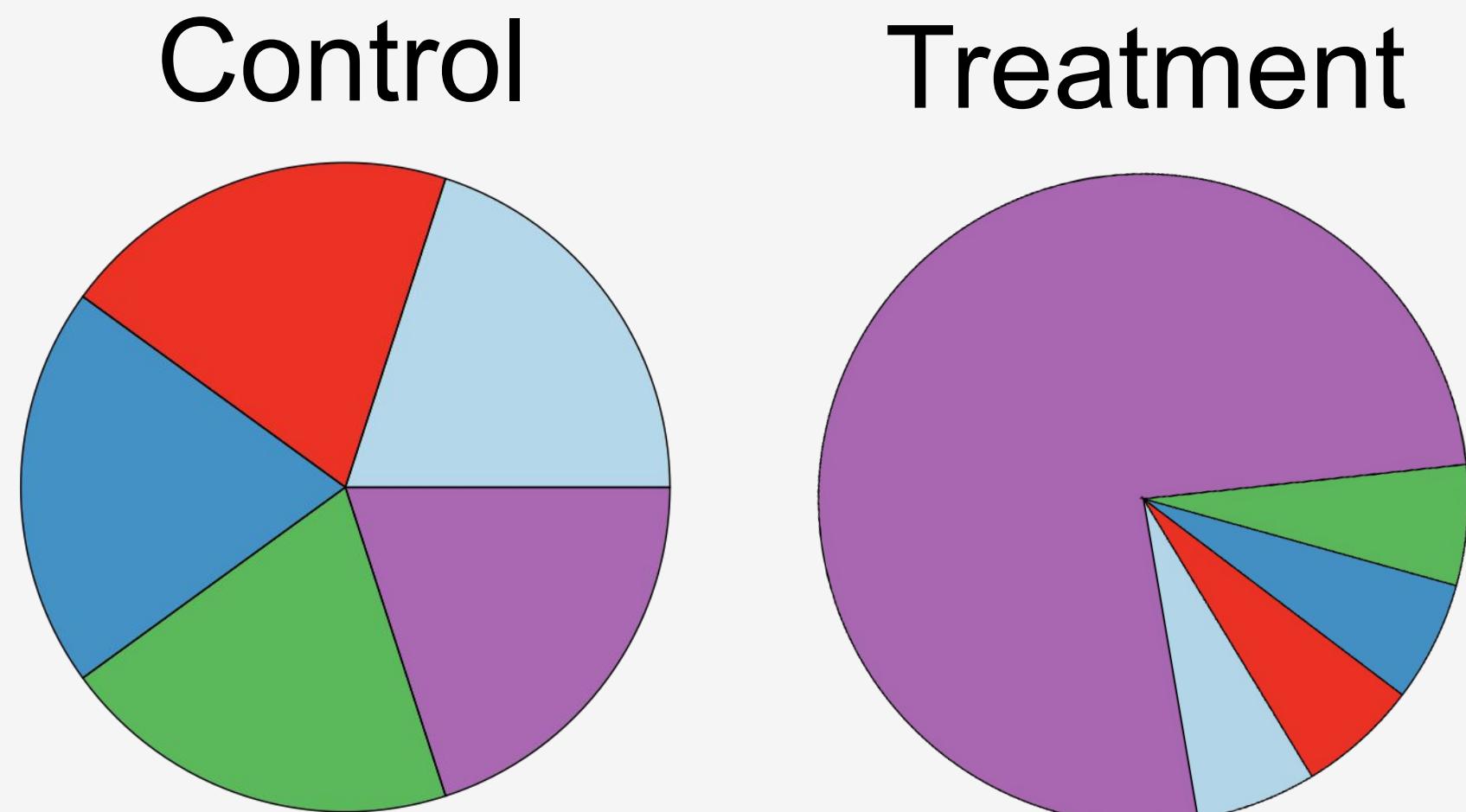
Calculate TPM

	Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	total RPK
Sample 1	281690	70423	84507	211268	352113	1000000
Sample 2	18181	36363	36363	909090	909090	1000000

most current normalization methods attempt to find a set of genes with minimal variance across samples for normalization

# Further normalization steps need to be taken before one can compare between samples

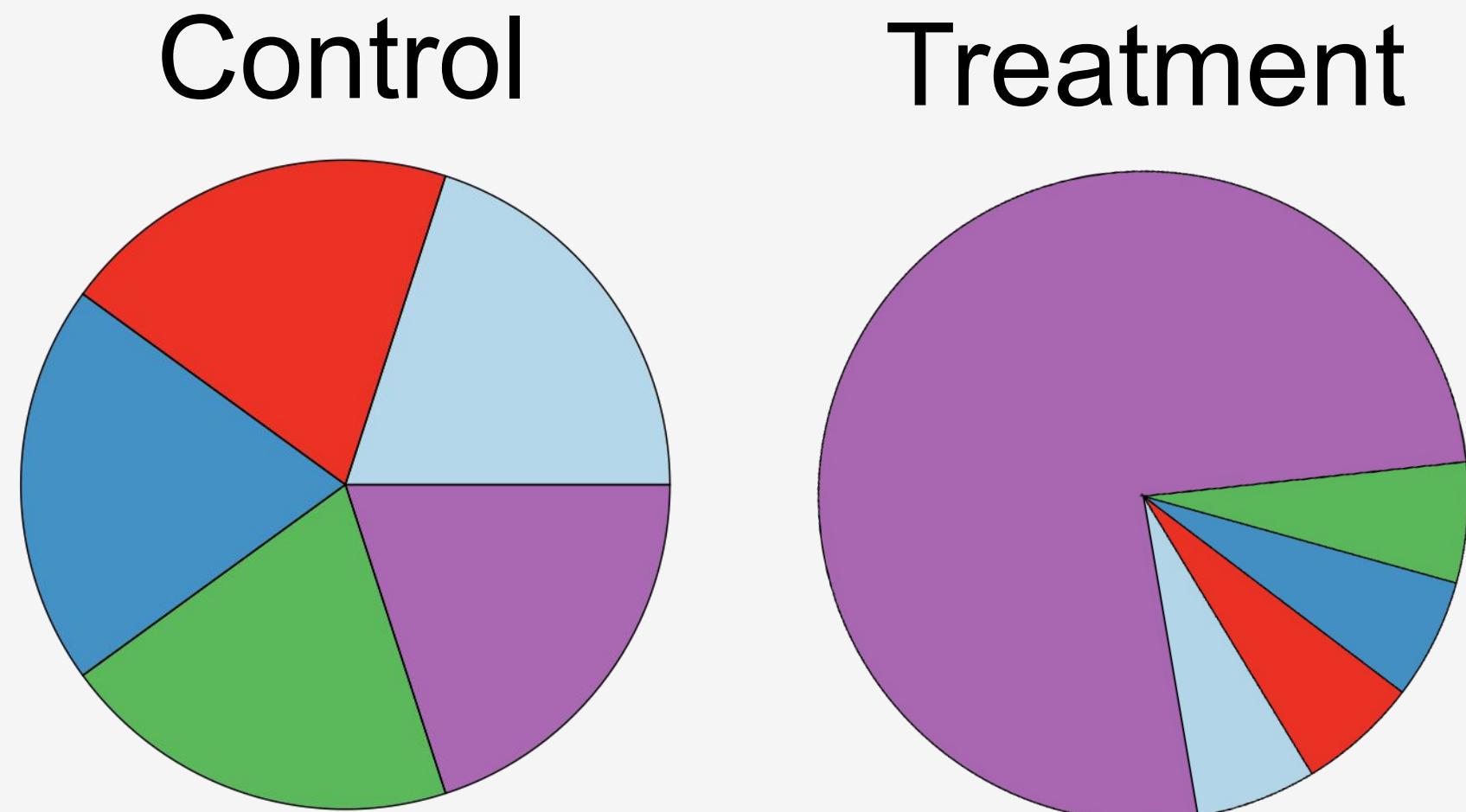
	Gene A	Gene B	Gene C	Gene D	Gene E	total counts
Control	0.2	0.2	0.2	0.2	0.2	10
Treatment	0.06	0.06	0.06	0.06	0.76	100



every gene appears to be differentially expressed!

# Further normalization steps need to be taken before one can compare between samples

		Gene A 100kb	Gene B 50kb	Gene C 25kb	Gene D 5kb	Gene E 1kb	Total RPK
	control	0.25	0.25	0.25	0.25	0.25	8
	treatment	0.25	0.25	0.25	0.25	3.17	24*



most current normalization methods attempt to find a set of genes with minimal variance across samples for normalization

# **Statistical tools for normalization and differential expression analysis make certain assumptions about your data**

## **Key assumptions**

1. Most genes are not differentially expressed (implications for comparing very different treatments/conditions)
2. Approx. equivalent numbers of up and down regulated genes

**End of Part 1**

# Setting up our analysis project



R version 3.6.3 (2020-02-29) -- "Holding the Windsock"  
Copyright (C) 2020 The R Foundation for Statistical Computing  
Platform: x86\_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> |

RStudio

Console Terminal x ~/

Environment History Connections

Import Dataset Global Environment

Files Plots Packages Help Viewer

Install Update

Name Description

System Library

Name	Description	Version	Action
abind	Combine Multiple Matrices into a Single Array	1.8-1	Install
acepack	ACE and AV Methods for Linear and Logistic Regression	1.3-3	Install
ade4	Analysis of Ecological Data: Multivariate Methods in Environmental Sciences	1.7-12	Install
annotate	Annotation for microarrays	1.64.0	Install
AnnotationDbi	Manipulation of SQLite-based annotations in Bioconductor	1.48.0	Install
AnnotationFilter	Facilities for Filtering Bioconductor Annotation Resources	1.10.0	Install
AnnotationHub	Client to access AnnotationHub resources	2.18.0	Install
ape	Analyses of Phylogenetics and Evolution	5.3	Install
askpass	Safe Password Entry for R, Git, and SSH	1.1	Install
assertthat	Easy Pre and Post Assertions	0.2.1	Install
ATACseqQC	ATAC-seq Quality Control	1.10.4	Install
audio	Audio Interface for R	0.1-7	Install
available	Check if the Title of a Package is Available, Appropriate and Interesting	1.0.4	Install
backports	Reimplementations of Functions Introduced Since R-3.0.0	1.1.6	Install
base	The R Base Package	3.6.3	Install
base64enc	Tools for base64 encoding	0.1-3	Install
batman	Convert Categorical Representations of Logicals to Actual Logicals	0.1.0	Install
bayesm	Bayesian Inference for Marketing/Micro-Econometrics	3.1-4	Install
beepR	Easily Play Notification Sounds on any Platform	1.3	Install
BH	Boost C++ Header Files	1.72.0-3	Install

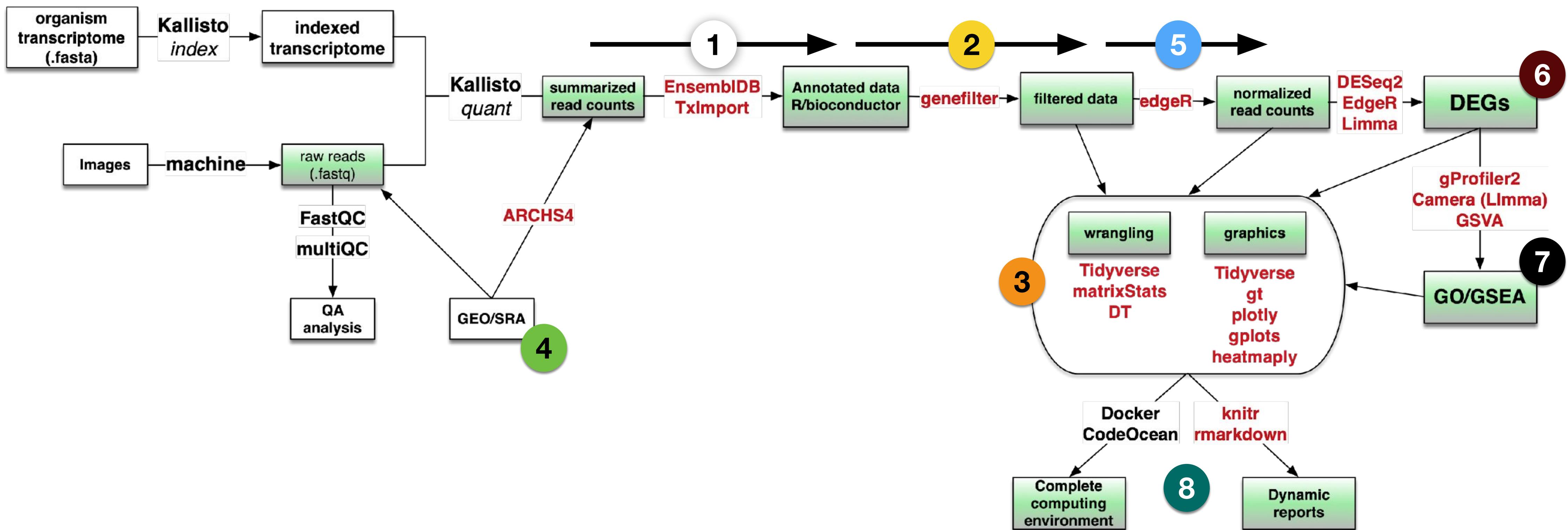
New Project... Open Project... Open Project in New Session... Close Project Clear Project List Project Options...

# R projects

- A new R session (process) is started
- The **.Rprofile** file in the project's main directory (if any) is sourced by R
- The **.RData** file in the project's main directory is loaded (if project options indicate that it should be loaded).
- The **.Rhistory** file in the project's main directory is loaded into the RStudio History pane (and used for Console Up/Down arrow command history).
- The **current working directory** is set to the project directory.
- Previously edited **source documents** are restored into editor tabs
- Other **RStudio settings** (e.g. active tabs, splitter positions, etc.) are restored to where they were the last time the project was closed.

**R script ≠ RStudio project**

# We'll break up our analysis into 'step' scripts



# A homework assignment

The screenshot shows the DataCamp Learn platform interface. At the top, there's a navigation bar with links for Home, Learn (which is highlighted in blue), Workspace, Certification, Jobs, Join Groups, a search bar, an Upgrade button, and a user profile icon. The main content area features a dark-themed card for an 'INTERACTIVE COURSE' titled 'Introduction to the Tidyverse'. Below the title are two buttons: 'Continue Course' (green) and 'Bookmark' (white). Below these buttons, course statistics are displayed: 4 hours, 16 Videos, 50 Exercises, 274,309 Participants, and 4,150 XP. To the left of the main card is a sidebar with a dark background containing various navigation links: Progress, Bookmarks, Leaderboard, Catalog (with sub-links for Tracks, Courses, Practice, Projects, Competitions - marked as BETA, Assessments, and Live Events). At the bottom of the sidebar, a purple button says 'Getting Started (1/4)'. On the right side of the main card, there's a purple box listing tracks this course is part of: Data Analyst with R, Data Scientist with R, R Programmer, and Tidyverse Fundamentals with R. Below this, a section titled 'DATASETS' shows a link to 'Gapminder' with a 'Explore datasets' button.

<https://app.datacamp.com/learn/courses/introduction-to-the-tidyverse>

[https://github.com/Lukalgnjatovic/DataCamp\\_-\\_Track\\_-\\_Data\\_Scientist\\_with\\_R\\_-\\_Course\\_03\\_-\\_Introduction\\_to\\_the\\_Tidyverse](https://github.com/Lukalgnjatovic/DataCamp_-_Track_-_Data_Scientist_with_R_-_Course_03_-_Introduction_to_the_Tidyverse)