

# Promoting Statistical Rigor and Reproducibility: Awareness, Action, and Advocacy

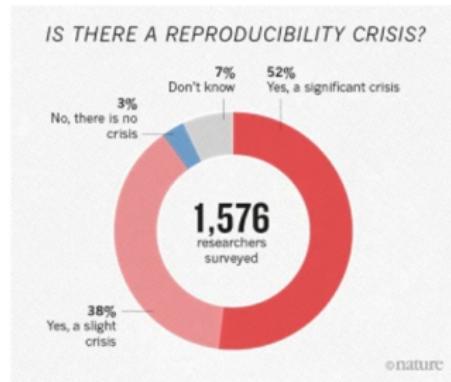
Zhaoxia Yu

Professor, Department of Statistics, University of California, Irvine

UCI, February 2025

# Reproducibility/Replication

- A cornerstone of scientific research
- An existing crisis is reproducibility/replication crisis
- Consequences of reproducibility crisis: wasting public resources, loosing public confidence, threatening people's health



[Nature](#) volume 533, pages452–454 (2016)

# Examples

Science

Current Issue First release papers Archive About Sub

HOME > SCIENCE > VOL. 346, NO. 6218 > FIXING PROBLEMS WITH CELL LINES

POLICY FORUM CELL BIOLOGY

## Fixing problems with cell lines

Technologies and policies can improve authentication

JON R. LORSCH<sup>1</sup>, FRANCIS S. COLLINS<sup>2</sup> AND JENNIFER LIPPINCOTT-Schwartz<sup>3</sup> Authors Info & Affiliations

SCIENCE • 19 Dec 2014 • Vol 346, Issue 6218 • pp. 1450-1453 • DOI: 10.1126/science.1258110

736 99 135



ILLUSTRATION: PETER AND MARIA HOEI; WWW.PETERHOEI.COM

Since the 1960's, more than 400 widely used cell lines worldwide have been shown to have been misidentified.

Studies using just two misidentified cell lines were included in three grants funded by NIH, two critical trials, 11 patents, and >100 papers.

Nature

## LETTERS

### A specific amyloid- $\beta$ protein assembly in the brain impairs memory

Sylvain Lamant<sup>1</sup>, Ming Tang Koh<sup>1</sup>, Linda Kettiliusk<sup>1</sup>, Rakez Kayed<sup>2</sup>, Charles G. Glabe<sup>2</sup>, Austin Yang<sup>2</sup>, Michael Gallagher<sup>2</sup> & Karen H. Ashe<sup>1,2,3\*</sup>

Memory function often declines with age<sup>1</sup>, and is believed to result from subtle changes in synaptic transmission or reduction in the number of synapses. Some memory loss is also associated with Alzheimer's disease with neurodegeneration. Here we use Tg2576 mice, which express a human amyloid- $\beta$  precursor protein (APP) variant linked to Alzheimer's disease, to show that memory loss due to the absence of neurofibrillary or amyloid-like protein assemblies. Young Tg2576 mice (6 months old) have memory deficits and older Tg2576 mice (18 months old) have more severe memory deficits without neuronal loss, and old aged (18 months old) Tg2576 mice show neuritic plaques and neurofibrillary tangles (Fig. 1). We found that memory deficit in middle-aged Tg2576 mice is caused by the truncated neurofibrillary protein assemblies, whereas the memory deficit in old Tg2576 mice disrupts memory when administered to young APP<sup>3E4</sup> (Aβ<sup>3E4</sup>). Aβ<sup>3E4</sup> protein from the brain-diseased Tg2576 mice disrupts memory when administered to young APP<sup>3E4</sup> mice, suggesting that the memory deficits of plaques or neurofibrils may contribute to cognitive deficits associated with Alzheimer's disease.

To further investigate the mechanism of memory loss, 15 years before diagnosis<sup>2</sup>, non-demented individuals at risk genetically for Alzheimer's disease show abnormalities in relational memory, which is a key feature of Alzheimer's disease. Alzheimer's disease has an insidious onset, which makes it particularly difficult to associate memory impairment with Alzheimer's disease<sup>3</sup>. Tg2576 mice exhibit memory loss and some of the histopathological features of Alzheimer's disease, such as neuritic plaques, neurofibrillary tangles, and inflammatory changes. However, Tg2576 mice lack neuritic filaments, suggesting that they have low and gross structure<sup>4</sup>. They may not be able to form the same type of memory as Alzheimer's disease, which depends on the degree of dementia or the onset of memory loss<sup>5</sup>.

In Tg2576 mice, at a lower APP transgenic value, there is strong evidence that truncated APP is unnecessary for age-related memory

for memory loss, as associated with either disease or memory function. Our solution to this conundrum was to use a mixture of soluble APP and truncated APP, which is a mixture of truncated APP (Aβ<sup>3E4</sup>) and soluble APP (Aβ<sup>3E4</sup>ΔC-terminal).

A challenge in analysis of Aβ<sup>3E4</sup>ΔC-terminal lies in the detection of the truncated protein because it is a complex intracellular molecule, membrane-bound, and insoluble. We overcome this obstacle by developing a high-fidelity extraction procedure that separates the truncated protein from the soluble protein (Fig. 1). Our extraction method allowed us to quantify and compare the relative proportion of truncated-derived soluble APP to total APP in the hippocampus of young, middle-aged, and old Tg2576 mice. To extract soluble APP, we used an extraction procedure to search for Aβ<sup>3E4</sup>ΔC-terminal molecules. We used antibodies that recognize Aβ<sup>3E4</sup>ΔC-terminal molecules to identify its effects. First, this truncated protein coincides with memory loss in 6 months. Second, this truncated protein coincides with memory loss in 18 months. Third, this truncated protein coincides with memory loss in APP<sup>3E4</sup> mice. Finally, this truncated protein coincides with memory loss in APP<sup>3E4</sup>ΔC-terminal mice. When we analyzed immunoblotting immunoreactive-depleted fractions extracts, we found a set of apparent assemblies of Aβ in the soluble, truncated, and full-length APP-immunoreactive proteins. In addition, to a faint 40-kDa band corresponding to Aβ monomers, 60-kDa- and 42-kDa-immunoreactive proteins (see Methods) were found in the truncated APP-immunoreactive proteins (10 kDa), hemispheric (27 kDa), hemispheric (40 kDa) and dendrite (56 kDa) Aβ<sub>42</sub>-immunoreactive bands. These species represent multiple forms of Aβ<sup>3E4</sup>ΔC-terminal, including the truncated Aβ<sup>3E4</sup>ΔC-terminal (40 kDa) appearing in mice older than 6 months. The detection of similar bands using 60-kDa and 42-kDa antibodies excludes the possibility that these bands are truncated Aβ<sup>3E4</sup>ΔC-terminal, which lacks the mid-domain (Ala<sub>40</sub>–Ala<sub>42</sub>) recognized by 42-kDa (Supplementary Fig. 1a). The bands were not recognized by 200 kDa (Supplementary Fig. 1a). The bands were neither Aβ<sup>3E4</sup> nor APP<sup>3E4</sup> degradation products (data not shown).

Although our results indicate that aging induces all tritons to bind to the high molecular weight Aβ<sup>3E4</sup>ΔC-terminal, we cannot exclude the possibility that this might account for all observed

Nature 440, 352–357 (2006) | Cite this article

76k Accesses | 2377 Citations | 1823 Altmetric | Metrics

This article was retracted on 24 June 2024

This article has been updated

# Factors Contributing to Reproducibility Crisis

- Publication bias
- Hyper-competitive research environment
- Lack of adherence to good scientific practices
- Bad statistical practice

# Factors Contributing to Reproducibility Crisis

- Bad statistical practice
  - underdeveloped experimental designs
  - over-interpretation of statistically marginal differences
  - lack of transparent presentation of methods
  - misuse of statistical methods

The screenshot shows a navigation bar with links for 'NEW TO NIH', 'FUNDING', 'GRANTS PROCESS', 'POLICY & COMPLIANCE', and 'NEWS & EVENTS'. Below the navigation is a 'ABOUT US' link. The main content area has a breadcrumb trail: 'Home > Policy & Compliance > Policy Topics > Enhancing Reproducibility through Rigor and Transparency > Principles and Guidelines for Reporting Preclinical Research'. A 'Policy & Compliance' section is visible. The main title is 'Principles and Guidelines for Reporting Preclinical Research'. Below the title, it says: 'NIH held a joint workshop in June 2014 with the Nature Publishing Group and Science or the issue of reproducibility and rigor of research findings, with journal editors represent over 30 basic/preclinical science journals in which NIH-funded investigators have most often published. The workshop focused on identifying the common opportunities in the scientific publishing arena to enhance rigor and further support research that is reproducible, robust, and transparent.' A 'Rigorous Statistical Analysis' section is present with the note: 'A section outlining the journal's policies for statistical analysis should be included in the information for Authors, and the journal should have a mechanism to check the statistical accuracy of submissions.' A 'Transparency in reporting' section is also mentioned.

# Awareness

Statistic analysis for a large number of cells from a few animals [Inbox x](#)

Xu, Xiangmin <xiangmix@hs.uci.edu>  
to me, Lujia, Haiz2 ▾

Thu, Mar 26, 2020, 4:37 PM

◀ Messages BiologyCollabor... Details

Hello! My data include  
1000 cells from 10  
animals. How should I  
analyze them?

The data are clustered.  
LME should work

What is LME? We often  
use t-test, ANOVA, and  
Wilcoxon when the data  
are not normal

You mean analyzing the  
1000 data points without  
using the animal IDs?

Yes. This is the common  
practice

Are you serious ...

◀ Messages BiologyCollabor... Details

Are you serious ...

I can show you how  
wrong that way is

Let's submit a letter to a  
journal

many weeks' teamwork ...

It is too long for a letter.  
Let's make it a primer



iFakeTextMessage.com



iFakeTextMessage.com



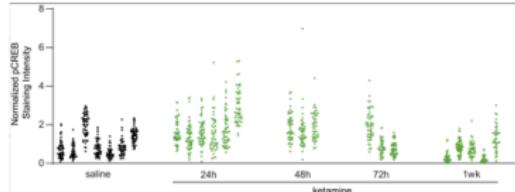
# Better Statistical Practice?

- How is the situation compared to ten years ago?
- Aarts et al. (2014 Nature Neuroscience) examined 314 articles published Jan 2012 – June 2013 in five top journals: Science, Nature, Cell, Nature Neuroscience, Neuron: 53% involved nested (clustered) data but “conventional” analysis methods were used in all the articles.
- We surveyed articles published in prestigious journals over the past few years. Among >100 articles in which multiple observations were measured from each individual animal, <50% accounted for data dependencies in any meaningful way.

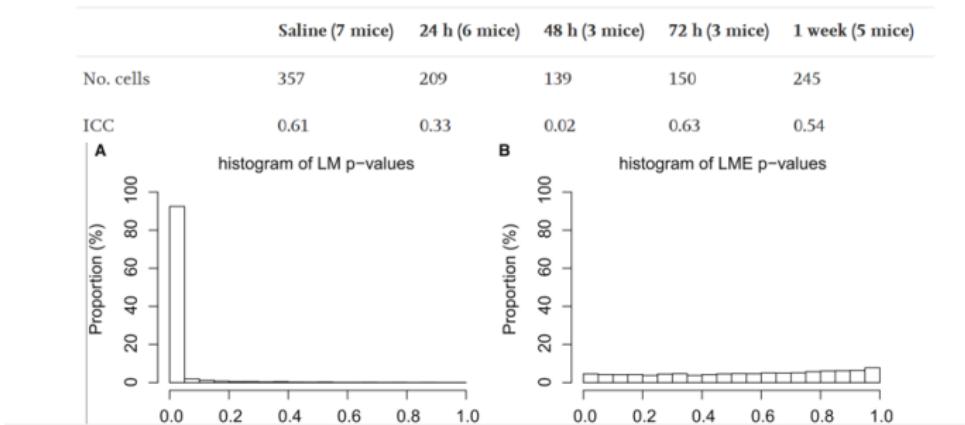
# Examples of Wrong Analyses

- “ $t(28656) = 314$  with  $p < 10^{-10}$  over a total of  $n=28657$  neurons pooled across six mice”
- “ $n = 377$  neurons from four mice, two-sided Wilcoxon signed rank test”
- “610 A cells, 987 B cells and 2584 C cells from 10 mice, one-way ANOVA and Kruskal–Wallis test”
- “two-sided paired t test,  $n=1597$  neurons from 11 animals, d.f. = 1596”

# Attention Please!



1200 neurons/measurements from 24 mice

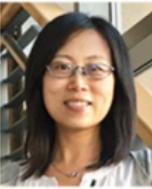


# Our First Action



Steven Grieco

Postdoctoral Scholar



Lujia Chen  
PhD Candidate

UCI's multidisciplinary team of researchers included (from left) Xiangmin Xu, PhD, professor and Chancellor's Fellow in the Department of Anatomy and Neurobiology, and director for the Center for Neural Circuit Mapping at the UCI School of Medicine; Zhaoxia Yu, PhD, and Michele Guindani, PhD, professors from the Department of Statistics UCI Donald Bren School of Information & Computer Sciences; and Todd Holmes, PhD, professor and vice chair of the UCI School of Medicine Department of Physiology & Biophysics.

UCI School of Medicine



PlumX Metrics



Beyond t test and ANOVA: applications of mixed-effects models for more rigorous statistical analysis in neuroscience research

Journal of Neuroscience, 2020, 40(16), 3603–3613, Issue 16, April 21, 2020

Publication Date: April 21, 2020



James (Jim) Dent, LSSBB, DTM on LinkedIn



Dr Jennifer Deem  
@Deemdeemlab · Follow

I have a love/hate relationship with statistics. This article increases the love side of the relationship.  
[cell.com/neuron/fulltext/10.1016/j.jneurosci.2020.03.038](http://cell.com/neuron/fulltext/10.1016/j.jneurosci.2020.03.038)

11:41 AM · Nov 16, 2021

Metric Details

CITATIONS: 16

CHINESE INDEXES: 36

SCOPUS: 36

CROSSREF: 2

CAPTURES: 368

MENTIONS: 1

MEET MEETERS: 1

MENTIONED BY: 1

SOCIAL MEDIA: 301

TWEETS: 267

SHARING: 267

REPLIES: 34

Most Recent Tweet

Review Description

In basic neuroscience research, data are often clustered or collected with repeated measures, rarely correlated. The most common statistical approach is to use t tests and ANOVAs to detect dependence into account and thus are often rejected. This Primer introduces linear and generalized mixed-effects models that enable researchers to analyze such data. It also provides new insight into recognizing when they are needed and how to apply them. The appropriate use of mixed effects models will help researchers improve their statistical design and will lead to data analysis with greater validity and higher reproducibility of the experimental findings.

Bibliographic Details

DOI: 10.1016/j.jneurosci.2020.03.038

PMID: 32547008

URL ID: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7694271/

DOI: 10.1016/j.jneurosci.2020.03.038

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7694271/

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7694271/

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7694271/

Provide Feedback

How ideas for a new metric? Would you like to see something else here? Let us know!



RRR

Awareness

Zhaoxia Yu (UCI)

Action 1

Action 2

Advocacy

Discussion

Appendix

Promoting Statistical Rigor and Reproducibility

UCI, February 2025

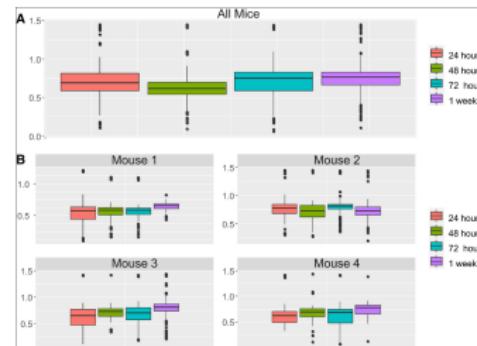
10 / 58

# Accessible Materials

- Start from the methods that everyone understands
- T-test  $\Rightarrow$  ANOVA  $\Rightarrow$  LM  $\Rightarrow$  LME  $\Rightarrow$  GLMM  $\Rightarrow$  Bayesian
- Representative examples and data
- Detailed supplemental materials
- Step-by-step guidance and reproducible code (including fixed random seeds)

# Representative Examples

- Example 1: unacceptable type I error rate due to ignoring data dependency.
- Example 2: when cells are pooled naively, the results can be dominated by the animal contributing to a large proportion of cells.
- Example 3: nested random effects. P-values became more significant after accounting for data dependency.
- Example 4: non-normal outcomes.



# Accessible Resources

The screenshot shows a web browser with two tabs open. The left tab is titled "A Introduction of Linear Mixed-Effects and Generalized Linear Mixed-effects Models" and contains R code and its output. The right tab is titled "Mixed-effects model and its applications" and displays the "Supplemental data download" section of a paper by Zhaoxia Yu et al. in Neuron.

**A Introduction of Linear Mixed-Effects and Generalized Linear Mixed-effects Models**

R code and output:

```
library(lme4)
icc.analysis<-lmer(Reaction ~ 1 | Subject + Condition, data = sleepstudy)
icc.analysis
```

Output:

```
Run Lmer (Est+Error)
Run lmer has not been executed in this session
executed at unknown time
"Saline", "24h"
icc.analysis
```

	n	icc	design_effect	ef
	cdbs	cdbs	cdbs	cdbs
Saline	7	0.62094868	32.047454	
24h	6	0.33006327	17.668195	
48h	3	0.01780304	1.807071	
72h	3	0.62810904	31.777343	
1wk	5	0.53694579	26.773598	

Essential results are summarized in the following table

	stats	Saline (7 mice)	24h (6 mice)
# of cells	357	209	
ICC	0.61	0.33	

The ICC indicates that the dependency due to cluster not be treated as 1,200 independent cells. When dep rate can be much higher than the pre-chosen level of examine the false positives based on the dependency below to generate 1,000 data sets, each of which will follow the five conditions. Surprisingly, the type I error observations is over 90% at the significance level of 0.05.

```
[1] source("http://xulab.anat.ucl.edu/Downloads/Type_I_error_rate_of_LMM_at_significance_level_0.05.R")
[1] "Type_I_error_rate_of_LMM_at_significance_level_0.05.R"
```

**Mixed-effects model and its applications**

Application of Mixed-effects Model in...

1 Supplemental data download

2 Introduction

3 Mixed-effects model analysis

4 Conduct LME in R: Example 1

4.1 Wrong analysis with lm or anova

4.2 LME: estimation methods

4.3 On the degrees of freedom an...

4.4 The overall p-value for the tre...

4.5 P-value adjustment for multip...  
4.6 Robust methods with paramet...

4.7 Additional tools

5 Conduct LME in R: Example 2

5.1 Wrong analysis

5.2 LME

5.3 Why pooling data naively is n...

5.4 Remark: on the minimum num...

6 Conduct LME in R: Example 3

6.1 Wrong analysis with lm, pare...

6.2 A note on "nested" random eff...

6.3 On models with more random ...

6.4 A note on the testing of random ...

7 Generalized Linear Mixed-Effects ...

7.1 Fit a GLMM in R

7.2 Use nonparametric methods t...

8 A Bayesian Analysis of Example 4

Zhaoxia Yu,<sup>1,2\*</sup> Michele Guidani,<sup>3</sup> Steven F. Greico,<sup>2</sup> Luja Chen,<sup>2</sup> Todd C. Holmes,<sup>2</sup> Xiangmin Xu<sup>1,2,4,5</sup>

This is the online supplementary to accompany Yu Z., Guidani M., Greico SF., Chen L., Holmes TC., Xu X. (2022) Beyond t-Test and ANOVA: applications of mixed-effects models for more rigorous statistical analysis in neuroscience research. *Neuron*. 110: 21-23. <https://doi.org/10.1016/j.neuron.2021.10.030>.

The data used in this supplementary can be downloaded from the following website: <https://www.wics.ucl.edu/~zhaoxia/Data/BeyondTandANOVA/>.

```
knitr:::include_graphics("Fig8.png")
```

## Chapter 1 Supplemental data download

This is the online supplementary to accompany Yu Z., Guidani M., Greico SF., Chen L., Holmes TC., Xu X. (2022) Beyond t-Test and ANOVA: applications of mixed-effects models for more rigorous statistical analysis in neuroscience research. *Neuron*. 110: 21-23. <https://doi.org/10.1016/j.neuron.2021.10.030>.

The data used in this supplementary can be downloaded from the following website:

<https://www.wics.ucl.edu/~zhaoxia/Data/BeyondTandANOVA/>.

## Neuron

CelPress

### Beyond t test and ANOVA: applications of mixed-effects models for more rigorous statistical analysis in neuroscience research

Zhaoxia Yu,<sup>1,2\*</sup> Michele Guidani,<sup>3</sup> Steven F. Greico,<sup>2</sup> Luja Chen,<sup>2</sup> Todd C. Holmes,<sup>2</sup> and Xiangmin Xu<sup>1,2,4,5</sup>

In basic neuroscience research, data are often clustered or collected with repeated measures, hence correlation. The most widely used measures such as t-test and ANOVA do not take data dependence into account and are often inappropriate. This Primer provides an introduction to mixed-effects models that consider data dependence and provides clear instruction on how to recognize when they are needed and how to apply them. The appropriate use of mixed-effects models will help researchers improve their experimental design and will lead to data analyses with greater validity and higher reproducibility of the experimental findings.

# Deliver Messages

- “Rigorous statistical analysis is not a hunt for the smallest p value (commonly known as p-hacking or significance chasing).”
- Overuse of p-values: “It should be noted that the modeling decisions should not be based on tests and p values alone, as the result may be significant even . . . or be insignificant . . . . Rather, the modeling decision should always be guided by the combined information provided by the study design, scientific reasoning, and previous evidence.”
- “Due to the limited space, it is overambitious to cover all of the practical issues . . . We refer the interested reader to specialized research articles ( . . . ) or to consult with experienced statisticians.”

# The Misuse of A Robust Test



## Confronting false discoveries in single-cell differential expression

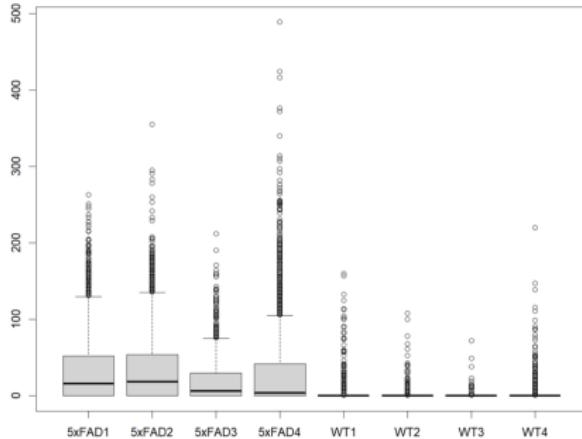
[Jordan W. Squair](#), [Matthieu Gautier](#), [Claudia Kathe](#), [Mark A. Anderson](#), [Nicholas D. James](#), [Thomas H. Hutson](#), [Rémi Hudelle](#), [Taha Qaiser](#), [Kaya J. E. Matson](#), [Quentin Barraud](#), [Ariel J. Levine](#), [Gioele La Manno](#), [Michael A. Skinnider](#)✉ & [Grégoire Courtine](#)✉

*Nature Communications* **12**, Article number: 5692 (2021) | [Cite this article](#)

**127k** Accesses | **197** Altmetric | [Metrics](#)

endothelial cells using representative single-cell and pseudobulk methods. We selected the Wilcoxon rank-sum test as a single-cell method, since this test has been the most widely used approach in the field of single-cell transcriptomics (Fig. 1b), and edgeR-LRT<sup>29</sup> as a pseudobulk

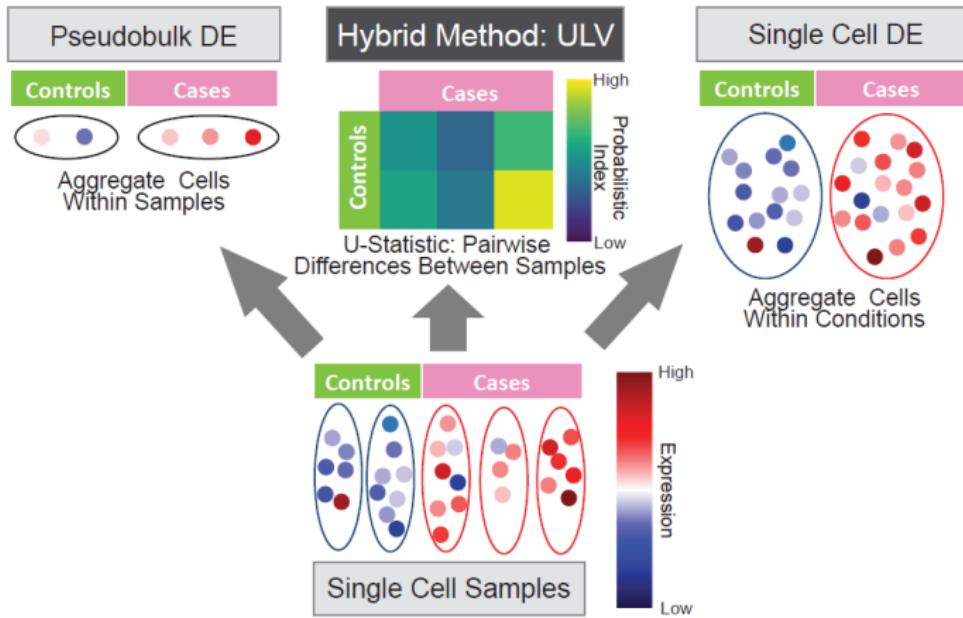
# Why Nonparametric Tests are Appealing?



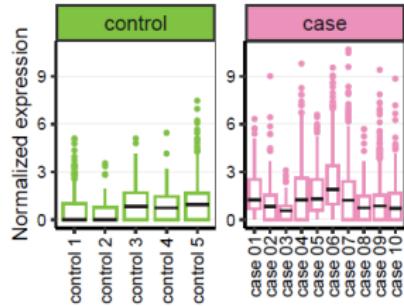
- Non-normal distributions, outliers, excess zeros
- Challenges: small sample sizes, large repeated measurements

# Existing and Proposed Approaches

<https://github.com/yu-zhaoxia/ULV>



# Goal: Robust, Valid, and Flexible



Subject Level Covariate Information

Subject ID	Age	Gender	...
case 1	63	Male	...
case 2	80	Female	...
...	...	...	...
control 1	31	Female	...
control 2	45	Male	...
...	...	...	...

#case m = 10  
#control n = 5

## Covariates

# Existing Approaches

- Parametric methods: MAST-GLMM with zero-inflated Poisson or negative binomial, pseudobulk DeSeq2.
- Wilcoxon rank-sum test for cell data is the most widely used (Squair et al. 2021). But ...
- Wilcoxon rank-sum test for clustered data
  - use cluster-level summaries: when  $m=n=4$ , the smallest p-value is  $1/70 \approx 0.014$
  - sample one observation per cluster at a time (Datta–Satten), large-sample approximations (Rosner et al.)
    - Large-sample sizes are required.
    - The number of repeated measurements per cluster is often small.
    - Better suited for longitudinal studies.
    - Not easy to extend to multiple groups and adjust for covariates

# U Statistic and Rank Sum for Two Samples

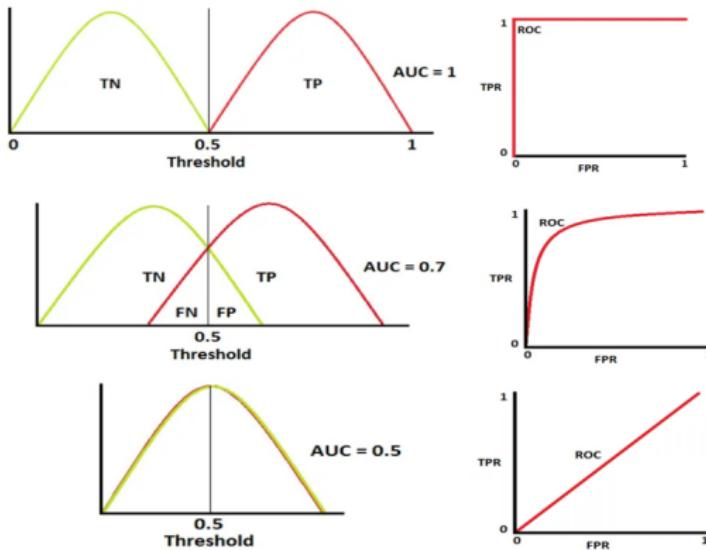
- Consider  $X_1, \dots, X_m$  and  $Y_1, \dots, Y_n$ .
- Mann-Whitney U test is based on

$$U = \sum_{i=1}^m \sum_{j=1}^n I(X_i > Y_j)$$

- Wilcoxon rank sum test is based on  $R_X$ , the sum of ranks in group  $X$ .
- Mann-Whitney U and Wilcoxon rank sum are equivalent because

$$U = R_X - \frac{m(m+1)}{2}$$

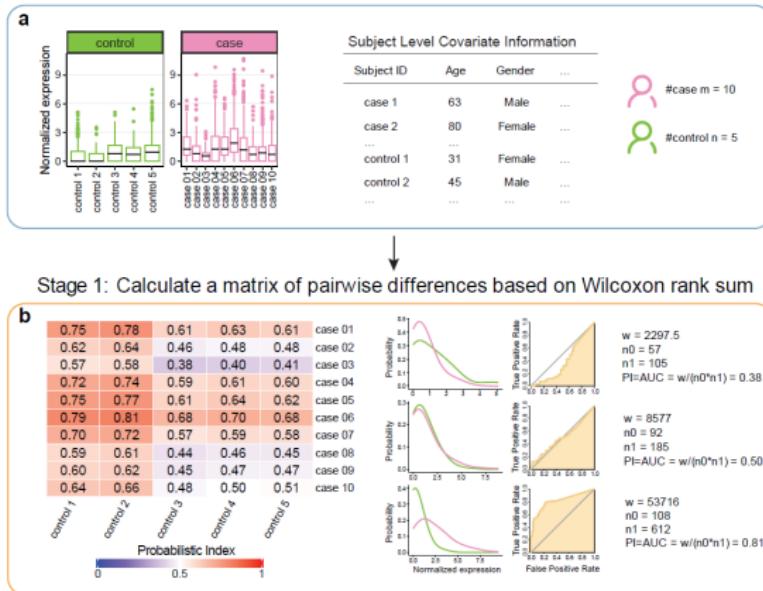
# AUC as A Metric for Prediction / Separation



- $AUC = \frac{U}{mn}$

# Quantify the Difference between Two Subjects

- AUC is also known as the Probabilistic Index (PI)



## Two-Group Comparisons: Notations

- Consider  $m$  cases and  $n$  controls, each with multiple measurements.
- The  $i$ th case is denoted by

$$\vec{y}_{1i} = (y_{1,i,1}, \dots, y_{1,i,K_{1i}})^T,$$

where  $K_{1i}$  is the number repeated measurements.

- Similarly, the  $j$ th control is denoted by

$$\vec{y}_{0j} = (y_{0,j,1}, \dots, y_{0,j,K_{0j}})^T$$

- The comparison between the  $i$ th case and  $j$ th control is given by  $h(\vec{y}_{1i}, \vec{y}_{0j})$ , which is a U statistic.

## Two-Group Comparisons

- For example, one robust choice of  $h$  is

$$h(\vec{y}_{1i}, \vec{y}_{0j}) = \frac{1}{K_{1i} K_{0j}} \sum_{k=1}^{K_{1i}} \sum_{k'=1}^{K_{0j}} \left[ I(y_{1ik} > y_{0jk'}) + \frac{1}{2} I(y_{1ik} = y_{0jk'}) \right].$$

- It is a consistent, unbiased, and minimum variance estimator of the probabilistic index for quantifying "greater".
- A reasonable statistic is  $U_{cluster} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n h(\vec{y}_{1i}, \vec{y}_{0j})$
- $U_{cluster}$  can be used for inference. However, its large-sample distribution requires estimating correlations, so sample sizes need to be large.

# A Latent Variable Model

- Let  $d_{ij} = h(\vec{y}_{1i}, \vec{y}_{0j})$ ,  $i = 1, \dots, m; j = 1, \dots, n$
- We model  $d_{ij}$  by

$$d_{ij} = a_i - b_j + \epsilon_{ij}, i = 1, \dots, m; j = 1, \dots, n$$
$$a_i \sim N(\mu, \sigma_1^2), b_j \sim N(0, \sigma_0^2), \epsilon_{ij} \sim N(0, \sigma^2)$$

- $a_i$  represents the latent (relative and not directly observed) level of the  $i$ th case
- $b_j$  is the latent level of the  $j$ th control
- $\mu$  is the difference in the population levels of the two groups, which is the parameter of interest.

$$H_0 : \mu = 0.5$$

# Justification of the Latent Variable Model

- The model includes the two-sample t-test as a special case when  $d(i,j) = \bar{x}_i - \bar{y}_j$ .
- Connection to Rasch (1 parameter Item Response Model) Model
  - Proposed by Rasch in 1960 for modeling responses of test or survey items.
  - The Rasch model and extensions are widely used in
    - psychometrics (van der Linden & Hambleton, 2013)
    - educational tests (Lan et al., 2016)
    - public health (Cappelleri et al., 2014)
    - market and financial research (Schellhorn & Sharma, 2013; Brzezinska, 2016).

# Justification of the Latent Variable Model

- Let  $Y_{ij} \in \{0, 1\}$  denote test-taker i's response to item j, with 1 for a correct response and 0 otherwise.
- The Rasch model assumes that

$$\log \left( \frac{\Pr(Y_{ih} = 1)}{\Pr(Y_{ih} = 0)} \right) = \theta_i - \beta_j$$

- The original Rasch model assumes that the ability scores ( $\theta_i$ ) are random and the item difficulties ( $\beta_j$ ) are fixed.

# Extensions

- Adjust for covariates
- Compare multiple groups
- Weighted analysis to account for unbalanced cluster sizes

# Adjust for Covariates

- Suppose there are  $p$  covariates for each subject,  $x_i = (x_{i1}, \dots, x_{ip})^T$ ,

$$\begin{aligned}d_{ij} &= (a_{0i} + \beta^T x_i) - (b_{0j} + \beta^T x_j) + \epsilon_{ij} \\&= a_{0i} - b_{0j} + \beta^T (x_i - x_j) + \epsilon_{ij},\end{aligned}$$

- $a_{0i}$  and  $b_{0j}$  represent the covariate-adjusted latent expression levels.

$$a_{0i} \sim N(\mu, \sigma_1^2), b_{0j} \sim N(0, \sigma_0^2), \epsilon_{ij} \sim N(0, \sigma^2)$$

- Testing  $H_0 : \mu = 0.5$ .

# Compare Multiple Groups

- Suppose there are a total number of  $M + 1$  conditions.
- The pairwise difference between the observations from the  $i$ th subject under condition  $m$  and the  $j$ th subject in the reference group is

$$d_{ij}^m = a_{mi} - b_{0j} + \epsilon_{ij}, i = 1, \dots, N_m; j = 1, \dots, N_0; m = 1, \dots, M,$$

- $a_{mi}$  represents the latent level of the  $i$ th subject in condition  $m$ , and  $b_{0j}$  is the latent level of the  $j$ th subject in the reference group.

$$a_{mi} \sim N(\mu_m, \sigma_m^2), b_{0j} \sim N(0, \sigma_0^2), \epsilon_{ij} \sim N(0, \sigma_\epsilon^2).$$

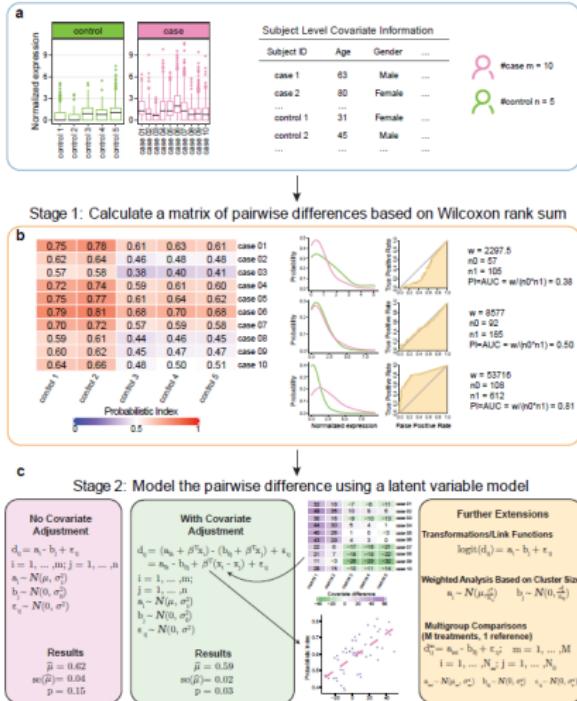
- Testing  $H_0 : \mu_1 = \dots = \mu_M = 0.5$ .

# Weighted Analysis

- Let  $K_{0i}$  and  $K_{1j}$  denote the cluster sizes of control  $i$  and patient  $j$ , respectively.
- $d_{ij} \propto (1/K_{0i} + 1/K_{1j})$
- Weighted latent variable model:

$$a_i \sim N(\mu, \frac{\sigma_1^2}{K_{1i}}), b_j \sim N(0, \frac{\sigma_0^2}{K_{0j}}).$$

# A Schematic Summary of ULV



# Simulation Results

- We simulated single-cell expression data based on parameters estimated from real sc RNA-seq data
- ULV has excellent control of type I error rates at different significance levels
- ULV has comparable power with competing parametric or pseudobulk methods
- ULV is much more powerful than the pseudobulk version of Wilcoxon rank sum test

# Application 1: ULV discovers DE proteins in AML patients

nature cancer

Article

<https://doi.org/10.1038/s43018-022-00480-0>

## An inflammatory state remodels the immune microenvironment and improves risk stratification in acute myeloid leukemia

Received: 4 February 2022

Audrey Lasry<sup>1,2\*</sup>, Bettina Nadorm<sup>1,2,3,4</sup>, Maarten Fornesrod<sup>5</sup>, Deedra Nicoll<sup>1,6</sup>,

Accepted: 4 November 2022

Huiyuan Wu<sup>7</sup>, Christopher J. Walker<sup>4,8</sup>, Zhengxi Sun<sup>1,2</sup>, Matthew T. Witkowski<sup>1,2</sup>,

Published online: 29 December 2022

Anastassia N. Tikhonova<sup>1,2</sup>, Maria Guittamont-Rusano<sup>1,2</sup>, Geraldine Cayanan<sup>1,2</sup>,

 Check for updates

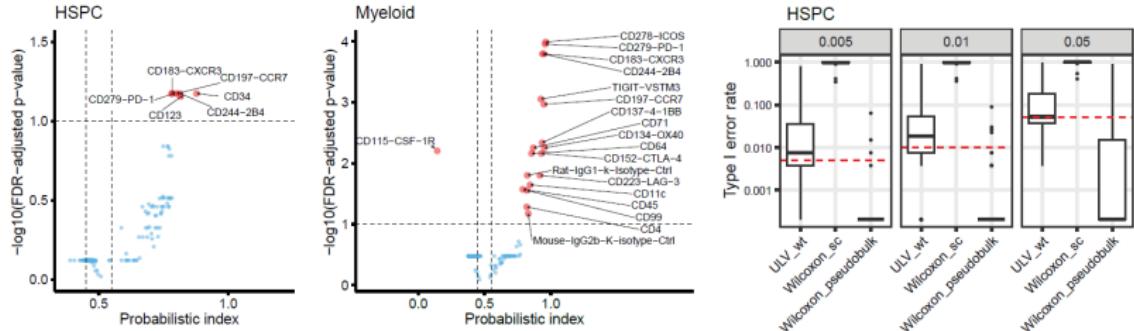
Anna Yeston<sup>1,2</sup>, Gabriel Robbins<sup>1,2</sup>, Esther A. Obeng<sup>1,2</sup>,

Aristoteles Tsirigos<sup>1,2</sup>, Richard M. Stone<sup>9</sup>, John C. Byrd<sup>10</sup>, Stanley Pounds<sup>1,2</sup>,

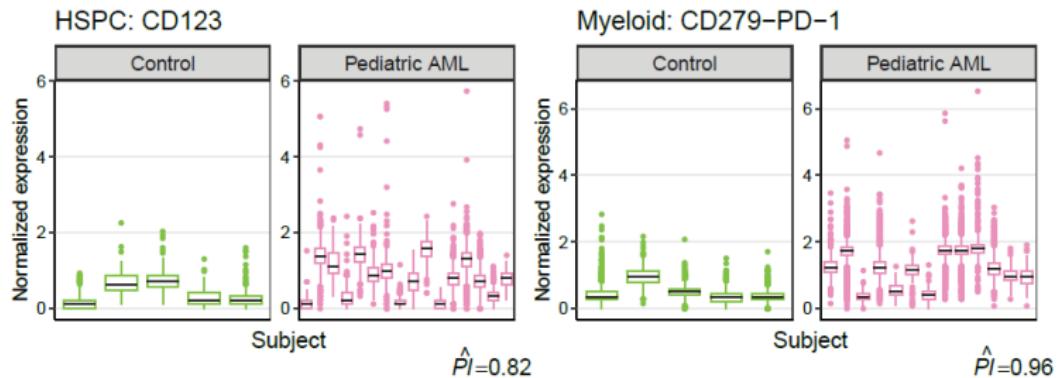
William L. Carroll<sup>1,2</sup>, Tanja A. Gruber<sup>1,2,11</sup>, Ann-Kathrin Eisfeld<sup>1,2,10,12</sup>, &

Iannis Alifantis<sup>1,2,13</sup>

- We focused on two cell types:
  - 56,764 malignant and healthy hematopoietic stem and progenitor cells (HSPCs): 269 proteins, 6 pediatric AML patients and 5 age-matched controls
  - 65,510 myeloid cells: 270 proteins, 13 pediatric AML patients and 5 age-matched controls
- Protein expression data were pre-processed using centered log-ratio normalization

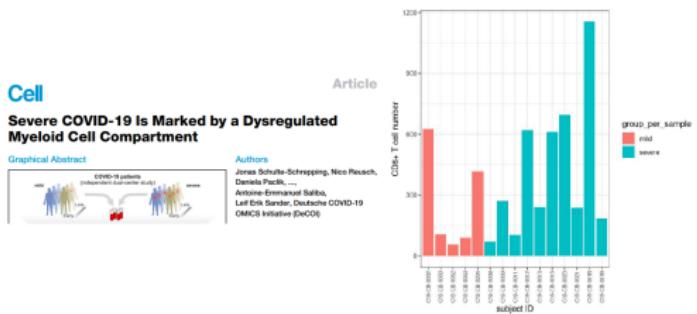


- Pseudobulk Wilcoxon didn't find any DE protein

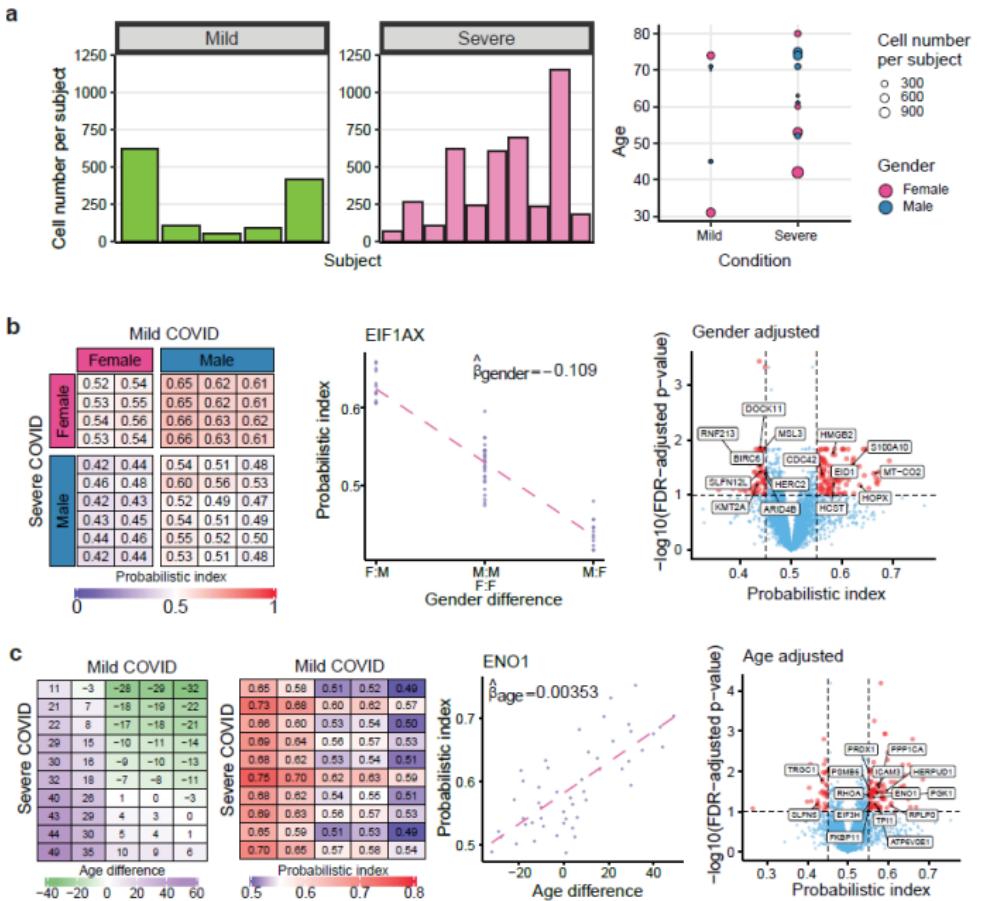


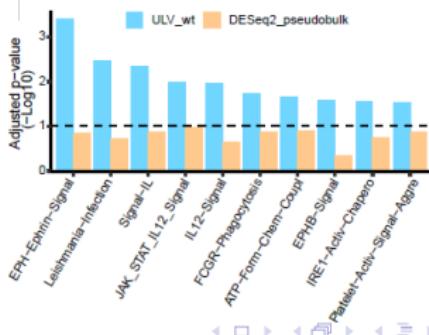
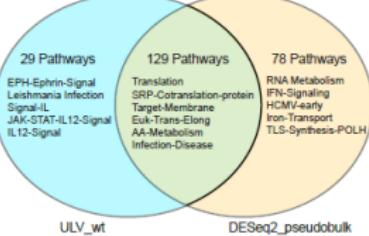
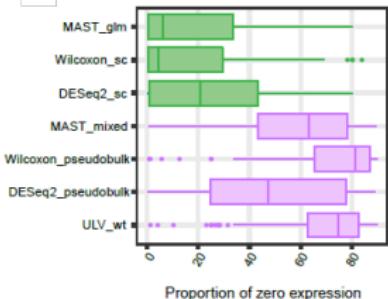
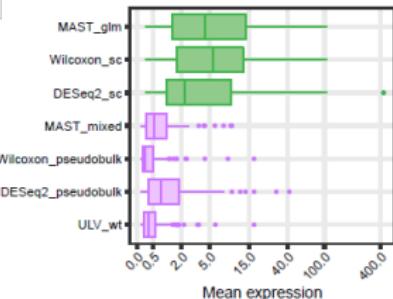
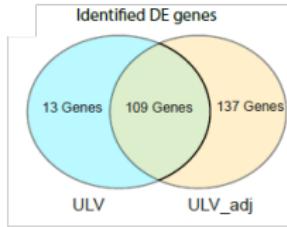
- The interleukin 3 receptor alpha (CD123) was previously found as a putative marker in AML blasts compared with normal HSCs
- The programmed cell death protein 1 (CD279-PD-1) has been linked to the suppression of immune responses during tumor development and is also associated with AML pathology

# Application 2: COVID-19 Data



- Blood samples from Germany COVID-19 patients.
- scRNA-seq data were measured using a droplet-based single-cell platform (10x Chromium), which contains the peripheral blood mononuclear cells (PBMC).
- We focused on patients with an adequate number of CD8+ T cells: 5 mild and 10 severe patients.





# Misuse of Statistical Methods

- It is easy and sometimes attempting to misuse methods
- Can happen in many journals, including the most prestigious journals
- High impact journals have an enormous amount readers
- High impact journals should be more responsible to set good examples

# Misuse of Statistical Methods

- Used the conventional rank sum test as a “robust” method
- No control of multiple comparisons:
  - reported 107 out 3350 genes were “significant” (nominal  $p < 0.05$ , fold change  $> 1.75$ )
  - But  $5\% \times 3350 = 167$

nature

Explore content ▾ About the journal ▾ Publish with us ▾

nature > articles > article

Article | Open access | Published: 07 February 2024

## Spatial transcriptomics reveal neuron–astrocyte synergy in long-term memory

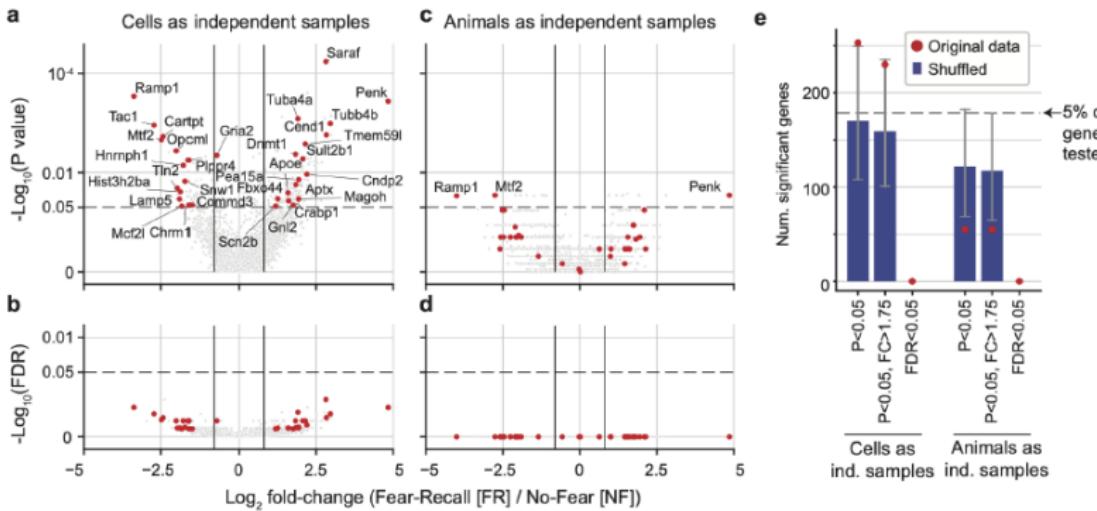
Wenfei Sun, Zhihui Liu, Xian Jiang, Michelle B. Chen, Hua Dong, Jonathan Liu, Thomas C. Südhof & Stephen R. Quake

Nature 627, 374–381 (2024) | Cite this article

56k Accesses | 222 Altmetric | Metrics

# Misuse of Statistical Methods

- Failing to adjust for multiple comparisons or correlated data leads to spurious DE genes



# Discussion

- Misuse of statistical methods has contributed to the reproducibility crisis.
- Raise awareness, take action, and continue advocacy efforts.
- Collaborate with domain scientists from diverse backgrounds.
- Provide accessible materials and resources
- Understand the challenges and develop new methods and tools

# Acknowledgements

- Student: Mingyu Du
- Collaborators:
  - Veronica Berrocal
  - Michele Guindani
  - Kevin Johnston (Utah Tech University), Steven F Grieco, Todd C Holmes, Wei Li, and Xingmin Xu, Eran Mukamel (UCSD)

# Appendices

- Appendix 1: Justification of the ULV
- Appendix 2: A cautionary note
- Appendix 3: Simulation methods and results

# Appendix 1: Justification

- As a two-way random-effects model
- Closed-form
- Least square

## As a Two-Way Random-Effects Model

- This is also a two-way random-effects model with a single observation per cell.
- $\bar{d}_{..}$  is the LME and is unbiased for  $\mu$ .
- The variance of  $\hat{\mu} = \bar{d}_{..}$  can be estimated using

$$\hat{var}(\hat{\mu}) = \frac{1}{m}\hat{\sigma}_1^2 + \frac{1}{n}\hat{\sigma}_0^2 + \frac{1}{mn}\hat{\sigma}^2,$$

where

$$\hat{\sigma}_1^2 = \frac{\sum_{i=1}^m (\bar{d}_{i\cdot} - \bar{d}_{..})^2}{m-1}, \quad \hat{\sigma}_0^2 = \frac{\sum_{j=1}^n (\bar{d}_{\cdot j} - \bar{d}_{..})^2}{n-1},$$

$$\hat{\sigma}^2 = \frac{1}{(m-1)(n-1)} \sum_{i=1}^m \sum_{j=1}^n (d_{ij} - \bar{d}_{i\cdot} - \bar{d}_{\cdot j} + \bar{d}_{..})^2.$$

- The third item is small and can be ignored.

## Closed-form

**Theorem:** Let  $\bar{d}_{i\cdot}, i = 1, \dots, m$  and  $\bar{d}_{\cdot j}, j = 1, \dots, n$  be the marginal means. Then

- $\bar{d}_{i\cdot}, i = 1, \dots, m$  are identically distributed; so are  $\bar{d}_{\cdot j}, j = 1, \dots, n$ .
- All the above means are asymptotically uncorrelated.
- Under the null hypothesis of no group difference

$$\frac{\bar{d}_{..} - 0.5}{\sqrt{\hat{var}(\bar{d}_{..})}} \rightarrow N(0, 1)$$

where

$$\hat{var}(\bar{d}_{..}) = \frac{1}{m}\hat{\sigma}_1^2 + \frac{1}{n}\hat{\sigma}_0^2, \hat{\sigma}_1^2 = \frac{\sum_{i=1}^m (\bar{d}_{i\cdot} - \bar{d}_{..})^2}{m-1}, \hat{\sigma}_0^2 = \frac{\sum_{j=1}^n (\bar{d}_{\cdot j} - \bar{d}_{..})^2}{n-1}$$

# Least Squares

- **Remark:** the above t-statistic is identical to the two-sample t-statistic based on  $(\hat{a}_1, \dots, \hat{a}_m)$  and  $(\hat{b}_1, \dots, \hat{b}_n)$ , where

$$\hat{a}_i = \bar{d}_{i..}, \hat{b}_j = \bar{d}_{.j} - \bar{d}_{..}, i = 1, \dots, m; j = 1, \dots, n,$$

which is the least squares solution to

$d_{ij} = a_i - b_j + \epsilon_{ij}, i = 1, \dots, m; j = 1, \dots, n$  under the constraint that  $\sum_{j=1}^b b_j = 0$ .

- Least squares solutions exist and are unique up to a constant.
- These closed forms can be used when there is no need to adjust for covariates.

## Appendix 2: A Note

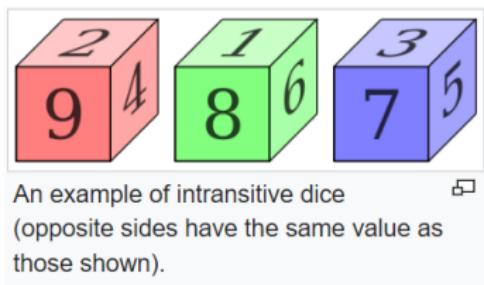
- Suppose  $Pr(G_2 > G_0) > \frac{1}{2}$  and  $Pr(G_0 > G_1) > \frac{1}{2}$
- What can we say about  $Pr(G_2 > G_1)$ ?

# A Cautionary Note About Rank-based Methods for Multiple Groups

- It turns out that

$$\Pr(G_2 > G_0) > \frac{1}{2} \\ \Pr(G_0 > G_1) > \frac{1}{2} \quad \left. \right\} \implies \Pr(G_2 > G_1) > \frac{1}{2}$$

- Nontransitive dice

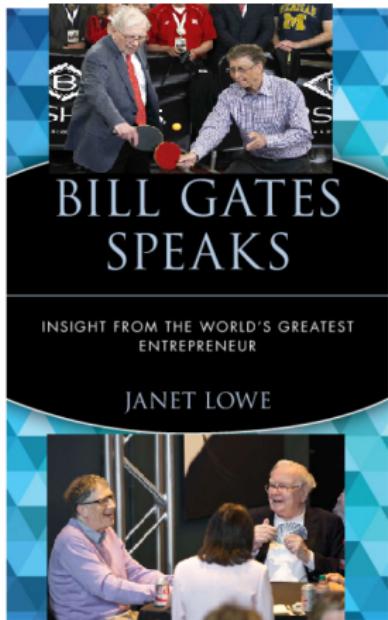


- A: 2, 2, 6, 6, 7, 7
- B: 1, 1, 5, 5, 9, 9
- C: 3, 3, 4, 4, 8, 8

$$\Pr(A > B) = \Pr(B > C) = \Pr(C > A) = \frac{5}{9}$$

# A Cautionary Note About Rank-based Methods for Multiple Groups

- A story between Warren Buffet and Bill Gates



Buffett once challenged Gates to a game of dice, using a set of four unusual dice with a combination of numbers from 0 to 12 on the sides. Buffett suggested that each of them choose one of the dice, then discard the other two. They would bet on who would roll the highest number most often. Buffett offered to let Gates pick his die first. This suggestion instantly aroused Gates's curiosity. He asked to examine the dice, after which he demanded that Buffett choose first.

"It wasn't immediately evident that because of the clever selection of numbers for the dice they were nontransitive," Gates said. "The mathematical principle of transitivity, that if A beats B and B beats C, then A beats C, did not apply. Assuming ties were

# A Cautionary Note About Rank-based Methods for Multiple Groups

- Nontransitivity could cause inconsistencies when comparing multiple conditions / treatments.
- The inference about overall effects is still valid.
- When two specific groups need to be compared, compare them directly.

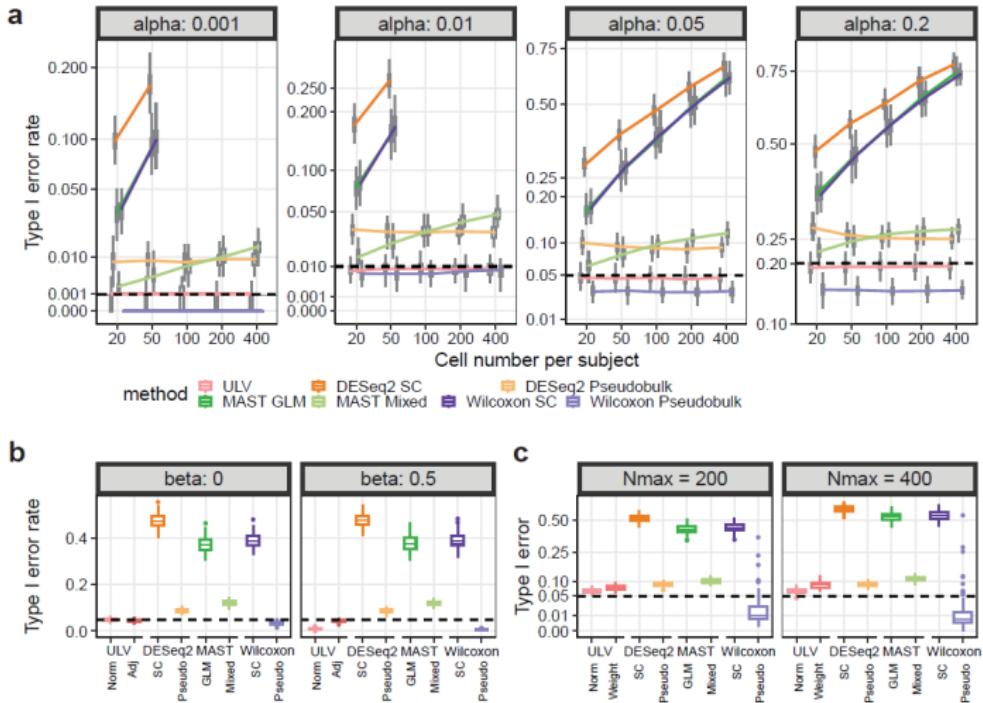
## Appendix 3: Simulations

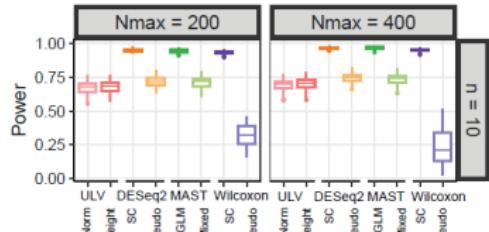
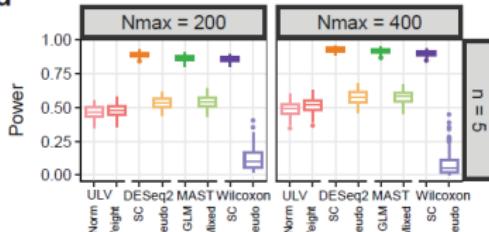
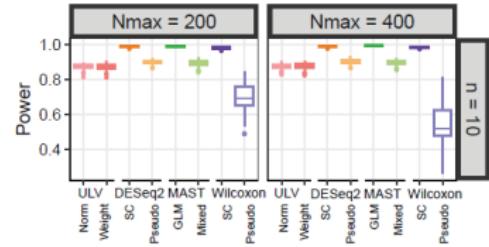
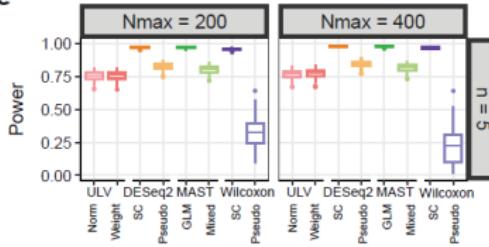
- Reference gene expression levels, inter/intraclass correlations, and correlations between genes were obtained using real data
- Simulation parameters
  - Number of cells per subject: 10-400
  - Number of subjects per group: 3-10 per group
  - Fold change: 1 to evaluate type I error rates, 4 to evaluate power
- Whether to adjust for covariates or unbalanced cluster sizes

# Methods to Evaluate

- ULV
- Two DeSeq2 methods: assumes negative binomial distribution.
  - Method 1: treat cells as independent observations;
  - Method 2: pseudobulk: treat each subject as a bulk
- Two MAST methods: based on two-part, generalized linear models adapted for zero-inflated data.
  - Method 1: GLM;
  - Method 2: GLMM
- Two Wilcoxon methods

# Simulation Results



**d****e**

## Appendix 4: Discussion of ULV

- Website: <https://github.com/yu-zhaoxia/ULV>
- Evaluate whether other difference metrics, such as the Hodges-Lehmann statistic, are better choices
- Extend the current work to model longitudinal data
- Develop robust methods for spatial single-cell data
- Extend the current work to perform causal inference