

Responsible Research

making experiments better

Tobias Straub

tobias.straub@lmu.de

Biomedizinisches Centrum, LMU München

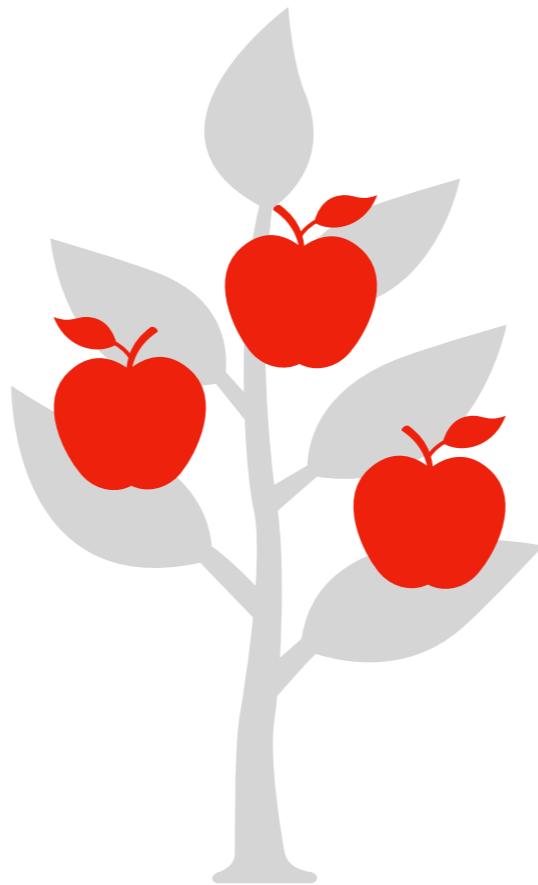
Why am I doing this?



- Purpose: do sth, measure/observe sth, then **predict!**
- Types of experiments
 - Descriptive versus **controlled**
 - Explorative versus **confirmatory**
 - Obscure versus **transparent**
 - Anecdotal versus **representative**
 - Context: Isolated versus **ensemble**
- Types of outcomes
 - **Convincing** versus unconvincing
 - **Reproducible** versus irreproducible
 - Career-changing breakthrough versus failed
- Matters: Data management, documentation, visualisation, interpretation

Taking measurements

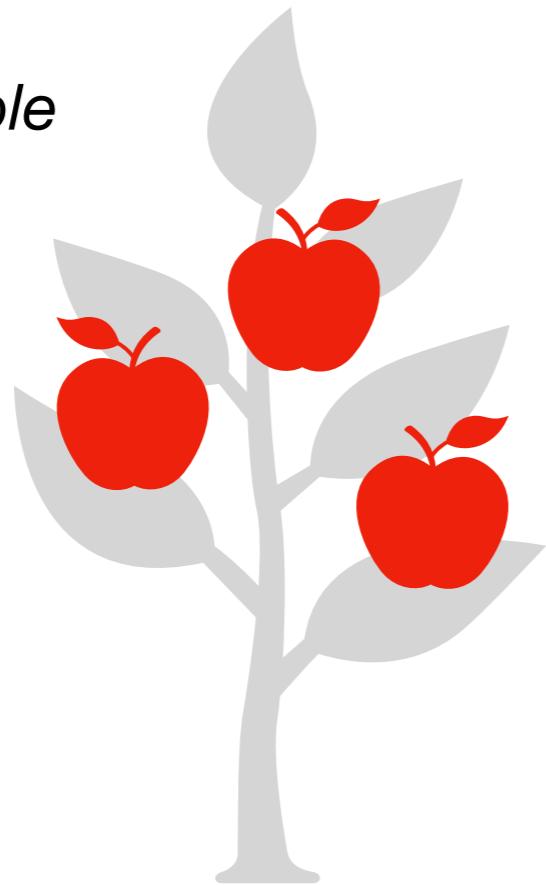
descriptive versus controlled



Number of apples
endpoint

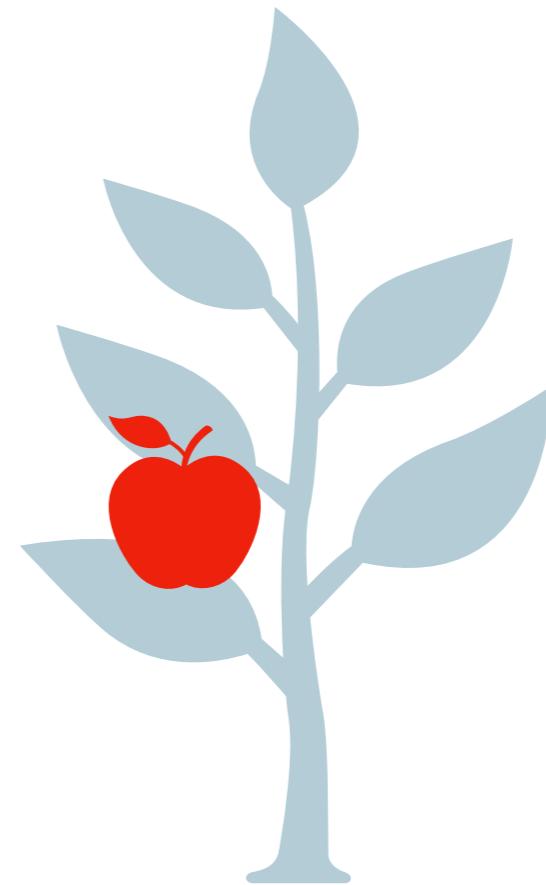
descriptive versus controlled

Number of apples
endpoint
dependent variable



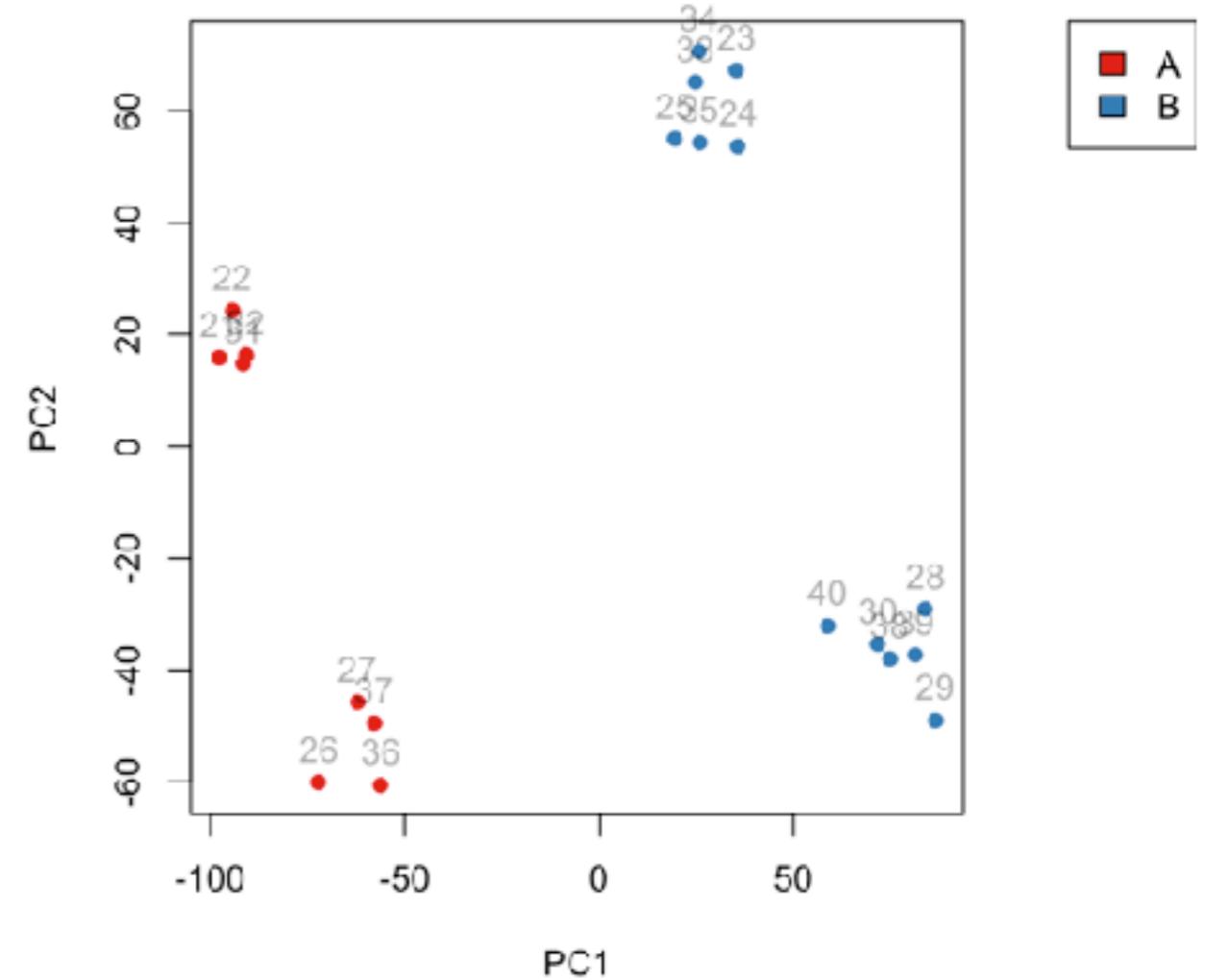
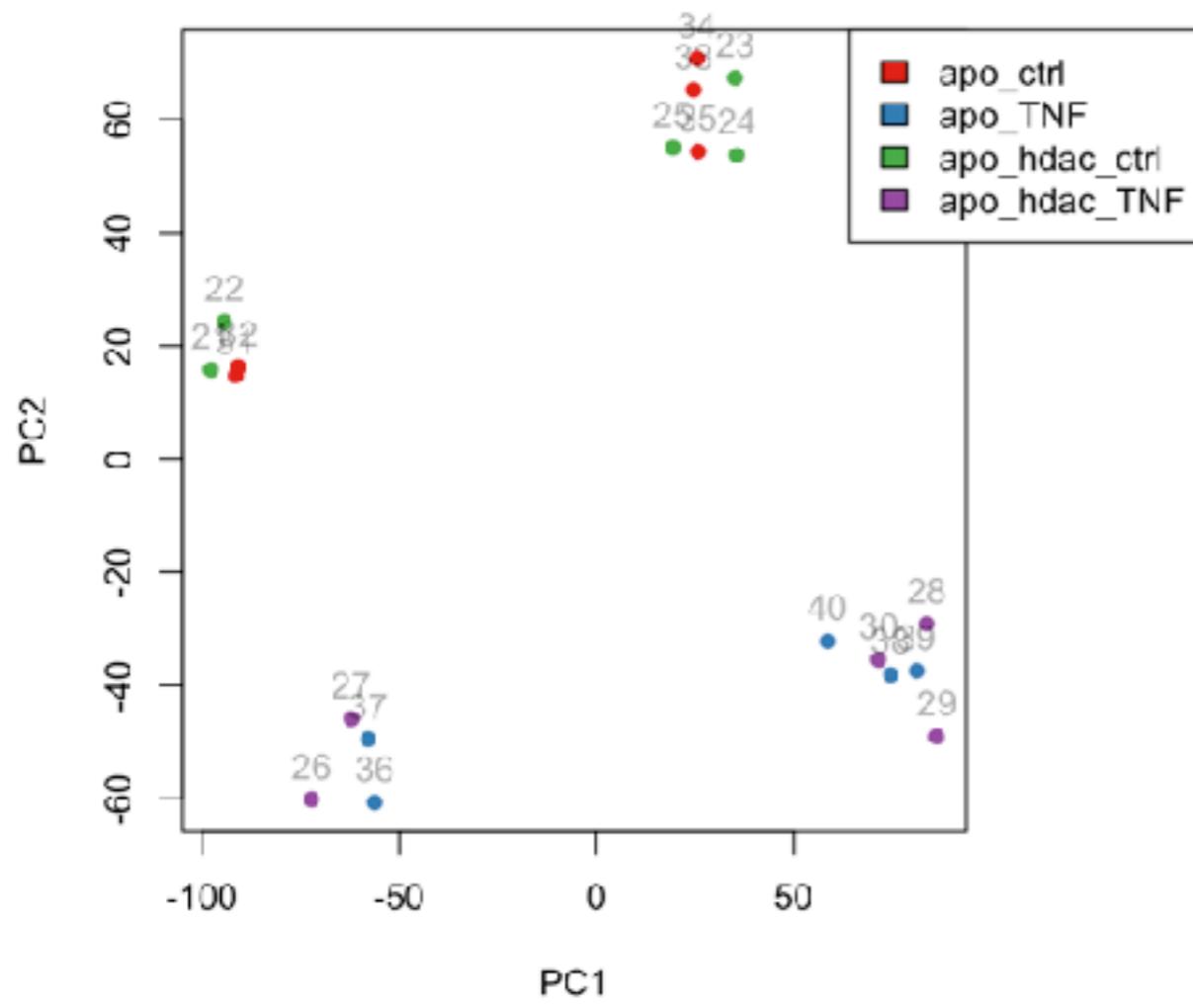
Difference in number of apples
The Effect

Genotype
independent variable



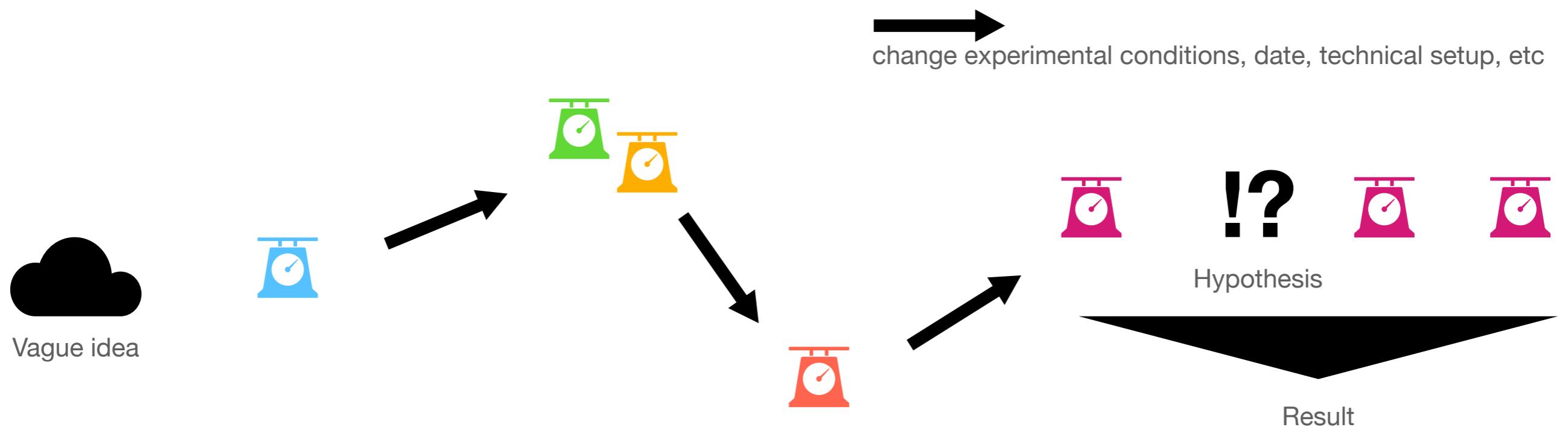
systematic influences on measurement are eliminated
relate the effect observed to the treatment, causal inference

Batch effects



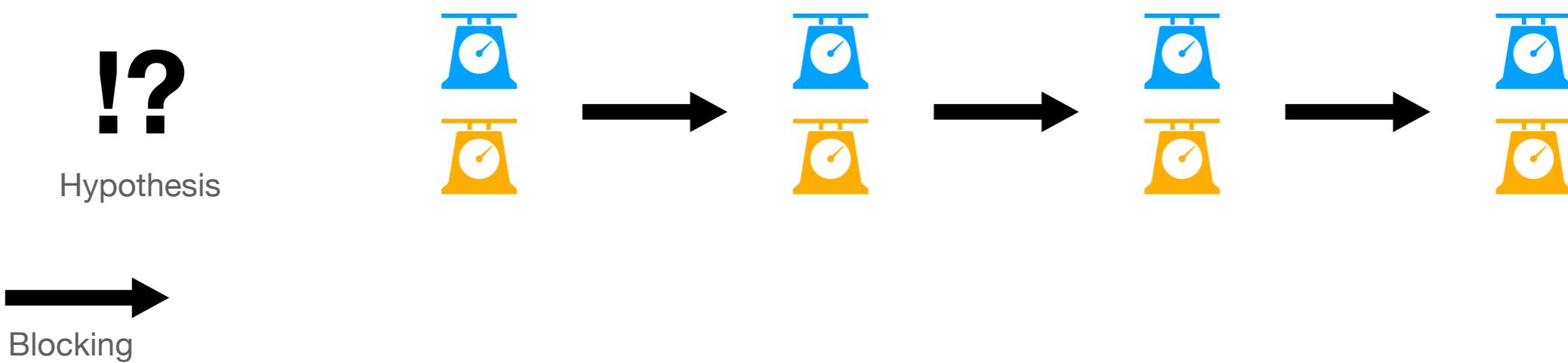
Planing
measurements

Exploratory Research



- Tends to be highly **irreproducible**
- Statistical evaluation highly problematic (**sampling bias, power**)
- Data organisation and documentation potentially **compromised** (lack of structure, retrospective collection & assembly)

Confirmatory Research

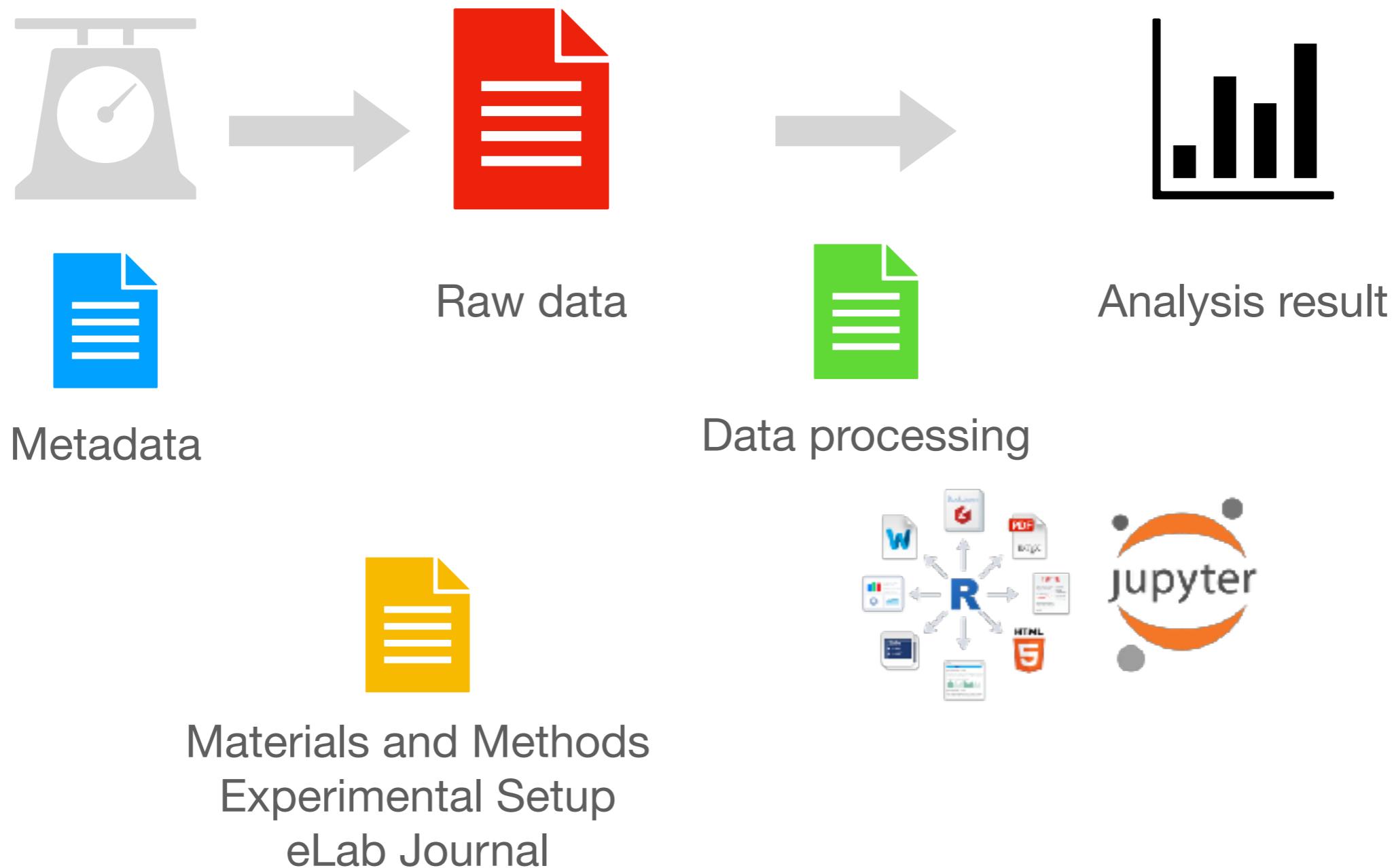


- Statistical evaluation possible
- Standardised measurements/data acquisition
- Estimation of robustness (reproducibility) possible
- Requires Experimental Design
- Prior definition of Raw Data and Metadata. Open Science practices facilitated.

Open Science

**State-of-the-art
Documentation and Data Management**

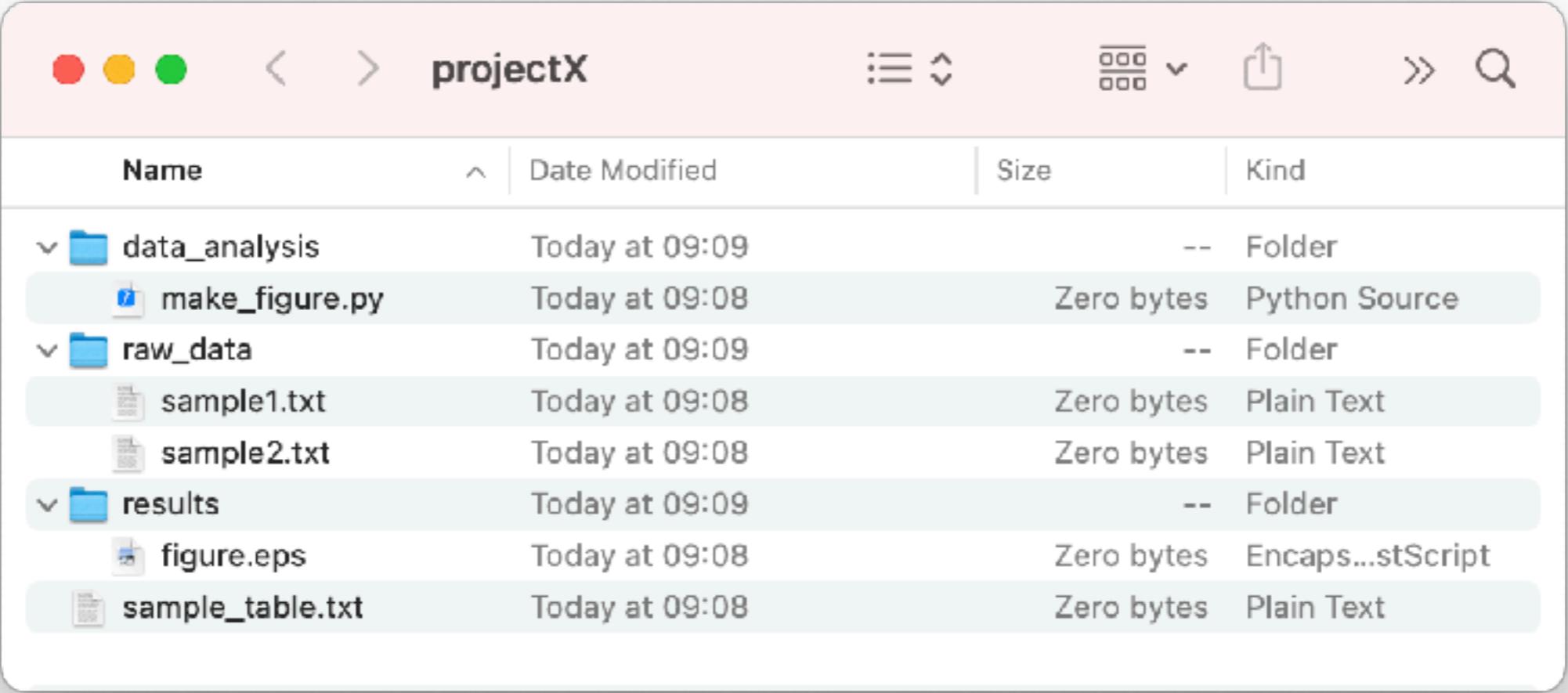
digital documentation



Open science practices (in short)

- All files centrally stored in universal format accessible/findable for group/peers/world (**FAIR**)
- Safely **backed up**
- Avoid redundancy
- Never edit the **raw** data files
- All relevant information provided in **meta data** table.
- Data analyses and environment fully documented to ensure **analytical reproducibility**

Data organisation

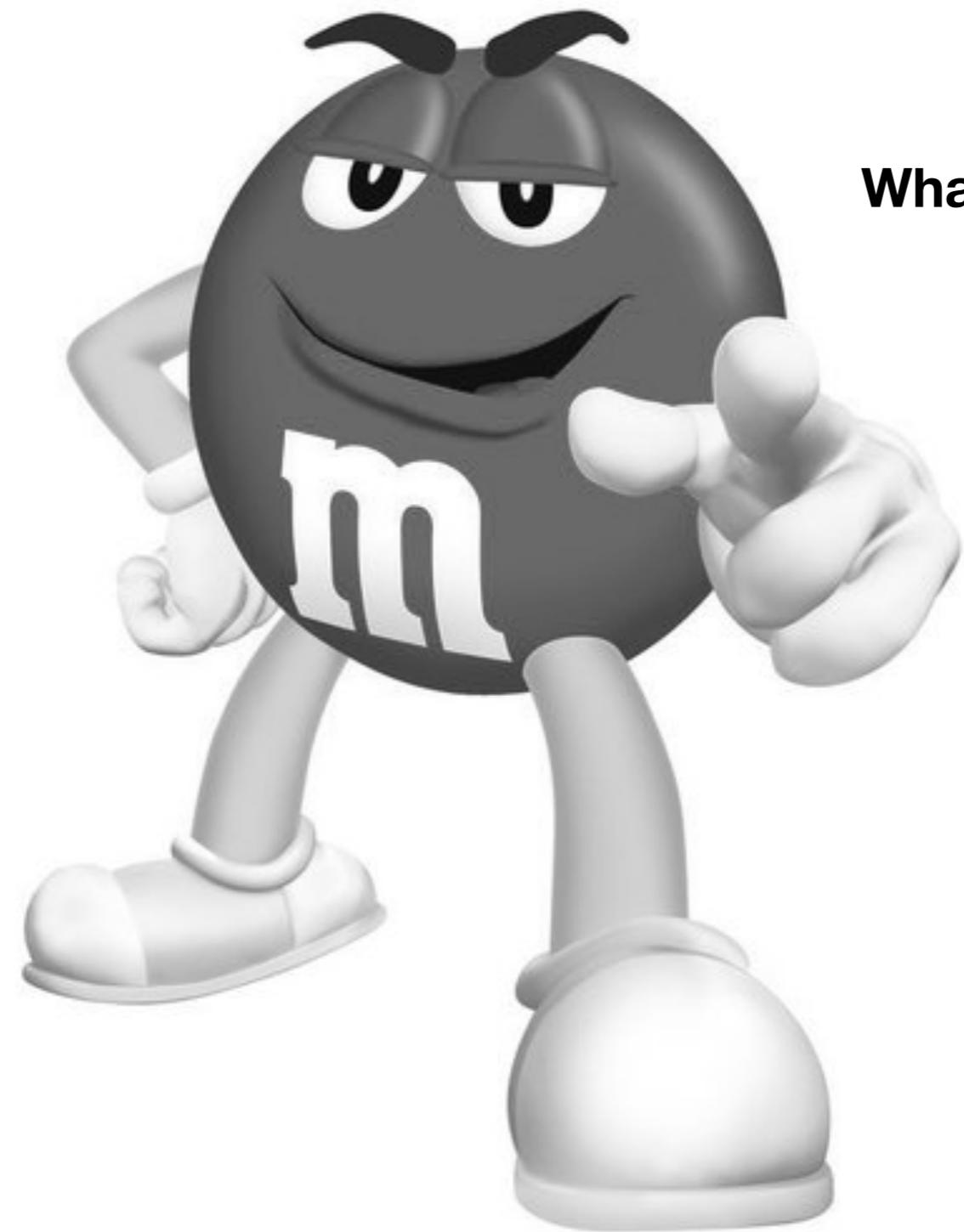


A screenshot of a file explorer window titled "projectX". The window has a light gray header bar with standard OS X-style controls (red, yellow, green buttons, back/forward arrows, title bar, and various icons for file operations). Below the header is a table-based file list.

Name	Date Modified	Size	Kind
▼ data_analysis	Today at 09:09	--	Folder
make_figure.py	Today at 09:08	Zero bytes	Python Source
▼ raw_data	Today at 09:09	--	Folder
sample1.txt	Today at 09:08	Zero bytes	Plain Text
sample2.txt	Today at 09:08	Zero bytes	Plain Text
▼ results	Today at 09:09	--	Folder
figure.eps	Today at 09:08	Zero bytes	Encapsulated PostScript
sample_table.txt	Today at 09:08	Zero bytes	Plain Text

Sampling & Representativity

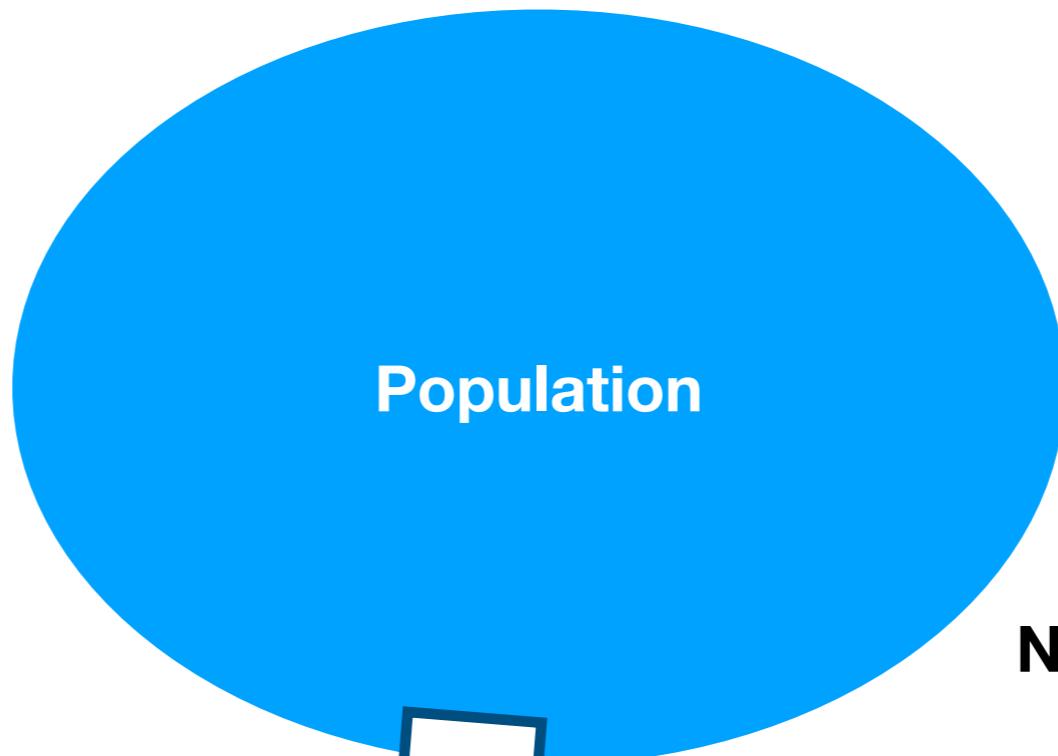
N is SMALL



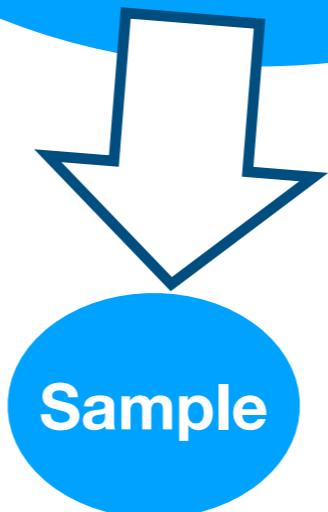
What's my colours?

Any expectation?

Experiment



N = usually infinite



N = 3



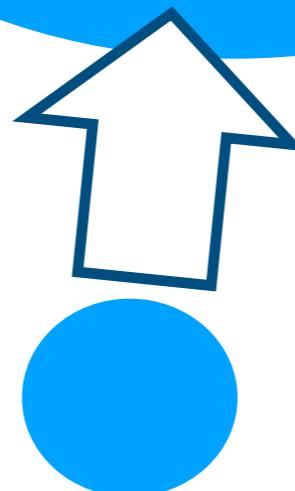
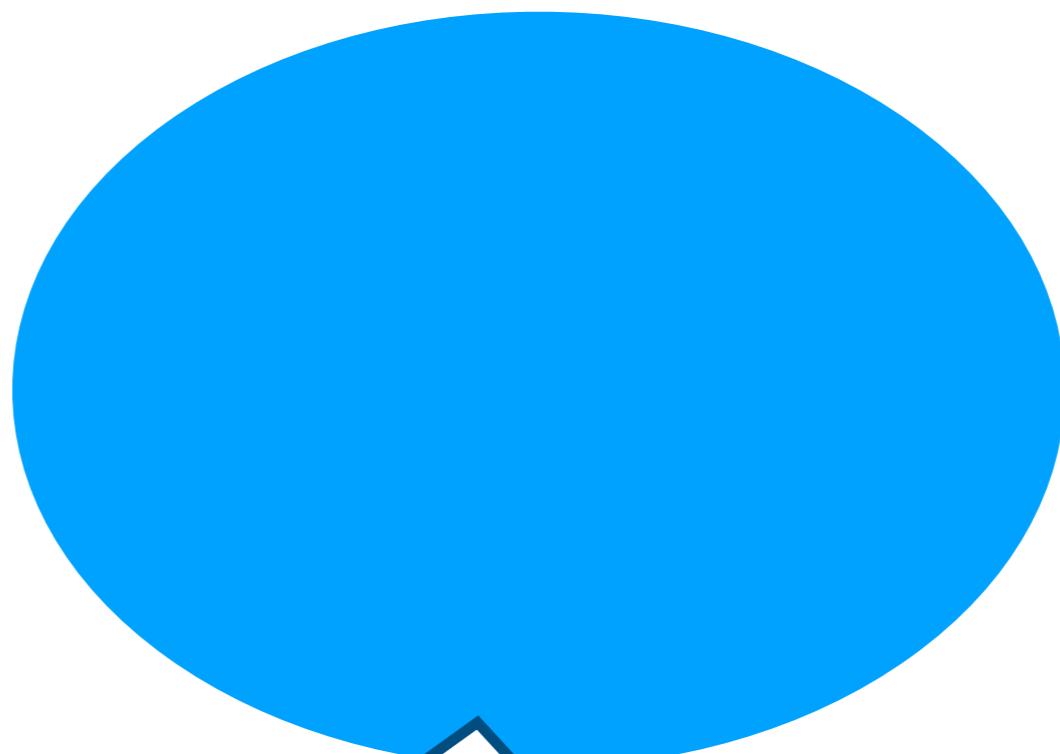






Inference

Estimate the colours of the population from the colours of the sample



sweets



All M&Ms
are yellow!

Reporting Bias

“**My** M&Ms are yellow”

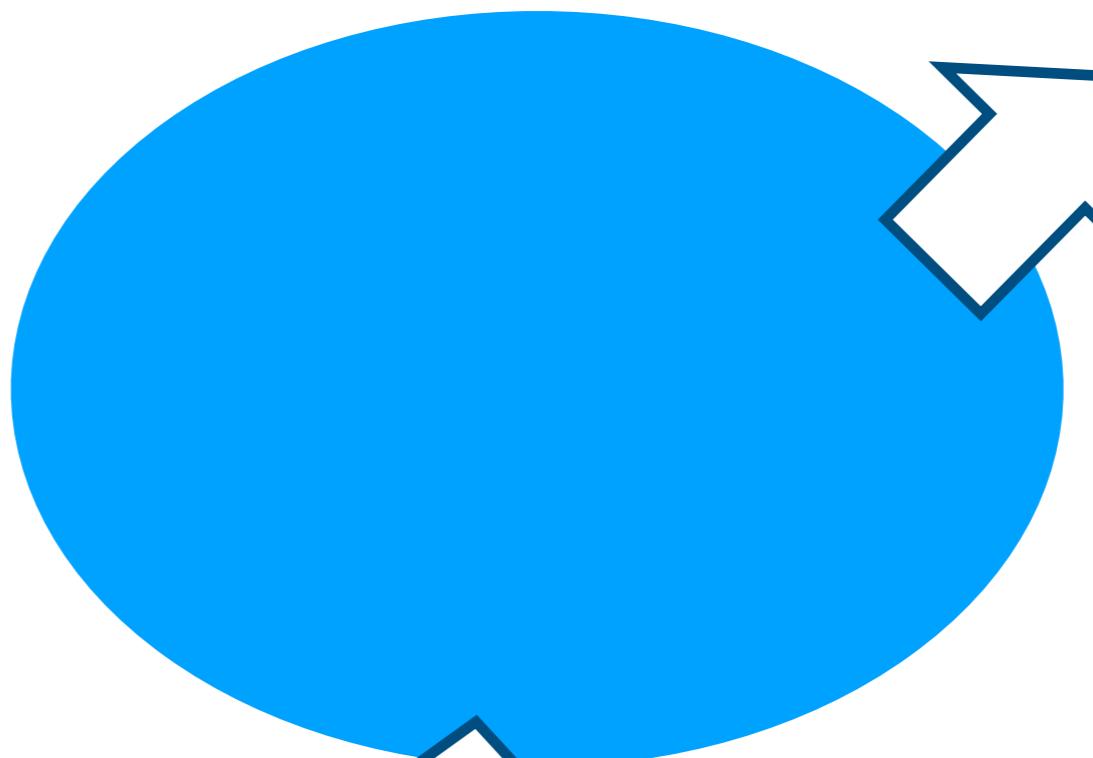
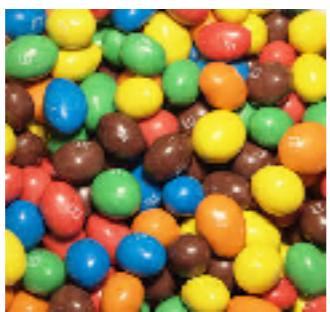
does not equal

“**All** M&Ms are yellow”

There is no stronger force in science but the love of a scientist for her/his own results

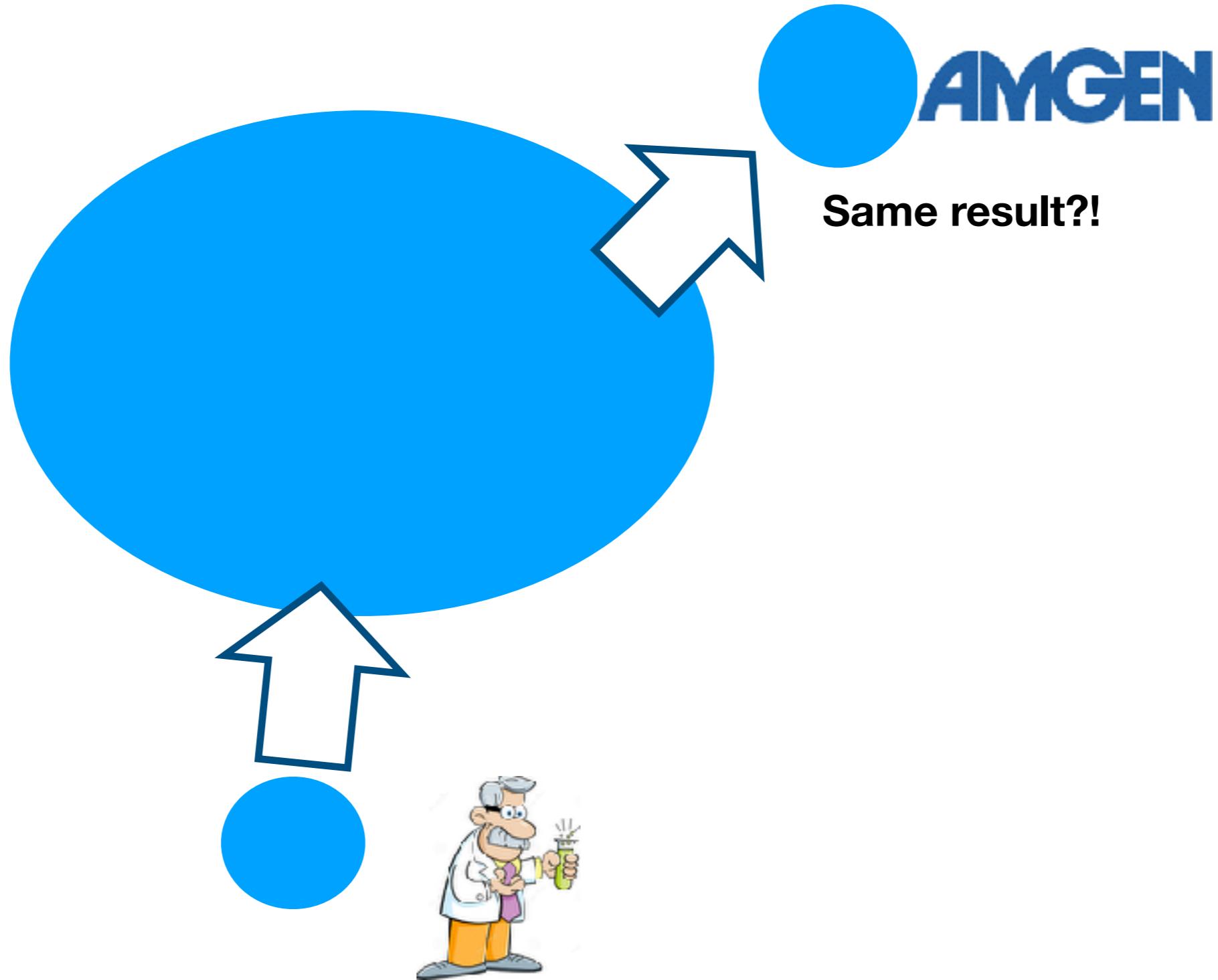
Reproducibility

Direct, Conceptual, Analytical ...

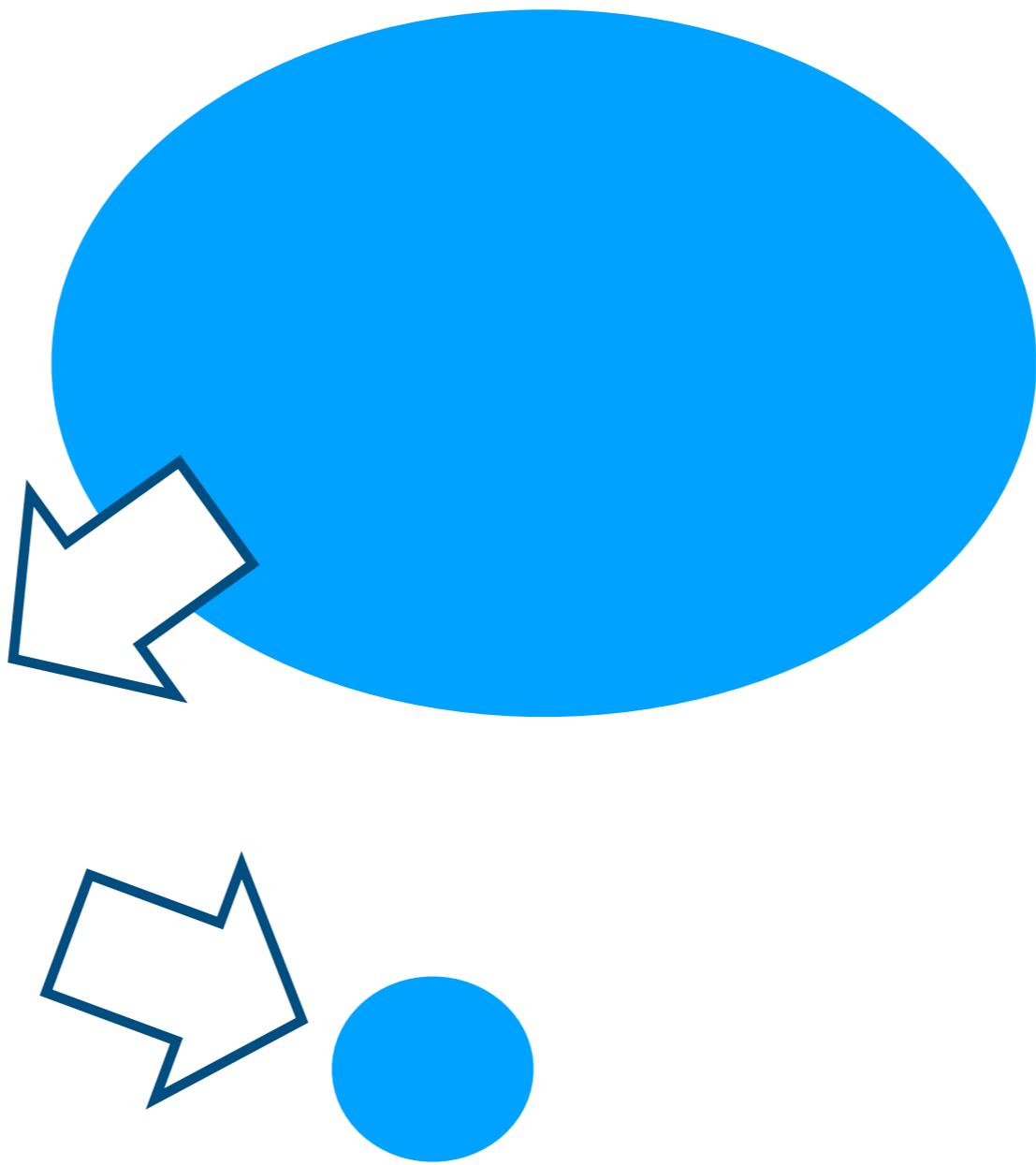


Same result?!

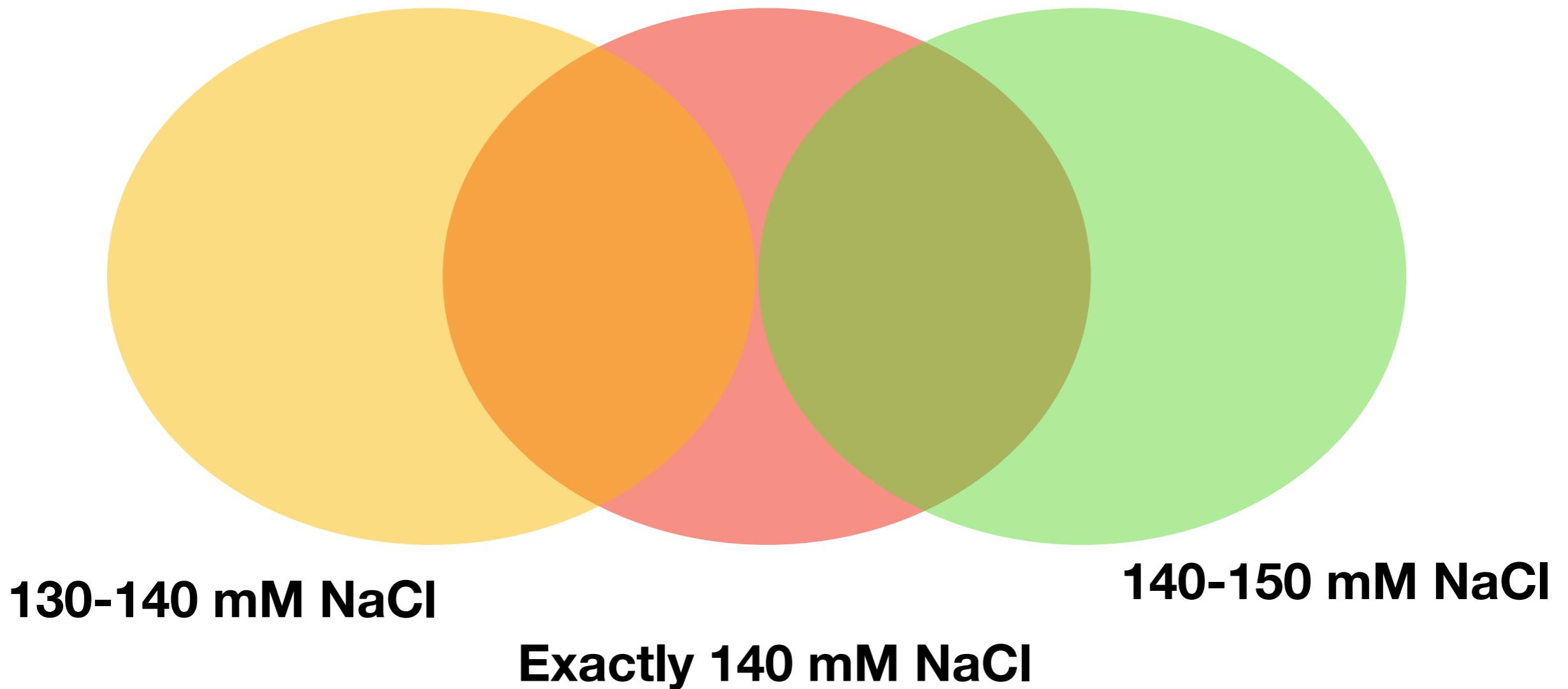
Roadblock to Translation



Sampling space = population?



What defines the population?



Reproduction of the original results using the same protocol/reagents/tools

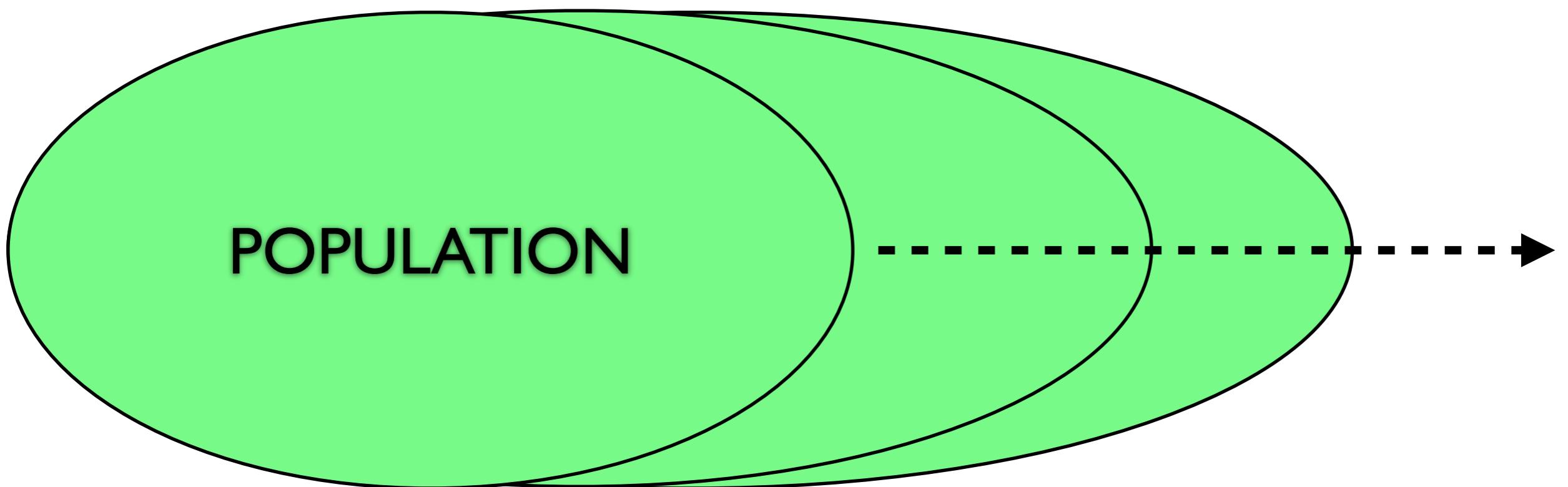
by the same person

by a different person in the lab

by a different person outside the lab

Reproduction using different reagents/tools but the same protocol by a different person outside the lab

Reproduction just based on text description



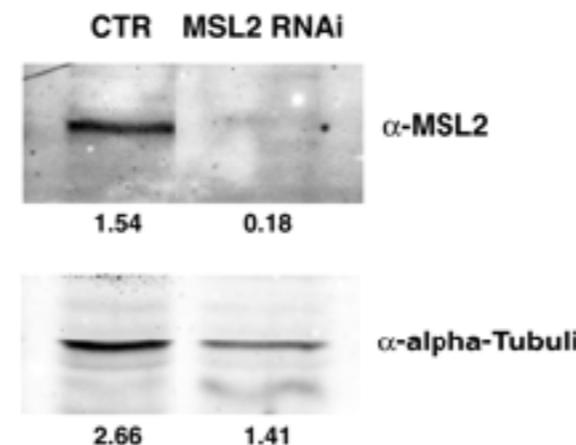
Parameters to consider

- Environmental conditions (animal house, food, health state, mood, season, time of day ...)
- Cell type
- Reagents (serum!, antibodies!!!, ...)
- Technical (detection limits, ...)
- Genomic background, secondary (off-target) mutations
- ...

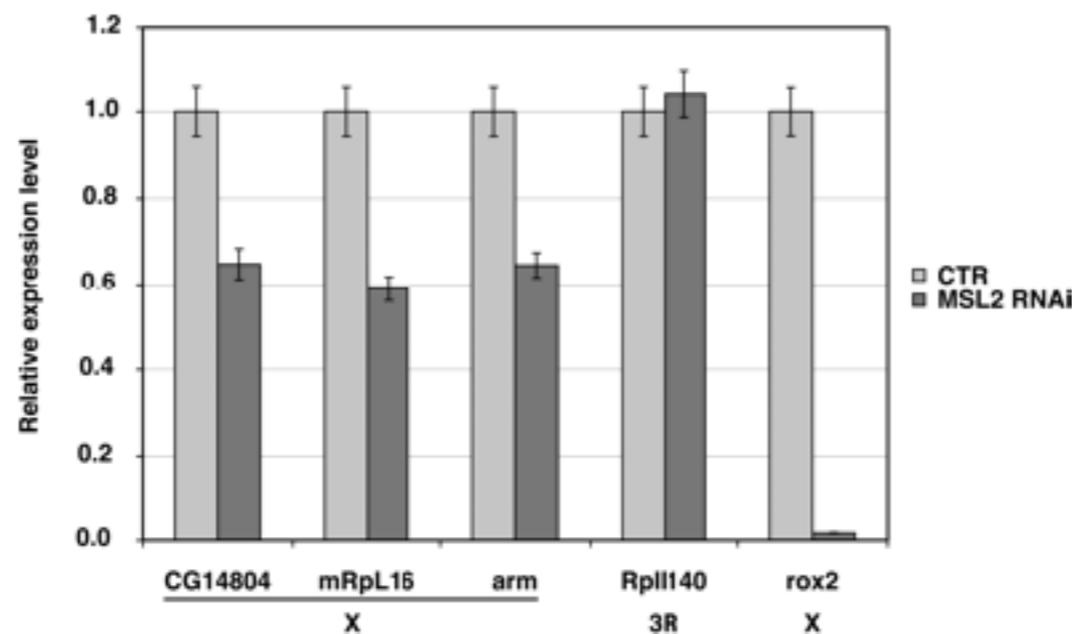
Certainty

Is This True?

A

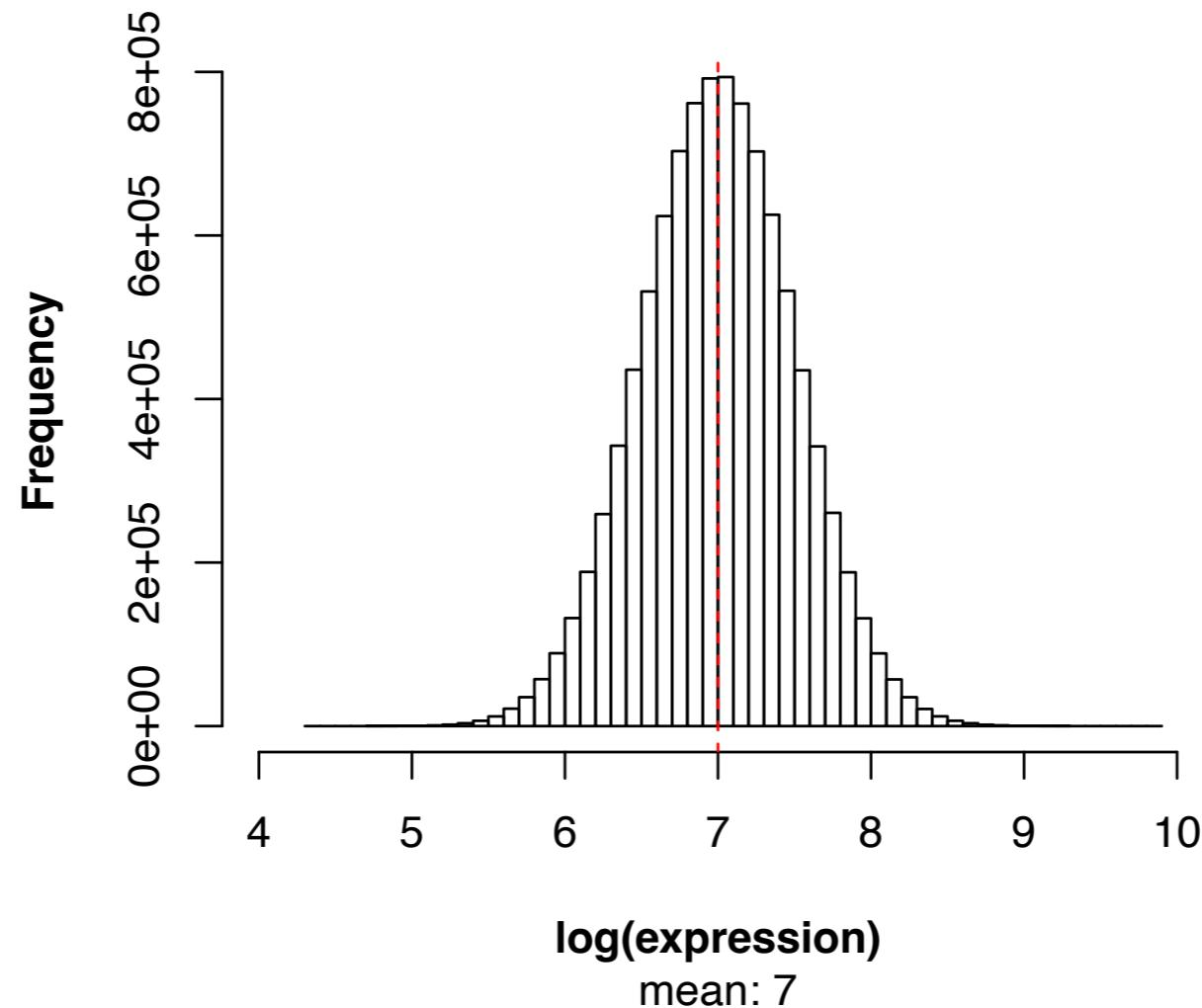


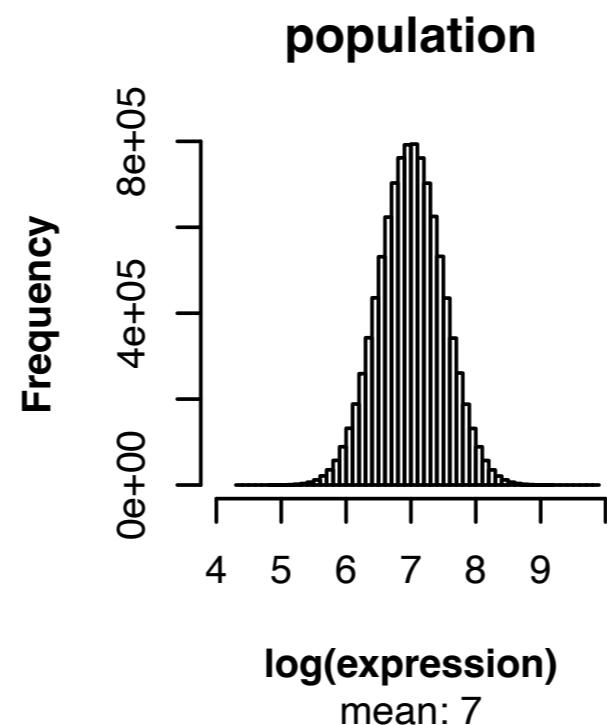
B



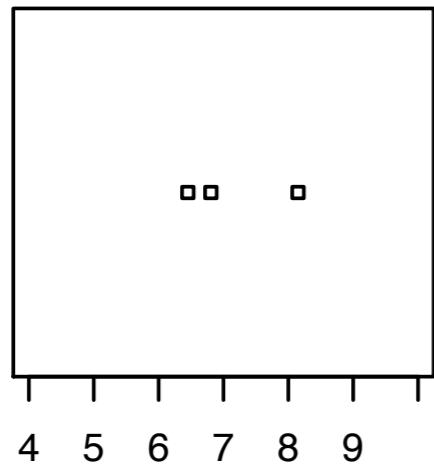
Objective Inference

gene X expression (population)



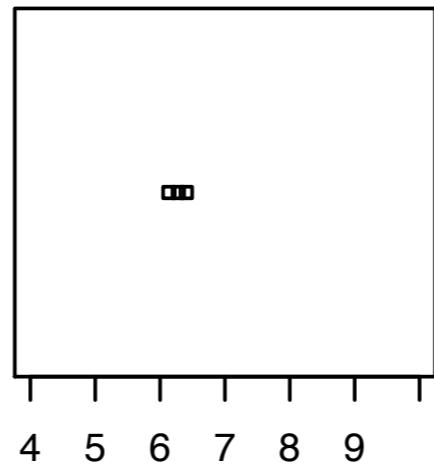


expt 1



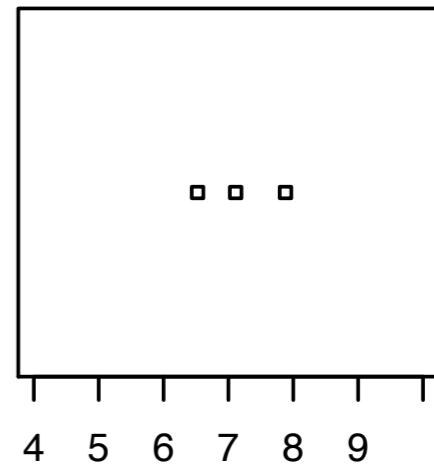
log(expression)
mean: 7.135

expt 2

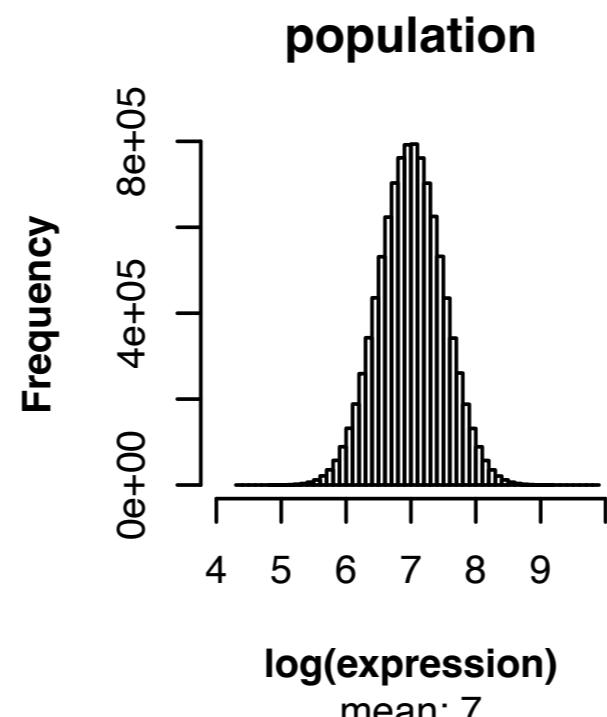


log(expression)
mean: 6.274

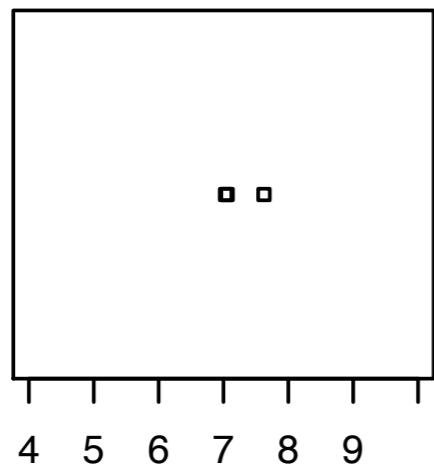
expt 3



log(expression)
mean: 7.172

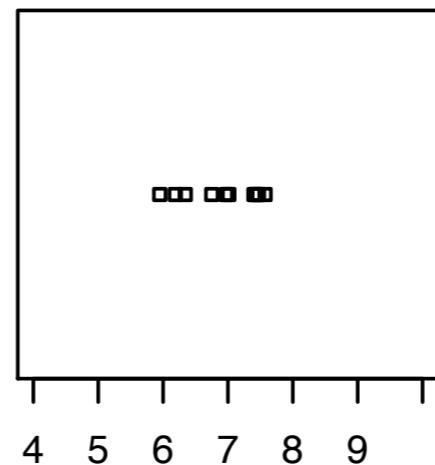


n=3



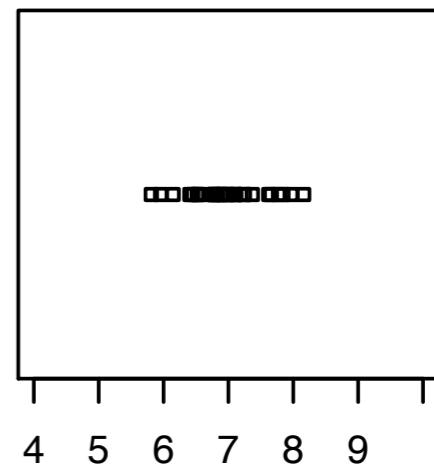
log(expression)
mean: 7.238

n=10



log(expression)
mean: 6.908

n=30



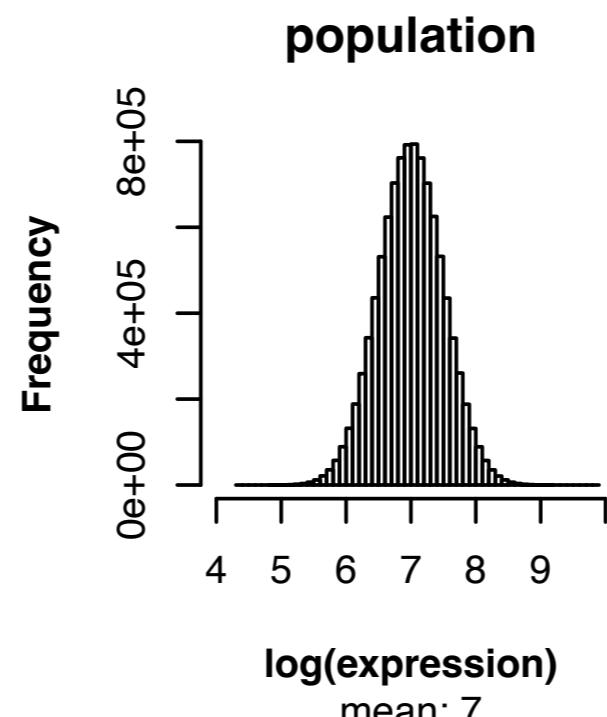
log(expression)
mean: 6.972

Standard Error (of the mean)

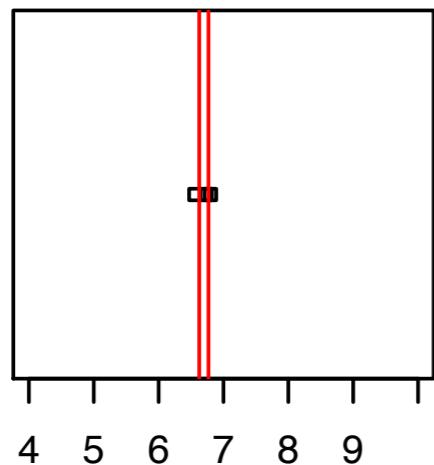
- The standard error of the mean (SEM) is the standard deviation of the sample mean estimate of a population mean.

$$\text{SEM} = \text{standard deviation}/\sqrt{n}$$

- a small SEM indicates that the sample mean is likely to be quite close to the true population mean
- a large SEM indicates that the sample mean is likely to be far from the true population mean

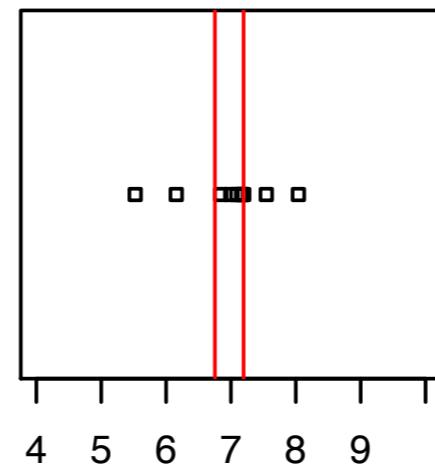


n=3



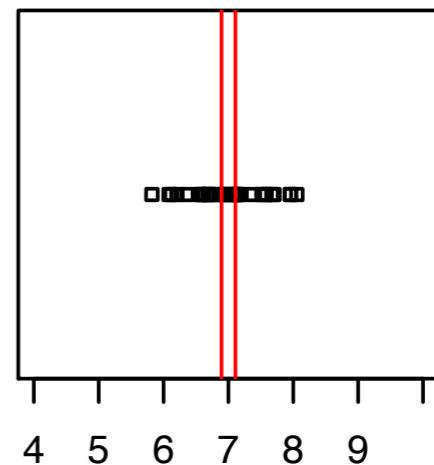
log(expression)
mean: 6.698

n=10



log(expression)
mean: 6.974

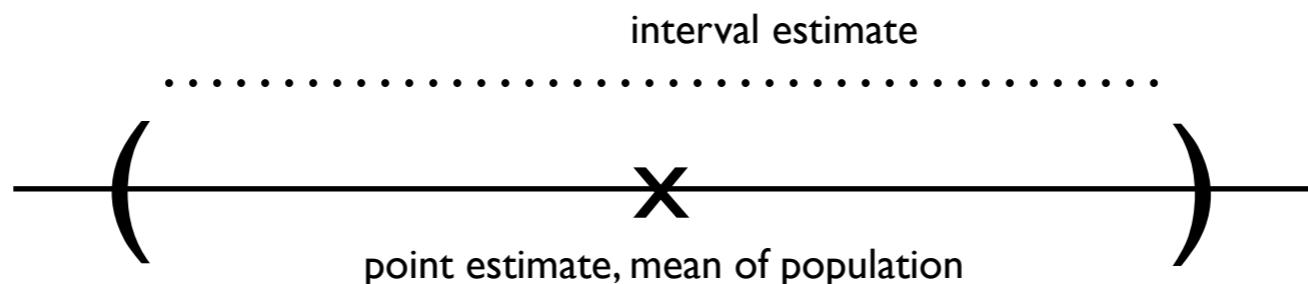
n=30



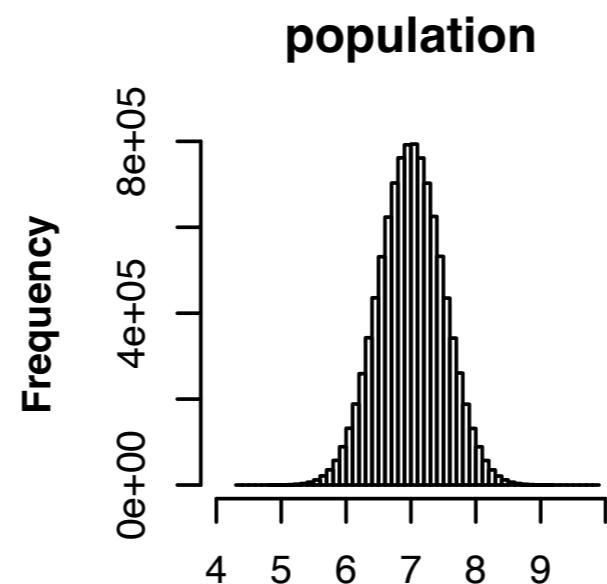
log(expression)
mean: 7.001

Confidence Intervals

- 95%-confidence interval: An estimated interval which contains the „true value“ of a quantity with a probability of 95%.

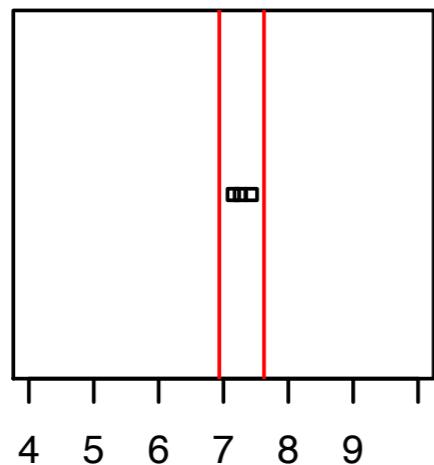


- $(1 - \alpha)$ -confidence interval: An estimated interval which contains the „true value“ of a quantity with a probability of $(1 - \alpha)$.
 $1 - \alpha$ = confidence level, α = error probability



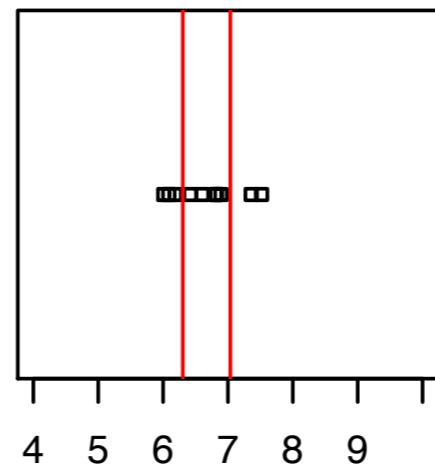
log(expression)
mean: 7

n=3



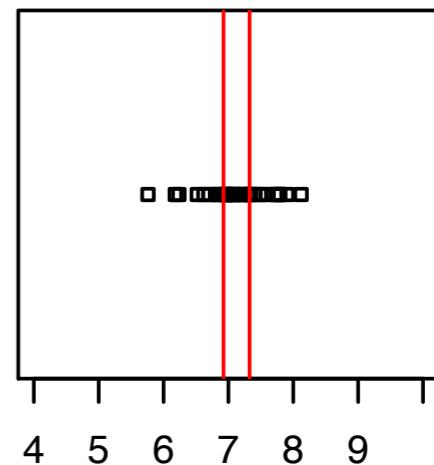
log(expression)
mean: 7.28

n=10



log(expression)
mean: 6.671

n=30



log(expression)
mean: 7.125

Practical example

Someone asks: “how many dead cells are in your culture?”

You use a hemocytometer to determine the viability of cells stained with trypan blue.
You count 94 unstained cells and 6 stained.

How can the data be reported?

95% CI=0.02-0.13

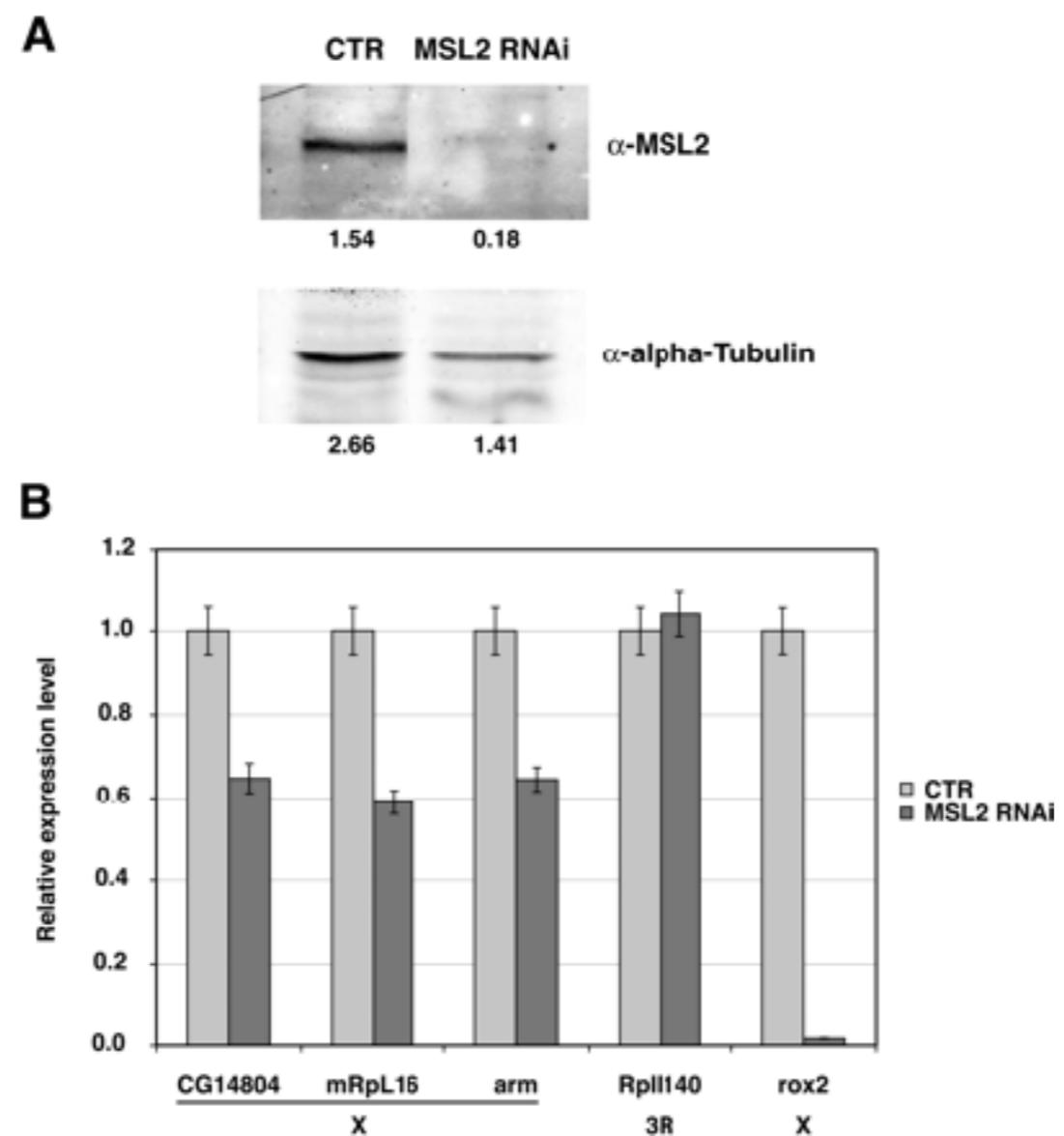
Prerequisites for inference

- the sample has to be representative
- how is representativity achieved?
 - large sample number (N)
 - independent sampling/
random recruitment of samples
(BIOLOGICAL replicates)

Unbiased Sampling

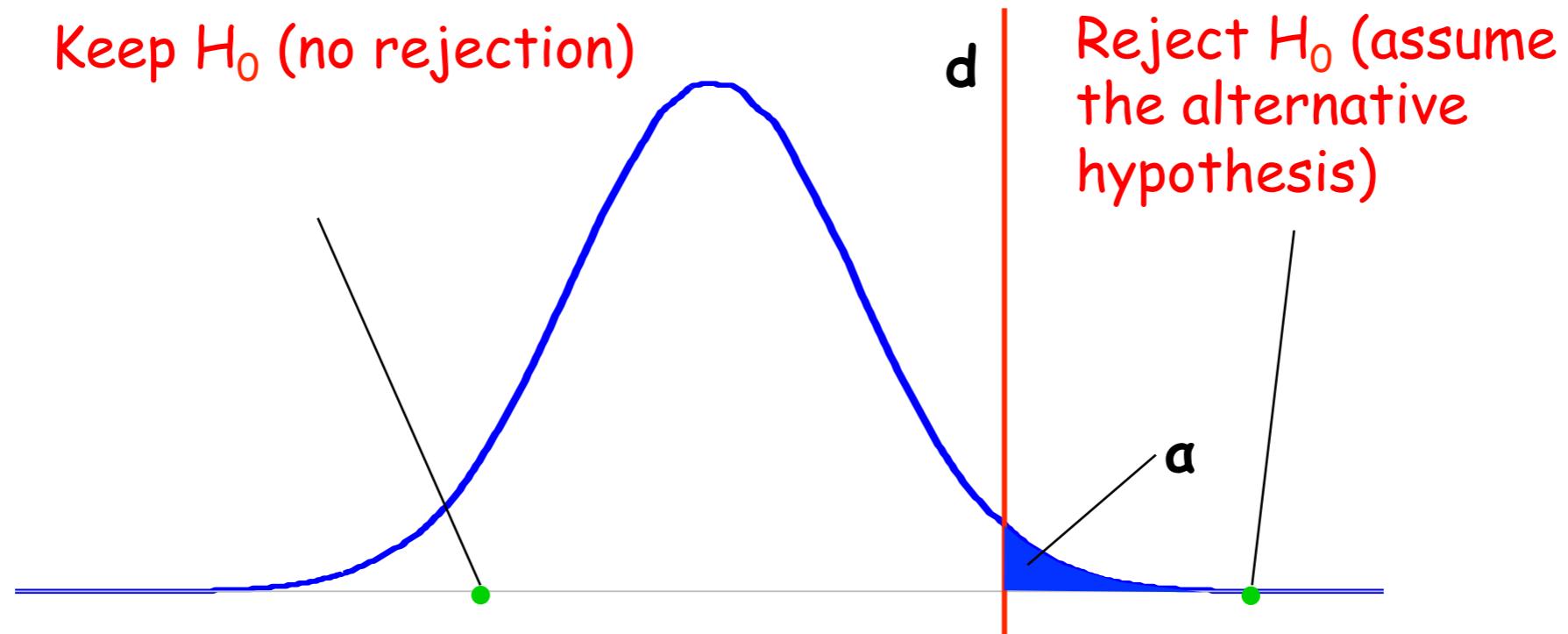
- Discard results only if:
 - objective **exclusion criteria** have been defined beforehand. Use positive/negative **controls**.
 - made transparent.
- **Randomisation** measures and blinding help fight the bias
- In “gonzo” designs, selection of sample objects defines the population.

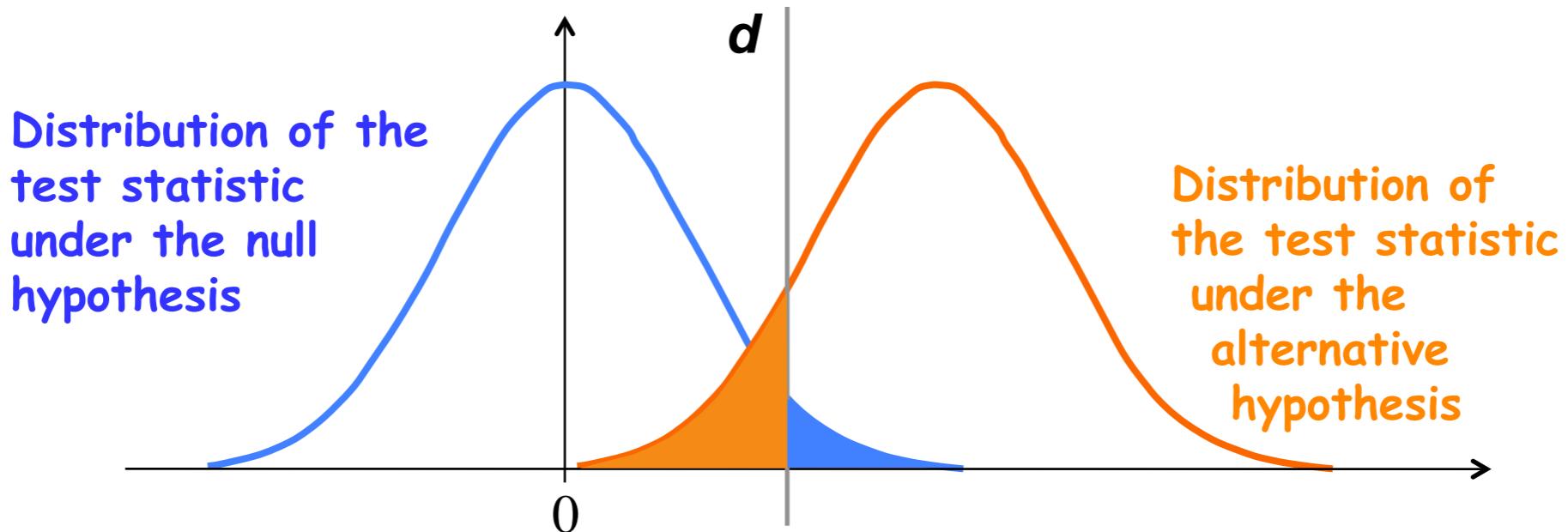
Truth is one, Paths are many @lakens



Truth is one, Paths are many @lakens

- Frequentist (Neyman Pearson):
 - evaluation of ongoing data collection. Says nothing concrete about the current test
- Likelihood
 - Compare the likelihood of different hypotheses given the data
- Bayesian Statistics
 - Update prior beliefs in light of new data



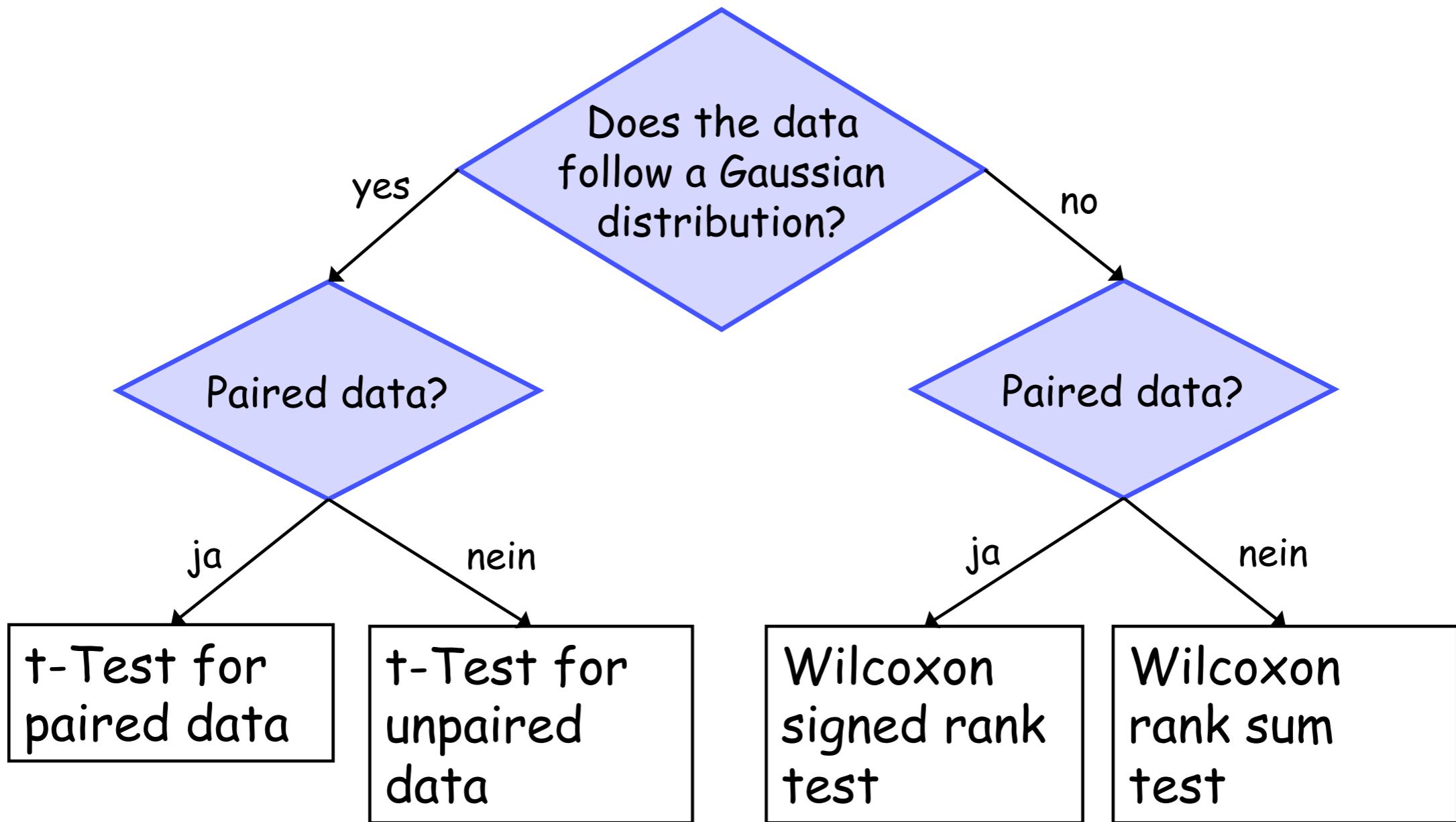


	Accept null hypothesis	Reject null hypothesis
null hypothesis is TRUE	correct decision	Type I Error “False Positive”
alternative hypothesis is TRUE	Type II Error “False Negative”	correct decision

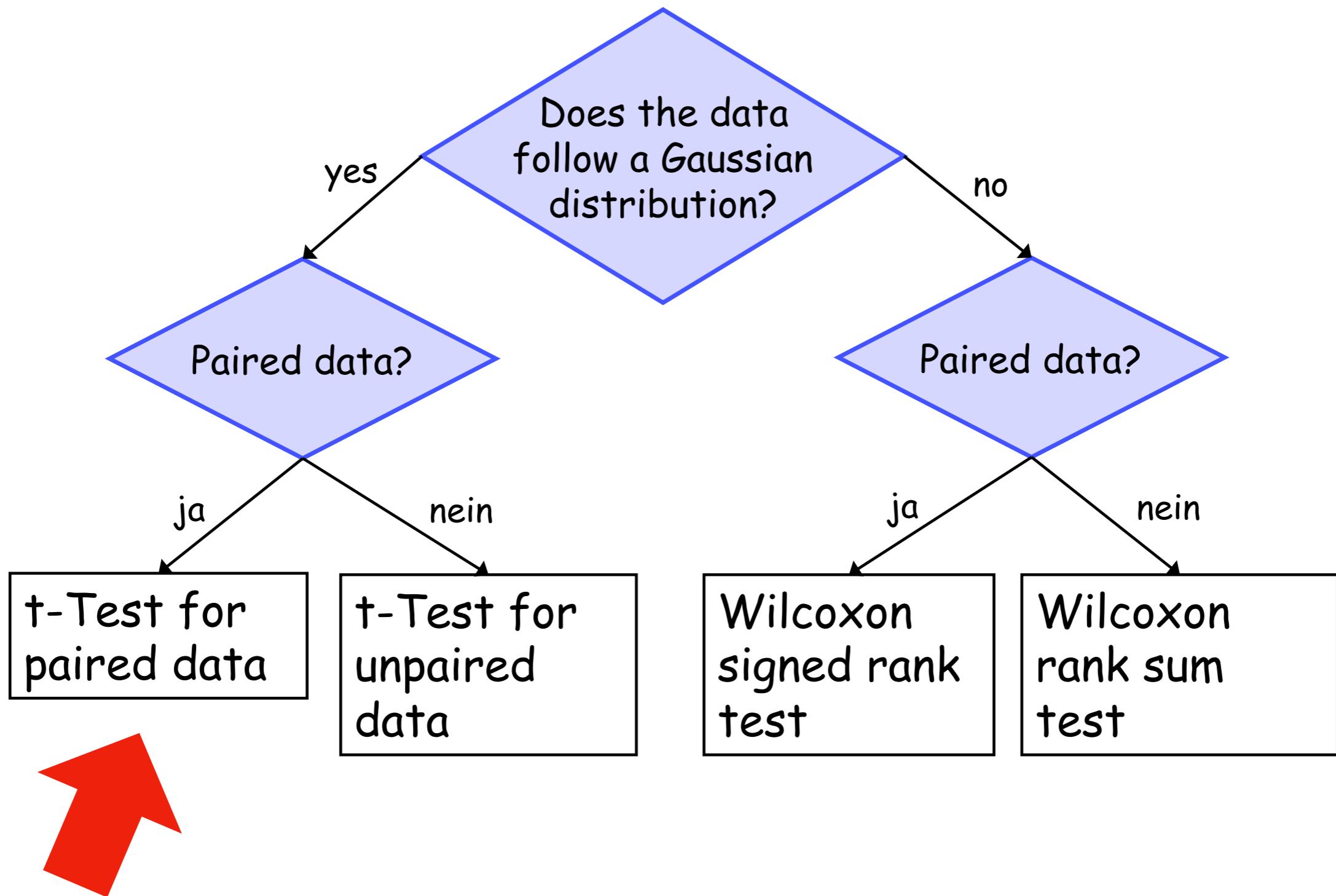
Statistical Testing

- Statistical tests are **inferential** procedures.
- Statistical tests only address the robustness of effects if applied over independent **biological replicates**.
(watch out X-omics!)
- No testing without education otherwise you are wrong.
Understand the **prerequisites** for statistical testing procedure.
- Understand the meaning of **statistical errors** and the concept of **multiple testing**.
- Understand the meaning of **statistical power**.

Question: Are group 1 and group 2 identical with respect to the distribution of the endpoint?

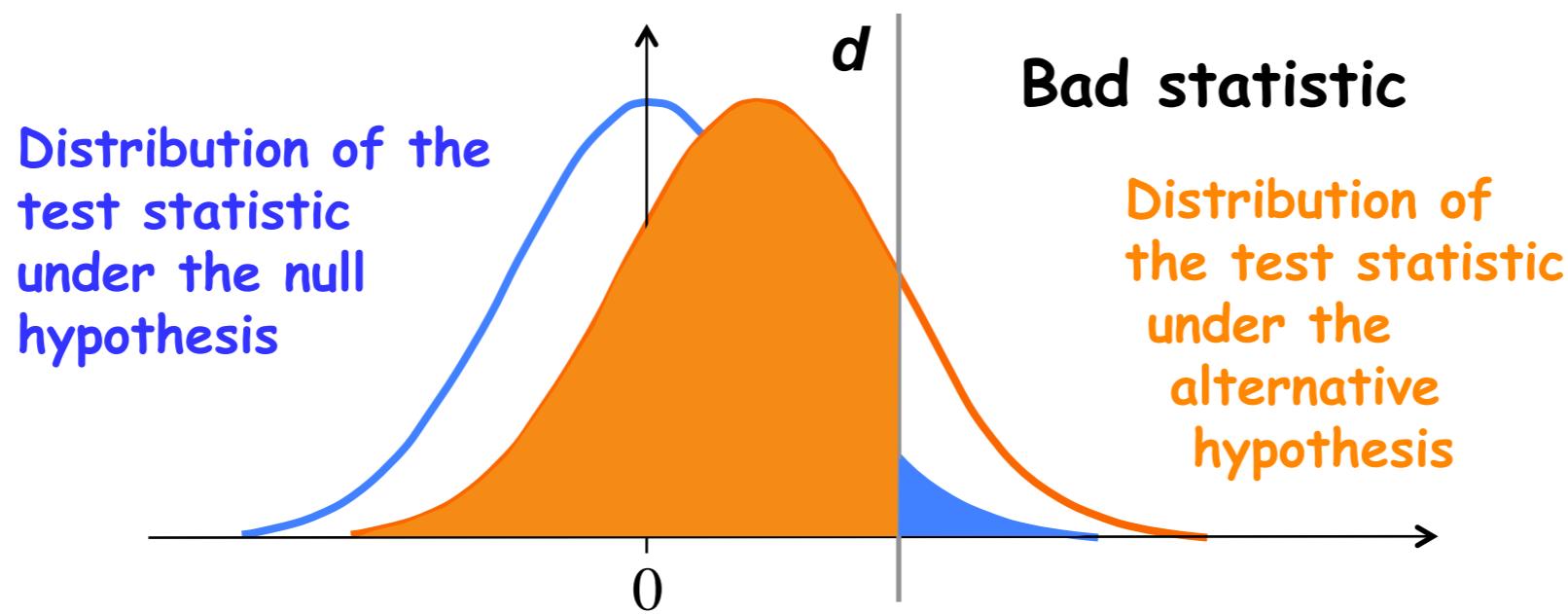
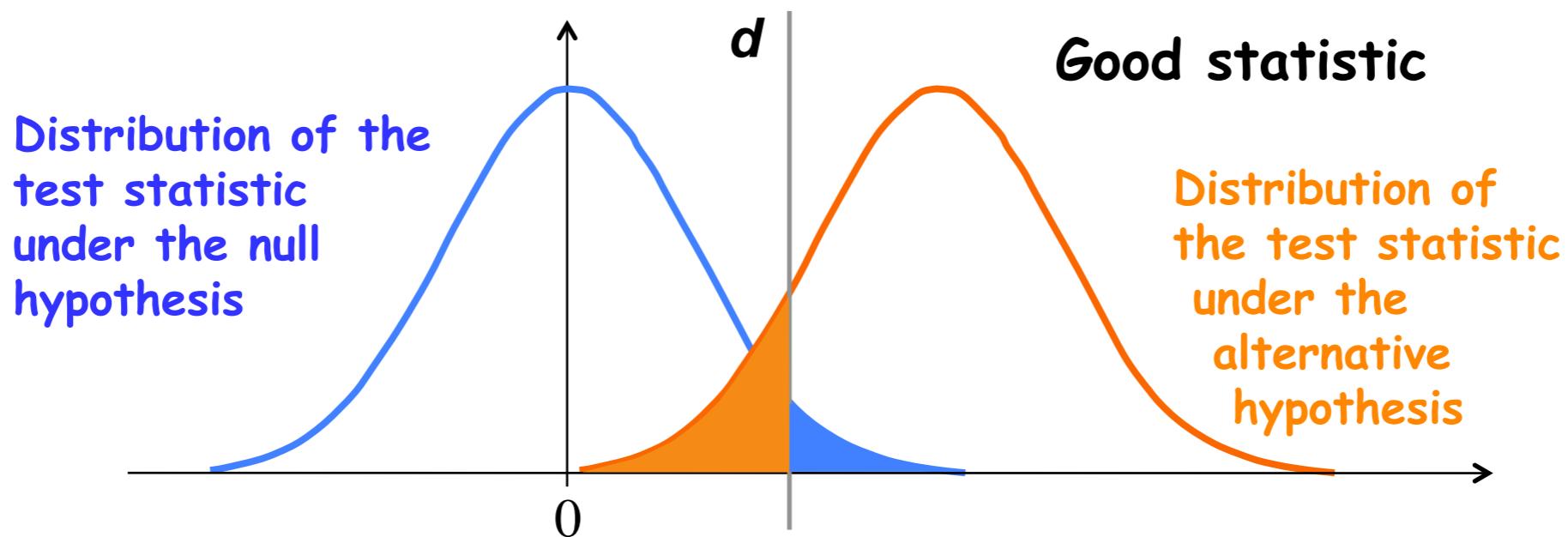


Question: Are group 1 and group 2 identical with respect to the distribution of the endpoint?



Statistical Power

As has been shown previously, the probability that a research finding is indeed true depends on the prior probability of it being true (before doing the study), the statistical power of the study, and the level of statistical significance [10,11]. Consider a 2×2

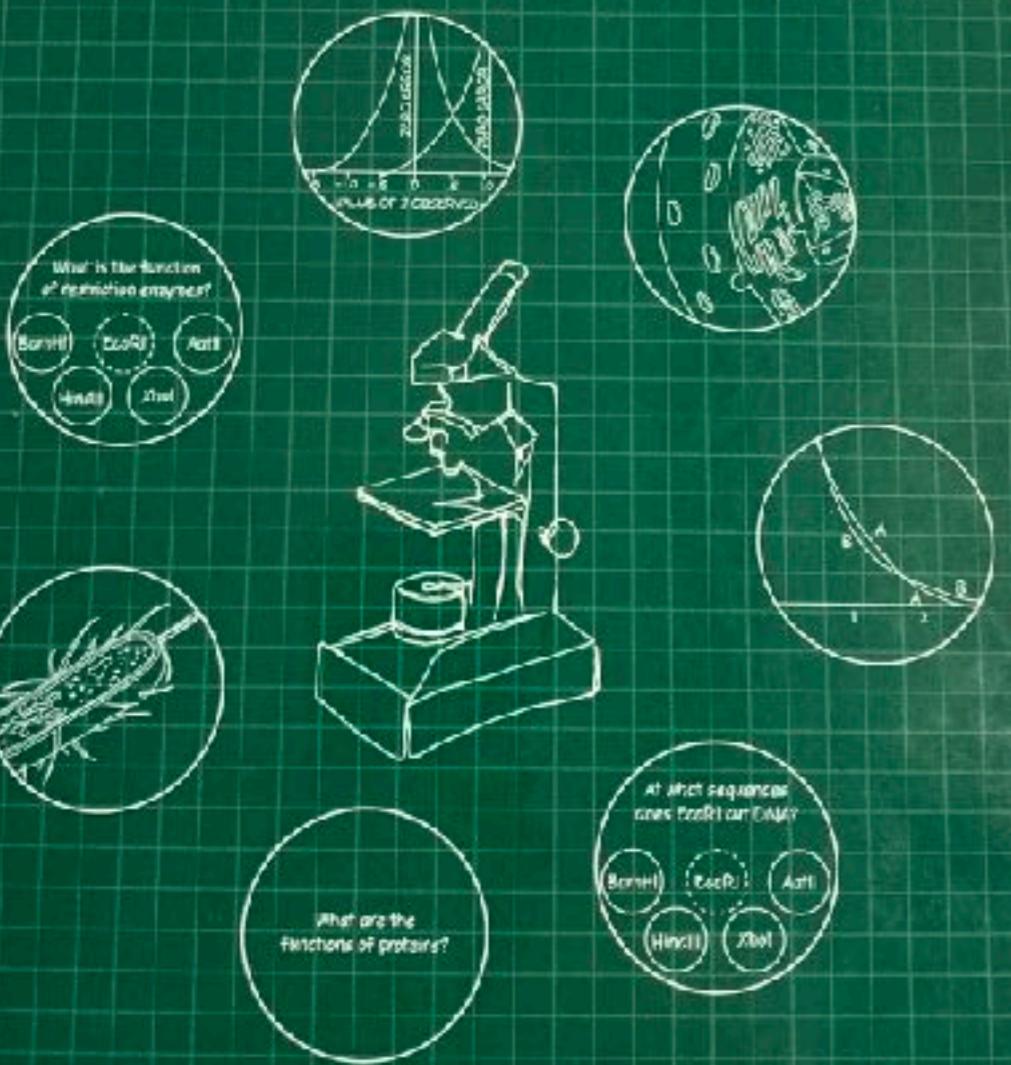


- Depending on the magnitude and variability of the effect (effect size) being investigated an experiment requires an **adequate sample size** (N) for reliable evaluation by statistical testing procedures.
- Very often N is too small, study **underpowered**, true effects not detected.
- If N is very large, small, potentially irrelevant effects become significant. (i.e. **overpowered** study)

Experimental Design

EXPERIMENTAL DESIGN FOR BIOLOGISTS

SECOND EDITION



DAVID J. GLASS

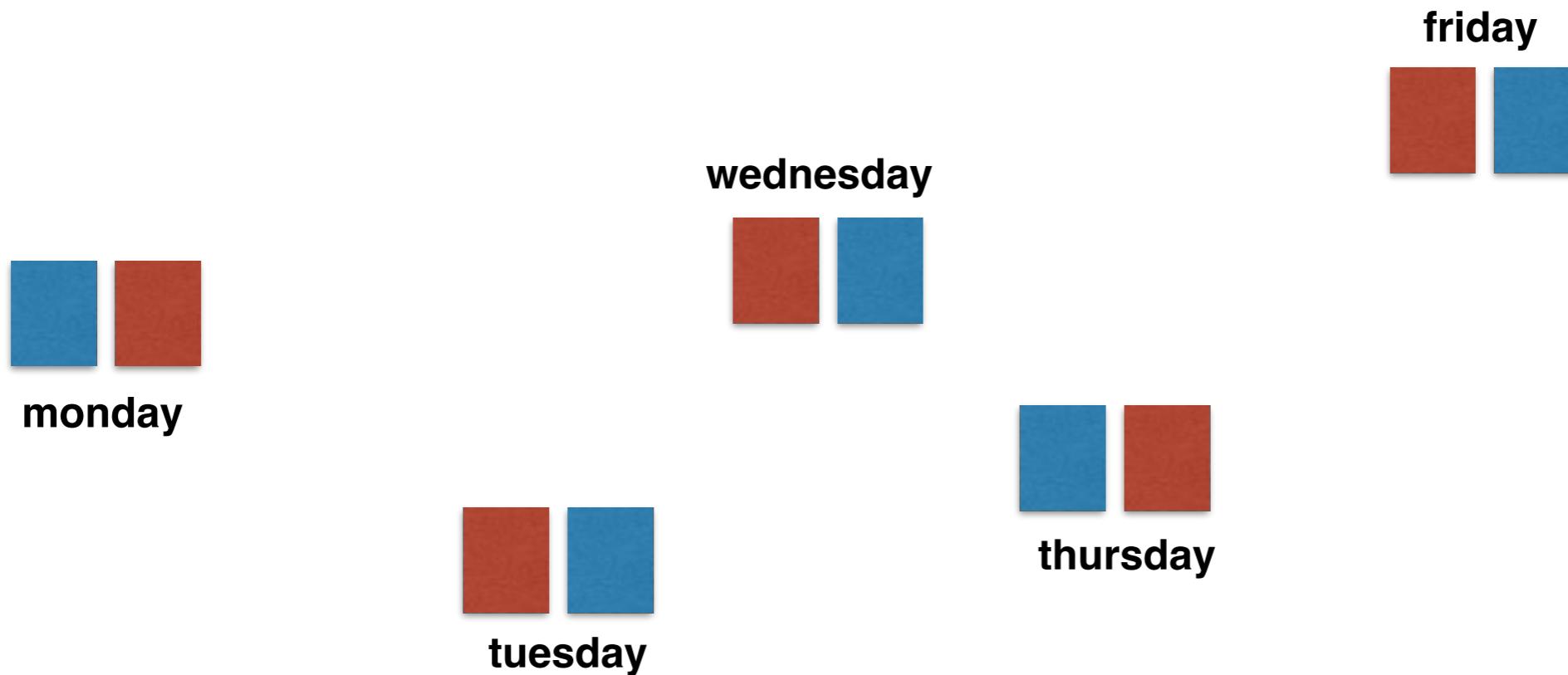
STANLEY E. LAZIC

Experimental Design for Laboratory Biologists

Maximising Information and
Improving Reproducibility

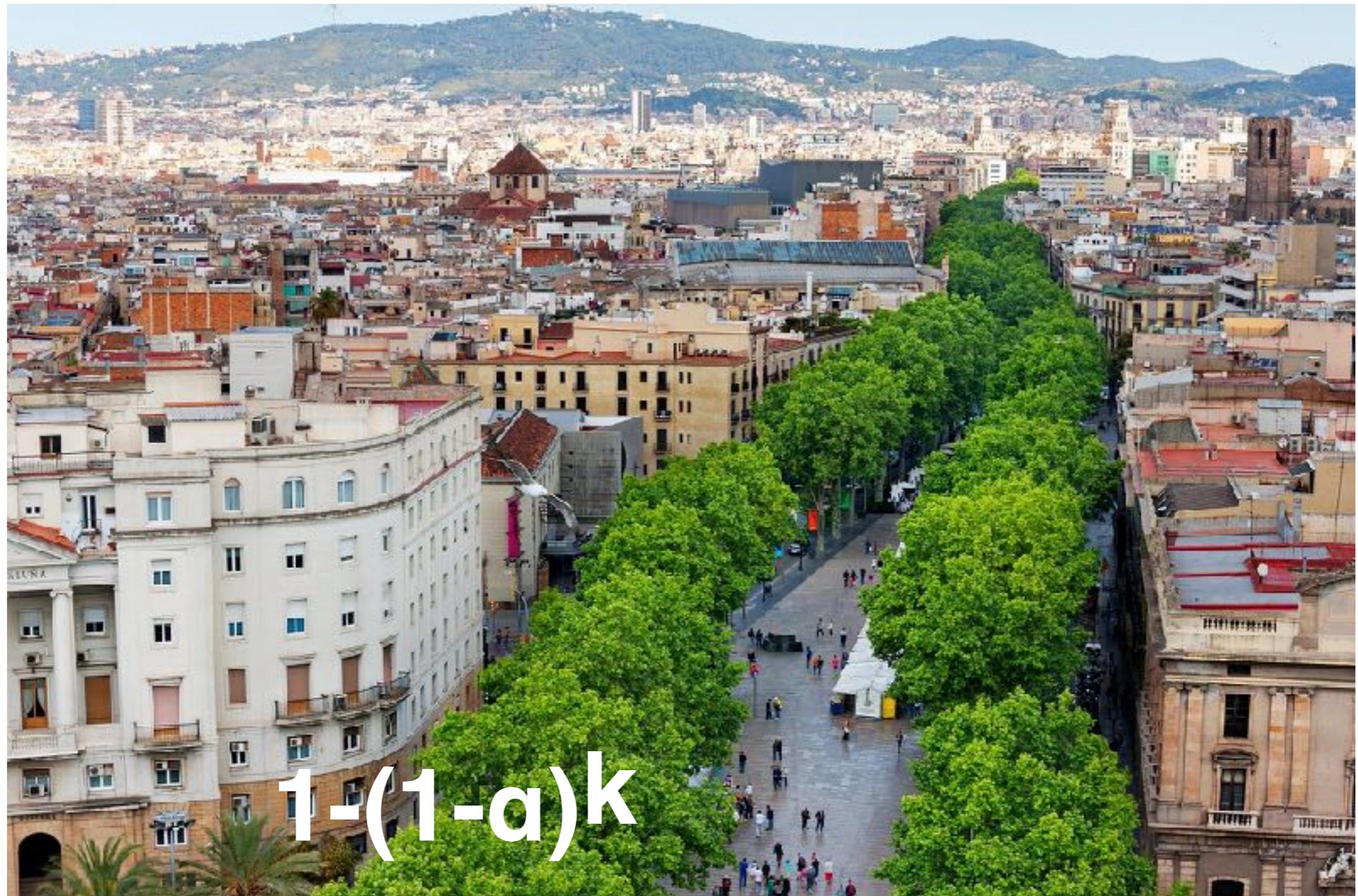


- removes systematic and reduces random errors
(blocking / replication)
- eliminates biases
(randomisation / blinding)
- Maximises sensitivity for the biological effect, (i.e. reveals interesting phenomena that are blurred in gonzo designs)
- can help to **increase statistical power** and/or reduce N



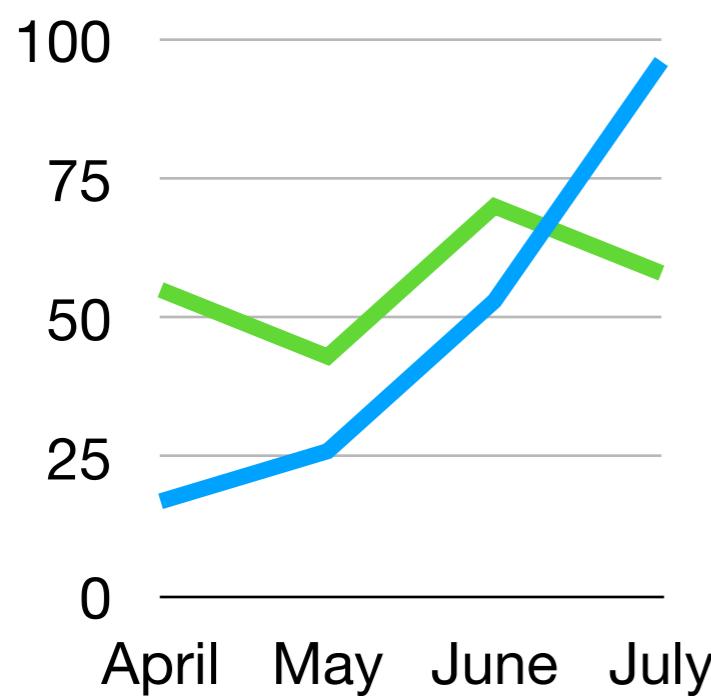
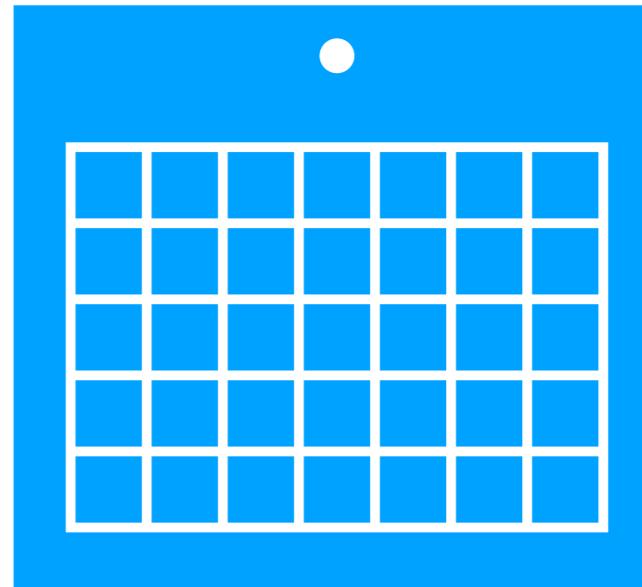
- randomised block design, only 2 factor levels (control, treatment)
- suited to control for day-to-day fluctuations which are very common. Also change reagents, batches of cells etc. between the blocks as well. Every block a new batch, every block new reagents.
- paired t-test (very powerful!)

Multiple Testing



$1-(1-a)^k$

Many endpoints, more than one independent variable



- In practice, number of endpoints/factors and number of replicates are typically inversely correlated (economic considerations)
- For rather small effects such experiments are typically underpowered (on the level of endpoint)
- Often inconclusive
- But perfect for narrowing the search space
- Figure 1.

Experiments in the context of a project

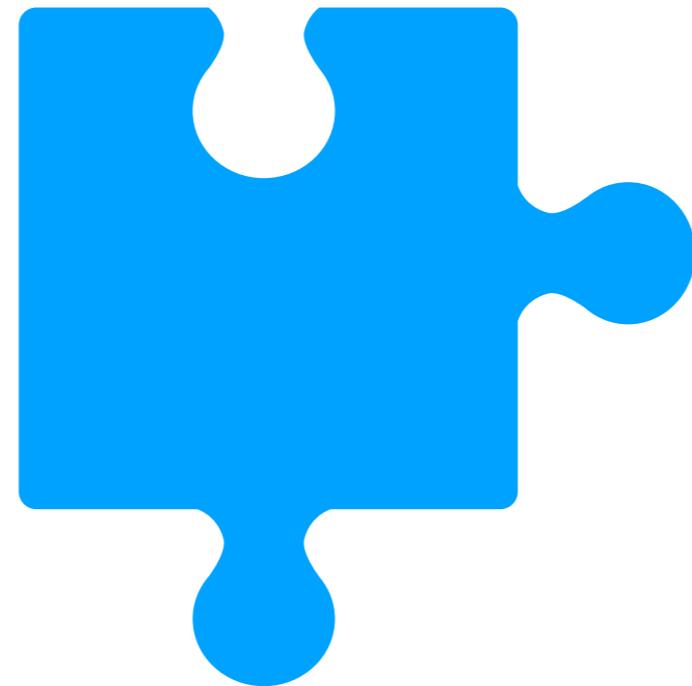
How to build the story for publication

Single Experiment

medical phase III/IV trial



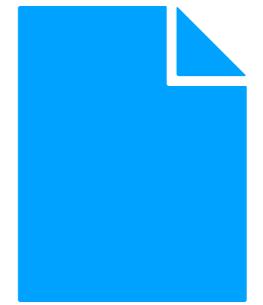
**Strong
Hypothesis**



**RCT
 $N >> 100$**



**Strong
Conclusion**



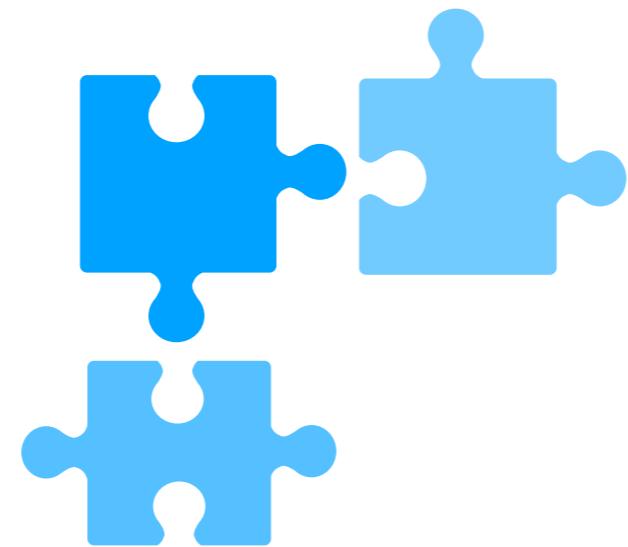
**Strong
Paper**

Preregistration

Triangulation



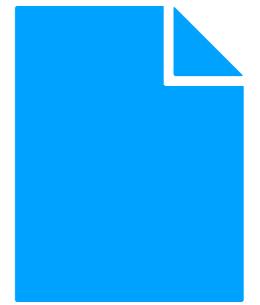
Vague idea
weak Hypothesis



Ensemble experiments
Of varying evidence strengths



Conclusion

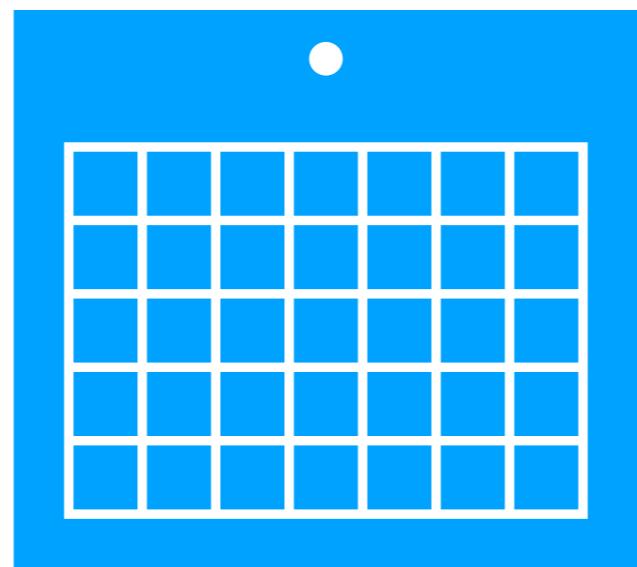


Paper

Screening

?

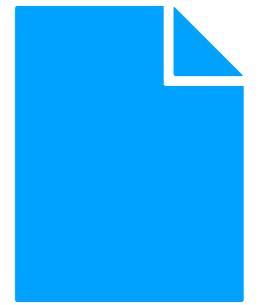
Vague idea
weak Hypothesis



High throughput
Omics
 $N < 3$

!

Inconclusive

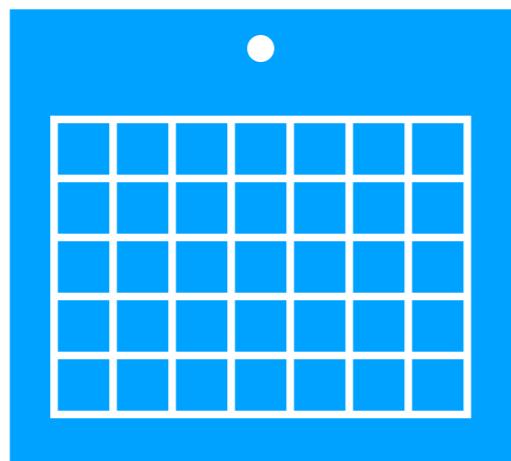


Paper

Screening+

**Ensemble experiments
Of varying evidence strengths**

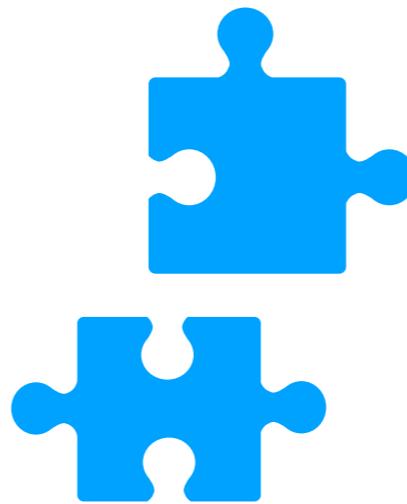
?



**High throughput
Omics
 $N < 3$**

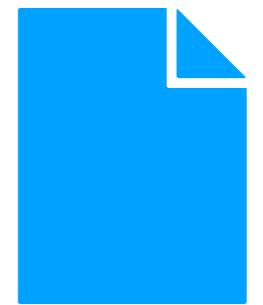
**Vague idea
weak Hypothesis**

?



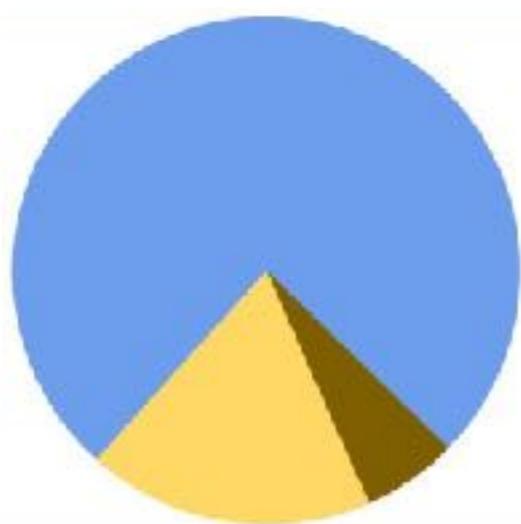
Strong Hypothesis

!



**Strong
Paper**

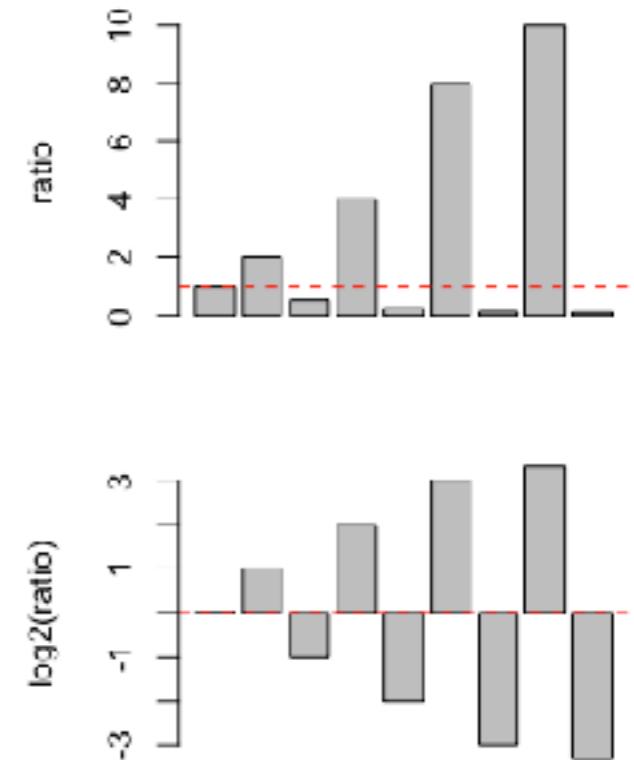
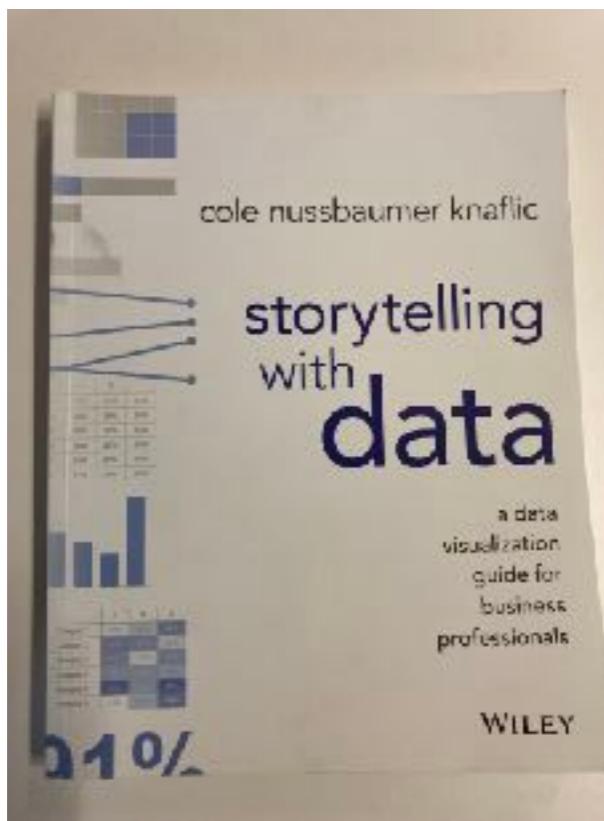
Conclusion



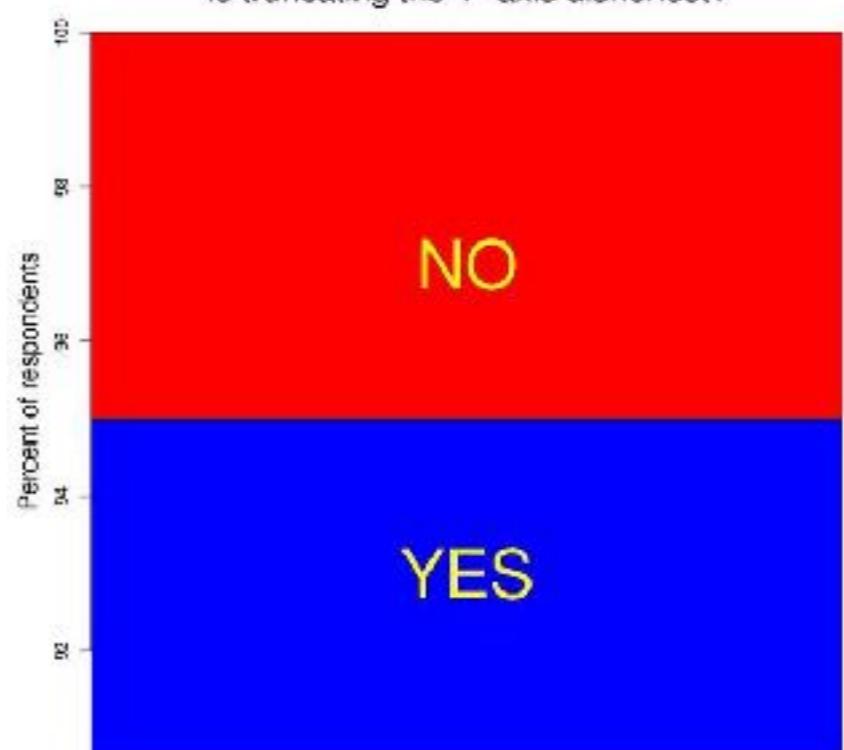
- Sky
- Sunny side of pyramid
- Shady side of pyramid

Visualisation

Another day



Is truncating the Y-axis dishonest?



Publication

Journals now have check lists



- How was the **sample size** chosen to ensure adequate **power** to detect a pre-specified effect size?
- For animal studies, include a statement about sample size estimate even if no statistical methods were used.
- Describe **inclusion/exclusion criteria** if samples or animals were excluded from the analysis. Were the criteria pre-established?
- Were any steps taken to minimise the effects of subjective bias when allocating animals/samples to treatment (e.g. **randomisation** procedure)? If yes, please describe.
For animal studies, include a statement about randomisation even if no randomisation was used.
- Were any steps taken to minimise the effects of subjective bias during group allocation or/and when assessing results (e.g. **blinding** of the investigator)? If yes please describe.
For animal studies, include a statement about blinding even if no blinding was done
- For every figure, are statistical tests justified as appropriate?
Do the data meet the assumptions of the tests (e.g., normal distribution)? Describe any methods used to assess it. Is there an **estimate of variation within each group** of data? Is the variance similar between the groups that are being statistically compared?

Making
better Experiments

- Think ahead, define strategy
- Be representative
- Aim for conclusion: Triangulation and/or confirmatory experiments
- Don't overdo screening, KISS towards the end.
- Experimental Design
- Open Science practices
(Transparency, Storage, Analytical Reproducibility)
- Less time, less money, more reliable/sustainable results.

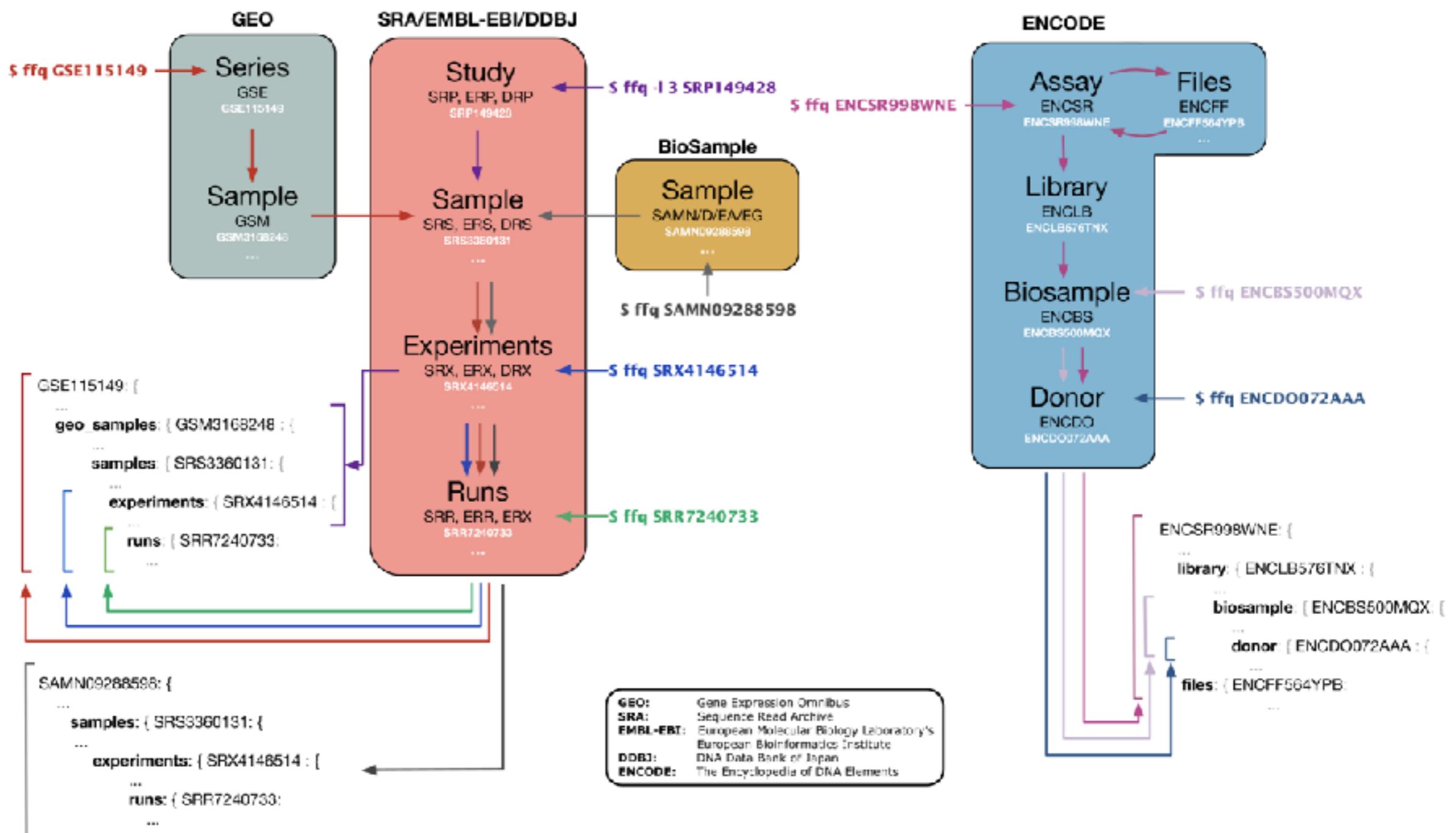
Analysis of Public Transcriptomic Data Sets

Tobias Straub

tobias.straub@lmu.de

Biomedizinisches Centrum, LMU München

Where the data @



ncbinlm.nih.gov

NCBI GEO Accession Display

GEO help: Mouse over screen elements for information.

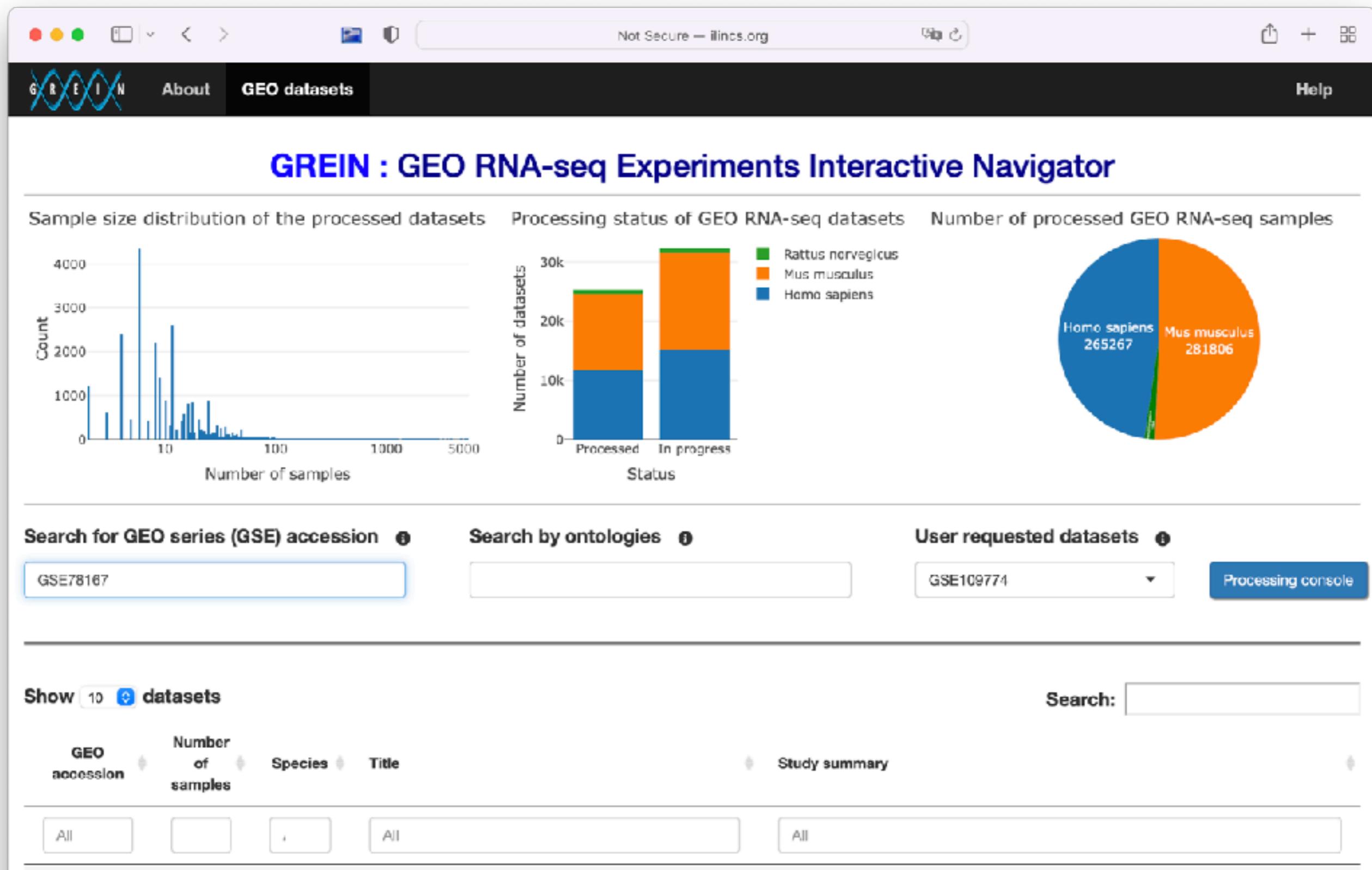
Scope: Self Format: HTML Amount: Quick GEO accession: GSE78167

Series GSE78167 Query DataSets for GSE78167

Status	Public on Aug 19, 2016
Title	An integrative transcriptomics approach identifies miR-503 as a candidate master regulator of the estrogen response [RNA-seq]
Organism	Homo sapiens
Experiment type	Expression profiling by high throughput sequencing
Summary	<p>Estrogen receptor α (ERα) is an important biomarker of breast cancer severity and a common therapeutic target. Recent studies have demonstrated that in addition to its role in promoting proliferation, ERα also protects tumors against metastatic transformation. Current therapeutics antagonize ERα and interfere with both beneficial and detrimental signalling pathways stimulated by ERα. The goal of this study is to uncover the dynamics of coding and non-coding RNA (microRNA) expression in response to estrogen stimulation and identify potential therapeutic targets that more specifically inhibit ERα-stimulated growth and survival pathways without interfering with its protective features. To achieve this, we exposed MCF7 cells (an estrogen receptor positive model cell line for breast cancer) to estrogen and prepared a time course of paired mRNA and miRNA sequencing libraries at ten time points throughout the first 24 hours of the response to estrogen. From these data, we identified three primary expression trends—transient, induced, and repressed—that were each enriched for genes with distinct cellular functions. Integrative analysis of paired mRNA and microRNA temporal expression profiles identified miR-503 as the strongest candidate master regulator of the estrogen response, in part through suppression of ZNF217—an oncogene that is frequently amplified in cancer. We confirmed experimentally that miR-503 directly targets ZNF217 and that over-expression of miR-503 suppresses breast cancer cell proliferation. Overall, these data indicate that miR-503 acts as a potent estrogen-induced tumor suppressor microRNA that opposes cellular proliferation and has promise as a therapeutic for breast cancer. More generally, our work provides a systems-level framework for identifying functional interactions that shape the temporal dynamics of gene expression.</p>
Overall design	Quantification of mRNAs in MCF7 cells responding to estrogen following a period of estrogen starvation. Three independent biological replicates (30 samples: 3 replicates x 10 time points) of MCF7 cells were exposed to 10nM Estradiol for 0, 1, 2, 3, 4, 5, 6, 8, 12, or 24 hours, and total RNA was extracted from the samples. Total RNA was used to generate paired RNA and miRNA sequencing. RNA libraries were prepared using an Illumina TruSeq stranded mRNA library preparation kit.
Contributor(s)	Baran-Gale J , Sethupathy P , Purvis J
Citation(s)	Baran-Gale J, Purvis JE, Sethupathy P. An integrative transcriptomics approach identifies miR-503 as a candidate master regulator of the estrogen response in MCF-7 breast cancer cells. <i>RNA</i> 2016 Oct;22(10):1592-603. PMID: 27539763
Submission date	Feb 22, 2016
Last update date	May 15, 2019
Contact name	Jeanette Baran-Gale
E-mail(s)	jbaran@email.unc.edu
Organization name	University of North Carolina at Chapel Hill
Lab	Sethupathy & Purvis Labs
Street address	120 Mason Farm Rd.
City	Chapel Hill

- The easy way - using convenience tools
 - Convenient
 - Fixed genome/annotation/tools
 - Potentially compromised by bad metadata
 - Reproducibility questionable (documentation)
 - Black box
 - Accessibility/availability not guaranteed
- The hard way - running analyses from scratch
 - Steep learning curve (but huge reward)
 - Flexible
 - Deposition flaws can be fixed
 - Reproducible, re-usable, time-efficient
 - Publication quality figures

<http://www.ilincs.org/apps/grein/>



Not Secure — illncs.org

processed. Please see the following table.

Show 10 datasets

Search:

GEO accession	Number of samples	Species	Title	Study summary
All	All	All	All	All
GSE78167	30	<i>Homo sapiens</i>	An integrative transcriptomics approach identifies miR-503 as a candidate master regulator of the estrogen response [RNA-seq]	Estrogen receptor α (ER α) is an important biomarker of breast cancer severity and a common therapeutic target. Recent studies have demonstrated that in addition to its role in promoting proliferation, ER α also protects tumors against metastatic transformation. Current therapeutics antagonize ER α and interfere with both beneficial and detrimental signaling pathways stimulated by ER α . The goal of this study is to uncover the dynamics of coding and non-coding RNA (microRNA) expression in response to estrogen stimulation and identify potential therapeutic targets that more specifically inhibit ER α -stimulated growth and survival pathways without interfering with its protective features. To achieve this, we exposed MCF7 cells (an estrogen receptor positive model cell line for breast cancer) to estrogen and prepared a time course of paired mRNA and miRNA sequencing libraries at ten time points throughout the first 24 hours of the response to estrogen. From these data, we identified three primary expression trends—transient, induced, and repressed—that were each enriched for genes with distinct cellular functions. Integrative analysis of paired mRNA and microRNA temporal expression profiles identified miR-503 as the strongest candidate master regulator of the estrogen response, in part through suppression of ZNF217—an oncogene that is frequently amplified in cancer. We confirmed experimentally that miR-503 directly targets ZNF217 and that over-expression of miR-503 suppresses breast cancer cell proliferation. Overall, these data indicate that miR-503 acts as a potent estrogen-induced tumor suppressor microRNA that opposes cellular proliferation and has promise as a therapeutic for breast cancer. More generally, our work provides a systems-level framework for identifying functional interactions that shape the temporal dynamics of gene expression.

Showing 1 to 1 of 1 datasets

Previous 1 Next

Not Secure — ilincs.org

G R E I N About GEO datasets Explore dataset Analyze dataset Help

Selected study

GSE78167

	Description	Metadata	Counts table	QC report	Visualization
Study link	GSE78167				
No. of GEO samples	30				
No. of SRA runs	30				
Species	Homo sapiens				
Title	An integrative transcriptomics approach identifies miR-503 as a candidate master regulator of the estrogen response [RNA-seq]				
Summary	<p>Estrogen receptor α (ERα) is an important biomarker of breast cancer severity and a common therapeutic target. Recent studies have demonstrated that in addition to its role in promoting proliferation, ERα also protects tumors against metastatic transformation. Current therapeutics antagonize ERα and interfere with both beneficial and detrimental signaling pathways stimulated by ERα. The goal of this study is to uncover the dynamics of coding and non-coding RNA (microRNA) expression in response to estrogen stimulation and identify potential therapeutic targets that more specifically inhibit ERα-stimulated growth and survival pathways without interfering with its protective features. To achieve this, we exposed MCF7 cells (an estrogen receptor positive model cell line for breast cancer) to estrogen and prepared a time course of paired mRNA and miRNA sequencing libraries at ten time points throughout the first 24 hours of the response to estrogen. From these data, we identified three primary expression trends—transient, induced, and repressed—that were each enriched for genes with distinct cellular functions. Integrative analysis of paired mRNA and microRNA temporal expression profiles identified miR-503 as the strongest candidate master regulator of the estrogen response, in part through suppression of ZNF217—an oncogene that is frequently amplified in cancer. We confirmed experimentally that miR-503 directly targets ZNF217 and that over-expression of miR-503 suppresses breast cancer cell proliferation. Overall, these data indicate that miR-503 acts as a potent estrogen-induced tumor suppressor microRNA that opposes cellular proliferation and has promise as a therapeutic for breast cancer. More generally, our work provides a systems-level framework for identifying functional interactions that shape the temporal dynamics of gene expression.</p>				

Not Secure — lincs.org

GREIN About GEO datasets Explore dataset Analyze dataset Help

Selected study

GSE78167

Data type
 Raw Normalized

Number of samples to show
30

Show counts table

Download data

Gene level Transcript level

Description Metadata Counts table QC report Visualization

Select any row to see boxplot of the selected gene below the table.

Show 8 genes Search:

Gene symbol	GSM2068643	GSM2068644	GSM2068645	GSM2068646	GSM2068647	GSM2068648
ENSG000000000003	TSPAN6	799	568	806	727	610
ENSG000000000005	TNMD	0	0	0	0	0
ENSG00000000419	DPM1	2230	1579	2423	2246	2134
ENSG00000000457	SCYL3	816	410	533	521	595
ENSG00000000460	C1orf112	375	205	345	291	456
ENSG00000000938	FGR	1	0	35	28	6
ENSG00000000971	CFH	1	0	1	0	0
ENSG00000001036	FUCA2	2716	1832	2888	2567	2909

Showing 1 to 8 of 27,990 genes

Previous 1 2 3 4 5 ... 3499 Next

Not Secure — ilincs.org

G R E I N About GEO datasets Explore dataset Analyze dataset Help

Selected study GSE78167

Description Metadata Counts table QC report Visualization

[Download QC report](#)

MultiQC v1.2

General Stats

Salmon

FastQC

Sequence Quality Histograms

Per Sequence Quality Scores

Per Base Sequence Content

Per Sequence GC Content

Per Base N Content

Sequence Length Distribution

Sequence Duplication Levels

Overrepresented sequences

Adapter Content

MultiQC

- /GSE78167/fastqc
- /GSE78167/salmon

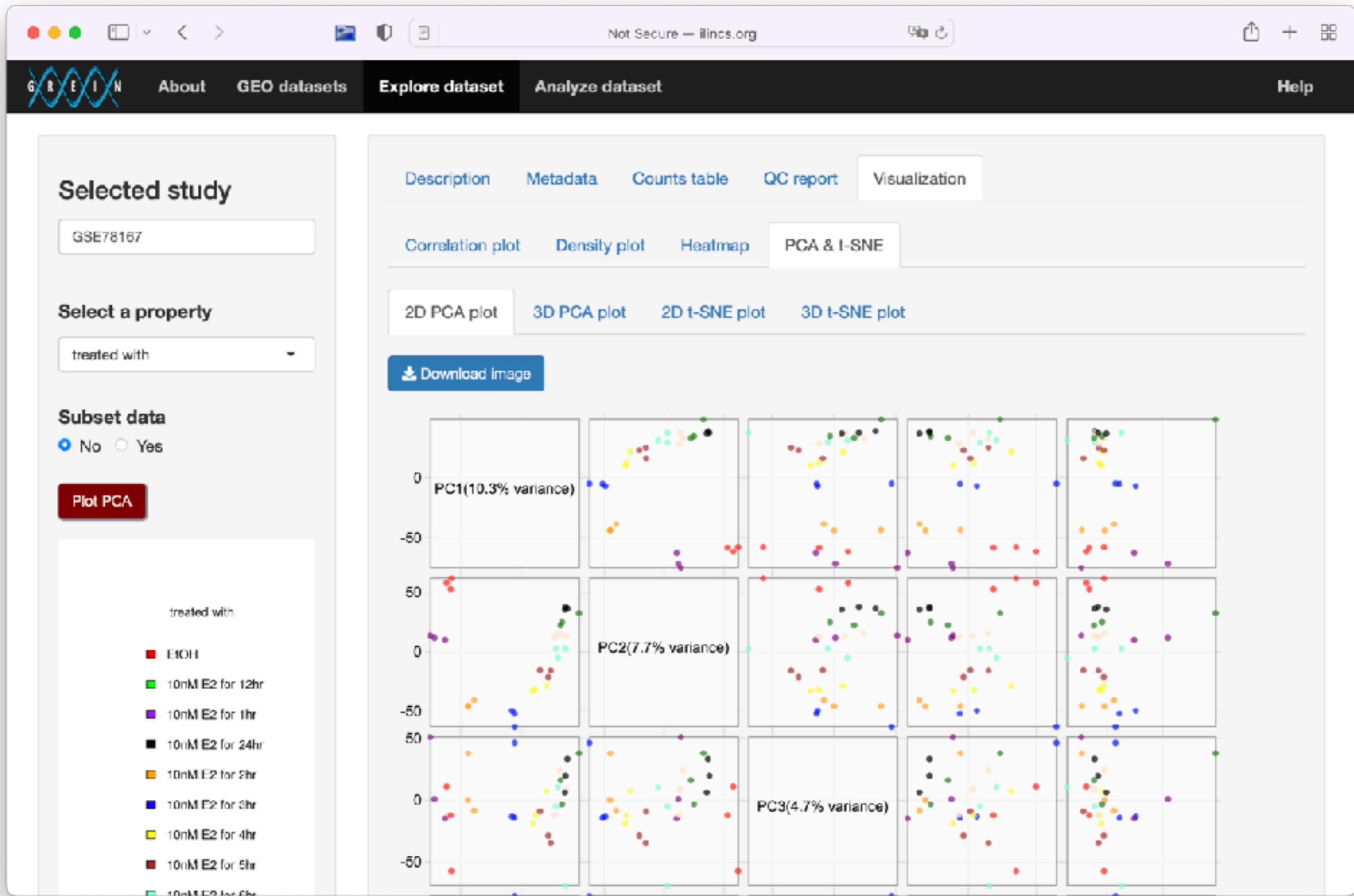
• Welcome! Not sure where to start? [Watch a tutorial video \(6:06\)](#) [don't show again](#)

General Statistics

Copy table | Configure Columns | Plot | Showing 10/93 rows and 5/7 columns.

Sample Name	% Aligned	M Aligned	%
SRR3182145_GSM2068643_transcripts_quant	92.0%	33.6	66
SRR3182145_pass_1			64
SRR3182145_pass_2			62
SRR3182146_GSM2068644_transcripts_quant	90.5%	22.7	60
SRR3182146_pass_1			58
SRR3182146_pass_2			56
SRR3182147_GSM2068645_transcripts_quant	90.0%	22.8	54

Toolbox



Selected study

GSE78167

Factor of interest

treated with

Sample selection

All samples

Specific samples

Experimental group

10nM E2 for 2hr

Control group

EtOH

Type of comparison

Two group without covariate

Subset samples

No Yes

Generate signature

Not Secure — illncs.org

Create a signature Power analysis

Metadata

Show 6 samples

Search:

	Selected groups	treated with	characteristics	passage
GSM2068643	Control	EtOH	MCF7_0hrs_post_E2	p9
GSM2068645	Experimental	10nM E2 for 2hr	MCF7_2hrs_post_E2	p9
GSM2068653	Control	EtOH	MCF7_0hrs_post_E2	p9
GSM2068655	Experimental	10nM E2 for 2hr	MCF7_2hrs_post_E2	p9
GSM2068663	Control	EtOH	MCF7_0hrs_post_E2	p10
GSM2068665	Experimental	10nM E2 for 2hr	MCF7_2hrs_post_E2	p10

Showing 1 to 6 of 6 samples

Previous 1 Next

Not Secure — ilincs.org

Selected study

GSE78167

[Signature visualization](#)

[Download signature](#)

Upload signature to iLINCS

[Upload all genes](#)

[Upload](#)

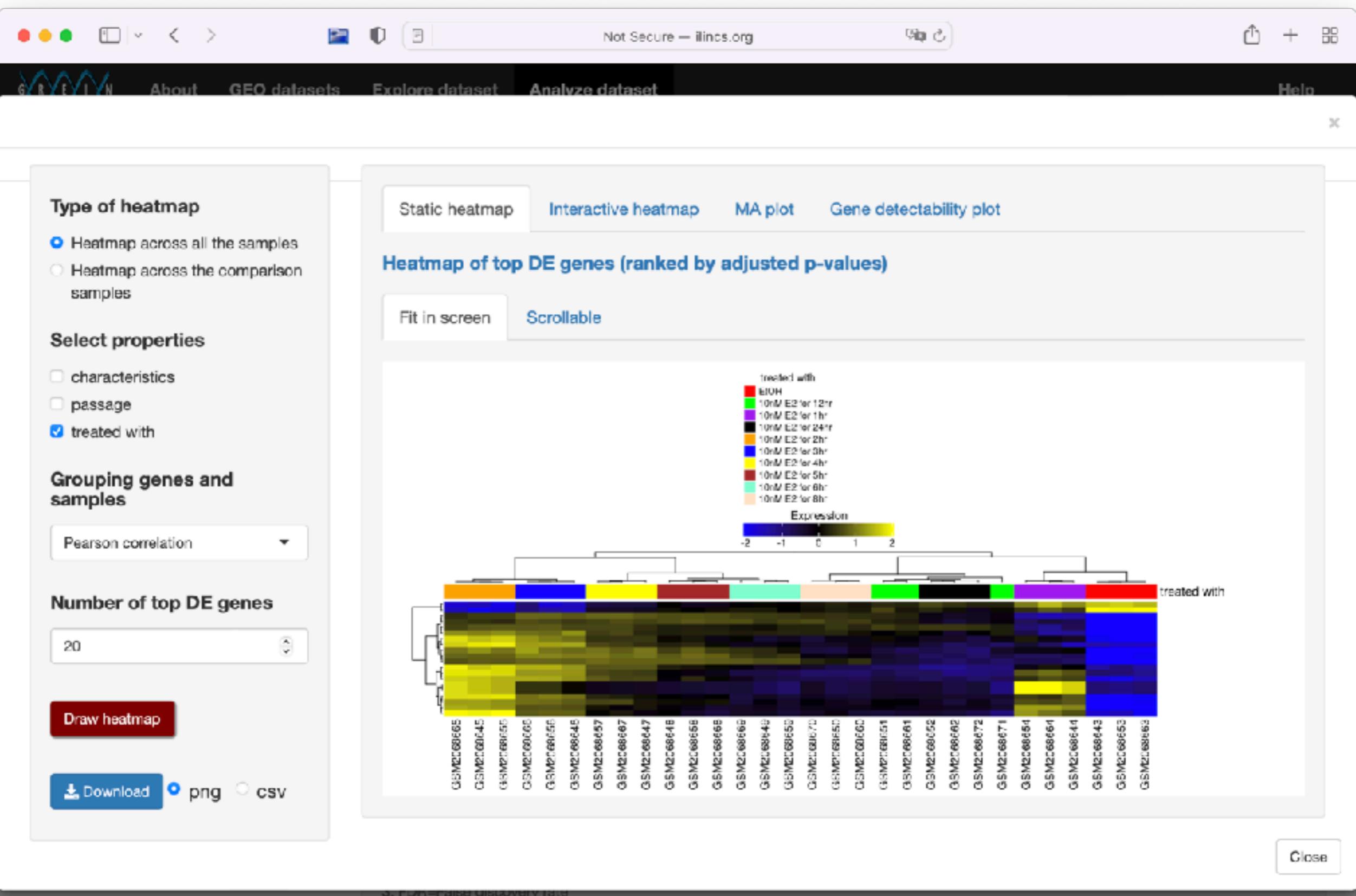
Create a signature [Power analysis](#)

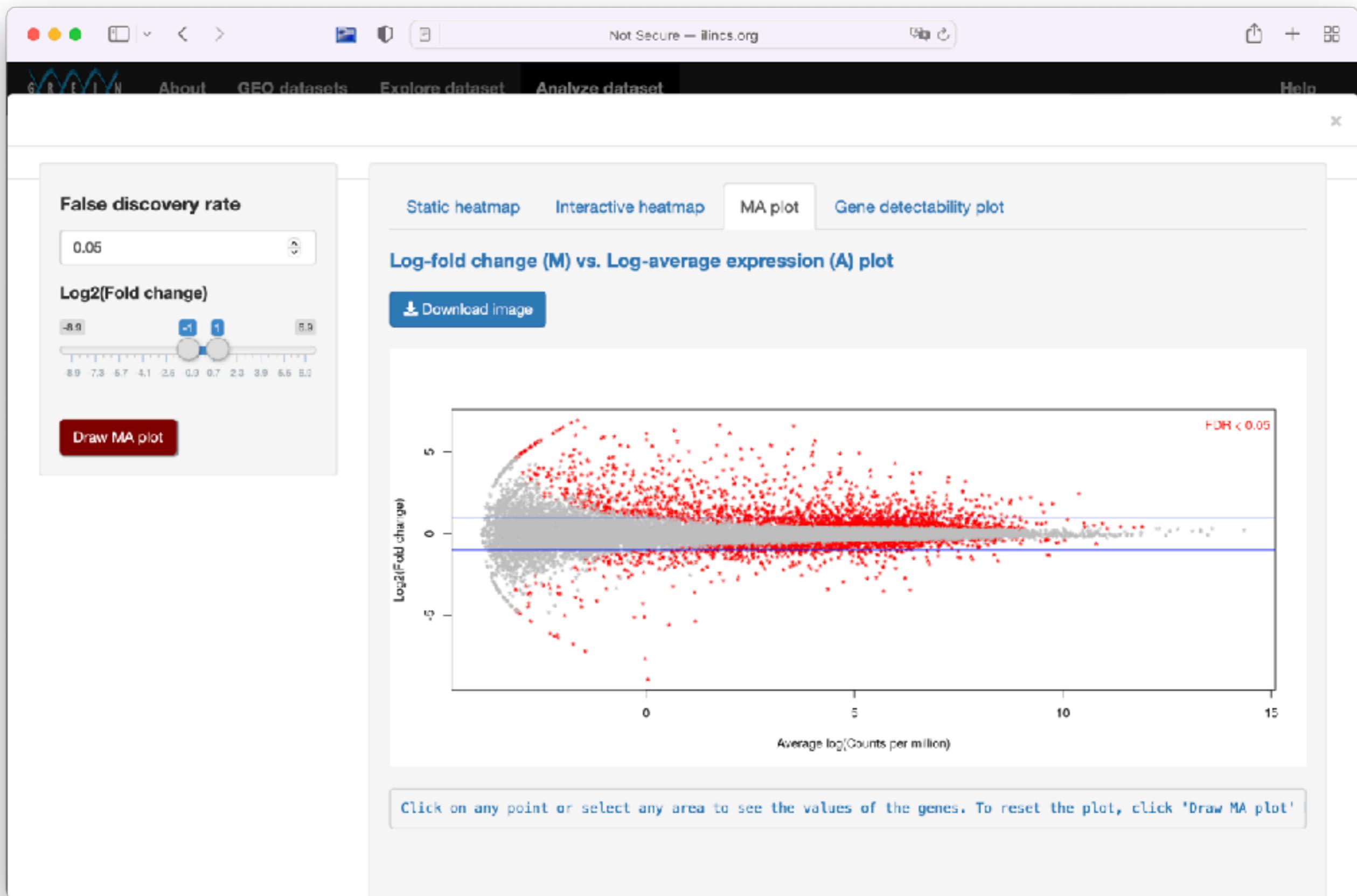
Metadata [Signature](#)

Show 10 genes [Search:](#)

Ensembl_ID	Gene_symbol	logFC	logCPM	PValue	FDR
All	All	All	All	All	All
ENSG00000148344	PTGES	4.933	5.122	5.24e-195	1.087e-190
ENSG00000165731	RET	3.176	7.235	2.615e-172	2.711e-168
ENSG00000175591	P2RY2	3.451	7.259	2.268e-165	1.568e-161
ENSG00000099337	KCNK6	3.203	7.657	5.967e-165	3.093e-161
ENSG00000198959	TGM2	5.485	4.008	8.701e-155	3.609e-151
ENSG00000108352	RAPGEFL1	3.123	6.394	1.133e-152	3.917e-149
ENSG00000164626	KCNK5	3.687	6.402	1.114e-146	3.299e-143
ENSG00000188176	SMTNL2	4.908	4.658	2.755e-144	6.402e-141
ENSG00000108375	RNF43	-3.506	5.701	2.778e-144	6.402e-141
ENSG00000115738	ID2	-3.447	6.34	1.667e-143	3.457e-140

1. logFC=Log fold change
2. logCPM=Log counts per million
3. FDR=False discovery rate





Not Secure — lincs.org

XXX GREEN

DCIC integrative LINCS genomics data portal | LINCS

iLINCS Signatures Datasets Genes iLINCS Paper new

Search for signatures / Upload a signature / Uploaded Signature

Uploaded Signature

Signature analysis

Modify the list of selected genes >

Other analyses with selected genes >

Signature Info

Session ID: Sat_May_28_09_46_50_2022_580356

File name: GSE78167_signatureData_13:43:43_2022-05-28_10124_up.txt

Genes not Found: C6ORF132, C14ORF132, FCMR, C1ORF228, C10ORF2, C2ORF54, DISP3, C15ORF69, C1ORF109, C1ORF111, PHF24, C8ORF46, C4ORF19, C11... [More](#)

Found 9471 out of 9873 submitted entries.

Complete signature (9471) Selected genes (100)

Download

Signature Analysis Tools ▶ Signature Data ▶ Connected Signatures ⓘ ▶ Connected Perturbations ⓘ ▶

Pathway Analysis ▶

Enrichr DAVID ToppFun Reactome

Network Analysis ▶

reactome.org

XXX GREIN DCIC integrative LINCS genomics data portal | LINCS PB | Estrogen-dependent nuclear events downstream of ESR-membr...

Pathways for: Homo sapiens Citation: Analysis: Tour: Layout:

Event Hierarchy:

- ESTG binds ESR2-chaperone complex
- HSP90-dependent ATP hydrolysis
- ESR dimerizes
- 27-hydroxysterol binds ESR1, ESR2
- + Estrogen-dependent gene expression
- + Extra-nuclear estrogen signaling
 - PRMT1 methylates ESRs
 - HSBP1 oligomer binds ESRs
 - ZDHHC7, ZDHHC21 palmitoylates ESRs
 - Palms-ESRs bind CAVs
 - Palms-ESRs:CAVs translocate to nucleus
 - Estrogen stimulates dimerization
 - ESR binds STRN
 - ESTG binding induces ESR dimerization
 - Heterotrimeric G protein (I) binds ESR
 - G-proteins dissociate from plasma membrane
 - Membrane estrogen receptor
 - ESR-associated SRC autophosphorylation
 - PI3K binds membrane-associated ESR
 - cNOS synthesizes NO
 - PTK2 is recruited to methylated ESR
 - + Estrogen-stimulated signaling
 - RUNX3 regulates WNT signaling
 - Signaling by Nuclear Receptors
 - ATF6 (ATF6-alpha) activates chaperone genes
 - ATF6 (ATF6-alpha) activates chaperones
 - Unfolded Protein Response (UPR)
 - Transcriptional Regulation by VENTX
 - + Estrogen-dependent nuclear events downstream of ESR-membrane signaling
 - MMPs cleave HB-EGF
 - PTK2 binds activated EGFR
 - PTK2 autophosphorylates downstream of ESR

Search for a term, e.g. pten ...

Signal Transduction

4.56E0

1.71E0

Description Molecules Structures Expression Analyses (34) Downloads

Expression analysis results for ENSEMBL [Data: Genes]

Pathway name	Entities found	Entities Total	Entities ratio	pValue	FDR	Reactions found	Reactions total
Estrogen-dependent nuclear events downstream of ESR-membrane signaling	2	3	0.002	6.02E-3	8.54E-2	4	5
Extra-nuclear estrogen signaling	2	3	0.002	5.02E-3	8.54E-2	4	5
ESR-mediated signaling	4	23	0.016	8.5E-3	9.35E-2	13	66
Estrogen-dependent gene expression	3	22	0.015	4.16E-2	3.33E-1	9	61
RUNX3 regulates WNT signaling	1	2	0.001	6.75E-2	4.06E-1	2	4
Signaling by Nuclear Receptors	4	49	0.084	9.12E-2	4.21E-1	13	126
ATF6 (ATF6-alpha) activates chaperone genes	1	5	0.003	1.6E-1	4.21E-1	1	5
ATF6 (ATF6-alpha) activates chaperones	1	5	0.003	1.6E-1	4.21E-1	1	5
Unfolded Protein Response (UPR)	4	61	0.043	1.63E-1	4.21E-1	4	61
Transcriptional Regulation by VENTX	1	7	0.005	2.17E-1	4.21E-1	2	13

1-20 of 34

Not Secure — lincs.org

GREIN DCIC integrative LINCS genomics data portal | LINCS

iLINCS Signatures Datasets Genes iLINCS Paper new

[Search for signatures](#) / [Upload a signature](#) / [Uploaded Signature](#) / [Pathway Analysis](#)

SPIA Functional Pathway Analysis

KEGG pathway name	KEGG pathway ID	Genes in Pathway	DE Genes in Pathway	Topology Score	KEGG link					
	ID	Pathway	Pathway	ORA pval	Top pval	SPIA pval	SPIA adj pval	Status	Link	
Intestinal immune network for IgA production	hsa04672	5	2	0.0024	-0.26	0.818	0.0144	0.3998	Inhibited	View
Osteoclast differentiation	hsa04380	52	3	0.0195	-13.9124	0.075	0.0215	0.3998	Inhibited	View
Amoebiasis	hsa05146	14	1	0.2017	-4.904	0.03	0.037	0.3998	Inhibited	View
cGMP-PKG signaling pathway	hsa04022	67	2	0.2697	19.8674	0.028	0.0472	0.3998	Activated	View
Hepatitis B	hsa05161	56	1	0.5951	20.3431	0.014	0.0482	0.3998	Activated	View
Taste transduction	hsa04742	4	1	0.0523	3.587	0.153	0.0539	0.3998	Activated	View
RNA degradation	hsa03018	5	1	0.0772	0	NA	0.0772	0.3998	Inhibited	View
Sphingolipid signalling pathway	hsa04071	40	2	0.1331	15.4392	0.112	0.0776	0.3998	Activated	View
Axon guidance	hsa04360	89	3	0.0971	16.8184	0.166	0.0782	0.3998	Activated	View
T cell receptor signalling pathway	hsa04660	41	1	0.4638	20.93	0.041	0.0976	0.4489	Activated	View

5 10 50 First < 1 2 3 4 5 > Last Page 1 of 5, of 43 entries

[Download data](#)

The
hard but rewarding
way

```
@BS-DSU-ELLAC_0:4:1:6214:930
NNTCCTGGCTGGTAGCTTAAATAATAGAGCTTAA
+
#****(*)+(@@@@1211557755.,55575775@
@BS-DSU-ELLAC_0:4:1:8611:931
NNACAAATGAGCGTGAGCTTCTGCCATCTTATGGG
+
#####
@BS-DSU-ELLAC_0:4:1:9150:948
NNTTTATATTCTAATTACATATGTACAAAAGTT
```

Raw data
+ Metadata

+ Genomic data

UNIX/LINUX

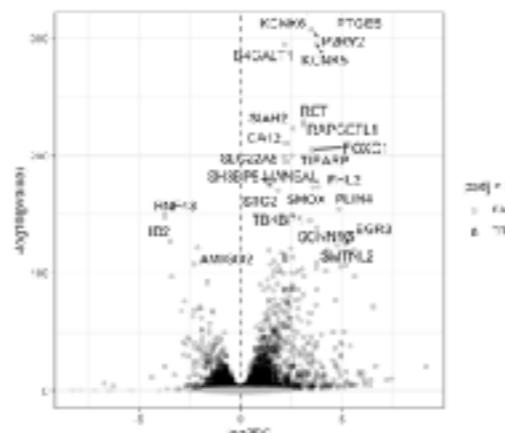
Pre-processing (aka the engineering part)

Processed data

SummarizedExperiment

Data analysis (the scientific part)

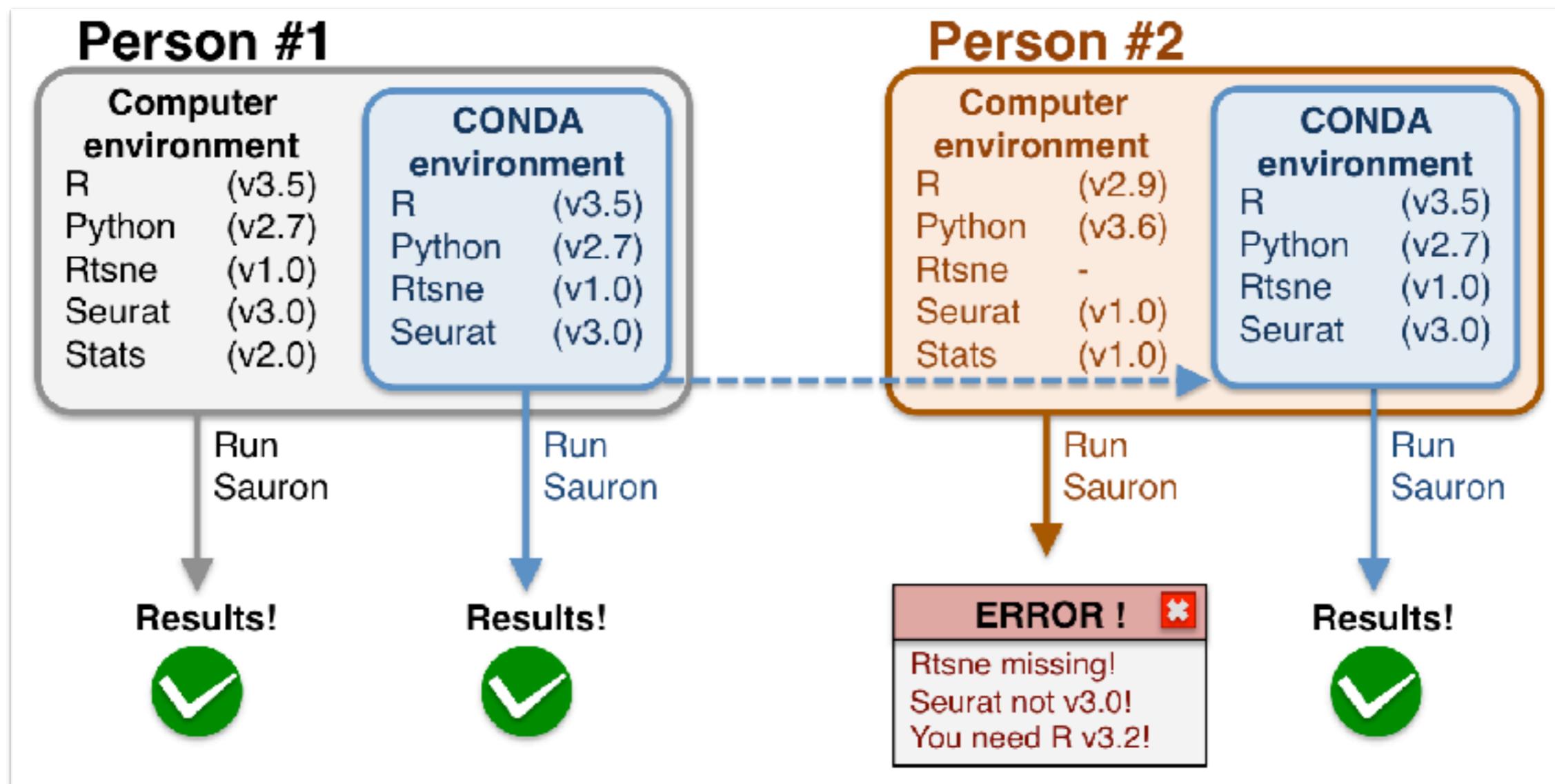
R/python



Requirements

- Desire, commitment
- Computer hardware / good Network connection
- Skills:
 - Google, copy, paste skills
 - (Very basic) Unix
 - Software mangagement with conda
 - Snakemake
 - R/bioconductor

Conda



Snakemake

Tell Snakemake what files you want to be created

Produce the files you want to have from some intermediate result

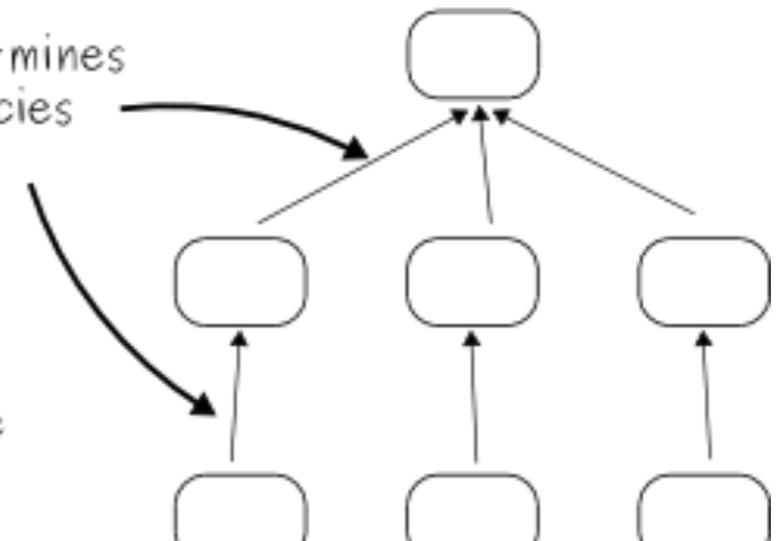
Create a needed intermediate result

```
rule:  
    input: "A.txt", "B.txt", "C.txt"
```

```
rule:  
    input: "{sample}.inter"  
    output: "{sample}.txt"  
    shell: "somecommand {input} {output}"
```

```
rule:  
    input: "{sample}.in"  
    output: "{sample}.inter"  
    run:  
        somepythoncode()
```

Snakemake determines the dependencies for you



Use wildcards to write general rules for all samples



ncbi.nlm.nih.gov

Library preparation kit with 50 ng of depleted RNA; Strand orientated RNA-Seq

Library strategy: RNA-Seq
Library source: transcriptomic
Library selection: cDNA
Instrument model: Illumina HiSeq 2500

Description: growth in M63 minimal media plus 0.4% glucose plus 50 µM NiCl2
Data processing: Basecalls performed using HCS 2.0.5 and RTA 1.17.20
RNA-seq reads were aligned to the W3110 genome using CASAVA 1.8.2
Gene expression (based on known genes) was determined using Cufflinks 2.0.2, only PF reads were retained
Genome_build: W3110 (NC_007779)
Supplementary_files_format_and_content: excel file with coverage and FPKM measurements

Submission date: Dec 18, 2015
Last update date: May 15, 2019
Contact name: Agnès Rodriguez
E-mail(s): agnes.rodrigue@insa-lyon.fr
Organization name: INSA Lyon CNRS
Lab: MAP UMR 5240
Street address: 10 rue Dubois
City: Villeurbanne
ZIP/Postal code: 69621
Country: France

Platform ID: GPL18133
Series (1): [GSE76167](#) Mechanisms of nickel toxicity in bacteria

Relations
BioSample: SAMN04351201
SRA: SRK1492019

Supplementary data files not provided

SRA Run Selector [?](#)
Raw data are available in SRA
Processed data are available on Series record

| NLM | NIH | GEO Help | Disclaimer | Accessibility |
HHS Vulnerability Disclosure

NCBI SRA Run Selector <https://www.ncbi.nlm.nih.gov/Tracesearch/Recs/PRJNA312817?max=1000>

Filter List

- 1 Bioproject
- 2 BYTES
- 3 Passage
- 4 source_name
- 5 treated_with

Accession PRJNA312817 [Search](#)

Common Fields

BioProject	PRJNA312817
Consort	PUBLIC
Assay Type	RNA-Seq
Avg. Bases	96
CellLine	MCF7
Center Name	GEO
DATA SOURCE TYPE	PASID: SRA
DATASTORE provider	OS NCBI SRA
DATASTORE region	gs.US.ncbi-public.s3.us-east-1

Select

	Runs	Bytes	Bases	Download	Cloud DataDelivery	Copying
Total	30	5159 Gb	9MBC	Metadata or Accession List		
Selected	0	0	0	Metadata or Accession List or JWT Cart	Deliver Data	Gallery

X Found 30 items [Search with file results](#)

#	Run	# Bases	# Bytes	# Experiment	GEO_Accession	Passage	# Sample Name	source_name	treated_with
1	SRR3132145	340G	2.18 Gb	SRK1556372	GSM2068542	p9	GSM2068543	MCF7_0hrs_port_E2	E0H
2	SRR3132146	252G	1.51 Gb	SRK1556373	GSM2068544	p9	GSM2068544	MCF7_1hrs_port_E2	10mE2 for 1hr
3	SRR3132147	392G	2.35 Gb	SRK1556374	GSM2068545	p9	GSM2068545	MCF7_2hrs_port_E2	10mE2 for 2hr
4	SRR3132148	301G	1.80 Gb	SRK1556375	GSM2068546	p9	GSM2068546	MCF7_3hrs_port_E2	10mE2 for 3hr
5	SRR3132149	314G	1.88 Gb	SRK1556376	GSM2068547	p9	GSM2068547	MCF7_4hrs_port_E2	10mE2 for 4hr
6	SRR3132150	340G	2.16 Gb	SRK1556377	GSM2068548	p9	GSM2068548	MCF7_5hrs_port_E2	10mE2 for 5hr
7	SRR3132151	300G	1.81 Gb	SRK1556378	GSM2068549	p9	GSM2068549	MCF7_6hrs_port_E2	10mE2 for 6hr
8	SRR3132152	264G	1.57 Gb	SRK1556379	GSM2068550	p9	GSM2068550	MCF7_7hrs_port_E2	10mE2 for 7hr
9	SRR3132153	330G	2.09 Gb	SRK1556380	GSM2068551	p9	GSM2068551	MCF7_12hrs_port_E2	10mE2 for 12hr
10	SRR3132154	330G	1.78 Gb	SRK1556381	GSM2068552	p9	GSM2068552	MCF7_24hrs_port_E2	10mE2 for 24hr
11	SRR3132155	319G	1.74 Gb	SRK1556382	GSM2068553	p9	GSM2068553	MCF7_0hrs_port_E2	E0H
12	SRR3132156	299G	1.75 Gb	SRK1556383	GSM2068554	p9	GSM2068554	MCF7_1hrs_port_E2	10mE2 for 1hr
13	SRR3132157	341G	2.06 Gb	SRK1556384	GSM2068555	p9	GSM2068555	MCF7_2hrs_port_E2	10mE2 for 2hr
14	SRR3132158	309G	1.87 Gb	SRK1556385	GSM2068556	p9	GSM2068556	MCF7_3hrs_port_E2	10mE2 for 3hr
15	SRR3132159	379G	2.25 Gb	SRK1556386	GSM2068557	p9	GSM2068557	MCF7_4hrs_port_E2	10mE2 for 4hr
16	SRR3132160	349G	2.11 Gb	SRK1556387	GSM2068558	p9	GSM2068558	MCF7_5hrs_port_E2	10mE2 for 5hr
17	SRR3132161	275G	1.69 Gb	SRK1556388	GSM2068559	p9	GSM2068559	MCF7_6hrs_port_E2	10mE2 for 6hr
18	SRR3132162	280G	1.67 Gb	SRK1556389	GSM2068560	p9	GSM2068560	MCF7_7hrs_port_E2	10mE2 for 7hr
19	SRR3132163	249G	1.60 Gb	SRK1556390	GSM2068561	p9	GSM2068561	MCF7_8hrs_port_E2	10mE2 for 8hr

Project — ~/Desktop/mount/work/project/tobias/GSE78167

Project

- GSE78167
 - .snakemake
 - GSM2068643_out
 - GSM2068644_out
 - GSM2068645_out
 - GSM2068646_out
 - GSM2068647_out
 - GSM2068648_out
 - GSM2068649_out
 - GSM2068650_out
 - GSM2068651_out
 - GSM2068652_out
 - GSM2068653_out
 - GSM2068654_out
 - GSM2068655_out
 - GSM2068656_out
 - GSM2068657_out
 - GSM2068658_out
 - GSM2068659_out
 - GSM2068660_out
 - GSM2068661_out
 - GSM2068662_out
 - GSM2068663_out
 - GSM2068664_out
 - GSM2068665_out
 - GSM2068666_out
 - GSM2068667_out
 - GSM2068668_out
 - GSM2068669_out
 - GSM2068670_out
 - GSM2068671_out
 - GSM2068672_out
 - logs
 - scripts
 - toSe.R
 - slurm
 - config.yaml
 - run.sh
 - se.rds
- Snakefile
- sra_kallisto_paired.yaml
- SraRunTable.txt

Snakefile

```

1 #!
2 configfile: "config.yaml"
3
4 import pandas as pd
5
6 sample_table = pd.read_table(config["runtable"], sep=",")
7 samples = list(sample_table['Sample Name'].unique())
8
9 # everything that need network connection for localrules: prefetch, get_index_files, create_se
10
11 rule all:
12     input:
13         "se.rds"
14
15 rule create_se:
16     input:
17         expand("{SRS}_out/{SRS}.abundance.tsv", SRS=samples),
18         "genome.gtf.gz"
19     output:
20         "se.rds"
21     script:
22         "scripts/toSe.R"
23
24 rule kallisto:
25     input:
26         files1=lambda wildcards: expand("raw/{sample}_1.fastq.gz",
27             sample=list(sample_table[sample_table['Sample Name']==wildcards.SRS].Run)),
28         files2=lambda wildcards: expand("raw/{sample}_2.fastq.gz",
29             sample=list(sample_table[sample_table['Sample Name']==wildcards.SRS].Run)),
30         gtf="genome.gtf.gz",
31         index="kallisto_index/genome"
32     output:
33         "{SRS}_out/{SRS}.abundance.h5",
34         "{SRS}_out/{SRS}.abundance.tsv"
35     threads: 32
36     params:
37         stranded=config['strandedness']
38     resources:
39         mem_mb=10000
40     shell:
41         """
42             if [params.stranded] == "reverse" :
43                 then
44                     kallisto quant --index={input.index} --gtf {input.gtf} --output-dir={wildcards.SRS}_out --rf-stranded --threads={threads} -b 30 {input.files1} {input.files2}
45             elif [params.stranded] == "forward" :
46                 then
47                     kallisto quant --index={input.index} --gtf {input.gtf} --output-dir={wildcards.SRS}_out --fr-stranded --threads={threads} -b 30 {input.files1} {input.files2}
48             else
49                 kallisto quant --index={input.index} --gtf {input.gtf} --output-dir={wildcards.SRS}_out --threads={threads} -b 30 {input.files1} {input.files2}
50             fi
51             mv {wildcards.SRS}_out/abundance.h5 {wildcards.SRS}_out/{wildcards.SRS}.abundance.h5
52             mv {wildcards.SRS}_out/abundance.tsv {wildcards.SRS}_out/{wildcards.SRS}.abundance.tsv
53         """
54
55
56 rule fastq_dump:
57     input:
58         "sra/{SRR}.sra"
59     output:
60         temp("raw/{SRR}_1.fastq.gz"),
61         temp("raw/{SRR}_2.fastq.gz")
62     shell:
63         """

```

Project — ~/Desktop/mount/work/project/tobias/GSE78167

Project

- GSE78167
 - > ■ .snakemake
 - > ■ GSM2068843_out
 - > ■ GSM2068844_out
 - > ■ GSM2068845_out
 - > ■ GSM2068846_out
 - > ■ GSM2068847_out
 - > ■ GSM2068848_out
 - > ■ GSM2068849_out
 - > ■ GSM2068850_out
 - > ■ GSM2068851_out
 - > ■ GSM2068852_out
 - > ■ GSM2068853_out
 - > ■ GSM2068854_out
 - > ■ GSM2068855_out
 - > ■ GSM2068856_out
 - > ■ GSM2068857_out
 - > ■ GSM2068858_out
 - > ■ GSM2068859_out
 - > ■ GSM2068860_out
 - > ■ GSM2068861_out
 - > ■ GSM2068862_out
 - > ■ GSM2068863_out
 - > ■ GSM2068864_out
 - > ■ GSM2068865_out
 - > ■ GSM2068866_out
 - > ■ GSM2068867_out
 - > ■ GSM2068868_out
 - > ■ GSM2068869_out
 - > ■ GSM2068870_out
 - > ■ GSM2068871_out
 - > ■ GSM2068872_out
 - > ■ logs
 - ✓ ■ scripts
 - toSe.R
 - > ■ slurm
 - config.yaml
 - run.sh
 - se.rds
 - Snakefile
- sra_kallisto_paired.yaml
- SraRunTable.txt

+ X sra_kallisto_paired.yaml 1:1 LF UTF-8 YAML GitHub

```
1 channels:
2   -- bioconda
3   -- conda-forge
4 dependencies:
5   -- snakemake-minimal =7.2
6   -- python
7   -- pandas
8   -- sra-tools =2.11
9   -- kallisto =0.48
10  -- r >=4.0
11  -- bioconductor-tximport
12  -- bioconductor-SummarizedExperiment
13  -- bioconductor-biomarR
14  -- bioconductor-GenomicFeatures
15
```

Project — ~/Desktop/mount/work/project/tobias/GSE78167

Project

- GSE78167
 - .snakemake
 - logs
 - series
 - sra_kallisto_paired.yaml
 - slurm
 - config.yaml
 - run.sh
 - Snakefile
 - sra_kallisto_paired.yaml
 - SraRunTable.txt

```

channels:
  - bioconda
  - conda-forge
dependencies:
  - snakemake-minimal >=7.2
  - python
  - pandas
  - sra-tools >=2.11
  - kallisto >=0.48
  - r >=4.0
  - bioconductor-tximport
  - bioconductor-SummarizedExperiment
  - bioconductor-biomaRt
  - bioconductor-GenomicFeatures

```

```

then
    kallisto quant --index=kallisto_index/genome --gtf genome.gtf.gz --output-dir=GSM2068659_out --fr-stranded --threads=1 -b 30 raw/SRR3182161_1.fastq.gz raw/SRR3182161_2.fastq.gz
else
    kallisto quant --index=kallisto_index/genome --gtf genome.gtf.gz --output-dir=GSM2068659_out --threads=1 -b 30 raw/SRR3182161_1.fastq.gz raw/SRR3182161_2.fastq.gz
fi
mv GSM2068659_out/abundance.h5 GSM2068659_out/GSM2068659.abundance.h5
mv GSM2068659_out/abundance.tsv GSM2068659_out/GSM2068659.abundance.tsv

```

[Sun May 29 15:54:10 2022]

```

localrule create_se:
    input: GSM2068643_out/GSM2068643.abundance.tsv, GSM2068644_out/GSM2068644.abundance.tsv, GSM2068645_out/GSM2068645.abundance.tsv, GSM2068646_out/GSM2068646.abundance.tsv, GSM2068647_out/GSM2068647.abundance.tsv, GSM2068648_out/GSM2068648.abundance.tsv, GSM2068649_out/GSM2068649.abundance.tsv, GSM2068650_out/GSM2068650.abundance.tsv, GSM2068651_out/GSM2068651.abundance.tsv, GSM2068652_out/GSM2068652.abundance.tsv, GSM2068653_out/GSM2068653.abundance.tsv, GSM2068654_out/GSM2068654.abundance.tsv, GSM2068655_out/GSM2068655.abundance.tsv, GSM2068656_out/GSM2068656.abundance.tsv, GSM2068657_out/GSM2068657.abundance.tsv, GSM2068658_out/GSM2068658.abundance.tsv, GSM2068659_out/GSM2068659.abundance.tsv, GSM2068660_out/GSM2068660.abundance.tsv, GSM2068661_out/GSM2068661.abundance.tsv, GSM2068662_out/GSM2068662.abundance.tsv, GSM2068663_out/GSM2068663.abundance.tsv, GSM2068664_out/GSM2068664.abundance.tsv, GSM2068665_out/GSM2068665.abundance.tsv, GSM2068666_out/GSM2068666.abundance.tsv, GSM2068667_out/GSM2068667.abundance.tsv, GSM2068668_out/GSM2068668.abundance.tsv, GSM2068669_out/GSM2068669.abundance.tsv, GSM2068670_out/GSM2068670.abundance.tsv, GSM2068671_out/GSM2068671.abundance.tsv, GSM2068672_out/GSM2068672.abundance.tsv, genome.gtf.gz
    output: se.rds
    jobid: 1
    resources: tmpdir=/tmp

```

[Sun May 29 15:54:10 2022]

```

localrule all:
    input: se.rds
    jobid: 0
    resources: tmpdir=/tmp

```

Job stats:

job	count	min threads	max threads
all	1	1	1
create_se	1	1	1
fastq_dump	30	1	1
get_index_files	1	1	1
kallisto	30	1	1
kallisto_index	1	1	1
prefetch	30	1	1
total	94	1	1

This was a dry-run (flag -n). The order of jobs does not reflect the order of execution.

(sra_kallisto_paired) [tobiass@master GSE78167]\$

Project — ~/Desktop/mount/work/project/tobias/GSE78167

Project

GSE78167

- .snakemake
- logs
- series
- slurm

config.yaml

```
1 runtable: SraRunTable.txt
2 transcriptsFTP: http://ftp.ensembl.org/pub/release-106/fasta/homo_sapiens/cdna/Homo_sapiens.GRCh38.cdna.all.fa.gz
3 gtfFTP: http://ftp.ensembl.org/pub/release-106/gtf/homo_sapiens/Homo_sapiens.GRCh38.106.gtf.gz
4 strandedness: reverse
5 organism: human
6 annotationVersion: 106
7
```

run.sh

Snakefile

sra_kallisto_paired.yaml

SraRunTable.txt

Kallisto_index 1 1 1
prefetch 30 1 1
total 94 1 1

This was a dry-run (flag -n). The order of jobs does not reflect the order of execution.

(sra_kallisto_paired) [tobiass@master GSE78167]\$

+ config.yaml 1:1 LF UTF-8 YAML GitHub

RNA_Seq_reference - RStudio

D1_project_overview.Rmd x

Source Visual B I Normal Format Insert Table Outline

```

1 ---  

2 title: "Project Overview"  

3 author: "tobiasst"  

4 date: `r format(Sys.time(), "%d %B, %Y")`  

5 output:  

6   html_document:  

7     df_print: paged  

8     code_folding: hide  

9     toc: true  

10 ---  

11  

12 ```{r libs}  

13 suppressPackageStartupMessages({  

14   library(SummarizedExperiment)  

15   library(RColorBrewer)  

16   library(ggplot2)  

17   library(ggrepel)  

18 })  

19 ...  

20  

21 **Series GSE78167**  

22  

23 An integrative transcriptomics approach identifies miR-503 as a candidate master regulator of the estrogen respon  

24  

25 Quantification of mRNAs in MCF7 cells responding to estrogen following a period of estrogen starvation. Three in  

26  

27 ```{r init_vars}  

28 cluster.dir <- "/Users/tobiasst/Desktop/mount/work/project/tobias/GSE78167/"  

29 ...  

30  

31 group variable  

32 ...  

33 ```{r}  

34 group_var <- "treated_with"  

35 group_var  

36 ...  

15:24 [1] "Chunk 1: libs"
  
```

Environment History Connections Tutorial

RNA_Seq_reference — Documents

Outline

Collected data from cluster.dir
 Read distribution
 Genes detected
 PCA
 sample-based
 PCA
 gene-based

R Global Environment

Data

se	Large SummarizedExperiment (260..)
Values	
batch_var	"Passage"
cluster.dir	"/Users/tobiasst/Desktop/mount/work/project/tobias/GSE78167"
group_var	"treated_with"

Files Plots Packages Help Viewer

com-apple-CloudDocs / Documents / RNA_Seq_reference

Name	Size	Modified
..		
.Rhistory	15.4 KB	May 29, 2018
01_project_overview.html	2.1 MB	May 28, 2018
01_project_overview.Rmd	3.9 KB	May 28, 2018
02_DGE.html	1 MB	May 6, 2018
02_DGE.Rmd	1.8 KB	May 9, 2018
03_EnrichR.html	5.3 MB	May 29, 2018
03_EnrichR.Rmd	2.3 KB	May 29, 2018
03_Heatmaps.html	1.1 MB	May 28, 2018
03_Heatmaps.Rmd	2 KB	May 28, 2018
03_MA_Vulcano.html	1.6 MB	May 28, 2018
03_MA_Vulcano.Rmd	2.1 KB	May 28, 2018
03_pathways_2h_alternative.Rmd	8.1 KB	May 29, 2018
03_pathways_2h.html	2.4 MB	May 28, 2018
03_pathways_2h.Rmd	7.2 KB	May 28, 2018
03_pathways_24h.html	3.3 MB	May 28, 2018

Console Terminal

R 4.2.0 - ~/Library/Mobile Documents/com-apple-CloudDocs/Documents/RNA_Seq_reference

```

> cluster.dir <- "/Users/tobiasst/Desktop/mount/work/project/tobias/GSE78167/"
> group_var <- "treated_with"
> group_var
[1] "treated_with"
> batch_var <- "Passage"
> batch_var
[1] "Passage"
>
  
```

PCA

sample based

colored for group

[Hide](#)

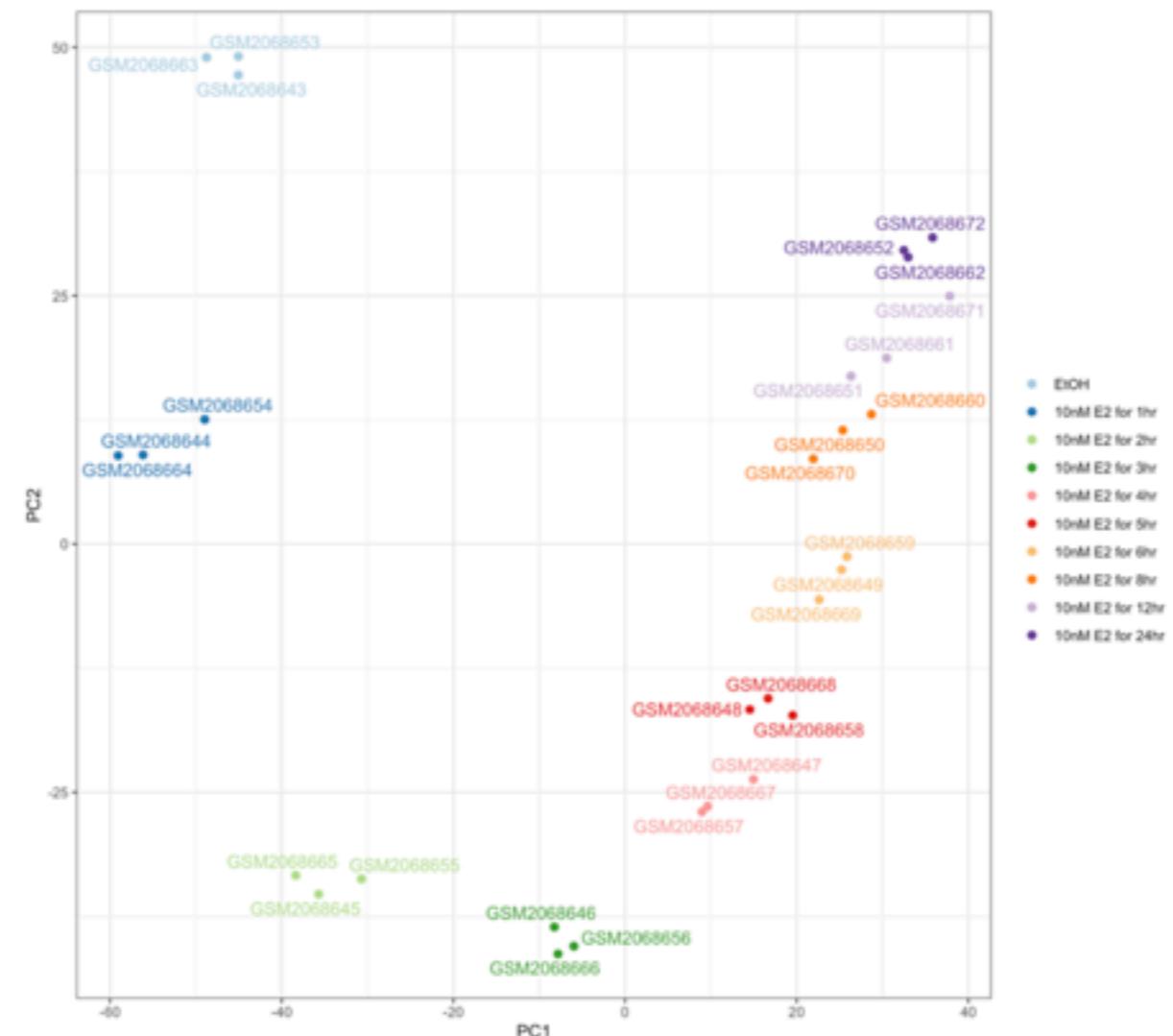
```
se <- se[!is.na(rowData(se)$entrez),]

mat <- log2(assays(se)$tpms+0.01)
rownames(mat) <- rowData(se)$symbol
mat <- mat[apply(mat,1,var)>0,]

pca <- prcomp(t(mat))

ggdf <- data.frame(pca$x[,1:2])
ggdf$class <- colData(se)[,group_var]
ggdf$sample_id <- colnames(se)

ggplot(ggdf, aes(x=PC1, y=PC2, col=class)) +
  geom_point(size=2) +
  geom_text_repel(aes_string(label = "sample_id"), show.legend = FALSE, max.overlaps = 10) +
  scale_color_manual(values = class.cols) +
  theme_bw() + theme(legend.title = element_blank())
```



immediate early

[Hide](#)

```
res <- results(dds, contrast=c("treated_with","10nM.E2.for.2hr","EtOH"))
summary(res)
```

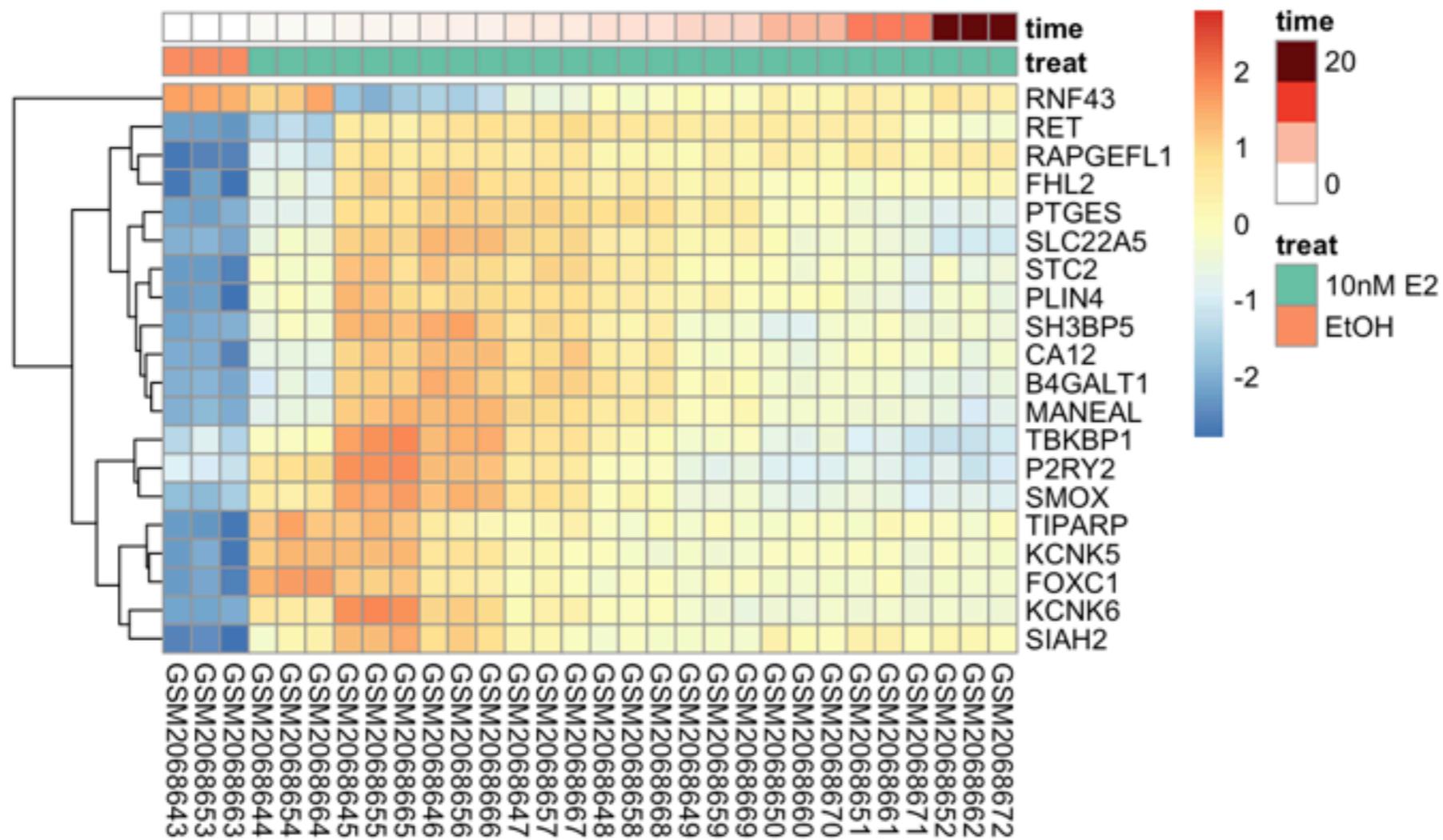
```
##
## out of 26098 with nonzero total read count
## adjusted p-value < 0.1
## LFC > 0 (up)      : 2502, 9.6%
## LFC < 0 (down)    : 2127, 8.2%
## outliers [1]       : 2, 0.0077%
## low counts [2]     : 10120, 39%
## (mean count < 4)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results
```

[Code](#)

2h

top 20 de-regulated genes

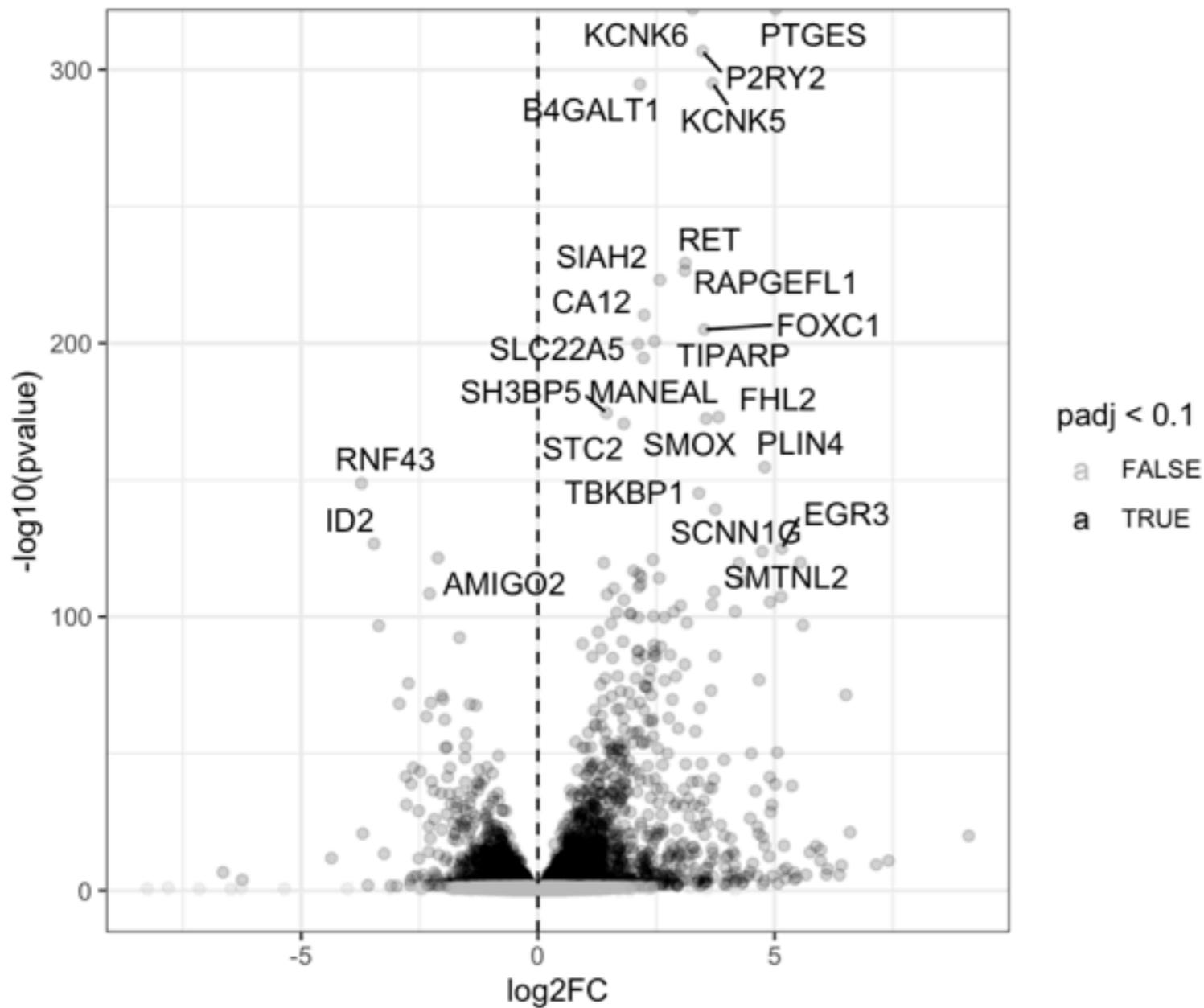
Code
Code



Vulcano

[Code](#)[Hide](#)

```
ggplot(ggdf, aes(x=log2FC, y=-log10(pvalue), label=gene, col=padj<0.1)) + geom_point(alpha=0.2) +  
  theme_bw() + geom_vline(xintercept = 0, linetype=2) +  
  geom_text_repel() +  
  scale_color_manual(values=c("grey","black"))
```



padj < 0.1

a FALSE

a TRUE

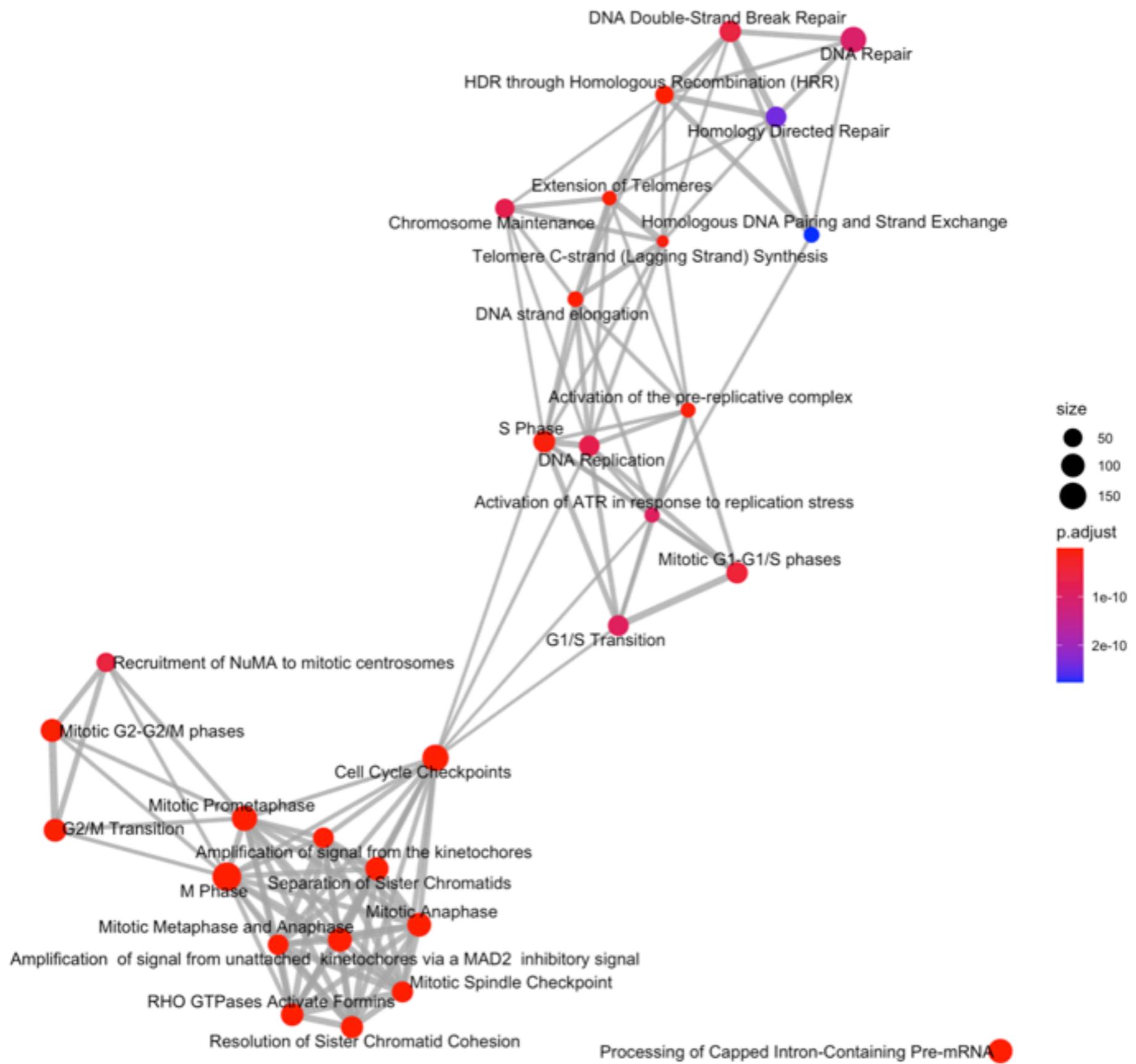
GSEA

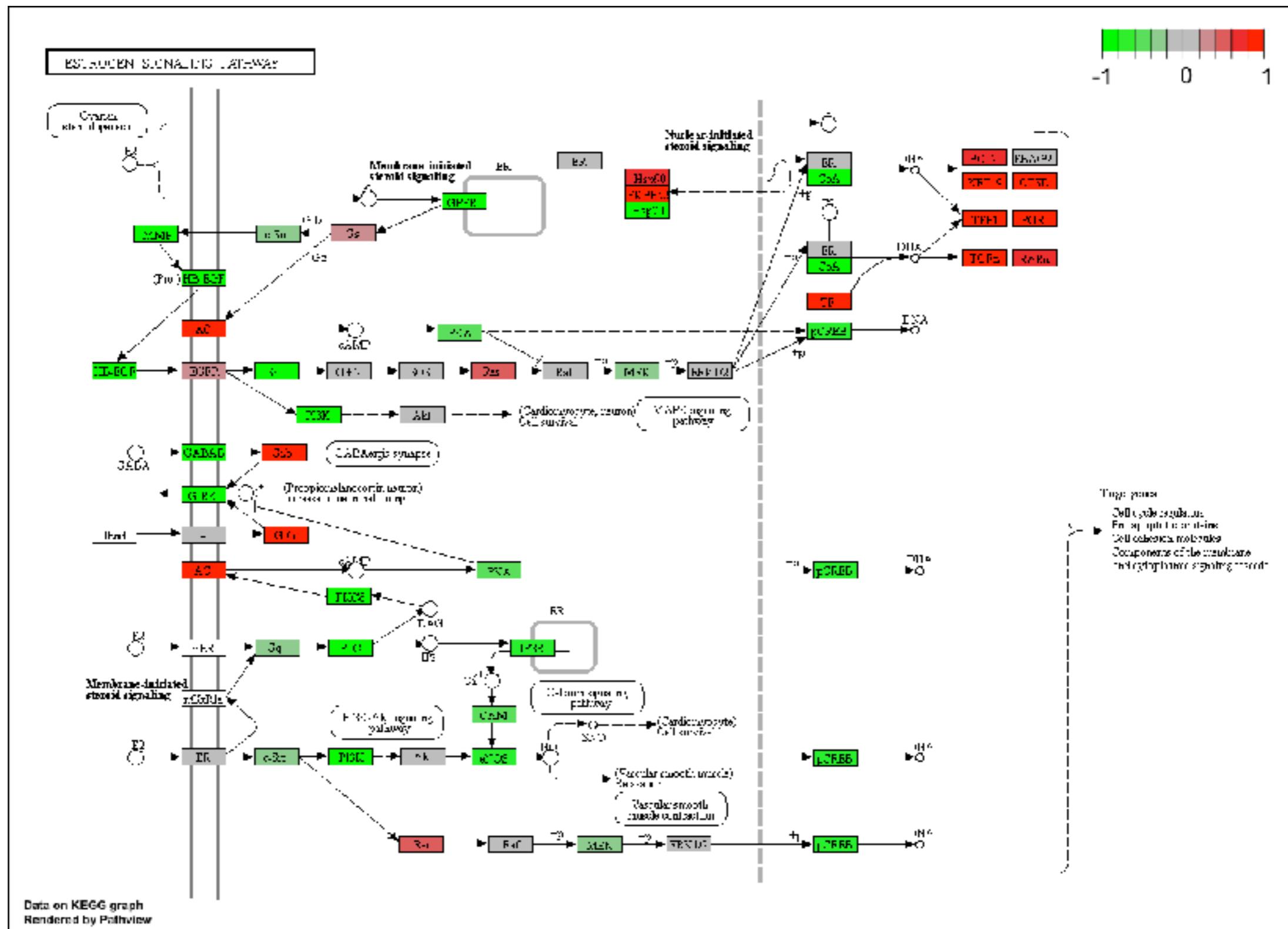
[Code](#)

Hallmark gene sets

[Code](#)

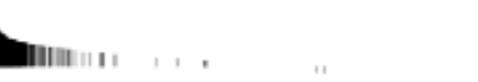
Pathway	Gene ranks	NES	pval	padj
HALLMARK_ESTROGEN_RESPONSE_EARLY		2.71	1.0e-10	2.5e-09
HALLMARK_ESTROGEN_RESPONSE_LATE		2.52	1.0e-10	2.5e-09
HALLMARK_UNFOLDED_PROTEIN_RESPONSE		1.81	1.5e-05	2.5e-04
HALLMARK_MYC_TARGETS_V2		1.76	1.4e-03	1.0e-02
HALLMARK_HYPOXIA		1.63	1.4e-04	1.7e-03
HALLMARK_WNT_BETA_CATENIN_SIGNALING		-1.27	1.2e-01	2.4e-01
HALLMARK_INTERFERON_GAMMA_RESPONSE		-1.28	1.3e-02	3.9e-02
HALLMARK_INTERFERON_ALPHA_RESPONSE		-1.40	4.3e-02	1.1e-01
HALLMARK_E2F_TARGETS		-1.42	3.5e-03	1.4e-02
HALLMARK_LIV_RESPONSE_DN		-1.48	6.4e-03	2.3e-02





curated gene sets

[Code](#)

Pathway	Gene ranks	NES	pval	padj
DUTERTRE_ESTRADIOL_RESPONSE_24HR_UP		3.13	1.0e-10	3.4e-09
DUTERTRE_ESTRADIOL_RESPONSE_6HR_UP		2.98	1.0e-10	3.4e-09
KOBAYASHI_EGFR_SIGNALING_24HR_DN		2.88	1.0e-10	3.4e-09
FLORIO_NEOCortex_BASAL_RADIAL_GLIA_DN		2.85	1.0e-10	3.4e-09
ROSTY_CERVICAL_CANCER_PROLIFERATION_CLUSTER		2.85	1.0e-10	3.4e-09
BHAT_ESR1_TARGETS_VIA_AKT1_DN		-2.62	1.0e-10	3.4e-09
BHAT_ESR1_TARGETS_NOT_VIA_AKT1_DN		-2.77	1.0e-10	3.4e-09
FRASOR_RESPONSE_TO_ESTRADIOL_DN		-2.79	1.0e-10	3.4e-09
DUTERTRE_ESTRADIOL_RESPONSE_6HR_DN		-3.05	1.0e-10	3.4e-09
DUTERTRE_ESTRADIOL_RESPONSE_24HR_DN		-3.10	1.0e-10	3.4e-09

0 4000 8000 12000 16000