

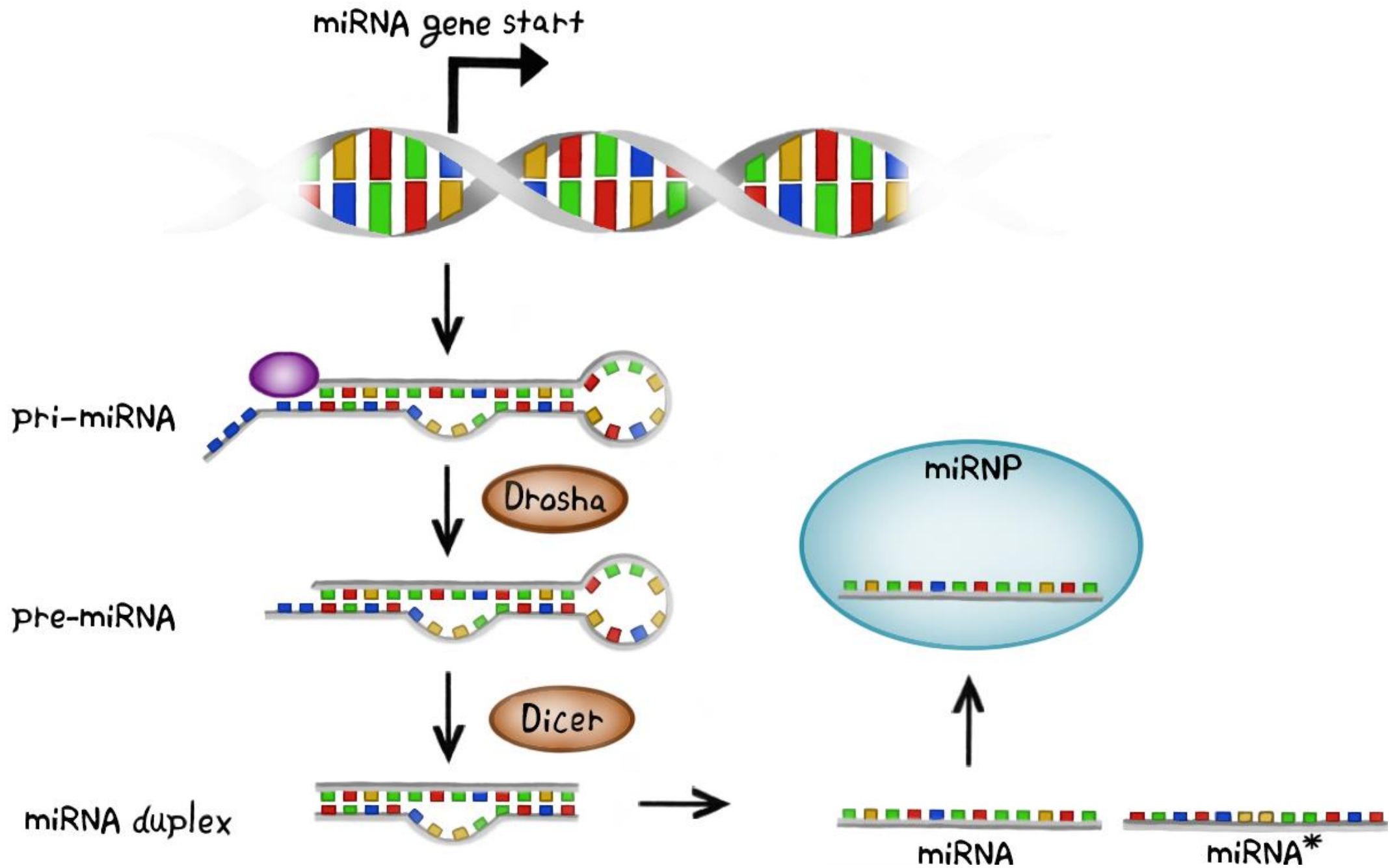
A new statistical approach to identify differential expression in small RNA-Seq data

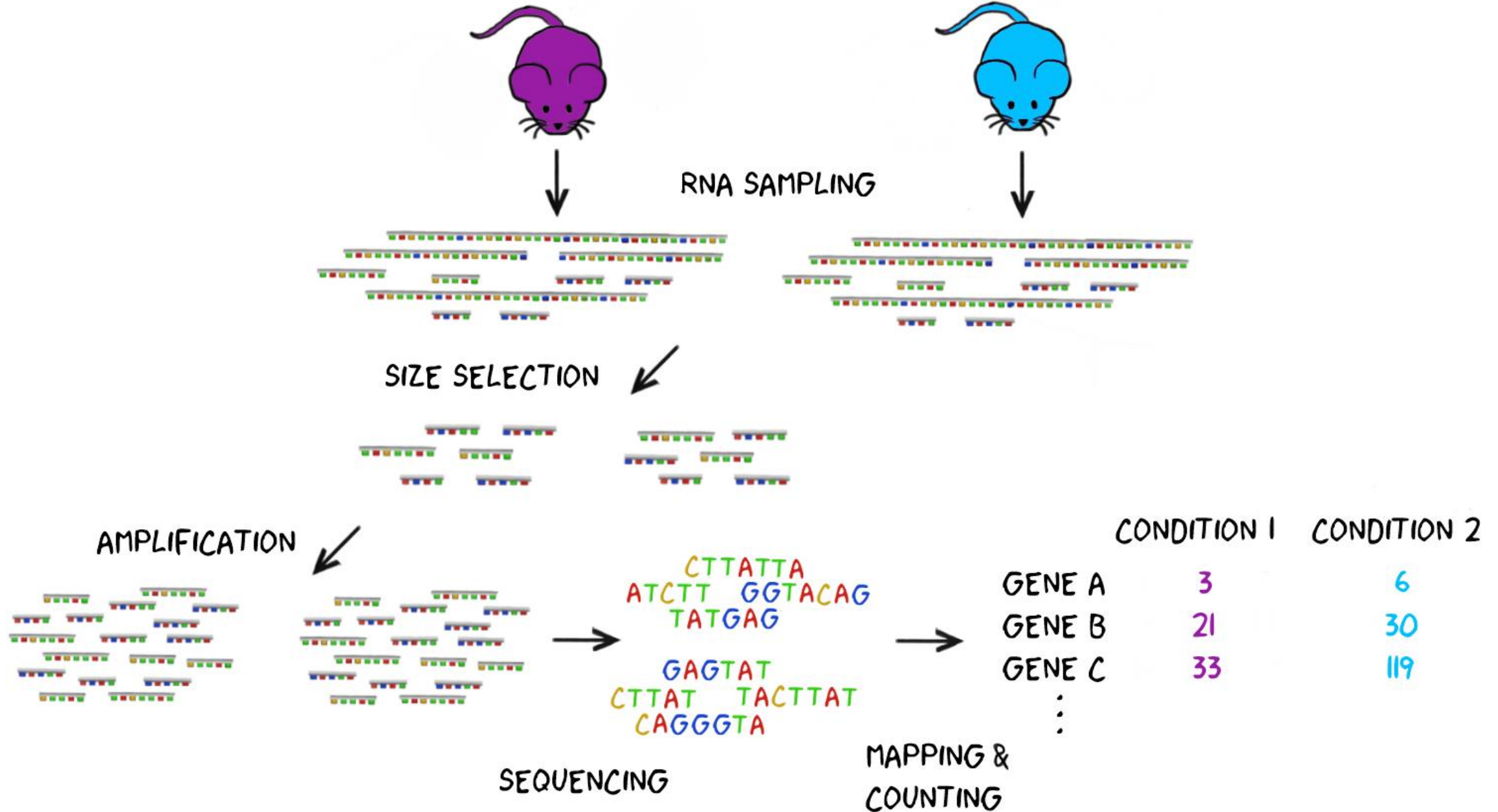
32nd TBI Winterseminar, 2017

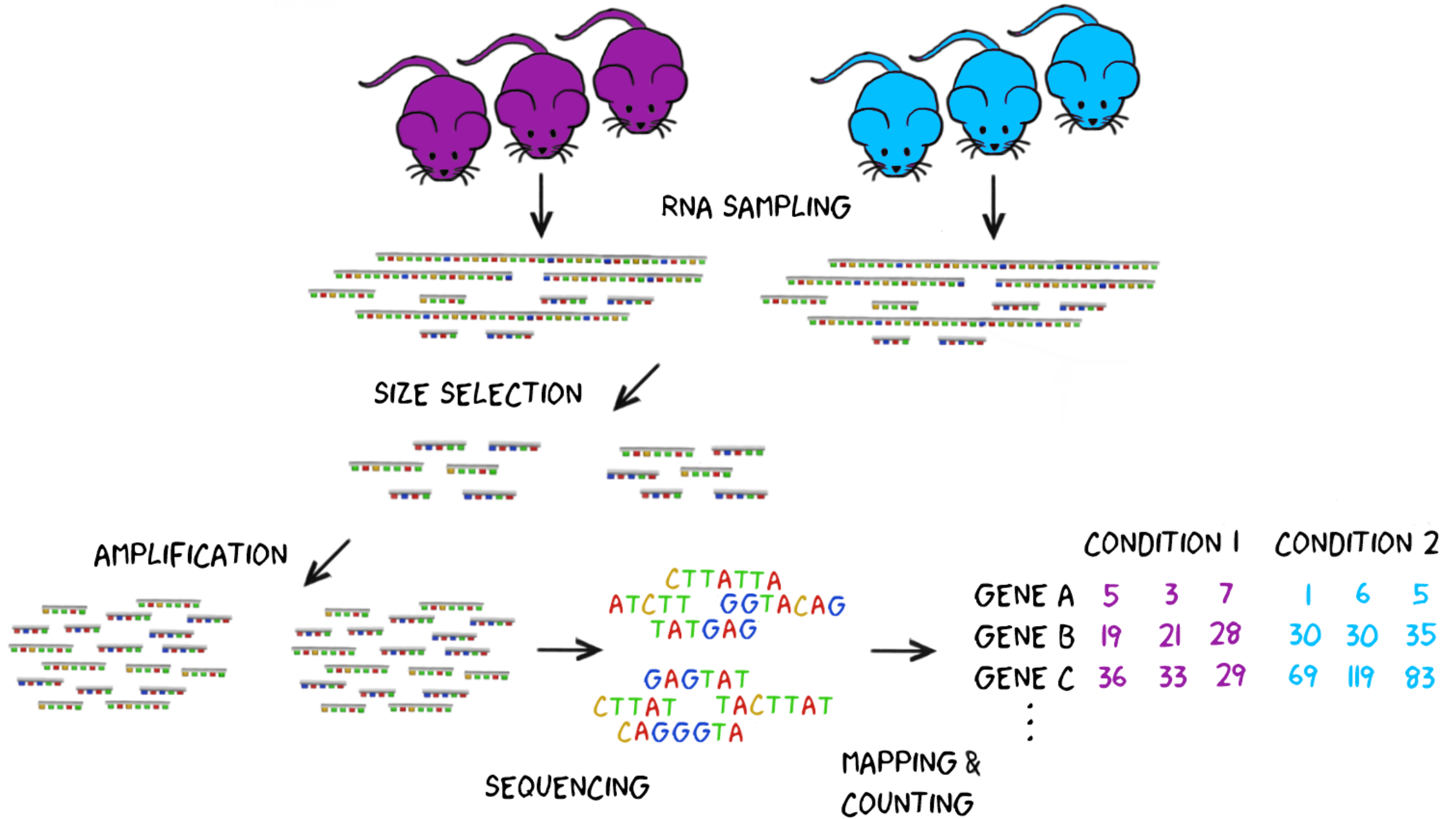
Emanuel Barth



seit 1558







SCIENTIFIC REPORTS

OPEN

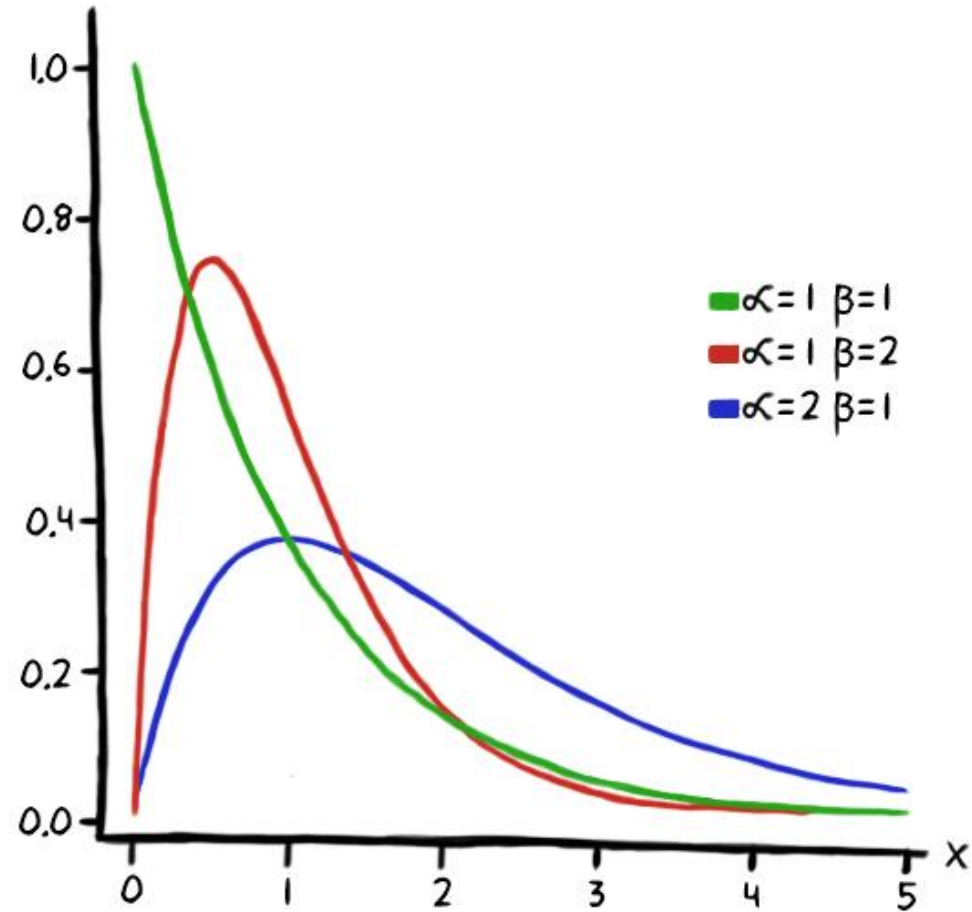
PUBLISHED
07 APRIL 2016

Empirical insights into the
stochasticity of small RNA
sequencing

Li-Xuan Qin, Thomas Tuschl and Samuel Singer

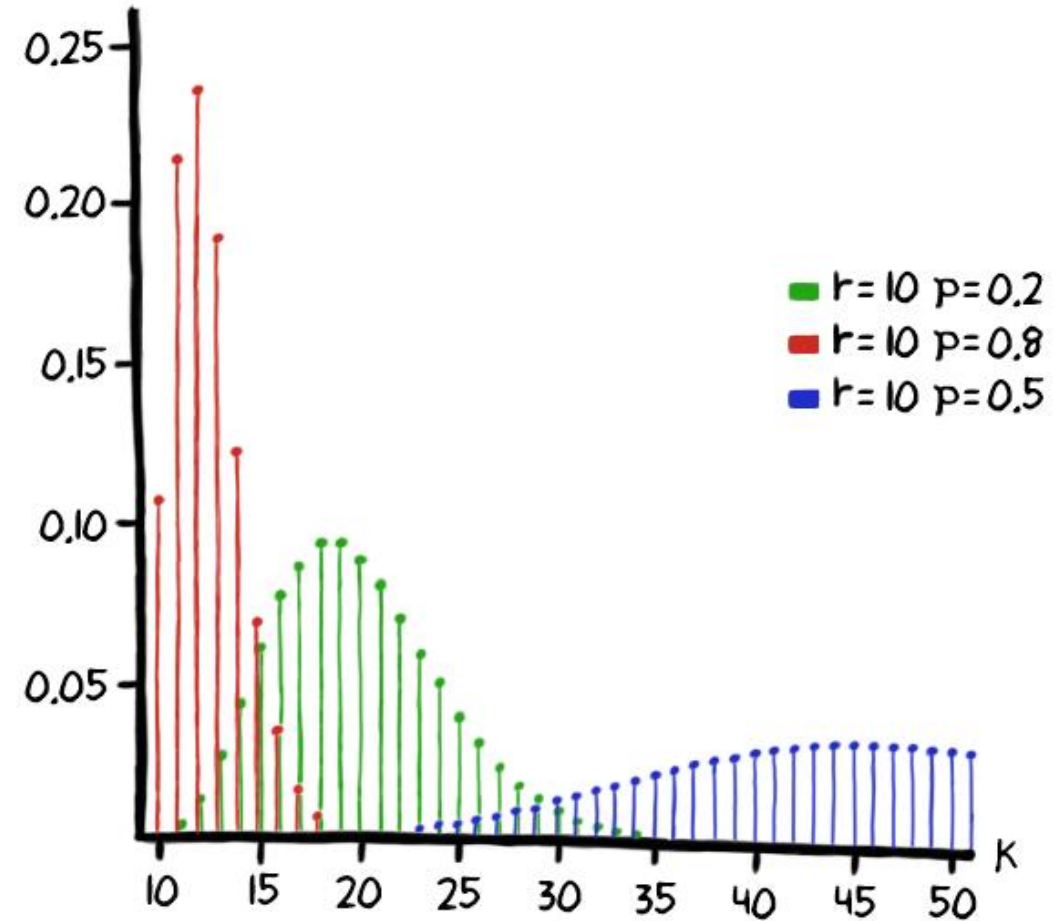
- Reads coming from small RNA-Seq are more likely to follow a Gamma distribution rather than a Poisson distribution
- Through a simple transformation of the count data, we can make use of much simpler statistical models for differential expression analysis

GAMMA



$$f(x) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} & x > 0 \\ 0 & x \leq 0 \end{cases}$$

NEGATIVE BINOMIAL



$$f(k) = \sum_{k=r}^{\infty} \binom{k-1}{r-1} p^r (1-p)^{k-r}$$

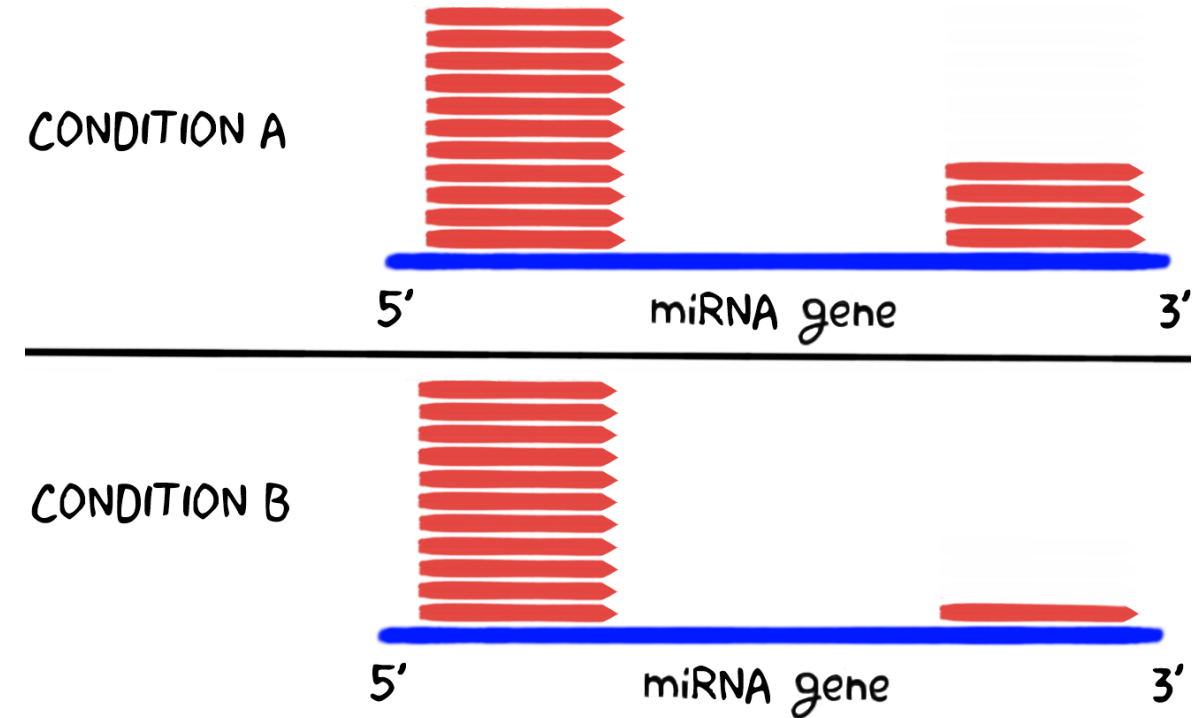
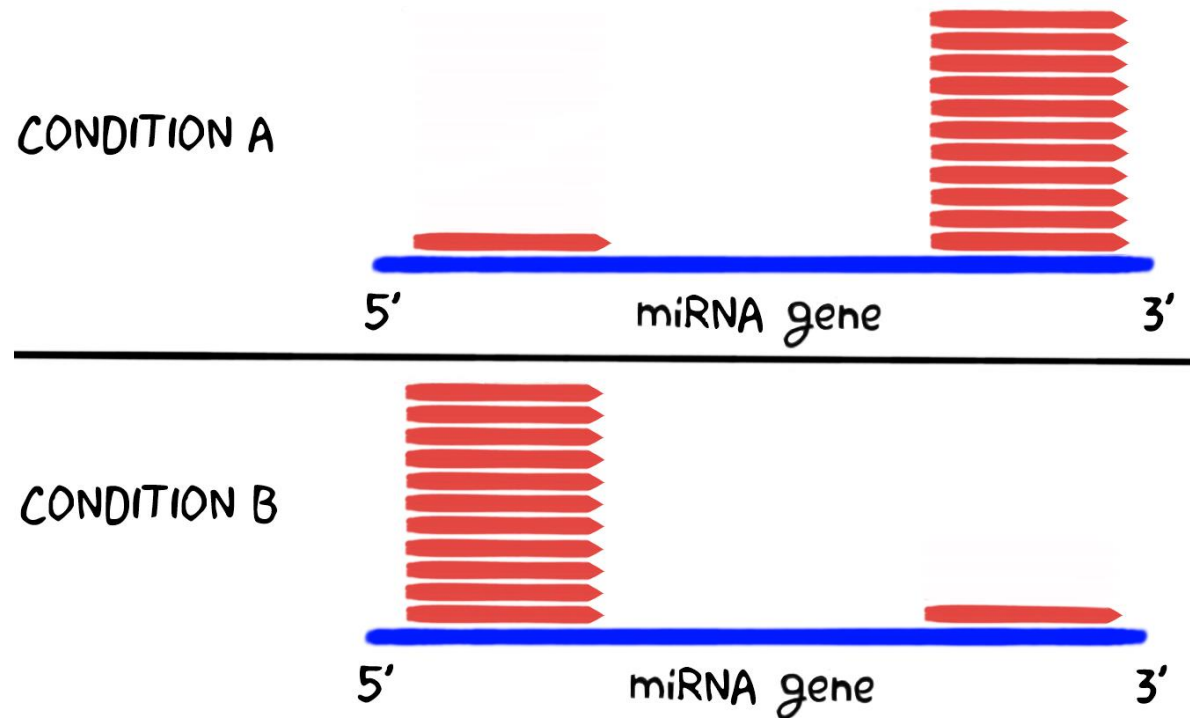
MicroRNA Differential Expression analysis (MeRDE)

- Input: sam/bam mapping files or raw count file
- Output: differentially expressed miRNAs (and some nice pictures)
- Count normalization based on pseudo-reference scaling¹
- Count transformation using the cubic root function
- Testing for differential expression with Welch's t-test

¹Anders et al. *Differential expression analysis for sequence count data*, G Bio, 2010, 11:R106

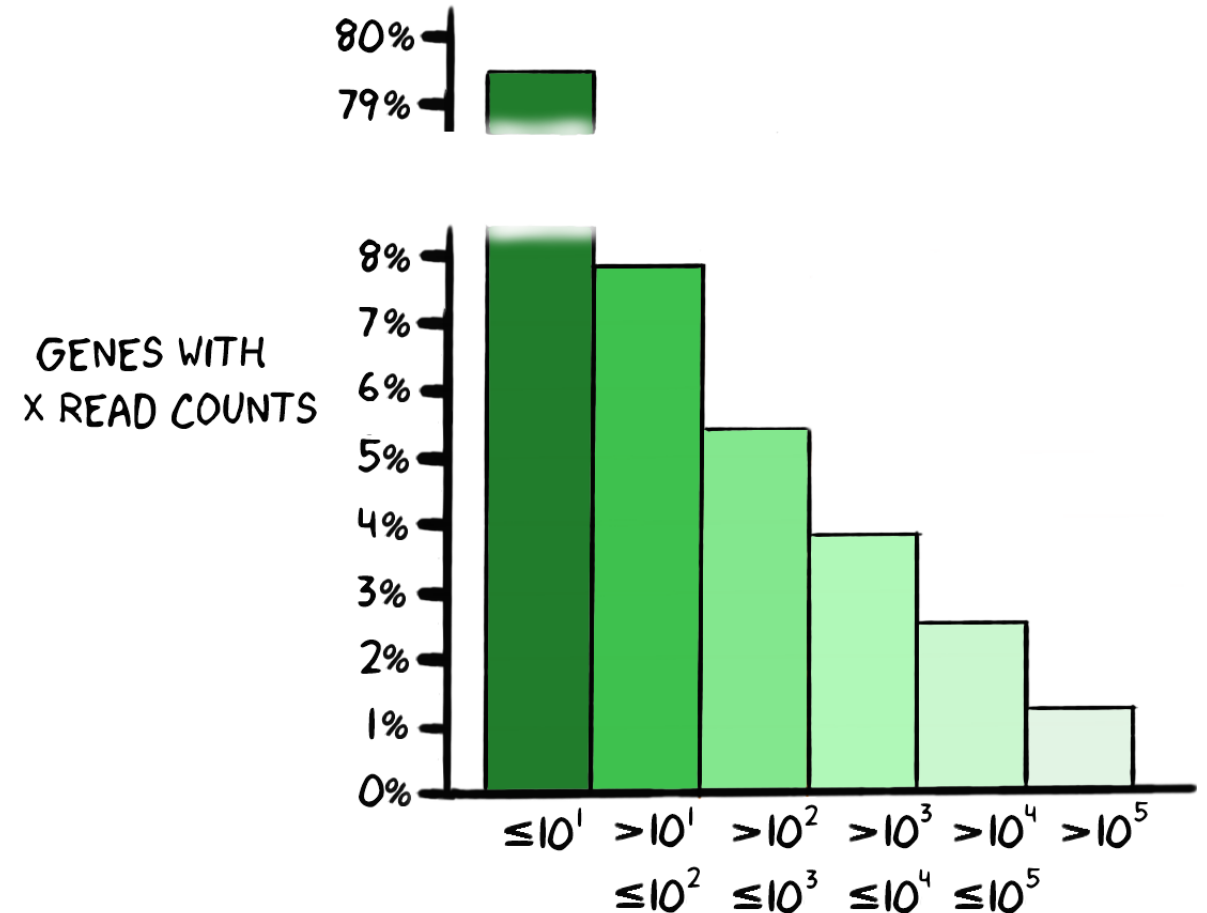
MicroRNA Differential Expression analysis (MeRDE)

- Treats 5'- and 3'-ends independently

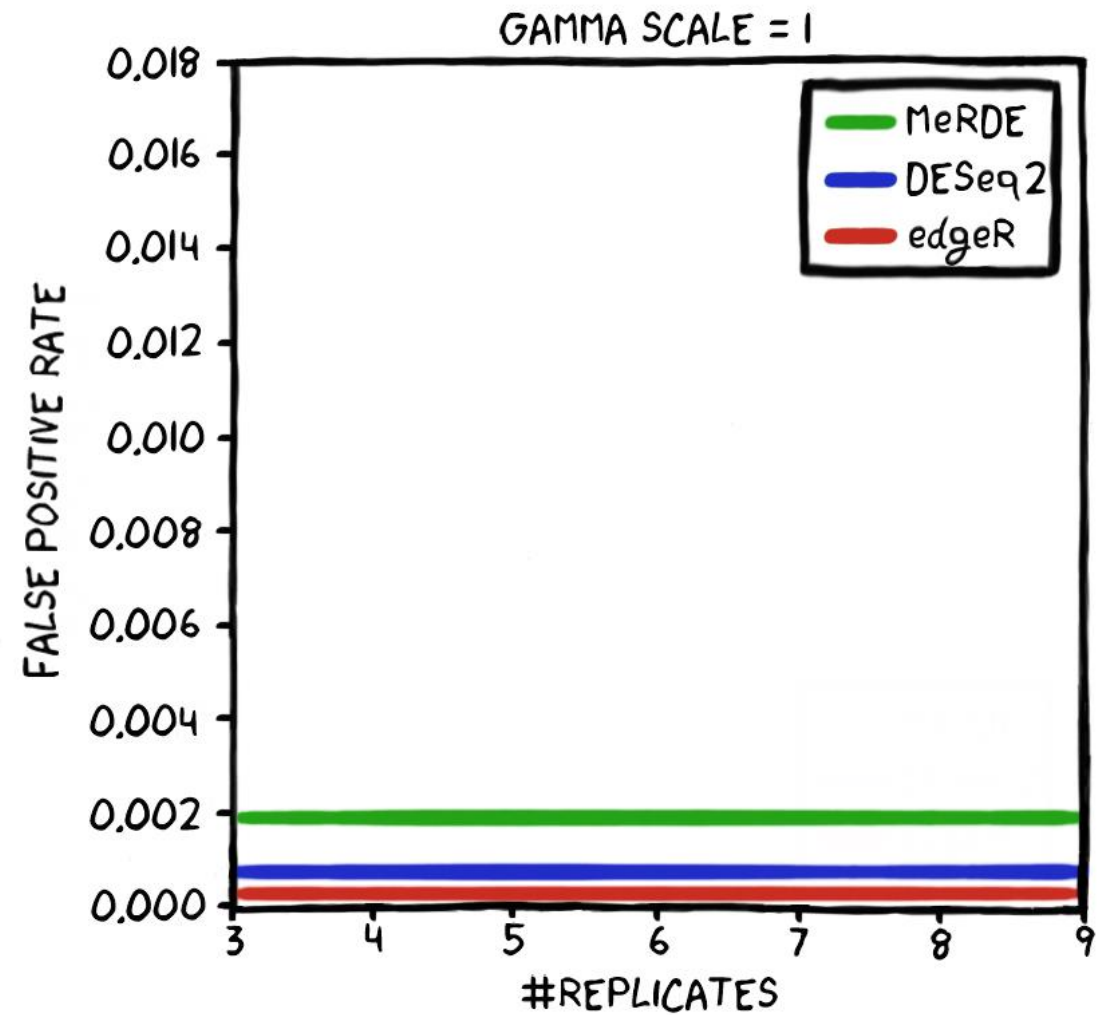
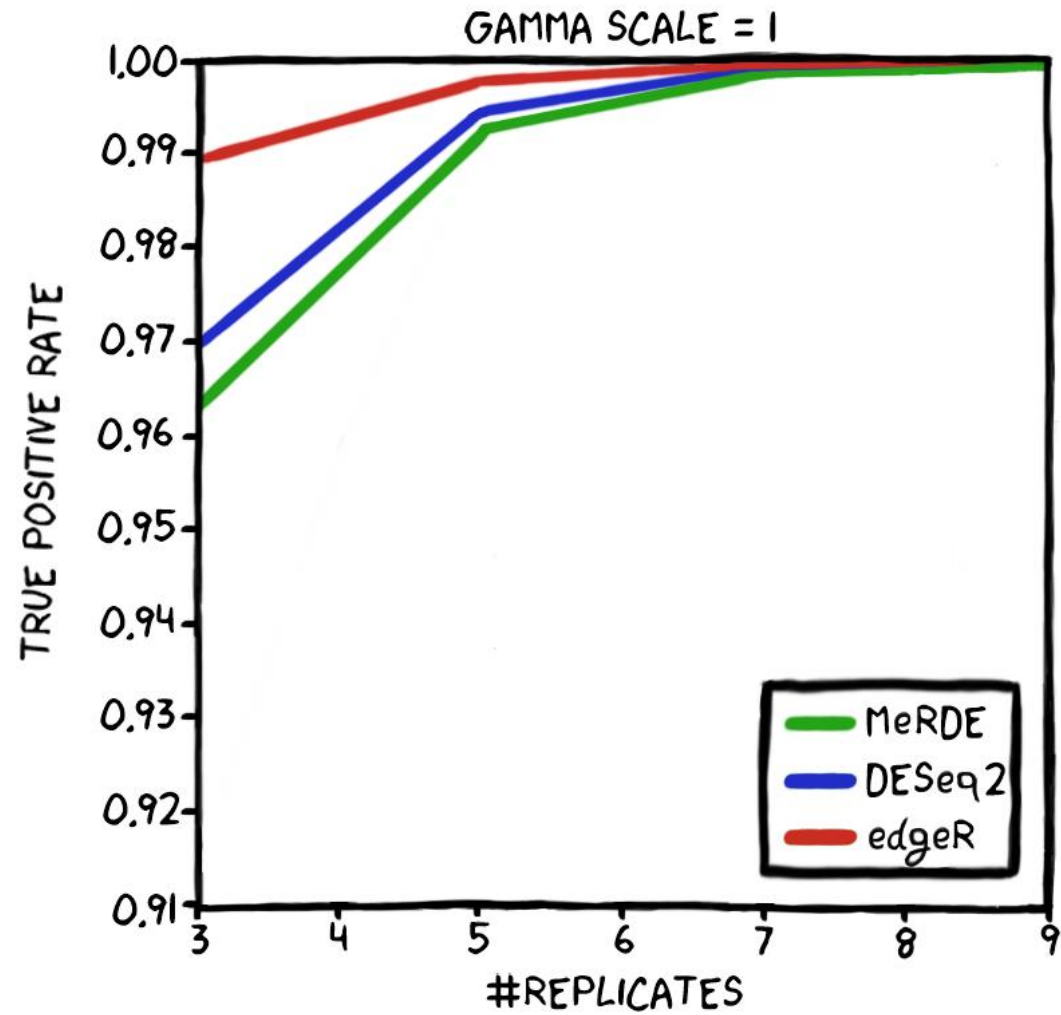


- small RNA-Seq data sets of three different species of different tissues between different ages

- Human: 60 Datasets
- Mouse: 74 Datasets
- *Notho. furzeri*: 160 Datasets

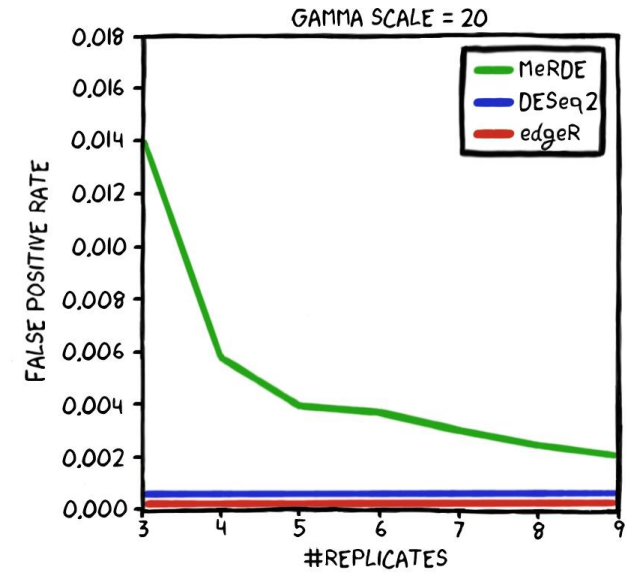
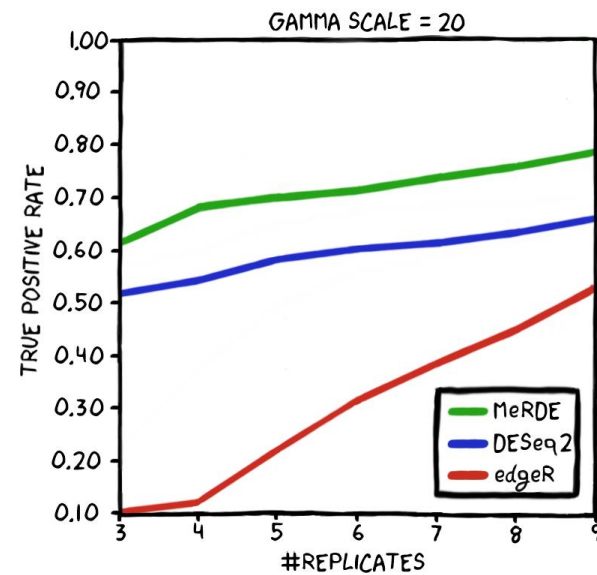
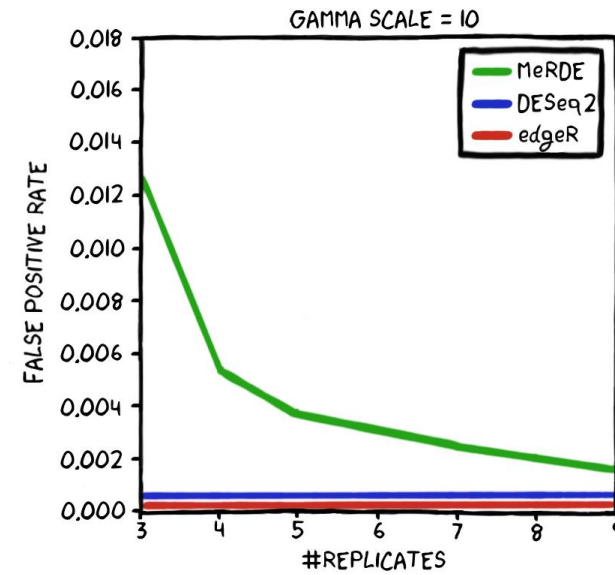
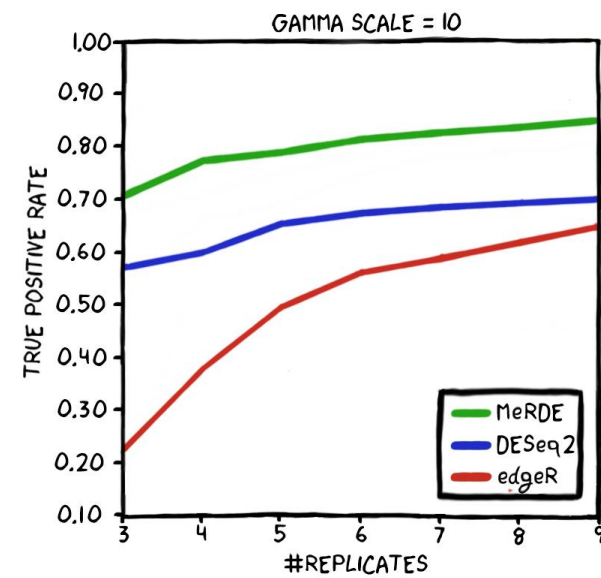
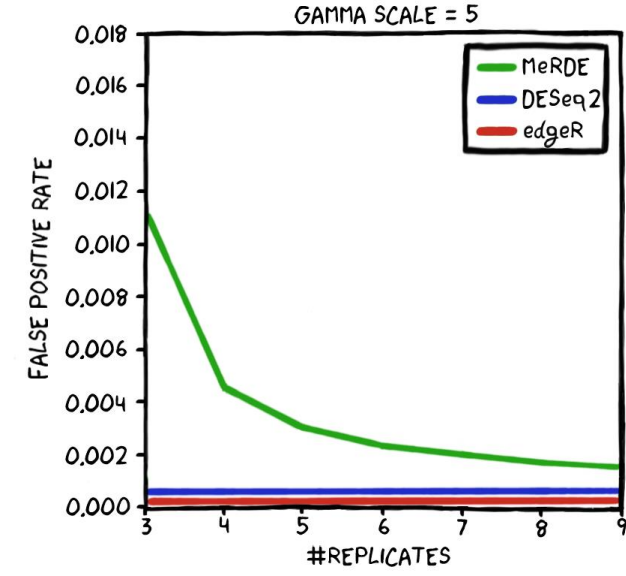
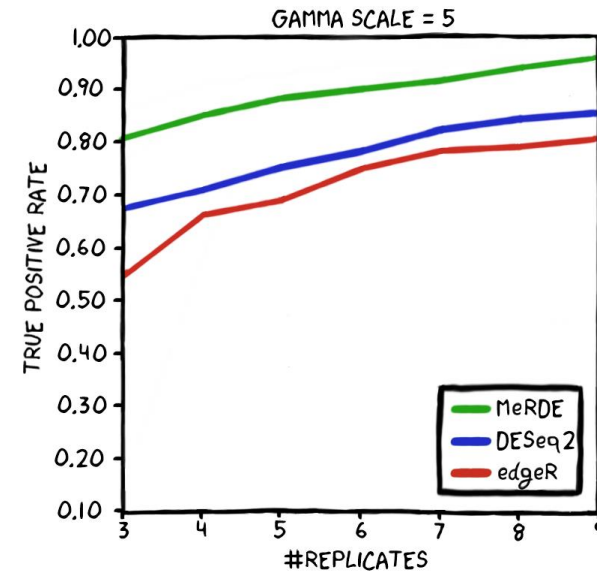
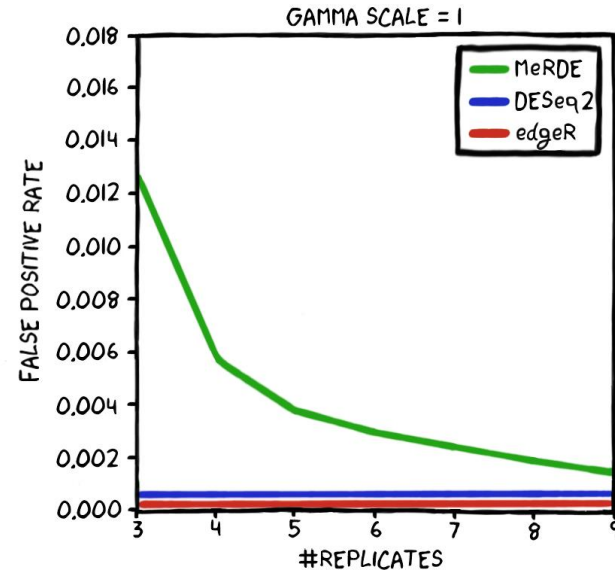
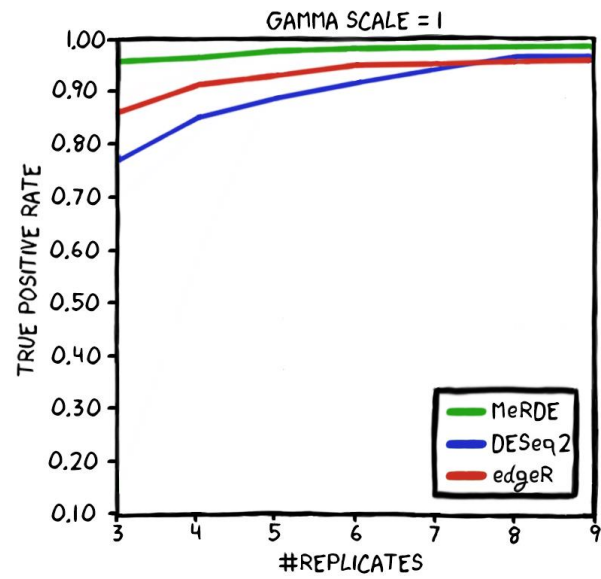


- Artificial data setup:
 - Gamma scaling factor $\beta = 1$
 - two conditions with 3 to 9 replicates
 - 1000 count files each containing 5000 simulated miRNAs
 - 1% - 5% of the genes are DE
 - DE factor between 2-fold and 5-fold up-/downregulated



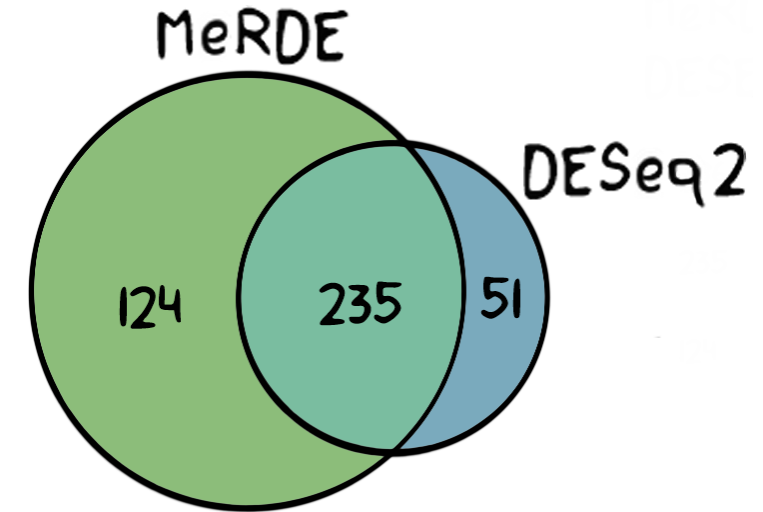
- Improved artificial data setup:
 - Gamma scaling factor $\beta = 1, 5, 10, 20$
 - two conditions with 3 to 9 replicates
 - 1000 count files each containing 5000 miRNAs
 - 1% - 5% of the genes are DE
 - DE factor between $f(x)$ and 2-fold up-/downregulated
 - outlier rate of 2%

$$f(x) = \frac{1.322}{\log_{4.8} x} + 1.1 \text{ where } x \text{ is the base mean expression}$$

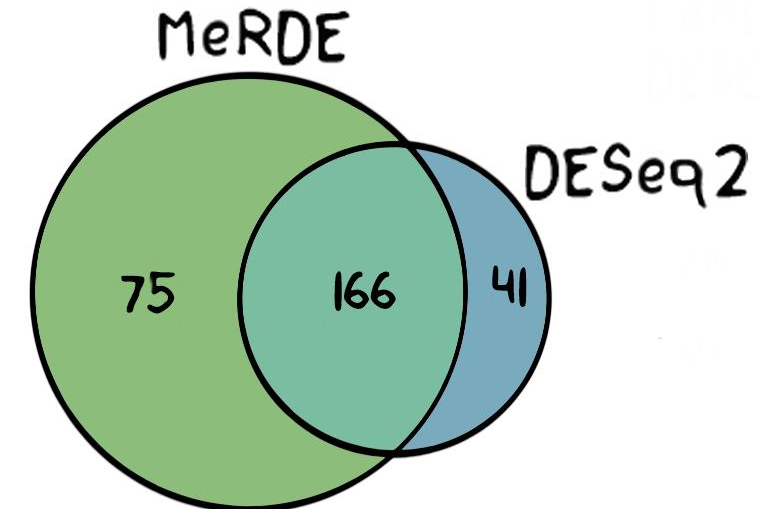


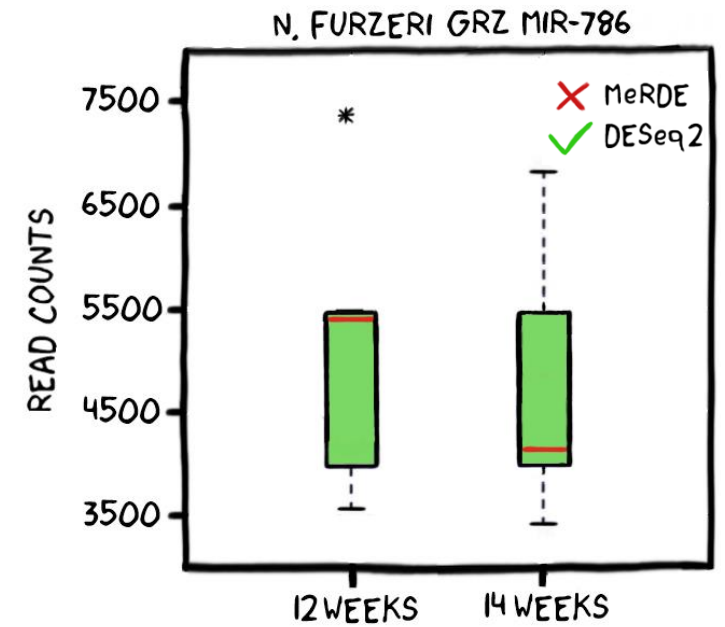
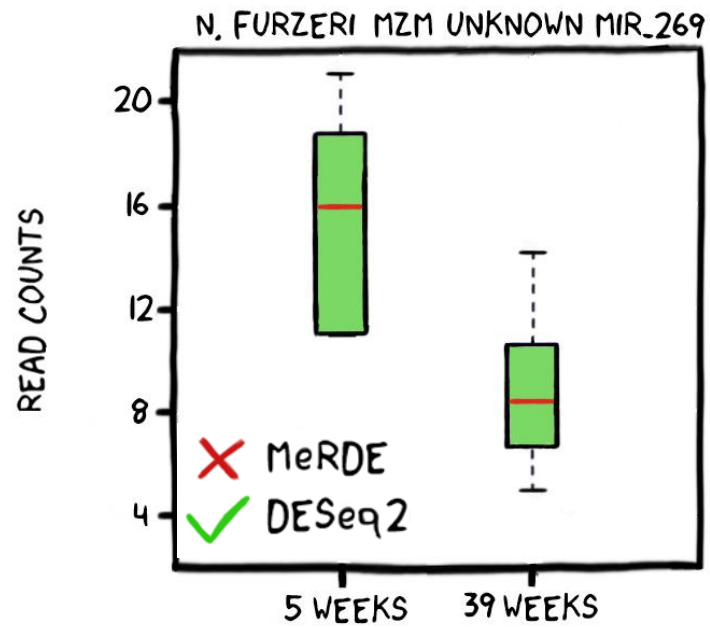
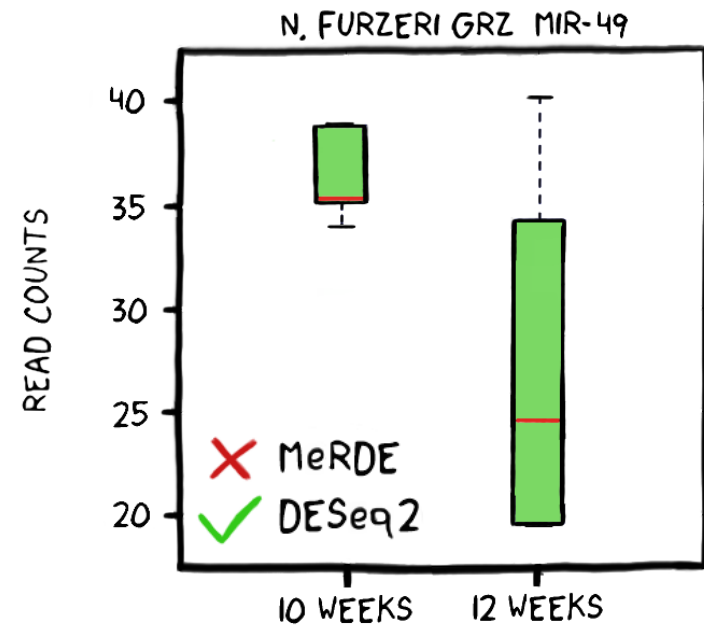
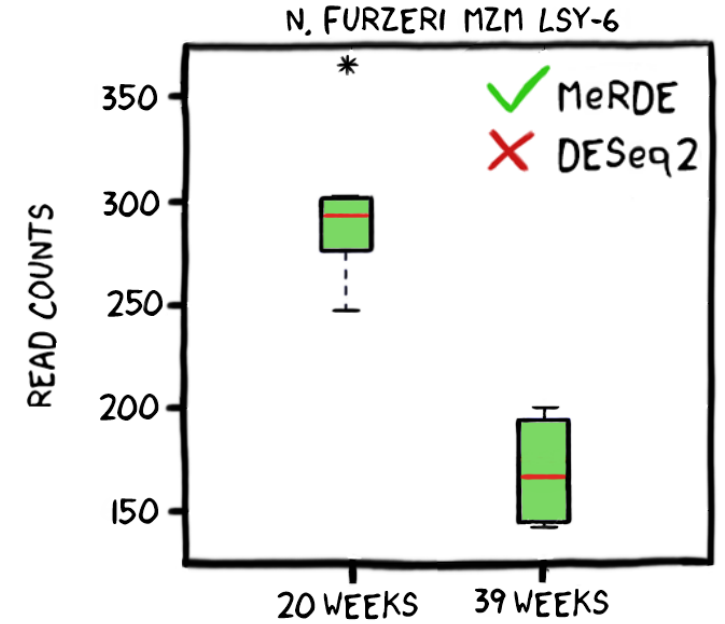
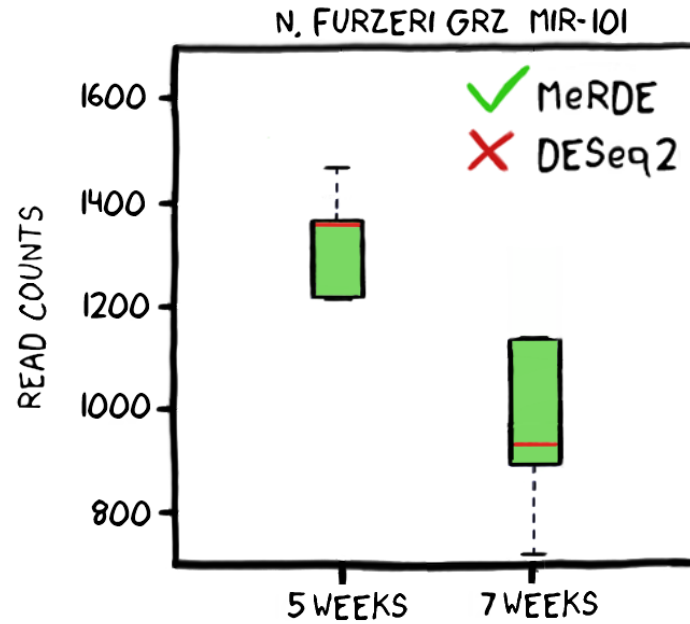
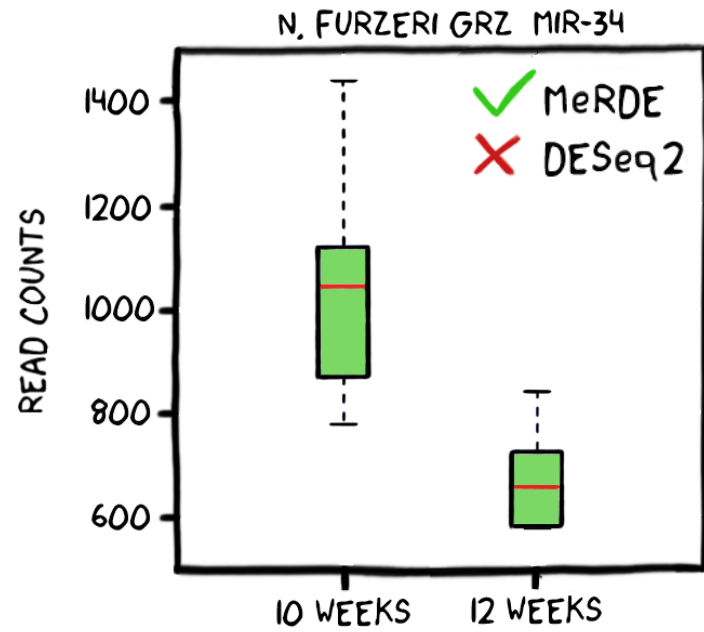
Real data results:

- Pairwise age comparison of two short-lived fish
 - *Nothobranchius furzeri* GRZ: 5, 7, 10, 12 and 14 weeks of age
 - $\approx 20\%$ more DEGs compared to DESeq2



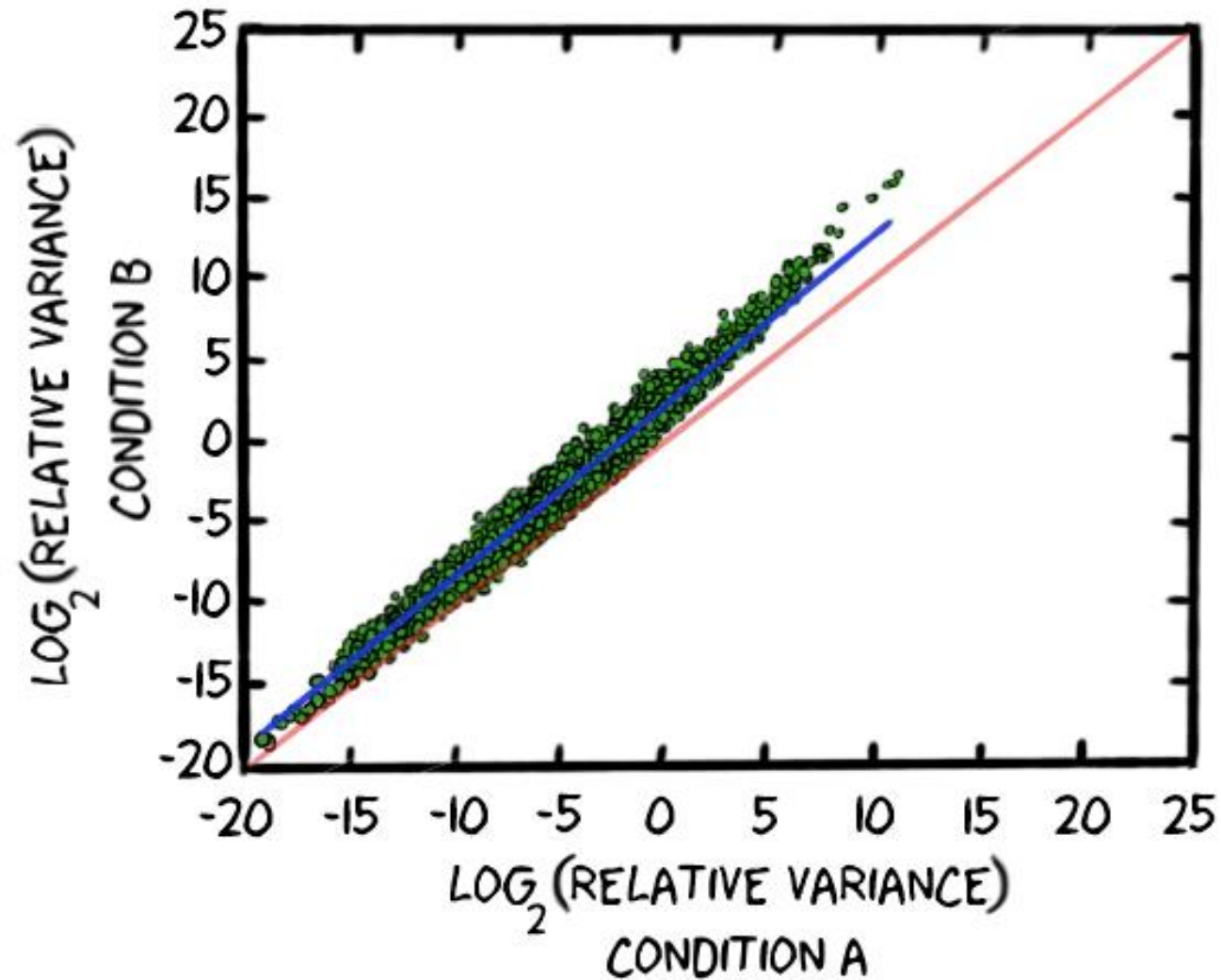
- *Nothobranchius furzeri* MZM 5, 12, 20, 27 and 39 weeks of age
- $\approx 14\%$ more DEGs compared to DESeq2





Still to do:

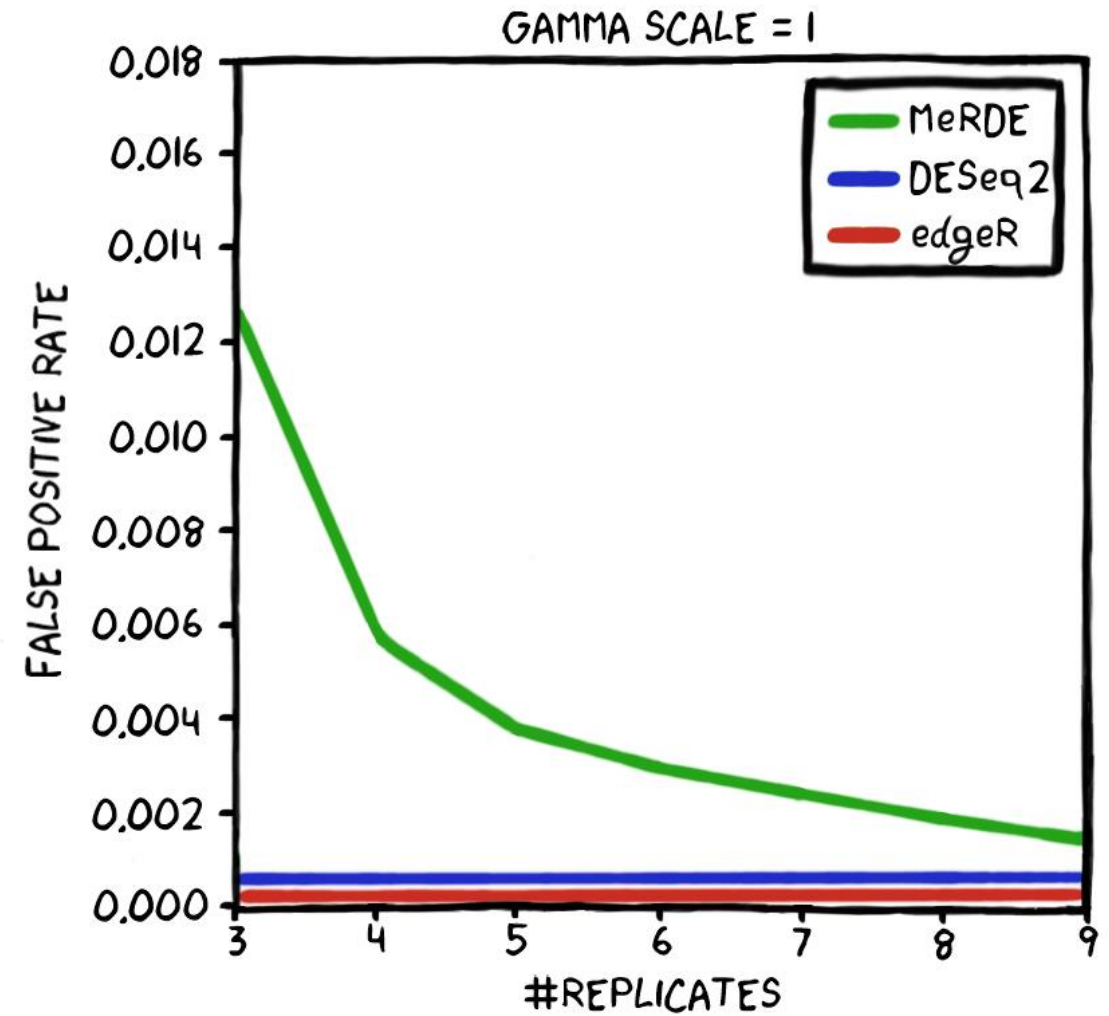
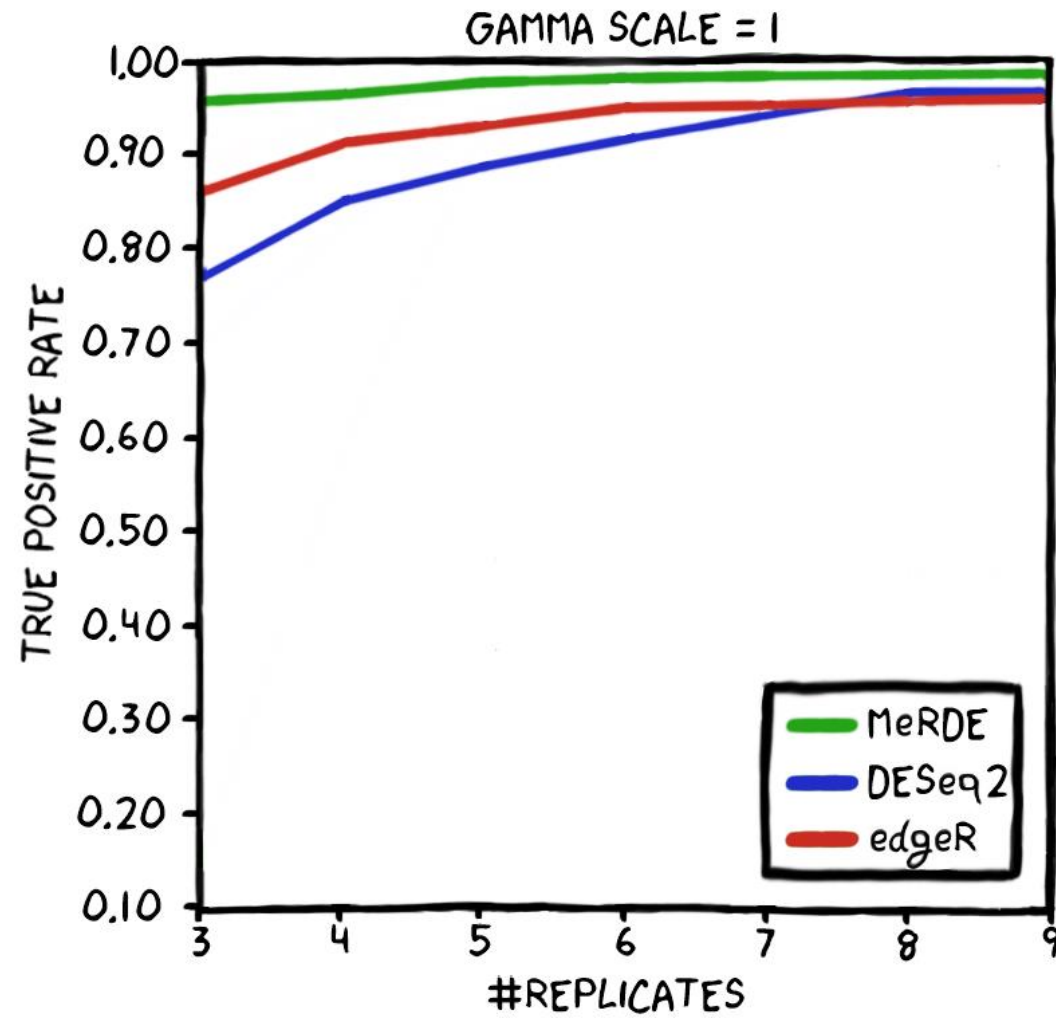
- Detecting and handling outliers

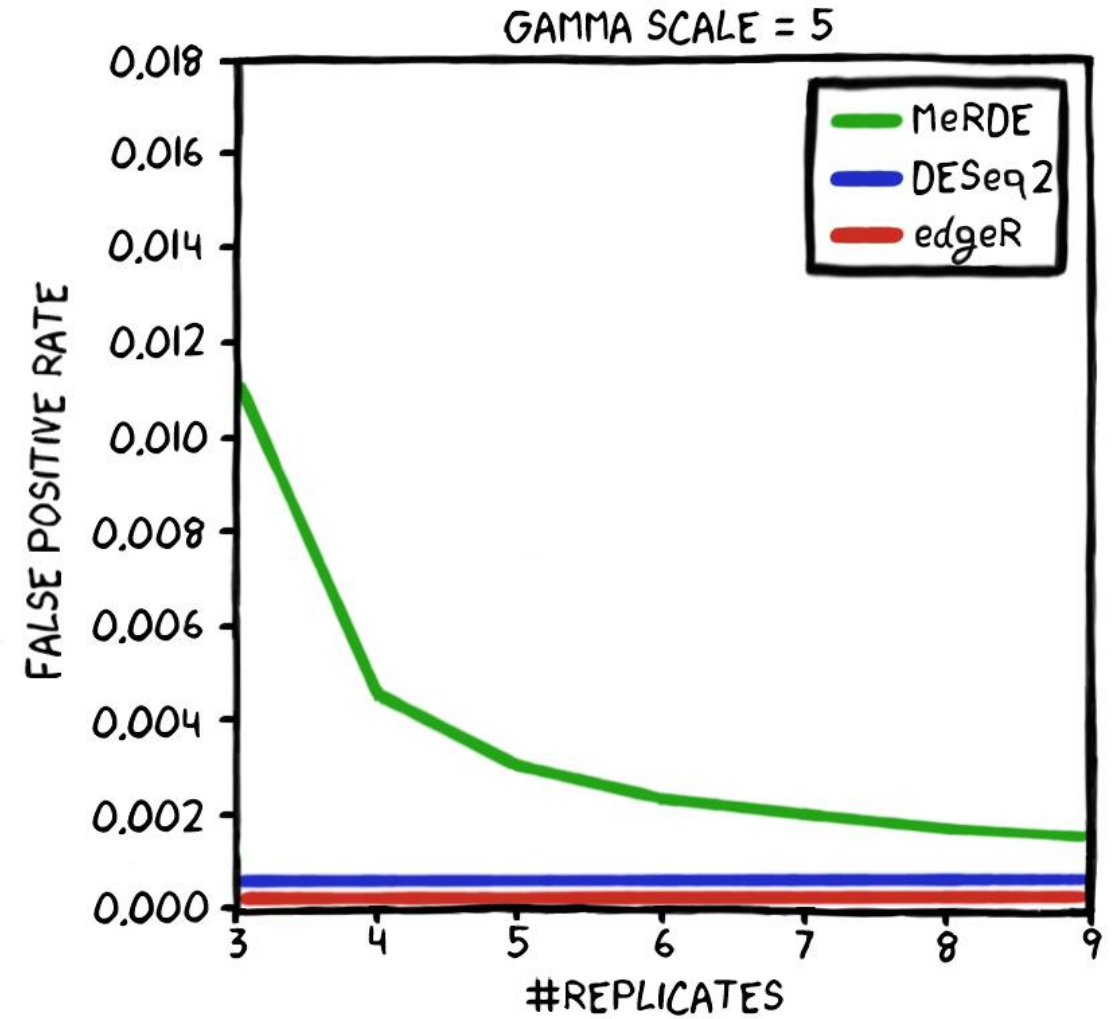
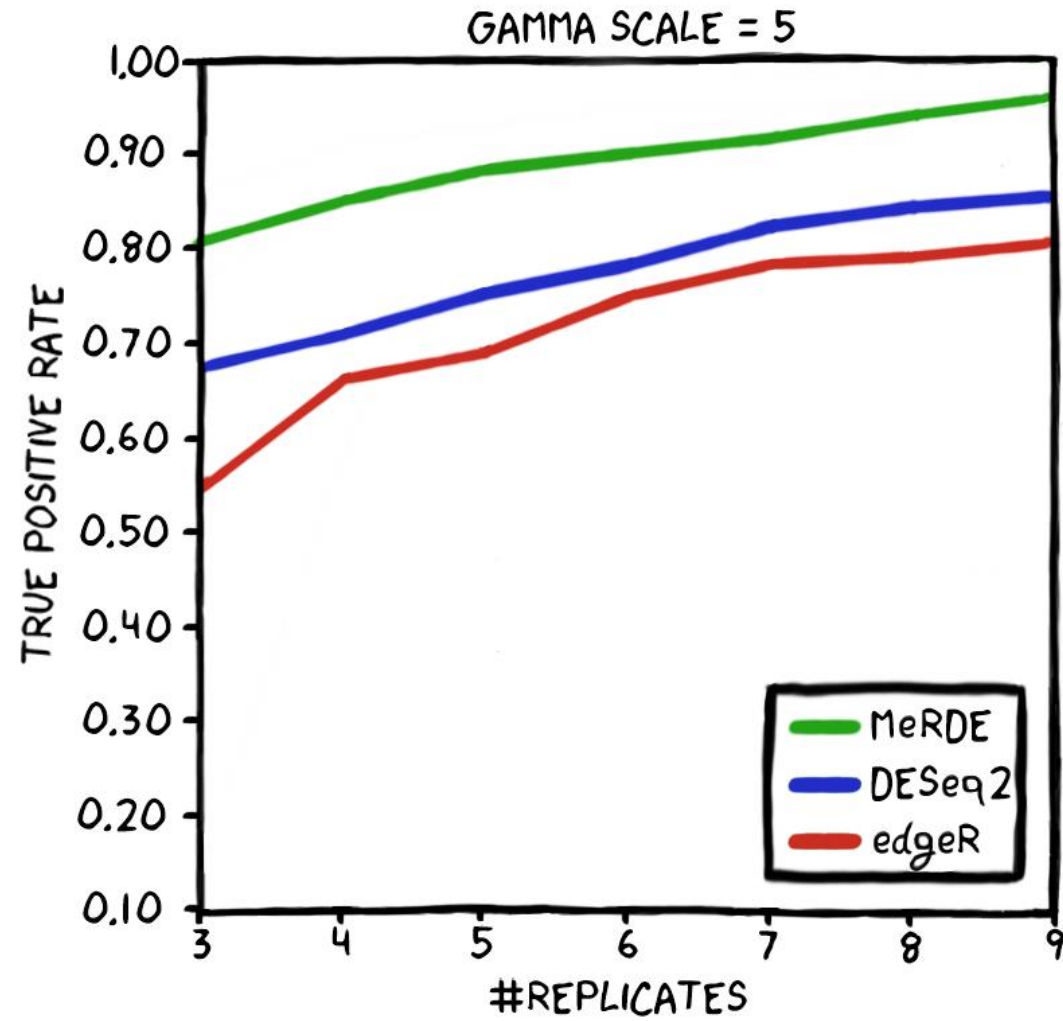


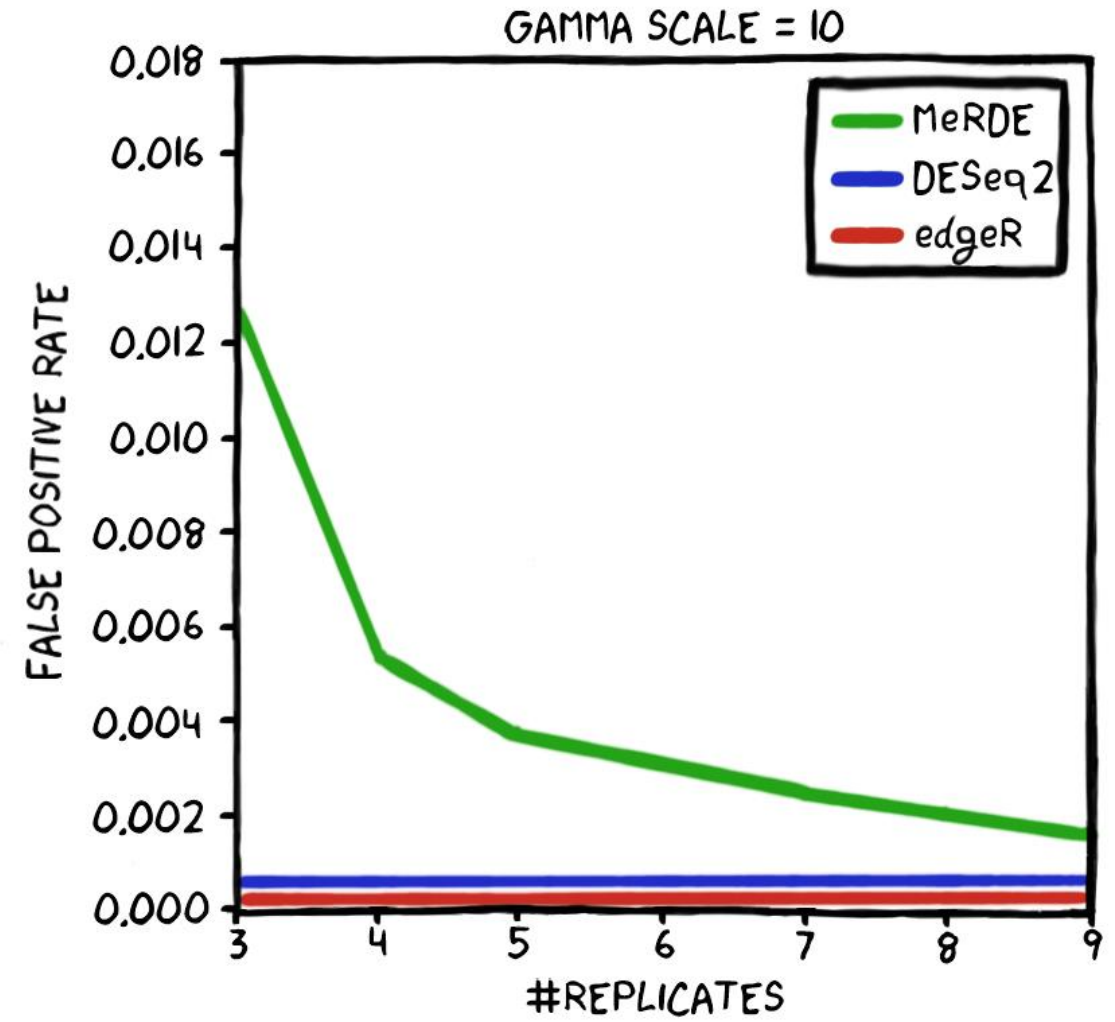
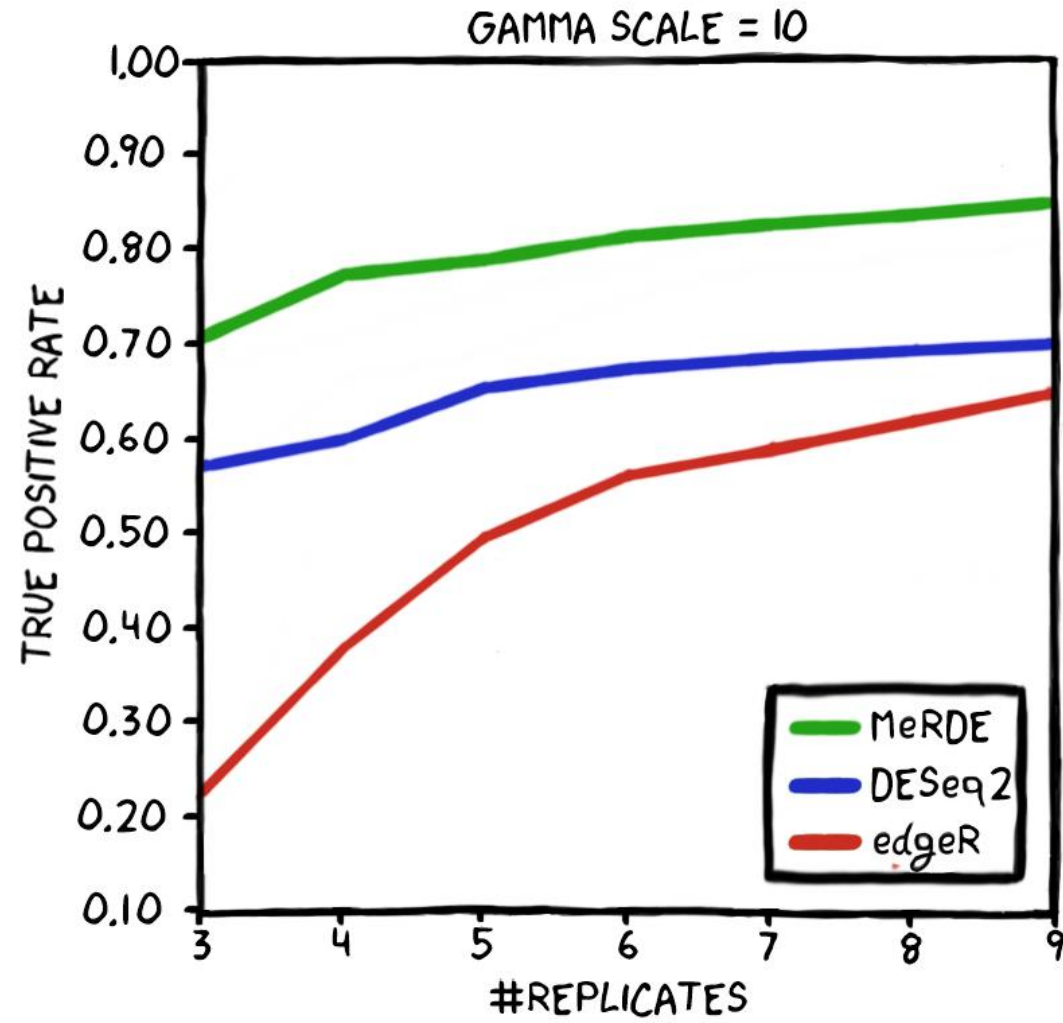
Still to do:

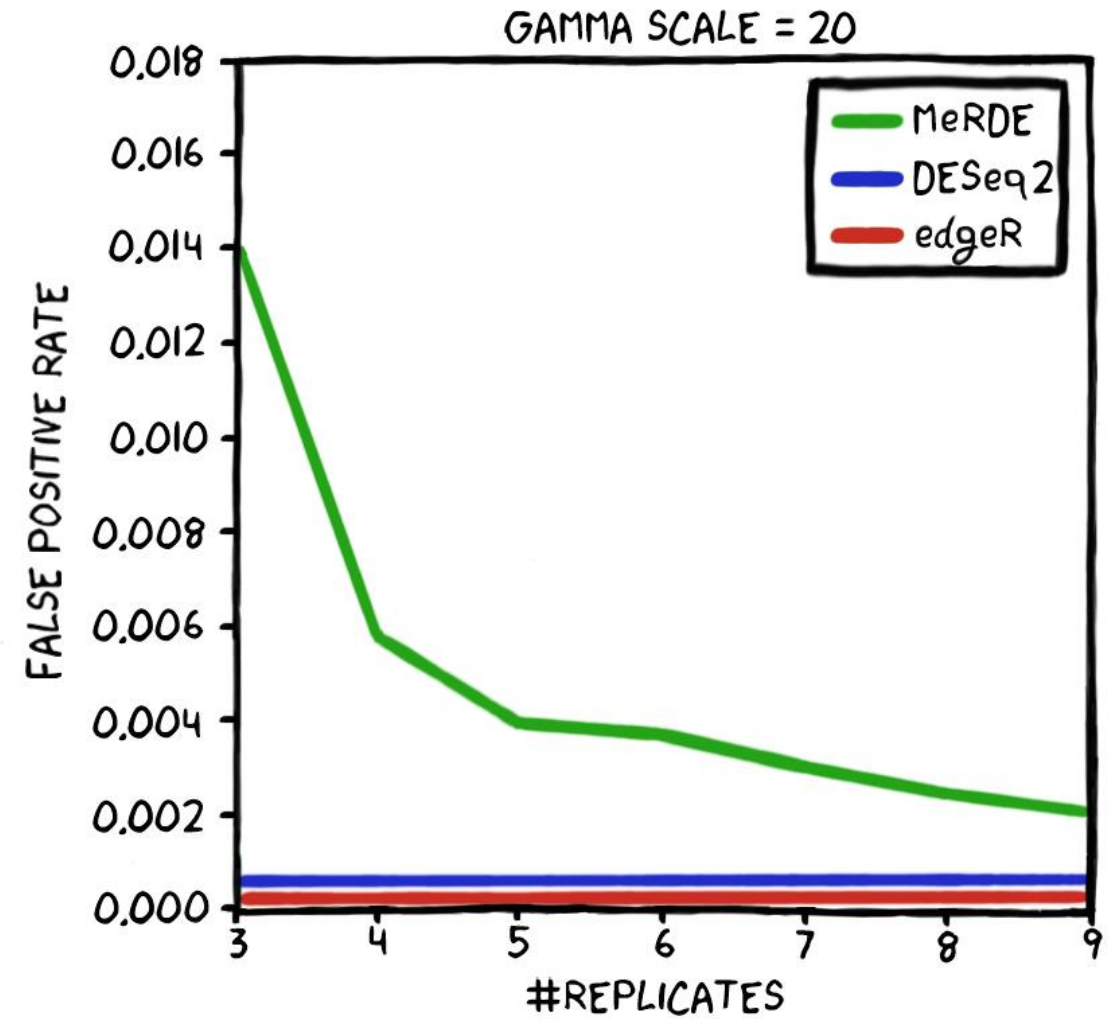
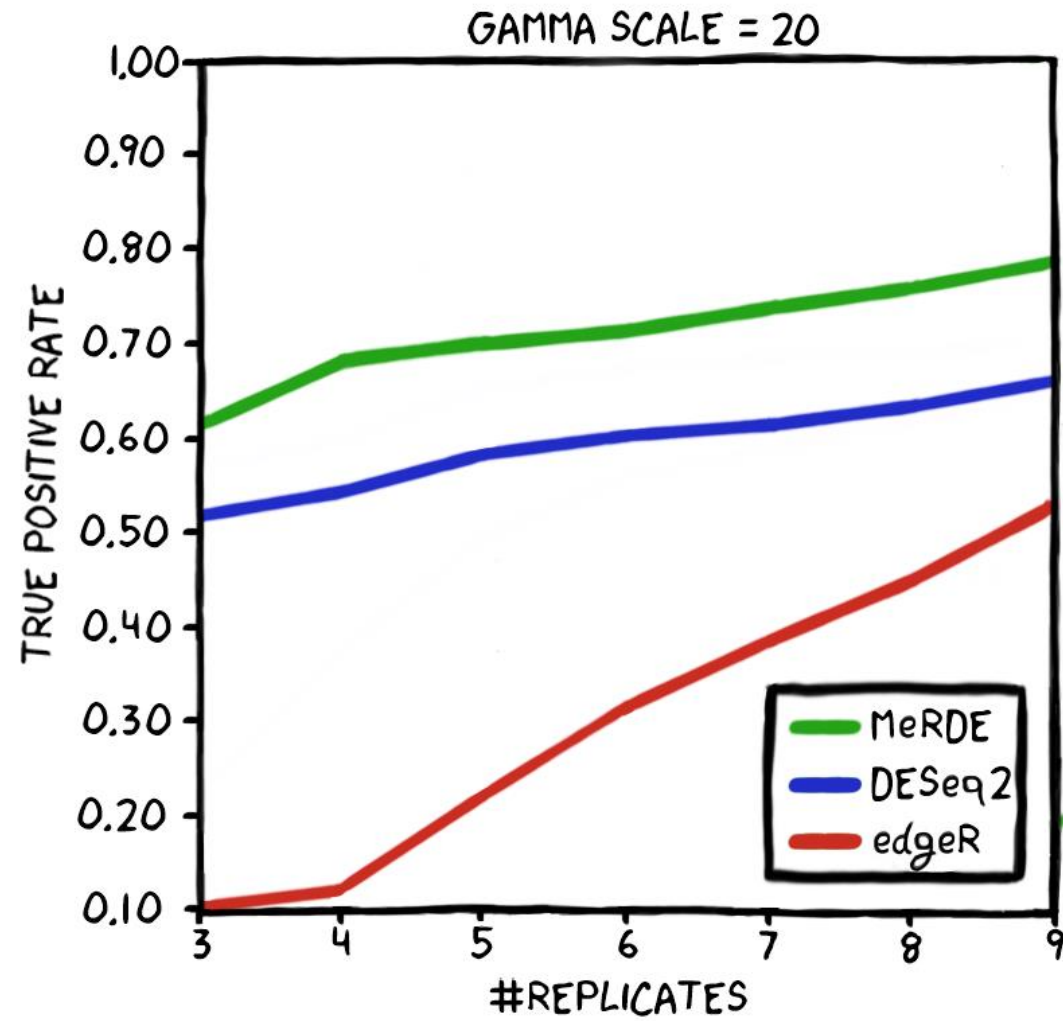
- Adjust transformation based on the genes estimated means and variances by applying a regularized cubic root function
- Instead of count transformation, create a new model and perform hypothesis tests assuming gamma-distributed random variables











- small RNA-Seq data sets of three different species of different tissues between different ages

- Human: 60 Datasets
- Mouse: 74 Datasets
- *Notho. furzeri*: 160 Datasets

