

## Response to Murphy, Fancy, and Skene, *bioRxiv*, 2023 from Mathys, et al, *Nature*, 2019 authors: Recommendations and pitfalls for computational method comparisons in single-cell genomics.

Manolis Kellis<sup>1,2,8</sup>, Jose Davila-Velderrain<sup>1,2,5</sup>, Liang He<sup>1,2,6</sup>, Hansruedi Mathys<sup>3,4,7</sup>, Li-Huei Tsai<sup>3,4,8</sup>

<sup>1</sup> MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, MA, USA.

<sup>2</sup> Broad Institute of MIT and Harvard, Cambridge, MA, USA.

<sup>3</sup> Picower Institute for Learning and Memory, Massachusetts Institute of Technology, Cambridge, MA, USA.

<sup>4</sup> Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA.

<sup>5</sup> Currently: Research Group Leader, Neurogenomics Research Centre, Human Technopole, Milan, Italy.

<sup>6</sup> Currently: Assistant Professor, Department of Public Health, U of Southern Denmark, Odense, Denmark.

<sup>7</sup> Currently: Assistant Professor, Neurobiology, University of Pittsburgh, Pittsburgh, PA, USA

<sup>8</sup> These authors jointly supervised this work: Manolis Kellis, Li-Huei Tsai.

\* e-mail: manoli@mit.edu; lhtsai@mit.edu

### Abstract

Single-cell gene expression can reveal molecular changes underlying disease, with three broadly-used computational paradigms for differentially-expressed gene (DEG) identification: cell-level analyses that treat cells as independent samples, pseudobulk aggregate analyses that average all cells of the same type in each individual; and hierarchical mixed-model analyses that model data as first sampling individuals and then sampling cells. Our Mathys, Davila-Velderrain *et al.* *Nature* 2019 paper combined mixed-model and cell-level approaches and reported 1,031 DEGs significant in both. Murphy, Fancy, and Skene (*bioRxiv* 2023) advocate instead for a pseudobulk-only approach using our study as an example, but several issues make their commentary fall short, which we discuss here. First, they exclude our mixed-model results from their comparison, thus misrepresenting our work as cell-level-only. Second, they interpret any differences between their results and ours as improvements, thus conflating false positives with true predictions and lack of power with specificity, without providing any independent support – by contrast, our paper validated our DEGs using both targeted experiments and comparison with previous human and animal studies. Third, their 16 DEGs show unrealistically-high ~1200% median changes, and only implicate the rarest cell types, suggesting unaddressed statistical artifacts stemming from outliers and small cell counts. They also fail to generalize their claims beyond a single dataset, lack any benchmarks, lack comparison to current state-of-the-art methods, and misuse statistical terminology. Lastly, they overstate their previous results from a similar Matters Arising commentary, which also advocated for pseudobulk using the same pattern of only analyzing a single dataset, methodological shortcomings, overstated conclusions, and not addressing corrections suggested by the original authors. Overall, this is a helpful example to discuss potential pitfalls in comparing bioinformatic pipelines, and to recommend perhaps more constructive strategies for such comparisons.

### Summary of our original 2019 paper

Briefly, our Mathys\*, Davila-Velderrain\* *et al.* *Nature* 2019 paper<sup>1</sup> reported the first single-cell analysis of post-mortem human brain samples for Alzheimer's Disease. We profiled 80,000 cells from 48 individuals, including 24 unaffected controls, and 24 individuals with early and late AD pathology. We identified 1,031 differentially-expressed genes (DEGs) across 6 major cell types using a combination of (a) cell-level differential analysis, and (b) a Poisson mixed model accounting for the individual of origin for nuclei and removing unwanted variation. Our methods wrote:

***Differential gene-expression analysis.*** *snRNAseq-based differential expression analysis was assessed using two tests. First, a cell-level analysis was performed using the Wilcoxon rank-sum test and FDR multiple-testing correction. Second, a Poisson mixed model accounting for the individual of origin for nuclei and for unwanted sources of variability was performed using the R packages lme4 and RUV-seq, respectively. The consistency of DEGs detected using the cell-level analysis model with those obtained with the Poisson mixed model was assessed by comparing the directionality and rank of DEGs in the two models. Consistency in directionality for all cell types was measured by counting the fraction of the top 1,000 DEGs (ranked by FDR scores) detected in cell-level analysis that showed consistent direction in the mixed model. High consistency was found, with a median fraction of 0.99 (Supplementary Table 2). Global consistency between the two models was assessed statistically using a resampling test. We tested whether the differential P value and z-score ranks corresponding to genes detected as upregulated or downregulated in the cell-level analysis were significantly higher or lower than those that were expected by chance when computed using the mixed model. Expected scores were estimated by randomly sampling*

*same-sized gene sets (n = 1,000 replicates). Both significance rank and direction deviated significantly from expectation, with directionality consistent for up- or downregulated genes. Results from the consistency tests and from both the cell-level and mixed-model differential tests are included in Supplementary Table 2. For analyses involving DEG counts, only genes that were significantly supported by both models using the criteria FDR-corrected  $P < 0.01$  in a two sided Wilcoxon-rank sum test, absolute  $\log_2(\text{mean gene expression in AD category } x / \text{mean gene expression in AD category } y) > 0.25$ , and FDR-corrected  $P < 0.05$  in a Poisson mixed model were considered. Owing to their low cell counts, pericyte and endothelial cell populations were excluded from differential analyses.*

To provide some context, there are three approaches broadly used in the field for analyzing single-cell data:

- **#1: Cell-level analysis** treats each cell as an independent sample, annotated with the phenotypes of the individual that each cell came from. This approach has the advantages of simplicity and of leveraging all the single-cell measurements, and the disadvantage of inflated statistics because cells from the same individual are not independent of each other.
- **#2: Pseudobulk analysis** treats all the cells of a given cell type from the same individual as a single measurement, typically by averaging their expression patterns. This approach has the advantages of simplicity and of not counting the multiple measurements from the same individual as independent, but it also has several disadvantages, including: reduced statistical power, because it is not leveraging all the richness of single-cell data; treating each individual equally, regardless of cell count or quality; results that depend on the aggregation method (mean, sum, median, etc); losing information about multimodal distributions or outliers; and losing information about expression variance between cells of the same individual, thus making it harder to capture statistically-significant but subtle expression differences.
- **#3: Hierarchical (or mixed-model) analysis** treats the sampling procedure as a two-step sampling, first of individuals, and second of cells from those individuals. This approach has the advantage of more accurately modeling the underlying process through which the data was sampled as a nested generative process, but the disadvantage of increased statistical and computational complexity and potential model misspecifications, as it makes more assumptions about the data sampling and distribution processes.

In our [Nature](#) 2019 paper, we combined approaches #3 and #1 and reported 1,031 genes that were captured by both approaches. We combined the two approaches to leverage their complementary nature: the statistical rigor of approach #3, a Poisson mixed model, accounted for the individual of origin for nuclei and for unwanted sources of variability, to ensure that stringent statistical significance thresholds are applied; and the nonparametric nature of approach #1, a cell-level Wilcoxon rank-sum test, to ensure that we did not introduce statistical artifacts in cases where approach #3 alone may have not modeled the data distribution correctly and may have resulted in false positives.

We have since formalized approach #3 into NEBULA (a NEgative Binomial mixed model Using a Large-sample Approximation)<sup>2</sup>, by Liang He, Jose Davila-Velderrain *et al*, [Communications Biology](#), 2021. In addition, many new statistical, computational, and experimental innovations have been made, across cell-level, pseudobulk, and hierarchical approaches and their combinations, and formal benchmarks have been developed for comparing them<sup>3,4</sup>. We believe that different strategies have distinct advantages and disadvantages in different settings, which can depend on a gene's level of expression, the number of cells it's expressed in, the variability of expression across cells, the variability across individuals, how bursty the expression is, whether it is bimodal or even multimodal, and other factors. We encourage students and practitioners of single-cell methods for differential expression analysis to explore this rapidly-expanding, and very exciting field.

### [Evaluation of Murphy \*et al\* commentary](#)

It seems that the Murphy, Fancy, and Skene, [bioRxiv](#) 2023 commentary<sup>5</sup> seeks to advocate for a pseudobulk-only approach (#2) for single-cell datasets, using our paper as an example where such a method might outperform previous approaches.

We want to start by thanking the authors for their attention to our manuscript, and for their efforts in reprocessing and analyzing our data. We also want to acknowledge that there can be advantages to pseudobulk methods, and that no single method is likely to outperform all others in all settings. We also want to clearly state that the technologies and methods used in our paper are now over four years old, in a rapidly moving field, so we expect many newer analyses will be able to clearly outperform previous papers in the field, including ours. Lastly, we encourage students and practitioners to revisit previous papers as a general practice

of open science, and to share the results in open, respectful, and transparent ways, which can move the field forward collaboratively and constructively.

That said, there are some issues with the Murphy *et al.* commentary that make it fall short of their goal of demonstrating a clear advantage of pseudobulk methods. Briefly: they only analyze partial results from our paper, not our full methodology, misrepresenting it in their comparisons as cell-level-only instead of joint hierarchical and cell-level; they conflate false positives with true predictions and lack of power with specificity for their method; they only find unrealistically-high ~12-fold median changes across their 16 DEGs, and they only find DEGs for the rarest cell types, likely due to noisy small counts; they lack any independent support for the differences observed; they fail to generalize their claims beyond a single paper; they lack proper benchmarks or comparison to current state-of-the-art methods; they make strong but weakly-supported claims; they use incorrect terminology; and they do not consider that even their few DEGs are potentially unaddressed statistical artifacts that simple filters could exclude.

We detail these issues below, which we hope can be helpful to the authors and to the broader community:

1. Murphy *et al.* do not analyze our results from approach #3 in their comparison, thus misrepresenting our work as using only approach #1 in their comparison. Any potential claims of methodological superiority are dubious, as they are not actually using our results, but instead an intermediate methodological step of our paper. We recommend comparing their method to our actual results.
2. Murphy *et al.* fail to consider the possibility that their method is simply underpowered due to its pseudobulk nature (#2) and lack of removal of unwanted variation, and instead interpret the fact that it misses our 1,031 DEGs as high specificity. We show that our 1,031 DEGs (detected by both our hierarchical model #3 and our cell-level analysis #1) are biologically meaningful through direct experimental validation and comparison with previous studies in human and in animal models. Murphy *et al.* fail to provide any independent support or validation of either their positive or their negative results.
3. Murphy *et al.* report a strong negative correlation ( $R=-0.48$ ) between the number of cells and the number of DEGs called by their method, when in fact power should increase, not decrease, with the number of cells. More specifically, the only cell types where any DEGs are called are the ones with the smallest numbers of cells, i.e. 1 DEG in OPCs, 1 in Astrocytes, and 14 in Microglia (the cell type with the fewest numbers of cells by far). This inverse relationship suggests that the genes detected by Murphy *et al.*'s pseudobulk method might represent unaddressed statistical artifacts stemming from outlier genes in low cell counts, as expected for these low-abundance cell types, and we encourage the authors to investigate further along these lines to understand their technical (or biological) sources.
4. Murphy *et al.* write: “*the absolute log<sub>2</sub> fold change of our DEGs is over 20 times larger*”, with a “median LFC of 3.33”. We note that a median fold change of  $2^{3.33}=11.9$ -fold, or 1190%, across their 16 DEGs seems biologically unlikely in a complex polygenic trait. An alternative interpretation of their results is that the only genes detected by their weakly-powered pseudobulk method are in fact false positives, stemming from unaddressed statistical artifacts of very small cell counts. Perhaps using a threshold on the minimum number of cells required to call DEGs with pseudobulk methods, or using a hierarchical model similar to ours, which considers the single-cell nature of our data, could help avoid such artifacts.
5. More generally, Murphy *et al.* interpret any disagreement with prior methods as an improvement, in the absence of any further computational or experimental evidence: they interpret the genes they miss as “false positives” of previous methods, and the new genes they find as “true positives” of their method, when in fact the opposite could be true. Disagreement does not always mean improvement, and we encourage Murphy *et al.* to critically evaluate and validate both positive and negative results from their method, with an open mind that their method might sometimes be incorrect, in either direction.
6. Murphy *et al.* seek to demonstrate the superiority of pseudobulk using a single paper, which makes it hard to draw broad conclusions. There are many published single-cell datasets besides ours, which could be used to demonstrate generalizability of the method. Applying their method to more than one dataset would be our recommended way to advocate for its broad applicability and usefulness (being cautious to use the actual reported results of those papers, and to not misinterpret their methods).
7. Murphy *et al.* do not compare their approach to any recent methods for cell-level (#1), pseudobulk (#2), or hierarchical (#3) analysis developed since our paper, or to any independent benchmark datasets established since then. Without such comparisons and benchmarks, it is difficult to determine whether

their approach is an improvement, and whether it should be applied more broadly.

8. Murphy *et al.* write that: “*Pseudobulk differential expression (DE) analysis has recently been proven to give optimal performance compared to both mixed models and pseudoreplication approaches*”, citing in support their own previous commentary<sup>6</sup>, similar to the one they write here. This previous commentary was also based on a single, and in fact simulated in that case, dataset, and also contained several methodological, statistical, and logical shortcomings which were pointed out in the same issue<sup>7</sup> by the authors of the original publication, and thus should be known to Murphy and Skene. We encourage Murphy *et al.* to avoid broad conclusions without broad benchmarks and to address the shortcomings pointed out by the authors.
9. In terms of nomenclature, Murphy *et al* incorrectly write “*this cell-level differential expression approach, also known as pseudoreplication...*”. This is a misuse of an otherwise well-defined statistical term, which conflates a property of method #1 (pseudoreplication bias) with the name of the method. Murphy and Skene first used the term incorrectly in their previous commentary<sup>6</sup> on Zimmerman *et al*<sup>8</sup>, who had used the term correctly. The terms “pseudoreplication” and “cell-level analysis” are not equivalent or interchangeable.
10. Murphy *et al.* conclude by writing: “*We would like to highlight that the use of [cell-level analysis] in neuroscience research is not isolated to the author’s work, others have used this approach<sup>[refs 11,12,13]</sup> and their results should be similarly scrutinized*”. It is unclear if the authors are suggesting that every paper that uses cell-level analysis warrants a similar “Matters Arising” commentary by Murphy and Skene. Similarly, when a new method appears, going back and writing “Matters Arising” commentaries on each of the previous papers that can be potentially improved by these methods does not seem to be the best strategy, as we discuss below.

Overall, while pseudobulk methods and other state-of-the-art tools can outperform previous studies, including ours, we find that Murphy *et al.* fail to confidently demonstrate that their pseudobulk-only approach is superior due to several shortcomings in their study, and the actual results of the Murphy *et al.* commentary fall short of their claims. This should not be interpreted as a general shortcoming of pseudobulk methods, as the results might stem from how Murphy *et al.* utilized them on our particular dataset. We hope our suggestions are useful, and perhaps help improve their methodology on this and on future datasets by avoiding these pitfalls.

### Contrasting strategies for methodological benchmark analyses

The two commentaries<sup>5-7</sup> by Murphy and Skene bear several similarities, and seem to share the goal of demonstrating that pseudobulk methods are superior, and a strategy of individual commentaries on individual papers. Both commentaries: (1) seek to utilize the “Matters Arising” mechanism of Nature-family journals to advocate for pseudobulk methods; (2) analyze a single dataset; (3) position their work as correcting flaws in the literature; (4) do not provide independent validation for their findings; (5) only provide marginal improvements or minor clarifications, if any, but have several flaws of their own; (6) do not address these flaws when they are identified by the authors of the original paper; and (7) are ultimately cited by the commenters as independent papers that draw broad conclusions in the field instead of narrow dataset-specific results.

We believe that this strategy might not be optimal for several reasons: (1) it’s hard to generalize from performance on a single dataset, and without comparison with state-of-the-art tools or systematic benchmarks; (2) such commentaries can spread confusion when they are only glanced at quickly by busy researchers who might take the criticism at face value; (3) reviewing these commentaries in isolation puts undue burden on the editors and authors of the original papers to write detailed responses over email or via journal-mediated correspondence, when the issues are broader than any one individual paper; (4) the many suggestions and corrections made by the authors of the original papers in these exchanges can sometimes go unaddressed, ignored, or dismissed by the commenters, leading to wasted time and energy; (5) these commentaries can put disproportionate focus on individual papers and potentially impact the careers of young scientists, when hundreds of papers use similar methodologies; and lastly (6) with hundreds of articles published each year, and with constantly-improving methods, going back and writing a “Methods Arising” commentary for each of these papers every time improved methods appear seems inefficient and unwarranted.

A different strategy that Murphy and Skene could consider, should they continue seeking to establish that pseudobulk methods are superior, could be: (1) applying their pseudobulk approach to multiple datasets at once; (2) using well-established independent and rigorous benchmark metrics and statistics for evaluation and performance comparison; (3) including other state-of-the-art methods against these benchmarks; (4) evaluating observed differences biologically using independent datasets and directed biological validation

experiments; (5) seeking to understand the biological and technical drivers of these differences; (6) scrutinizing their own predictions for potential confounders such as small cell counts and outliers; and (7) presenting the results in a stand-alone paper, which can then be peer-reviewed for conference or journal publication, as is typical both for papers introducing new methodologies and for papers benchmarking previous methods.

If the goal of Murphy and Skene is to provide general methodological insights and recommendations that can help the field move towards broader adoption of robust tools by the community, this second strategy could be more constructive and impactful, and we encourage them to pursue it instead.

A third strategy, which is increasingly taken in several fields, is to coalesce around comprehensive databases and portals, encompassing raw and reprocessed datasets across hundreds of studies, which can allow re-running of primary data processing, QC filtering, normalization, and secondary analyses more systematically. Several such efforts are underway in single-cell genomics as the field is maturing. These efforts can ensure that the most accurate cell-level, pseudobulk, and hierarchical methods can be re-run systematically on multiple datasets and their impact evaluated. We do not believe that only one method should be chosen, as each method has distinct advantages and disadvantages. Such efforts can help make more informed choices and enable their broad applicability.

### References:

1. Mathys H, Davila-Velderrain J, Peng Z, Gao F, Mohammadi S, Young JZ, Menon M, He L, Abdurrob F, Jiang X, Martorell AJ, Ransohoff RM, Hafler BP, Bennett DA, Kellis M, Tsai LH. Single-cell transcriptomic analysis of Alzheimer's disease. *Nature*. 2019 Jun;570(7761):332–337. PMID: PMC6865822
2. He L, Davila-Velderrain J, Sumida TS, Hafler DA, Kellis M, Kulminski AM. NEBULA is a fast negative binomial mixed model for differential or co-expression analysis of large-scale multi-subject single-cell data. *Commun Biol*. 2021 May 26;4(1):629. PMID: PMC8155058
3. Junttila S, Smolander J, Elo LL. Benchmarking methods for detecting differential states between conditions from multi-subject single-cell RNA-seq data. *Brief Bioinform [Internet]*. 2022 Sep 20;23(5). Available from: <http://dx.doi.org/10.1093/bib/bbac286> PMID: PMC9487674
4. Gagnon J, Pi L, Ryals M, Wan Q, Hu W, Ouyang Z, Zhang B, Li K. Recommendations of scRNA-seq Differential Gene Expression Analysis Based on Comprehensive Benchmarking. *Life [Internet]*. 2022 Jun 7;12(6). Available from: <http://dx.doi.org/10.3390/life12060850> PMID: PMC9225332
5. Murphy AE, Fancy NN, Skene NG. Avoiding false discoveries: Revisiting an Alzheimer's disease snRNA-Seq dataset [Internet]. *bioRxiv*. 2023 [cited 2023 Apr 17]. p. 2023.04.01.535040. Available from: <https://www.biorxiv.org/content/10.1101/2023.04.01.535040v1>
6. Murphy AE, Skene NG. A balanced measure shows superior performance of pseudobulk methods in single-cell RNA-sequencing analysis. *Nature communications*. 2022. p. 7851. PMID: PMC9780232
7. Zimmerman KD, Evans C, Langefeld CD. Reply to: A balanced measure shows superior performance of pseudobulk methods in single-cell RNA-sequencing analysis. *Nature communications*. 2022. p. 7852. PMID: PMC9780219
8. Zimmerman KD, Espeland MA, Langefeld CD. A practical solution to pseudoreplication bias in single-cell studies. *Nat Commun*. 2021 Feb 2;12(1):738. PMID: PMC7854630