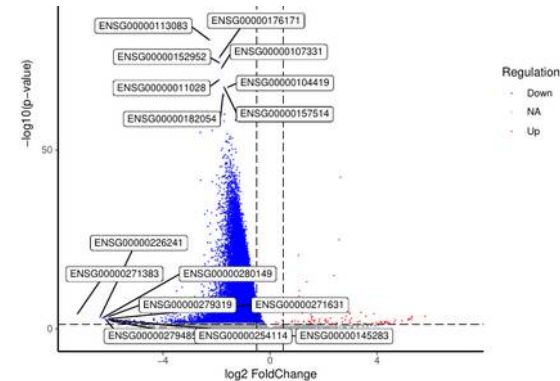


**fli**Leibniz Institute on Aging –  
Fritz Lipmann Institute

## Differential gene expression analyses under global transcriptional shifts



# Bulk RNAseq in a nutshell

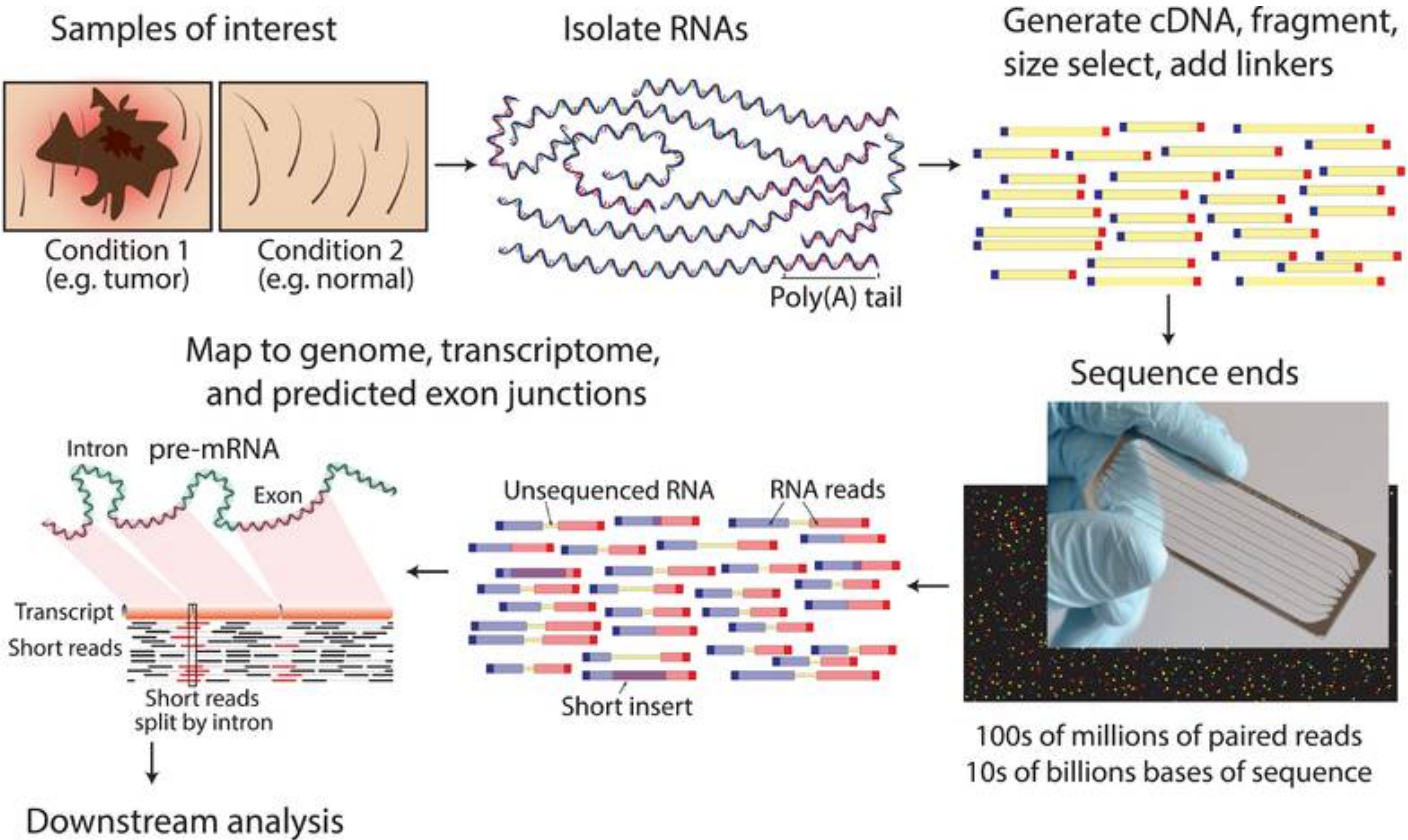


Image from <https://en.wikipedia.org/wiki/File:Journal.pcbi.1004393.g002.png>

# Global transcriptional shutdown or amplification

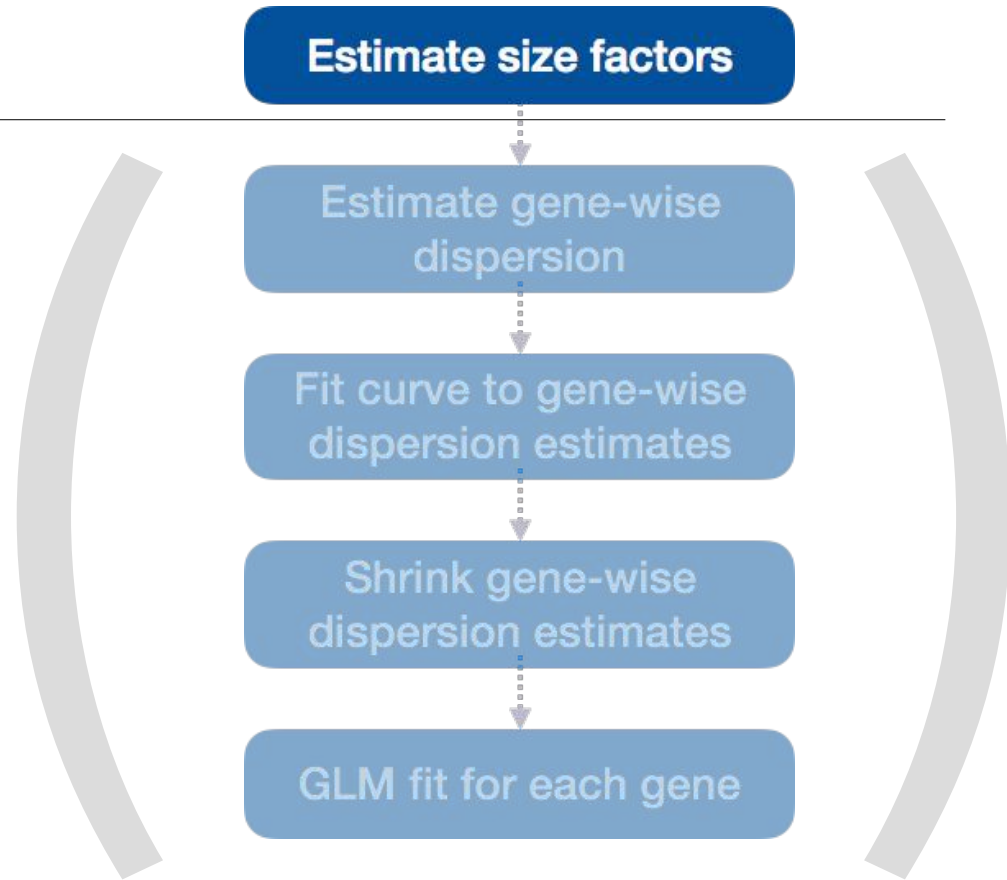
<b>Shutdown</b>	<b>Amplification</b>
Heat-shock	Cancer
Starvation	Embryonic development
Viral infections	Organ/tissue regeneration

See e.g. also: Percharde M, Bulut-Karslioglu A, Ramalho-Santos M. Hypertranscription in Development, Stem Cells, and Regeneration. Dev Cell. 2017 Jan 9;40(1):9-21. doi: 10.1016/j.devcel.2016.11.010. Epub 2016 Dec 15. PMID: 27989554; PMCID: PMC5225143.

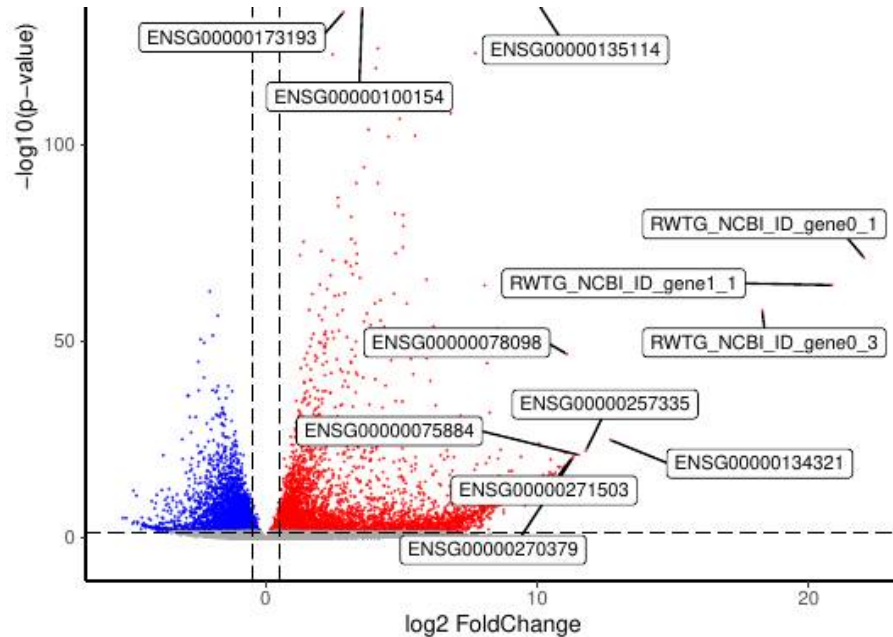
# Roadmap

DESeq2's workflow.

Today, the roadtrip ends here



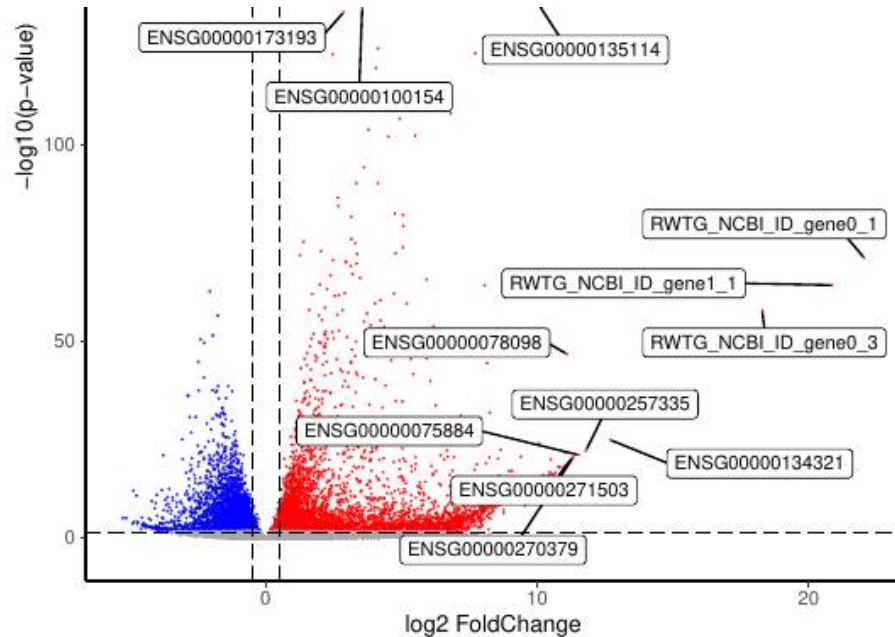
# Control vs Rift valley fever virus infected cells - observed vs expected



Observed: 3458 down- and 4534 upregulated genes

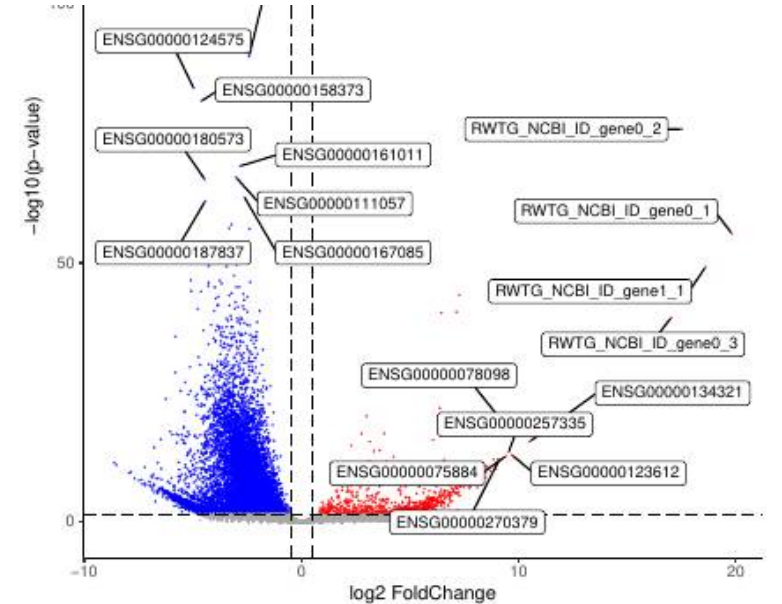
Pinkham et al. (2017): “*The large number of upregulated genes identified in our RNA-seq study was somewhat surprising, given the ability of NSs to suppress host transcription*”

# Control vs Rift valley fever virus infected cells - observed vs expected



Observed: 3458 down- and 4534 upregulated genes

DESeq2's default sizefactor estimation



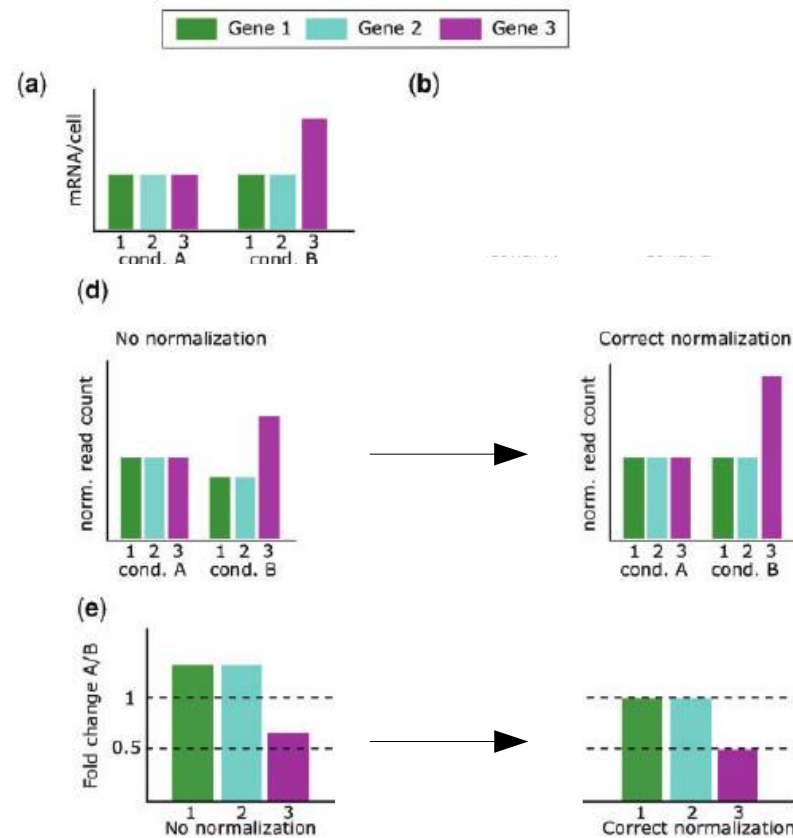
12851 down- and 778 upregulated genes

New sizefactor estimation method



# From mRNA/cell to fold changes

## Sizefactor estimation



# DESeq2's *median of ratios* method

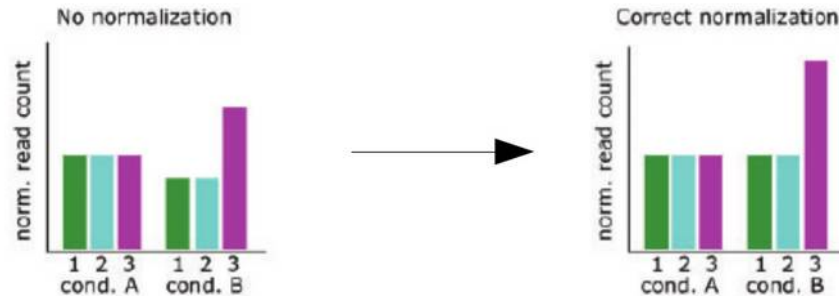
- Each sample is assigned **one** single sizefactor  $\hat{s}_j$
- The sizefactors are determined by the *median of ratios* to the geometric mean of all samples.
- (The transcript length shortens out, so DESeq2 doesn't need it)

$$\hat{s}_j = \underset{i}{\text{median}} \frac{k_{ij}}{\left( \prod_{v=1}^m k_{iv} \right)^{1/m}}$$



# DESeq2's *median of ratios* method

Sample A	Sample B	Pseudo-reference (geometric mean)	Ratio of reference to Sample A	Ratio of reference to Sample B	Sample A scaled	Sample B scaled
1	4	2	2/1	2/4	3	$1, \bar{3}$
1	9	3	3/1	3/9	3	3
1	16	4	4/1	4/16	3	$5, \bar{3}$



**But:** this only works, if most genes are not DE

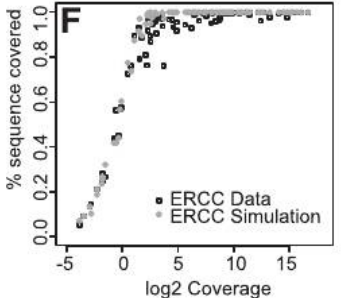
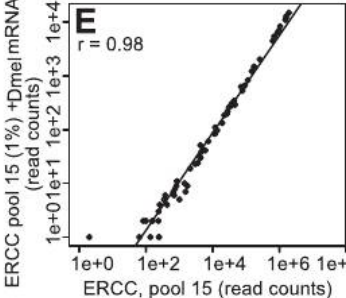
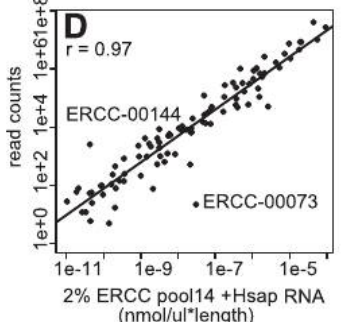
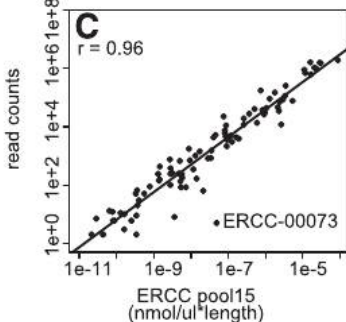
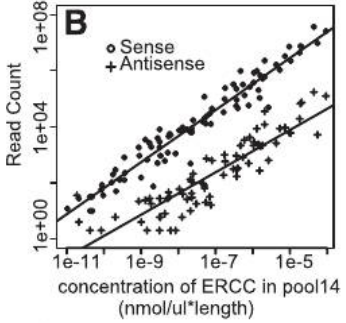
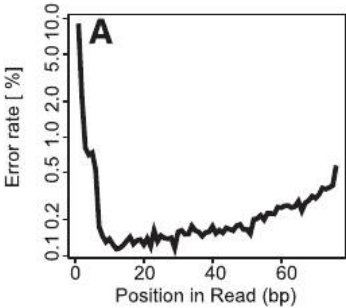
# RNA spike-ins

- The External RNA Control Consortium (ERCC) provides 96 synthetic RNAs with various lengths, and GC content covering a  $2^{20}$  concentration range
- Same amount of spike-ins applied to each sample
- DESeq2 estimates sizefactors based on the spike-ins
- **But:** cells have to be counted 😞
- Spike-ins often unavailable 😞
- *Housekeeping genes* potentially unreliable (inhibition of Pol II) 😞



Invitrogen™  
**ERCC RNA Spike-In Mix**  
 Catalog number: 4456740  
 Related applications: [Real Time PCR \(qPCR\)](#) | [RNA Sequencing](#) | [SOLID® Next-Generation Sequencing](#)  
[Technical Support](#) | [Customer Service](#)

Catalog Number	Unit Size	Price (EUR)	Availability ⓘ	Quantity
✓ 4456740	1 kit	Price: 1.392,00 Online Offer: 1.268,65 ⓘ (ends 31-Dec-2023) Your Price: <a href="#">Sign In ⓘ</a>		<input type="text"/>



# Existing tools - qsmooth

Qsmooth quantile-normalizes the counts between the reference quantile of biological replicates or towards the overall reference quantile, depending on variability within and between conditions.

The screenshot shows the top portion of a journal article page. At the top, the journal title "Biostatistics" is displayed in a large, white serif font against a blue background. Below this is a dark blue navigation bar with white text for "Issues", "Advance articles", "Submit" (with a dropdown arrow), "Purchase", "Alerts", and "About" (with a dropdown arrow). The main content area has a light blue background. It starts with an "Article Navigation" section. Below that is a "JOURNAL ARTICLE" label. The article title "Smooth quantile normalization" is in a bold, black serif font, followed by a green "FREE" badge. The authors' names "Stephanie C Hicks, Kwame Okrah, Joseph N Paulson, John Quackenbush, Rafael A Irizarry, Héctor Corrada Bravo" are listed in a smaller blue font. Below the authors is the journal information: "Biostatistics, Volume 19, Issue 2, April 2018, Pages 185–198, <https://doi.org/10.1093/biostatistics/kxx028>". The publication date "Published: 10 July 2017" and a link to "Article history" are also present. At the bottom of this section are icons for "PDF", "Split View", "Cite", "Permissions", and "Share". The "SUMMARY" section begins with the text: "Between-sample normalization is a critical step in genomic data analysis to remove systematic bias and unwanted variation in high-throughput data. Global normalization methods are based on the assumption that observed global properties is due to technical reasons and are unrelated to the biology of interest. For example, some me".

# Existing tools - moose<sup>2</sup>

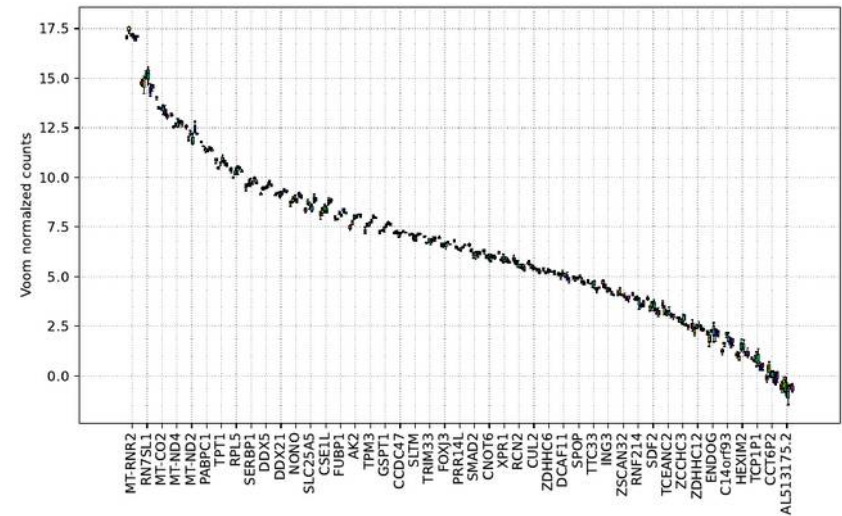
Research | [Open access](#) | [Published: 05 September 2017](#)

## RNA-sequence data normalization through in silico prediction of reference genes: the bacterial response to DNA damage as case study

[Bork A. Berghoff](#), [Torgny Karlsson](#), [Thomas Källman](#), [E. Gerhart H. Wagner](#) & [Manfred G. Grabherr](#) 

[BioData Mining](#) **10**, Article number: 30 (2017) | [Cite this article](#)

**4126** Accesses | **13** Citations | **5** Altmetric | [Metrics](#)



# Benchmark against spiked datasets with transcriptional shifts

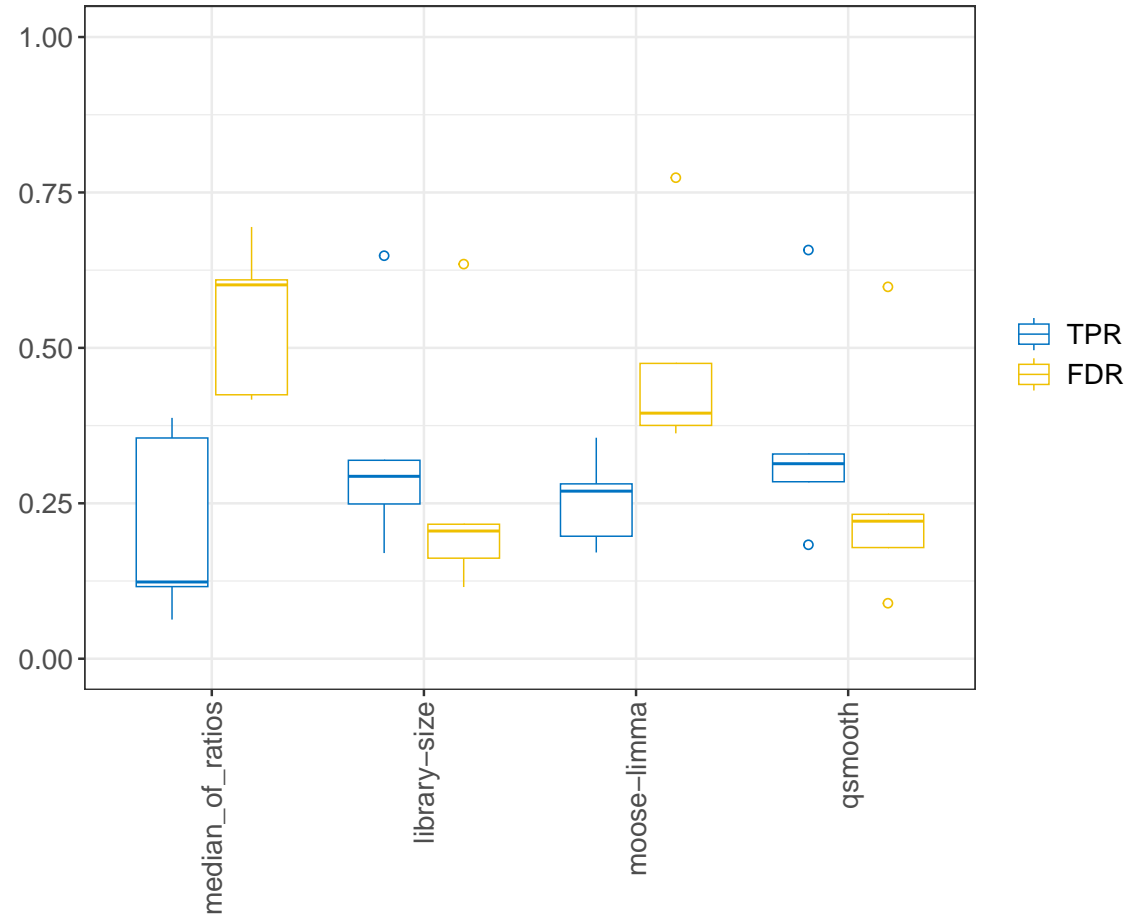
- Buschle et al. (2021) examined Raji-cells expressing the Epstein-Barr virus (**EBV**)
- Lau et al. (2023) found a large downregulation in topoisomerase I (**TOP1**) overexpressing Human embryonic kidney 293 cells. TOP1 is known for its role in relieving DNA supercoils for enabling transcription
- Bruno et al. (2020) found that **Che-1** depletion induces a global transcription shutoff by reducing histone acetylation. Che-1 is an interactor of RNA polymerase II

*In all publications, cells were actually **counted** (many other authors don't, but use spikes as control for technical bias).*

1. Buschle, A., Mrozek-Gorska, P., Cernilogar, F.M., Ettinger, A., Pich, D., Krebs, S., Mocanu, B., Blum, H., Schotta, G., Straub, T., et al. (2021). Epstein-Barr virus inactivates the transcriptome and disrupts the chromatin architecture of its host cell in the first phase of lytic reactivation. *Nucleic Acids Res* 49, 3217–3241. [10.1093/nar/gkab099](https://doi.org/10.1093/nar/gkab099).
2. Lau, M.S., Hu, Z., Zhao, X., Tan, Y.S., Liu, J., Huang, H., Yeo, C.J., Leong, H.F., Grinchuk, O.V., Chan, J.K., et al. (2023). Transcriptional repression by a secondary DNA binding surface of DNA topoisomerase I safeguards against hypertranscription. *Nat Commun* 14, 6464. [10.1038/s41467-023-42078-9](https://doi.org/10.1038/s41467-023-42078-9).
3. Bruno, T., De Nicola, F., Corleone, G., Catena, V., Goeman, F., Pallocca, M., Sorino, C., Bossi, G., Amadio, B., Cigliana, G., et al. (2020). Che-1/AATF-induced transcriptionally active chromatin promotes cell proliferation in multiple myeloma. *Blood Advances* 4, 5616–5630. [10.1182/bloodadvances.2020002566](https://doi.org/10.1182/bloodadvances.2020002566).

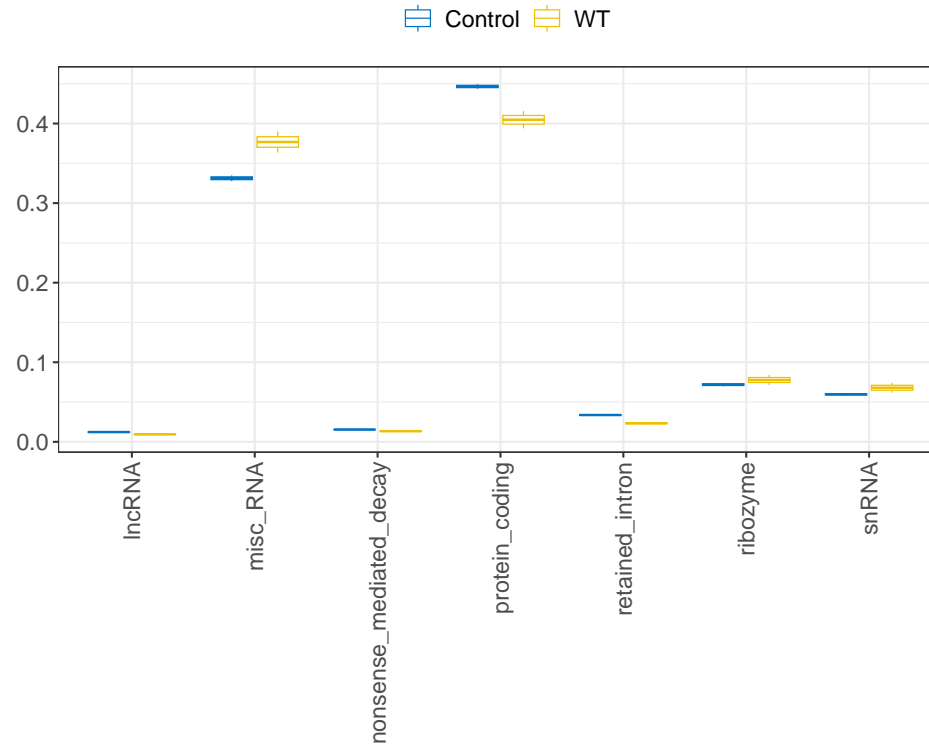
# DEGs detection performance by normalization method

Control vs treated  
compared to spike-ins  
( $p_{\text{adjust}} < 0.05$ )



# Fragment percentages by ensembl biotype

TOP1



- lncRNAs transcribed by Pol I, II or III <sup>1</sup>
- snRNAs transcribed by Pol II or III

Fragment := mapped read divided by exon length



## Normalization by biotype — basic assumptions

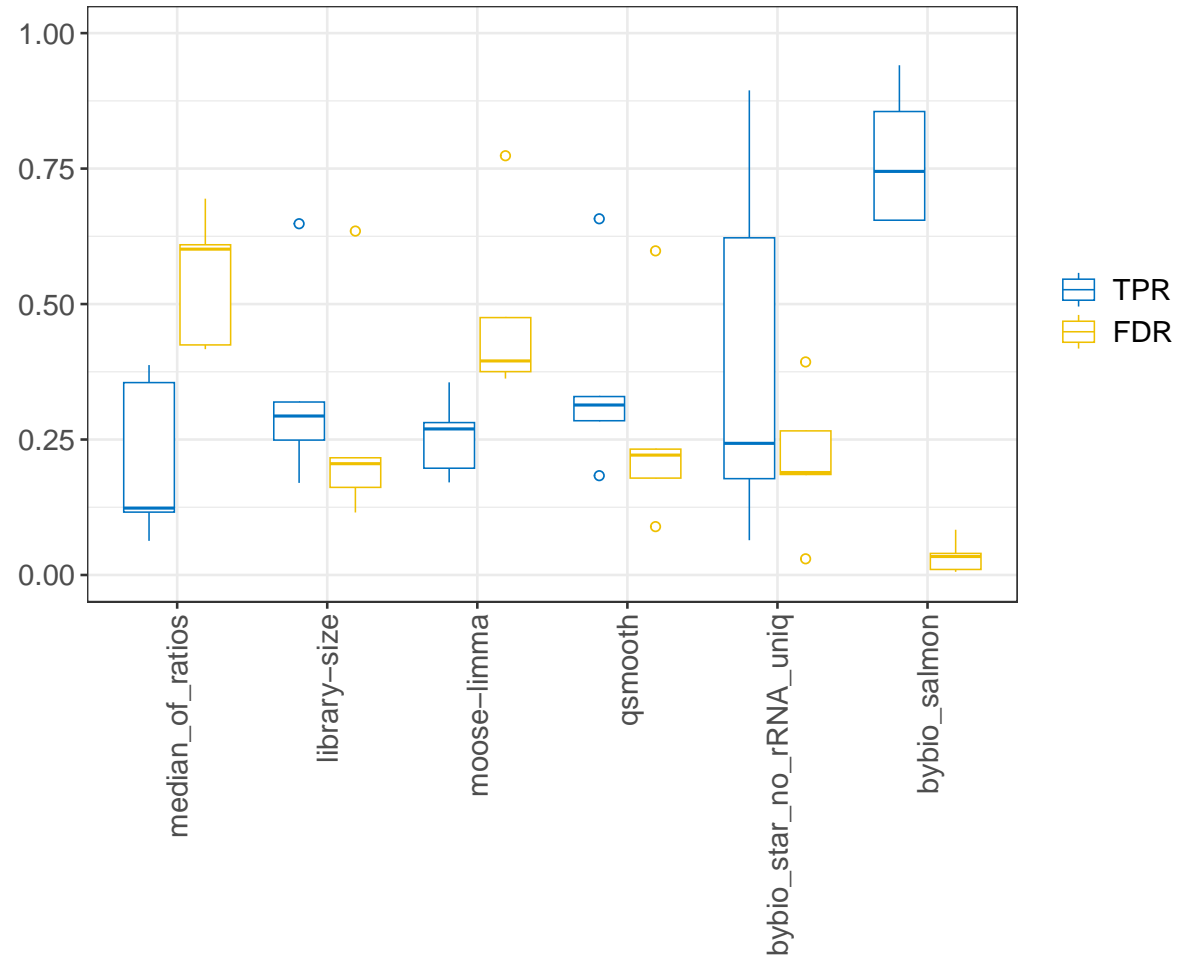
- **The total number of read ncRNA fragments should be *roughly* equal across conditions. If it's higher, in case of a shutdown: likely simply more cells were sequenced**
- At least ***some abundant*** ncRNAs are likely unaffected by the global shift
- ncRNA is presumably less affected by treatment, e.g. it's transcribed by POL I, II and III as opposed to cRNA (Pol II)
- some ncRNAs are quite abundant and stable, such as circular RNAs (mRNA-levels!)<sup>1</sup>
- a greater share of ncRNA is likely non-functional as compared to cRNA<sup>1</sup>

# DEGs detection performance by normalization method

Control vs treated  
compared to spike-ins

*bybio\_salmon*:

- mapped to whole transcriptome (including rRNA)
- salmon keeps multi-mapping reads
- transcript-level length used instead of, e.g., canonical length
- sizefactors based on the salmon-counts were applied to STAR uniquely mapped and rRNA-filtered counts



# Perspectives

- Test on more datasets (awaiting ethics approval for a large spiked “cancer dataset”)
- Test on polyA RNAseq (currently only rRNA depleted total-RNA)
- Check performance on spiked data without transcriptional shifts → preliminary benchmarks seem to perform well
- Evaluate implications of *slight* overall changes in global transcription levels.

# Acknowledgments

Thanks to the whole Hoffman  
group, especially

Konstantin

Robert

Steve

And thanks for your attention!



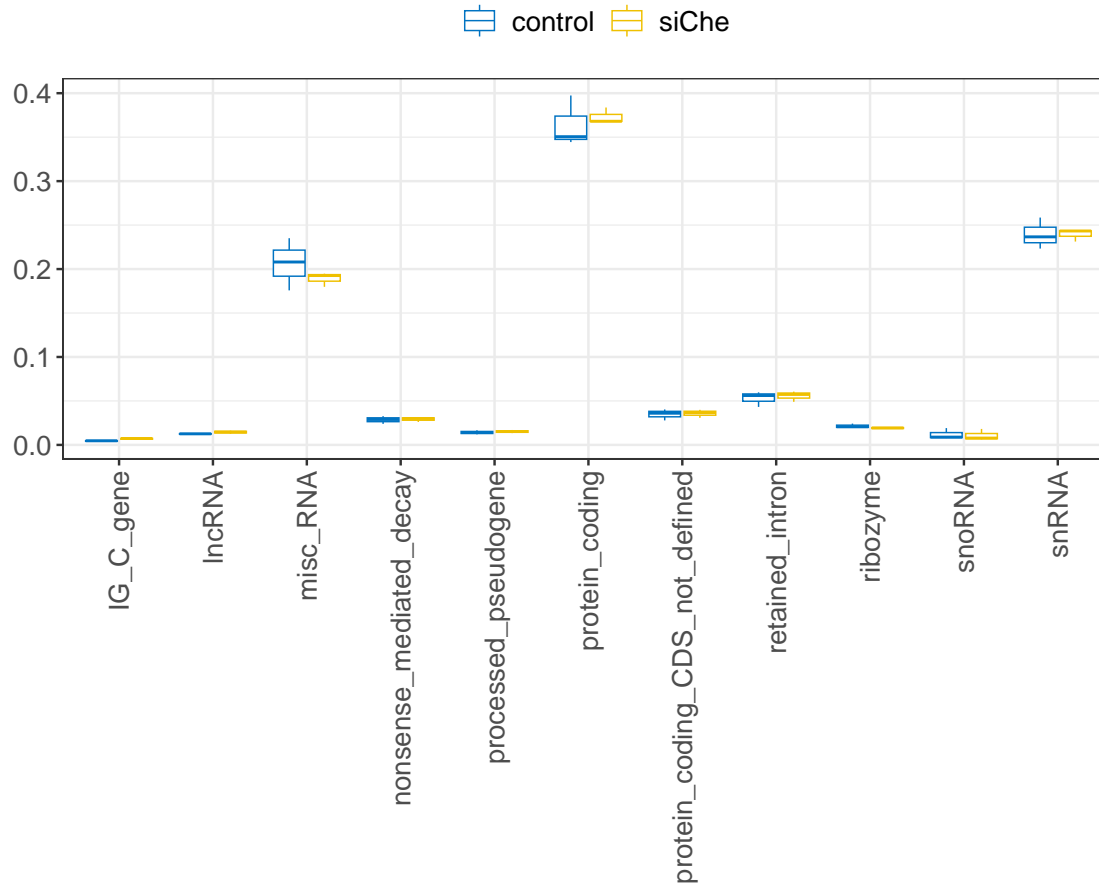
## Appendix – Normalization algorithm for **different** ncRNA fragmentcounts across conditions

- Normalize by transcript length and scale all samples to the same total number of fragments
- Ignore cRNAs (as well as pseudogenes and some other protein-associated biotypes)
- Remove the most variant ncRNA's **within** replicates. As the variance is dependent on the mean, these will be predominantly highly expressed genes which are already unstable in one condition
- Calculate a per-sample pseudo-reference by geometric mean of the remaining sum of fragments
- Calculate scaling factors for each sample to that pseudo-reference
- Apply the scaling factors to the whole sample (including cRNA's) and run DESeq2 with disabled sizefactor estimation

## Appendix – Normalization algorithm for **equal** ncRNA fragmentcounts across conditions

- If total ncRNA fragmentcounts are already similar across conditions, scaling by these makes no sense
- Instead, perform the same procedure as for the different ncRNA fragmentcounts on the slide before, but instead of a pseudo-reference by geometric mean, fall back to DESeq2's *median of ratios* method — but only for the ncRNAs.

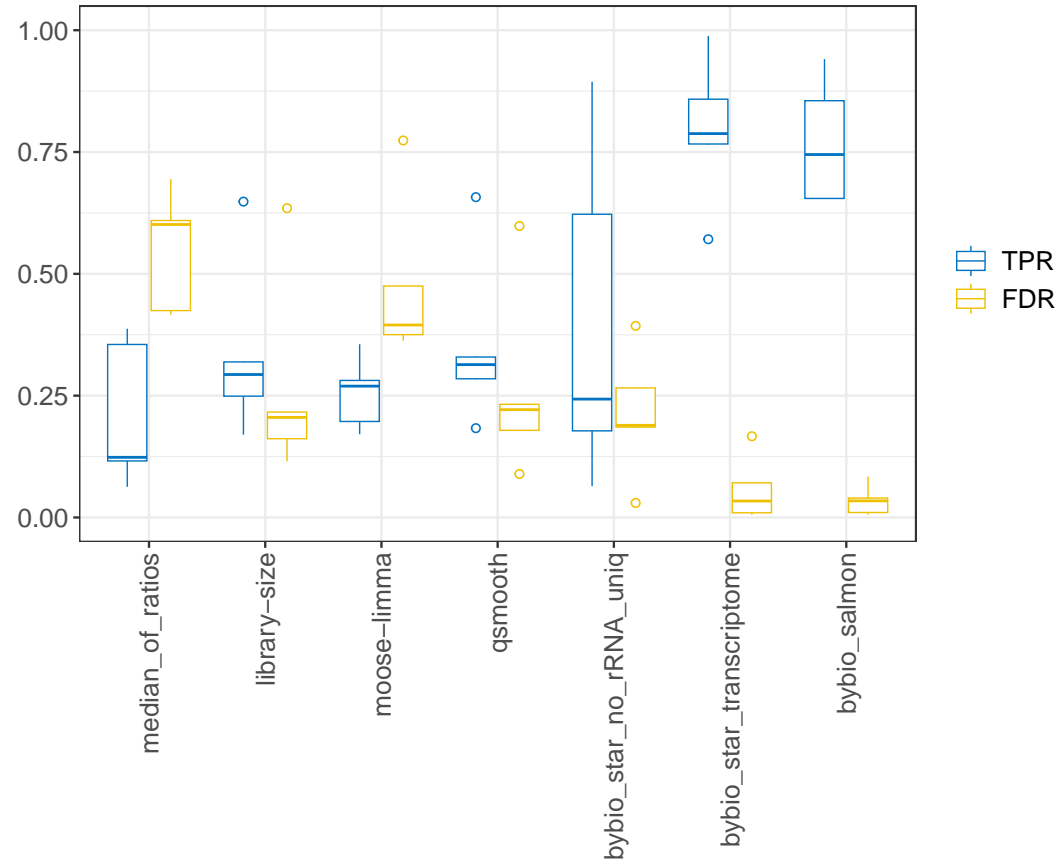
# Appendix – No global fragment differences in Che1



→ Thus, we fallback to  
*median of ratios*  
normalization of ncRNA



# Appendix – star on transcriptome with canonical transcript length



# Appendix – RNA molecules per cell

Type	Percent of total RNA by mass	Molecules per cell	Average size (kb)	Total weight picograms/cell	Notes	Reference
rRNAs	80 to 90	3–10 × 10 <sup>6</sup> (ribosomes)	6.9	10 to 30		Blobel and Potter (1967), Wolf and Schlessinger (1977), Duncan and Hershey (1983)
tRNA	10 to 15	3–10 × 10 <sup>7</sup>	<0.1	1.5 to 5	About 10 tRNA molecules /ribosome	Waldron and Lacroute (1975)
mRNA	3 to 7	3–10 × 10 <sup>5</sup>	1.7	0.25 to 0.9		Hastie and Bishop (1976), Carter et al. (2005)
hnRNA (pre-mRNA)	0.06 to 0.2	1–10 × 10 <sup>3</sup>	10*	0.004 to 0.03	Estimated at 2–4% of mRNA by weight	Mortazavi et al. (2008), Menet et al. (2012)
Circular RNA	0.002 to 0.03	3–20 × 10 <sup>3</sup>	~0.5	0.0007 to 0.005	Estimated at 0.1–0.2% of mRNA**	Salzman et al. (2012), Guo et al. (2014)
snRNA	0.02 to 0.3	1–5 × 10 <sup>5</sup>	0.1–0.2	0.008 to 0.04		Kiss and Filipowicz (1992), Castle et al. (2010)
snoRNA	0.04 to 0.2	2–3 × 10 <sup>5</sup>	0.2	0.02 to 0.03		Kiss and Filipowicz (1992), Cooper (2000), Castle et al. (2010)
miRNA	0.003 to 0.02	1–3 × 10 <sup>5</sup>	0.02	0.001 to 0.003	About 10 <sup>5</sup> molecules per 10 pg total RNA	Bissels et al. (2009)
7SL	0.01 to 0.2	3–20 × 10 <sup>4</sup>	0.3	0.005 to 0.03	About 1–2 SRP molecules/100 ribosomes	Raue et al. (2007), Castle et al. (2010)
Xist	0.0003 to 0.02	0.1–2 × 10 <sup>3</sup>	2.8	0.0001 to 0.003		Buzin et al. (1994), Castle et al. (2010)
Other lncRNA	0.03 to 0.2	3–50 × 10 <sup>3</sup>	1	0.002 to 0.03	Estimated at 1–4% of mRNA by weight	Mortazavi et al. (2008), Ramsköld et al. (2009), Menet et al. (2012)

\*The size for the average unspliced pre-mRNA is 17 kb; however, most pre-mRNAs are partially spliced at any given time, and the average size of hnRNA is estimated at 10 kb (Salditt-Georgieff et al., 1976).

\*\*Based on the finding that 1–2% of all mRNA species generate circular RNA, which is present at 10% of the level of the parental mRNA.

## Appendix – qPCR as confirmation?

### OAS/RNASE L: SENSING VIRAL PAMP TRIGGERS GLOBAL RNA DEGRADATION AND TRANSLATIONAL ARREST

Degrading viral genomes presents one potent method of antiviral activity; digesting viral genetic material ensures that no further steps in replication can occur. However, the challenge lies within being able to control RNA degradation to ensure cellular survival or limit destruction within the host. The OAS/RNase L pathway is activated upon sensing the PAMP of dsRNA, serving two functions: sensing viral intruders and inhibiting viral replication by degrading RNA almost indiscriminately, inducing global translational arrest.

**Note:** Pinkham et al. “confirmed” some of the RNAseq results using the  $\Delta\Delta C_t$  method and 18 S ribosomal RNA!

Activation of RNase L quickly arrests global translation. This rapid translational arrest is traditionally attributed to degradation of transcripts involved in host translation machinery, as evidenced by degradation of 28S and 18S rRNA upon RNase L activation. However, closer examination revealed

1. Yang, E., and Li, M. (2020). All About the RNA: Interferon-Stimulated Genes That Interfere With Viral RNA Processes. *Frontiers in Immunology* 11. [10.3389/fimmu.2020.605024](https://doi.org/10.3389/fimmu.2020.605024).