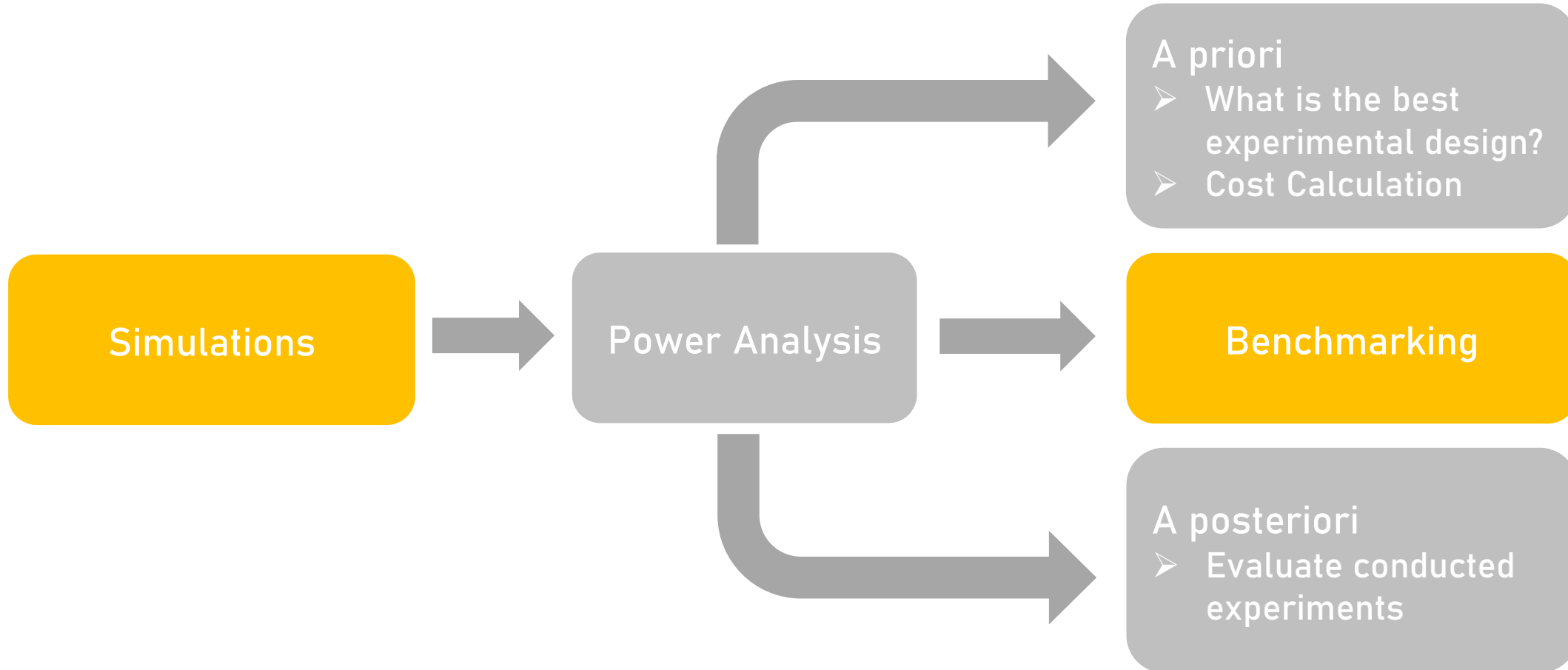




Power simulations for cell type identification from single cell RNA-seq data



Observed Expression

=

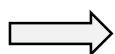
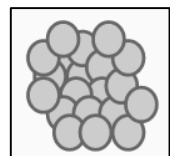
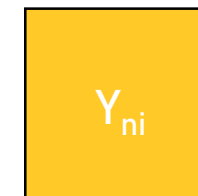
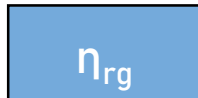
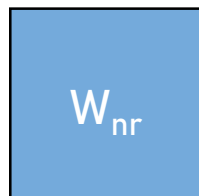
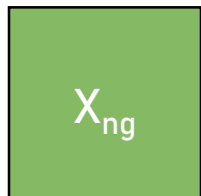
Unwanted Technical Factors

+

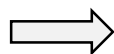
Unwanted Biological Factors

+

Biological Factors



X_{gn}	Cell 1	Cell 2	Cell 3	Cell n
Gene 1	0	80	0	10
Gene 2	10	0	5	0
Gene 3	25	65	39	0
Gene g	0	0	54	24



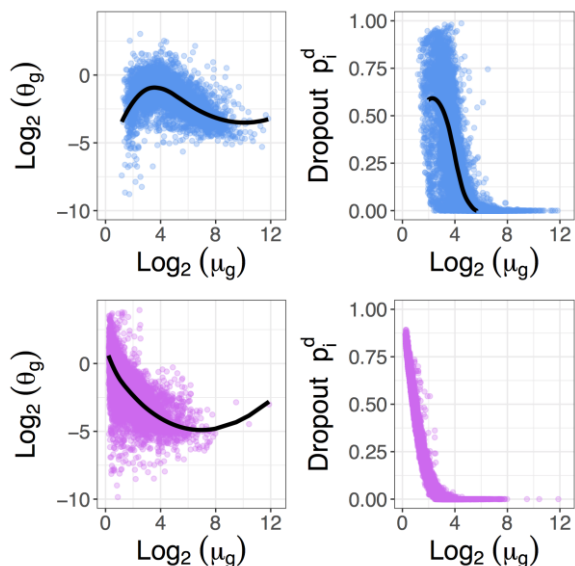
non-UMI Protocols

$$X_{ng} \sim \text{ZINB}(\mu_{g+}, \theta_{g+}, \pi_g)$$

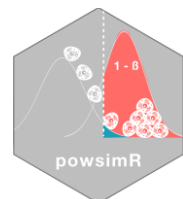
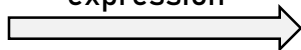
UMI Protocols

$$X_{ng} \sim \text{NB}(\mu_g, \theta_g)$$

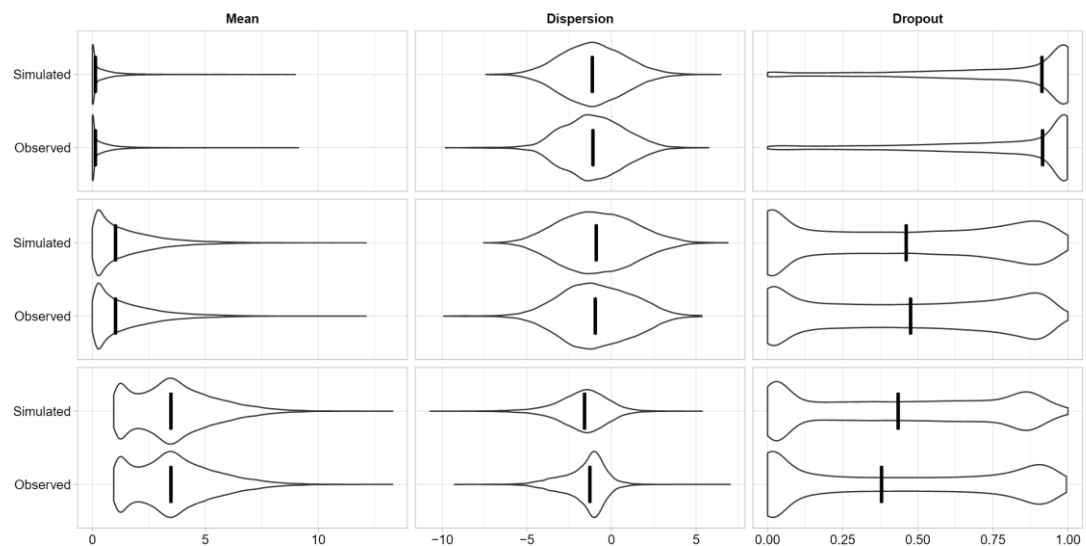
Fitting mean-variance and mean-dropout relationship:



Simulate expression

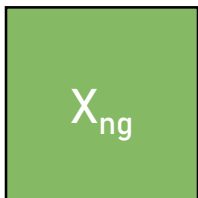


Vieth et al. 2017 Bioinformatics



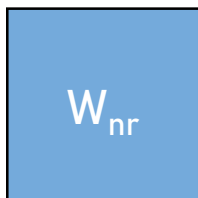
Vieth et al. 2019 Nat. Comm.

Observed Expression



=

Unwanted Technical Factors



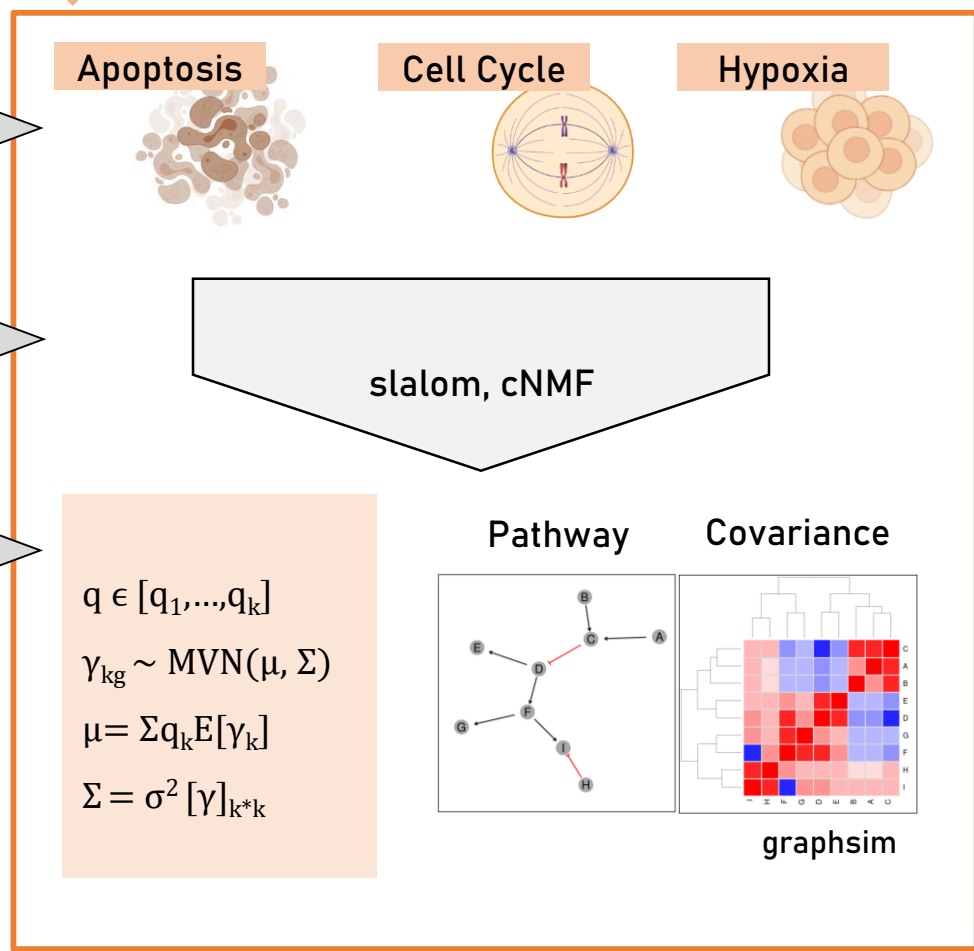
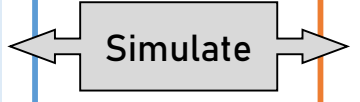
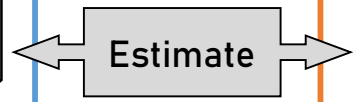
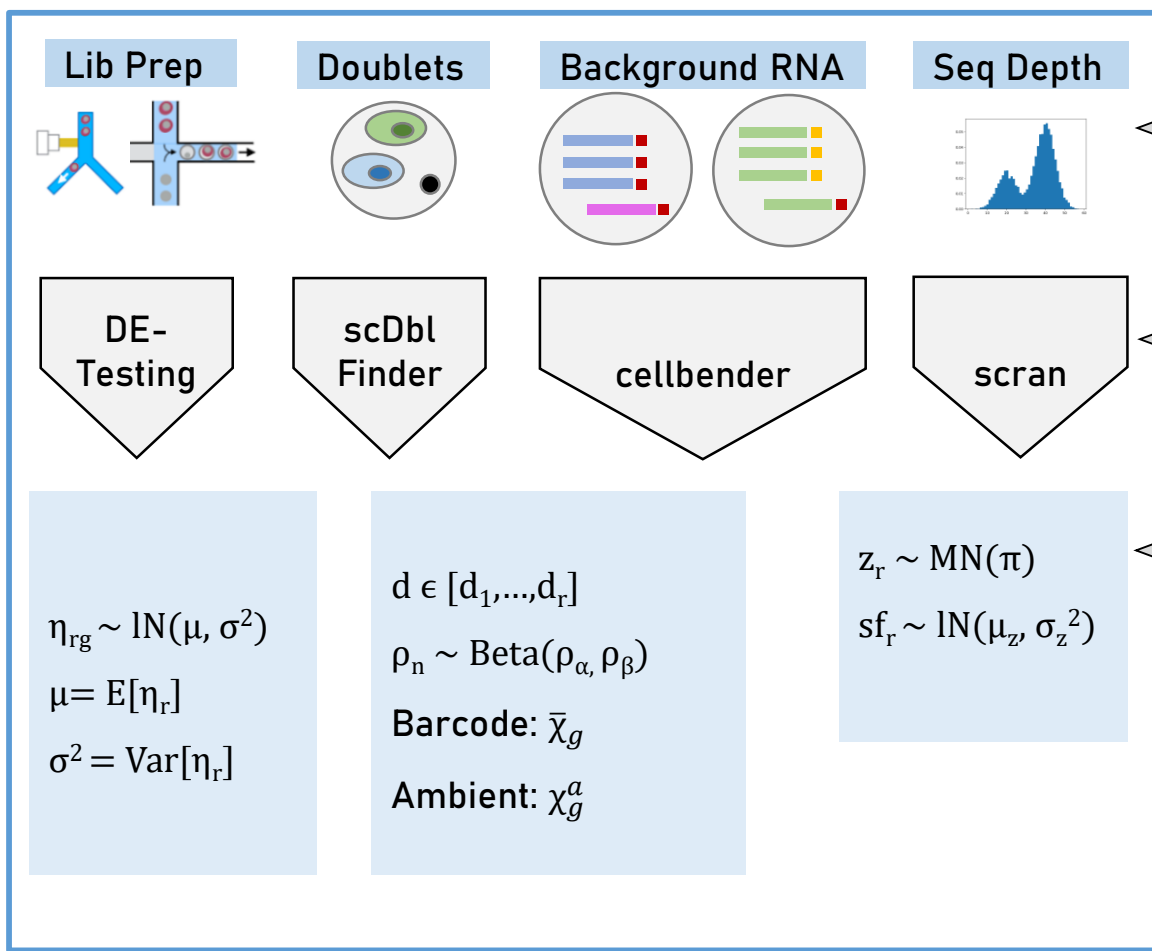
+

Unwanted Biological Factors



+

Biological Factors



Observed Expression

=

Unwanted Technical Factors

+

Unwanted Biological Factors

+

Biological Factors

$$X_{ng}$$

$$W_{nr}$$

$$\eta_{rg}$$

$$Z_{nk}$$

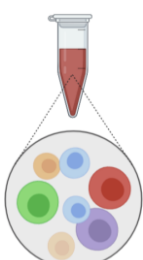
$$Y_{kg}$$

$$Y_{ni}$$

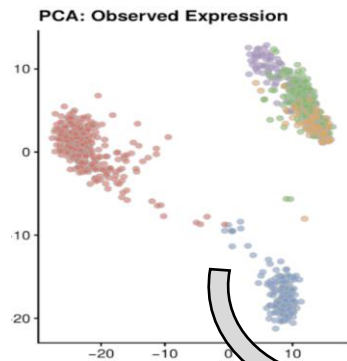
$$\beta_{ig}$$



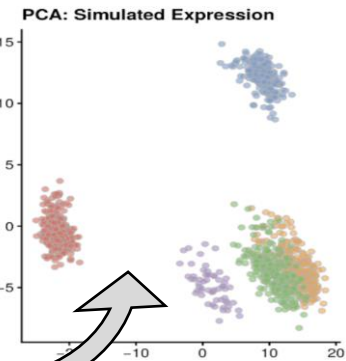
Blood



PCA: Observed Expression



PCA: Simulated Expression



$\beta_{ig} \sim \text{MVN}(\mu, \Sigma)$

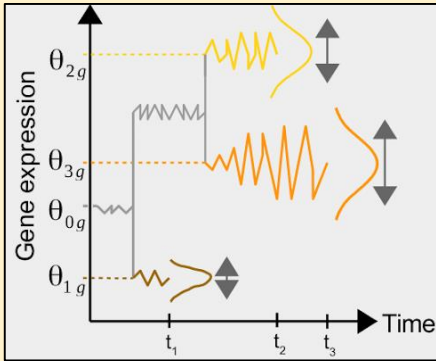
Ornstein Uhlenbeck Process:

$$dX_{ig}(t) = \alpha_g[\theta_{ig} - X_{ig}(t)]dt + \sigma_g dB(t)$$

$$\beta_{ig} \sim F(f_{ig}, \theta_{0g}, \alpha_g, \sigma_g, \phi)$$

$$\phi \text{ Tree} \quad \alpha_g \sim \Gamma(\alpha, \beta) \quad \sigma_g \sim N(\mu, \sigma^2)$$

$$f_{ig} \sim \Gamma(\alpha, \beta)$$

$$\theta_{0g} \sim N(\text{Par}(\alpha, \beta), \frac{\sigma_g^2}{2\alpha_g})$$



Gene expression

Time

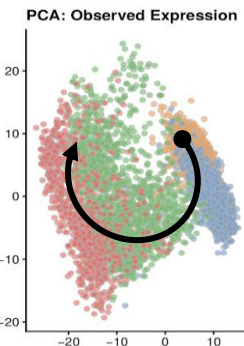
t_1, t_2, t_3

$\theta_{1g}, \theta_{0g}, \theta_{3g}, \theta_{2g}$

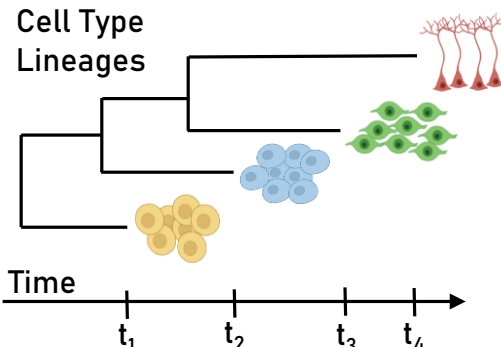
Differentiation



PCA: Observed Expression



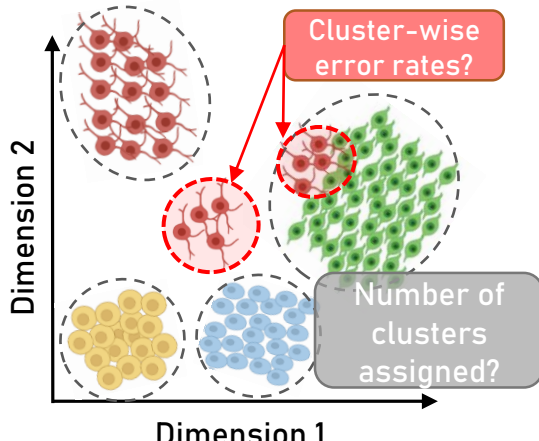
Cell Type Lineages



Time

t_1, t_2, t_3, t_4

Evaluate classification results



Dimension 2

Dimension 1

Cluster-wise error rates?

Number of clusters assigned?