Detecting Rhythmic Gene Expression in Single Cell Transcriptomics

Bingxian Xu^{1,2}, Rosemary Braun^{1,2,3,4,5*}

¹Department of Molecular Biosciences, Northwestern University, Evanston, IL 60208, USA; ²NSF-Simons Center for Quantitative Biology, Northwestern University, Evanston, IL 60208, USA; ³Department of Engineering Sciences and Applied Mathematics, Northwestern University, Evanston, IL 60208, USA;

⁴Department of Physics and Astronomy, Northwestern University, Evanston, IL 60208, USA; ⁵Northwestern Institute on Complex Systems, Northwestern University, Evanston, IL 60208, USA; *To whom correspondence should be addressed. Email: rbraun@northwestern.edu.

Abstract

An autonomous, environmentally-synchronizable circadian rhythm is a ubiquitous feature of life on Earth. In multicellular organisms, this rhythm is generated by a transcription—translation feedback loop present in nearly every cell that drives daily expression of thousands of genes in a tissue—dependent manner. Identifying the genes that are under circadian control can elucidate the mechanisms by which physiological processes are coordinated in multicellular organisms. Today, transcriptomic profiling at the single-cell level provides an unprecedented opportunity to understand the function of cell-level clocks. However, while many cycling detection algorithms have been developed to identify genes under circadian control in bulk transcriptomic data, it is not known how best to adapt these algorithms to single-cell RNAseq data. Here, we benchmark commonly used circadian detection methods on their reliability and efficiency when applied to single cell RNAseq data. Our results provide guidance on adapting existing cycling detection methods to the single-cell domain, and elucidate opportunities for more robust and efficient rhythm detection in single-cell data. We also propose a subsampling procedure combined with harmonic regression as an efficient, reliable strategy to detect circadian genes in the single-cell setting.

1 Introduction

The circadian rhythm is an approximate 24-hour oscillation of physiology, metabolism and behaviour that enables organisms to adapt to a changing daily environment [1, 2, 3]. At the cellular level, autonomous molecular oscillations are generated by a transcription-translation feedback loop. In drosophila, for example, the CLOCK/CYCLE (CLK/CYC) heterodimer binds to the E-box to activate gene expression of period (per) and timeless (tim); the PER and TIM proteins then dimerize in the cytoplasm and translocate to the nucleus to inhibit the DNA binding activity of CLK/CYC [4]. This core clock circuit in turn drives circadian oscillations of hundreds of downstream targets. It has been reported that $\sim 40\%$ of all protein coding genes may exhibit circadian oscillations in a tissue–specific manner [5], underscoring the role of the circadian rhythm in orchestrating physiological processes in multicellular organisms.

Proper temporal coordination requires synchronization of tissue–specific rhythms, which in turn requires synchronization of cell–level clocks. Loss of synchrony between organ systems and dampening of circadian rhythms has been implicated in aging and disease [6]. In older mice, for example, it has been reported that the flattened amplitude of circadian oscillation is related to deficits in long-term spatial memory due to impaired sleep [7]. Even within the same tissue, it was recently observed that the number of detectably cycling genes decreases as the animal gets older [8]. Such dampening may be due to an overall loss of amplitude, or a loss of synchrony between cells that leads to no discernible rhythm at the tissue level.

Today, single cell RNAseq profiling provides an opportunity to examine this coordination by investigating cell type—specific transcriptomic cycling. Circadian transcriptomic timeseries data (both bulk and single-cell) typically comprises samples taken every 2–4 hours over a 24–48 hour period, potentially with replicates [9]. The goal of such studies is to identify genes that are under circadian control in different conditions, providing insights into how the circadian rhythm orchestrates cellular physiology. Because these data are noisy, sparsely sampled in time, and high dimensional in the number genes, statistical tests for evidence of cycling and differential cycling remain an area of active research [10, 11, 12]. To date, however, all proposed methods were developed in the context of bulk RNAseq data, leaving open the question of how best to analyze circadian single-cell data. Single cell RNAseq measurements

yield many more observations than bulk data (one per cell), but are much noisier due to stochastic gene expression and drop-outs; in consequence, methods designed for bulk RNAseq may not translate well to the single cell context. Below, we briefly review state-of-the-art methods for cycling detection in bulk transcriptomic data, and highlight the opportunities and challenges posed by single-cell RNAseq profiling.

Perhaps the simplest form of cycling detection is harmonic regression, in which a 24 h sinusoidal curve is fit to the data and goodness-of-fit statistics serve as an assessment of cycling. However, the noisiness of the data and the low number of replicates per timepoint in bulk transcriptomic data mean that a single outlying observation can produce a significant "cycling" component, leading to a large number of false positives when applied to bulk transcriptomic data. To overcome this challenge, AR-SER [13] removes linear trends in the data and smooths the data using a fourth-order Savitzky-Golay filter before the harmonic regression. In addition, there are methods that rely on Fourier analysis [14, 15], essentially conducting a discrete Fourier transform and assessing the statistical significance of the peak corresponding to a period close to 24 hours [16]. Like harmonic regression, Fourier-based methods also have inaccuracies [17, 18] due to insufficient temporal sampling (frequency and duration).

Because it has been observed that cycling genes may exhibit sharply peaked or asymmetric waveforms, non-parametric approaches have been developed to test for evidence of cycling without assuming sinusoidality [17, 18, 19, 20]. JTK-cycle [17] employs the Jonckheere-Terpstra trend test/Kendall's τ rank correlation to look for monotonic patterns of rising and falling in a 24-hour window. Here, the observed pattern of gene expression is compared to template waveforms of different phases and asymmetries (e.g. a longer rising interval than falling interval). Because corrections need to be made for the multiplicity of templates considered, JTK-cycle can lose power as the number of patterns of interest increases. To circumvent the need for predetermined waveforms, RAIN [18] was developed using the general umbrella test, a statistical test for an umbrella shape with a flexible inflection point, allowing it to accommodate to asymmetric waveforms automatically. Two alternative methods, SW1PerS [21] and TimeCycle [12], use techniques from topological data analysis to nonparametrically quantify the cyclicity of observed expression profiles. In the context of bulk transcriptomic data, these nonparametric approaches generally outperform parametric ones. Nevertheless, it has also been shown that no cycling detection method is universally optimal, and that the best analysis method depends on both the study design (i.e. the frequency of sampling, number of periods sampled, and replicates per timepoint) and the shape of the waveforms of interest [9].

Analysis of single-cell RNAseq data presents new opportunities and challenges. The large number of cells assayed at each timepoint opens the possibility of regarding them as many replicate samples (in contrast to bulk RNAseq studies, which often have only two or three), enabling the variance between cells to inform the analysis. However, the number of cells is very large, and may vary from timepoint to timepoint. As a result, analysis methods that require a regular number of samples per timepoint (such as ARSER [13]) or those that scale poorly with the number of replicates may not be directly applicable to single-cell data. As an alternative, one can consider averaging expression levels across all cells of the same type. Methods developed for bulk transcriptomic data may then be applied to the resulting "pseudo-bulk" data. In support of this approach, it was reported that considering cells as replicates for differential expression analysis can result in a large number of false positives due to systematic correlations between the cells of a given sample [22]. There, the authors found that pseudobulking yields more reproducible results in tests of differential expression than treating cells as independent observations. However, it is not known to what degree pseudobulking may enhance the specificity of cycling detection (by reducing false positives) vs. reducing its sensitivity (by lowering the number of replicates), and there remains no guidance on the best way to analyze single-cell circadian timeseries data.

The goal of this work is to evaluate the performance of various existing approaches applied to a circadian timecourse single-cell RNAseq dataset of *drosophila* brain tissue [23]. We compared a number of different methods (including JTK-cycle and RAIN) as well as a variety of application approaches (i.e., treating cells as distinct samples or averaging to create a pseudobulk) and assessed the computational efficiency, reproducibility, and robustness to noise of the algorithms. Our results suggest that methods designed for circadian detection in bulk RNAseq data may not be optimal for single-cell data, and suggest opportunities for new approaches.

2 Materials and Methods

2.1 Data acquisition

We use publicly available single-cell RNAseq data from Ma et al.[23] in our study. The data comprise four circadian time-series of fly brains, with two time-series collected under each of two conditions (12h light/dark [LD] and dark/dark [DD]). Each series comprises six samples collected every four hours. In total, 4671 cells were assayed. Transcriptomic data was obtained as a Seurat object [24, 25] containing all relevant metadata. The original authors annotated these into 39 clusters, of which 17 were assigned to known cell types. In the analyses that follow, we used the normalized counts (kindly provided by the authors) of the LD cells as input to the cycling detection algorithms. Details of the preprocessing may be found in the Supplement.

2.2 Creation of pseudo-bulk expression profiles

In addition to considering single cells as individual observations, we also constructed pseudo-bulk expression profiles for each cell type (cluster). Pseudo-bulk profiles were constructed from the normalized expression matrix $X \in \mathbf{R}^{g \times c}$ where g denotes the number of genes and c the number of cells for the cell-type of interest. Each column in X corresponds to a cell with a time stamp $t_i \in \{t\}$, where $\{t\}$ is the set of sampling times. The pseudo-bulk expression profile for gene g at time t is the average across the cells in a given cell type with timestamp $t_i = t$,

$$\hat{x}_g(t) = \frac{1}{c_t} \sum_{k:t_k = t} X_{gk} \,, \tag{1}$$

where c_t denotes the number of cells sampled at time t. Pseudobulk temporal profiles were calculated per Eq. 1 for the four largest clusters (Figure 4).

2.3 Application of cycling detection algorithms

JTK-cycle [17], RAIN [18], and harmonic regression were applied to single cell data using the recommended parameter settings (see Supplement), treating each cell as an independent sample. Additionally, we applied JTK-cycle, RAIN, harmonic regression, and ARSER [13] to the pseudobulk data. ARSER, harmonic regression, RAIN and JTK-cycle were all implemented using their respective R packages [17, 26, 18, 27]. We excluded from our analysis SW1PerS [21], which does not return a p-value, and TimeCycle [12], which requires at least 18 time points to construct the cycle.

2.4 Synthetic data contamination

In addition to testing the performance of the algorithms against the original data, we wanted to explore whether a systematic artefact affecting a single time-point would generate false-positive results. In [22], it was demonstrated that differential expression analysis is prone to false positives when cells are considered replicates due to the fact that a small artefact would be amplified by the large number of cells. Would the same hold true for cycling detection?

To examine this possibility, we identified genes that showed no evidence of cycling, and systematically contaminated the data by elevating the expression level of genes at specific timepoints to explore whether the algorithms would erroneously identify it as cycling. We restricted this analysis to the largest cell cluster (1:DN1p_CNMa) collected under the LD condition. At a given timepoint, we increased the expression level for each gene by an amount proportional to its mean expression over all cells (in that cluster at all timepoints. For each round of analysis, we only increased the expression level of cells collected in a single time point and replicate, mimicking an artifact affecting a single scRNA-seq collection. That is, for each gene X measured in cell c(t) at timepoint t, we increase its expression to

$$\tilde{X}_{c(t)} = X_{c(t)} + \lambda \langle X_{c(t)} \rangle_{c,t}, \qquad (2)$$

where $X_{c(t)}$ ($X_{c(t)}$) denotes the contaminated (original) normalized count for gene X in cell c, $\langle \cdot \rangle_{c,t}$ denotes an average over all cells (from the same cell type) at all timepoints, and λ is a parameter that controls the extent of the noise contamination. Non-cycling genes selected for synthetic contamination

were chosen as those that were detectably expressed in at least half of the cells and had a harmonic regression p > 0.9 (treating cells as replicates). This yielded 59 "non-cycling" genes.

3 Results

3.1 A framework for evaluating cycling detection reproducibility in singlecell RNAseq datasets

The large number of observations (cells) obtained at each timepoint in a single-cell RNAseq experiment permits a novel approach for evaluating the reliability of cycling detection methods. The approach is based on the following conjecture: if a gene is truly under circadian control in a certain cell type, it should be detected as cycling reliably even when considering only a subset of the cells at each timepoint. We note that for most cell-type clusters, even taking half of the observed cells yields more replicates than are typically found in bulk circadian transcriptomics studies (where there are commonly <3 observations per timepoint). We thus compute the reliability of a cycling detection method as follows. For each cell-type at each timepoint, we select at random half of the cells as "sub-experiment" 1, and consider the rest as "sub-experiment" 2. We then apply the cycling detection algorithm to both subexperiments, yielding for each sub-experiment the genes detected as cycling. Ideally, the intersection of those sets should be complete, with the same genes detected as cycling (or not) in each sub-experiment. We use the size of the intersection N_{\cap} over the number of cyclers detected when all cells are used N_{full} as a measure of reliability. The choice of developing this measure rather than using existing measures such as the Jaccard index (N_{\cap}/N_{\cup}) , where N_{\cup} is the union of the cyclers in the two sub-experiments) is that the union of the cyclers detected in the two sub-experiments N_{\cup} are only a small fraction of the cyclers detected when all cells are used $(N_{\cup} \ll N_{\text{full}})$. Since typical analysis of single-cell data would use all cells, we wanted $N_{\rm full}$ to be the basis for comparison. The same procedure can also be used with pseudo-bulk data; in this case, division into the two sub-experiments is done before Eq.1 is calculated.

In our tests, we performed 10 random splits of the data into the two sub-experiments to gather statistics for N_{\cap} . Higher N_{\cap}/N_{full} indicates greater reproducibility of the cycling detection, even when the data are subsampled.

3.2 Computational efficiency

We first tested how run-time scales with the total sample size, ranging from ~ 10 to 200 cells, when treating each cell as an independent observation (i.e. without pseudo-bulking). We focused on three methods—JTK-cycle [17], RAIN [18], and harmonic regression [27]—chosen because of their popularity and their ability to handle uneven replicates. The latter is a crucial feature, since the number of cells is likely to vary per time-point in single-cell data; as a result, we exclude methods (such as ARSER) that require a consistent number of observations at each timepoint. We then considered the computational cost of running these algorithms without pseudo-bulking.

In contrast to what was previously reported [20], we observed that RAIN can actually be faster than JTK-cycle when the sample size (i.e. number of cells) is small (Figure S1). However, as sample size increases, RAIN becomes slower than JTK-cycle as expected. These differences in efficiency are likely due to implementation details, rather than the complexity of the underlying test; we derived the theoretical scaling relationship as having $\mathcal{O}(m^2)$ computational complexity in the number of samples m in both cases (see Supplementary Information). However, we observe empirically that the measured run time for both JTK-cycle and RAIN deviates considerably from the estimated quadratic complexity (Figure 1). Our empirical benchmarking suggests that for a moderate single cell dataset with only 1000 cells, JTK-cycle (RAIN) will take at least seven hours (sixteen days) to complete on a MacBook Pro with 2.9GHz dual-core Intel Core i5 and 8G memory. In particular, both implementations appear to scale exponentially, which may make them unfeasible for application to single-cell data without pseudo-bulking. The discrepancy between the theoretical computational complexity and the empirical performance also suggests that these algorithms efficiency may still be improved. Harmonic regression, as expected, scales linearly in the number of observations, and takes < 10s to compute even for very large clusters.

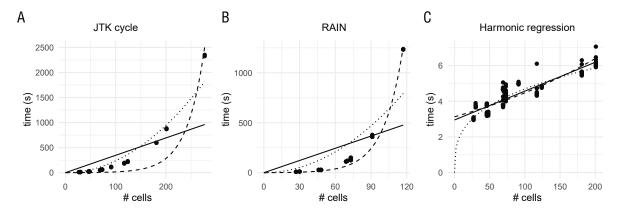


Figure 1: Scaling relationship between the number of samples and run time for (A) JTK-cycle, (B) RAIN, and (C) harmonic regression. Trend-lines were fitted for linear (solid), quadratic (dotted) and exponential (dashed) computational complexities. Note the difference in y-axis scales.

3.3 Reliability of cycling detection in single-cell data

We next explored the reliability of cycling detection in sc-RNAseq data using the concordance framework described above. As in the previous section, we applied JTK-cycle, RAIN and harmonic regression to eight cell types within the Ma data that contained at least two cells at each time point, and which were sufficiently small to analyze given the efficiency attributes noted above. Our choice of methods is motivated both by their popularity [28, 8, 23, 29] and their ability to handle replicates and run efficiently when sample size is low [20]. For all methods tested, we consider a gene to be cycling if the p-value is less than 0.05. We compared how these three methods differed in the number of cycling genes detected and in their subsample concordance. Figure 2 shows the number of genes detected as cycling when no sub-sampling is done (N_{full} , panel A); the number of genes in the intersection of two sub-experiments when subsampling is done (N_{n} , panel B); and their ratio (panel C).

Interestingly, while it has been reported that RAIN is less conservative than JTK-cycle in bulk data [9, 18], we observe that the number of detected cycling genes by these two methods tend to be similar when cells were considered as replicates (Figure 2). We then conducted subsampling (as described above) to examine the reproducibility of genes detected as cycling in both subsamples. Interestingly, though harmonic regression considered more genes to be cycling compared to JTK-cycle and RAIN, all the obtained N_{\cap} tended to be similar (Figure 2). Lastly, we looked at the ratio between N_{\cap} and N_{full} and observed that RAIN had the highest subsample concordance in all tested cases, around two times greater than that from both JTK-cycle and harmonic regression.

We note that the above analysis used unadjusted p-values. To examine if adjusting p-values for multiple hypotheses increases subsample concordance, we repeated above analysis using FDR values computed with the Benjamini-Hochberg method, setting FDR < 0.05 as the cycling detection threshold. Interestingly, we observed a general decrease of subsample concordance (Figure S2) with this more stringent criterion, suggesting that using the FDR is insufficient to robustly identify circadian genes.

We then examined whether the various methods identified the same genes as cycling both across subsamples and across methods (Figure 3 and Supplemental figures S5, S6). Harmonic regression, which is known to have a high false–positive rate, yielded many subsample-specific cyclers (i.e. genes detected as cycling in one but not the other of N1, N2, shown as the first two sets in the upset plot), confirming that the harmonic regression p-value alone may generate a large amount of false positive cyclers. However, considering the overlap in the sub-experiments mitigates this issue. For the genes detected by harmonic regression in both sub-experiments N1 and N2 (ie, sets with overlap in HR-N1, HR-N2 in Figure 3), the majority are also identified by RAIN and JTK-cycle as well. Additionally, we again observed that adjusting p-values does not increase the agreement between each experiment (Figure 3). Importantly, this suggests that looking for genes consistently detected as cycling across two sub-samples using a fast algorithm (harmonic regression) may be a reliable and efficient alternative to more computationally intensive methods.



Figure 2: A: The number of cycling genes detected using all cells, $N_{\rm full}$. B: The average number of cyclers detected in both subsamples across 10 trials, N_{\cap} . C: The ratio of N_{\cap} to N_{full} , the proportion of consistently–detected sub-sample cyclers relative to those found using all cells. Tests were conducted in the largest eight labeled cell clusters that RAIN can handle. JTK-cycle and harmonic regression were run in parallel. Error bars indicate standard deviation across the 10 subsamplings.

3.4 Cycling detection using pseudo-bulk data

The problems of large run-times and uneven replicates can also be circumvented by first averaging across cells of a given cell type within each timepoint, yielding a single "pseudo-bulk" gene expression value for each cell type at each timepoint. Cycling detection may then be applied to the resulting pseudo-bulk data, as was done in a dataset of mouse SCN neurons [29]. To examine how pseudobulking affects cycling detection, we applied JTK-cycle and harmonic regression to pseudo-bulk data and compared it against the results obtained using cells as replicates. As an additional comparison, we also conducted cycling detection on pseudo-bulk data using ARSER (which cannot handle uneven replicates) and RAIN (which is computationally inefficient in the non-pseudo-bulk setting, Figure 1B). We observed that the number of cycling genes (N_{full}) detected using pseudo-bulk data was similar for JTK-cycle and ARSER. RAIN identified the largest amount of rhythmic transcripts, approximately twice that of JTK-cycle and ARSER, and slightly more than that of harmonic regression (Figure 4A). This is expected from previous studies that also noted RAIN's permissiveness. For JTK-cycle and harmonic regression, where we can directly compare the outcome of treating cells as replicates or as a pseudo-bulk tissue, we observed that both methods detected more cycling genes when cells were treated as replicates, which points out the intuitive fact that a larger number of replicates can help the identification of less obvious cyclers(Figure 4A, Figure S3, S4).

In addition, we split our data into two subsamples as described previously, both containing cells collected at all of the time points as required by JTK-cycle and RAIN. With this, we conducted both the pseudo-bulking and the subsequent cycling detection using all samples ("full" data), as well as the first and second subsample separately. We then counted the number of gene in the intersection of those considered cycling in both subsamples (N_{\cap}) , and quantified N_{\cap}/N_{full} as a measure of the concordance. From our analysis, we observed the lowest subsample concordance when JTK-cycle was applied to the pseudo-bulk (Figure 4B,C).

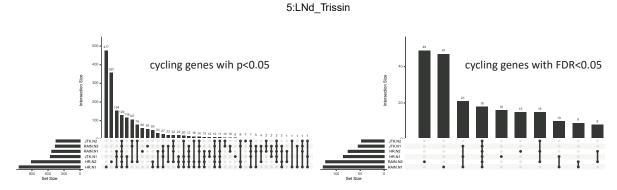


Figure 3: UpSet plots showing the intersection of detected cycling genes, with either p-values or FDR-adjusted p-values less than 0.05, by harmonic regression, RAIN and JTK-cycle in one instance of the two sub-samples. Here, HR indicates harmonic regression; JTK = JTK-cycle; N1 = sub-experiment 1; N2 = sub-experiment 2.

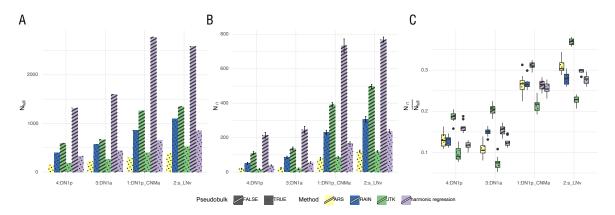


Figure 4: A: the number of cycling genes detected using each method. B: the number of genes detected as cycling in both samples in each method. C: subsample concordance of each method.

3.5 Impact of systematic contamination

While we have shown that considering cells as replicates can reliably detect more cyclers than pseudobulking, it has been reported that methods designed to conduct single cell differential expression analysis where cells were considered as replicates tend to overestimate the number of differentially expressed genes, due to the fact that the cells from a given sample are not truly independent observations [22]. To test whether the same idea holds true for circadian detection algorithms, we first synthetically increased the expression of genes that are considered not cyclic at individual timepoints (see Materials and Methods) and applied harmonic regression on data with and without pseudo-bulking. As expected, the estimated phase of these non-cyclic gene converged to the time of elevated expression (Figure 5A). Additionally, we observed that p-value decreases as the magnitude of gene expression elevation increases; consequently, a gene will have a higher chance of being a false positive when noise is high. Though p-value decreases with or without pseudo-bulking, we observed zero false positive cyclers when pseudobulked data is used, suggesting that pseudo-bulking is robust to elevating expression at a single time point (Figure 5A).

We next examined the effect of elevating the expression of all genes including those considered cycling using the DN1p_CNMa cells as an example. In general, we observed that the number of detected cyclers increases as a function of the magnitude of contamination for methods that consider each cell as a single replicate (JTK-cycle and harmonic regression) (Figure 5C). As expected, the number of detected cyclers did not increase indefinitely for JTK-cycle as it only depends on the rank order of gene expression, which eventually stops changing and does not impact the outcome of the cycling detection algorithm any further. For methods where pseudo-bulk data was used as input, we observed a decreasing number of detected cyclers as a function of contamination, where ARSER and

harmonic regression experienced the greatest amount of loss of detected cyclers (Figure 5 A). Closer inspection revealed interesting connection between the methods. When looking at how the number of total detected cycling genes changed as a function of noise magnitude and the affected time-point, we observed that JTK-cycle and RAIN, both using the Jonckheere-Terpstra test, behaved similarly when applied to pseudo-bulked data (Figure S7, S8). Similar observations were also made when harmonic regression and ARSER were applied to pseudo-bulked data (Figure S9, S10) and when harmonic regression and JTK-cycle were applied considering each cell as a replicate (Figure S11, S12).

In total, we observe that considering cells as replicates can give rise to false positive cyclers in the presence of strong experimental artifacts. In contrast, pseudo-bulking affords protection from this effect, but at the cost of detecting fewer true positives.

4 Discussion

We evaluated the performance of four circadian detection algorithms on single cell transcriptomic timeseries data, with and without pseudobulking. Our analysis indicates that existing cycling detection algorithms scale poorly with replicate size, and showed low subsample concordance in general.

While the problem of long run times can be circumvented by pseudo-bulking, we find that cycling detection on pseudo-bulked data identifies fewer cycling genes in comparison to treating each cell as a replicate. To test whether this increase in the number of cyclers is meaningful, we proposed and conducted subsample concordance analysis by splitting the data into two subsamples to see how the number of cycling genes detected in the full data compares to that of the intersection of the two subsamples. By conducting this analysis on four major clock neuron clusters, we observed that pseudobulking significantly reduced $\frac{N_{\Box}}{N_{full}}$, with greater N_{\Box} when using cells-as-replicates. This suggests that the cyclers detected after pseudobulking may be less reproducible (more differences between sub-experiments) than those detected when treating cells as replicates.

Our analysis also points to a simple method to address the problem of long run times while treating cell as replicates to maintain high $\frac{N_{\cap}}{N_{full}}$: specifically, using the fast (but suboptimal) harmonic regression on two randomly selected subsets of cells, and then considering the overlap of the results. Considering the overlap of the subsets mitigates the false–positive issues that affect harmonic regression. In Fig 3, we demonstrate that using this overlap yields the same genes as would be detected via other methods.

However, care should still be taken when treating cells as replicates. In single cell differential expression analysis, it has been shown that treating cells as replicate can lead to false positives. To investigate whether the same happens in circadian detection algorithms, we artificially contaminated our data by elevating the expression at a single time point, mimicking what might happen if a single sample in a circadian time–series was affected by a systematic artefact. We then compared how pseudobulking impacts the outcome of commonly used circadian detection algorithms. When all genes (both cycling and non-cycling in the original data) were contaminated by a systematic offset of one time–point, we found that the number of detected cyclers increased with the amount of contamination when cells were considered as replicates, suggesting that treating cells as replicates may result in false positives due to experimental artefacts. On the other hand, the number of cycling genes decreased with the amount of contamination when pseudobulk data was used, suggesting that pseudobulking may provide protection from false detection of cyclers, albeit at the expense of detecting fewer true cyclers.

Finally, our analysis highlights the importance of computing subsample concordance. By looking at the agreement between detected cyclers from two subsamples, we noted poor subsample concordance from all tested methods. This provides an opportunity for designing and evaluating new methods for cycling detection designed specifically for single cell datasets.

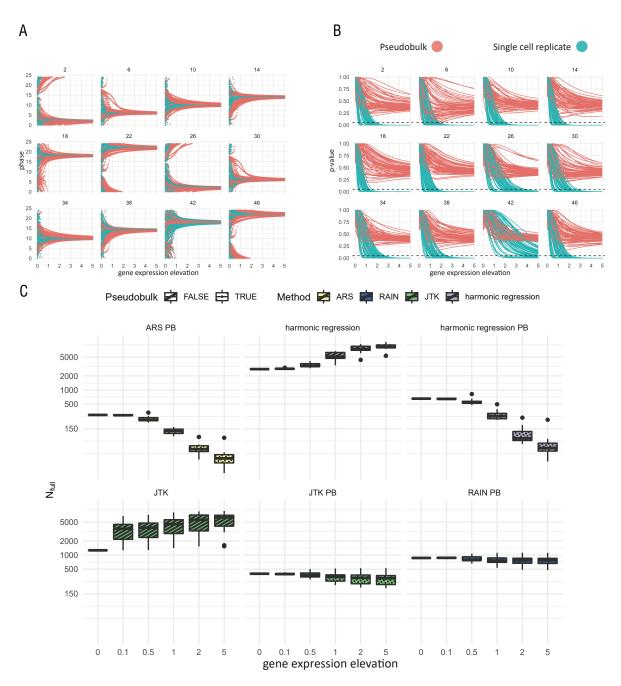


Figure 5: A: Estimated phase of each non-cyclic gene as a function of gene expression elevation. Number of each subplot indicates the timepoint which was subject to gene expression elevation. B: p-value of each gene as a function of gene expression elevation. The dashed line indicates p = 0.05. C: The number of cyclers detected at each level of gene expression elevation.

5 Acknowledgements

This work was supported by NSF grant DMS-1764421, Simons Foundation grant 597491, and NIH grant R01AG068579.

The authors thank Dingbang Ma, Katherine Abruzzi, and Michael Rosbash for processed (normalized counts) single-cell data from [23] as well as for helpful discussions.

References

- [1] Francis Levi and Ueli Schibler. Circadian Rhythms: Mechanisms and Therapeutic Implications. *Annual Review of Pharmacology and Toxicology*, 47(1):593–628, February 2007.
- [2] Michael H Hastings and Michel Goedert. Circadian clocks and neurodegenerative diseases: time to aggregate? Current Opinion in Neurobiology, 23(5):880–887, October 2013.
- [3] Biliana Marcheva, Kathryn Moynihan Ramsey, Ethan D. Buhr, Yumiko Kobayashi, Hong Su, Caroline H. Ko, Ganka Ivanova, Chiaki Omura, Shelley Mo, Martha H. Vitaterna, James P. Lopez, Louis H. Philipson, Christopher A. Bradfield, Seth D. Crosby, Lellean JeBailey, Xiaozhong Wang, Joseph S. Takahashi, and Joseph Bass. Disruption of the clock components CLOCK and BMAL1 leads to hypoinsulinaemia and diabetes. *Nature*, 466(7306):627–631, July 2010.
- [4] Ravi Allada and Brain Y Chung. Circadian organization of behavior and physiology in drosophila. *Annual Review of Physiology*, 72:605–624, 2010.
- [5] Ray Zhang, Nicholas F. Lahens, Heather I. Ballance, Michael E. Hughes, and John B. Hogenesch. A circadian gene expression atlas in mammals: Implications for biology and medicine. *Proceedings of the National Academy of Sciences*, 111(45):16219–16224, November 2014.
- [6] Yitong Huang, Yuanzhao Zhang, and Rosemary Braun. A minimal model of peripheral clocks reveals differential circadian re-entrainment in aging. Chaos: An Interdisciplinary Journal of Nonlinear Science, 33(9), 2023.
- [7] Anna A. Kondratova and Roman V. Kondratov. The circadian clock and pathology of the ageing brain. *Nature Reviews Neuroscience*, 13(5):325–335, May 2012.
- [8] Christopher A. Wolff, Miguel A. Gutierrez-Monreal, Lingsong Meng, Xiping Zhang, Lauren G. Douma, Hannah M. Costello, Collin M. Douglas, Elnaz Ebrahimi, Bryan R. Alava, Andrew R. Morris, Mehari M. Endale, G. Ryan Crislip, Kit-yan Cheng, Elizabeth A. Schroder, Brian P. Delisle, Andrew J. Bryant, Michelle L. Gumz, Zhiguang Huo, Andrew C. Liu, and Karyn A. Esser. Defining the age-dependent and tissue-specific circadian transcriptome in male mice. preprint, Genomics, April 2022.
- [9] Elan Ness-Cohn, Marta Iwanaszko, William L. Kath, Ravi Allada, and Rosemary Braun. Time-Trial: An Interactive Application for Optimizing the Design and Analysis of Transcriptomic Time-Series Data in Circadian Biology Research. *Journal of Biological Rhythms*, 35(5):439–451, October 2020.
- [10] Jordan M. Singer and Jacob J. Hughey. LimoRhyde: A Flexible Approach for Differential Analysis of Rhythmic Transcriptome Data. *Journal of Biological Rhythms*, 34(1):5–18, February 2019.
- [11] Jake Hughey, Dora Obodo, and Elliot Outland. limorhyde2: Quantify Rhythmicity and Differential Rhythmicity in Genomic Data, 2023. https://limorhyde2.hugheylab.org, https://github.com/hugheylab/limorhyde2.
- [12] Elan Ness-Cohn and Rosemary Braun. TimeCycle: topology inspired method for the detection of cycling transcripts in circadian time-series data. *Bioinformatics*, 37(23):4405–4413, December 2021.
- [13] R. Yang and Z. Su. Analyzing circadian expression data by harmonic regression based on autoregressive spectral estimation. *Bioinformatics*, 26(12):i168–i174, June 2010.

- [14] Joel D Levine, Pablo Funes, Harold B Dowse, and Jeffrey C Hall. Signal analysis of behavioral and molecular cycles. *BMC Neuroscience*, 3(1):1, 2002.
- [15] Martin Straume. DNA Microarray Time Series Analysis: Automated Statistical Assessment of Circadian Rhythms in Gene Expression Patterning. In *Methods in Enzymology*, volume 383, pages 149–166. Elsevier, 2004.
- [16] S. Wichert, K. Fokianos, and K. Strimmer. Identifying periodically expressed transcripts in microarray time series data. *Bioinformatics*, 20(1):5–20, January 2004.
- [17] Michael E. Hughes, John B. Hogenesch, and Karl Kornacker. JTK_cycle: An Efficient Nonparametric Algorithm for Detecting Rhythmic Components in Genome-Scale Data Sets. *Journal of Biological Rhythms*, 25(5):372–380, October 2010.
- [18] Paul F. Thaben and Pål O. Westermark. Detecting Rhythms in Time Series with RAIN. *Journal of Biological Rhythms*, 29(6):391–400, December 2014.
- [19] David Laloum and Marc Robinson-Rechavi. Methods detecting rhythmic gene expression are biologically relevant only for strong signal. PLOS Computational Biology, 16(3):e1007666, March 2020.
- [20] Wenwen Mei, Zhiwen Jiang, Yang Chen, Li Chen, Aziz Sancar, and Yuchao Jiang. Genome-wide circadian rhythm detection methods: systematic evaluations and practical guidelines. *Briefings in Bioinformatics*, 22(3):bbaa135, May 2021.
- [21] Jose A Perea, Anastasia Deckard, Steve B Haase, and John Harer. Sw1pers: Sliding windows and 1-persistence scoring; discovering periodicity in gene expression time series data. *BMC bioinformatics*, 16(1):1–12, 2015.
- [22] 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, and Grégoire Courtine. Confronting false discoveries in single-cell differential expression. *Nature Communications*, 12(1):5692, September 2021.
- [23] Dingbang Ma, Dariusz Przybylski, Katharine C Abruzzi, Matthias Schlichting, Qunlong Li, Xi Long, and Michael Rosbash. A transcriptomic taxonomy of Drosophila circadian neurons around the clock. *eLife*, 10:e63056, January 2021.
- [24] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi, William M Mauck III, Yuhan Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija. Comprehensive integration of single-cell data. Cell, 177:1888–1902, 2019.
- [25] Yuhan Hao, Stephanie Hao, Erica Andersen-Nissen, William M. Mauck III, Shiwei Zheng, Andrew Butler, Maddie J. Lee, Aaron J. Wilk, Charlotte Darby, Michael Zagar, Paul Hoffman, Marlon Stoeckius, Efthymia Papalexi, Eleni P. Mimitou, Jaison Jain, Avi Srivastava, Tim Stuart, Lamar B. Fleming, Bertrand Yeung, Angela J. Rogers, Juliana M. McElrath, Catherine A. Blish, Raphael Gottardo, Peter Smibert, and Rahul Satija. Integrated analysis of multimodal single-cell data. Cell, 2021.
- [26] Gang Wu, Ron C. Anafi, Michael E. Hughes, Karl Kornacker, and John B. Hogenesch. MetaCycle: an integrated R package to evaluate periodicity in large scale data. *Bioinformatics*, 32(21):3351–3353, November 2016.
- [27] Sarah Lueck, Kevin Thurley, Paul F. Thaben, and Paul O. Westermark. Rhythmic degradation explains and unifies circadian transcriptome and proteome data. *Cell Reports*, 9:741–751, 2014.
- [28] Sandipan Ray, Utham K. Valekunja, Alessandra Stangherlin, Steven A. Howell, Ambrosius P. Snijders, Gopinath Damodaran, and Akhilesh B. Reddy. Circadian rhythms in the absence of the clock gene *Bmal1*. *Science*, 367(6479):800–806, February 2020.
- [29] Shao'ang Wen, Danyi Ma, Meng Zhao, Lucheng Xie, Qingqin Wu, Lingfeng Gou, Chuanzhen Zhu, Yuqi Fan, Haifang Wang, and Jun Yan. Spatiotemporal single-cell analysis of gene expression in the mouse suprachiasmatic nucleus. *Nature Neuroscience*, 23(3):456–467, March 2020.

6 Supplemental material

6.1 Data preprocessing and clustering

Preprocessed data and cell type information was contained within the Seurat object provided by Ma et al [23]. Single cell data from six time points (two experimental replicates, hence twelve in total) were processed by Ma et al., who integrated the data on a per time point basis using Seurat with genes that identified by Seurat as highly variable across all time points. Integrated data were then reduced to 2D using t-SNE, from which cluster identities were assigned. For all of our analysis, we used normalized counts kindly shared my Ma et al. to ensure correspondence with results reported in [23].

6.2 Cycling detection

Given the two experimental replicates from Ma et al., were collected consecutively, the data were concatenated to a time series with twelve time points. Consequently, the pseudobulked time series also consist of twelve data points, and similarly for the subsampled data (constructed such that each subexperiment has at least one sample from each one of the twelve time points).

We applied four cycling detection methods:

- JTK-cycle is implemented using JTK_CYCLEv3.1.R. The period was set to 6 and the sampling interval was set to 4 to match the experimental procedure. When considering single cells as replicates, we also supplied the number of samples for each time point.
- Harmonic regression is implemented using the harmonicRegression package in R [27] inputting only time and the expression matrix.
- RAIN is implemented in the R package [18]. It takes as input the expression matrix ordered by sampling time and the sampling interval *deltat*, set to 4 to match the experimental design.
- ARSER is implemented using the metacycle package in R [26], whose only required input is the expression matrix and sampling time.

For the purpose of our analysis, a gene is considered to be cycling if its raw p-value is less than 0.05. We also repeat the analysis using FDR-adjusted p-values.

6.3 Derivation of theoretical computational complexity

6.3.1 Mann-Whitney U Statistics

Both JTK cycle and RAIN are non-parametric methods that are built upon the Mann-Whitney U statistics. Let $(X_{11}, \ldots, X_{1m}), \ldots (X_{T1}, \ldots, X_{Tm})$ be a set of T time-associated observations following probability distributions $P_1(x), \ldots, P_T(x)$ for timepoints 1-T. For the purpose of deriving computational complexity, we have assumed that each time point contains the same number of cells m without loss of generality. The Mann-Whitney test tests the null hypothesis that $P_i = P_j (i \neq j)$ by computing the U statistic between sample i and sample j as

$$U_{i,j} = \sum_{k=1}^{m} \sum_{l=1}^{m} \mathbb{I}(X_{ik} < X_{jl})$$
(3)

It is easy enough to see that $U_{i,j}$ should have a computational complexity of $\mathcal{O}(m^2)$ and therefore the computational complexity of $\mathbf{U} = (U_{1,2}, \dots, U_{T-1,T})$ should be approximately $\mathcal{O}(T^2 \times m^2)$. Since in real datasets $T \ll m$ and the number of time points do not change from one cluster to another, our expected observed computational complexity, when changing sample size, should also be $\mathcal{O}(m^2)$.

6.3.2 Jonckheere-Terpstra Test

JTK cycle employs the Jonckheere-Terpstra test, which is an extension of the Mann-Whitney U statistics to the case of having more than two samples. It tests for the presence of monotonic trend with the test statistics s, defined as

$$s = \sum_{i=1}^{T-1} \sum_{j=i+1}^{T} U_{i,j} \tag{4}$$

Again, it is clear that the computational complexity for s should be $\mathcal{O}(m^2 + T^2)$. When the number of time points is fixed, the computational complexity scales quadratically with sample size.

In JTK cycle, the rising and falling part of the oscillation is tested for each a pre-determined waveform (e.g., sine waves of differing phases), resulting in a complexity of $\mathcal{O}(km^2)$, where k is the number of templates. As k is fixed and $k \ll m$, the complexity will scale approximately $\mathcal{O}(m^2)$.

6.3.3 General Umbrella

To overcome the JTK-cycle's loss of power when considering multiple waveforms, RAIN uses a variation of the general umbrella which tests for the presence of an umbrella shape, which can be formally written as

$$H1: P_1(x) < P_2(x) < \dots P_t(x) > \dots > P_T(x)$$
 (5)

Here, t is a predetermined inflection point. The test statistics of the general umbrella is the sum of two Jonckheere-Terpstra statistics

$$s = \sum_{i=1}^{t-1} \sum_{j=i+1}^{t} U_{i,j} + \sum_{i=t}^{T-1} \sum_{j=i+1}^{T} U_{j,i}$$
(6)

RAIN tests for all potential inflection points T, yielding a computational complexity $\mathcal{O}(Tm^2)$. Since T is fixed by the experimental design, the complexity scales with the squared number of cells.

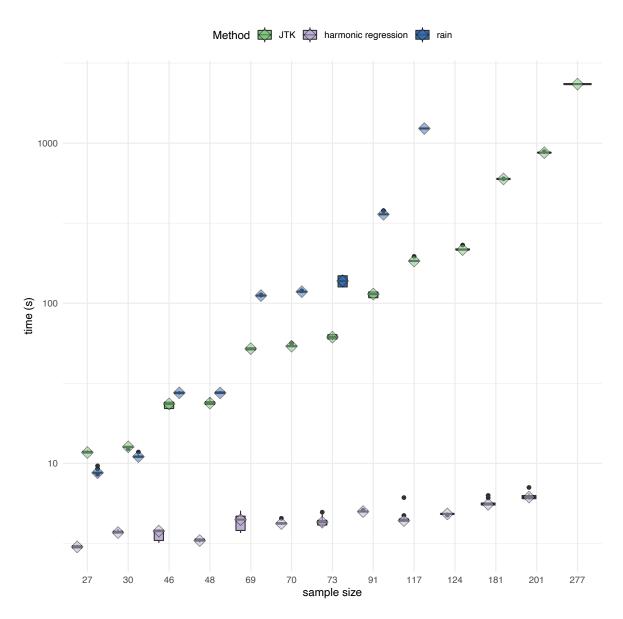


Figure S1: Run time for JTK cycle, RAIN and harmonic regression on multiple cell types.

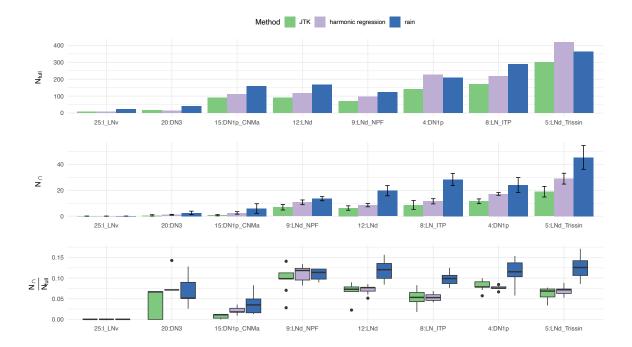


Figure S2: Top: The number of cycling genes detected using all cells, $N_{\rm full}$. Middle: The average number of cyclers detected in both subsamples across 10 trials, N_{\cap} . Bottom: The ratio of N_{\cap} to N_{full} , the proportion of consistently–detected sub-sample cyclers relative to those found using all cells. Error bars indicate standard deviation across the 10 subsamplings. A gene is considered cycling if its BH corrected p-value (i.e., FDR) is less than 0.05.

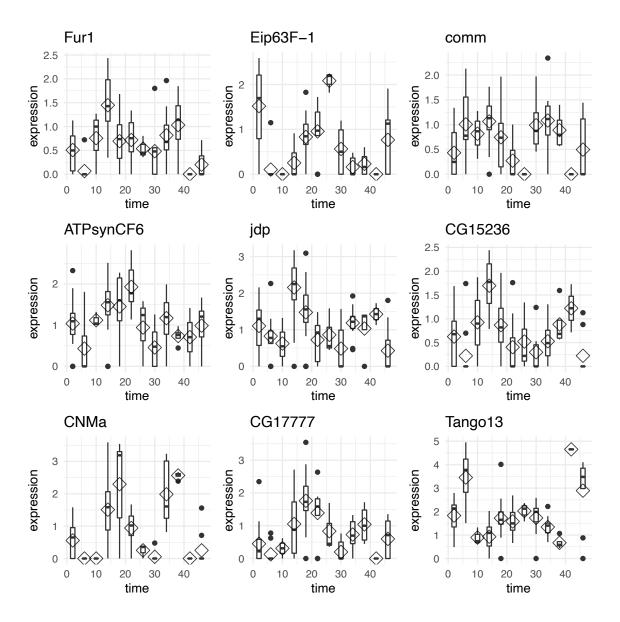


Figure S3: Example genes considered cycling by JTK-cycle when all cells were used as replicates but not when pseudobulk expression was used.

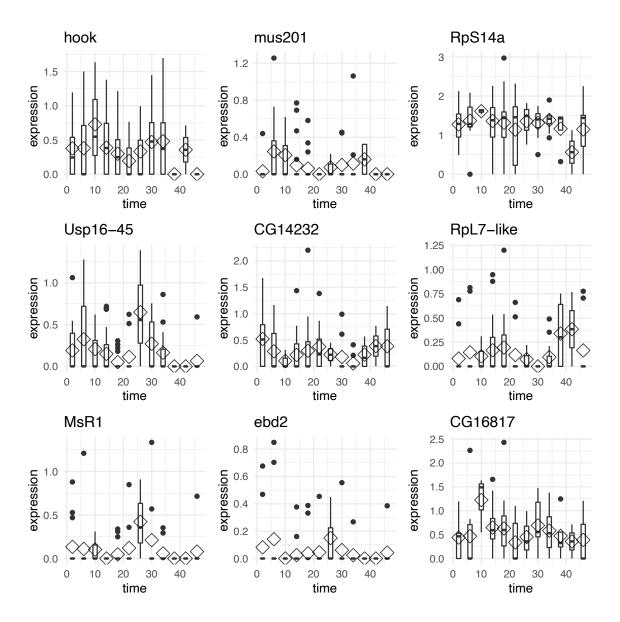


Figure S4: Example genes considered cycling by JTK cycle using pseudobulk expression but not when cells were treated as replicates

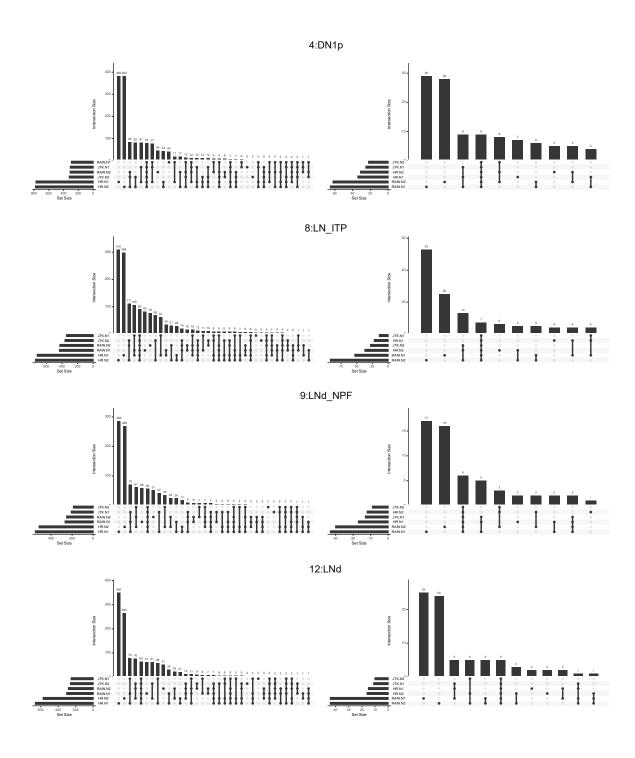


Figure S5: Upset plot of cluster 4, 8, 9, and 12. The left (right) column uses p-values (adjusted p-values) to identify cycling genes.

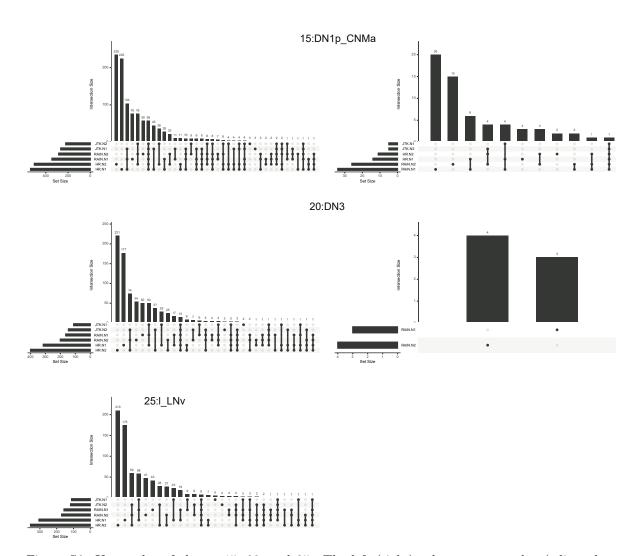


Figure S6: Upset plot of cluster 15, 20, and 25. The left (right) column uses p-value (adjusted p-values) to identify cycling genes. Right column for cluster 25 is empty for that the average intersection between all sets were less than one.

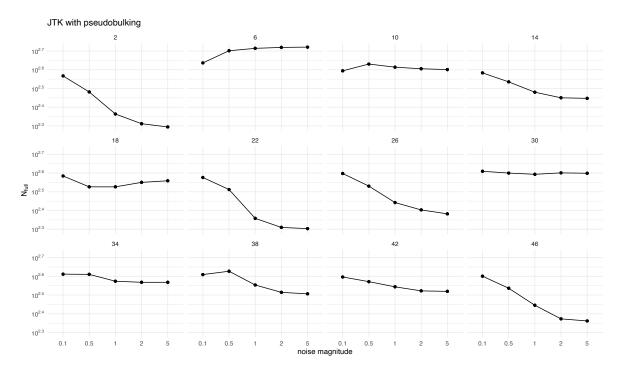


Figure S7: Effect of systematic contamination when using JTK-cycle with pseudobulking. Each panel shows how the number of detected cyclers $N_{\rm full}$ varies as function of the magnitude of noise injected at a single time-point.

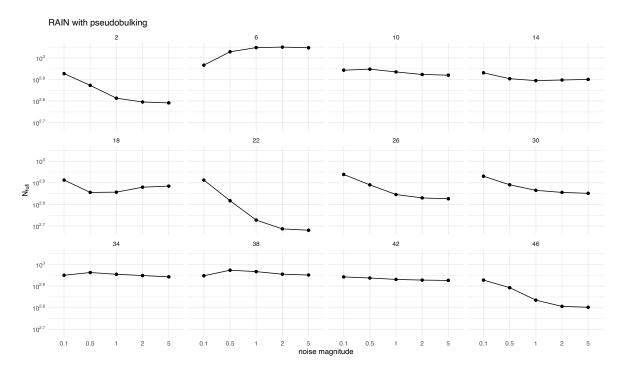


Figure S8: Effect of systematic contamination when using RAIN with pseudobulking. Each panel shows how the number of detected cyclers $N_{\rm full}$ varies as function of the magnitude of noise injected at a single time-point.

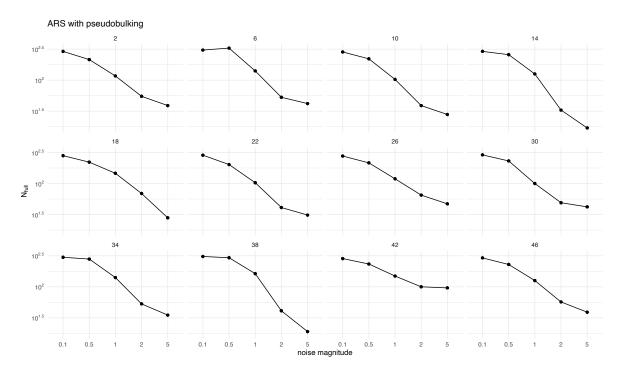


Figure S9: Effect of systematic contamination when using ARSER with pseudobulking. Each panel shows how the number of detected cyclers $N_{\rm full}$ varies as function of the magnitude of noise injected at a single time-point. The time-point contaminated by noise is given at the top of each panel.

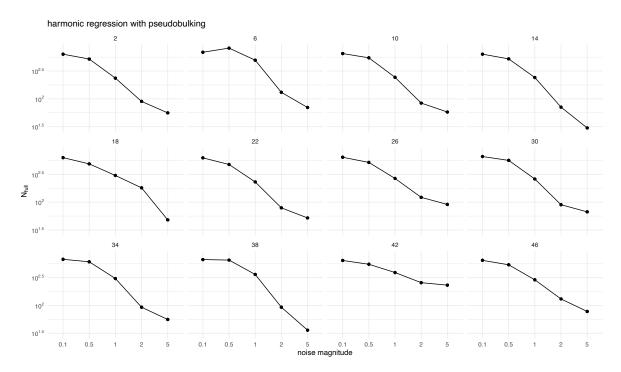


Figure S10: Effect of systematic contamination when using harmonic regression with pseudobulking. Each panel shows how the number of detected cyclers $N_{\rm full}$ varies as function of the magnitude of noise injected at a single time-point. The time-point contaminated by noise is given at the top of each panel.

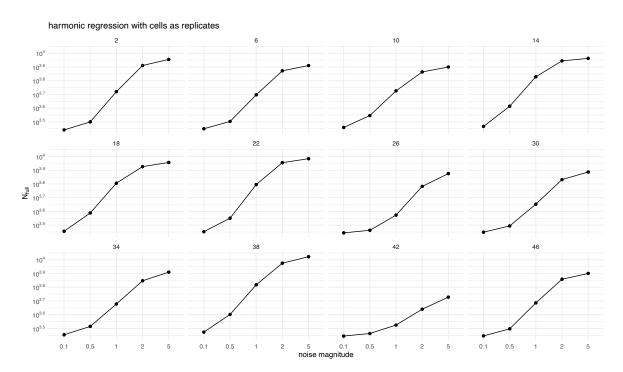


Figure S11: Effect of systematic contamination when using harmonic regression with cells as replicates. Each panel shows how the number of detected cyclers $N_{\rm full}$ varies as function of the magnitude of noise injected at a single time-point. The time-point contaminated by noise is given at the top of each panel.

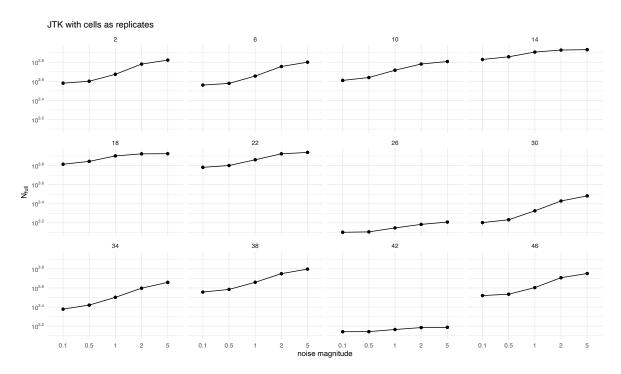


Figure S12: Effect of systematic contamination when using JTK-cycle with cells as replicates. Each panel shows how the number of detected cyclers $N_{\rm full}$ varies as function of the magnitude of noise injected at a single time-point. The time-point contaminated by noise is given at the top of each panel.