

Power and Sample Size for RNA-seq Experiments

OCTRI BERD Research Forum

Jessica Minnier

OHSU

11/1/23

Goals

- Review power and sample size statistical definitions
 - Power, sample size, effect size, type I error, type II error
- Overview of RNA-seq data generation
- Required components for power calculation
- Overview of tools available
- Examples of power calculations
- Brief overview of more complex designs and scRNA-seq issues

Prerequisites

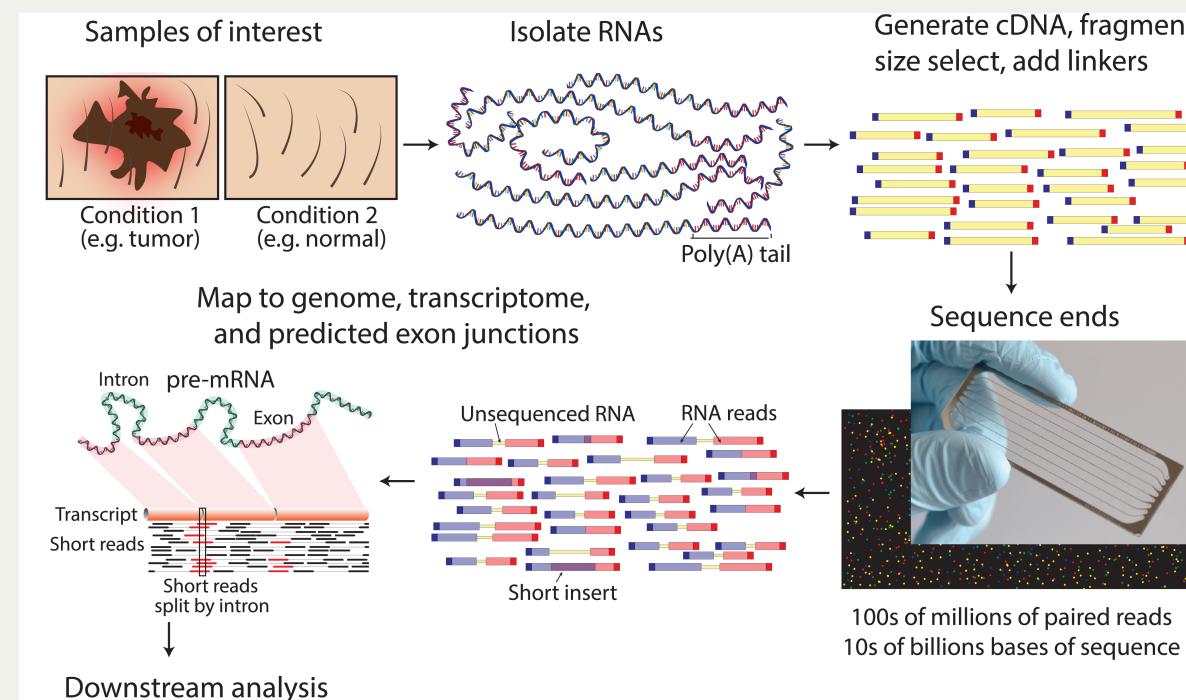
- Some knowledge of basic power and sample size calculations and concepts
 - PSS 101 from BERD, worth reviewing and watching!
- Familiarity with RNA-seq experiments; good overview: ([Conesa et al. 2016](#))
- Familiarity with statistical models and concepts such as regression, p-values, probability distributions (some review in context of RNA-seq: [Harvard Chan Bioinformatics Core \(HBC\) training's DGE and hypothesis testing lessons](#))
- Preferred: some experience with R/Bioconductor

RNA-seq

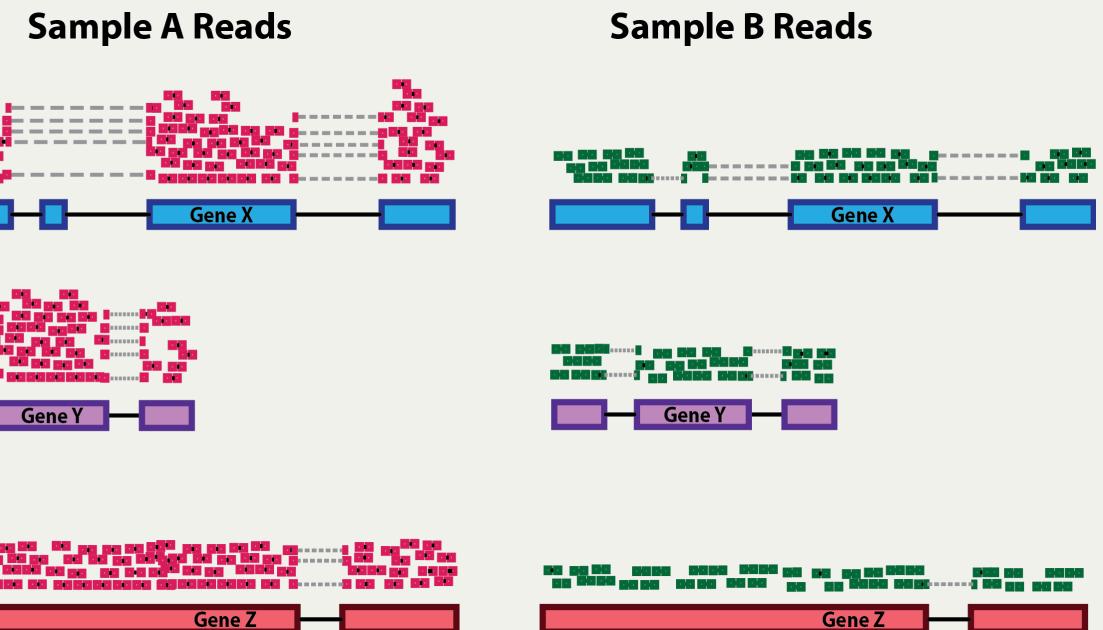
- Typically, transcriptomics: **gene expression profiling** of samples, using next-generation sequencing methods, reference genome
- RNA-seq can also examine alternative **gene spliced transcripts, mutations, levels of other RNA besides messenger RNA, etc**, but common aim is to evaluate gene expression measured by mRNA observed in sample tissue
- “**Bulk RNA-seq**” to distinguish from single cell RNA-seq, large populations of cells, mixed cell types
- **Alternative method: microarrays**, some of the following concepts still apply but distribution of expression measures is not count-based, but more related to continuous normal distribution ([Lee and Whitmore 2002](#))

RNA-seq data generation

NGS or massively parallel sequencing: sample preparation, mRNA fragmentation, reverse transcription to complementary DNA, map cDNA to reference genome



(Griffith et al. 2015)



(Mistry et al. 2021)

What does the data look like?

Raw read sequencing data (large, fastq files) → Count data (gene ID x sample, matrix in .csv)

From a downstream statistical analysis standpoint:

- Level of measurement is the gene (or transcript) or “feature”
- Measurement/outcome is “gene count”: how many reads aligned to that gene’s RNA sequence after alignment to a reference genome? → count matrix
- Need a statistical model for **count data** that models the variability appropriately

samples: want to see if differences across condition are significant (w.r.t. biological and technical variation)

features (e.g. genes)

	SRR1039508	SRR1039509	SRR1039512	SRR1039513	SRR1039516
ENSG00000000003	679	448	873	408	1138
ENSG00000000005	0	0	0	0	0
ENSG00000000419	467	515	621	365	587
ENSG00000000457	260	211	263	164	245
ENSG00000000460	60	55	40	35	78

Simplest Study: Two groups

Suppose we are studying two treatment groups (treatment vs. control), and we want to know which genes are differentially expressed between these two groups.

Two groups: Two experimental groups, multiple *biological* replicates within each group. The two groups contain different samples (i.e., not paired, not the same samples over time).

Differentially expressed genes (DEGs): A gene is “differentially expressed” between the treatment vs. control if there is a difference observed in read counts or expression levels.

Question of interest: How many samples (biological replicates) in each group do we need to “detect a difference”?

- ... with some level of confidence (power), restricting false positives (type I error, FDR)?

Simplest Study: Two groups

Question of interest: How many samples do we need to “detect a difference”?

(Some) Information we need:

- Definition of “how different” (i.e. effect size, fold change)
- Desired level of power and type I error/p-value expectations (“significance level”)
- Information about sequencing depth, average read count
- Gene expression info: expected variability/dispersion of gene expression levels in each group, number of genes we expect to measure/observe
- Statistical model and test, method for multiple comparison adjustment
- “How much signal?” % genes differentially expressed above some fold change

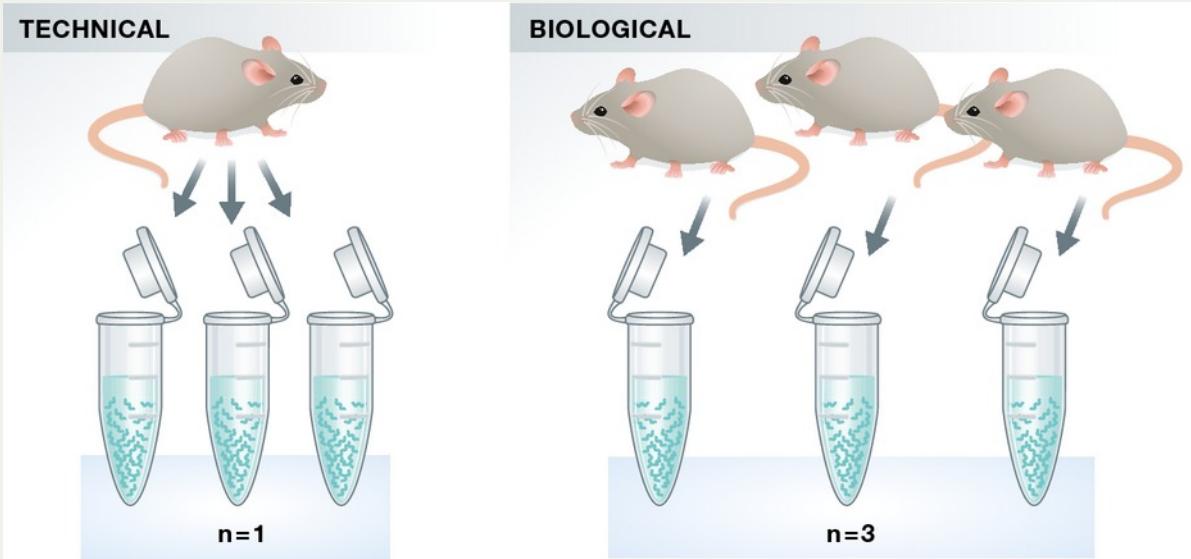
Reminder: Statistical Power Concepts

https://bit.ly/berd_pss_rnaseq

Sample Size

More biological replicates leads to:

- better estimates of variation (gene, biological, sample-to-sample)
- identify outliers or possible sources of technical variation (batch effects)
- improve precision of estimates
- observe low abundance genes



Klaus (2015) and Mistry et al. (2021)

Components of Sample Size

In general, need to know 3 of the 4 to determine the 4th:

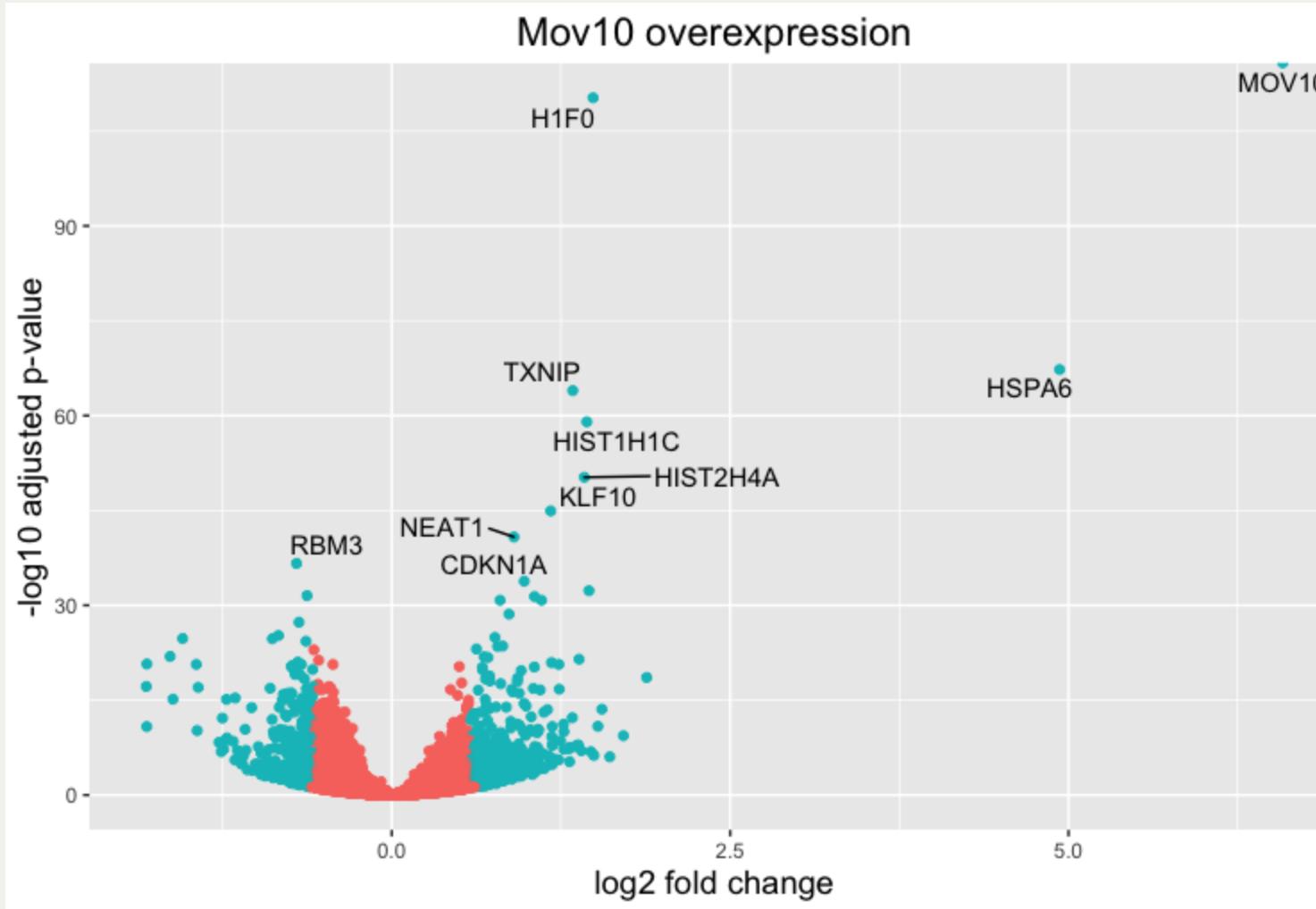
Do We Know?	Measure	Definition
??	Effect Size	Magnitude of difference or fold change
??	Sample Size	N total, n per group
0.05, 0.01	Type I Error / Significance level	α = probability of rejecting null hypothesis when it is true <i>need to adjust for multiple comparisons!</i>
0.9, 0.8	Power	$1 - \beta$ = 1 - Type II error = probability of rejecting null hypothesis when it is false <i>need to consider multiple comparisons!</i>

Effect size = Mean Expression Ratio (FC)

In a two group comparison, effect size is usually *fold change* (or logFC) for an individual gene:

$$FC_{\text{gene } X} = (\text{mean expression of gene } X \text{ in trt}) / (\text{mean expression of gene } X \text{ in control})$$

- May observe a large range of fold changes, or fold changes may be close to 1 for all genes
- FC measures the average size of the difference between groups, not variability
- Alternative effect size: biological coefficient of variation (BCV)



Volcano plot (hbctraining)

(Mistry et al. 2021)

Type I error (with multiple comparisons)

“The worst error” Truth = No difference, Test conclusion = there is a difference

- Typical power calculations are for one test, we could have tens of thousands of tests for each gene
- Multiple comparisons/multiple testing problem: If we use $p\text{-value} < 0.05$, for one test there is a 5% chance we have a false positive, if we test 20,000 genes and use $p < 0.05$, we would expect to detect 1000 false positives by chance; if we found 3000 DEGs total, one third are false positives
- Use “False Discovery Rate” = FDR as “adjusted” p-values
 - FDR = proportion of expected false positives in our set of DEGs, often controlled at 5% or 10%
 - Benjamini-Hochberg method common, alternative Storey/Q-value method, see ([Mistry et al. 2021](#)) for summary

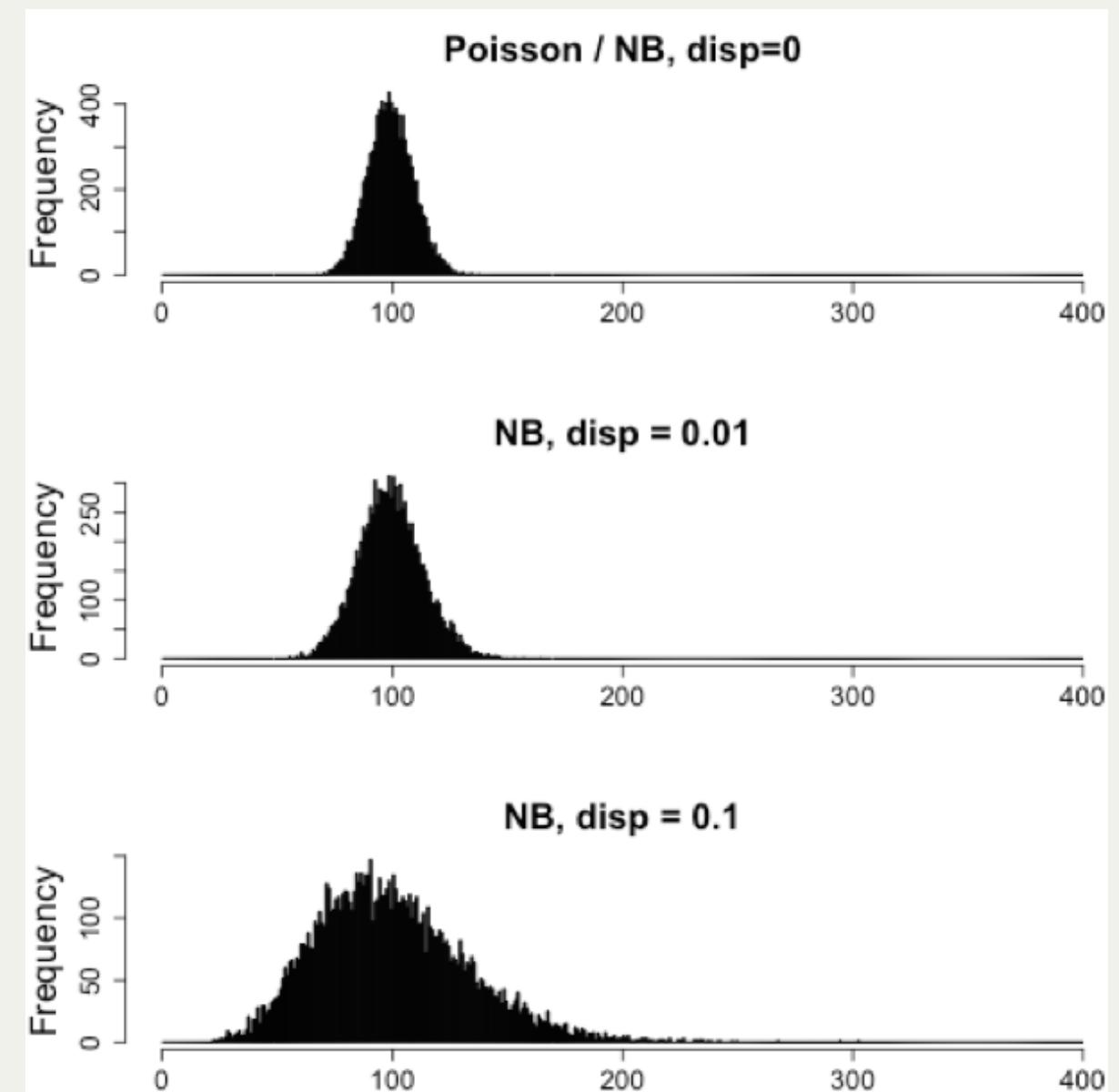
Power (with multiple comparisons)

- With one test, power = $1 - \text{type II error}$ for a given effect size
 - Power = probability that our test is “significant” when the truth is that there is a difference
- With multiple tests/genes, could have many definitions
 - Probability we detect at least X% of genes that are truly different (at least some FC)
 - Probability that we detect *all* genes that are truly different (at least some FC)
 - **Probability that we detect *one specific gene* that is truly different by some minimum FC**
- Define “significant” as FDR < cutoff

Gene Expression/Abundance Distributions

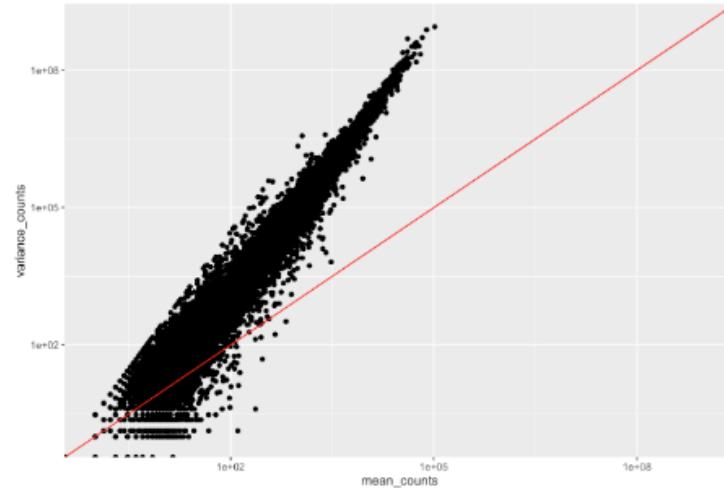
https://bit.ly/berd_pss_rnaseq

- Count data = discrete data (0, 1, 2...)
- Need **probability distribution** to match the type of data, such as: Poisson distribution, Negative Binomial distribution, Poisson-gamma mixture distribution
- Poisson model for count data assumes mean = variance, but RNA-seq has other **sources of variation** than the counting process → NB with dispersion parameter, other models
- Lower count genes can be harder to detect/measure/test
- Sequencing depth (total number of reads, i.e. 5-200 million) influences **average gene count**

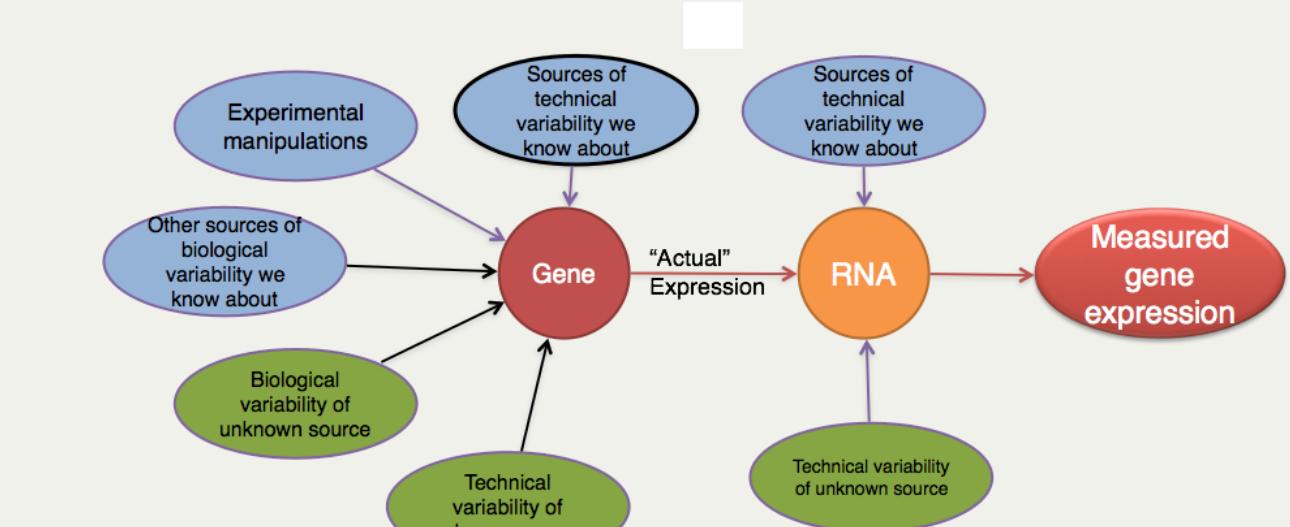


(Mistry et al. 2021)

Variability



1. The mean is not equal to the variance (the scatter of data points does not fall on the diagonal).
2. For the genes with high mean expression, the variance across replicates tends to be greater than the mean (scatter is above the red line).
3. For the genes with low mean expression we see quite a bit of scatter. We usually refer to this as "heteroscedasticity". That is, for a given expression level in the low range we observe a lot of variability in the variance values.



Courtesy of Paul Pavlidis, UBC

(Mistry et al. 2021)

Factors affecting gene expression + power

- *Biological differences and variation*
- Sequencing depth/coverage
- Gene length
- Total sample RNA output, variation between samples within biological groups
- Variance / dispersion of gene abundance distributions
 - dispersion: parameter that defines how far we expect observed count will deviate from mean value (estimate from a model)

Sequencing depth vs. N

https://bit.ly/berd_pss_rnaseq

Typical analysis

- Common tools: `DESeq2`, `edgeR`, `baySeq`, and `limma / voom` packages in R/Bioconductor
- Normalization for gene length, sequencing depth (TMM, RPKM)
- `DESeq2`, `edgeR`, `baySeq` use generalized linear models (GLMs) with negative binomial (NB) distribution to fit the data:

raw count for gene i , sample j

The mean is taken as “normalized counts” scaled by a normalization factor

one dispersion per gene

$$K_{ij} \sim \text{NB}(s_{ij}q_{ij}, \alpha_i)$$

- `Voom` models mean-variance relationships to use normal based models in `limma` (linear model with Empirical Bayes variance smoothing)
- Hypothesis test options: Wald Test, Likelihood Ratio Test, and similar

Power calculation software/tools

- See review in Jeon et al. (2023) and comparisons in Poplawski and Binder (2018)
- RNASeqPower Bioconductor/R package
 - closed form equation based on Score statistic, NB model
- RnaSeqSampleSize Bioconductor/R package and web/Shiny app
 - NB model, gene-specific power function, can overestimate sample size
- ssizRNA R package
 - NB model with normal-based test statistic via voom
- PROPER R package
 - NB model simulation based

Based on above reviews and simulation studies, ssizRNA and PROPER were most recommended

Example R code, ssizeRNA

Adapted from [ssizeRNA vignette](#) and Suppl Material S1 Jeon et al. (2023)

Arguments for `ssize_single()` function, assuming all genes have the same characteristics:

- Total number of genes: `G = 10000`
- Proportion of non-DE “null” genes: `pi0 = 0.8`
- FDR level to control: `fdr = 0.05`
- Desired average power to achieve: `power = 0.8`
- Average read count for each gene in control group: `mu = 10`
- Dispersion parameter for each gene: `disp = 0.1`
- Fold change for each gene: `fc = 2`

```
1 ## Install ssizeRNA package if needed
2 # install.packages("ssizeRNA")
3
4 ## Load package
5 library(ssizeRNA)
6
7 size_out <- ssizeRNA_single(nGenes = 10000,
8                               pi0 = 0.8,
9                               fdr = 0.05,
```

https://bit.ly/berd_pss_rnaseq

```

10          power = 0.8,
11          mu = 10,
12          disp = 0.1,
13          fc = 2)

1 # Output variables:
2 # ssize: sample sizes (for each treatment) at which desired power is first reached. size$ssize
3 size_out$ssize

  pi0 ssize      power
[1,] 0.8    14 0.8343262

1 # power: power calculations with corresponding sample sizes.
2 size_out$power

  n      0.8
[1,] 3 0.0000000
[2,] 4 0.0000000
[3,] 5 0.0056349
[4,] 6 0.0953837
[5,] 7 0.2332468
[6,] 8 0.3697507
[7,] 9 0.4895573
[8,] 10 0.5898277
[9,] 11 0.6719140
[10,] 12 0.7383110
[11,] 13 0.7916434
[12,] 14 0.8343262
[13,] 15 0.8683667
[14,] 16 0.8954845
[15,] 17 0.9170581
[16,] 18 0.9341977
[17,] 19 0.9478082
[18,] 20 0.9586136
[19,] 21 0.9671815
[20,] 22 0.9739771
[21,] 23 0.9793624

```

Typical inputs

- Dispersion often depends on animal model, might be 0.1 for mouse data, 0.2-0.5 for human data
- Read count depends on sequencing depth and distribution of abundance, if using average can be low, typically 5-30, but in reality will vary widely by gene
- Proportion of non-DEG, varies greatly based on experiment/biological difference expectations
- Fold change, the magic number :) 1.5, 2, 2.5...
- Power 0.9 and FDR 0.05, more “relaxed” adjust power 0.8-0.85 and/or FDR 0.1

Interactive web app

- RnaSeqSampleSize Bioconductor/R package and web/Shiny app

n: Sample Size
15

Sample Size Estimation by single parameter Sample Size Estimation by prior data Generate Power Curves Parameters Optimization

f: FDR level
0.05

w: Ratio of normalization factors between two groups
1

m: Total number of genes for testing
10000

m1: Expected number of prognostic genes
200

rho: Minimum fold changes for prognostic genes between two groups
2

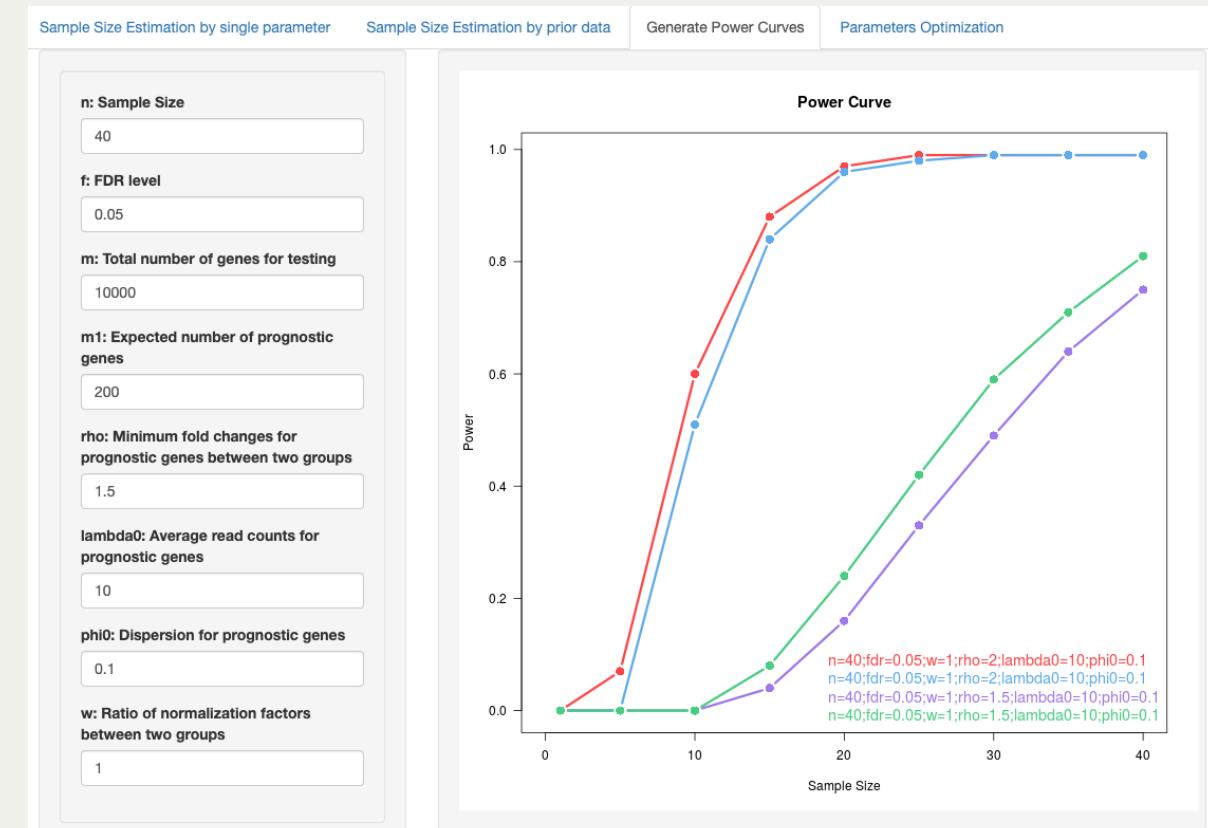
lambda0: Average read counts for prognostic genes
10

phi0: Dispersion for prognostic genes
0.1

0.88

Description:

We are planning a RNA sequencing experiment with 15 experimental subjects in each group to identify differential gene expression between two groups. Prior data indicates that the minimum average read counts among the prognostic genes in the control group is 10, the maximum dispersion is 0.1, and the ratio of the geometric mean of normalization factors is 1. Suppose that the total number of genes for testing is 10000 and the top 200 genes are prognostic. If the desired minimum fold change is 2, we will be able to reject the null hypothesis that the population means of the two groups are equal with probability (power) 0.88 using exact test. The FDR associated with this test is 0.05.



Pilot data or public data

- Ideally, pilot data specific to the experiment would be used to determine expected variability / dispersion and effect sizes
- `ssizeRNA` gives examples on using existing data to estimate `mu` and `disp` across all the genes, then use `ssizeRNA_vary()` function, and see below from Suppl Material S1 Jeon et al. (2023) for R code
- `RnaSeqSampleSize` has TCGA data examples, see `vignette`, and also `web app`

```
1 ## Install other packages if needed.
2 # install.packages("BiocManager")
3 # BiocManager::install("edgeR") # cf. https://bioconductor.org/packages/release/bioc/html/edge_R.html
4 # BiocManager::install("Biobase") # cf. https://bioconductor.org/packages/release/bioc/html/Biobase.html
5
6 library(edgeR)
7 library(Biobase)
8 library(ssizeRNA)
9
10 ## Example data saved in ssizerNA package:
11 ## Step 2-1. load hammer dataset (Hammer, P. et al., 2010)
12 ## two group rat data
13 data(hammer.eset)
14 counts <- exprs(hammer.eset)[,phenoData(hammer.eset)$Time=="2 weeks"]
15 counts <- counts[rowSums(counts) > 0, ]
16 dim(counts)
[1] 18463      4
1 trt <- hammer.eset$protocol[which(hammer.eset$Time=="2 weeks")]
2 https://bit.ly/berd_pss_rnaseq
```

```
3 ## After generating count data with column names of control and treatment, you may estimate the parameters of mu and disp using the f
4 ## mu: average read count in the control group
5 ## The following apply function averages the count values for each gene in the control group.
6 mu <- apply(counts[, trt == "control"], 1, mean)
7 summary(mu)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.0	9.5	146.5	952.0	662.5	386345.0

```
1 ## disp: dispersion parameters estimates using the edgeR package with count data.
2 d <- DGEList(counts)
3 d <- calcNormFactors(d)
4 d <- estimateCommonDisp(d)
5 d <- estimateTagwiseDisp(d)
6 disp <- d$tagwise.dispersion
7 summary(disp)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.05566	0.06991	0.11093	0.13949	0.18397	1.01143

```
1 size_out <- ssizeRNA_vary(nGenes = 10000,
2                               mu = mu, # estimated above
3                               disp = disp, # estimated above
4                               fc = 2,
5                               up = 0.5,
6                               fdr = 0.05,
7                               power = 0.8,
8                               maxN = 35)
```

```
1 # Output variables:
2 # ssize: sample sizes (for each treatment) at which desired power is first reached. size$ssize
3 size_out$ssize
```

```
pi0 ssize      power
[1,] 0.8     9 0.8434739
```

```
1 # power: power calculations with corresponding sample sizes.
2 size_out$power
```

```
n      0.8
[1,] 3 0.0000000
[2,] 4 0.1495363
[3,] 5 0.4275414
```

```
[4,] 6 0.6056247  
[5,] 7 0.7184589  
[6,] 8 0.7927465  
[7,] 9 0.8434739  
[8,] 10 0.8792088  
[9,] 11 0.9050462  
[10,] 12 0.9241417  
[11,] 13 0.9385395  
[12,] 14 0.9495731  
[13,] 15 0.9581587  
[14,] 16 0.9649286  
[15,] 17 0.9703306  
[16,] 18 0.9746886  
[17,] 19 0.9782380  
[18,] 20 0.9811584  
[19,] 21 0.9835790  
[20,] 22 0.9856014  
[21,] 23 0.9873035
```

More complex studies

What are the additional experimental factors?

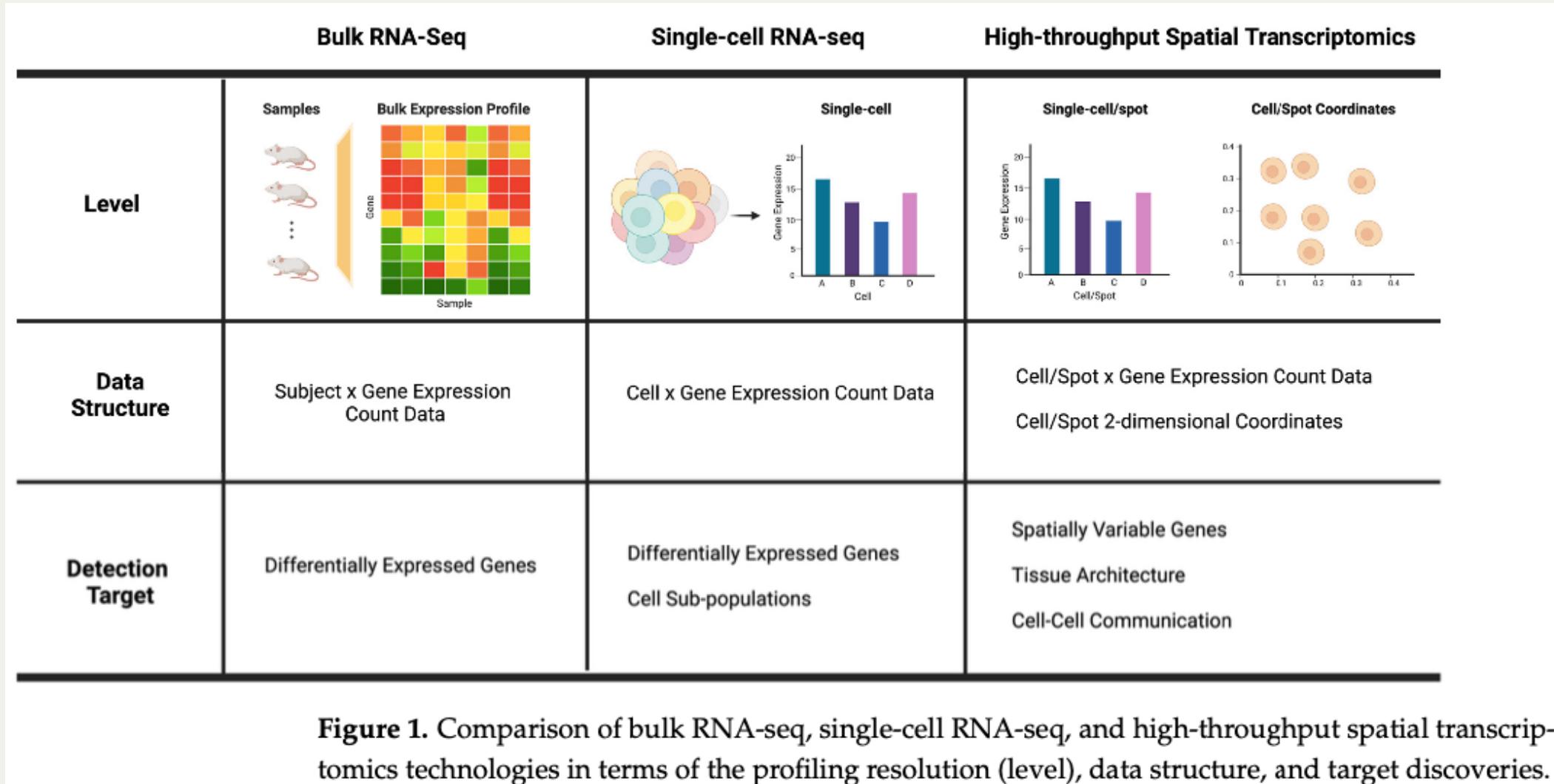
- Multiple groups, multiple treatment combinations
- Interaction effects
- Continuous predictor rather than groups
- Time or repeated measure component
- Technical replicates
- Batch effects or technical variation

→ Likely need to use **simulation based methods** for power / sample size

Pathway analyses?

Even more complex to calculate power with simulations and many assumptions

Single-cell RNA-seq or spatial transcriptomics



https://bit.ly/berd_pss_rnaseq

scRNA-seq

https://bit.ly/berd_pss_rnaseq

- scRNA-seq have more 0 counts than bulk which affects dispersions/variation, require “zero-inflated” model
- Power analysis for DEG (1) between biological groups/experimental conditions for a specific cell type or (2) between cell types or sub-populations within a sample
- Power analyses for identifying cell sub-populations, either within a single tissue, or proportional differences across experimental conditions

Table 2. A table with information about different software tools for scRNA-seq power analysis with two distinct detection targets. Experimental Factors: cell number (1), individual number (2), Sequencing depth (3).

Detection Target	# of Samples	Tool Name	Experimental Factor	Software	Model	Power Assessment
Cell sub-population	Single sample	'SCOPIT' [37]	(1)	R package & Web application	Multinomial	Analytical
		'howmanycells'		Web application	Negative Binomial	
	'Sensei' [38]			Web application	Beta Binomial	
		'scPOST' [39]	(1), (2)	R package	Linear mixed model	Simulation-based
DEG	Multi sample	'scPower' [40]		R package & Web server		Pseudobulk
		'hierarchicell' [41]	(1), (2), (3)		Negative Binomial	
	Single sample	'powsimR' [42]	(1)	R package	A mixture of zero-inflated Poisson and log-normal Poisson distributions	Simulation-based
		'POWSC' [43]				
		'scDesign' [44]	(1), (3)		Gamma-Normal mixture model	
				https://bit.ly/berd_pss_rnaseq		

Conclusions

Important Factors (bulk RNA-seq DEG)

- Distributions: average (or distribution of) read counts, as well as average (or distribution of, or maximum) dispersion
- Proportion of differentially expressed genes
- Fold change effect size
- False discovery rate and power
- Experimental design complexity (beyond two groups, consult statistician)

Primary resources

- Jeon et al. (2023) and Mistry et al. (2021) tutorial

Thank you

- Recording will be available at OCTRI BERD Research Forum Website
- email: minnier [at] ohsu.edu
- Thanks for consultation / tips from ONPRC BBC Director Suzi Fei and all I've learned from IGL/MPSSR Directors Bob Searles and Chris Harrington
- Please fill out the survey that Amy Laird will send you!
- More "omics" PSS series to come
- Need stats help? Contact BDP or Knight BSR

References

- Conesa, Ana, Pedro Madrigal, Sonia Tarazona, David Gomez-Cabrero, Alejandra Cervera, Andrew McPherson, Michał Wojciech Szcześniak, et al. 2016. “A Survey of Best Practices for RNA-Seq Data Analysis.” *Genome Biology* 17 (1): 1–19.
- Griffith, Malachi, Jason R Walker, Nicholas C Spies, Benjamin J Ainscough, and Obi L Griffith. 2015. “Informatics for RNA Sequencing: A Web Resource for Analysis on the Cloud.” *PLoS Computational Biology* 11 (8): e1004393.

- Jeon, Hyeongseon, Juan Xie, Yeseul Jeon, Kyeong Joo Jung, Arkobrato Gupta, Won Chang, and Dongjun Chung. 2023. "Statistical Power Analysis for Designing Bulk, Single-Cell, and Spatial Transcriptomics Experiments: Review, Tutorial, and Perspectives." *Biomolecules* 13 (2): 221.
- Klaus, Bernd. 2015. "Statistical Relevance—Relevant Statistics, Part i." *The EMBO Journal* 34 (22): 2727–30.
- Lee, Mei-Ling Ting, and George Alex Whitmore. 2002. "Power and Sample Size for DNA Microarray Studies." *Statistics in Medicine* 21 (23): 3543–70.
- Liu, Yuwen, Jie Zhou, and Kevin P White. 2014. "RNA-Seq Differential Expression Studies: More Sequence or More Replication?" *Bioinformatics* 30 (3): 301–4.
- Mistry, Meeta, Mary Piper, Jihe Liu, and Radhika Khetani. 2021. "hbctraining/DGE_workshop_salmon_online: Differential Gene Expression Workshop Lessons from HCBC (first release)." Zenodo.
- Poplawski, Alicia, and Harald Binder. 2018. "Feasibility of Sample Size Calculation for RNA-Seq Studies." *Briefings in Bioinformatics* 19 (4): 713–20.

https://bit.ly/berd_pss_rnaseq