

1 **FRONT MATTER**

2

3 **Title**

- 4 • **noisyR**: Enhancing biological signal in sequencing datasets by characterising random
5 technical noise
6 • Noise removal unveils robust biological signal

7 **Authors**

8 I. Moutsopoulos,¹ L. Maischak,² E. Lauzikaite,¹ S. A. Vasquez Urbina,² E. C. Williams,¹
9 H. G. Drost,² I. I. Mohorianu^{1*}

10 **Affiliations**

11 ¹ Wellcome-MRC Cambridge Stem Cell Institute, Jeffrey Cheah Biomedical Centre,
12 Cambridge Biomedical Campus, University of Cambridge, CB2 0AW, UK.

13 ² Computational Biology Group, Department of Molecular Biology, Max Planck Institute
14 for Developmental Biology, Max-Planck Ring 1, 72076 Tübingen, Germany.

15 * Corresponding author: I. I. Mohorianu, iim22@cam.ac.uk

16 **Abstract**

17 High-throughput sequencing enables an unprecedented resolution in transcript
18 quantification, at the cost of magnifying the impact of technical noise. The consistent
19 reduction of random background noise to capture functionally meaningful biological signals
20 is still challenging. Intrinsic sequencing variability introducing low-level expression
21 variations can obscure patterns in downstream analyses.

22 We introduce **noisyR**, a comprehensive noise filter to assess the variation in signal
23 distribution and achieve an optimal information-consistency across replicates and samples;
24 this selection also facilitates meaningful pattern recognition outside the background-noise
25 range. **noisyR** is applicable to count matrices and sequencing data; it outputs sample-
26 specific signal/noise thresholds and filtered expression matrices.

27 We exemplify the effects of minimising technical noise on several datasets, across various
28 sequencing assays: coding, non-coding RNAs and interactions, at bulk and single-cell level.
29 An immediate consequence of filtering out noise is the convergence of predictions
30 (differential-expression calls, enrichment analyses and inference of gene regulatory
31 networks) across different approaches.

32 Keywords: next generation sequencing, noise, bulk sequencing, single-cell sequencing,
33 count matrix, expression profile, differential expression, enrichment analysis, gene
34 regulatory network

35 **Teaser**

36 Noise removal from sequencing quantification improves the convergence of downstream
37 tools and robustness of conclusions.
38

39 MAIN TEXT

40

41 Introduction

42 High-throughput sequencing (HTS) became a new standard in most life science studies
43 yielding unprecedented insights into the complexity of biological processes. The increase
44 in sequencing depth and number of samples, across both bulk and single cell experiments,
45 facilitated a greater diversity in biological questions (1), at the same time allowing a higher
46 sensitivity for the detection of perturbations in gene expression levels between samples (2).
47 This increased accuracy greatly assists with the biological interpretation of results such as
48 identification and characterisation of differential expression (DE) at tissue and cellular
49 levels (3) or the inference and characterisation of gene regulatory networks (4). However,
50 HTS may exhibit high background noise levels resulting from non-biological/technical
51 variation, introduced at different stages of the RNA-seq library preparation, or from
52 amplification/sequencing bias (5) to random hexamer priming during the sequencing
53 reaction (6). These technical alterations of signal can affect the accuracy of the downstream
54 DE results or create spurious patterns biasing downstream interpretations. Statistical
55 methods developed to date (7-9), focused mainly on batch/background correction,
56 normalisation, and evaluation of DE, have been designed to mitigate the impact of these
57 biases on DE analyses (10). A noise filter for pre-processing of the data before these steps
58 would ensure a reduction of further amplification of these biases. Here, we introduce a new
59 high-throughput noise filter to remove random technical noise from sequencing data and
60 illustrate the downstream information consistency that is achieved.

61 While technologies may exhibit different technical biases, the sequencing bias across an
62 experiment was expected to be uniform. This expectation was based on the assumption that
63 sequencing reads would uniformly cover the expressed transcripts, with the algebraic sum
64 of reads from each gene being proportional to the expression of that gene (11). However, in
65 practice we observe a reproducible, yet uneven distribution of signal across transcripts (11);
66 moreover, highly abundant genes show a higher consistency of transcript-coverage than
67 lower abundance genes. This coverage bias of lower abundance genes is one of the main
68 origins of technical noise (12). The latter can be attributed to the stochasticity of the
69 sequencing process, the limits of sequencing depth, and alignment inaccuracies during the
70 mapping procedure. To further explore the coverage bias of lower abundance genes, we
71 define genes whose quantification is characterised by such a lack of coverage-uniformity as
72 “noisy”.

73 The presence of noise in HTS data has been widely acknowledged, and there have been
74 several attempts to understand and quantify it. A recent study (13) presented a variety of
75 common experimental errors that may increase sequencing noise and proposed ways to
76 alleviate their effect such as using a mild acoustic shearing condition to minimise the
77 occurrence of DNA damage. Fischer-Hwang and colleagues (14) presented a denoising tool
78 that can be applied on aligned genomic data with high fold-coverage of the genome to
79 improve variant calling performance. The recent prevalence of single-cell sequencing
80 technologies has further highlighted the issue of noise, as the lower sequencing depth per
81 cell leads to more uncertainty of the quantification of (low abundance) genes. Efforts have
82 been made to reduce the noise levels experimentally, such as by utilising a different
83 barcoding approach (15).

84 On the computational side, several imputation and denoising algorithms have been
85 proposed, e.g. a machine learning (ML) based deep count autoencoder (16). Other tools
86 focus on DE analysis, such as TASC (17), which uses a hierarchical mixture model of the
87 biological variation. However, successful methods usually rely on assumptions specific to
88 the biological experiment and are tailored to particular settings or model systems, thus
89 leaving most large-scale sequencing efforts, lacking such specific experimental design,
90 exposed to random technical noise. To our knowledge, there is little focus on bulk
91 experiments, where technical noise still exists at low abundances, independent of biological
92 assumptions; for these experiments the low number of replicates hinders imputation-based
93 approaches.

94 Existing approaches for calling DE genes mitigate to various extents the presence of noise,
95 however these are not designed to identify and assess the impact of genes showing random,
96 low-level variation. As a result, some of these are detected by the DE analyses, biasing the
97 biological interpretation of the results. In addition, the choice of tools used for pre-
98 processing steps may influence the relative transcript expression estimation accuracy (18).
99 These analytical biases mainly arise from differences in the detection and handling of
100 transcript isoforms or processing of unmapped and multi-mapping reads (3). Such variation
101 in abundance estimation can in turn strongly affect the downstream analyses (19).

102 We developed **noisyR**, a denoising pipeline to quantify and exclude technical noise from
103 downstream analyses, in a robust and data-driven way. The approach underlines consistency
104 of signal over a user-defined threshold. **noisyR** is applicable on either the original, un-
105 normalised count matrix, or alignment data (BAM format). Noise is quantified based on the
106 correlation of expression across subsets of genes for the former, or distribution of signal
107 across the transcripts for the latter, in different samples/replicates and across all gene
108 abundances (Methods). We illustrate the approach on bulk and single cell RNA-seq datasets
109 and highlight the impact of the noise removal on refining the biological interpretation of
110 results.

111 Results

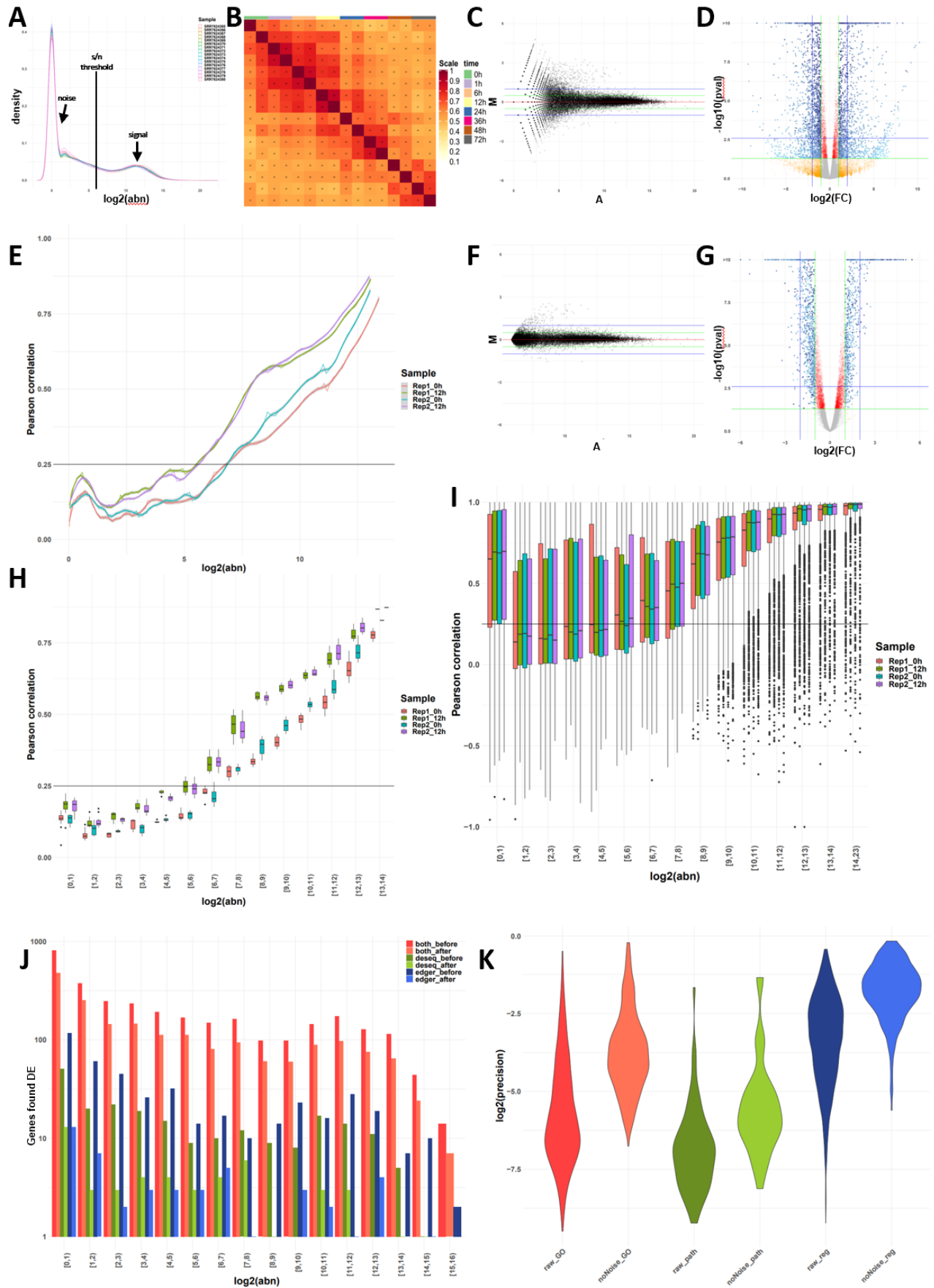
112 Noise quantification in bulk RNA-seq data

113 To exemplify the impact of denoising on the biological interpretations from bulk RNA-seq
114 experiments, we applied **noisyR** on mRNA-seq and smallRNA-seq data. First, we illustrated
115 the advantages of using the pipeline on a subset of mRNA-seq samples from a 2019 study
116 by Yang and colleagues (20). To assess the distributions of signal we used density plots
117 (Fig. 1A) and summaries of Jaccard similarity indices (Fig. 1B) across all samples. For the
118 former, we observed a multi-modal distribution that suggests a signal to noise transition
119 range between [3,7] on log₂ scale; for the latter, the high similarity along the diagonal
120 mirrors the temporal component of the time series. To reduce the number of low abundance,
121 high fold change DE calls (Fig. 1C, fig.S1A for sample similarity and the secondary DE
122 distribution visible in Fig. 1D and fig. S1C), we used first the **noisyR** count-based pipeline,
123 on default parameters: window length = 10% x #genes and sliding step = 5% x window
124 length (Fig. 1, E and H, fig. S1E). We used a correlation threshold of 0.25 and the boxplot
125 median method, a combination of hyper-parameters producing the smallest coefficient of
126 variation across abundance thresholds for the considered samples (Methods); the
127 interquartile ranges (IQRs) of noise thresholds for the different samples ranged between 39
128 and 63, with an average of 58, for sequencing depths varying between 58M and 82M (Fig.
129 1E, fig. S1E). We detected an outlier with a low threshold of 18 (corresponding to a

sequencing depth of ~77M) and three with values of over 100, corresponding to sequencing depths of 73M, 71M and 96M respectively. Next, we applied the transcript approach focusing on the correlation of the expression profiles across exons/transcripts (Methods); despite the higher runtime compared to the count-based approach, the transcript-approach was more robust, as illustrated by the lower variance in signal/noise thresholds across samples (Fig. 1I). The parameters that minimised the coefficient of variation were: correlation threshold = 0.26 and the boxplot median method; the resulting noise threshold IQRs ranged between 64 and 79, with an average of 75 and one outlier at 104. The signal/noise thresholds were similar for the two options, with an increased level of detail for the transcript-based approach.

These thresholds were used to exclude noisy genes from the count matrix (~44k genes were excluded out of ~56k genes expressed); the number of retained genes were 19.7k and 15.6k for the counts and transcript approaches, respectively. As a DE pre-processing step, the averaged noise threshold was added to all entries in the count matrix (Methods). The effect of the noise removal is illustrated by the narrower distribution in the MA plots (Fig. 1F, fig. S1B). Next, we performed a DE analysis between the 0h and 12h samples of the Yang dataset using the denoised matrix. Following the noise correction, we saw a 46% reduction in the number of DE genes - from 3,607 to 1,952. A large number of low abundance genes with spuriously high fold-changes were no longer called DE (12). Moreover, when comparing the outputs of two standard DE pipelines, edgeR (8) and DESeq2 (7), we noticed that the number of genes identified as DE by both methods only marginally decreased when the noise corrected input is used, whereas the number of DE genes called only with edgeR or only with DeSeq2 decreased significantly (Fig. 1J, fig. S1F); therefore we observed an increase in output consistency across methods when the noise filtered inputs were used. Moreover, the fold-changes and p-values of denoised genes correlated better and we no longer saw a large set of DE genes with (adjusted) p-values marginally below the DE threshold (Fig. 1D vs G. fig. S1C vs D). This step was followed by a functional enrichment analysis focusing on the DE genes, with the genes expressed (post filtering) as background set (21). The number of enriched terms was lower in the denoised data, 1,108 vs 4,671 in the original analysis; ~24% of the terms were retained and the terms found with the denoised dataset were approximately a subset of the ones found without the noise correction (~99.6% of terms found after denoising were also found prior to noise removal). In addition, the noise-correction terms corresponded to a higher percentage of genes assigned per pathway (Fig. 1K). Thus, applying *noisyR* focused the interpretation of results on the enrichment terms with highest confidence, ensuring biological relevance.

The *noisyR* transcript approach was also applied on two small RNA (sRNA) datasets, from plants (*A. thaliana*) and animals (*M. musculus*), respectively. In contrast to the mRNAseq data, sRNAs samples had different correlation vs abundance distributions. Overall low abundance sRNA transcripts/loci contained more noisy entries (22). Also, we observed a sharper increase to high correlation entries highlighting the transition from degraded transcripts to precisely excised sRNAs (23, 24). For both model organisms, miRNA hairpins and transposable elements (TEs) were analysed separately. For the former, we observed overall higher correlations than for mRNAs, likely because of the precise cleavage of the mature duplex and the lack of signal outside the duplex region (25); this characteristic is stronger for the animal case (fig. S2C). For both animals and plants, the increasing distribution was clearly detectable (fig. S2, A and C). The TE distributions also reflected the characteristics of the underlying sRNAs; for the animal example (fig. S2D) we saw a sharper increase along the abundance bins, specific for the piRNAs (26), whereas in plants (fig. S2B), the distribution of signal (expressed siRNAs) mirrored the biogenesis of heterochromatin siRNAs (27).



181 **Fig. 1 Overview of QC measures and original vs denoised outputs on standard components of**
182 **an mRNA-seq pipeline.**

183 (A) Distributions of gene abundances by sample; the RHS distribution corresponds to the biological
184 signal, the LHS distribution to the technical noise; the aim of *noisyR* is the identification of
185 biologically meaningful values for the signal/noise threshold in between. (B) JSI on the 100 most
186 abundant genes per sample; the replicates, and consecutive time points share a larger proportion of
187 abundant genes. (C) MA plot of the raw abundances for the two 12h biological replicates; a larger
188 proportion of low abundance genes exhibit high fold-changes, potentially biasing the DE calls. (D)
189 Volcano plot of differentially expressed genes on the original, normalised count matrix; the colour
190 gradient is proportional to the gene abundance. (E) Line plot of the PCC calculated on windows of
191 increasing average abundance for the count matrix-based noise removal approach. (F) MA plot of
192 the de-noised abundances for the two 12h biological replicates; the low-level variation is
193 significantly reduced. (G) Volcano plot of differentially expressed genes on the denoised count
194 matrix. (H) Box plot of the PCC binned by abundance for the count matrix-based noise removal
195 approach. (I) Box plot of the PCC binned by abundance for the transcript-based noise removal
196 approach. (J) Histogram of the differentially expressed genes found by applying DESeq and edgeR
197 on the original and denoised count matrix respectively, binned by abundance; counts are on a log-
198 scale for visualization. (K) Violin plot of the precision (intersection size divided by the query size)
199 for the results of the enrichment analysis performed on the differentially expressed genes found for
200 the original (*raw*) and denoised (*noNoise*) matrices (log-scale). In the Gene Ontology set (*GO*) the
201 terms from Biological Processed, Cellular Component and Molecular Function were grouped; in
202 the Pathway set (*path*) the *Kegg* and *Reactome* terms were grouped; in the Regulatory terms (*reg*)
203 the enriched Transcription Factors and microRNA entries were grouped.

204 Effect of noise on single cell (smartSeq) data

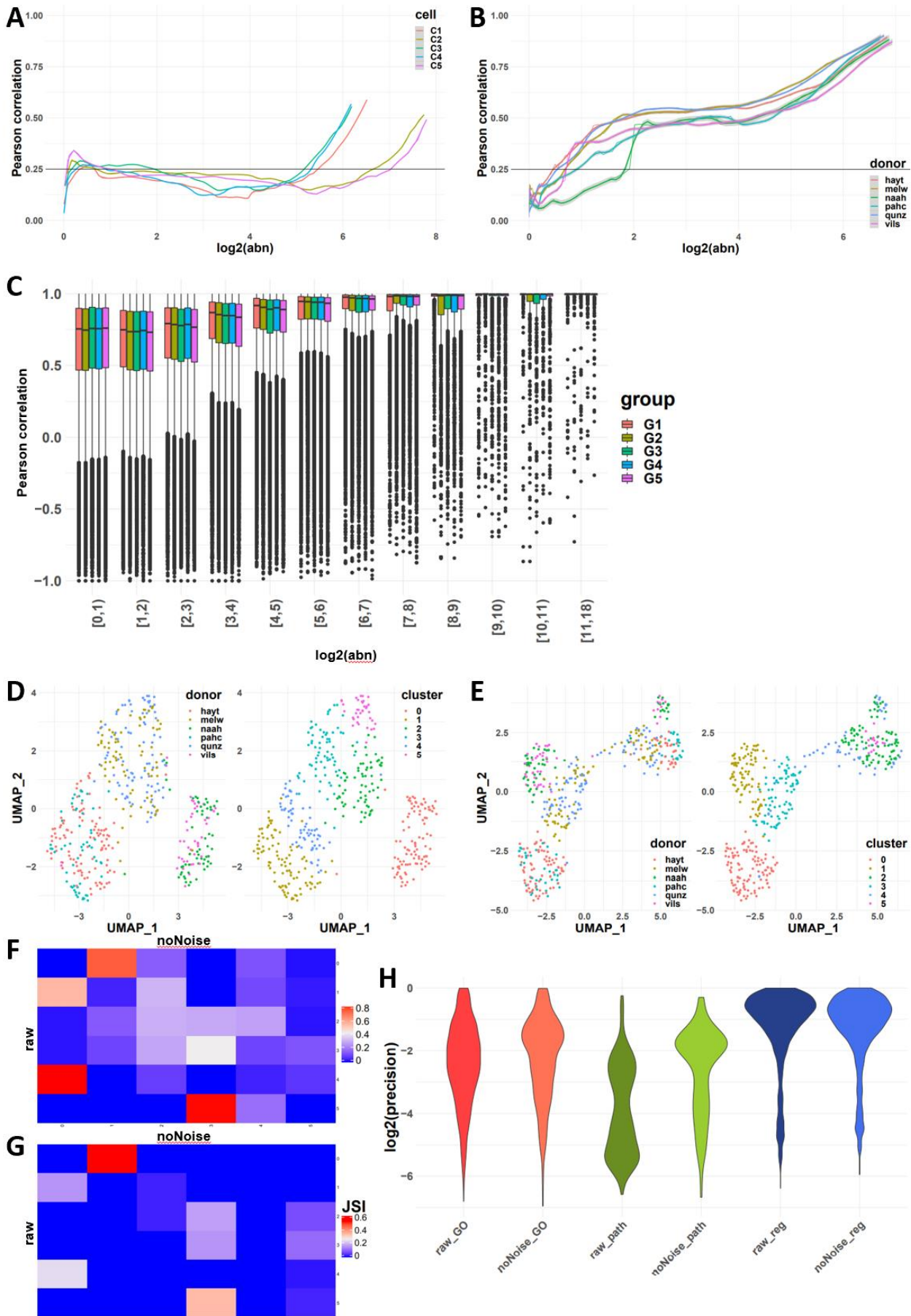
205 To illustrate the broad applicability of *noisyR* on different HTS data, we present its output
206 on single cell (smartSeq2) sequencing output focusing on a subset of samples from the
207 dataset presented by Cuomo and colleagues (28); we focused on 6 donors, and one time-
208 point, the number of cells per donor varied between 45 and 107. A common difficulty in
209 single-cell experiments is that due to the higher number of samples/cells, the runtime is
210 much higher if the pipeline is applied without modification, making the transcript approach
211 intractable in practice, for higher number of cells; we also assessed whether the inferred
212 signal/noise threshold was informative.

213 First, we applied *noisyR* using the count matrix approach on all cells with default
214 parameters; we observed that correlation values rose to a weakly positive plateau (0.2-0.4)
215 and remained stable over a wide range of abundances (Fig. 2A). Our interpretation of this
216 result is that lower sequencing depths and higher resolution of smart-seq compared to bulk
217 data induces more dissimilarity across medium abundance values. To alleviate this effect,
218 we grouped cells into a small number of “pseudo-samples”, both randomly and according
219 to the sample origin (i.e. donor). For each pseudo-sample, we applied the count-based
220 approach on the averaged expression of genes. In the resulting *noisyR* output, we observed
221 a clearer step in the abundance-correlation plot (Fig. 2B, Fig. S3A), especially when the
222 cells were grouped by donor. This indicates that an effect of the summarisation is a reduction
223 in cell-to-cell variability which also focuses the noise identification procedure. The
224 thresholds obtained via pseudo-sample summarisation and count-based noise identification
225 varied between 2 and 4 with an average of 2.6 (corresponding to a sequencing depth per
226 pseudo-sample between 590K and 689K, representative of the average sequencing depth
227 per cell of 640K); these were used in a similar manner as for the bulk data, to produce a
228 denoised count matrix.

229 As the transcript approach is more computationally intensive, we applied it on a subsampled
230 set of 25 cells. The subsamples were chosen randomly, and the process was reiterated five
231 times, with the requirement that the summarised cells originate from the same donor.
232 Formatting the data for *noisyR* was achieved by concatenating the BAM files for the
233 selected cells and treating them as one sample. Whereas for the count approach, the results
234 were highly variable between the cells, with several instances of low or negative correlations
235 observed even at high abundances (Fig. 2A), the results obtained using the transcript
236 approach with the concatenated BAM files were more consistent, with an expected
237 increasing trend in the distribution of correlations (Fig. 2C). The correlation distributions
238 were high, even at low abundances, which may be a consequence of the summarisation; a
239 suitable threshold may be selected on the median, IQR, or 5-95% range to infer a signal to
240 noise threshold, as the distributions are stable for low values and increase as the abundance
241 increases above ~ 2 on a \log_2 -scale.

242 To assess the impact of *noisyR* on the biological interpretation of results, we performed
243 some downstream analyses before and after the noise removal and compared the results. In
244 this study, we focus on the structure and mathematical characteristics of the outputs, rather
245 than specific biological interpretations. The gene abundances were normalised and the cells
246 were clustered using the Seurat R package (29) (see Methods). The different clusterings
247 were visualised using the UMAP (non-linear) dimensionality reduction (30) (Fig. 2, D and
248 E, fig. S3, B and C). We observed that cells clustered into three groups of two donors each
249 when original data was used, suggesting a batch effect; However, cells corresponding to
250 the four donors are mixed across clusters, when the denoised data was analysed, suggesting

251 that some of the putative initial batch effect may have been alleviated with the noise
252 correction. We also observed a better separation of clusters in the denoised data, especially
253 on the first UMAP component, which may be an indication of robustness. We further
254 assessed the similarity of the two clustering results using a cell-centred contingency table
255 (Fig. 2F). We observed a good correspondence between the original and denoised matrices
256 i.e. clusters 1 and 4 largely merged into cluster 0, and cluster 0 remains intact and turns into
257 cluster 1. While the total number of clusters remained the same (under default parameters),
258 the partitioning of cells was altered, which led us to believe that the results obtained with
259 the original and denoised matrices may be qualitatively different, potentially affecting the
260 downstream biological interpretations. To evaluate the changes in interpretation, we
261 compared the clusters obtained prior to and post noise filtering by identifying the (positive)
262 markers and computing the JSI between the top 50 markers of each cluster (Fig. 2G, fig. S3,
263 D, E). Similarly as for the contingency table, the JSI heatmap shows an analogous
264 correspondence between clusters, albeit weaker. Finally, we performed a functional
265 enrichment analysis of the markers identified pre/post noise filtering. Similarly to the bulk
266 results, there were fewer DE genes (markers per cluster) identified in the denoised dataset,
267 with the precision being higher on average across the different GO terms, pathways, and
268 regulatory terms (Fig. 2H). This strengthens our conclusion that the noise filtering process
269 can add focus to the downstream biological analysis without significantly altering the
270 overall composition of the data.



272 **Fig. 2 Overview of noise filtering on smartSeq data and impact on biological interpretation**
273 **of results.**

274 (A) PCC calculated on windows of increasing average abundance for the count-matrix based noise
275 removal approach applied to the full count matrix of all cells (four cells shown). (B) PCC calculated
276 on windows of increasing average abundance for the count-matrix based noise removal approach
277 applied to the “pseudo-samples” formed by grouping all cells from each donor. (C) Box plot of the
278 PCC binned by abundance for the transcript-based noise removal approach applied to five groups
279 of five cells each obtained by concatenating the corresponding BAM files. (D) UMAP
280 representation of the cells using the raw count matrix grouped by donor (left) and by inferred cluster
281 (right). (E) UMAP representation of the cells using the denoised count matrix grouped by donor
282 (left) and by inferred cluster (right) (F) Contingency matrix of the clusters formed before and after
283 the noise removal; the shade of each tile represents the proportion of the cluster from the raw matrix
284 (row) that belongs to the corresponding cluster of the denoised matrix (column). (G) Heatmap of
285 the Jaccard similarity index between the 50 most significant markers identified for each cluster on
286 the raw matrix (rows) and denoised matrix (columns). (H) Violin plot of the precision (intersection
287 size divided by the query size) for the results of the enrichment analysis performed on the marker
288 genes found for each cluster of the raw and denoised matrix respectively (log-scale).

289 Effects of noise filtering on the biological interpretation of regulatory interactions

290 One of the main aims of high-throughput sequencing projects, besides the identification of
291 differentially expressed genes (the effect), is to infer the complex interactions of genes that
292 lead to biological functions, the cause (e.g. disease, development or stress response).
293 Understanding these interactions between genes and the corresponding regulatory elements
294 (at transcriptional level, such as transcription factors (31, 32), or post-transcriptional, small
295 RNAs (33)) allows us to unveil the molecular mechanisms encoding phenotypic outcomes,
296 including causes of diseases.

297 *Effect on PARE data on predicting regulatory miRNA/mRNA interactions*

298 First, we sought to understand the effect of noise removal on the identification of
299 miRNA/mRNA interactions. We applied the *noisyR* transcript approach to a Parallel
300 Analysis of RNA Ends Sequencing (PAREseq) dataset (34). The distribution of degraded
301 fragments across transcripts showed the same distribution of correlation vs abundance as
302 we earlier observed for the bulk RNAseq data (Fig. 3A). Using a correlation threshold of
303 0.25, we determined a signal/noise threshold of 60 for this dataset. Next, we matched the
304 highly abundant reads to known miRNAs (Methods, Fig. 3B) and illustrated that by
305 removing the noisy reads, having abundance less than the noise threshold (Fig. 3, C-D), the
306 prediction of interactions is simplified (35) i.e. for most genes only a few peaks were left.
307 In some cases (e.g. Fig. 3C), only a very clear peak was retained after the noise removal,
308 while for other transcripts some secondary interactions were kept. These results illustrate
309 that noise-filtering is a crucial step for producing biologically meaningful mRNA/targeting
310 predictions.

311 *Effect on the inference and interpretation of Gene Regulatory Networks*

312 Characterising direct interactions between regulatory elements and their targets is only
313 feasible for a limited set of interactions (such as the miRNA/mRNA interaction in plants,
314 leading to mRNA degradation (34). To capture more of the vast complexity of gene
315 interactions, for thousands of genes in tandem, Gene Regulatory Networks (GRNs) have
316 been proposed as a systems biology tool to infer (direct and indirect) regulatory interactions
317 from high-throughput sequencing data (expression data). In a gene regulatory network,
318 nodes represent individual genes (e.g. transcription factors) and edges denote the regulatory
319 interaction between connected genes. When edge-weights are considered, they encode the
320 relative strength of the modeled interaction between two genes. After the network inference
321 step, the resulting topology of GRNs can be used as a proxy for capturing the underlying
322 biological and regulatory complexity of the studied process which in combination with
323 enrichment analyses based on various Gene Ontologies generates a comprehensive model
324 of the investigated process.

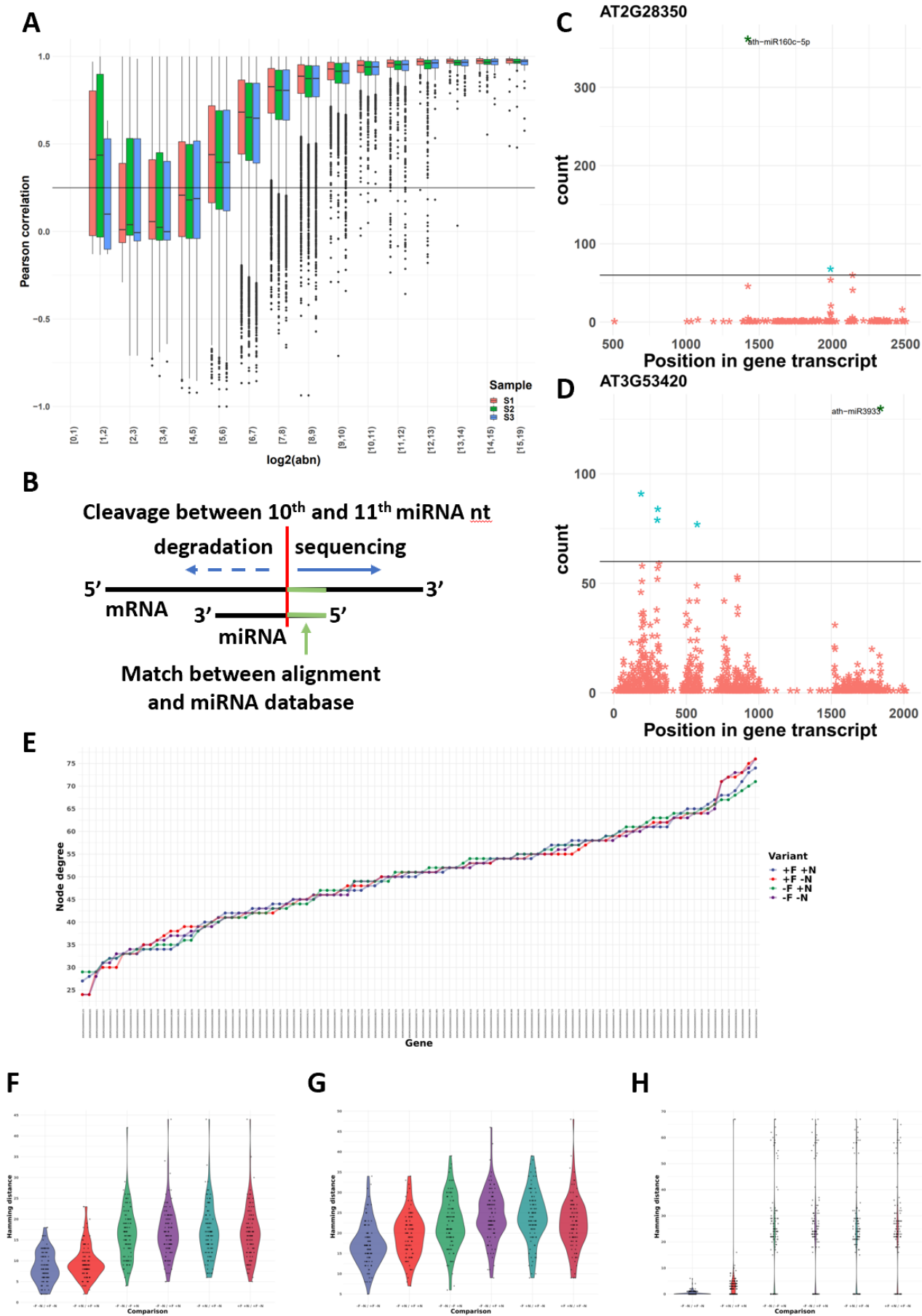
325 We evaluate the impact of noise-filtering on the inference of GRNs on particular network
326 modules (subnetworks), associated with annotated pathways; we quantify the impact of
327 random noise in altering network topologies and subsequent biological interpretations. To
328 achieve this, we run our Network Inference Pipeline (NIP) and *edgynode* network analytics
329 package (Methods) on bulk RNA-seq datasets using non-noise-filtered original, non-noise-
330 filtered normalised, and noise-filtered normalised count matrices. Bulk RNAseq data has
331 been widely used despite its well-known effect to dilute expression signals of individual
332 cells or tissue types. However, in the context of technical noise, the averaging across cells
333 and tissues may buffer the noise effect on general patterns while reducing the possibility to
334 detect weak, but biologically meaningful, expression signals (e.g. transcription factor (36)
335 or transposable element expression (37)).

336 Using the Yang dataset (20) in four different setups (original, -F(iltered) -N(normalised);
337 noise-filtered but not normalised, +F -N; not filtered but normalised, -F +N; and noise-
338 filtered and normalised, +F +N) and subsetting on five biological pathways (Placenta
339 development, 46 genes; Neuron differentiation, 102 genes; Cell differentiation, 249 genes;
340 Phosphorus metabolic process, 493 genes; and Multicellular organism development 996
341 genes), we ran NIP to infer GRNs using three inference approaches GENIE3, GRNBoost2,
342 and PIDC, detailed in Methods. The inferred weighted correlation networks were imported
343 into *edgynode* and rescaled to the range [0,100] to allow comparisons across inference tools.

344 Next, all rescaled weight matrices (fig. S4, E and F) were converted to binary format, using
345 the median value over the entire weight matrix as threshold; a zero was assigned if the
346 weight was below the median value, and a one, if the weight was above the median value.
347 The resulting binary adjacency matrices were then used as input to compute the gene-
348 specific node degrees and to calculate the pairwise Hamming distances for each gene
349 between combinations of original, noise-filtered, and normalised datasets (Fig. 3F, fig. S4,
350 A-D) (Methods). This per-gene Hamming distance is a direct assessment of the number of
351 edges that differ between inferences and captures both edge gain and loss. A low Hamming
352 distance illustrates a robust network, whereas a high Hamming distance is proportional to
353 large changes in the overall GRN topology. Panels Fig. 3, F-H illustrate pairwise
354 comparisons between all combinations of input datasets: 1) original -F -N; 2) not noise-
355 filtered but normalised -F +N; 3) noise-filtered but not normalised +F -N; and 4) noise-
356 filtered and normalised +F +N exemplified for 102 genes corresponding to the neuron
357 differentiation pathway and shown for all three network inference tools (GENIE3, Fig. 3F;
358 GRNBoost2, Fig. 3G; and PIDC Fig. 3H). For all network inference tools, a common pattern
359 is the refining effect of noise-filtering on the overall network topologies. Interestingly, the
360 normalisation step has, in most cases, much greater impact on the network topology than
361 noise-filtering. This result implies that the filtering procedure can detect and remove
362 technical noise without disrupting the global network topology.

363 In addition, (fig. S4, E and F) shows a comparison between rescaled weight matrix
364 distributions for an original and a noise-filtered and normalised network inferred with
365 GENIE3. In this analysis, most genes had a large number of low-weight values within their
366 edge-weight distributions that would result in thousands of biologically meaningless,
367 weakly supported, connections with other genes. Noise-filtering in this bulk RNAseq
368 dataset allows the exclusion of noisy genes as these fall below the median-threshold level
369 which results in a more refined and biologically meaningful network topology after
370 binarisation was applied (Methods).

371 Together, these results suggest that across network inference tools noise-filtering has
372 refining effects on the inferred network topologies in original or normalised data, further
373 illustrating the advantages of noise-filtering to magnify biological signals by reducing
374 technical noise (35).



376 **Fig. 3 Effect of *noisyR* on PARE-Seq and GRN inference**

377 (A) Box plot of the PCC binned by abundance for the transcript-based noise removal approach
378 applied to PARE-Seq data. (B) Schematic overview of the microRNA/mRNA interaction; cleavage
379 of the mRNA transcript occurs between the 10th and 11th nucleotide of the microRNA; (C, D)
380 PARE t-plot illustrating the distribution of degradation products (each point) across the transcripts
381 AT2G28350 and AT3G53420, respectively. All reads with summarised abundance less than the
382 signal/noise thresholds are represented in red; degradation products corresponding to the signal,
383 consistently identified across replicates, are represented in blue. The ones potentially generated by
384 miRNAs are labelled. (E) node degree distributions (total number of edges connected to a
385 node/gene) of 102 genes assigned to the neuron differentiation pathway from the Yang et al dataset.
386 The four input data variants are shown: original (-F -N, purple); not noise-filtered but normalised
387 (-F +N, green); noise-filtered but not normalised (+F -N, red); and noise-filtered and normalised
388 (+F +N, blue) sorted by increasing values using -F -N as sorting key. (F-H) Pairwise hamming
389 distance comparisons for each gene between all combinations of original (-F -N), noise-filtered
390 (+F), and normalised (+N) input datasets using 102 Neuron differentiation genes from the bulk
391 RNAseq (Yang et al.) dataset (Methods) show a comparable pattern across different gene regulatory
392 network inference tools: (F) GENIE3; (G) GRNBoost2; (H) PIDC. The results consistently show
393 that across network inference tools, noise-filtering has refining effects on the inferred network
394 topologies in original or normalised data, further illustrating the advantages of noise-filtering to
395 magnify biological signals by reducing technical noise.

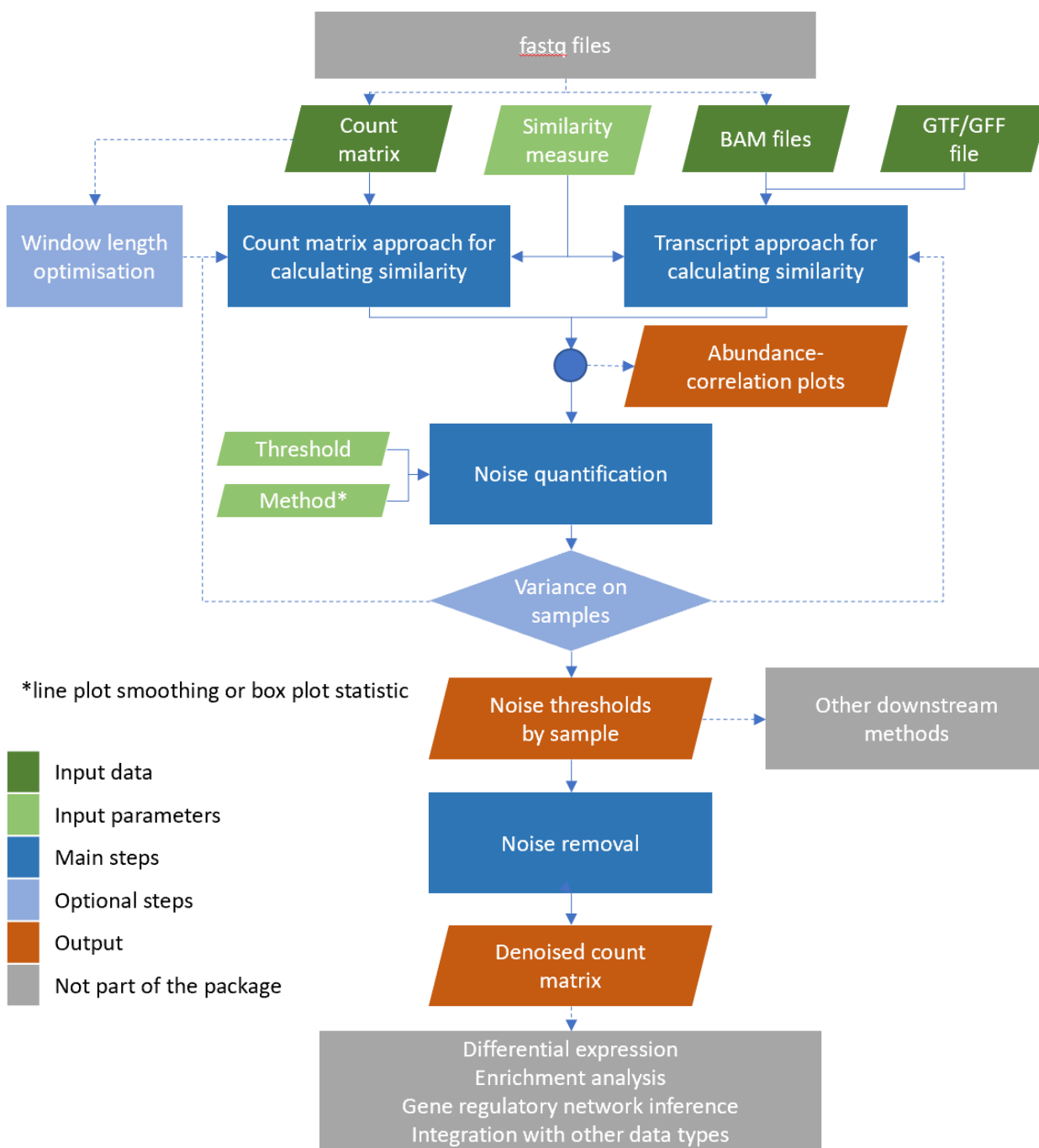
396 **noisyR** package

397 The **noisyR** package is available on CRAN (<https://CRAN.R-project.org/package=noisyR>)
398 and comprises an end-to-end pipeline for quantifying and removing technical noise from
399 HTS datasets. The three main pipeline steps are [i] similarity calculation across samples,
400 [ii] noise quantification, and [iii] noise removal; each step can be finely tuned using hyper-
401 parameters; optimal, data-driven values for these parameters are also determined. The
402 package is written in the R (version 4.0.3) programming language and is actively maintained
403 on <https://github.com/Core-Bioinformatics/noisyR>.

404 For the sample-similarity calculation, two approaches are available. The **count matrix**
405 **approach** uses the original, un-normalised count matrix, as provided after alignment and
406 feature quantification; each sample is processed individually, only the relative expressions
407 across samples are compared. Relying on the hypothesis that the majority of genes are not
408 DE, most of the evaluations are expected to point towards a high similarity across samples.
409 Choosing from a collection of >45 similarity metrics (38), users can select a measure to
410 assess the localised consistency in expression across samples (12). A sliding window
411 approach is used to compare the similarity of ranks or abundances for the selected features
412 between samples. The window length is a hyperparameter, which can be user-defined or
413 inferred from the data (supplementary methods 1). The **transcript approach** uses as input
414 the alignment files derived from read-mappers (in BAM format). For each sample and each
415 exon, the point-to-point similarity of expression across the transcript is calculated across
416 samples in a pairwise all-versus-all comparison. The output formats for the two approaches
417 are the same; the number of entries varies, since the count approach focuses on windows,
418 whereas for the transcript approach we calculate a similarity measure for each transcript.

419 The noise quantification step uses the abundance-correlation (or other similarity measure)
420 relation calculated in **step i** to determine the noise threshold, representing the abundance
421 level below which the gene expression is considered noisy e.g. if a correlation threshold is
422 used as input then the corresponding abundance from a (smoothed) abundance-correlation
423 line plot is selected as the noise threshold for each sample. The shape of the distribution can
424 vary across experiments; we provide functionality for different thresholds and recommend
425 the choice of the one that results in the lowest variance in the noise thresholds across
426 samples. Options for smoothing, or summarising the observations in a box plot and selecting
427 the minimum abundance for which the interquartile range (or median) is consistently above
428 the correlation threshold are also available. Depending on the number of observations, we
429 recommend using the smoothing with the count matrix approach, and the boxplot
430 representation with the transcript option.

431 The third step uses the noise threshold calculated in **step ii** to remove noise from the count
432 matrix (and/or BAM file). The count matrix can be calculated by exon or by gene; if the
433 transcript approach is used, the exon approach is employed. Genes/exons whose expression
434 is below the noise thresholds for every sample are removed from the count matrix. The
435 average noise threshold is calculated and added to every entry in the count matrix. This
436 ensures that the fold-changes observed by downstream analyses are not biased by low
437 expression, while still preserving the structure and relative expression levels in the data. If
438 downstream analysis does not involve the count matrix, the thresholds obtained in **step ii**
439 can be used to inform further processing and potential exclusion of some genes/exons from
440 the analysis.



441

442 **Fig. 4 Workflow diagram of the noisyR pipeline.**

443 Workflow diagram describing the series of steps comprising the *noisyR* pipeline. Individual
 444 algorithms, finely tuned through hyper-parameters, are highlighted in blue. Optional steps
 445 are indicated through higher transparency. Common data pre- and post- processing steps not included
 446 in the package are indicated in grey.

447

448 Discussion

449 User-defined or data-driven options for the hyperparameters

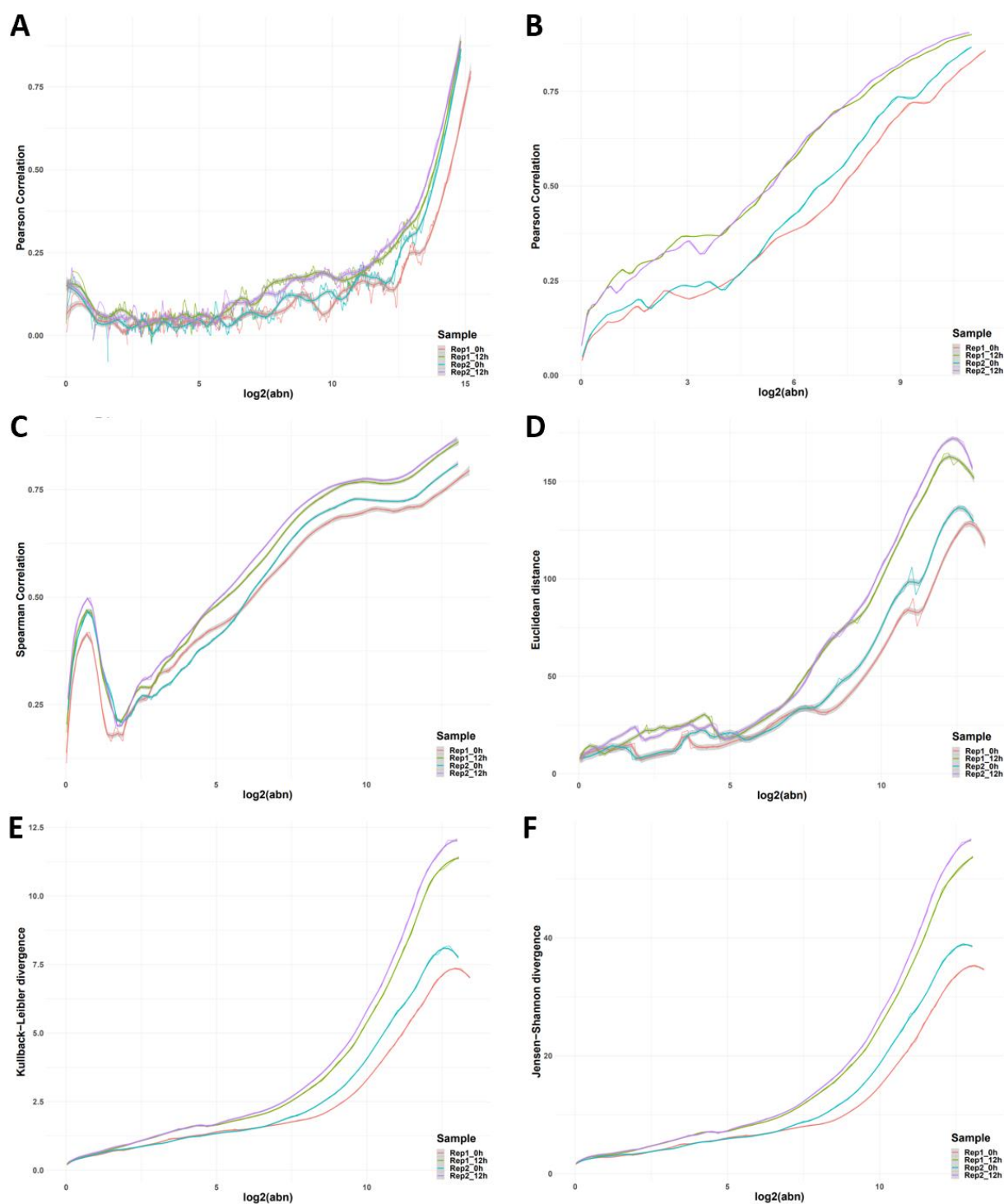
450 **noisyR** hyperparameters can be used to finely tune the identification of the signal/noise
451 thresholds. To optimise the noise filtering procedure and dampen the stochastically induced
452 differences between samples (e.g. derived from variation in sequencing depth or sample
453 read-complexity) the noise removal step is performed by adding the average of the
454 signal/noise thresholds across samples, on the raw count matrix. Nevertheless, comparable
455 thresholds across the dataset are essential for a meaningful filtering; we recommend the use
456 of consistency and robustness checks throughout the pipeline to ensure that the input
457 samples are comparable, coupled with the data-driven selection of threshold values for
458 setting hyper-parameters. The option of user-defined values is available, however the
459 selected values should be based on observations from the input dataset, rather than
460 exclusively following default recommendations. Next, we discuss in detail the options
461 available for selecting the hyperparameters for a more adaptive noise-filtering based on the
462 structure of the input data.

463 For the count matrix approach, the length of the sliding windows plays a significant role in
464 assessing the similarity across samples. Smaller windows require more computational time;
465 however the intended level of detail may not always be preferable, as small gene expression
466 fluctuations, from sample to sample, would reduce the across-sample similarity if the
467 abundance range is not wide enough (Fig. 5A). Even for medium-high abundances,
468 expression or rank inconsistencies characterise smaller windows, indirectly leading to
469 higher (and more variable across samples) signal/noise thresholds. If the window size is too
470 large, less information is captured by the similarity measure and the accuracy of the noise
471 threshold identification is also reduced (Fig. 5B). We recommend medium-sized windows
472 that cover the abundance range in small incremental steps as larger overlaps between
473 windows result in a more robust estimation of similarity-variation. An intuitive approach
474 for determining an informative window size for a dataset relies on monotony changes of the
475 similarity measure, quantified as the number of times the derivative of the correlation (as a
476 function of abundance) changes sign. On several datasets, this resulted in a window length
477 of 1/10th of the total number of expressed genes and a sliding window step size of 1/20th
478 of the total gene number. A different tactic, also implemented in **noisyR**, tackles this task
479 from a different direction; it relies on optimising the window length using an entropy-based
480 approach with the Jensen-Shannon divergence to assess the stability achieved as the window
481 length is increased (supplementary methods 1). The shape of the distribution of correlations
482 changes as the window length increases; however the change is less significant (evaluated
483 using a t-test) for larger windows. The first point of stability is selected as the optimal
484 window length, as it provides the largest possible granularity while maintaining robustness.
485 The results from this approach are also consistent with earlier, empirical findings when
486 applied to the Yang dataset (20).

487 Yet another hyperparameter is the similarity measure; we compared the results for different
488 correlation and distance metrics. We aim to achieve a high consistency in quantifying the
489 signal/noise thresholds that is independent of the similarity measure. We tested the standard
490 parametric and non-parametric correlation measures as well as the ones implemented in the
491 *philentropy* package (38), which provides a variety of >45 distance measures. Dissimilarity
492 measures are being inverted for comparison purposes (Fig. 5C-F illustrates the Spearman
493 correlation, Euclidean distance, Kulbeck-Leibler divergence, and Jensen-Shannon
494 divergence). Some measures have fixed ranges (e.g. the correlation coefficients), while
495 others are semi- or unbounded. This raises the question of how to choose a similarity
496 threshold when the range of values resulting from the similarity measure is unknown.

497 Inspired by the correlation threshold, which provides a good separation at 0.25 for many
498 datasets, we focus, as a starting point, on the naive assumption to use a quarter of the full
499 range of the observed similarity values as a first cut-off approximation. Picking a threshold
500 in a data-driven manner is, however, preferable and, in this case, achievable. Selecting from
501 a variety of threshold values that minimise the coefficient of variation (standard deviation
502 divided by the mean) of the corresponding noise thresholds in different samples is an
503 empirical approach that works in practice. If the samples are semantically grouped e.g.
504 replicates or time points, it may be better to minimise the variation in each individual group
505 rather than across the full experimental design.
506

507



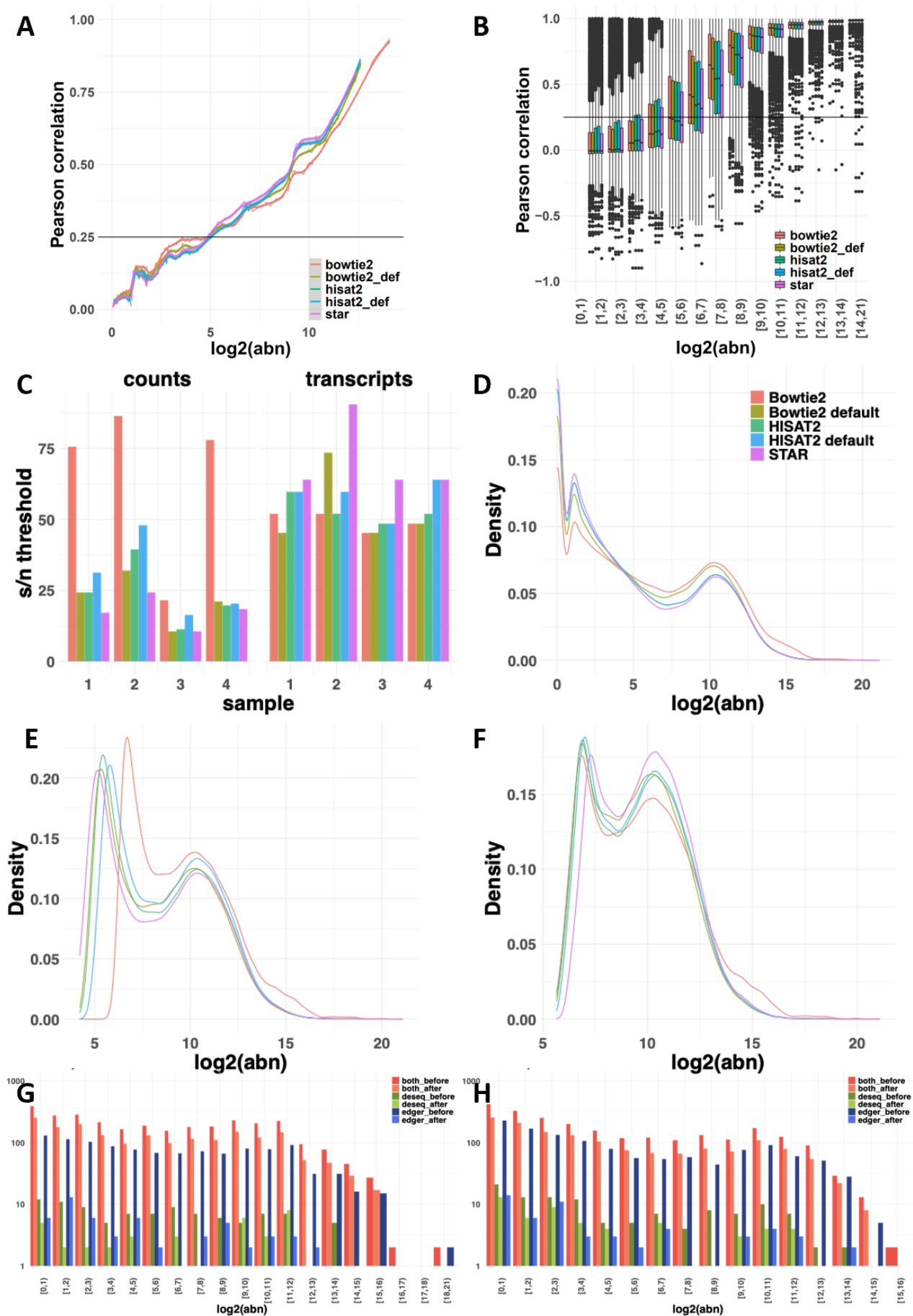
508

509 **Fig. 5 Effects of hyperparameter selection on noise quantification.**

510 (A) PCC-abundance plot for a window length of 1,000 genes, $\sim 1/5$ th of the default. (B) PCC-
 511 abundance plot for a window length of 20,000 genes, ~ 4 times the default. (C) Spearman
 512 correlation plotted against abundance for the default window length of $\sim 5,500$. (D) Inverse of the
 513 Euclidean distance plotted against abundance for the default window length of $\sim 5,500$. (E)
 514 Inverse of the Kulbeck-Leibler divergence plotted against abundance for the default window
 515 length of $\sim 5,500$. (F) Inverse of the Jensen-Shannon divergence plotted against abundance for the
 516 default window length of $\sim 5,500$.

517 Effect of aligner choice on noise quantification

518 The choice of the read-aligner was shown to influence the downstream DE analyses when
519 the same quantification model was applied (18). To assess the effect of different alignment
520 approaches on the quantification and observed levels of noise, mRNA quantification using
521 featureCounts (39) was performed on reads aligned with STAR (40), HISAT2 (41) and
522 Bowtie2 (42). The latter two were run both using their default parameters and with
523 parameters set to match STAR functionality. For the count-based approach, the distribution
524 of the Pearson Correlation Coefficients across abundance bins (Fig. 6A) shows that noise
525 levels were relatively consistent regardless of the applied alignment algorithm. Similarly,
526 for the transcript-based approach, the correlation distributions across abundance bins (Fig.
527 6B) illustrate little variation across aligners (additional examples in fig. S6, A and B). The
528 estimated signal/noise thresholds were also comparable between the datasets generated by
529 different aligners (Fig. 6C), with transcripts-based noise results being less variable. Once
530 the noise correction was applied, the substantial peak in the abundance distributions around
531 zero (Fig. 6D) was removed or significantly diminished and a second peak corresponding
532 to the true signal was revealed around $\log_2(\text{abundance})$ of five using both counts and
533 transcripts based approaches (Fig. 6,E and F respectively). The similarity of the abundance
534 distributions across the datasets produced by the different aligners was observable both
535 before and after the noise correction. This demonstrates that the proposed correction
536 approaches are non-destructive and preserve the underlying biological signal. To further
537 validate this point, the overlap between edgeR and DESeq2 analyses was investigated. The
538 differentially expressed (DE) genes (adjusted p-value < 0.05 and $|\log_2(\text{FC})| > 1$) detected
539 by the two methods were compared for outputs produced using STAR (Fig. 1J), Bowtie2
540 (Fig. 6G) and HISAT2 (Fig. 6H). In all cases, there were fewer DE genes in total after noise
541 correction was applied, and the specific differences for each DE method were reduced. The
542 same conclusions were reached for the processing with Bowtie2 and HISAT2 applied with
543 their default parameters (fig. S6C).



545 **Fig. 6 Assessment of aligner choice on noise quantification**

546 **(A)** The distribution of PCC across abundance bins in datasets for a single mRNAseq sample
547 obtained by STAR, Bowtie2 and HISAT2 alignment followed by featureCounts quantification
548 using a counts-based noise removal approach. **(B)** The distribution of PCC across abundance bins
549 in aligned read counts obtained by the five aligners for the same sample in the transcript-based
550 noise correction approach. **(C)** The detected signal-to-noise thresholds in the four mRNAseq
551 samples varied when the counts or transcripts-based noise correction methods were applied. **(D)**
552 The distribution of abundance of reads aligned by the five algorithms and quantified by
553 featureCounts. **(E)** The distribution of abundance of the quantified counts after counts-based noise
554 correction **(F)** The distribution of abundance of the quantified counts after transcripts-based noise
555 correction. **(G)** The number of the differentially expressed genes found by applying DESeq and
556 edgeR on the original and denoised (using transcripts-based approach) count matrices obtained by
557 Bowtie2 alignment. **(H)** The overlap between the DESeq and edgeR analyses performed on the
558 original and denoised counts matrices obtained by HISAT2.
559

560 The effects of noise-filtering on GRN inference for single cell RNAseq data

561 The recent emergence of single-cell sequencing technologies enabled the simultaneous
562 assessment of expression variation between individual cells across thousands of cell-
563 lineages. Although conceptually powerful, sequencing depths remain constrained by cost
564 and in comparison to bulk RNAseq experiments the total number of reads is now shared
565 among these (hundreds-) thousands of individual cells expressing thousands of genes each.
566 This limit on the sequencing depth per cell underlines, yet again, the technical noise,
567 whereby the quantification of low-abundance transcripts can be the result of either low
568 biological expression or due to stochastic effects (likelihood) of read capturing. The
569 requirement of an adaptive noise-filtering pipeline is fulfilled by *noisyR*; the retained gene
570 expression levels increases the robustness of quantification of single-cell data.

571 Analogous to the Yang et al. dataset, we used the Cuomo (28) dataset in four different setups
572 (original, -F -N; noise-filtered but not normalised, +F -N; not filtered but normalised, -F
573 +N; and noise-filtered and normalised, +F +N) and subsampled into three distinct biological
574 pathways (Metabolism, 57 genes; Catalytic activity, 133 genes; Cellular metabolic process,
575 246 genes), we ran the Network Inference Pipeline to infer GRNs using the same three
576 inference methods GENIE3, GRNBoost2, and PIDC (Methods) as used for bulk RNAseq
577 data. The inferred weighted correlation networks were imported into *edgynode* and rescaled
578 (fig. S5, C and D) analogous to the bulk RNAseq data shown in Results and Methods. The
579 resulting pairwise Hamming distances for each gene between combinations of original,
580 noise-filtered, and normalised datasets and for genes corresponding to various biological
581 pathways (Fig. 7, A-C, fig. S5, A and B) show that total Hamming distances over all genes
582 are larger in single-cell data. This implied that noise-filtering had a more significant/
583 refining impact on the inference and biological interpretations drawn from single-cell data
584 when compared with analogous bulk RNA data.

585 Fig. 7 D-F illustrates such analogous pairwise comparisons between all combinations of
586 input datasets: 1) original -F -N; 2) not noise-filtered but normalised -F +N; 3) noise-filtered
587 but not normalised +F -N; and 4) noise-filtered and normalised +F +N exemplified for 133
588 genes corresponding to catalytic activity pathways derived from single-cell RNAseq data
589 (Cuomo et al.) and also shown for all three network inference tools (GENIE3, Fig. 7A;
590 GRNBoost2, Fig. 7B; and PIDC Fig. 7C). Analogous to the bulk RNAseq results, noise-
591 filtering has smaller effects on changes in network topologies than the normalisation step.
592 Interestingly, it seems that the overall effect of noise-filtering in single-cell data has a
593 stronger impact than in bulk RNAseq data (Fig. 3, F-H). Together, these conclusions hint
594 toward a more useful effect of noise-filtering in single-cell data as is particularly expected
595 for datasets with limited sequencing depth, but high individual cell numbers.

596 These highlight the positive effects of noise-filtering on magnifying meaningful biological
597 signals in single-cell RNAseq data, with more significant effects in single-cell data due to
598 the nature of technical noise induced by sequencing depth-constraints in combination with
599 technical variation.

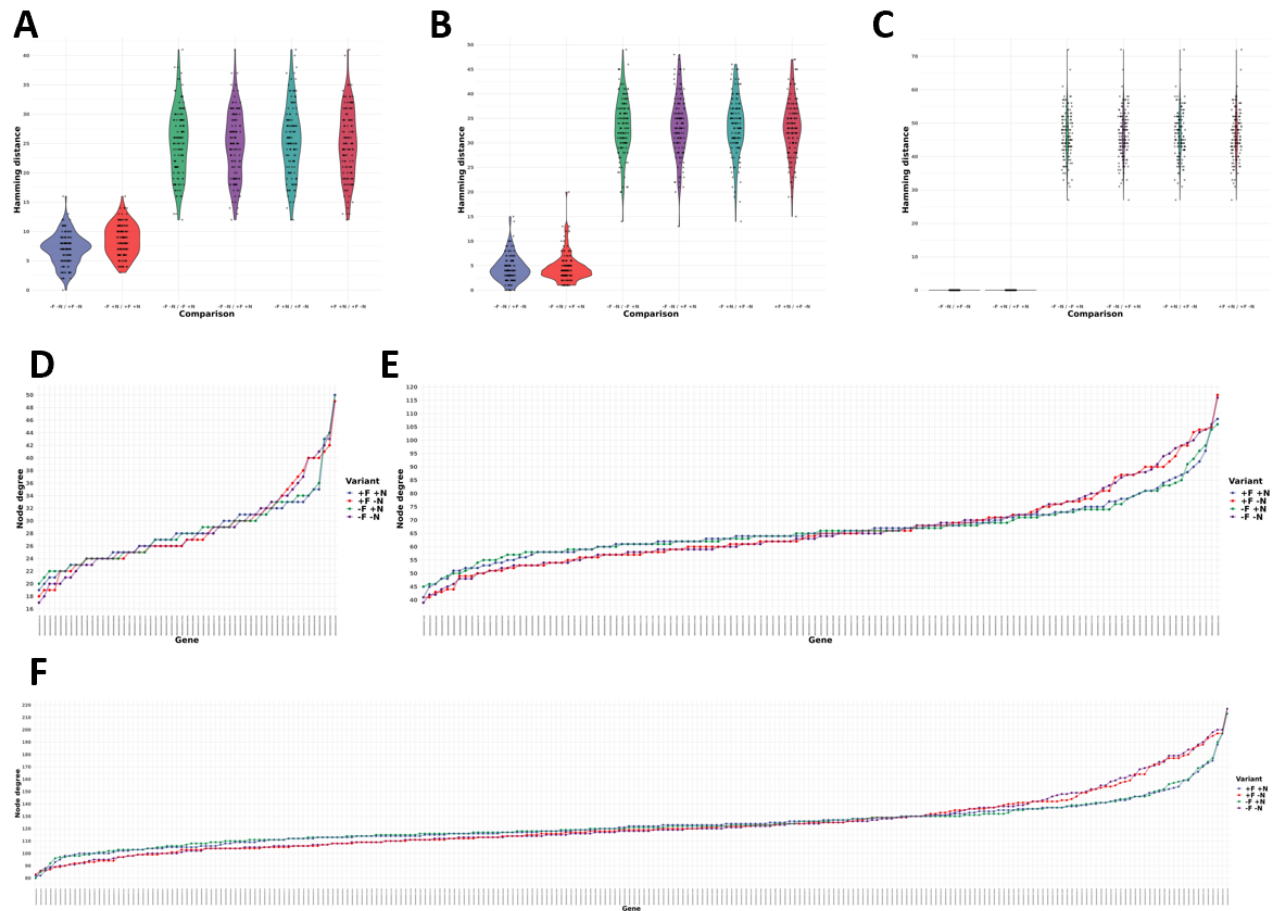


Fig. 7 Node degree distributions and pairwise Hamming distance distributions between combinations of original, noise-filtered, and normalised input smartSeq datasets

(A-C) Pairwise Hamming distance comparisons for each gene between all combinations of original (-F -N), noise-filtered (+F), and normalised (+N) input datasets using 133 genes associated with catalytic activity pathways from the smartSeq (Cuomo et al.) dataset (Methods) show a comparable pattern across different gene regulatory network inference tools: (A) GENIE3; (B) GRNBoost2; (C) PIDC. The results consistently show that across network inference tools noise-filtering has refining effects on the inferred network topologies in original or normalised data, further illustrating the advantages of noise-filtering to magnify biological signals by reducing technical noise. (D-F). For the smartSeq (Cuomo et al.) dataset the node degree distributions (total number of edges connected to a node/gene) of three sets of genes corresponding to different pathways are shown: (D) 57 genes associated with metabolism pathways; (E) 133 genes associated with catalytic activity pathways; (F) 246 genes associated with the cellular metabolic process. All four input data variants are shown: original (-F -N, purple); not noise-filtered but normalised (-F +N, green); noise-filtered but not normalised (+F -N, red); and noise-filtered and normalised (+F +N, blue) sorted by increasing values using -F -N as sorting key.

619 Using *noisyR*, we demonstrate that unfiltered RNAseq quantifications can cause spurious
620 false positive effects in various standard expression-based data analysis steps. To overcome
621 these limitations, we introduce an R package to equip life scientists with a flexible solution,
622 applicable across different bulk and single cell datasets, for excluding inconsistent transcript
623 quantifications that would otherwise introduce stochastic variability in processed datasets.
624 A comprehensive selection of automatic and semi-automatic threshold detection options
625 provided by *noisyR* allows the robust inference of noise-thresholds to exclude low-
626 confidence transcripts from processed RNAseq data. We illustrate the importance of such a
627 noise-filtering procedure by assessing the convergence of DE identification and by inferring
628 and comparing gene regulatory networks from various biological pathways, across gold-
629 standard network inference tools. As a result, we find that noise-filtering is indeed able to
630 significantly reduce stochastic effects magnifying underlying biological signals, thereby
631 yielding more robust biological interpretations.

632 **Materials and Methods**

633 **Materials**

634 The bulk mRNA-seq used to illustrate *noisyR* was generated by Yang et al (20). The dataset
635 comprises 16 samples across 8 time points [0-72 hours post stem cell induction. The raw
636 data (fastq files and metadata) were downloaded from GEO (accession numbers
637 GSE117896, GSM3314677 - GSM3314692).

638 Next, sRNA data was retrieved from Paicu et al (43) for the plant dataset (2 samples, a
639 wildtype and DCL1 knockdown, with 3 biological replicates each, in *A. thaliana*,
640 GSM2412286 - GSM2412291) and from Wallach et al (44) for the animal dataset, 6
641 samples generated for the identification of microRNAs as TLR-activating molecules in *M.*
642 *musculus* (PMID: 31940779, GSE138532, GSM4110737 - GSM4110742). For both
643 datasets, the reads were aligned to mature and hairpin miRNAs, downloaded from miRBase
644 (45) and TEs, downloaded from TAIR and Ensembl, for *M. musculus*.

645 For assessing the impact of noise on direct biological interpretations and predictions, such
646 as the interaction of miRNAs and mRNAs, we selected a PARE (parallel analysis of RNA
647 ends, also known as degradome sequencing) dataset, consisting of 3 biological replicates
648 (GSE113958) presented in Thody et al (34).

649 The single-cell mRNA-seq dataset used to illustrate *noisyR* was generated by Cuomo et al
650 (study of stem cell differentiation) (28). The data is available on ENA, ERP016000 -
651 PRJEB14362. The six donors with the highest number of cells (hayt, naah, vils, pahc, melw,
652 qunz) were selected, cells in time point 3 were included.

653 The reference genomes used for alignment were: Homo_sapiens.GRCh38.98 (Ensembl
654 version 98), Mus_musculus.GRCm38.98 (Ensembl version 98) and *A. thaliana* (46).

655 **Methods, bulk mRNAseq data**

656 *Data pre-processing and quality checking*

657 Initial quality checks were performed using fastQC (version 0.11.8), summarised with
658 multiQC (version 1.9) (47). Alignments to reference genomes were performed using STAR
659 (version 2.7.0a) with default parameters (40); the count matrices were generated using
660 featureCounts (version 2.0.0) (39) against the *M. musculus* exon annotations obtained from
661 the Ensembl database (genome assembly GRCm38.p6). Additional quality checks included
662 density plots, (comparable distributions are a necessary but not sufficient condition for
663 comparability), MA plots for the sufficiency check (expected to have a funnelling shape;
664 observed outliers are candidates for differentially expressed transcripts), incremental
665 dendrograms and PCA plots to evaluate the similarity of distributions (12, 48).

666 *Data post-processing and biological interpretation of results*

667 The differential expression analysis was performed after quantile normalisation of the count
668 matrix using the standard functions from edgeR, version 3.28.0 (8) and DESeq2, version
669 1.26.0 (7). The thresholds for DE were $|\log_2(\text{FC})| > 1$ and adjusted p-value < 0.05
670 (Benjamini-Hochberg multiple testing correction). The enrichment analysis was performed
671 using g:profiler (R package gprofiler2, version 0.2.0) (21), against the standard GO terms,
672 and the KEGG (49) and reactome (50) pathway databases. The observed set consisted of
673 the DE genes, the background set comprised all expressed genes, using the full or de-noised
674 count matrix respectively.

675 To assess the effect of noise correction across the multiple options of mRNA quantification,
676 the sequencing reads were aligned to the reference genome using Bowtie2 (version 2.4.2)

677 (42) and HISAT2 (version 2.1.0) (41). Aligners were run both with default parameters and
678 with parameters set to match the STAR functionality of searching for up to 10 distinct, valid
679 alignments for each read ("bowtie2 --end-to-end -k 10" and "hisat2 -q -k 10"). The transcript
680 expression was quantified using featureCounts. The robustness of the quantification was
681 assessed by investigating the overlap between edgeR and DESeq2 analyses. The genes with
682 adjusted p-value < 0.05 (Benjamini-Hochberg multiple testing correction) and $|\log_2(\text{FC})| >$
683 1 were considered before and after noise correction.

684 *Gene regulatory network inference*

685 To assess the implications of the noise filter on downstream biological interpretations, we
686 used the bulk and single-cell datasets as inputs for various gene regulatory network (GRN)
687 inference tools and compared the results for filtered and unfiltered inputs. For this purpose,
688 we selected several gene subsets, ranging in size from 49 to 996 genes for the bulk dataset
689 and from 57 to 246 genes for the single-cell dataset, based on enrichment analyses
690 performed on the DE genes according to their inclusion in annotated pathways.
691 (Supplementary table 1)

692 We chose a subset of the GRN inference tools benchmarked by BEELINE (51): GENIE3
693 (52), GRNBoost2 (53), and PIDC (54). We packaged the tools as Singularity containers
694 (<https://github.com/drostlab/network-inference-toolbox>) and then assembled them into a
695 custom pipeline (<https://github.com/drostlab/network-inference-pipeline>).

696 This pipeline extracts the subsets of genes corresponding to selected pathways and uses
697 them as inputs for the GRN inference tools. The results are rescaled, binarised and compared
698 using the *edgynode* package (v0.3.0, <https://github.com/drostlab/edgynode>). The edge
699 weights and node degree distributions for all genes across the selected subsets are then
700 visualised.

701 In detail, the similarity assessment of network topologies was performed using the *edgynode*
702 function `network_benchmark_noise_filtering()` and was visualized using
703 `plot_network_benchmark_noise_filtering()`. For this purpose, the inferred networks were
704 converted to a binary format (presence/absence of an edge) using the overall median edge
705 weight per network as a threshold. In `network_benchmark_noise_filtering()` four different
706 types of matrices are used as input: a weighted adjacency matrix returned by a network
707 inference tool where 1) no noise filter and no quantile normalisation (original) was
708 performed (denoted in the figures as -F -N), 2) a noise filtering but no quantile normalisation
709 was performed (+F -N), 3) no noise filtering but a quantile normalisation was performed (-
710 F +N), and 4) both, noise-filtering and quantile normalization were performed (+F +N).

711 In a pairwise all versus all comparison, for each gene, the Hamming distance over the binary
712 edge weight vectors was computed using the `hamming.distance()` function from the R
713 package `e1071 v1.7-4` (55), yielding a distribution of distances, which captures how many
714 genes gained or lost their connection with other genes. A Kruskal-Wallis Rank Sum Test
715 was performed using the `stats::kruskal.test()` function in R to assess whether comparisons
716 of Hamming distance distributions between original, noise-filtered, and normalized
717 combinations were statistically significantly different. Furthermore, visualising these
718 distributions across comparisons and for all network inference tools facilitated an evaluation
719 of the overall change of network topologies driven by the network inference tool or the
720 normalisation/noise-filtering that was applied. These visualizations were then used to assess
721 the impact and robustness of our noise-filter on the interpretation of biological network
722 topologies. We applied the pipeline, including *edgynode*, with the same parameter
723 configurations to both, bulk (Yang et al.) and single-cell (Cuomo et al.) data to retrieve
724 comparable results for direct comparisons. Computationally reproducible analysis scripts to

725 perform all inference steps, data transformations, and visualisations, including the ones used
726 in this study can be found at <https://github.com/drostlab/network-inference-pipeline>.

727 Methods, sRNAseq data

728 The 6 *A. thaliana* sRNA samples were assessed using multiQC version 1.9 (47). Next, the
729 sequencing adapters (both standard and HD) were trimmed using Cutadapt (version 3.2)
730 (56) and the UEA sRNA Workbench (57). The larger 3 samples were subsampled without
731 replacement to 8M reads (12); the smaller 3 samples were left unchanged. The read/sRNA-
732 length distributions were bimodal with peaks at 21nt and 24nt, corresponding to miRNAs
733 and TE- sRNAs, respectively. These sRNAs were aligned (using STAR (version 2.7.0a)
734 (40)) to both microRNA hairpins (miRBase Release 22.1) (45) and TEs (obtained from
735 TAIR10) (46).

736 The 6 *M. musculus* sRNA samples were processed in a similar way as the plant samples and
737 subsampled without replacement to 3.5M sequences (12). The distribution of read lengths
738 was bimodal with peaks at 22nt and 30nt corresponding to microRNAs and piRNAs
739 respectively. The sRNAs were aligned to microRNA hairpins (miRBase Release 22.1) (45)
740 and TEs (Ensembl release 101).

741 Methods, PARE data

742 The 3 *A. thaliana* PARE samples (GSE113958) were QCed (multiQC version 1.9) (47) and
743 the reads trimmed to 20nt; next, all samples were randomly subsampled without
744 replacement to 25M (12). The subsampled reads were aligned to the reference genome
745 (obtained from TAIR10 (46)) using STAR (using STAR (version 2.7.0a) (40)), with default
746 parameters. The reads aligned to each position along a transcript were grouped on sequence
747 and summarised by frequency. Each summarised fragment was matched (as reverse
748 complement) to *A. thaliana* miRNAs. To visualise the distribution of signal across
749 transcripts, t-plots were created, where each point corresponds to a summarised PARE
750 fragment; the points for which a corresponding miRNA was identified were highlighted
751 using the miRNA label (34).

752 Methods, single cell data

753 For the single cell SmartSeq2 data, the cellranger software version 3.0 (58) was used for
754 pre-processing, initial quality checks, and to generate the count matrix (it internally uses the
755 STAR aligner). Further quality checks included distribution plots for the number of features,
756 counts, mitochondrial and ribosomal reads per cell; significant outliers were removed during
757 pre-processing. Dimensionality reduction and clustering were performed with the Seurat R
758 package version 3.2 (29). The UMAP reduction method (30) was used for visualisation and
759 assessment of results.

760 Methods, noise quantification

761 Two approaches were implemented for the identification of noise. (1) The “count matrix
762 approach” is a simple, fast way to obtain a threshold utilising solely the un-normalised count
763 matrix (m genes x n samples). (2) The “transcript approach” is more refined, as it takes into
764 account the distribution of signal across the transcript obtained by summarising the aligned
765 reads from the BAM alignment files. For both approaches, a variety of correlation and
766 distance measures are used to assess the stability of signal across samples (38). Most results
767 were obtained using Pearson Correlation Coefficient (the default); similar results are
768 obtained with other similarity or inverted dissimilarity measures such as Spearman
769 Correlation, Euclidean distance, Kulbeck-Leibler divergence, and Jensen-Shannon
770 divergence.

771 *Count matrix approach*

772 For each sample in the count matrix, the genes are sorted, in descending order, by
773 abundance. A sliding window approach is used to scan the sorted genes (genes with similar
774 abundances are grouped into “windows”). The window length is a hyper-parameter that can
775 be user-defined or a single value inferred from the data using a Jensen-Shannon entropy
776 based approach (supplementary methods 1). The sliding step can be varied to reduce
777 computational time at the cost of reducing the number of data points and potentially losing
778 accuracy. For each window, the correlation of the abundances of the genes from the sample
779 of interest and all other samples is calculated and averaged using the arithmetic mean. Per
780 sample, the variation in correlation coefficient (y-axis) is represented vs the average window
781 abundance, x-axis. A correlation threshold (as a hyper-parameter) is used to determine a
782 corresponding abundance threshold as a cut-off - the noise threshold. The correlation
783 threshold is inferred from the data to minimise the variance of noise thresholds across the
784 different samples. Several available approaches are based on the (smoothed) line plot or a
785 binned boxplot of abundance against correlation (supplementary methods 2). Genes with
786 abundances below the sample specific noise thresholds across samples were excluded from
787 downstream analyses; the average of the thresholds were added to the count matrix, to avoid
788 further biases. By increasing the minimum values in the count matrix from zero to the noise
789 threshold, methods that are based on fold-changes will not emphasise small differences in
790 abundance at very low values, which becomes especially problematic for genes that are
791 seemingly absent in some samples but present and lowly expressed in others. This effect is
792 particularly striking in single-cell data.

793 *Transcript approach*

794 Using the transcript coordinates of the aligned reads as input, the expression profile for each
795 individual transcript was built as an algebraic point sum of the abundances of reads incident
796 to any given position (59); if the alignment was performed per read, the corresponding
797 abundance for every entry was set to +1. For each sample j , and for each transcript T , the
798 point-to-point Pearson Correlation between the expression profile in j and the one in all
799 other samples is calculated. The noise detection is based on the relative location of the
800 distribution of the point-to-point Pearson Correlation Coefficient (p2pPCC) versus the
801 abundances of genes and is specific for each individual sample. For low abundance
802 transcripts the stochastic distribution of reads across the transcript leads to a low p2pPCC;
803 the aim of the approach is to determine the range where the distribution of correlation
804 coefficients (used as proxy for the distribution of reads across a transcript) are above a user-
805 defined threshold; to approximate the signal-to-noise threshold a binning on the abundances
806 was performed. For all examples presented in this study, the binning was done on log2
807 ranges; the signal-to-noise thresholds were defined as the abundance above which the first
808 quartile of the p2pPCC distribution consistently remains above 0.25 (IQR method - see
809 supplementary methods 2). Once a noise threshold was determined for each sample, the
810 original count matrix was then filtered analogous to the count matrix approach. The BAM
811 files can also be filtered directly by removing all genes which fall below the noise threshold
812 in every sample. Downstream analysis that is not based on the count matrix, such as
813 alternative splicing analysis can also be informed by the noise threshold by setting a lower
814 bound of expression acceptance.

815 **References**

- 816
817 1. R. Stark, M. Grzelak, J. Hadfield, RNA sequencing: the teenage years. *Nature Reviews*
818 *Genetics* **20**, 631-656 (2019).

- 819 2. A. Oshlack, M. D. Robinson, M. D. Young, From RNA-seq reads to differential
820 expression results. *Genome Biology* **11**, 220 (2010).
- 821 3. A. Conesa *et al.*, A survey of best practices for RNA-seq data analysis. *Genome Biology*
822 **17**, 13 (2016).
- 823 4. M. Li, J. C. I. Belmonte, Ground rules of the pluripotency gene regulatory network.
824 *Nature Reviews Genetics* **18**, 180-191 (2017).
- 825 5. S. Parekh, C. Ziegenhain, B. Vieth, W. Enard, I. Hellmann, The impact of amplification
826 on differential expression analyses by RNA-seq. *Scientific Reports* **6**, 25533 (2016).
- 827 6. K. D. Hansen, S. E. Brenner, S. Dudoit, Biases in Illumina transcriptome sequencing
828 caused by random hexamer priming. *Nucleic Acids Research* **38**, e131-e131 (2010).
- 829 7. M. I. Love, W. Huber, S. Anders, Moderated estimation of fold change and dispersion for
830 RNA-seq data with DESeq2. *Genome Biology* **15**, (2014).
- 831 8. D. J. McCarthy, Y. Chen, G. K. Smyth, Differential expression analysis of multifactor
832 RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* **40**,
833 4288-4297 (2012).
- 834 9. T. Stuart, R. Satija, Integrative single-cell analysis. *Nature Reviews Genetics* **20**, 257-272
835 (2019).
- 836 10. F. Rapaport *et al.*, Comprehensive evaluation of differential gene expression analysis
837 methods for RNA-seq data. *Genome Biology* **14**, R95 (2013).
- 838 11. Z. Wang, M. Gerstein, M. Snyder, RNA-Seq: a revolutionary tool for transcriptomics.
839 *Nature Reviews Genetics* **10**, 57-63 (2009).
- 840 12. I. Mohorianu *et al.*, Comparison of alternative approaches for analysing multi-level RNA-
841 seq data. *PLOS ONE* **12**, e0182694 (2017).
- 842 13. G. Park *et al.*, Characterization of background noise in capture-based targeted sequencing
843 data. *Genome Biology* **18**, (2017).
- 844 14. I. Fischer-Hwang, I. Ochoa, T. Weissman, M. Hernaez, Denoising of Aligned Genomic
845 Data. *Scientific Reports* **9**, (2019).
- 846 15. K. Shiroguchi, T. Z. Jia, P. A. Sims, X. S. Xie, Digital RNA sequencing minimizes
847 sequence-dependent bias and amplification noise with optimized single-molecule
848 barcodes. *Proceedings of the National Academy of Sciences* **109**, 1347-1352 (2012).
- 849 16. G. Eraslan, L. M. Simon, M. Mircea, N. S. Mueller, F. J. Theis, Single-cell RNA-seq
850 denoising using a deep count autoencoder. *Nature Communications* **10**, (2019).
- 851 17. C. Jia *et al.*, Accounting for technical noise in differential expression analysis of single-
852 cell RNA sequencing data. *Nucleic Acids Research* **45**, 10978-10988 (2017).
- 853 18. A. Srivastava *et al.*, Alignment and mapping methodology influence transcript abundance
854 estimation. *Genome Biology* **21**, (2020).
- 855 19. L. A. Corchete *et al.*, Systematic comparison and assessment of RNA-seq procedures for
856 gene expression quantitative analysis. *Scientific Reports* **10**, (2020).
- 857 20. P. Yang *et al.*, Multi-omic Profiling Reveals Dynamics of the Phased Progression of
858 Pluripotency. *Cell Systems* **8**, 427-445.e410 (2019).
- 859 21. U. Raudvere *et al.*, g:Profiler: a web server for functional enrichment analysis and
860 conversions of gene lists (2019 update). *Nucleic Acids Research* **47**, W191-W198 (2019).
- 861 22. I. Mohorianu, M. B. Stocks, J. Wood, T. Dalmay, V. Moulton, CoLide. *RNA Biology* **10**,
862 1221-1230 (2013).
- 863 23. V. N. Kim, J. Han, M. C. Siomi, Biogenesis of small RNAs in animals. *Nature Reviews*
864 *Molecular Cell Biology* **10**, 126-139 (2009).
- 865 24. F. Borges, R. A. Martienssen, The expanding world of small RNAs in plants. *Nature*
866 *Reviews Molecular Cell Biology* **16**, 727-741 (2015).
- 867 25. M. Ha, V. N. Kim, Regulation of microRNA biogenesis. *Nature Reviews Molecular Cell*
868 *Biology* **15**, 509-524 (2014).

- 869 26. B. Czech *et al.*, piRNA-Guided Genome Defense: From Biogenesis to Silencing. *Annual*
870 *Review of Genetics* **52**, 131-157 (2018).
- 871 27. R. K. Papareddy *et al.*, Chromatin regulates expression of small RNAs to help maintain
872 transposon methylome homeostasis in Arabidopsis. *Genome Biology* **21**, (2020).
- 873 28. A. S. E. Cuomo *et al.*, Single-cell RNA-sequencing of differentiating iPS cells reveals
874 dynamic genetic effects on gene expression. *Nature Communications* **11**, (2020).
- 875 29. Y. Hao *et al.* (Cold Spring Harbor Laboratory, 2020).
- 876 30. E. Becht *et al.*, Dimensionality reduction for visualizing single-cell data using UMAP.
877 *Nature Biotechnology* **37**, 38-44 (2019).
- 878 31. R. Andersson, A. Sandelin, Determinants of enhancer and promoter activities of
879 regulatory elements. *Nature Reviews Genetics* **21**, 71-87 (2020).
- 880 32. M. Levo, E. Segal, In pursuit of design principles of regulatory sequences. *Nature Reviews*
881 *Genetics* **15**, 453-468 (2014).
- 882 33. D. Holoch, D. Moazed, RNA-mediated epigenetic regulation of gene expression. *Nature*
883 *Reviews Genetics* **16**, 71-84 (2015).
- 884 34. J. Thody, V. Moulton, I. Mohorianu, PAREameters: a tool for computational inference of
885 plant miRNA-mRNA targeting rules using small RNA and degradome sequencing data.
886 *Nucleic Acids Research* **48**, 2258-2270 (2020).
- 887 35. J. Thody *et al.*, PAREsnip2: a tool for high-throughput prediction of small RNA targets
888 from degradome sequencing data using configurable targeting rules. *Nucleic Acids*
889 *Research*, (2018).
- 890 36. C. E. Ang, M. Wernig, Profiling DNA-transcription factor interactions. *Nature*
891 *Biotechnology* **36**, 501-502 (2018).
- 892 37. S. Lanciano, G. Cristofari, Measuring and interpreting transposable element expression.
893 *Nature Reviews Genetics* **21**, 721-736 (2020).
- 894 38. H.-G. Drost, Philentropy: Information Theory and Distance Quantification with R.
895 *Journal of Open Source Software* **3**, 765 (2018).
- 896 39. Y. Liao, G. K. Smyth, W. Shi, featureCounts: an efficient general purpose program for
897 assigning sequence reads to genomic features. *Bioinformatics* **30**, 923-930 (2014).
- 898 40. A. Dobin *et al.*, STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21
899 (2013).
- 900 41. D. Kim, J. M. Paggi, C. Park, C. Bennett, S. L. Salzberg, Graph-based genome alignment
901 and genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology* **37**, 907-915
902 (2019).
- 903 42. B. Langmead, S. L. Salzberg, Fast gapped-read alignment with Bowtie 2. *Nature Methods*
904 **9**, 357-359 (2012).
- 905 43. C. Paicu *et al.*, miRCat2: accurate prediction of plant and animal microRNAs from next-
906 generation sequencing datasets. *Bioinformatics* **33**, 2446-2454 (2017).
- 907 44. T. Wallach *et al.*, Identification of CNS Injury-Related microRNAs as Novel Toll-Like
908 Receptor 7/8 Signaling Activators by Small RNA Sequencing. *Cells* **9**, 186 (2020).
- 909 45. A. Kozomara, M. Birgaoanu, S. Griffiths-Jones, miRBase: from microRNA sequences to
910 function. *Nucleic Acids Research* **47**, D155-D162 (2019).
- 911 46. T. Z. Berardini *et al.*, The arabidopsis information resource: Making and mining the “gold
912 standard” annotated reference plant genome. *genesis* **53**, 474-485 (2015).
- 913 47. P. Ewels, M. Magnusson, S. Lundin, M. Käller, MultiQC: summarize analysis results for
914 multiple tools and samples in a single report. *Bioinformatics* **32**, 3047-3048 (2016).
- 915 48. I. T. Jolliffe, J. Cadima, Principal component analysis: a review and recent developments.
916 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and*
917 *Engineering Sciences* **374**, 20150202 (2016).

- 918 49. M. Kanehisa, KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids*
919 *Research* **28**, 27-30 (2000).
- 920 50. G. Viteri *et al.*, Reactome and ORCID—fine-grained credit attribution for community
921 curation. *Database* **2019**, (2019).
- 922 51. A. Pratapa, A. P. Jalihal, J. N. Law, A. Bharadwaj, T. M. Murali, Benchmarking
923 algorithms for gene regulatory network inference from single-cell transcriptomic data.
924 *Nature Methods* **17**, 147-154 (2020).
- 925 52. V. A. Huynh-Thu, A. Irrthum, L. Wehenkel, P. Geurts, Inferring Regulatory Networks
926 from Expression Data Using Tree-Based Methods. *PLoS ONE* **5**, e12776 (2010).
- 927 53. T. Moerman *et al.*, GRNBoost2 and Arboreto: efficient and scalable inference of gene
928 regulatory networks. *Bioinformatics* **35**, 2159-2161 (2019).
- 929 54. T. E. Chan, M. P. H. Stumpf, A. C. Babbie, Gene Regulatory Network Inference from
930 Single-Cell Data Using Multivariate Information Measures. *Cell Systems* **5**, 251-267.e253
931 (2017).
- 932 55. E. Dimitriadou, K. Hornik, F. Leisch, D. Meyer, A. Weingessel. (2009), vol. 1.
- 933 56. M. Martin, Cutadapt removes adapter sequences from high-throughput sequencing reads.
934 *EMBnet.journal* **17**, 10 (2011).
- 935 57. M. B. Stocks *et al.*, The UEA sRNA workbench: a suite of tools for analysing and
936 visualizing next generation sequencing microRNA and small RNA datasets.
937 *Bioinformatics* **28**, 2059-2061 (2012).
- 938 58. G. X. Y. Zheng *et al.*, Massively parallel digital transcriptional profiling of single cells.
939 *Nature Communications* **8**, 14049 (2017).
- 940 59. W. J. Kent *et al.*, The Human Genome Browser at UCSC. *Genome Research* **12**, 996-1006
941 (2002).

942

943 **Acknowledgments**

944 IM, EL, EW and IIM acknowledge the constructive feedback from the Core Bioinformatics
945 group and the support by the Core grant awarded to the Wellcome-MRC Cambridge Stem
946 Cells Institute. SAVU, LM and HGD acknowledge that their work was supported by the
947 Max Planck Society. IM and IIM designed the study and implemented the R package
948 **noisyR**; SAVU, LM and HGD implemented the R package *edgynode*, the analyses were
949 performed by IM, IIM (the mRNA bulk and single cell analyses), EW, IIM (sRNA analysis),
950 LM and HGD (GRN inference), EL, IM, IIM (comparison across tools). IM, HGD and IIM
951 wrote the manuscript; all authors read and approved the submitted manuscript.

952