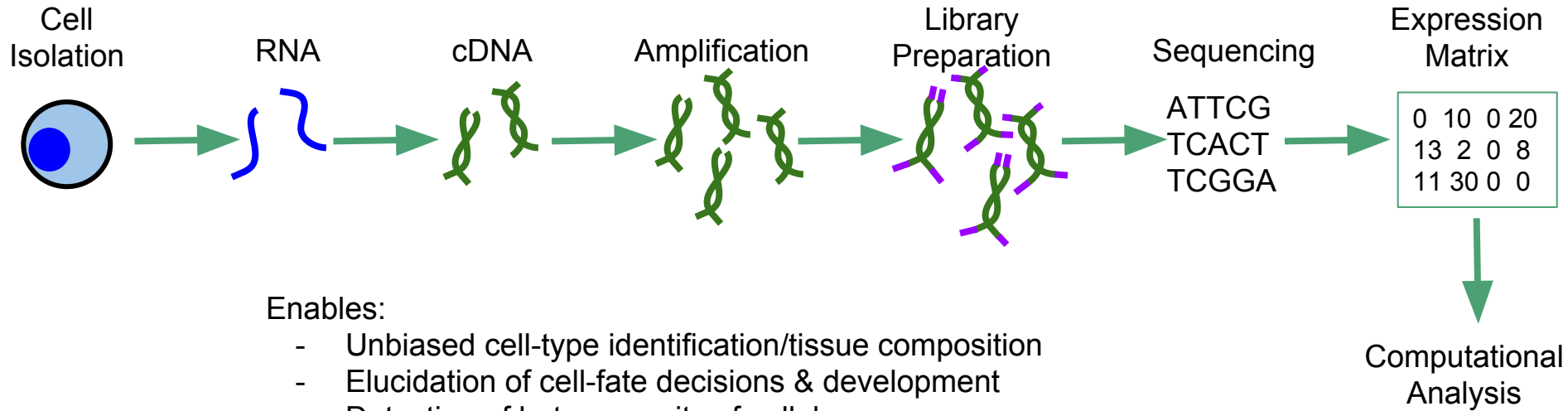


# Understanding Nothing: Zeros in scRNASeq

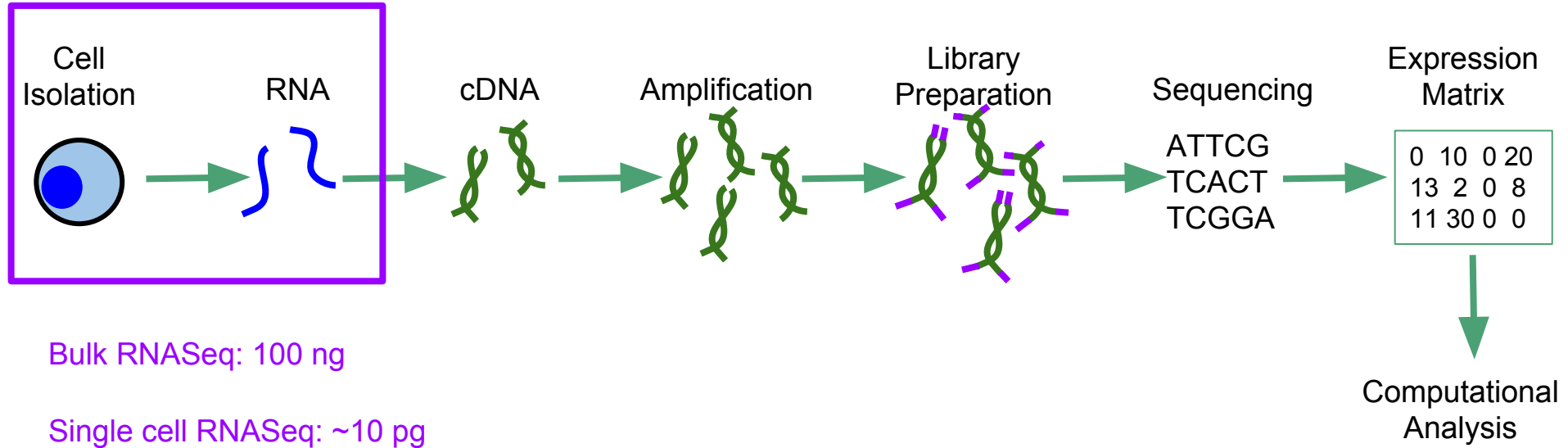
---

Tallulah Andrews, 27 Sept 2016

# Single-cell vs bulk RNASeq



# Single-cell vs bulk RNASeq

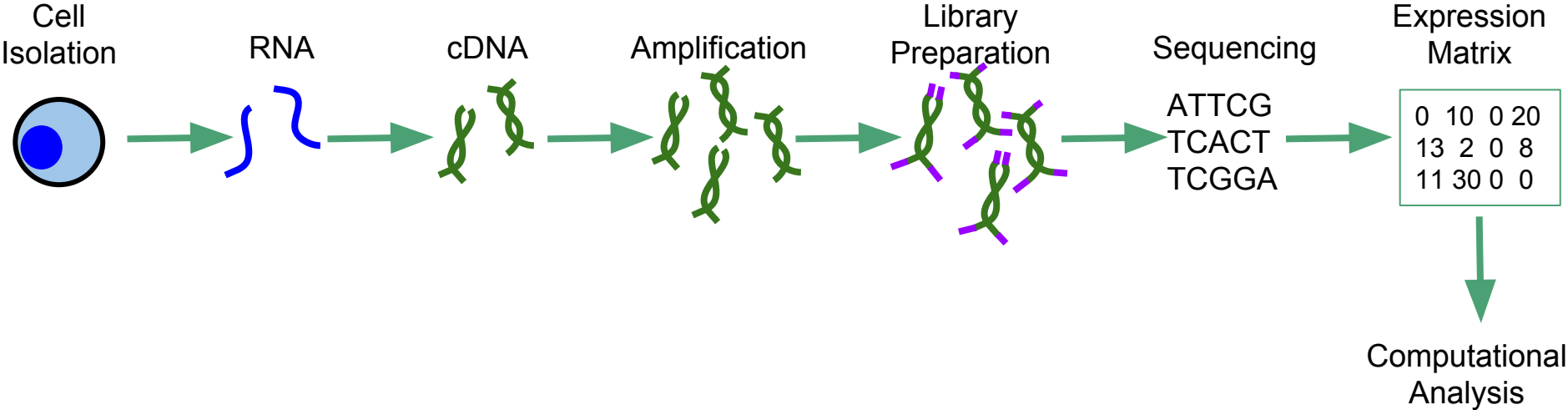


# Zeros Dominate scRNASeq

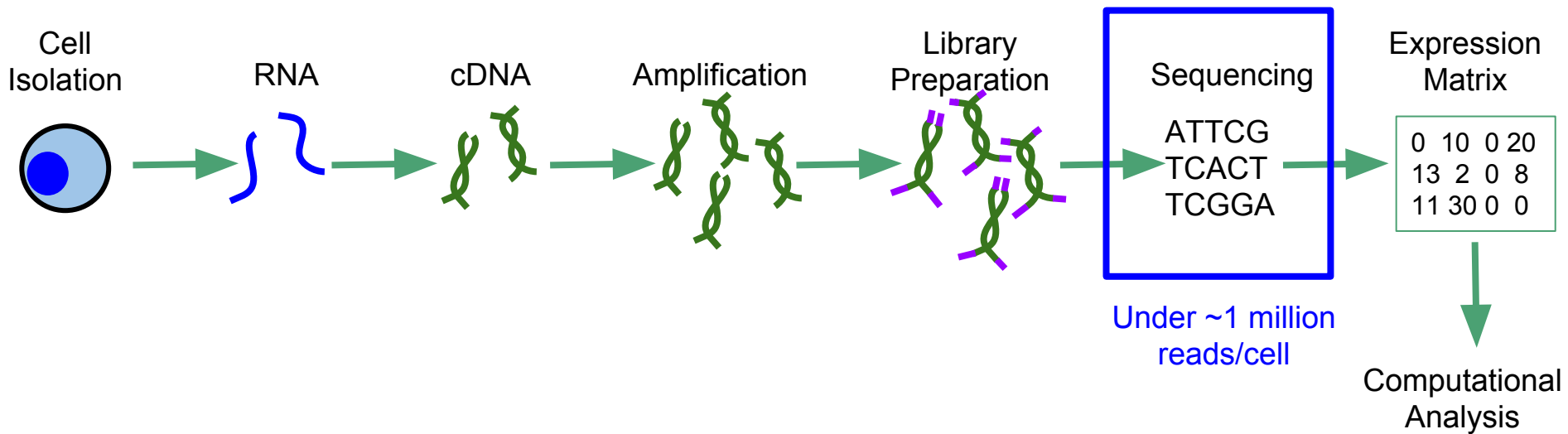
Dataset	Type	No. Cells	No. Genes	Prop Zero
Buettner	mouse ESCs	279	17,231	<b>51.2%</b>
Shalek	mouse bone marrow	324	12,474	<b>66.4%</b>
Deng	mouse embryo	255	17,406	<b>50.2%</b>
Usoskin	mouse neuron	530	15,585	<b>72.5%</b>
Kirschner	mouse ESCs	2,448	23,729	<b>62.5%</b>
Linnarsson	mouse brain	2,542	17,867	<b>76.9%</b>
Pollen	human neural	301	19,624	<b>60.3%</b>
Zhong	mouse embryo	49	20,558	<b>38.0%</b>

\*Cells with > 2,000 detected genes  
\*\*Genes seen in >3 cells

# Source of Zeros

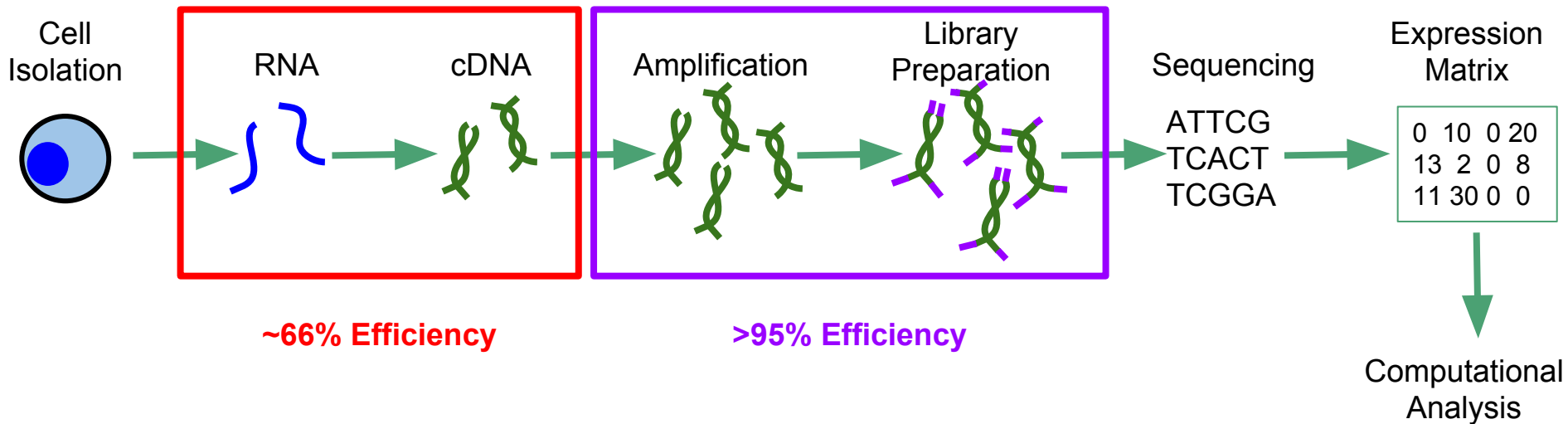


# Source of Zeros



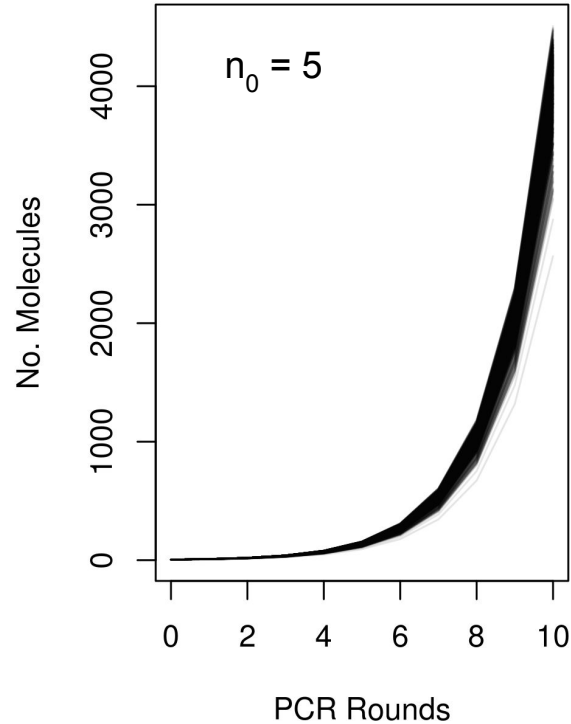
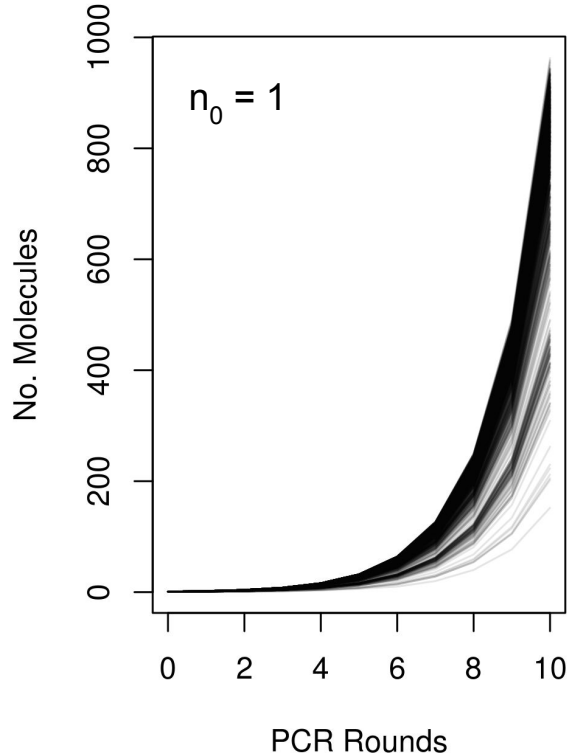
Svensson et al. (2016)

# Source of Zeros



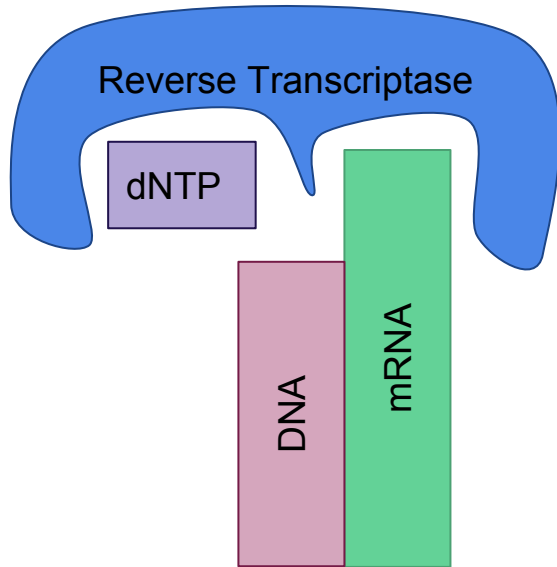
Reiter et al. (2011) & Bengtsson et al. (2008)

# RT failure propagates downstream



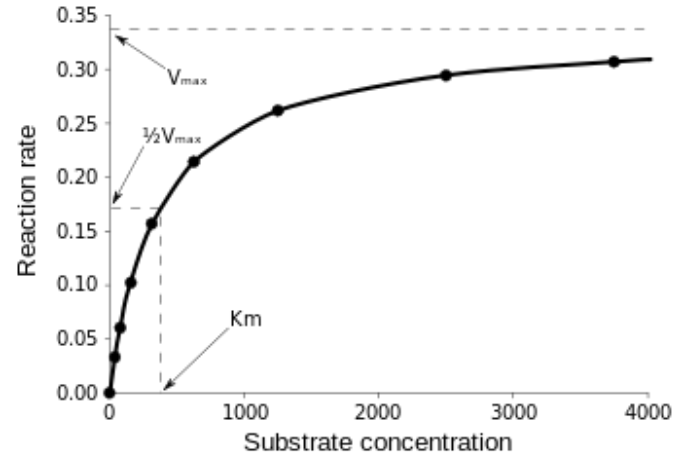


# Reverse Transcription = Michaelis-Menten



To model probability:  
 $V_{\max} = 1$

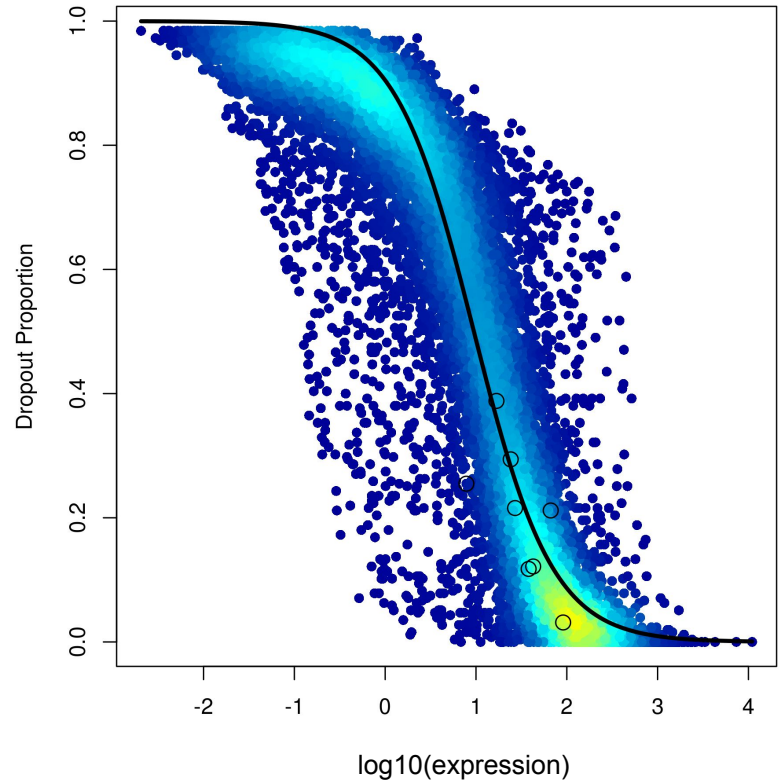
Detection probability = 
$$= \frac{V_{\max} [S]}{(K_M + [S])}$$



# MM vs Other Models

Michaelis-Menten Modelling of Dropouts (M3Drop)

- $P_{\text{dropout}} = 1 - [s]/(K+[s])$
- **For Deng: K = 9.5**



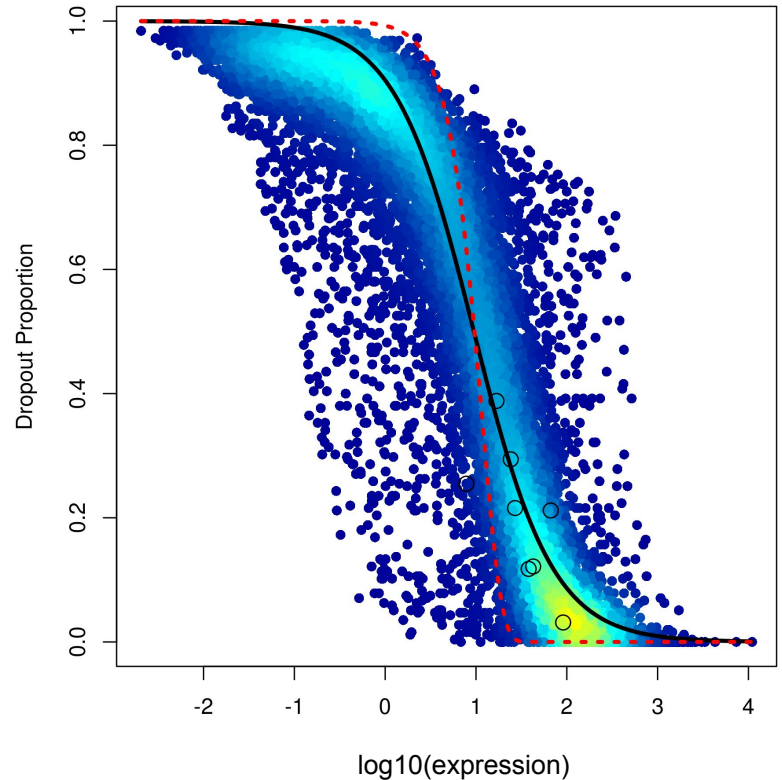
# MM vs Other Models

## Michaelis-Menten Modelling of Dropouts (M3Drop)

- $P_{\text{dropout}} = 1 - [s]/(K+[s])$
- **For Deng:  $K = 9.5$**

## Zero Inflated Factor Analysis (ZIFA)

- Dimensionality Reduction for scRNASeq
- $P_{\text{dropout}} = e^{-\lambda[s]}$
- **For Deng:  $\lambda = 0.0075$**



# MM vs Other Models

## Michaelis-Menten Modelling of Dropouts (M3Drop)

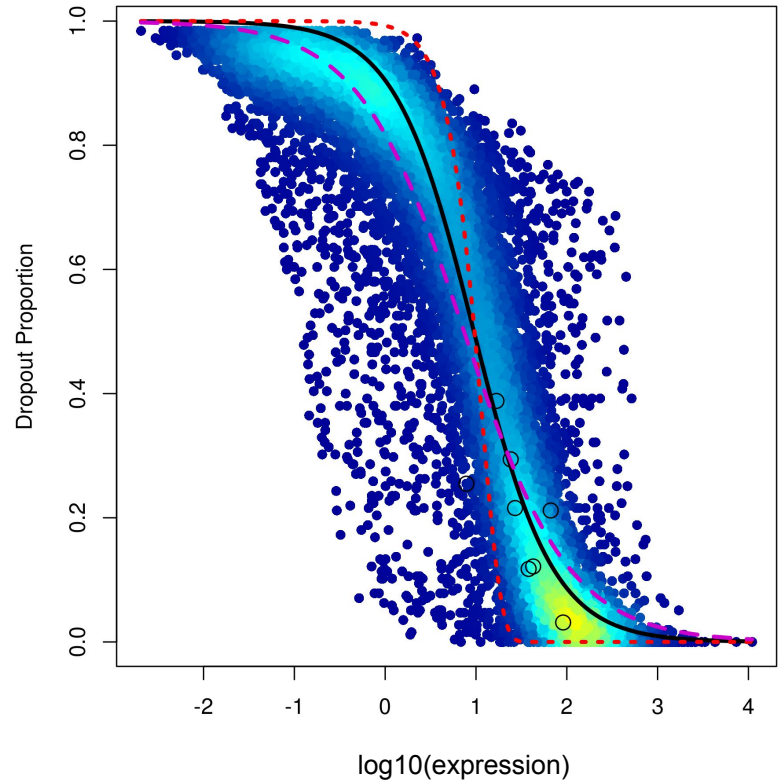
- $P_{\text{dropout}} = 1 - [s]/(K+[s])$
- **For Deng:  $K = 9.5$**

## Zero Inflated Factor Analysis (ZIFA)

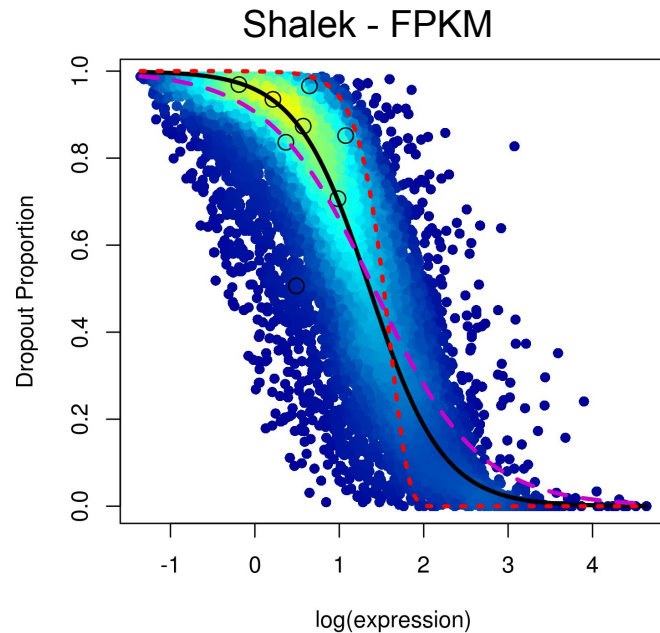
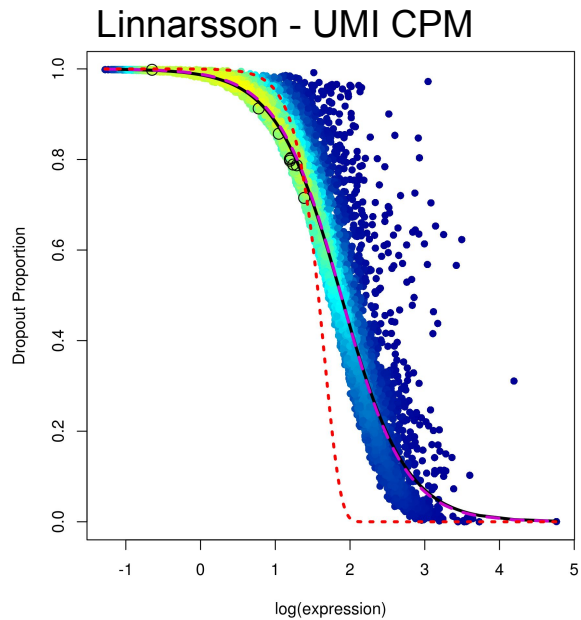
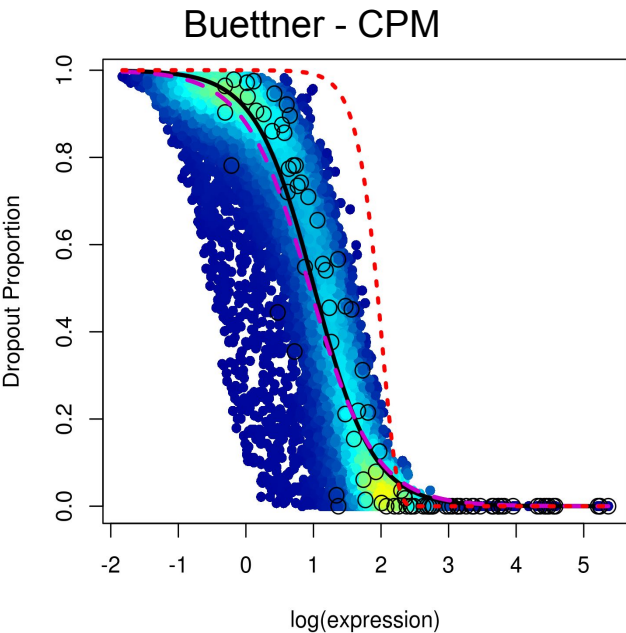
- Dimensionality Reduction for scRNASeq
- $P_{\text{dropout}} = e^{-\lambda[s]}$
- **For Deng:  $\lambda = 0.0075$**

## Single Cell Differential Expression (SCDE)

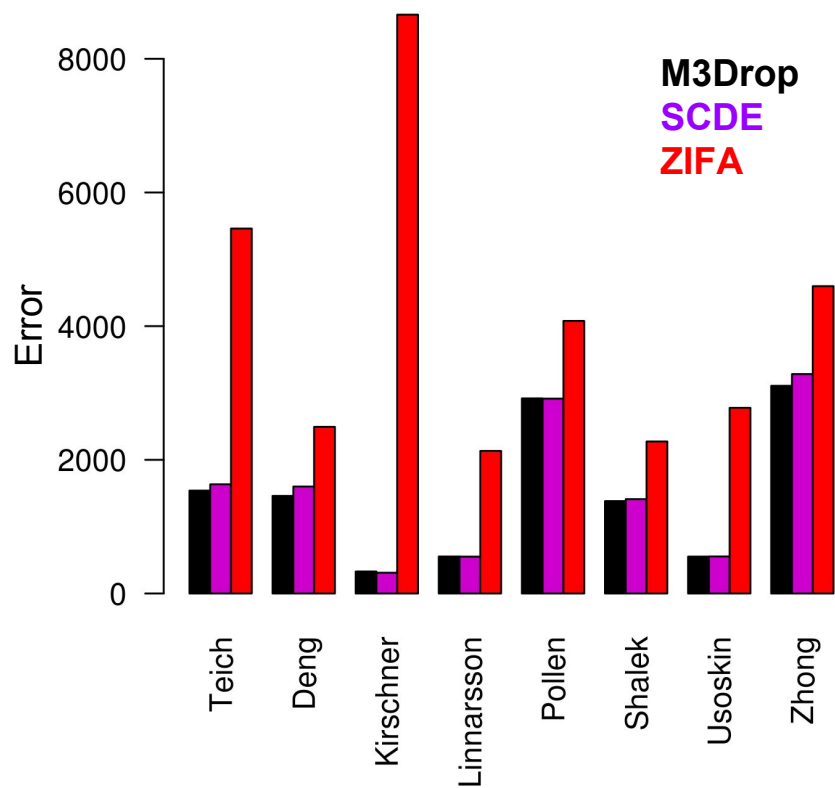
- $P_{\text{dropout}} = 1/(1+e^{-(a+b*\log([s]))})$
- **For Deng:  $a = 1.5, b = -0.75$**



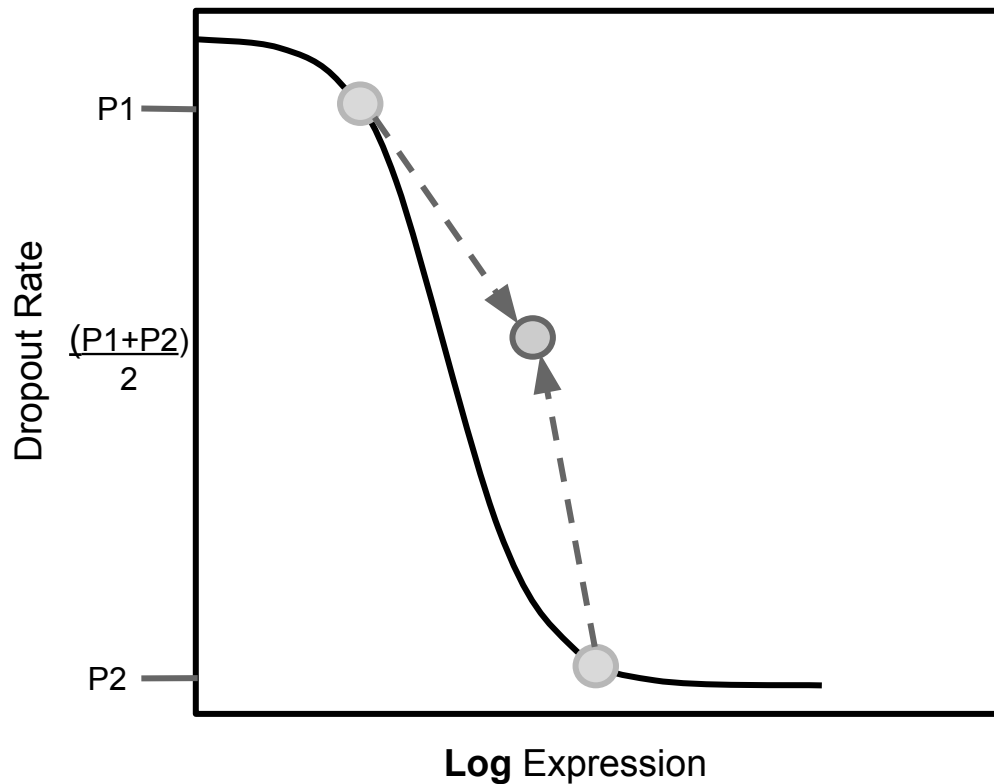
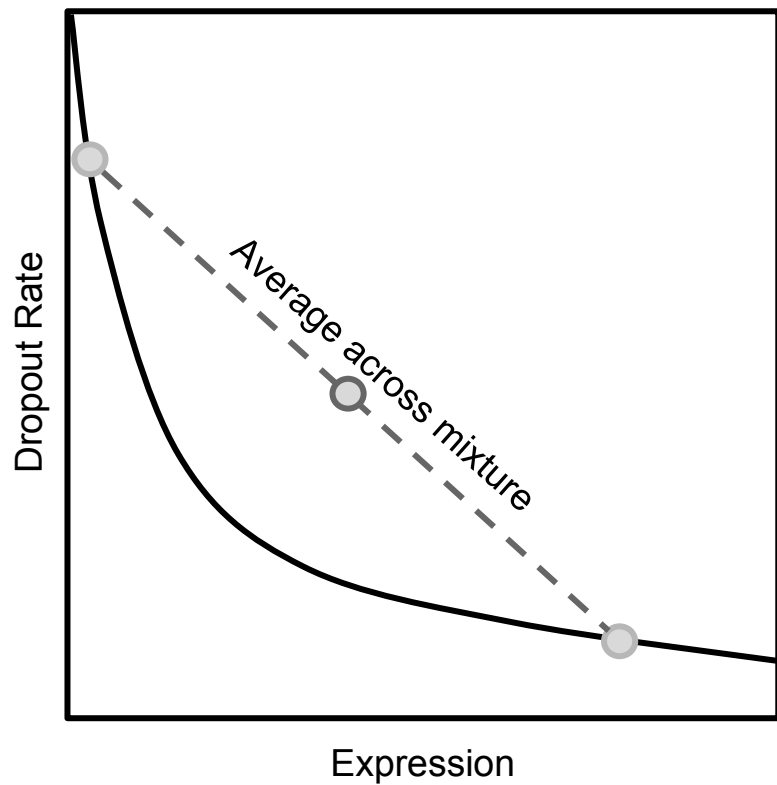
# Michaelis-Menten fits diverse datasets.



# Michaelis-Menten fits diverse datasets.



# Differentially Expressed Genes are Outliers



# Outlier/DE gene detection

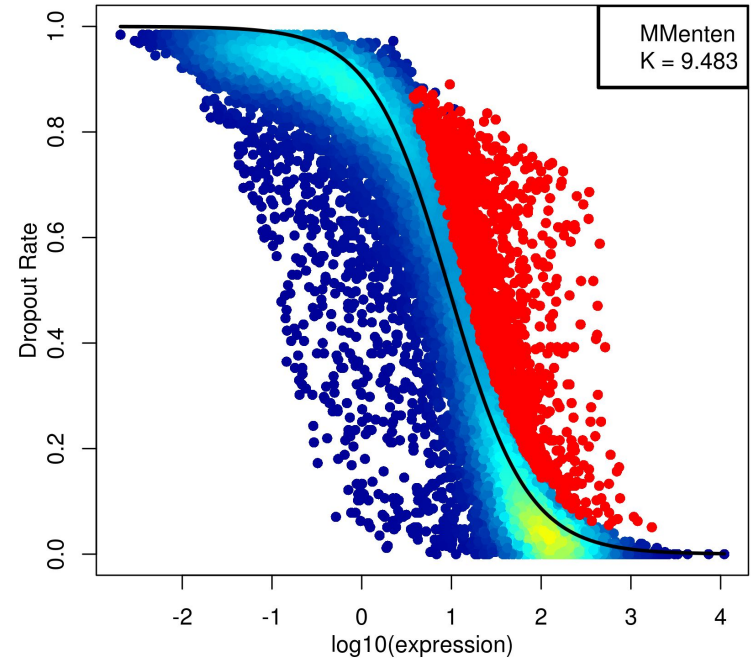
Michaelis-Menten:

$$P_{\text{dropout}} = 1 - S/(K+S)$$

Rearrange to solve for K:

$$K = P / (1-P) * S$$

1. Calculate  $K_j$  for each gene
2. Propagate errors in estimates for  $S$  (mean expression) and  $P$  (observed dropout rate) to get error for  $K_j$
3. Estimate error of global  $K_M$
4. Test whether  $K_j$  is significantly larger than  $K_M$  fit across all genes using a Z-test combining errors of (2) & (3)





# Highly Variable Genes

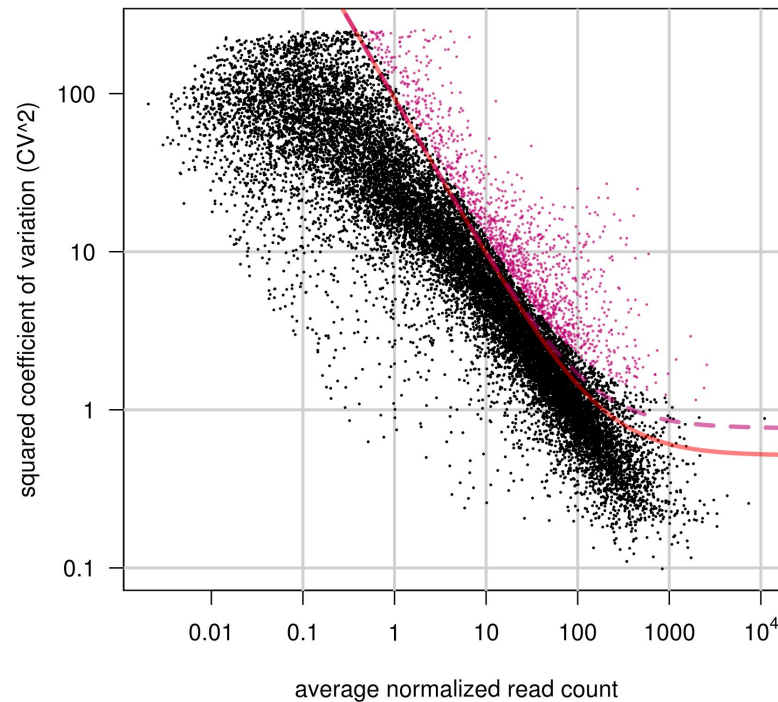
In general:

$$f(\text{variance}) = g(\text{mean})$$

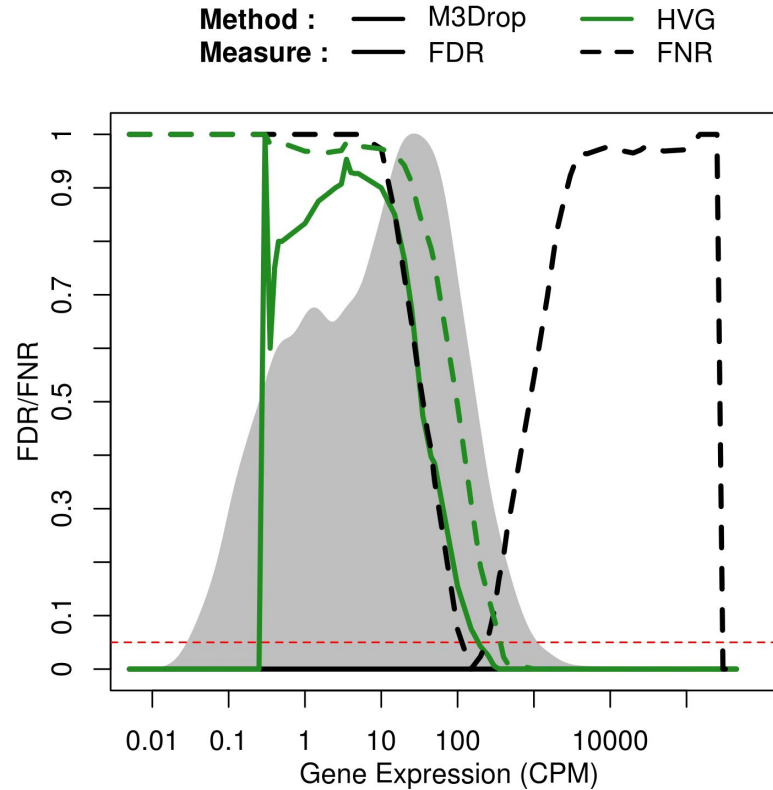
1. Fit a relationship between variance and mean expression
  - a. May use all genes or only spike-ins in fitting
2. Identify points above this relationship

Brennecke et al. (2013) :

1.  $CV^2 = a_1/\mu + \alpha_0$
2. Significant outliers detected using  $\chi^2$ -test

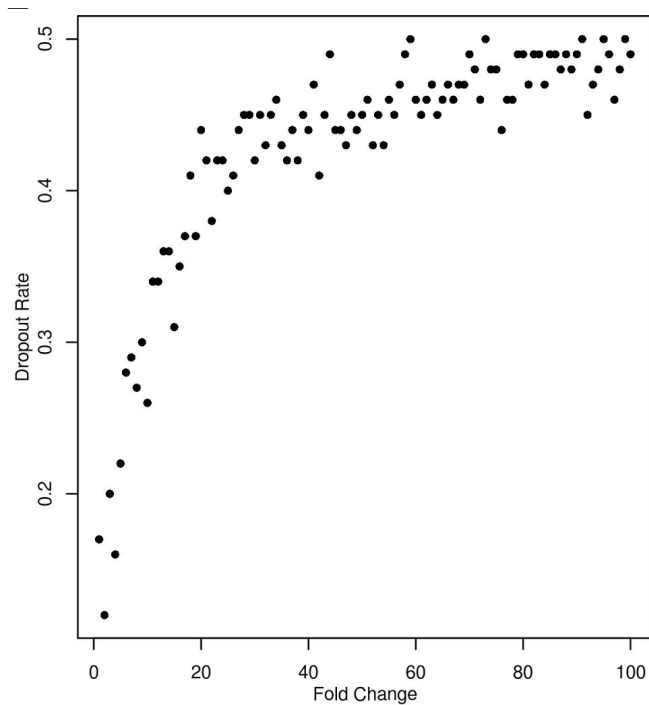
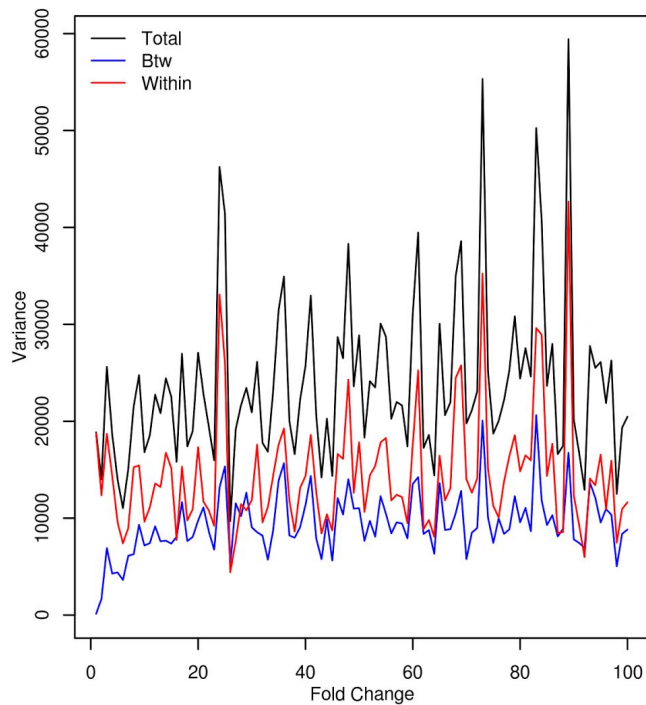


# DE Simulations - Dropouts vs Variance.

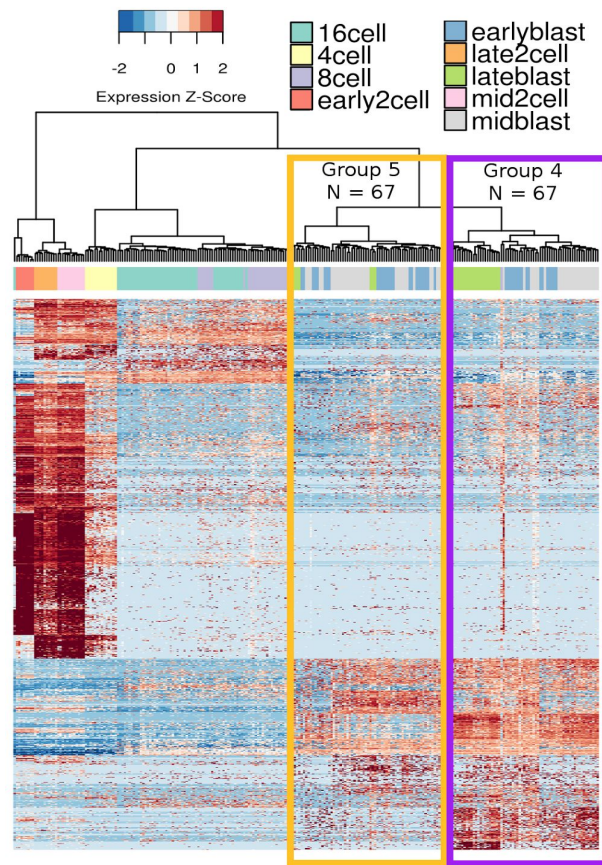
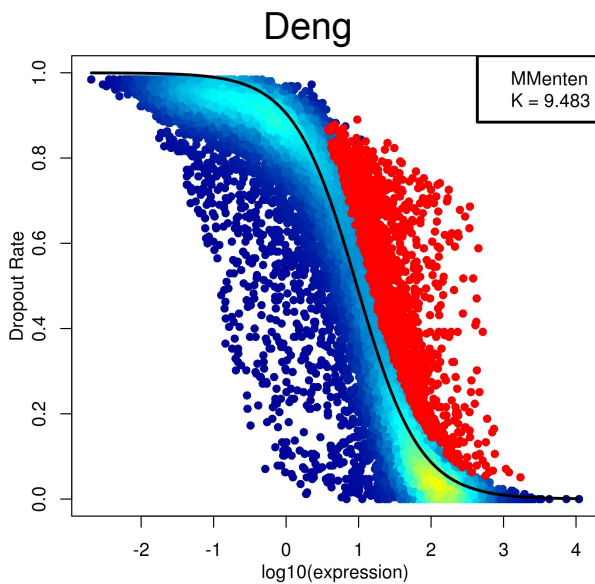
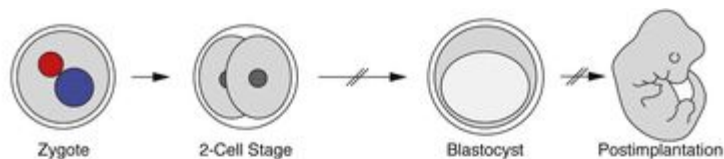


# DE Simulations - Dropouts vs Variance.

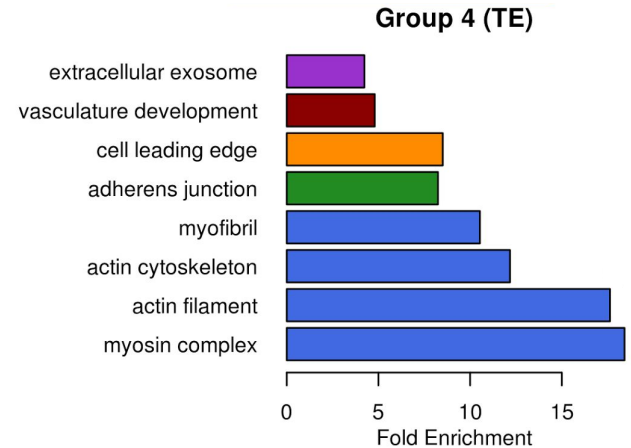
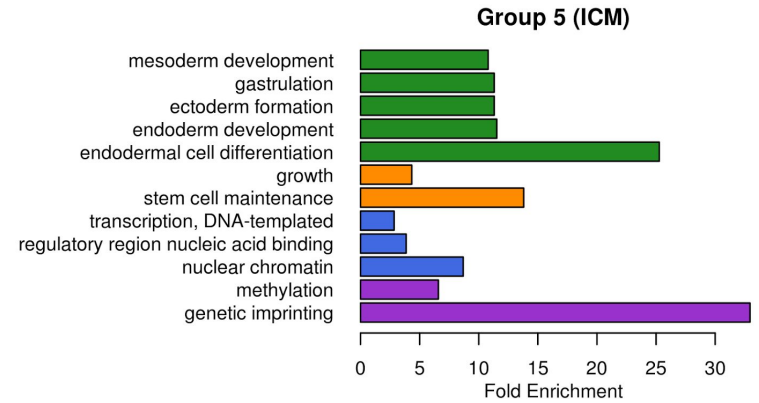
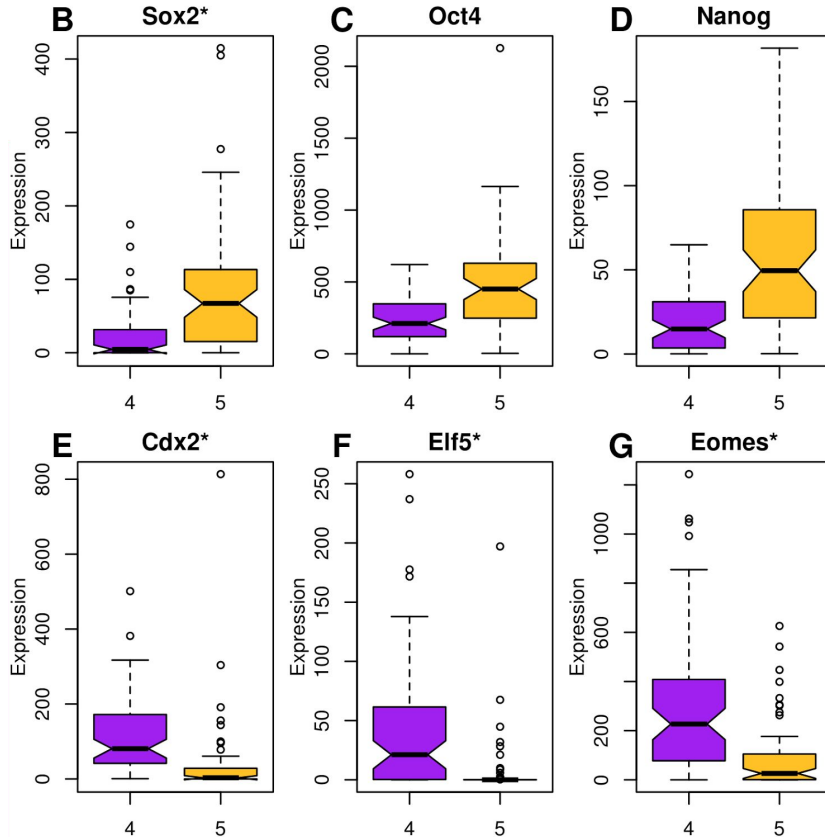
$\mu = 100, n = 100$



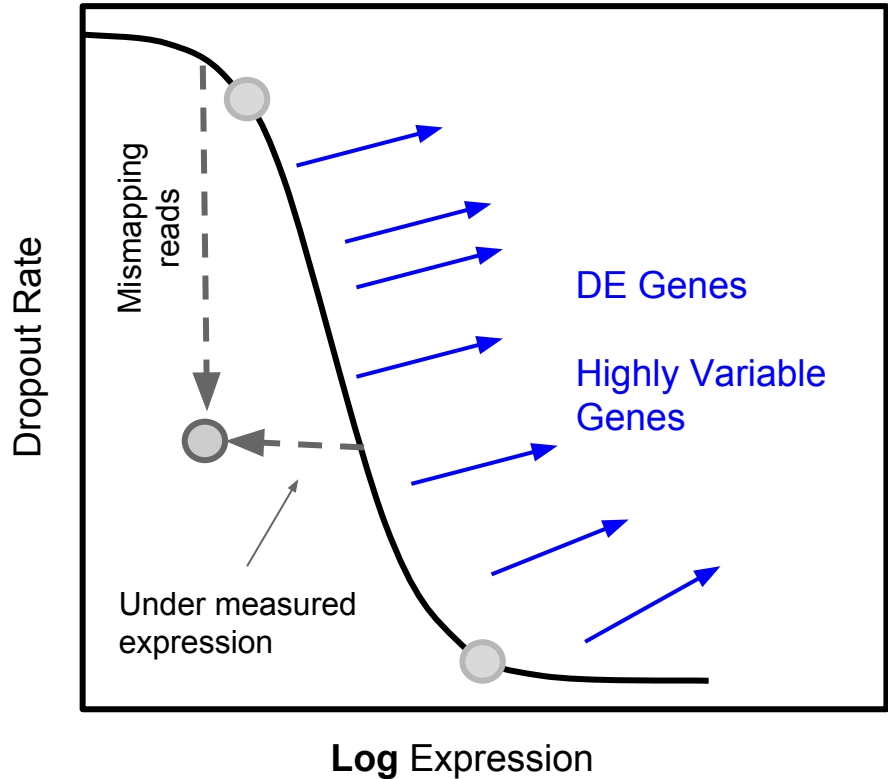
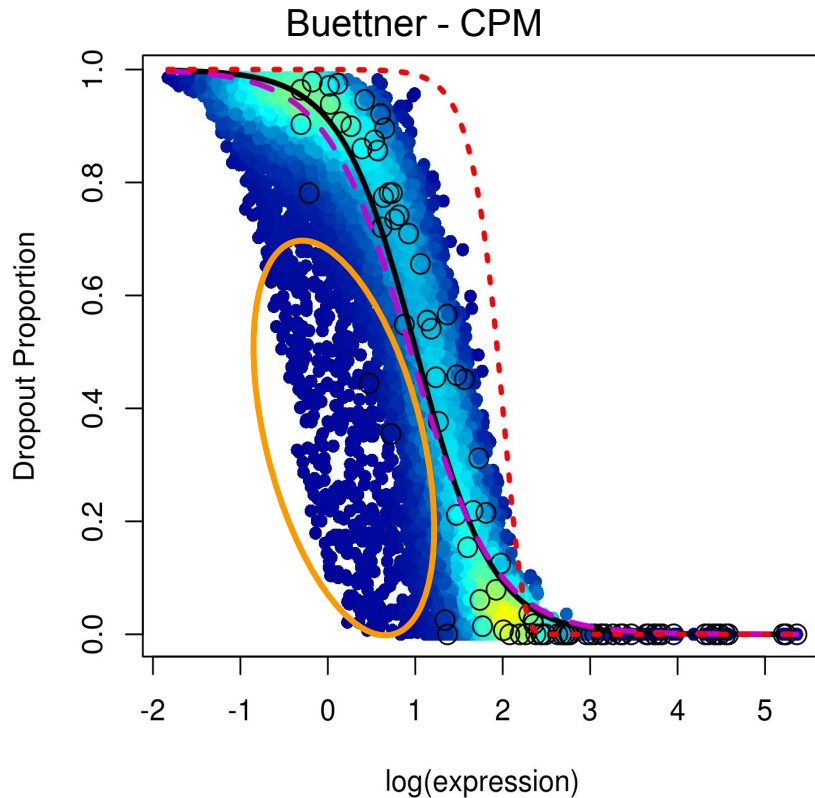
# Applying M3Drop to Early Mouse Development



# Identification of TE and ICM



# What are outliers to the left?

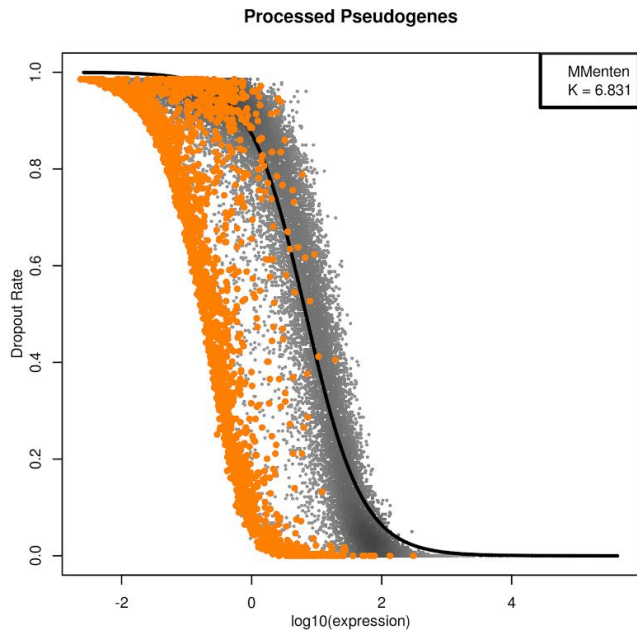


# Processed Pseudogenes = True Negatives

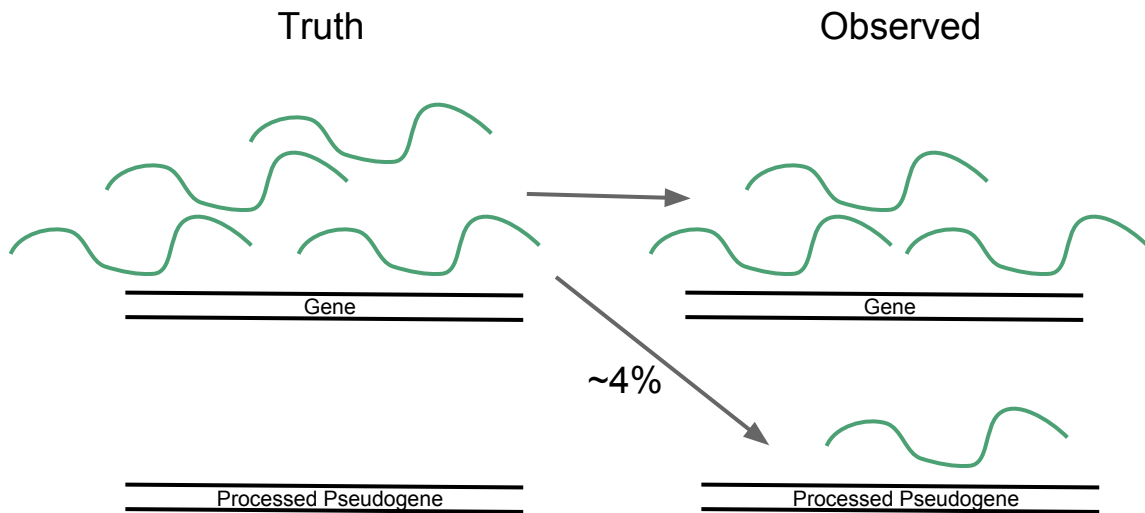


- Identical sequence to original transcript
- Lacks introns
- Lacks promoters & regulatory sequences
  - Assumed to not be transcribed
- >3,000 identified in the mouse genome
  - only 150 have confirmed expression

# Processed Pseudogenes - Mismapping Reads



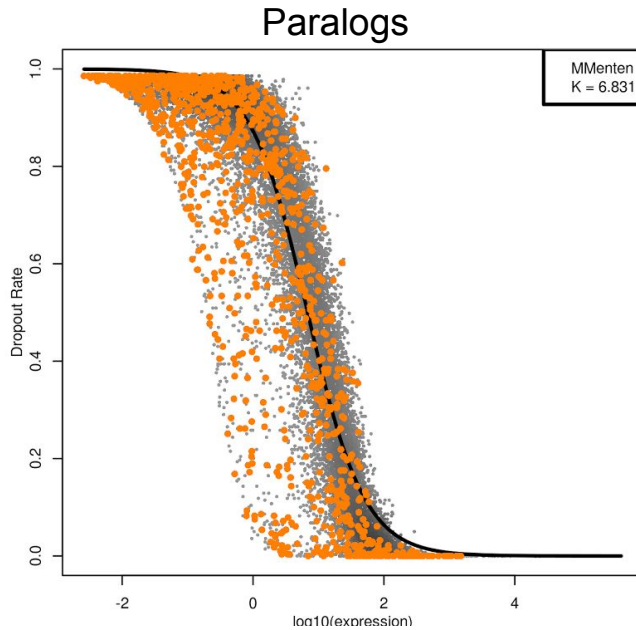
Processed Pseudogenes  
Left shifted by 1.4 ( $p \sim 0$ )



1% sequencing error rate x 100bp reads:  
4% of reads have 3+ sequencing errors

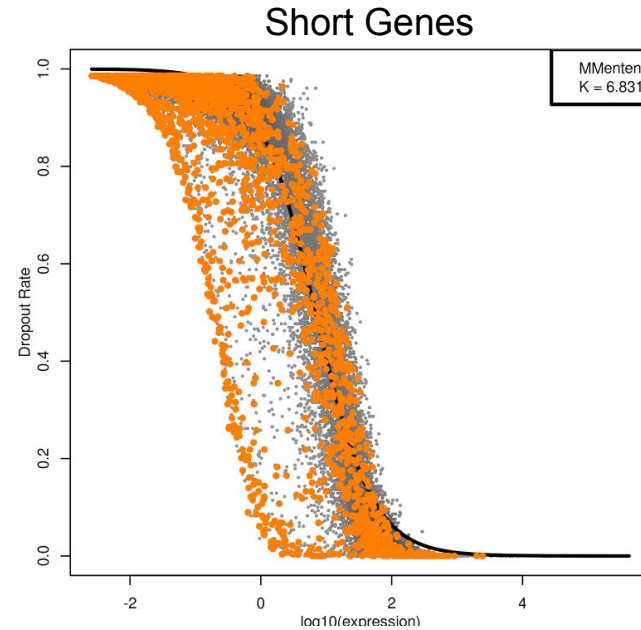


# Under-Measured Expression



Duplication node: Mus musculus  
Left shifted by 0.66 ( $p < 10^{-40}$ )

multimapping reads =  
under counting



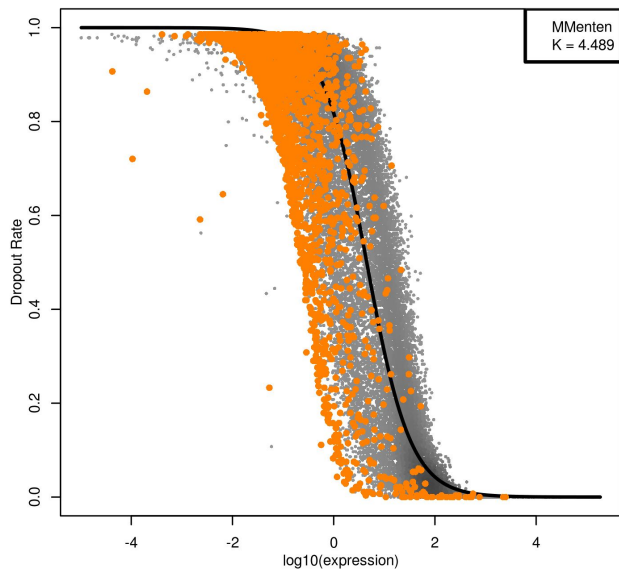
CDS < 300 n.t.  
Left shifted by 0.21 ( $p < 10^{-45}$ )

fewer unique fragments =  
fewer unique reads

# Tophat2 maps more reads to processed pseudogenes

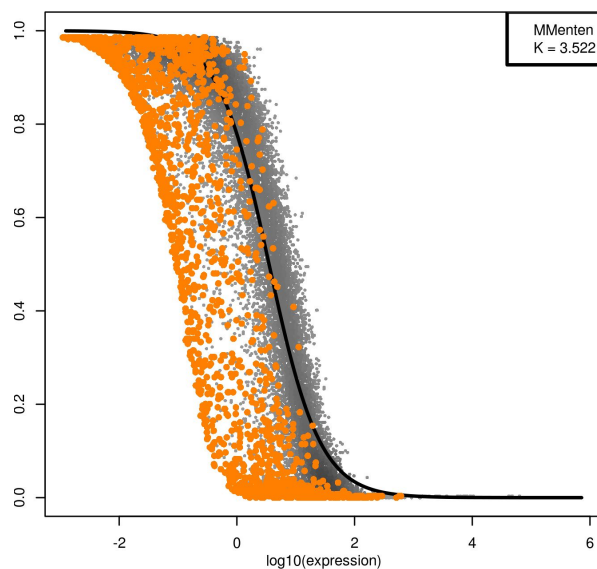
Kallisto

Processed Pseudogenes



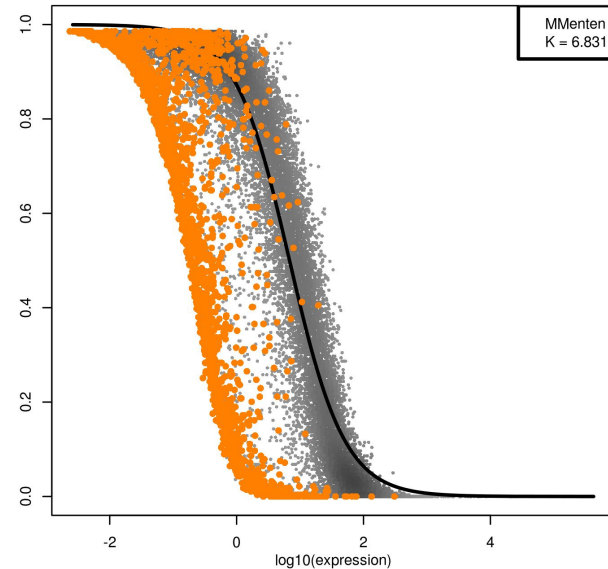
Tophat2

Processed Pseudogenes

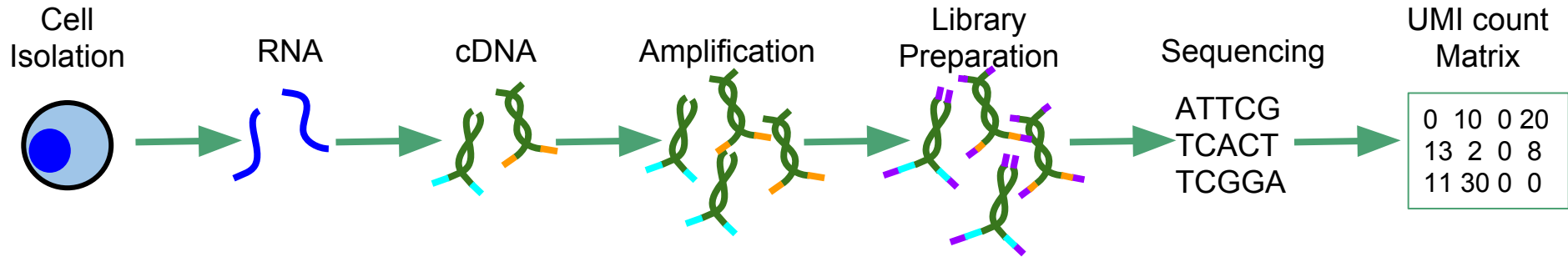


STAR

Processed Pseudogenes



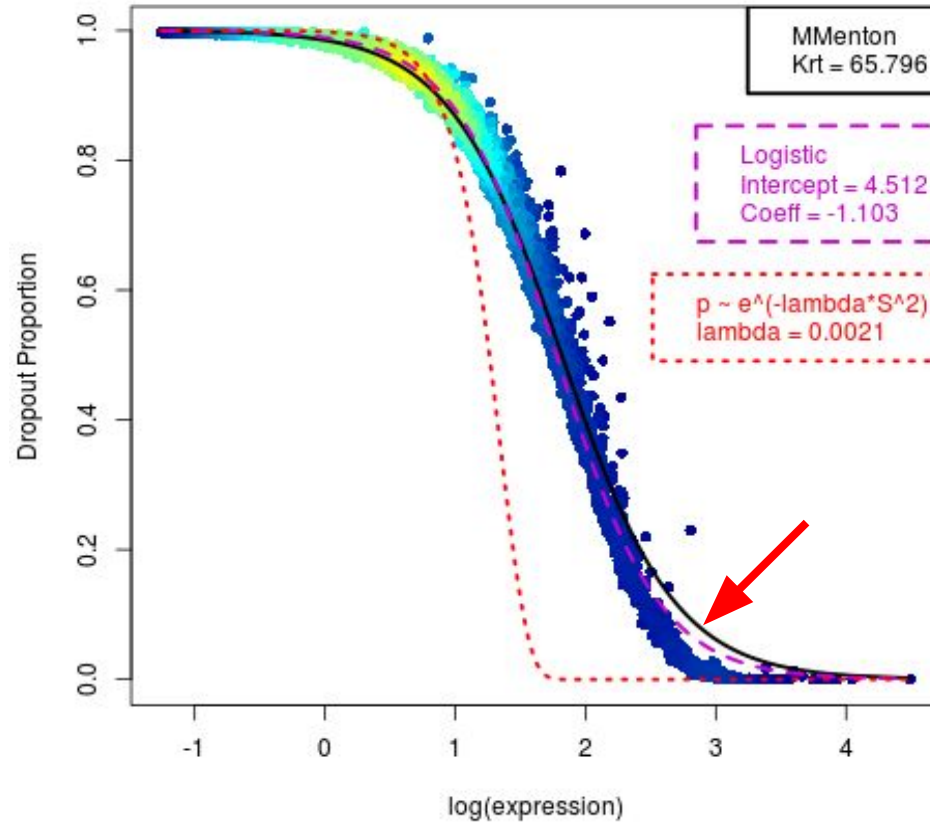
# Unique Molecular Identifiers (UMIs)



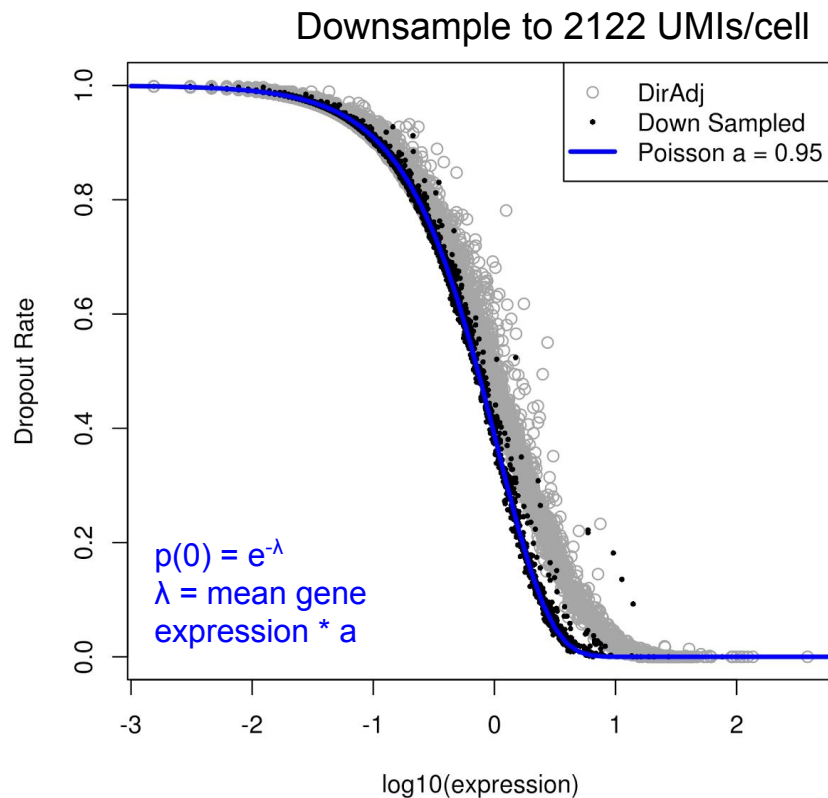
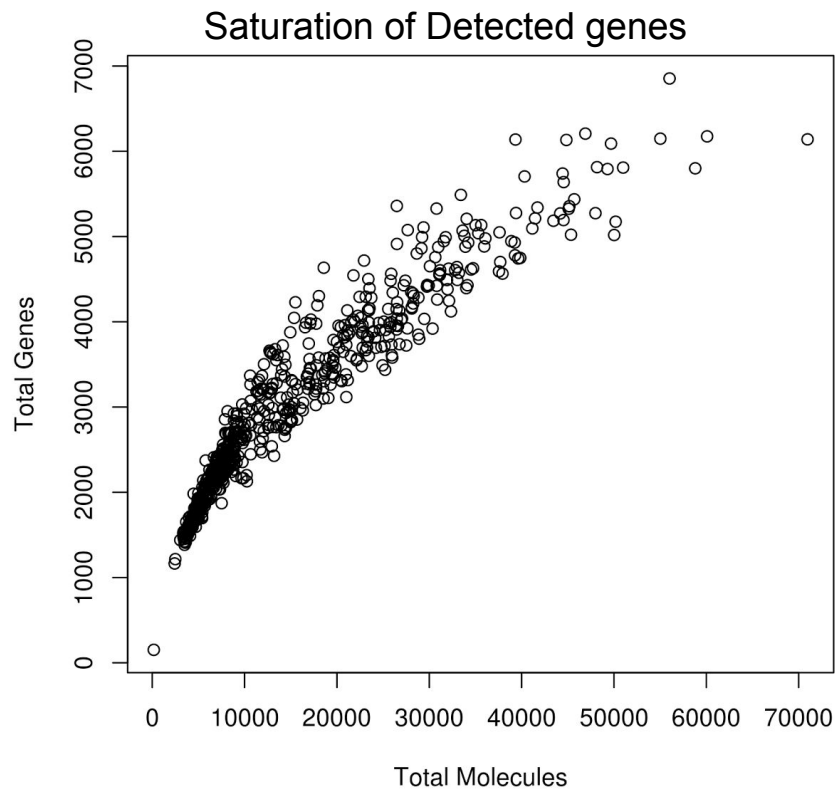
Enables:

- Correction for PCR duplicates (amplification noise)

# None of the proposed models fit corrected UMIs



# Cell-specific detection rates obscure true relationship

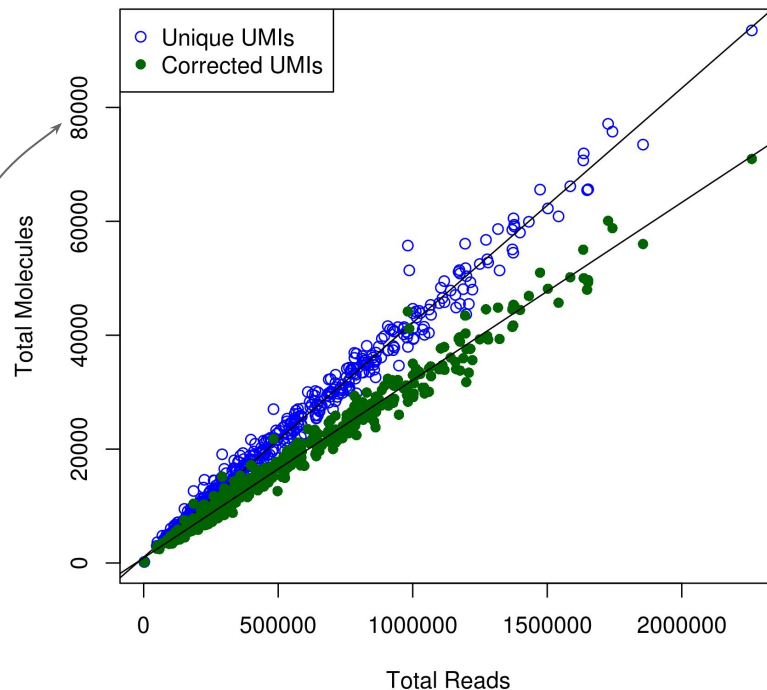


# The PoissonUMIs Model

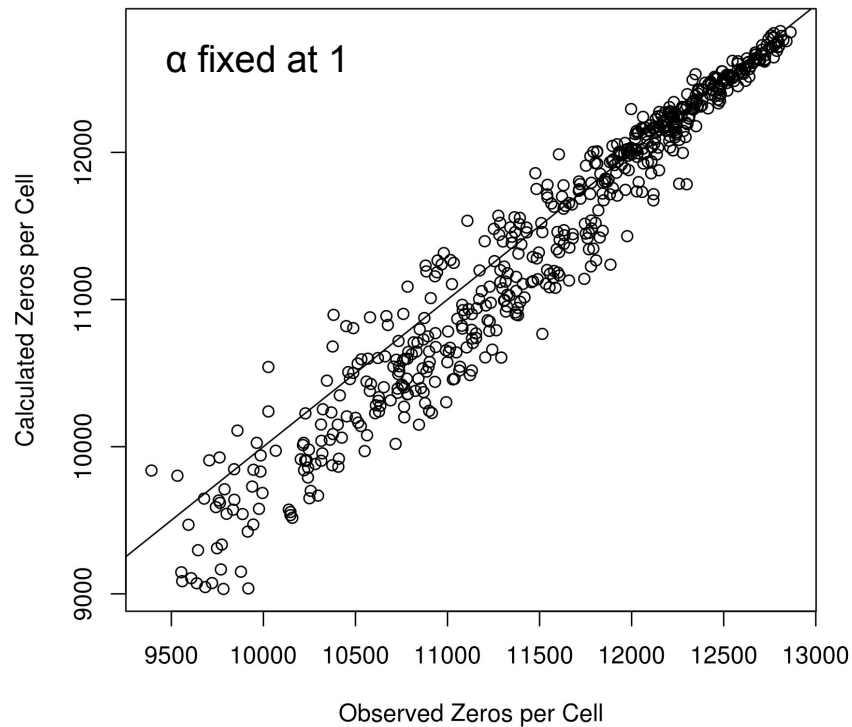
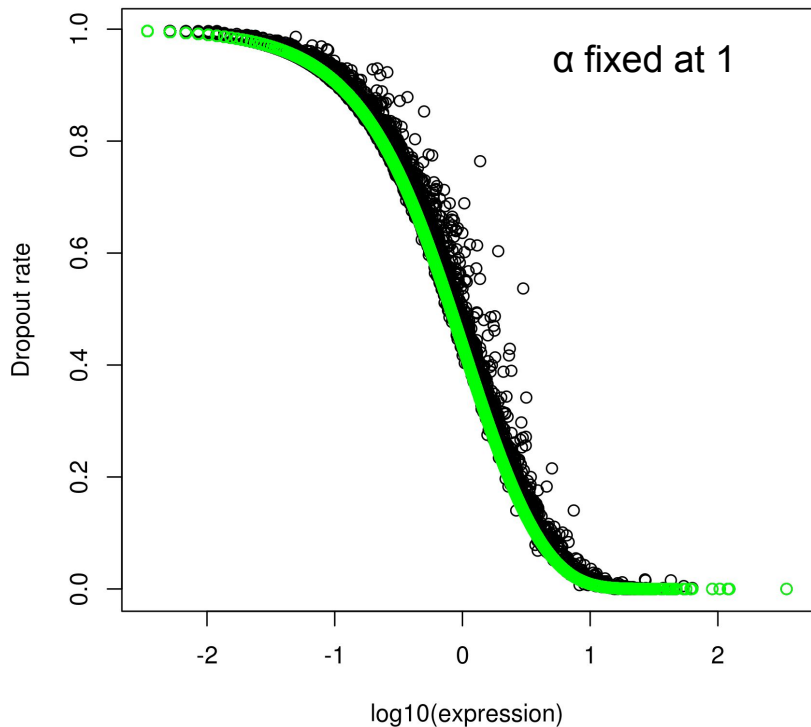
$$M_{ij} \sim \text{Poisson}(\lambda)$$
$$\lambda = m_i * m_j * \text{total} * \alpha$$

$M_{ij}$  = Molecules of gene  $j$  in cell  $i$   
 $m_i$  = proportion of molecules in cell  $i$   
 $m_j$  = proportion of molecules for gene  $j$   
total = total detected molecules  
 $\alpha$  = scaling factor

Account for different  
counting methods

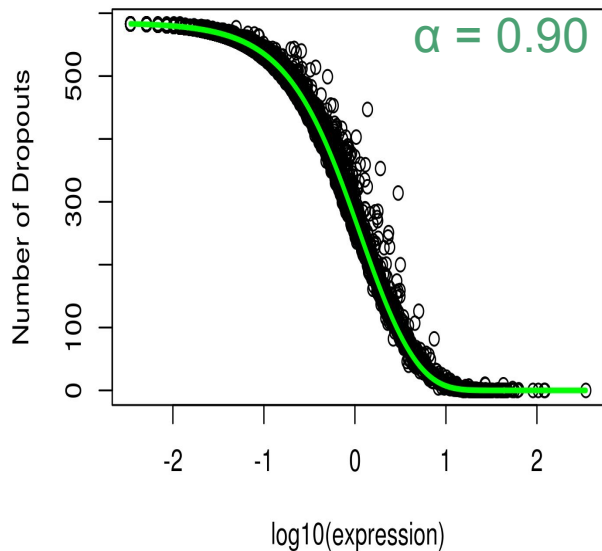


# Poisson model accounting for differences in read depth

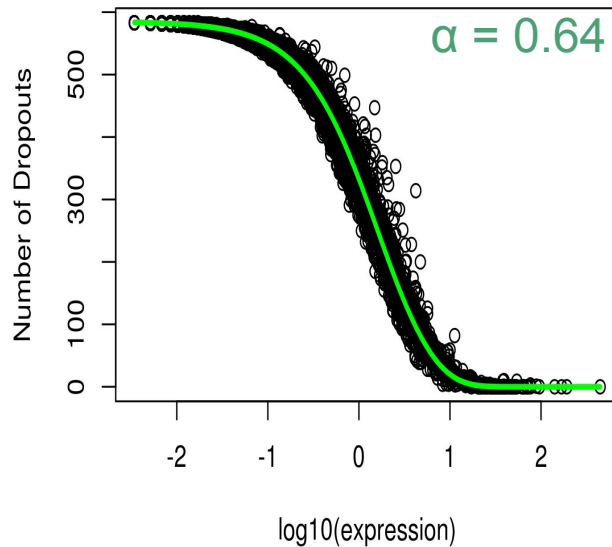


# Fitted alpha reflects quantification method

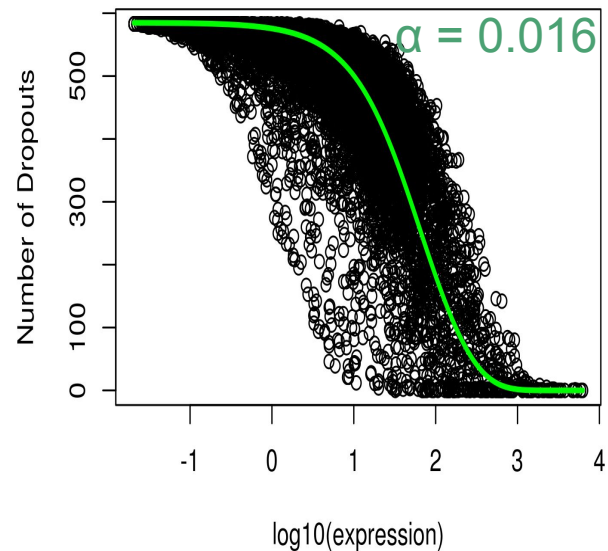
Corrected UMIs



Unique UMIs



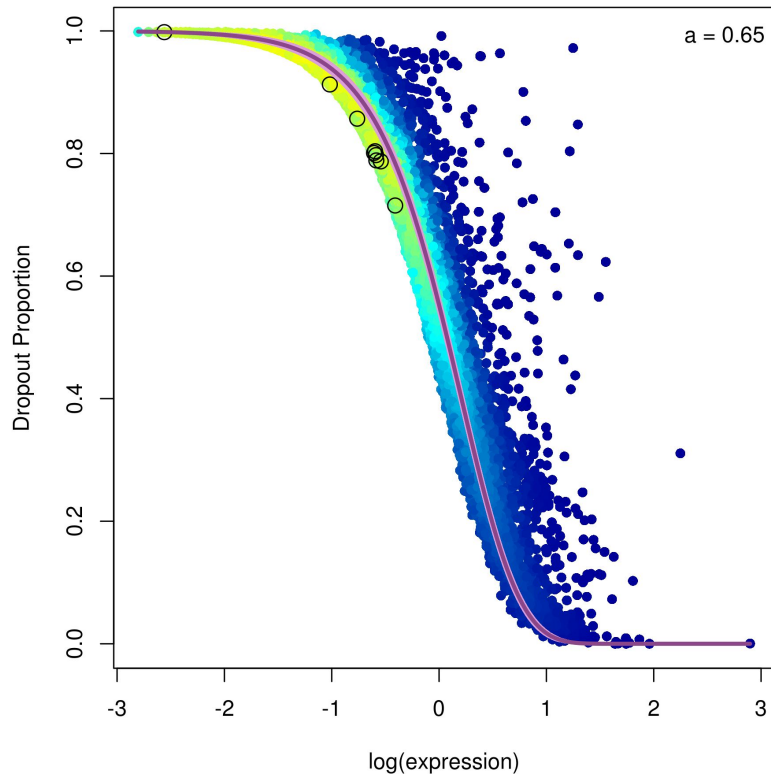
Reads



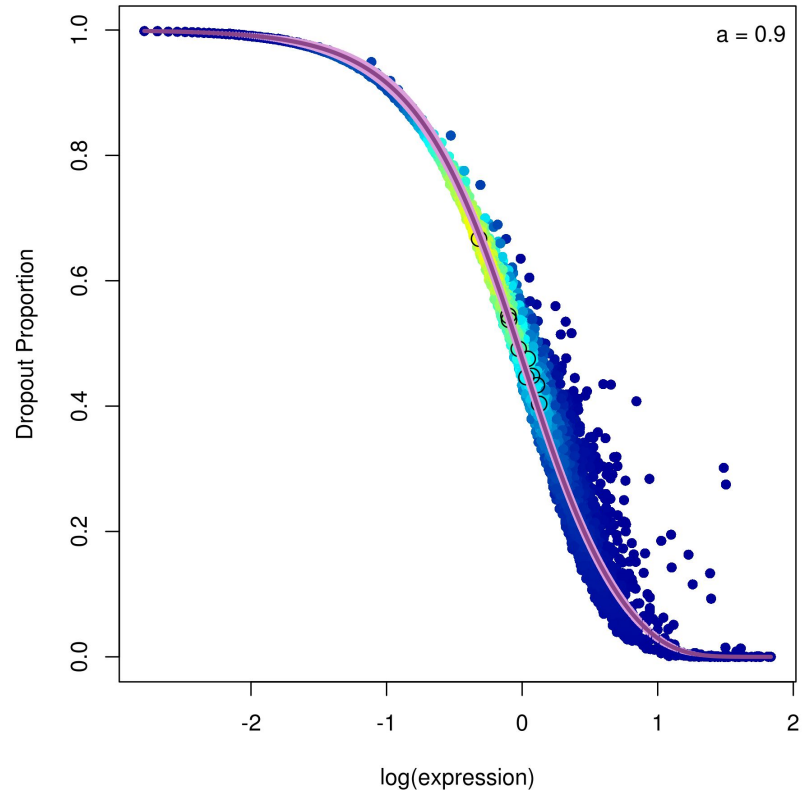


# Fitting the model to other UMI datasets

Linnarsson  $\alpha = 0.65$

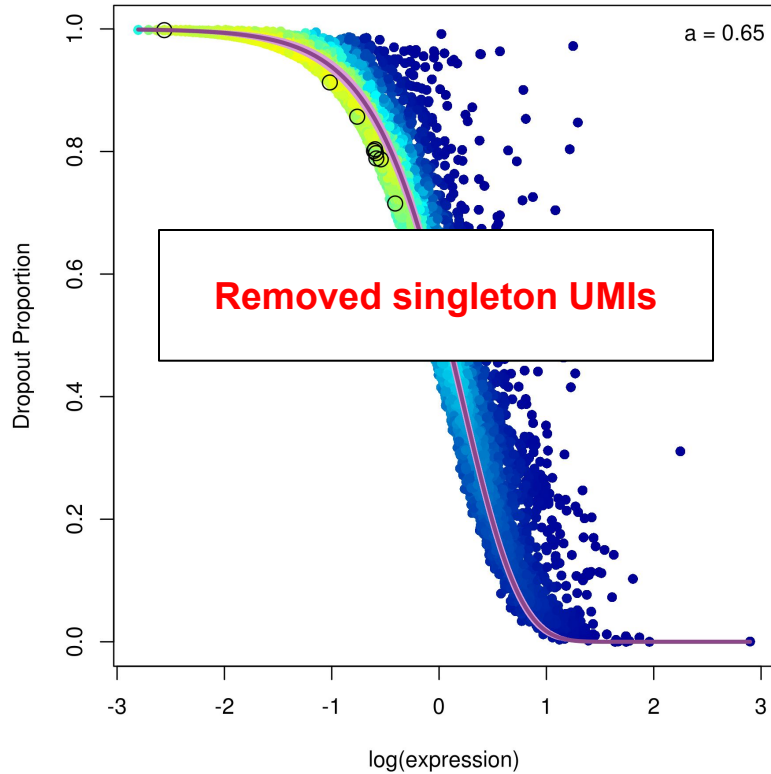


Kirschner  $\alpha = 0.90$

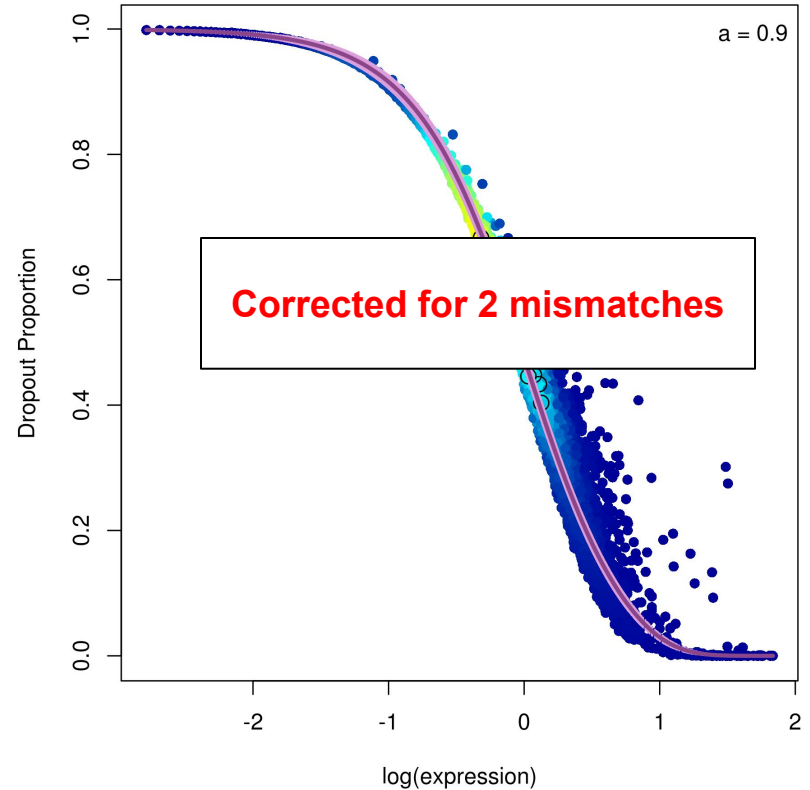


# Fitting the model to other UMI datasets

Linnarsson  $\alpha = 0.65$

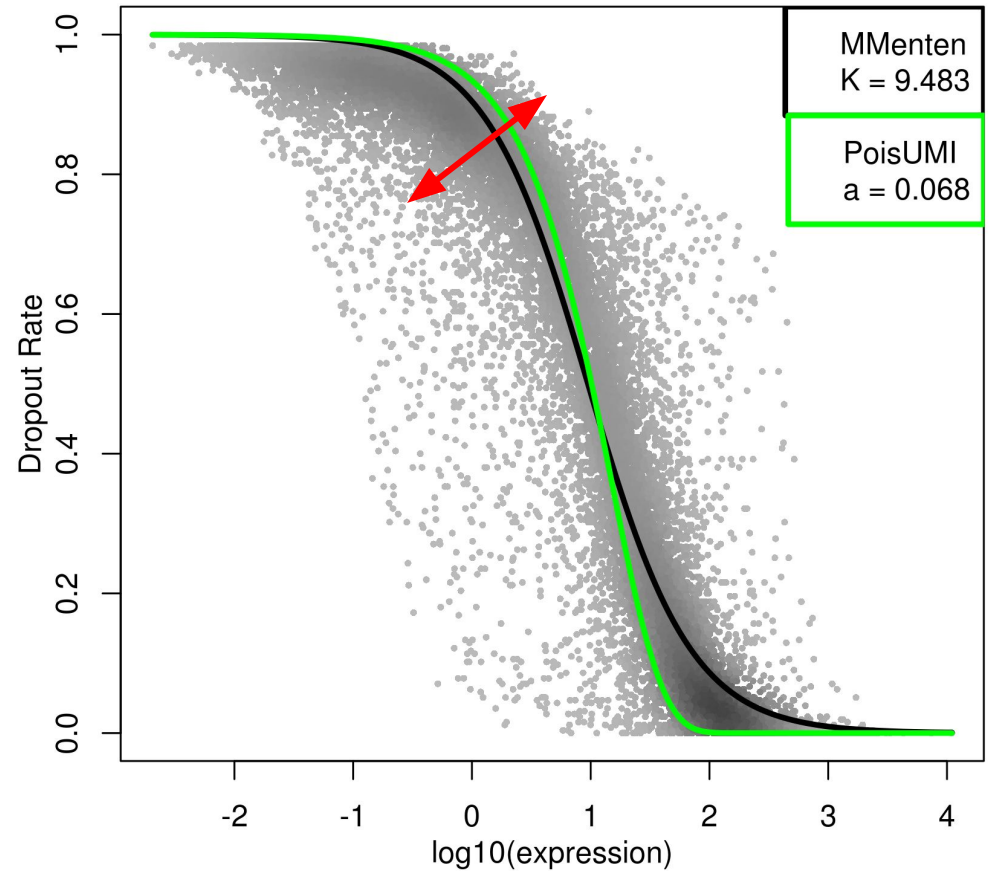


Kirschner  $\alpha = 0.90$



# Summary

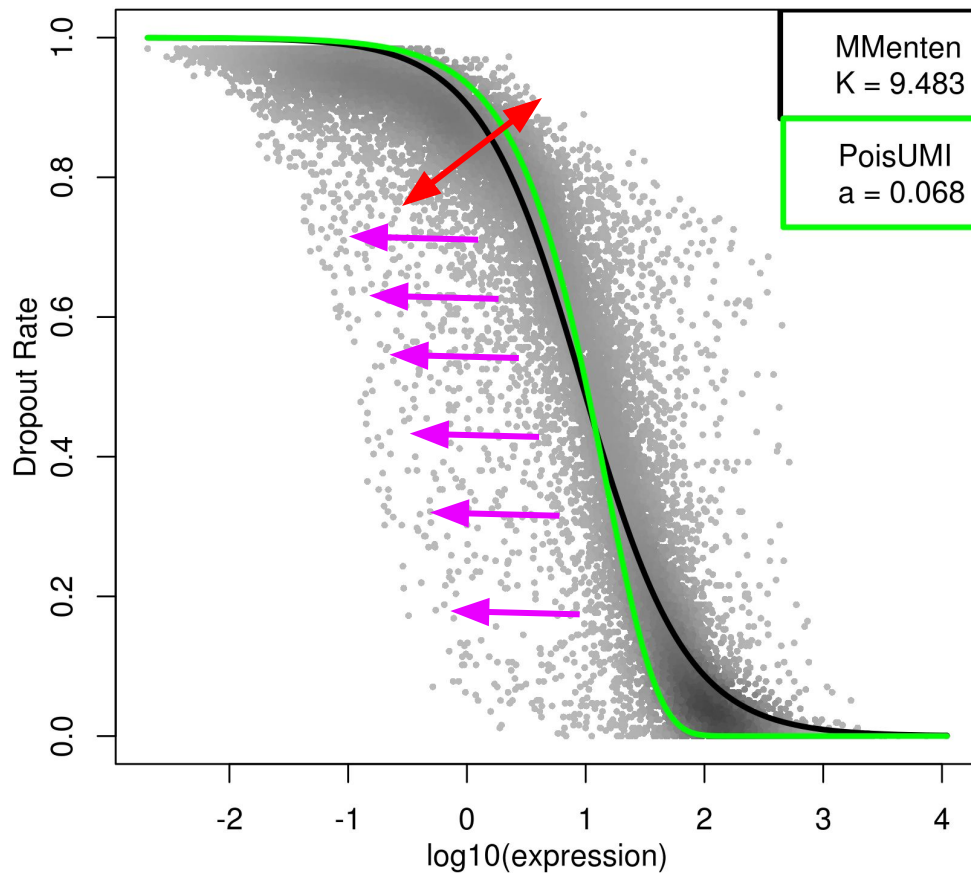
**Amplification noise**



# Summary

**Amplification noise**

**Mismapping /  
Miscounting**

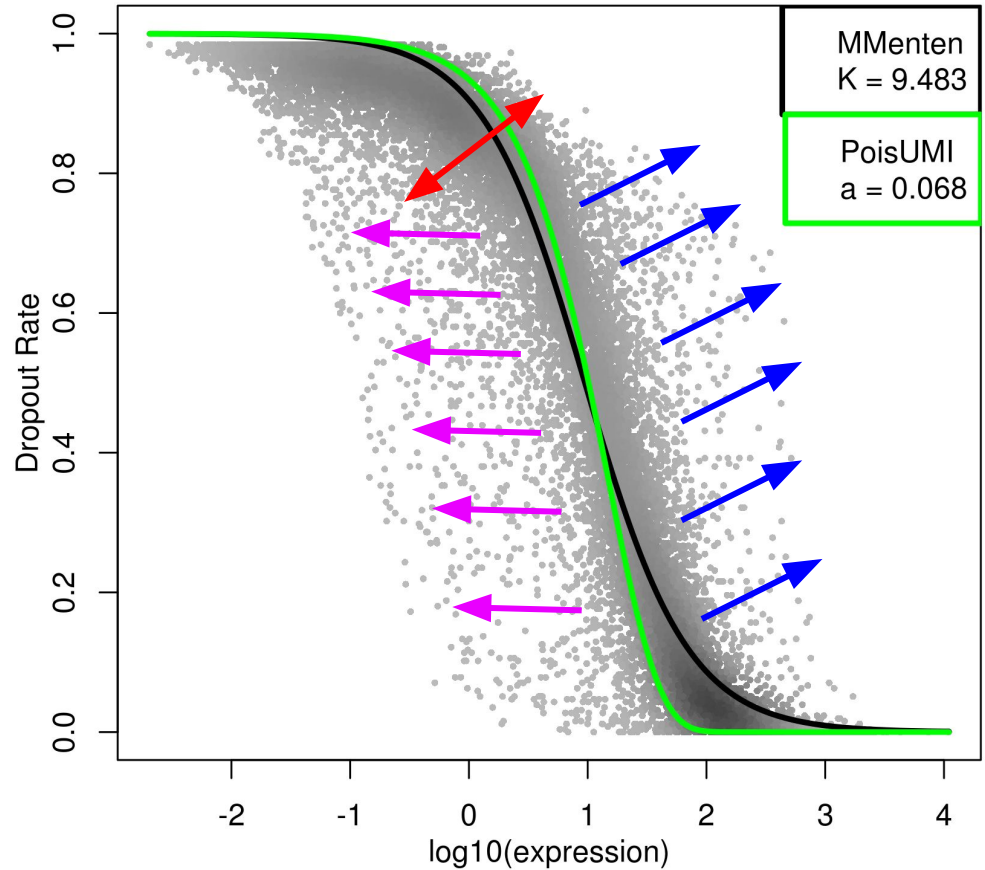


# Summary

**Amplification noise**

**Mismapping /  
Miscounting**

**Differential Expression**



# Acknowledgements

## Wellcome Trust Sanger Institute

Martin Hemberg

Vladimir Kiselev



## Availability

M3Drop : <https://github.com/tallulandrews/M3Drop>

PoissonUMIs: <https://github.com/tallulandrews/PoissonUMIs>

## EMBL Rome

Christophe Lancrin

Isabelle Bergiers

