Selecting genes for analysis using historically contingent progress

Farhaan Lalit¹, Antony M Jose^{1*}

Affiliations:

¹University of Maryland, College Park, MD, USA.

*Corresponding author. Email: amjose@umd.edu

Author Contributions: A. M. J. designed the study; F. L. and A. M. J. performed the analyses; and F. L. and A. M. J. wrote the paper.

Competing Interest Statement: The authors declare no competing interests.

Keywords: Mutual Information, RNA silencing, gene regulation, *C. elegans*.

This file includes:

Main Text

Figures 1 and 2

SUMMARY

Progress in biology has generated numerous lists of genes that share some property. But, advancing from the initial implication of a set of genes in a process to understanding their roles in the process is slow and unsystematic. Here we use RNA silencing in *C. elegans* to illustrate a general approach for comparing lists of data accumulated by a field to prioritize genes for detailed study given limited resources. The partially subjective relationships between genes forged by both functional relatedness of the genes and biased progress in the field was captured as historical mutual information (HMI) and used as a quantitative measure for clustering genes. These clusters suggest regulatory links connecting RNA silencing with other processes like the cell cycle and identify understudied regulated genes that could be used to sense perturbation or mediate feedback inhibition.

MAIN TEXT

Introduction

Genes and gene products are often collected as lists based on unifying characteristics or based on experiments. For example, genes that show enrichment of a chromatin modification, mRNAs that change in response to a mutation, proteins that interact with another protein, etc. After the initial identification of a set of genes as belonging to a list, multiple approaches [1] are needed to generate an explanatory model. Since single papers often analyze only one or a few genes, a wider view of genes with roles in a process could be gained by comparing lists generated by a group of studies or even all studies in a relatively young field. Such exploration could identify genes that are present in multiple lists but have not yet been selected for detailed study. Identifying these understudied genes is especially useful during the early stages of a field, when coherent models that provide explanations for most observed phenomena have not yet emerged. While this approach is also extensible to lists of anything that is used to characterize living

systems (changes in lipids, metabolites, localizations, etc.), here we focus on lists of mRNAs, proteins, and small RNAs generated by the field of RNA silencing in the nematode *C. elegans*.

A gene present in many lists could be regulated in multiple separable ways and/or be regulated in one or a few ways by connected sets of regulators (Fig. 1A). For example, mRNA levels could be regulated through changes in transcription, turnover, localization, small RNA production, etc. or all changes could occur because of turnover regulation by a connected set of regulators. Changes in such genes could alter specific regulatory outputs, making them integrators of inputs from many other regulators. Alternatively, they could have no measurable consequence, making them experimentally useful as general sensors of perturbation. Here we present an approach to identify these regulated but understudied genes in the field of RNA silencing in *C. elegans*.

Results

To determine if there are any understudied regulated genes that are relevant for RNA silencing in C. elegans, we examined data from past studies in the field. While complete replication of each study might be needed for direct comparisons, this ideal is impractical. Even beginning with the 'raw' data deposited to public resources (e.g., fastq files after RNA-seq) and repeating the analyses reported in a paper is not always feasible. Summary tables from previous analyses presented in papers provide a practical intermediate level of data to use for comparisons across studies. Therefore, we collated a total of 432 tables from 112 papers for comparison (see methods and Table S1 for list of studies) and joined the tables together after standardizing gene names to yield genes that can be compared for presence or absence across 432 lists (Fig. 1B). To identify a set of genes (g) that receive extensive regulatory input and/or that encode proteins that interact with many other proteins and are yet selectively regulated, we propose a metric r_g (Fig. 1B). Since the likelihood of including a gene from the lists increases with g, the metric is

specified with a subscript for each analysis (e.g., r_{25} refers to a regulation score when the top 25 genes that are most commonly present in lists are considered) and defined to be:

$$r_g \coloneqq \sum_{i=1}^n \frac{S_i}{T_i}$$

where g = size of gene set chosen for analysis, n = total number of lists with altered genes, S_i = number of genes from the ith list that is also present in the gene set g, and T_i = total number of genes in the ith list. The larger the set of genes (g) selected, the greater the chance of a dataset (with T_i genes) having at least one overlapping gene within the selected gene set (probability given by $P(S_i > 0)$ in Fig. 1C). The metric r_g is a decision aid that helps with choosing genes for experimental analysis and is not to be taken as an objective measure of the importance of the gene for the biological process under study.

The top 25 genes with the highest r_{25} values included the germline Argonaute proteins CSR-1 [2] and HRDE-1 [3], which have each been the subject of numerous studies (Fig. 1*D*). While most other genes are understudied (fewer than 10 publications on WormBase), among them is W09B7.2/sdg-1, which was recently reported to be regulated by the double-stranded RNA importer SID-1 and encodes a protein with a suggested role in feedback regulation of heritable RNA silencing by colocalizing with perinuclear germ granules [4]. This discovery suggests that the analysis of the additional genes with high r_{25} values could also be fruitful. Of the 16 understudied genes that encode proteins, 7 had high-confidence AlphaFold structures [5], which were then used to identify related protein domains using Foldseek [6] (Fig. 1*E*). Three more proteins have been proposed to be nucleocapsids encoded from genes within retrotransposons ([4, 7]; Fig. 1*E*). These candidates can be experimentally analyzed in the future for roles in RNA silencing, if any.

To explore the relationships between these genes with the highest r_{25} values (Fig. 1*F* and 1*G*), we clustered the genes and generated a dendrogram where genes present together in

different lists are closer together (see supplementary methods). The dendrogram revealed the gene *hil-4*, which encodes a Histone H1-like protein [8], as the understudied gene clustering closest to *hrde-1* and *csr-1*, making it a strong candidate for a role potentially downstream of RNA-mediated gene regulation. Another cluster (brown in Fig. 1*G*) included all four pseudogenes, suggesting that this method could capture functional relatedness despite the limitations and biases introduced by the available data.

To examine if the observations using r_{25} hold when analyzing a larger set of genes, we examined the top 100 genes with the highest r_{100} values. To quantify the correlated presence or absence of genes in different lists we used a measure of mutual information [9] named here as historical mutual information (HMI) to emphasize the subjective nature of this measure because it depends on both functional relatedness of the genes and biased availability of data (see supplementary methods). Using HMI to cluster these genes revealed three major clusters (43, 42, and 11 genes), another cluster with two genes and two other unconnected genes (Fig. 2A). Only one cluster (cluster 1 in Fig. 2A) had significant numbers of genes associated with gene ontology terms. Many of these genes encode proteins that bind and/or hydrolyze RNA (Fig. 2B, top), localize to cytoplasmic ribonucleoprotein granules (Fig. 2B, middle), and/or play roles in other processes such as cell division (Fig. 2B, bottom). Consistently, this cluster also had the greatest number of genes that have been described in multiple publications (Fig. 2C), including all the genes that have been featured in abstracts on RNA silencing (Fig. 2D). Therefore, the analysis of additional genes in this cluster could be relevant for RNA silencing and connect such regulation to other processes (e.g., the cell cycle). Since four of the five pseudogenes are in a small cluster (Fig. 2E, 4 of 11 genes in cluster 2), the other genes in this cluster could potentially be targets of regulation without specific downstream regulation or be co-regulated sensors of pseudogene RNA levels. Intriguingly, there is a large overlap between a set of genes that require HRDE-1 for downregulation (67 genes in both replicates from worms grown at 15°C [10]) and genes in a single cluster (Fig. 2F, 17 of 42 genes in cluster 3). One possible explanation for this abundance and

clustering could be that *hrde-1*-dependent gene lists are among the most numerous generated by the field and/or included in our analysis (44 of 298 lists with fewer than 2000 genes). Alternatively, genes that are subject to HRDE-1-dependent silencing could be extensively regulated by many other regulators and require this additional downregulation for fitness – i.e., overexpression of these genes is detrimental. Consistent with this possibility, loss of HRDE-1 results in progressive sterility that can be reversed by restoring HRDE-1 activity [10]. Also, as expected for the use of HRDE-1 downstream of SID-1, genes upregulated using *sid-1* (18 genes in animals with a deletion in *sid-1* [4]) overlap with genes in the same cluster (Fig. 2F, 4 of 42 in cluster 3). Future studies by labs working on RNA silencing in *C. elegans* have the potential to test and enrich the classification of regulated yet understudied genes revealed here.

Discussion

Our analysis has identified selectively regulated yet understudied genes in the field of RNA silencing in *C. elegans*. Clustering these genes reveal that better studied genes are together in one cluster. Many of these genes have known roles, providing regulatory links between RNA silencing and other processes. The other two larger clusters include genes and pseudogenes that have been described as targets of RNA regulation. These regulated genes, which are present in many lists, could be under the independent control of multiple regulators and/or be jointly regulated by a connected set of regulators, making them experimentally useful general sensors of perturbation in RNA silencing and/or regulators that mediate feedback inhibition of RNA silencing.

While future extensions of this work could automate the process of aggregating and comparing data, flexible inclusion of different lists in the analysis would be needed to enable customization based on the expertise, interests, and risk tolerance of individual labs. Furthermore, earlier studies that had to use older technologies with limitations could have led to conclusions that need revision. For example, when analyzed using multi-copy transgenes, the dsRNA-binding protein RDE-4 showed a cell non-autonomous effect [11, 12], but when analyzed using single-

copy transgenes, RDE-4 showed a cell autonomous effect [13]. Since different researchers could interpret such conflicting data differently (e.g., differences in levels of tissue-restricted expression versus differences in extent of misexpression in other tissues), it is useful to preserve customization for the lists included.

Different properties of a single protein or RNA could be important for different biological roles [14, 15], or the same properties could be important for different processes. Despite such plurality, a gene found in many lists could become associated with a single label because of the historical sequence of discovery (e.g., HRDE-1-dependent genes; many in cluster 3, Fig. 2F), thereby obscuring additional roles of that gene. With the expanding number of lists generated through large-scale experimental approaches in different fields, identifying selectively regulated yet understudied genes could aid the prioritization of genes for detailed mechanistic studies using the limited resources and time available for any lab.

Methods

Data tables from 112 studies on RNA silencing in *C. elegans* that were published between 2007 and 2023 were downloaded (Table S1), reformatted manually and/or using custom scripts, and filtered to generate lists that only include entries with reported p-values or adjusted p-values < 0.05, when such values were available. The top 'g' genes that occur in the greatest numbers of tables were culled as the most frequently identified genes. A measure for the extent of regulation of each gene (r_g) was used to aid prioritization for detailed study. Co-occurrence patterns of genes in different tables were captured using the Jaccard distance (d_J) [16] or as a symmetric measure of normalized mutual information [9], defined here as Historical Mutual Information (HMI). The d_J values were used to generate a dendrogram using the average linkage method (Fig. 1 G). HMI was used to group genes into clusters according to the Girvan-Newman algorithm [17] and different sets of genes were highlighted (Fig. 2). Gene ontology (GO) analyses were performed

using Gene Ontology Resource (https://geneontology.org/; [18, 19]) and significant GO terms were collected for visualization using REVIGO [20]. All programs used in this study are available at GitHub (AntonyJose-Lab/Lalit_Jose_2024).

Acknowledgements

We thank Tom Kocher and members of the Jose lab for comments on the manuscript. This work is supported in part by National Institutes of Health Grant R01GM124356 and National Science Foundation Grant 2120895 to A.M.J.

References

- 1. A. M. Jose, The analysis of living systems can generate both knowledge and illusions. *Elife* **9** (2020).
- 2. J. M. Claycomb *et al.*, The Argonaute CSR-1 and its 22G-RNA cofactors are required for holocentric chromosome segregation. *Cell* **139**, 123-134 (2009).
- 3. B. A. Buckley *et al.*, A nuclear Argonaute promotes multigenerational epigenetic inheritance and germline immortality. *Nature* **489**, 447-451 (2012).
- 4. N. Shugarts *et al.*, SID-1 regulates a retrotransposon-encoded gene to tune heritable RNA silencing. *bioRxiv* 10.1101/2021.10.05.463267 (2023).
- 5. J. Jumper *et al.*, Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583-589 (2021).
- 6. M. van Kempen *et al.*, Fast and accurate protein structure search with Foldseek. *Nature Biotechnology* 10.1038/s41587-023-01773-0 (2023).
- 7. S. E. J. Fischer, G. Ruvkun, Caenorhabditis elegans ADAR editing and the ERI-6/7/MOV10 RNAi pathway silence endogenous viral elements and LTR retrotransposons. *Proc Natl Acad Sci U S A* **117**, 5987-5996 (2020).
- 8. M. A. Jedrusik, E. Schulze, A single histone H1 isoform (H1.1) is essential for chromatin silencing and germline development in Caenorhabditis elegans. *Development* **128**, 1069-1080 (2001).
- 9. I. H. Witten, E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques* (Morgan Kaufmann, San Francisco, ed. 2nd, 2005).
- 10. J. Z. Ni *et al.*, A transgenerational role of the germline nuclear RNAi pathway in repressing heat stress-induced transcriptional activation in *C. elegans. Epigenetics Chromatin* **9**, 3 (2016).
- 11. A. M. Jose, G. A. Garcia, C. P. Hunter, Two classes of silencing RNAs move between Caenorhabditis elegans tissues. *Nat Struct Mol Biol* **18**, 1184-1188 (2011).
- 12. D. Blanchard *et al.*, On the nature of in vivo requirements for rde-4 in RNAi and developmental pathways in C. elegans. *RNA Biol* **8**, 458-467 (2011).
- 13. P. Raman, S. M. Zaghab, E. C. Traver, A. M. Jose, The double-stranded RNA binding protein RDE-4 can act cell autonomously during feeding RNAi in C. elegans. *Nucleic Acids Res* **45**, 8463-8473 (2017).
- 14. C. E. Chapple *et al.*, Extreme multifunctional proteins identified from a human protein interaction network. *Nat Commun* **6**, 7412 (2015).
- 15. D. M. Ribeiro, G. Briere, B. Bely, L. Spinelli, C. Brun, MoonDB 2.0: an updated database of extreme multifunctional and moonlighting proteins. *Nucleic Acids Res* **47**, D398-D402 (2019).
- 16. P. Jaccard, Distribution de la flore alpine dans le bassin des dranses et dans quelques régions voisines. *Bulletin de la Société vaudoise des sciences naturelles* **37**, 241-272 (1901).
- 17. M. Girvan, M. E. Newman, Community structure in social and biological networks. *Proc Natl Acad Sci U S A* **99**, 7821-7826 (2002).
- 18. M. Ashburner *et al.*, Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet* **25**, 25-29 (2000).

- 19. G. O. Consortium et al., The Gene Ontology knowledgebase in 2023. Genetics 224 (2023).
- 20. F. Supek, M. Bosnjak, N. Skunca, T. Smuc, REVIGO summarizes and visualizes long lists of gene ontology terms. *PLoS One* **6**, e21800 (2011).

Figures and Figure Legends

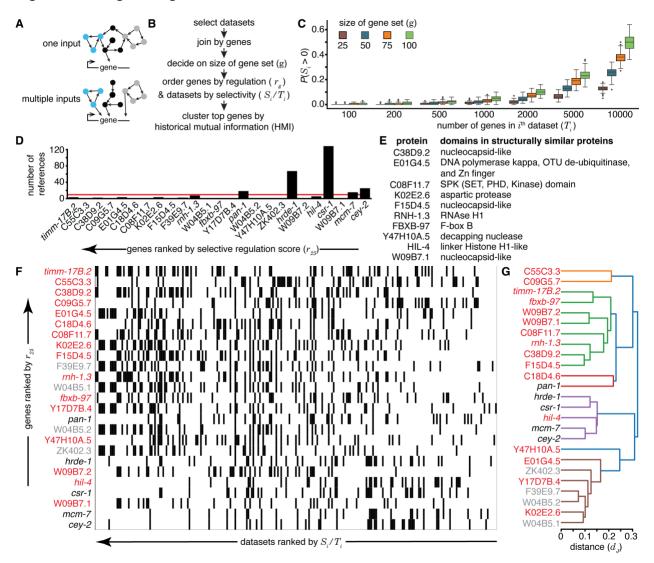


Figure 1. Some genes are selectively regulated, reported as part of many lists, and yet are understudied. (A) Schematics of possible regulatory architectures for genes found on multiple lists. (top) Gene receiving one input form a large network. (bottom) Gene receiving multiple inputs from separable networks. (B) Strategy for the identification of regulated genes. See Methods for details. (C) Relationship between S_i , T_i , and g obtained using simulated data from an organism with 20,000 genes. Distributions of 100 runs for each parameter combination are presented as box and whisker plots. (D) Numbers of publications listed on WormBase for the top 25 regulated genes as measured using r_{25} in the field of RNA silencing in C. elegans. Red line marks 10

publications. (*E*) Domains present in proteins encoded by understudied genes among the top 25 genes that are suggestive of function. Proteins with high-confidence AlphaFold structures [5] were used to identify related proteins using Foldseek [6] or based on the literature ([4, 7]; C38D9.2, F15D4.5, and W09B7.2). (*F*) Heat map showing the top 25 regulated genes. Presence (black) or absence (white) of each gene in each dataset is indicated. Relatively understudied (<10 references on WormBase) genes (red) or pseudogenes (grey) identified in (*D*) are indicated. (*G*) Hierarchical clustering of the top 25 genes based on co-occurrence in studies, where gene names colored as in (*D*) and 'distance (d_x)' indicates Jaccard distance.

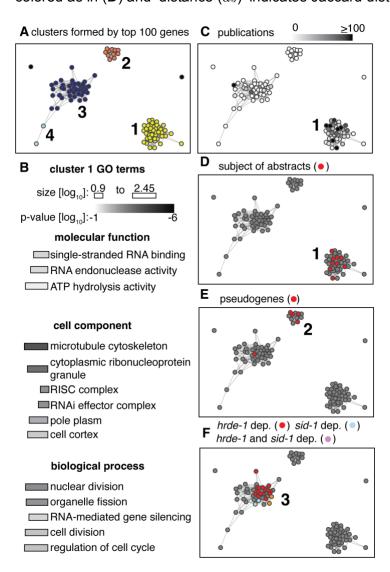


Figure 2. Clusters formed by regulated genes suggest priorities for detailed study. (A and B) Properties of the top 100 regulated genes as measured using r_{100} in the field of RNA silencing in C. elegans. (A) Clusters of genes based on their historical mutual information. Threshold for link: HMI > 0.9. Also see Table S2. (B) Molecular functions (top), cell components (middle), and biological processes (bottom) of genes in cluster 1 as in (A). Length of boxes near each term in (B) indicates log_{10} (annotations for GO term in C. elegans), with largest and smallest bars indicating ~285 and ~8 annotations, and shading indicates $-log_{10}$ (Bonferroni-corrected p-value), with black and white indicating a p-value of ~10-6 and ~10-1, respectively. (C-F) Network in (A) with nodes colored to show number of publications per gene (white, 0; black, ≥ 100) (C), genes that have been the main subject of abstracts on RNA silencing in C. elegans (D), pseudogenes (red) (E), and genes changed in hrde-1 mutants [10] (red), a sid-1 mutant [4] (blue), or both (orange) (F).

Supplementary Information

Selecting genes for analysis using historically contingent progress

Farhaan Lalit¹, Antony M Jose^{1*}

This file includes:

Supplementary Methods

SI References

2 Supplementary Tables

Supplementary Methods

To identify studies on RNA silencing in C. elegans that have data tables that can be compared across all studies, we used the term 'C. elegans RNA silencing' to search PubMed. After examining the more than 2000 studies that resulted from the search, the available data tables from 112 studies that were published between 2007 and 2023 were downloaded (Table S1), reformatted into 432 distinct tables manually and/or using custom scripts. Metadata if supplied by the authors for each table were retained as comments above each table. Gene names were unified using the Gene Name Sanitizer (https://wormbase.org/tools/mine/gene sanitizer.cgi) as on 26 April 2022. It is unclear how an exhaustive list of papers that is nevertheless fieldrestricted could ever be defined for any field. Accordingly, our list of RNA silencing studies in C. elegans is not exhaustive and we apologize to colleagues whose work is not included in our analysis. Nevertheless, this effort captured additional datasets compared with those available in other more unrestricted collections that attempt to collect tables from all studies on an organism (e.g. WormExp 2.0 [1]). Only 29 studies included in this study overlapped with the 461 included in WormExp 2.0 as on 27 Jan 2023, which was determined by comparing the paper IDs after downloading all datasets from https://wormexp.zoologie.uni-kiel.de/wormexp/ using a tool on WormBase (http://tazendra.caltech.edu/~azurebrd/cgibin/forms/generic.cgi?action=PapIdToWBPaper). Data tables that reported p-values or adjusted p-values were filtered to only include entries with p < 0.05. Since fold-changes were not always available, for every dataset, genes were scored as present or absent to generate a heatmap featuring the most frequently changed genes (highest values of r_g), where the number of genes considered (g) can be arbitrary (e.g., 25 in Fig. 1F and 100 in Fig. 2). The relationships between the parameters S_i , T_i , and g (Fig. 1C) were obtained using simulated data by sampling 100 random sets of genes as the top q genes from a total of 20,000 genes and similarly sampling the genes in datasets of various sizes (T_i) . For each gene, the number of references listed on

Wormbase (https://wormbase.org/) was used as a measure of the extent to which the gene has been studied. Genes with fewer than 10 references were defined as understudied (Fig. 1*D*). To generate the heatmap, genes were ordered in decreasing values of r_{25} (top to bottom in Fig. 1*F*) and datasets were ordered in decreasing values of $\frac{S_i}{T_i}$. (left to right in Fig. 1*F*). To determine the co-occurrence patterns of all pairs of genes, Jaccard distances ($d_J = 1 - \frac{|X \cap Y|}{|X \cup Y|}$) where X and Y are sets of lists containing genes x and y, respectively) were calculated for each pair and all genes were hierarchically clustered using the 'average' linkage method. Relationships between genes based on occurrence in datasets were also captured as normalized mutual information and defined as historical mutual information (HMI) to emphasize the dependence on the biased availability of data based on historical progress in addition to the functional relatedness of the genes. Specifically, it was defined to be a symmetric and normalized mutual information score [2] and was calculated using the function normalized_mutual_info_score from scikit-learn [3] for genes X and Y:

$$HMI(X;Y) := \frac{2.MI(X;Y)}{H(X) + H(Y)},$$

where $MI(X;Y) = \sum_{y} \sum_{x} P_{(X,Y)}(x,y) \log_2\left(\frac{P_{(X,Y)}(x,y)}{P_X(x)P_Y(y)}\right)$, $H(X) = -\sum_{x} P(x)\log_2\left(P(x)\right)$, and $H(Y) = -\sum_{y} P(y)\log_2\left(P(y)\right)$. Mutual information (MI) determines how different the joint distribution of the gene pair (X, Y) is from the product of the marginal distributions of each gene, H(X) and H(Y) are the entropies of the two genes, and P(...) indicates probabilities. Clusters of genes based on HMI values were identified using the Girvan-Newman algorithm [4]. Gene Ontology (GO) analysis was performed on all clusters using the Gene Ontology Resource ([5, 6]; https://geneontology.org/) and the significant terms (selected as having P < 0.05 after Bonferroni correction for multiple testing, associated with > 3 genes, and with a > 3-fold enrichment), if any, for each cluster were reduced for visualization using REVIGO ([7]; http://revigo.irb.hr/) with the organism set to Caenorhabditis elegans, the resulting list size set to 'small', and displaying only terms with

frequency < 3% (selects for more specific terms). The interactive graphical user interface (GUI) for visualizing clusters and genes of interest was created using Dash (Python).

Table S1. Published papers from which tables were used for this study

Table 51. Published papers from which tables were used for this study				
Paper	Pubmed			
2007 Welker et al RNA	https://pubmed.ncbi.nlm.nih.gov/17526642/			
2007 Zhang et al Mol Cell	https://pubmed.ncbi.nlm.nih.gov/18042455/			
2008 Batista et al Mol Cell	https://pubmed.ncbi.nlm.nih.gov/18571452/			
2008 Spike et al Development	https://pubmed.ncbi.nlm.nih.gov/18234720/			
2008 Wang et al Curr Biol	https://pubmed.ncbi.nlm.nih.gov/18501605/			
2009 Claycomb et al Cell	https://pubmed.ncbi.nlm.nih.gov/19804758/			
2009 Gent et al Genetics	https://pubmed.ncbi.nlm.nih.gov/19805814/			
2009 Gu et al Mol Cell	https://pubmed.ncbi.nlm.nih.gov/19800275/			
2009 Han et al PNAS	https://pubmed.ncbi.nlm.nih.gov/19846761/			
2009 vanWolfswinkel et al Cell	https://pubmed.ncbi.nlm.nih.gov/19804759/			
2010 Conine et al PNAS	https://pubmed.ncbi.nlm.nih.gov/20133686/			
2010 Correa et al PLoS Genet	https://pubmed.ncbi.nlm.nih.gov/20386745/			
2010 Vasale et al PNAS	https://pubmed.ncbi.nlm.nih.gov/20133583/			
2010 Welker et al RNA	https://pubmed.ncbi.nlm.nih.gov/20354150/			
2011 Fischer et al PLoS Genet	https://pubmed.ncbi.nlm.nih.gov/22102828/			
2011 Maniar et al Curr Biol	https://pubmed.ncbi.nlm.nih.gov/21396820/			
2011 Thivierge et al NSMB	https://pubmed.ncbi.nlm.nih.gov/22179787/			
2011 Wu et al NSMB	https://pubmed.ncbi.nlm.nih.gov/21909095/			
2011 Zhang et al PNAS	https://pubmed.ncbi.nlm.nih.gov/21245313/			
2012 Bagijn et al Science	https://pubmed.ncbi.nlm.nih.gov/22700655/			
2012 Buckley et al Nature	https://pubmed.ncbi.nlm.nih.gov/22810588/			
2012 Gu et al Cell	https://pubmed.ncbi.nlm.nih.gov/23260138/			
2012 Lee et al Cell	https://pubmed.ncbi.nlm.nih.gov/22738724/			
2012 Warf et al Genome Res	https://pubmed.ncbi.nlm.nih.gov/22673872/			
2012 Zhang et al Curr Biol	https://pubmed.ncbi.nlm.nih.gov/22542102/			
2013 Conine et al Cell	https://pubmed.ncbi.nlm.nih.gov/24360276/			
2013 Hall et al RNA	https://pubmed.ncbi.nlm.nih.gov/23329696/			
2013 Sarkies et al Genome Res	https://pubmed.ncbi.nlm.nih.gov/23811144/			
2014 Cecere et al NSMB	https://pubmed.ncbi.nlm.nih.gov/24681887/			
2014 Kasper et al Dev Cell	https://pubmed.ncbi.nlm.nih.gov/25373775/			
2014 Ni et al BMC Genomics	https://pubmed.ncbi.nlm.nih.gov/25534009/			
2014 Phillips et al Curr Biol	https://pubmed.ncbi.nlm.nih.gov/24684932/			
2014 Rechavi et al Cell	https://pubmed.ncbi.nlm.nih.gov/25018105/			

0044.0.11: 1.1.1.1.1.1.0	
2014 Sakaguchi et al PNAS	https://pubmed.ncbi.nlm.nih.gov/25258416/
2014 Stoeckius et al EMBO J	https://pubmed.ncbi.nlm.nih.gov/24957527/
2014 Weick et al Genes n Dev	https://pubmed.ncbi.nlm.nih.gov/24696457/
2014 Yang et al Curr Biol	https://pubmed.ncbi.nlm.nih.gov/24684930/
2014 Zhou et al Genetics	https://pubmed.ncbi.nlm.nih.gov/24532782/
2015 Albuquerque et al Dev Cell	https://pubmed.ncbi.nlm.nih.gov/26279485/
2015 Phillips et al Dev Cell	https://pubmed.ncbi.nlm.nih.gov/26279487/
2015 Tsai et al Cell	https://pubmed.ncbi.nlm.nih.gov/25635455/
2015 Tu et al Nucl Acids Res	https://pubmed.ncbi.nlm.nih.gov/25510497/
2015 Zinovyeva et al PNAS	https://pubmed.ncbi.nlm.nih.gov/26351692/
2016 Gerson-Gurwitz et al Cell	https://pubmed.ncbi.nlm.nih.gov/27020753/
2016 Houri-Zeevi et al Cell	https://pubmed.ncbi.nlm.nih.gov/27015309/
2016 Ni et al Epigenetics Chromatin	https://pubmed.ncbi.nlm.nih.gov/26779286/
2016 Tang et al Cell	https://pubmed.ncbi.nlm.nih.gov/26919432/
2017 Akay et al Dev Cell	https://pubmed.ncbi.nlm.nih.gov/28787591/
2017 Andralojc et al PLoS Genet	https://pubmed.ncbi.nlm.nih.gov/28182654/
2017 Brown et al Nucl Acids Res	https://pubmed.ncbi.nlm.nih.gov/28645154/
2017 Kalinava et al Epigenetics Chromatin	https://pubmed.ncbi.nlm.nih.gov/28228846/
2017 Lev et al Curr Biol	https://pubmed.ncbi.nlm.nih.gov/28343968/
2017 Tyc et al Dev Cell	https://pubmed.ncbi.nlm.nih.gov/28787592/
2017 Weiser et al Dev Cell	https://pubmed.ncbi.nlm.nih.gov/28535375/
2018 Almeida et al EMBO J	https://pubmed.ncbi.nlm.nih.gov/29769402/
2018 Davis et al elife	https://pubmed.ncbi.nlm.nih.gov/30575518/
2018 Newman et al Genes and Dev	https://pubmed.ncbi.nlm.nih.gov/29739806/
2018 Reich et al Genes and Dev	https://pubmed.ncbi.nlm.nih.gov/29483152/
2018 Uebel et al PLoS Genet	https://pubmed.ncbi.nlm.nih.gov/30036386/
2018 Xu et al Cell Rep	https://pubmed.ncbi.nlm.nih.gov/29791857/
2019 Almeida et al PLoS Genet	https://pubmed.ncbi.nlm.nih.gov/30759082/
2019 Bezler et al PLoS Genetics	https://pubmed.ncbi.nlm.nih.gov/30735500/
2019 Dodson et al Dev Cell	https://pubmed.ncbi.nlm.nih.gov/31402284/
2019 Gushchanskaia et al Neucl Acids Res	https://pubmed.ncbi.nlm.nih.gov/31216042/
2019 Lev et al Curr Biol	https://pubmed.ncbi.nlm.nih.gov/31378614/
2019 Marnik et al Genetics	https://pubmed.ncbi.nlm.nih.gov/31506335/
2019 Ouyang et al Dev Cell	https://pubmed.ncbi.nlm.nih.gov/31402283/
2019 Posner et al Cell	https://pubmed.ncbi.nlm.nih.gov/31178120/
2019 Svendsen et al Cell Reports	https://pubmed.ncbi.nlm.nih.gov/31801082/
2019 Zeng et al Cell Reports	https://pubmed.ncbi.nlm.nih.gov/31216475/
2020 Barucci et al Nat Cell Biol	https://pubmed.ncbi.nlm.nih.gov/32015436/
2020 Esse et al Cells	https://pubmed.ncbi.nlm.nih.gov/32781660/
2020 Fischer Ruvkun PNAS	https://pubmed.ncbi.nlm.nih.gov/32123111/
-	, ,

2020 Houri Zeevi et al Cell
2020 Lewis et al Mol Cell
2020 Manage et al eLife
2020 Mao et al PloS Biol
2020 Mattout et al Nat Cell Biol
2020 Reed et al Nucl Acids Res
2020 Schwartz-Orbach et al eLife
2020 Shukla et al Nature
2020 Suen et al Nat Commun
2020 Wan et al Genetics
2021 Charlesworth et al Nucl Acids Res
2021 Chaves et al Mol Cell
2021 Choi et al eLife
2021 Cornes et al Dev Cell
2021 Gudipati et al Mol Cell
2021 Houri-Zeevi et al eLife
2021 Kim et al eLife
2021 Kim et al eLife
2021 Montgomery et al Cell Reports
2021 Nguyen et al Nat Commun
2021 Placentino et al EMBO J
2021 Price et al eLife
2021 Qi et al Nat Commun
2021 Quarato et al Nat Commun
2021 Singh et al Nat Commun
2021 Spichal et al Nat Commun
2021 Wahba et al Dev Cell
2021 Wan et al EMBO J
2022 Chen et al Nat Commun
2022 Cornes et al Dev Cell
2022 Dai et al Cell Rep
2022 Efstathiou et al Nat Cell Biol
2022 Garrigues et al G3
2022 Hebbar et al Sci Rep
2022 Marnik et al PLoS Genet
2022 Wang et al Cell Rep
2023 Du et al Cell Rep
2023 Liontis et al BBA Adv
2023 Seroussi et al eLife

https://pubmed.ncbi.nlm.nih.gov/32841602/ https://pubmed.ncbi.nlm.nih.gov/32348780/ https://pubmed.ncbi.nlm.nih.gov/32338603/ https://pubmed.ncbi.nlm.nih.gov/33264285/ https://pubmed.ncbi.nlm.nih.gov/32251399/ https://pubmed.ncbi.nlm.nih.gov/31872227/ https://pubmed.ncbi.nlm.nih.gov/32804637/ https://pubmed.ncbi.nlm.nih.gov/32499657/ https://pubmed.ncbi.nlm.nih.gov/32843637/ https://pubmed.ncbi.nlm.nih.gov/33055090/ https://pubmed.ncbi.nlm.nih.gov/34329465/ https://pubmed.ncbi.nlm.nih.gov/33378643/ https://pubmed.ncbi.nlm.nih.gov/33587037/ https://pubmed.ncbi.nlm.nih.gov/34921763/ https://pubmed.ncbi.nlm.nih.gov/33852894/ https://pubmed.ncbi.nlm.nih.gov/33729152/ https://pubmed.ncbi.nlm.nih.gov/34003109/ https://pubmed.ncbi.nlm.nih.gov/34003111/ https://pubmed.ncbi.nlm.nih.gov/34879267/ https://pubmed.ncbi.nlm.nih.gov/34244496/ https://pubmed.ncbi.nlm.nih.gov/33231880/ https://pubmed.ncbi.nlm.nih.gov/34730513/ https://pubmed.ncbi.nlm.nih.gov/33627668/ https://pubmed.ncbi.nlm.nih.gov/33664268/ https://pubmed.ncbi.nlm.nih.gov/34108460/ https://pubmed.ncbi.nlm.nih.gov/33658512/ https://pubmed.ncbi.nlm.nih.gov/34388368/ https://pubmed.ncbi.nlm.nih.gov/33438773/ https://pubmed.ncbi.nlm.nih.gov/36085149/ https://pubmed.ncbi.nlm.nih.gov/34921763/ https://pubmed.ncbi.nlm.nih.gov/36070689/ https://pubmed.ncbi.nlm.nih.gov/36471127/ https://pubmed.ncbi.nlm.nih.gov/35088854/ https://pubmed.ncbi.nlm.nih.gov/35504914/ https://pubmed.ncbi.nlm.nih.gov/35657999/ https://pubmed.ncbi.nlm.nih.gov/36516753/ https://pubmed.ncbi.nlm.nih.gov/37505984/ https://pubmed.ncbi.nlm.nih.gov/37082252/ https://pubmed.ncbi.nlm.nih.gov/36790166/

Table S2. Clusters formed by $r_{\it 100}$ genes with HMI > 0.9

Cluster 1	Cluster 2	Cluster 3	Cluster 4	Unclustered
wago-4	Y37E11B.2	R03D7.2	C09G5.7	sea-2
par-5	H09G03.1	T02G5.4	C55C3.3	Y47H10A.5
eel-1	W04B5.2	fbxb-97		
egg-6	F39E9.7	pan-1		
mcm-7	ZK402.3	T20F7.1		
gfat-2	F39F10.4	fkb-8		
F34D10.4	Y17D7B.4	lin-15B		
pod-1	W05H12.2	bath-45		
ani-1	W04B5.1	W06A11.4		
spd-5	E01G4.5	timm-17B.2		
wago-1	K02E2.6	glit-1		
ima-3		elf-1		
ani-2		sdg-1		
mex-5		saeg-2		
mrp-4		ceh-20		
cdc-48.1		W09B7.1		
top-2		F40D4.13		
csr-1		F41G4.7		
hmg-12		C38C3.3		
tbb-2		rnh-1.3		
<u>simr-1</u>		C38D9.2		
idh-1		Y48G1BM.6		
hsp-90		F15D4.5		
pyk-1		citk-1		
cpg-1		Y20F4.4		
<u>rme-2</u>		F58H7.5		
puf-3		C04G6.6		
klp-15		R06C1.4		
hsp-4		saeg-1		
hrde-1		R03H10.6		
rpn-9		spe-41		
hsp-1		his-24		
tba-2		T16G12.4		
<u>pgl-3</u>		gly-13		
daf-18		clp-6		
mut-16		qdpr-1		
set-2		fbxa-192		
cey-2		C18D4.6		

 vig-1
 K09H9.7

 hil-4
 C08F11.7

 klp-7
 pdfr-1

 cdk-1
 scrm-4

deps-1

SI References

- 1. W. Yang, K. Dierking, H. Schulenburg, WormExp: a web-based application for a Caenorhabditis elegans-specific gene expression enrichment analysis. *Bioinformatics* **32**, 943-945 (2016).
- 2. I. H. Witten, E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques* (Morgan Kaufmann, San Francisco, ed. 2nd, 2005).
- 3. G. V. Fabian Pedregosa, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, Édouard Duchesnay, Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12, 2825-2830 (2011).
- 4. M. Girvan, M. E. Newman, Community structure in social and biological networks. *Proc Natl Acad Sci U S A* **99**, 7821-7826 (2002).
- 5. M. Ashburner *et al.*, Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet* **25**, 25-29 (2000).
- 6. G. O. Consortium *et al.*, The Gene Ontology knowledgebase in 2023. *Genetics* **224** (2023).
- 7. F. Supek, M. Bosnjak, N. Skunca, T. Smuc, REVIGO summarizes and visualizes long lists of gene ontology terms. *PLoS One* **6**, e21800 (2011).