

# A theory of genetic analysis using transcriptomic phenotypes

Thesis by  
David Angeles-Albores

In Partial Fulfillment of the Requirements for the  
degree of  
Doctor of Philosophy

The Caltech logo, consisting of the word "Caltech" in a bold, orange, sans-serif font.

CALIFORNIA INSTITUTE OF TECHNOLOGY  
Pasadena, California

[2018]  
Defended [18 September 2018]

© [2018]

David Angeles-Albores  
ORCID: 0000-0001-5497-8264

All rights reserved

## ACKNOWLEDGEMENTS

1

This thesis is possible thanks to the unwavering support of a long list of individuals. 2  
I would like to thank my advisor, Paul W. Sternberg, for providing a laboratory and 3  
an intellectual home for the past few years. I will always remember our conversations 4  
over coffee at 9am in the lab. I taught Systems Genetics with Paul as my project 5  
came into full maturity; these lectures on genetics gave us a chance to recast the 6  
classical interpretations with a genomics perspective. Those lectures are some of 7  
my fondest memories at Caltech. I also need to acknowledge my thesis committee, 8  
Dianne K. Newman, Elliot Meyerowitz and Matt Thomson. Without their advice, 9  
I would be considerably more confused than I am today. Dianne has been a major 10  
figure throughout my Ph.D., a great scientist with a heart to match, and I feel lucky 11  
to have had an opportunity to learn from her insights. I also need to thank Erich M. 12  
Schwarz, who taught me to argue, and taught me to write. Where others were content 13  
to say my work was fine, Erich found every possible loophole, every minor detail 14  
and every open question and pushed me to be complete without being redundant. I 15  
am sincerely grateful for his guidance and his mentorship. 16

I have been very lucky to have benefitted from a fantastic set of collaborators 17  
in the Sternberg lab. I have been lucky to work with Hillel Schwartz, a fantastic 18  
geneticist and good friend; Carmie Puckett Robinson, with whom I started to work on 19  
transcriptome genetics; Daniel Leighton, who taught me about worm pheromones 20  
and aging; Raymond Y. Lee and Juancarlos Chan, with whom I learned all the 21  
intricacies of WormBase and tool design. Throughout my time here, I have worked 22  
with three extremely talented undergraduates who made working in the lab much 23  
more exciting: Tiffany Tsou, Kyung Hoi (Joseph) Min and Vladimir Molchanov. 24  
Finally, I need to thank all the members of the Sternberg lab for making science 25  
come to life: Jon Liu, Han Wang, James Lee, Pei-Yin, Katie Brugman, Cynthia 26

Chai, Wen Chen, Sarah Cohen, Elizabeth Holman, Sandy Wong, Heenam Park, 27  
Daniel Jun Oh, Ravi Nath, Margaret Ho, Srimoyee Ghosh, Sophie Walton, Sarah 28  
Torres, Shahla Gharib, Barbara Perry, Animesh Ray and Elizabeth Glater. I cannot 29  
name all of the friends I have made here at Caltech; I hope they know how grateful 30  
I am for their friendship. 31

I would like to briefly acknowledge the programs and the people who brought me 32  
to Caltech. I would not be here without the EXtraordinary Research Opportunities 33  
Program (EXROP) from HHMI, where I met Andrew Quon and Christy Schultz. 34  
Through EXROP I met and had a chance to work for Susan Lindquist and her (then) 35  
postdoc Georgios Karras, who taught me the beauty of yeast genetics. I wish I could 36  
show Sue what I have done with the doors that she opened for me. At Cornell, I was 37  
incredibly lucky to be advised by Laurel Southard, who believed in my potential no 38  
matter what grades might say. 39

Throughout my time at Caltech, I have never been alone. My family has been 40  
a source of unconditional support. I thank my parents, Lilia and Josué, and my 41  
brother, Andrés, for always believing in me. 42

Finally, I would like to give my heartfelt thanks to Heather Curtis. 43  
  
Heather, you have been the sun, and the moon, and the stars in my life since I met 44  
you. You brought new colors into my world. Every day, I learn to think in new 45  
ways, I learn to see new things, thanks to you. I am a better person because I am 46  
with you and I am grateful that life brought us together. 47

*This thesis is for you.* 48

## ABSTRACT

49

This thesis deals with the conceptual and computational framework required to 50  
use transcriptomes as effective phenotypes for genetic analysis. I demonstrate 51  
that there are powerful theoretical reasons why Batesonian epistasis should feature 52  
prominently in transcriptional phenotypes. I also show how to compute and interpret 53  
the aggregate statistics for transcriptome-wide epistasis and transcriptome-wide 54  
dominance using whole-organism transcriptomic profiles of *C. elegans* mutants. 55  
Finally, I developed the WormBase Enrichment Suite for enrichment analysis of 56  
genomic data. 57

## PUBLISHED CONTENT AND CONTRIBUTIONS

58

Angeles-Albores, David, Raymond Y. N. Lee, et al. (2018). "Two new functions in the WormBase Enrichment Suite". In: <i>Micropublication: biology. Dataset</i> . doi: <a href="https://doi.org/10.17912/W25Q2N">https://doi.org/10.17912/W25Q2N</a> .	59 60 61
Angeles-Albores, David, Carmie Puckett Robinson, et al. (2018). "Reconstructing a metazoan genetic pathway with transcriptome-wide epistasis measurements". In: <i>Proceedings of the National Academy of Sciences</i> , p. 201712387. issn: 0027-8424. doi: <a href="https://doi.org/10.1073/pnas.1712387115">10.1073/pnas.1712387115</a> .	62 63 64 65
Angeles-Albores, David and Paul W Sternberg (2018). "Using Transcriptomes as Mutant Phenotypes Reveals Functional Regions of a Mediator Subunit in <i>Caenorhabditis elegans</i> ." In: <i>Genetics</i> , genetics.301133.2018. issn: 1943-2631. doi: <a href="https://doi.org/10.1534/genetics.118.301133">10.1534/genetics.118.301133</a> .	66 67 68 69
Angeles-Albores, David, Daniel H. W. Leighton, et al. (2017). "The <i>Caenorhabditis elegans</i> Female-Like State: Decoupling the Transcriptomic Effects of Aging and Sperm Status". In: <i>G3: Genes, Genomes, Genetics</i> 7.9.	70 71 72
Angeles-Albores, David, Raymond Y. N. Lee, et al. (2016). "Tissue enrichment analysis for <i>C. elegans</i> genomics". In: <i>BMC Bioinformatics</i> 17.1, p. 366. issn: 1471-2105. doi: <a href="https://doi.org/10.1186/s12859-016-1229-9">10.1186/s12859-016-1229-9</a> .	73 74 75

## TABLE OF CONTENTS

Acknowledgements . . . . .	iii	77
Abstract . . . . .	v	78
Published Content and Contributions . . . . .	vi	79
Table of Contents . . . . .	vii	80
List of Illustrations . . . . .	viii	81
List of Tables . . . . .	xxv	82
Preface . . . . .	1	83
Chapter I: Introduction . . . . .	3	84
Chapter II: A Statistical Mechanical Theory of Genetics using Gene Expression Phenotypes . . . . .	24	86
Chapter III: Reconstructing a metazoan genetic pathway with transcriptome-wide epistasis measurements . . . . .	43	88
Chapter IV: The <i>Caenorhabditis elegans</i> Female-Like State: Decoupling the Transcriptomic Effects of Aging and Sperm Status . . . . .	79	90
Chapter V: Using transcriptomes as mutant phenotypes reveals functional regions of a Mediator subunit in <i>C. elegans</i> . . . . .	105	92
Chapter VI: Tissue enrichment analysis for <i>C. elegans</i> genomics . . . . .	132	93
Chapter VII: Two new functions in the WormBase Enrichment Suite . . . . .	156	94
7.1 Description . . . . .	156	95
Conclusion . . . . .	159	96

## LIST OF ILLUSTRATIONS

97

<i>Number</i>		<i>Page</i>
11	Biologists work with two distinct types of epistasis. <b>A.</b> Batesonian, or classical, epistasis refers to those cases where the qualitative phenotype associated with one null mutation is masked completely by the presence of a second mutation at a distinct locus. <b>B.</b> Generalized epistasis is used for quantitative phenotypes and measures the systematic deviation in the phenotype of a double mutant relative to a statistical null model. Unlike Batesonian epistasis, generalized epistasis cannot be used to infer genetic pathways, since the choice of null model is arbitrary. The effects associated with allele $x$ are labelled $\beta_x$ , and the generalized epistasis is given the symbol $\Delta$ . . . .	98 99 100 101 102 103 104 105 106 107 6 108

12	Analysis methodology to infer genetic interactions using transcriptome data. <b>A.</b> After fitting all transcripts to a general linear model to calculate the individual and the epistatic components of null mutations in two distinct genes, the resulting parameters can be clustered and visualized in a heatmap. Each observed cluster can be grouped into one of 27 epistatic classes. All clusters are considered biologically relevant regardless of the number of transcripts they contain.	109 110 111 112 113 114 115
	A simple conclusion cannot be reached from these heatmaps. This approach was used in <b>Dixit2016; Adamson2016</b> <b>B.</b> Starting from the same statistical model, only transcripts that have all parameters different from zero are considered informative. These transcripts are plotted on a scatterplot, where the x-axis reflects the expected value of the double mutant under an additive or log-additive hypothesis, and the systematic deviation from additivity (generalized epistasis) is plotted on the y-axis. The resulting points form a ray on the plot. The slope of this ray is an aggregate statistic that can be interpreted in terms of a genetic pathway if the two genes exhibit Batesonian epistasis. This approach was used in <b>Angeles-Albores2017; Angeles-Albores2018</b> .	116 117 118 119 120 121 122 123 124 125 126
	12	126

13 Genes that are differentially expressed in genotypes containing mutant ( <i>a, b</i> ) alleles relative to a wild type homozygote can be categorized into phenotypic classes. Each phenotypic class can in turn be associated with a dominance behavior. The Venn diagram represents differentially expressed transcripts in each genotype relative to the wild-type control. Each of the possible 7 intersections is labelled with its dominance interpretation if the intersection is real. In this context, semi-recessiveness means that one allele is partially or completely dominant to the other along a continuous spectrum between 0 and 1. The dominance sign between an allele and the heterozygote genotype indicates heterosis or over-dominance. . . . .	127
	128
	129
	130
	131
	132
	133
	134
	135
	136
	137
16	137





31	Genetic and biochemical representation of the hypoxia pathway in <i>C. elegans</i> . Red arrows are arrows that lead to inhibition of HIF-1, and blue arrows are arrows that increase HIF-1 activity or are the result of HIF-1 activity. EGL-9 is known to exert VHL-1-dependent and independent repression on HIF-1 as shown in the genetic diagram. The VHL-1-independent repression of HIF-1 by EGL-9 is denoted by a dashed line and is not dependent on the hydroxylating activity of EGL-9. RHY-1 inhibits CYSL-1, which in turn inhibits EGL-9, but this interaction was abbreviated in the genetic diagram for clarity.	182 183 184 185 186 187 188 189 46 190
32	Analysis workflow. After sequencing, reads are quantified using Kallisto. Bars show estimated counts for each isoform. Differential expression is calculated using Sleuth, which outputs one $\beta$ coefficient per isoform per genotype. $\beta$ coefficients are analogous to the natural logarithm of the fold-change relative to a wild type control. Downstream analyses are performed with $\beta$ coefficients that are statistically significantly different from 0. $q$ -values less than 0.1 are considered statistically different from 0. . . . .	191 192 193 194 195 196 197 48 198
33	Principal component analysis of various <i>C. elegans</i> mutants. Genotypes that have an constitutive hypoxia response (i.e. <i>egl-9(lf)</i> ) cluster far from genotypes that do not have a hypoxic response (i.e. <i>hif-1(lf)</i> ) along the first principal component. The second principal component separates genotypes that do not participate hypoxic response pathway.	199 200 201 202 51 203

34	Interacting genes have correlated transcriptional signatures. The rank order of transcripts contained in the shared transcriptional phenotype is plotted for each pairwise combination of genotypes. Correlations between in-pathway genotypes are strong whereas comparisons with a <i>fog-2(lf)</i> genotype are dominated by noise. Comparisons between some genotypes show populations of transcripts that are anticorre- lated, possibly as a result of feedback loops. Plots are color-coded by row. Comparisons with genotypes with a constitutive hypoxia re- sponse are in blue; comparisons with genotypes negative for <i>hif-1(lf)</i> are black; and comparisons involving <i>fog-2(lf)</i> are red. X- and y-axes show the rank of each transcript within each genotype. . . . .	53	214
----	--	----	-----

- 35 (A) Schematic diagram of an epistasis plot. The X-axis on an epistasis plot is the expected coefficient for a double mutant under an log-additive model (null model). The Y-axis plots deviations from this model. Double mutants that deviate in a systematic manner from the null model exhibit transcriptome-wide epistasis ( $s$ ). To measure  $s$ , we find the line of best fit and determine its slope. Genes that act log-additively on a phenotype (**Ph**) will have  $s = 0$  (null hypothesis, orange line); whereas genes that act along an unbranched pathway will have  $s = -1/2$  (blue line). Strong repression is reflected by  $s = -1$  (red line), whereas  $s > 0$  correspond to synthetic interactions (purple line). (B) Epistasis plot showing that the *egl-9(lf)*; *vhl-1(lf)* transcriptome deviates significantly from a null additive. Points are colored qualitatively according to density (purple—low, yellow—high) and size is inversely proportional to the standard error (S.E.) of the y-axis. The green line is the line of best fit from an orthogonal distance regression. (C) Comparison of simulated epistatic coefficients against the observed coefficient. Green curve shows the bootstrapped observed transcriptome-wide epistasis coefficient for *egl-9* and *vhl-1*. Dashed green line shows the mean value of the data. Simulations use only the single mutant data to idealize what expression of the double mutant should look like.  $a > b$  means that the phenotype of  $a$  is observed in a double mutant  $a^-b^-$ .

- 36 Transcriptomes can be used to order genes in a pathway under certain assumptions. Arrows in the diagrams above are intended to show the direction of flow, and do not indicate valence. **A.** A linear pathway in which *rhy-1* is the only gene controlling *egl-9*, which in turn controls *hif-1* does not contain information to infer the order between genes. 237  
**B.** If *rhy-1* and *egl-9* have transcriptomic effects that are separable from *hif-1*, then the *rhy-1* transcriptome should contain contributions 238  
from *egl-9*, *hif-1* and *egl-9*- and *hif-1*-independent pathways. This 239  
pathway contains enough information to infer order. **C.** If a path- 240  
way is branched both upstream and downstream, transcriptomes will 241  
show even faster decorrelation. Nodes that are separated by many 242  
edges may begin to behave almost independently of each other with 243  
marginal transcriptomic overlap or correlation. **D.** The hypoxia path- 244  
way can be ordered. We hypothesize the rapid decay in correlation 245  
is due to a mixture of upstream and downstream branching that hap- 246  
pens along this pathway. Bars show the standard error of the weighted 247  
coefficient from the Monte Carlo Markov Chain computations. . . . 59 248  
253
- 37 **A.** 56 genes in *C. elegans* exhibit non-classical epistasis in the hy- 254  
poxia pathway, characterized by opposite effects on gene expression, 255  
relative to the wild type, of the *vhl-1(lf)* compared to *egl-9(lf)* (or 256  
*rhy-1(lf)*) mutants. Shown are a random selection of 15 out of 56 257  
genes for illustrative purposes. **B.** Genes that behave non-canonically 258  
have a consistent pattern. *vhl-1(lf)* mutants have an opposite effect 259  
to *egl-9(lf)*, but *egl-9* remains epistatic to *vhl-1* and loss-of-function 260  
mutations in *hif-1* suppress the *egl-9(lf)* phenotype. Asterisks show 261  
 $\beta$  values significantly different from 0 relative to wild type ( $q < 10^{-1}$ ). 62 262

38	A hypothetical model showing a mechanism where HIF-1-hydroxyl antagonizes HIF-1 in normoxia. <b>A.</b> Diagram showing that RHY-1 activates EGL-9. EGL-9 hydroxylates HIF-1 in an oxygen-dependent manner. HIF-1 is rapidly hydroxylated and the product, HIF-1-OH is rapidly degraded in a VHL-1-dependent fashion. EGL-9 can also inhibit HIF-1 in an oxygen-independent fashion. In our model, HIF-1 and HIF-1-OH have opposing effects on transcription. The width of the arrows represents rates in normoxic conditions. <b>B.</b> Table showing the effects of loss-of-function mutations on HIF-1 and HIF-1-OH activity, showing how this can potentially explain the <i>ftn-1</i> expression levels in each case. S.S = Steady-state. . . . .	65	273
41	Experimental design to identify genes associated with sperm loss and with aging. Studying the wild-type worm alone would measure time- and sperm-related changes at the same time, without allowing us to separate these changes. Studying the wild-type worm and a <i>fog-2(lf)</i> mutant would enable us to measure sperm-related changes but not time-related changes. By mixing both designs, we can measure and separate both modules. . . . .	82	280

42	A. Differentially expressed isoforms in the aging category. We identified a common aging expression signature between N2 and <i>fog-2(lf)</i> animals, consisting of 6,193 differentially expressed isoforms totaling 5,592 genes. The volcano plot is randomly down-sampled 30% for ease of viewing. Each point represents an individual isoform. $\beta_{\text{Aging}}$ is the regression coefficient. Larger magnitudes of $\beta$ indicate a larger log-fold change. The y-axis shows the negative logarithm of the q-values for each point. Green points are differentially expressed isoforms; orange points are differentially expressed isoforms of predicted transcription factor genes ( <b>Reece-Hoyes2005</b> ). An interactive version of this graph can be found on our website.	281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297
	B. Enriched tissues in aging-associated genes. Tissue Enrichment Analysis ( <b>Angeles-Albores2016</b> ) showed that genes associated with muscle tissues and the nervous system are enriched in aging-related genes. Only statistically significantly enriched tissues are shown. Enrichment Fold Change is defined as <i>Observed/Expected</i> . hmc stands for head mesodermal cell. . . . .	86

43	Explanation of linear regressions with and without interactions. <b>A.</b> A linear regression with two variables, age and genotype. The expression level of a hypothetical gene increases by the same amount as worms age regardless of genotype. However, <i>fog-2(lf)</i> has higher expression of this gene than the wild-type at all stages (blue arrow). <b>B.</b> A linear regression with two variables and an interaction term. In this example, the expression level of this hypothetical gene is different between wild-type worms and <i>fog-2(lf)</i> (blue arrow). Although the expression level of this gene increases with age, the slope is different between wild-type and <i>fog-2(lf)</i> . The difference in the slope can be accounted for through an interaction coefficient (red arrow). . . . .	298
		299
		300
		301
		302
		303
		304
		305
		306
		307
		308
87		
44	<i>fog-2(lf)</i> partially phenocopies early aging in <i>C. elegans</i> . The $\beta$ in each axes is the regression coefficient from the GLM, and can be loosely interpreted as an estimator of the log-fold change. Loss of <i>fog-2</i> is associated with a transcriptomic phenotype involving 1,881 genes. 1,040/1,881 of these genes are also altered in wild-type worms as they progress from young adulthood to old adulthood, and 905 change in the same direction. However, progression from young to old adulthood in a <i>fog-2(lf)</i> background results in no change in the expression level of these genes. <b>A.</b> We identified genes that change similarly during feminization and aging. The correlation between feminization and aging is almost 1:1. <b>B.</b> Epistasis plot of aging versus feminization. Epistasis plots indicate whether two genes (or perturbations) act on the same pathway. When two effects act on the same pathway, this is reflected by a slope of -0.5. The measured slope was $-0.51 \pm 0.01$ . . . . .	309
		310
		311
		312
		313
		314
		315
		316
		317
		318
		319
		320
		321
		322
93		
		323

45	Phenotype and GO enrichment of genes involved in the female-like state. <b>A.</b> Phenotype Enrichment Analysis. <b>B.</b> Gene Ontology Enrichment Analysis. Most of the terms enriched in PEA reflect the abundance of ribosomal subunits present in this gene set. . . . .	95	324
46	<b>A.</b> A substrate-dependent model showing how <i>fog-2</i> promotes sperm generation, whereas aging promotes sperm depletion, leading to entry to the female-like state. Such a model can explain why <i>fog-2</i> and aging appear epistatic to each other. <b>B.</b> The complete <i>C. elegans</i> life cycle. Recognized stages of <i>C. elegans</i> are marked by black arrows. States are marked by red arrows to emphasize that at the end of a state, the worm returns to the developmental timepoint it was at before entering the state. The L2d state is an exception. It is the only stage that does not return to the same developmental timepoint; rather, the L2d state is a permissive state that allows entry into either dauer or the L3 stage. We have presented evidence of a female-like state in <i>C. elegans</i> . At this point, it is unclear whether the difference between hermaphrodites and females is reversible by males. Therefore, it remains unclear whether it is a stage or a true state. . . . .	98	327
51	Protein sequence schematic for DPY-22. The positions of the non-sense mutations used are shown. . . . .	107	331
52	Principal component analysis of the analyzed genotypes. The analysis was performed using only those transcripts that were differentially expressed in at least one genotype. The plot shows that the <i>trans</i> -heterozygotes phenocopy the <i>dpn-22(bx93)</i> homozygotes along the first two principal dimensions. . . . .	114	348

53	Flowchart for an analysis of arbitrary allelic series. A set of alleles is selected, and the corresponding genotypes are sequenced. Independent phenotypic classes are identified, and classes that are the result of noise are discarded via a false hit analysis. For each phenotypic class, the alleles are ordered in a dominance/complementation hierarchy, which can then be used to infer functional regions (FR) within the genes in question. . . . .	116	355
54	Shared Transcriptomic Phenotypes amongst the <i>dpy-22</i> genotypes are regulated in the same direction. For each pairwise comparison, we found those transcripts that were commonly differentially expressed in both genotypes relative to the wild-type control and plotted the $\beta$ coefficients for each. We performed a linear regression on each plot to find the line of best fit (broken blue line). Only the comparison between <i>dpy-22(sy622)</i> and <i>dpy-22(bx93)</i> homozygotes was used to establish that the magnitude of the <i>dpy-22(sy622)</i> allele is greater than the magnitude of the <i>dpy-22(bx93)</i> allele. The other comparisons are shown for completeness. . . . .	120	365
55	<i>dpy-22</i> phenotypic classes are statistically significantly enriched for signatures of <i>let-60</i> (ras) and <i>bar-1</i> (wnt) signaling. We tested whether the overlap between the differentially expressed genes in <i>bar-1(ga80)</i> , <i>let-60(n1046gf)</i> or <i>let-60(n2021)</i> and the <i>dpy-22</i> phenotypic classes was statistically significant using a hypergeometric enrichment test. Since the hypergeometric enrichment test is very sensitive to deviations from random, and since we suspect that there may be a broad genotoxic response to all mutants, we used a statistical significance threshold of $p < 10^{-10}$ (dashed black line). . . . .	123	374

56	The functional regions associated with each phenotypic class can be mapped intragenically. The number of genes associated with each class is shown. The <i>dpy-22(bx93)</i> -associated class may be controlled by two functional regions. FR1 is a dosage-sensitive unit. FR2 and FR3 could be redundant if FR4 is a modifier of FR2 functionality at <i>dpy-22(bx93)</i> -associated loci. Note that the <i>dpy-22(bx93)</i> -associated phenotypic class is actually three classes merged together. Two of these classes are DE in <i>dpy-22(bx93)</i> homozygotes and one other genotype. Our analyses suggested that these two classes are likely the result of false negative hits and genes in these classes should be differentially expressed in all three genotypes, so we merged these three classes together (see Methods). . . . .	126	386
61	Schematic representation of trimming filters for an acyclical ontology.	387	
	<b>a.</b> The parent node (green) contains at least as many annotations as the union of the two sisters. These two sisters share annotations extensively, as expressed by the overlap in the Venn diagram, so they qualify for removal. <b>b.</b> Nodes with less than a threshold number of genes are trimmed (red) and discarded from the dictionary. Here, the example threshold is 25 genes. Nodes $\epsilon, \zeta, \eta$ , shown in red are removed. <b>c.</b> Parent nodes are removed recursively, starting from the root, if all their daughter nodes have more than the threshold number of annotations. Nodes in grey ( $\epsilon, \zeta, \eta$ ) were removed in the previous step. Nodes $\alpha, \beta$ shown in red are trimmed because each one has a complete daughter set. Only nodes $\gamma$ and $\delta$ will be used to generate the static dictionary. . . . .	136	399

62	Screenshot of results from the web GUI. After inputting a gene-list, the user is provided with the results. An HTML table is output with hyperlinks to the ontology terms. A publication-ready graph is provided below, which can be saved by dragging to the desktop. The graph is colored for better visualization; color is not intended to convey information. The graph and the table show anatomy terms in human-readable format, followed by their unique WBbt ID. Finally, lists of the genes used and discarded for the analysis are also presented.	400 401 402 403 404 405 406 407 140
63	TEA Workflow. The complete ontology is annotated continuously by WormBase curators. After each update, the ontology is processed to remove uninformative terms, and the remaining terms are used for statistical testing. Users can select a gene list and input it into our tool using our WormBase portal. The gene list is tested for enrichment using the trimmed ontology, and results are output in tabular and graphic formats for analysis.	408 409 410 411 412 413 141 414
64	Kernel density estimates (KDE) for 30 gold standard datasets. We ran TEA on 30 datasets we believed to be enriched in particular tissues and pooled all the results to observe the distribution of q-values. The mode of the distribution for dictionaries with annotation cut-offs of 100 and 50 genes are very similar; however, when the cut-off is lowered to 25 genes, the mode of the distribution shifts to the left, potentially signalling a decrease in measurement power.	415 416 417 418 419 420 143 421

65	Independently derived gene sets show similar results when tested with the same dictionary. <b>Set 1.</b> GABAergic gene set from Watson ( <b>Watson2008a</b> ). <b>Set 2.</b> GABAergic gene set from Spencer ( <b>Spencer2011</b> ) <sup>424</sup> Arrowheads highlight identical terms between both analyses. All terms refer to neurons or neuronal tissues and are GABA-associated.	422 423 424 425 426 427
	Dictionary with cutoff: 33; threshold: 0.95; method: ‘any’.	145
66	<i>D. coniospora</i> Gene Enrichment Analysis and Tissue Enrichment Analysis results. We compared and contrasted the results from a gene enrichment analysis program, pantherDB, with TEA by analyzing genes that were significantly down-regulated when <i>C. elegans</i> was exposed to <i>D. coniospora</i> in a previously published dataset by Engelmann <i>et al</i> ( <b>Engelmann2011</b> ) with both tools. <b>a.</b> pantherDB screenshot of results, sorted by p-value. Only top hits shown. <b>b.</b> TEA results, sorted by q-value (lowest on top) and fold-change. Both pantherDB and TEA identify terms associated with neurons (red square). The two analyses provide complementary, not redundant, information.	428 429 430 431 432 433 434 435 436 437 147
		438

## LIST OF TABLES

439

<i>Number</i>	<i>Page</i>
31 Number of differentially expressed genes in each mutant strain with respect to the wild type (N2). . . . .	50 442
51 The number of differentially expressed genes relative to the wild-type control for each genotype with a significance threshold of 0.1. . . . .	114 444
52 Dominance analysis for the <i>dpy-22/MDT12</i> allelic series. Dominance values closer to 1 indicate <i>dpy-22(bx93)</i> is dominant over <i>dpy-22(sy622)</i> , whereas 0 indicates <i>dpy-22(sy622)</i> is dominant over <i>dpy-22(bx93)</i> . . . . .	121 448
61 Parameter specifications and number of tissues for all dictionaries. The ‘Method’ column refers to the trimming criterion for the similarity metric. We used two such criteria, ‘any’ and ‘avg’.‘any’: For a given sister set, if any sister had a similarity exceeding the corresponding threshold, all sisters were removed from the final dictionary. ‘avg’: For a given sister set, if the average similarity across all the sisters in the set was greater than the corresponding threshold, all sisters were removed from the final dictionary. . . . .	138 456

62 Comparison of results for a GABAergic neuronal-enriched gene set from Watson ( <b>Watson2008a</b> ) showing that results are similar regardless of annotation cutoff. We ran the same gene list on a dictionary with a minimum annotation cutoff of 50, similarity threshold of 0.95 and similarity method ‘any’ versus another with a minimum annotation cutoff of 33, similarity threshold of 0.95 and similarity method ‘any’. In the table, columns are labeled with their significance value (Q-value) or enrichment fold change followed by a hyphen and a number which indicates which the cutoff for the dictionary that was used for testing. Not all tissues are present in either dictionary. Hyphens denote not-applicable values, which occurs when a particular tissue is not present in both dictionaries.	145	468
		457
		458
		459
		460
		461
		462
		463
		464
		465
		466
		467

## PREFACE

I have tried to organize this thesis in a way that makes sense. Briefly, the thesis can be 470  
viewed in three segments: Epistasis, complementation and software development. 471  
In doing so, I have broken the chronological order of my work, but I do not see this 472  
as a problem. Science is rarely linear, and it may often turn out that the last concepts 473  
to be found are actually those concepts that allow us to make sense of everything 474  
else. This has certainly been the case with my work. 475

In Chapter 1, the reader will find a brief overview of the problem facing transcript- 476  
tome genetics. This chapter encompasses a review of the relevant literature, but 477  
beyond that, I have tried to make arguments I think are important. First, transcript- 478  
tome genetics has obviated the chasm between statistical epistasis and classical, or 479  
Batesonian, epistasis. The confusion between the two (related) terms has been one 480  
of the great misfortunes in the field of genetics, since it has hampered a significant 481  
amount of work. I am glad to say that in this thesis I have achieved the unification 482  
of both concepts, such that no confusion should happen. Second, although we now 483  
know how to search signs of epistasis and dominance in transcriptomes, the issue of 484  
counting phenotypic classes or modules is becoming increasingly ominous. Unless 485  
and until we can confidently identify and purge spurious modules, we will not be 486  
able to use these phenotypes to their full extent. 487

In Chapter 2, I have written a theoretical argument that is the basis for the rest of 488  
the chapters dealing with epistasis. In this chapter, I prove that epistasis emerges 489  
from statistical mechanics, such that even genes that have enormously complex 490  
transcriptional mechanisms can in some cases exhibit Batesonian epistasis. This 491  
chapter establishes genetics as a variational method with which to probe an unknown 492  
partition function, enabling us to make statements about what values the partition 493

function is or is not allowed to take.

494

In Chapter 3, I develop the concept of transcriptome-wide epistasis and use it to 495  
reconstruct the well-studied hypoxia pathway. In Chapter 4, I use the concept 496  
of transcriptome-wide epistasis to identify a novel stage in the life cycle of the 497  
roundworm *C. elegans*. 498

Chapter 5 deals with the issue of complementation, and its study through expression 499  
profiles. In my opinion, this is the most complicated chapter in this thesis. I 500  
struggled with every aspect of this project, but the result is, to my mind, pleasing. 501

Chapter 6 and 7 deal with the creation of the WormBase Enrichment Suite. 502

Throughout this thesis, I have tried to be pedagogical. If we don't make efforts 503  
to explain the computational methods we are developing, biology will pass from a 504  
scientific discipline to an astrological pseudo-science, and we will fail to see the 505  
true beauty in the stars above and instead imbue them with our human desires and 506  
flaws, asking them to help us reach fame instead of helping us to solve the mysteries 507  
that abound in our universe. 508

## INTRODUCTION

510

**Abstract**

511

**Transcriptomes are microscopic phenotypes of enormous complexity. In spite** 512  
**of this complexity, it is becoming apparent that transcriptomes follow the same** 513  
**genetic rules as all other mesoscopic and macroscopic phenotypes. Due to** 514  
**their complexity, the genetic rules that bind transcriptomes appear more com-** 515  
**plicated. There is significant interest in developing statistical and biological** 516  
**methods that can deconvolute transcriptomes to extract the maximum amount** 517  
**of information encoded within them. Here, we review the basic concepts that un-** 518  
**derlie transcriptome genetics, identify confusions in the field and point towards** 519  
**the emerging challenges and opportunities associated with these intriguing new** 520  
**phenotypes.**

521

**Introduction**

522

The recent explosion in genomic technologies has provided us with unparalleled 523  
insight into the inner workings of cells. The cost of sequencing continues to 524  
drop, and new technologies are continuously increasing the number of samples that 525  
can be sequenced. In turn, these massive datasets have promoted the appearance of 526  
increasingly complex algorithms to make sense of them. A common tenet in these 527  
methods has been to reduce the dimensionality of these datasets (dimensionality 528  
refers to the number of measurements per sample) to look for trends in the data. 529  
Though sometimes these methods are rooted in biological principles, more often they 530  
come from algebraic methods that have no immediate connection to the underlying 531  
biology. This means that although these methods may be quite powerful, the results 532

may be hard to interpret in biological terms. Moreover, these methods may not utilize 533  
the rich structure inherent to biological systems that could place strong constraints 534  
on the problem under study to reduce the space of reasonable solutions. 535

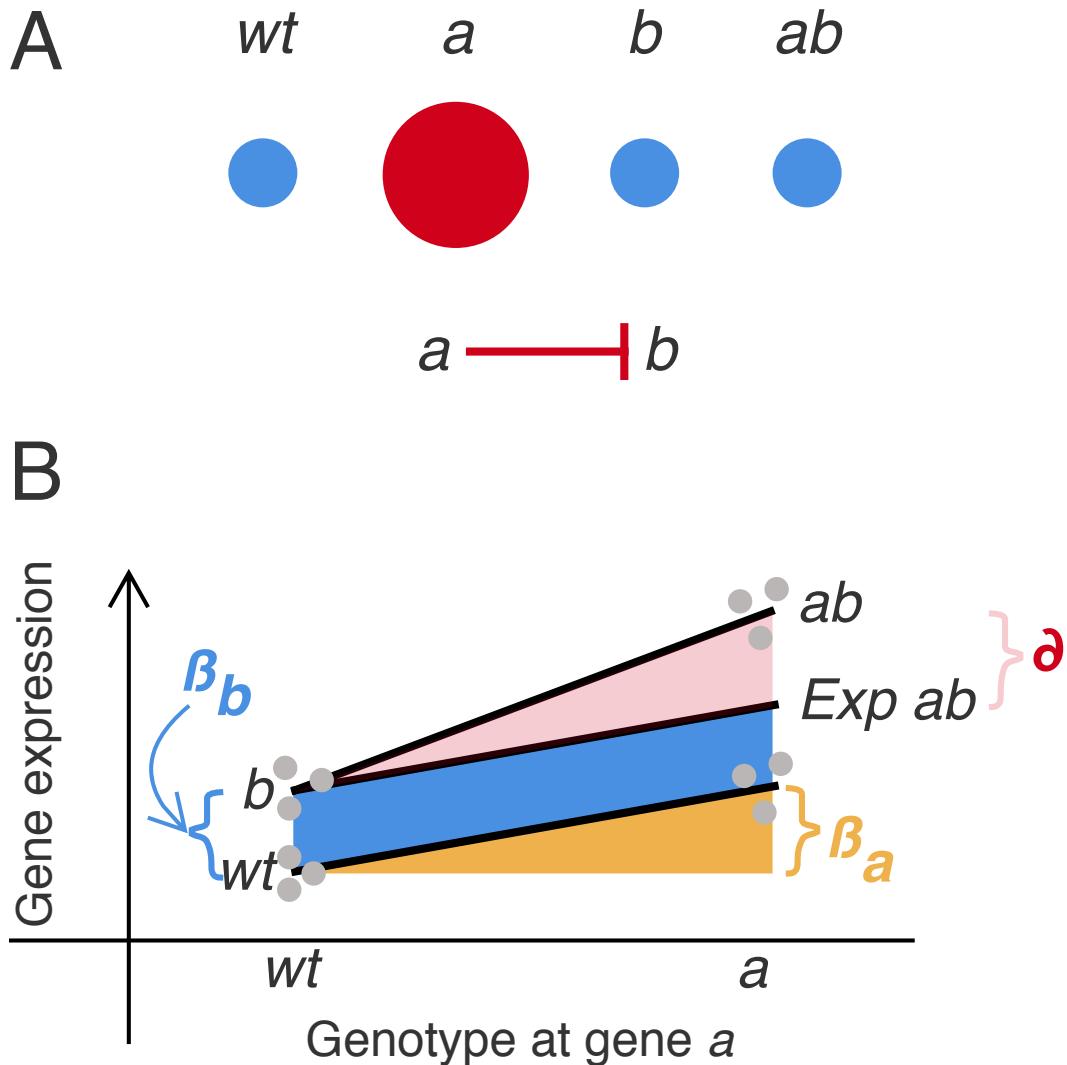
Biological systems can be daunting in their complexity. In general, there is no way 536  
to solve these complex systems from first principles, since the identity and activity 537  
of each component is in general not known. For this reason, biologists developed 538  
a set of methods, collectively referred to as genetic analyses, that do not make 539  
many assumptions regarding the underlying molecular details. Genetics is limited 540  
to making a limited set of true statements regarding certain kinds of molecular inter- 541  
actions. Due to its limited scope, genetics is robust to biological variation. A major 542  
goal over the last century has been to reconstruct the set of all genetic networks 543  
that result in a specific phenotype in a specific condition (the genotype-phenotype 544  
map). Though spectacular progress has been made in some cases (Costanzo et al., 545  
[2016](#)), we are still far from understanding genetic networks. Previously, generating 546  
sufficient genotypes in model organisms to analyze any network in detail was a 547  
bottleneck to perform thorough genetic reconstructions. However, with the advent 548  
of genome engineering, generating specific mutants is rapidly becoming easier. On 549  
the other hand, sensitive and fast phenotyping methods have lagged behind. A 550  
possible solution to this problem is bulk expression profiling, but the complexity of 551  
expression profiles had proved a daunting challenge for genetic analysis. Further- 552  
more, expression profiles have brought to the forefront a major source of confusion 553  
in genetics: The definition of genetic interactions. 554

Biologists identify genetic interactions between genes using a specific method called 555  
epistasis analysis. The term ‘epistasis’ was used for the first time over one hundred 556  
years ago by William Bateson (Bateson and Mendel, [2009](#)) to refer to the observation 557  
that the distribution of offspring phenotypes from a double heterozygote cross did 558  
not match the expected distribution prescribed by Mendelian segregation of two loci. 559

Under Mendelian laws, if two loci are associated with different phenotypes, crossing double heterozygotes of these two loci should generate animals with four phenotypic classes, with each class occurring in a 9:3:3:1 ratio. Bateson realized through segregation analyses that in certain cases, the phenotypic class associated with the double mutant was missing, and instead there was an excess of one phenotypic class typically associated with homozygotes of one mutant allele, an effect similar to Mendel's observations of allelic dominance. He coined the term epistasis to refer to the effect by which an allele at one locus, when present in two copies, can completely mask the phenotypic effect of another allele at a separate locus.

Since he coined the term, Batesonian or classical epistasis has become a popular tool amongst geneticists with which to identify genetic interactions. An important caveat is that in order to perform an epistasis analysis, geneticists must restrict themselves to alleles that are completely devoid of function. When this is the case, the phenotypic transformation of the double mutant is used to construct a genetic pathway (Avery and Wasserman, 1992; Huang and Sternberg, 2006) (see Fig. 11). Classical epistasis has become a cornerstone of biology.

Classical epistasis means that the phenotype of the double mutant is exactly the same as the phenotype of one of the single mutants. However, the problem can also be recast in quantitative terms. Statistical geneticists defined generalized epistasis as a systematic deviation between the observed values and a null model (usually additive or log-additive) that can be corrected by adding a second order interaction term (Fisher, 1919). In the terms of generalized genetics, epistasis in the heterozygote crosses is measured in the systematic excess of one phenotypic class and the systematic depletion of a second class. Notably, generalized epistasis is not constrained in the values it can take, and it is not constrained to measurements of population properties or properties of single individuals.



**Figure 11** Biologists work with two distinct types of epistasis. **A.** Batesonian, or classical, epistasis refers to those cases where the qualitative phenotype associated with one null mutation is masked completely by the presence of a second mutation at a distinct locus. **B.** Generalized epistasis is used for quantitative phenotypes and measures the systematic deviation in the phenotype of a double mutant relative to a statistical null model. Unlike Batesonian epistasis, generalized epistasis cannot be used to infer genetic pathways, since the choice of null model is arbitrary. The effects associated with allele  $x$  are labelled  $\beta_x$ , and the generalized epistasis is given the symbol  $\Delta$ .

As a result of its definition, the magnitude of generalized epistasis is completely 586 dependent on the null model selected by the researcher. Unlike physical models that 587 can be derived from first principles, statistical models of genetic interactions are 588 heuristic models that may or may not represent the molecular interactions underlying 589 the system accurately. In this sense, second order ‘interaction’ terms are *ad hoc* 590 corrections, technically useful for machine-learning, but not instructive in terms 591 of understanding the genetic mechanisms at play. The conceptual proof for this is 592 simple: Imagine two different statistical models that describe how two genes interact 593 along a phenotype. Both models perform equally well. One of the models has a 594 statistically significant interaction (generalized epistasis) term whereas the other 595 does not. It is not possible to select one model over the other based on statistical 596 properties. In fact, based on model simplicity, we may even prefer the model with 597 fewer parameters, which could rule out the model that includes an epistasis term. 598

Like classical epistasis, generalized epistasis has become a useful concept in many 599 areas of biology. Unlike classical epistasis, generalized epistasis measurements have 600 not been restricted to those generated by null alleles; instead, generalized epistasis, 601 particularly in human genetics, is measured between any two molecular variants 602 at different loci measured under a specific null model. As a result of the subtle 603 differences between classical and generalized epistasis, there has been considerable 604 concern about the apparent disagreement between these two concepts (Phillips, 605 2008; Cordell, 2002; Lehner, 2011). In this review, we will show how generalized 606 and classical epistasis can be successfully unified. Moreover, this unification has 607 important ramifications for our ability to detect genetic interactions between two 608 mutants using genome-wide studies. 609

**Motivation: A brief introduction to RNA-sequencing**

610

RNA-sequencing (Mortazavi et al., 2008) is a powerful method that can measure all 611 the gene expression levels in an organism simultaneously. These measurements can 612 be made in bulk, from homogenized tissues or even from whole-organisms. Recent 613 technological breakthroughs have made measuring expression levels from single 614 whole organisms (Serra et al., 2018; Chan, Rando, and Conine, 2018; Lott et al., 615 2011) or even single cells possible (Tang et al., 2009). As a result of its technical 616 advantages, RNA-seq has largely replaced microarrays as the method of choice to 617 monitor gene expression. 618

Since the advent of genome-wide measurement methods, the idea of a cell- or 619 organismal-state, defined by its gene expression levels, has drawn significant atten- 620 tion. Such states make sense in light of gene regulatory network theory, which posits 621 that the expression of many genes is coordinated by regulatory factors that, when 622 expressed, drive development forward (Britten and Davidson, 1969). A common 623 experimental design used to identify the genes that are controlled by a specific reg- 624 ulatory module is to measure a baseline (typically wild type) sample and a contrast 625 sample where the regulatory module has been perturbed (often through mutation). 626 These experimental designs identify differentially expressed genes between the wild 627 type and the mutant samples. These batteries can then be analyzed through ontolog- 628 ical enrichment analyses that attempt to integrate information from all the enriched 629 transcripts and identify the biological processes or signaling pathways contained 630 within this list (see for example Mi et al. (2009) and Angeles-Albores, N. Lee, et al. 631 (2016)). In spite of the enormous amount of quantitative information that RNA-seq 632 can provide about the genes that respond to a downstream perturbation, these single 633 factor experimental designs are generally used to select a small number of novel 634 downstream genes that can be studied to extend a pathway of interest. The problem 635 of how to analyze the rich datasets generated by RNA-seq has proved difficult, and 636

no one answer will be suitable for all problems. Analyses of these datasets rely on 637  
a combination of biological intuition, enrichment analyses or comparisons to other 638  
existing datasets. 639

If we are willing to sacrifice the requirement for interpretability, these datasets 640  
are still useful. Their practicality derives significantly from enormous advances 641  
in library preparation methods (Picelli et al., 2014) and improved quantification 642  
algorithms (Patro, Mount, and Kingsford, 2014; Patro, Duggal, et al., 2016; Bray 643  
et al., 2016) that have made RNA-seq an eminently replicable protocol that is 644  
fast to execute. As a result, transcriptomes can readily be used to compare the 645  
extent to which two perturbations are similar through clustering methods. Thus, 646  
transcriptomes could be thought of as extremely long barcodes that are associated 647  
with specific, potentially hidden, variables. If two barcodes are similar, then it is 648  
plausible to hypothesize that the perturbations applied to generate each barcode were 649  
also similar, even though we may not understand what these barcodes mean or how 650  
they were generated. However, it is not sufficient to develop algorithms that show 651  
two perturbations are similar on average. To use transcriptomes for genetic analysis, 652  
we need methods that quantitatively reveal what aspects of two transcriptomes are 653  
similar, by how much and that allow us to understand why they are similar. 654

## **Genetic interactions detection through sequencing** 655

### **A brief overview of the problem** 656

Expression profiles are vectors where each entry corresponds to the expression level 657  
of a single transcript. Conceptually, each entry could be treated as an independent 658  
continuous phenotypes. Since continuous phenotypes can be used to detect statistical 659  
epistasis, we could fit a statistical model to explain the expression level of this 660  
transcript in each genotype measured (wild type, single and double mutants). This 661  
statistical model will fit two parameters,  $\beta_a$  and  $\beta_b$ , that explain the *individual* effects 662

of each null mutation, and a third parameter,  $\Delta$ , quantifies the extent to which these 663 individual effects do not add when both null mutations are present at once (see 664 Fig. 11). Each parameter is associated with a  $p$ -value. These models are generated 665 for every measured transcript. The generated  $p$ -values should then be adjusted 666 for multiple comparisons (these adjusted values are referred to as  $q$ -values), and 667 parameters with  $q$ -values below a pre-specified threshold (often 0.1) are considered 668 statistically different from zero. 669

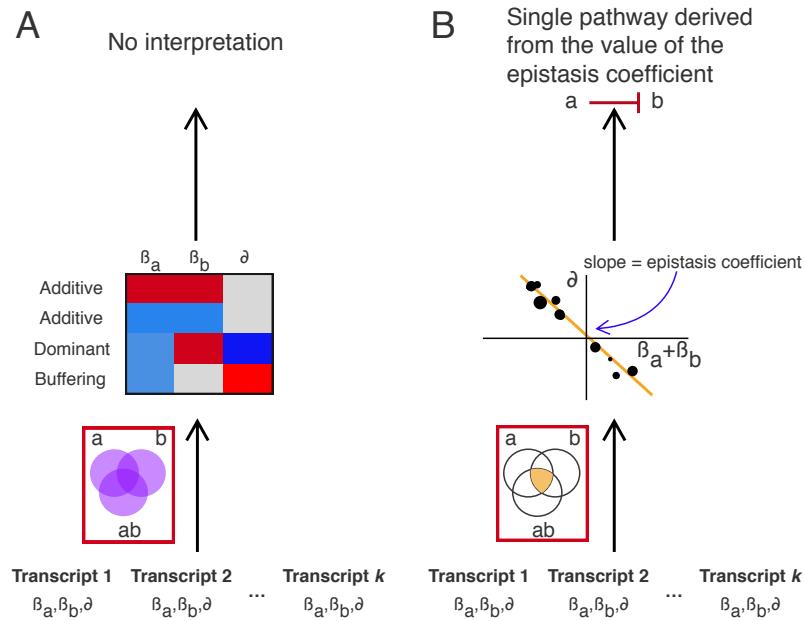
As a result, each transcript is associated with six values: the three model parameters 670 and their corresponding  $q$  values. At this point, the complexity of the problem is 671 obvious: Even if each parameter only acquires one of three values (-, 0, +), there are 672 27 possible parameter combinations (epistatic classes). Of these 27 classes, only 673 four classes could give rise to the classical epistasis regime (genetic suppression), so 674 only those 4 classes can give rise to a genetic diagram. One approach to visualize and 675 attempt to understand this complex space of epistatic combinations is to use heatmaps 676 to look for patterns and guide interpretation. This approach was used in single-celled 677 organisms (Capaldi et al., 2008; Van Driessche et al., 2005; Sameith et al., 2015; 678 Van De Peppel et al., 2005), and more recently has been used to perform high- 679 throughput analyses of genetic interactions in mammalian cells (Dixit et al., 2016). 680 Regression models with interactions have also been successfully implemented using 681 whole-organism transcriptomic measurements (Angeles-Albores, Leighton, et al., 682 2017). 683

The large number of parameter combinations is not the only (or major) drawback 684 to fitting models with interactions for every transcript. Another challenge is the 685 significant false positive and false negative rates for RNA-seq. RNA-seq studies 686 often accept an estimated false discovery rate of 10%, and, although false negative 687 rates are unknown, estimates are as high as 90% for mammalian cells (Pimentel et al., 688 2017). These rates seriously impair attempts to classify transcripts into any one of the 689

27 possible classes. If parameters are controlled at a rate of 10%, then the probability 690  
that at least one of the parameters in a dense class (classes where all 3 parameters 691  
are + or -, not 0) has been falsely accepted is almost 27%. Thus, almost one in three 692  
of the transcripts categorized into one of the 8 possible dense epistatic classes (+++, 693  
---, +-+, etc...) is misclassified and instead belongs to one of the twelve doublet 694  
epistatic classes (++0, 0--, etc...). The situation becomes considerably worse once 695  
we consider false negative rates, which are generally unknown but estimates range 696  
up to 90% in mammalian systems (Pimentel et al., 2017). In general, false rates 697  
greatly exacerbate the difficulties associated with analyzing transcriptomic datasets. 698  
If all transcripts actually belonged to a single epistasis class to begin with, the 699  
addition of statistical noise will split this class into many more classes that mimicry 700  
complex interactions. The situation is further worsened by the fact that interaction 701  
parameters can often be harder to measure than first order parameters. Classifying 702  
transcripts into epistatic classes is a major obstacle for successful epistatic analyses, 703  
and so far there has been little to no work done to assess which classes are real and 704  
which are artifactual (some work has been done in the context of allelic series, see 705  
page 14). Equally concerning is the fact that none of these epistatic classes can be 706  
translated into genetic diagrams. These epistatic classes do not provide a biological 707  
mechanism (genetic, biochemical or cellular) between the genes under study (see 708  
Fig. 12). 709

### Occam's razor, information pooling and constrained epistasis 710

To extract biological mechanisms from transcriptome data, we must apply simplifying- 711  
ing constraints. If transcripts are to be classified into 27 possible epistatic classes, 712  
we must develop methods to assess which of these classes have sufficient statistical 713  
leverage to accept their existence (in other words, we need a statistical test that exam- 714  
ines the null hypothesis that such a class could appear purely by chance). However, 715



**Figure 12** Analysis methodology to infer genetic interactions using transcriptome data. **A.** After fitting all transcripts to a general linear model to calculate the individual and the epistatic components of null mutations in two distinct genes, the resulting parameters can be clustered and visualized in a heatmap. Each observed cluster can be grouped into one of 27 epistatic classes. All clusters are considered biologically relevant regardless of the number of transcripts they contain. A simple conclusion cannot be reached from these heatmaps. This approach was used in Dixit et al. (2016) and Adamson et al. (2016) **B.** Starting from the same statistical model, only transcripts that have all parameters different from zero are considered informative. These transcripts are plotted on a scatterplot, where the x-axis reflects the expected value of the double mutant under an additive or log-additive hypothesis, and the systematic deviation from additivity (generalized epistasis) is plotted on the y-axis. The resulting points form a ray on the plot. The slope of this ray is an aggregate statistic that can be interpreted in terms of a genetic pathway if the two genes exhibit Batesonian epistasis. This approach was used in Angeles-Albores, Leighton, et al. (2017) and Angeles-Albores, Puckett Robinson, et al. (2018)

even if such a test were developed, we still require a method that allows us to summarize the information in these modules, and which lets us build a genetic pathway if the data suggests a pathway exists. A natural way to do this may be to use the natural structure of biological networks to pool the information from all transcripts, and test the interaction of this *structure* between the two mutants, instead of testing the individual transcripts. Information sharing is a powerful concept that allows us to incorporate more data points into a calculation, thus increasing our statistical power for any single test, but it requires the data to be drawn from a structure that permits sharing.

One such information sharing approach was implemented in Angeles-Albores, Puckett Robinson, et al. (2018) and Angeles-Albores, Leighton, et al. (2017). Briefly, these studies obtained whole-organism bulk RNA-seq transcriptome profiles for single and double perturbations and identified differentially expressed transcripts in each condition relative to the wild-type. Next, transcripts that were differentially expressed in all non-control conditions were aggregated and analyzed jointly for systematic deviations from a linear pathway. This systematic deviation was quantified in a single coefficient, called the transcriptome-wide epistasis coefficient. This coefficient can be interpreted in terms of simple genetic pathways because it can be used to test whether the perturbations result in a phenotypic transformation diagnostic of Batesonian epistasis. In this sense, the transcriptome-wide epistasis coefficient represents a unification of generalized epistasis and classical epistasis. This approach is powerful because it avoids multiple hypothesis testing (a single interaction coefficient is tested), and it doesn't rely on any one transcript to draw conclusions. A significant advantage of this method is that these studies were able to test and verify that the generalized epistasis measurements they made were equivalent to Batesonian epistasis (in other words, the double mutant had the same perturbations as one of the single mutants), culminating in a formal genetic

pathway. Both studies assumed that the genetic interaction between two genes is 743 unimodal, in other words, these two genes do not interact along multiple pathways 744 with different strengths and valences. This last assumption may not always hold. 745 This strongly simplifying assumption contrasts with the previously referenced work 746 that assumes unbounded complexity for all genetic interactions. Neither is correct, 747 though it is our opinion that biological interactions tend to be much simpler than is 748 often assumed in genomic studies. 749

### Beyond genetic interactions: Dominance studies to map gene functions

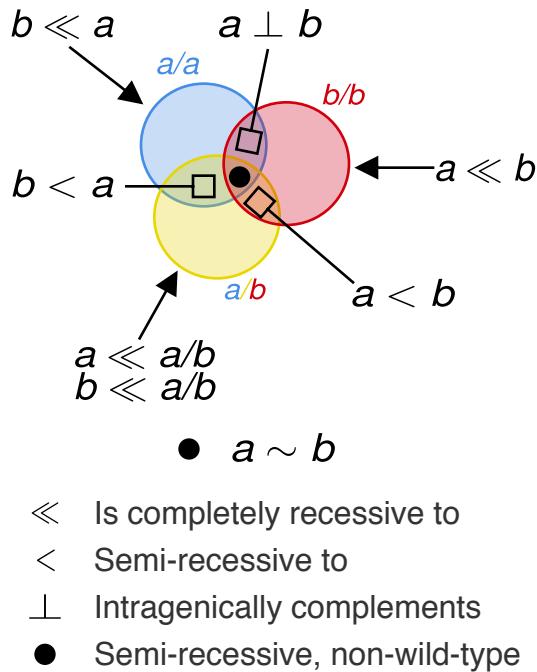
Although transcriptomes have been used as phenotypes for analysis of genetic in- 751 teractions for many years, their uses need not be restricted for theoretic analysis. 752 In population genetics, transcriptomes have been used as phenotypes with which 753 to identify expression quantitative trait loci in a number of organisms (Brem et al., 754 2002; DeCook et al., 2006; Kirst et al., 2004; Schadt et al., 2003). Transcriptomes 755 can also be used to compare the genetical properties of different alleles of a single 756 gene (Angeles-Albores and Sternberg, 2018). 757

Allelic series require considerably more analysis than tests for genetic interactions. 758 To infer functional units from the activity of multiple allelic variants, the phenotypes 759 associated with each variant must be carefully enumerated. Alleles must be ordered 760 according to the phenotypic severity they cause when animals are homozygotes for 761 each variant, with a separate hierarchy drawn for each phenotype. Alleles must 762 also be ordered according to their dominance hierarchy over other alleles along each 763 phenotype by measured the phenotypes of *trans*-heterozygotes. Particular care must 764 be taken to ensure that the phenotypes of the *trans*-heterozygotes are not the result 765 of maternal effects by testing progeny generated from a second, reciprocal, cross. 766 The overall results are examined and the most parsimonious explanation is accepted 767 to draw functional units and establish their sequence requirements. The the number 768

and resolution of the functional units that can be defined depends on the density of 769  
the allelic series that is tested. For a more thorough introduction to dominance and 770  
its role in allelic series, see Yook (2005). 771

As a result of the rigor required to analyze them, allelic series provides an excellent 772  
testing ground in which to explore the potential, but also the shortcomings, of 773  
transcriptomes as molecular phenotypes. To be successful, the analysis of even the 774  
smallest allelic series must order the tested variants. Angeles-Albores and Sternberg 775  
(2018) reported the first allelic series, to our knowledge, to be analyzed using 776  
expression profiles in any organism. In this analysis, the transcriptomic analogue 777  
of distinct phenotypes, phenotypic classes consisting of groups of differentially 778  
expressed genes, were identified by labelling each gene with the genotypes where 779  
it was differentially expressed. Subsequently, the expression level of these genes 780  
in *trans*-heterozygotes was approximated by a linear combination of the expression 781  
levels in each homozygote, with the weighting coefficients constrained to add to 782  
unity. The weighting coefficients, bounded in this manner, reflect the dominance 783  
of one allele over the other. These transcriptome-wide dominance coefficients are 784  
analogous to the transcriptome-wide epistasis aggregate statistics derived in previous 785  
studies (Angeles-Albores, Puckett Robinson, et al., 2018). The intersections from 786  
the Venn diagram (see Fig. 13) are understood to occur as a result of the activity of 787  
one or more functional units which may or may not have dosage-saturated activity 788  
(this is inferred from the dominance behavior of the given intersection). 789

This study highlighted the importance of recognizing and characterizing the sta- 790  
tistical artifacts that can occur in genomic datasets (see Fig. 14). The analyzed 791  
dataset had sufficiently large false positive and false negative rates to generate arti- 792  
ficial phenotypic classes that nevertheless could be identified and removed from the 793  
analysis. Unlike epistatic classes, for which we do not have a sense of what classes 794  
can most easily arise as a result of statistical artifacts, all the phenotypic classes 795

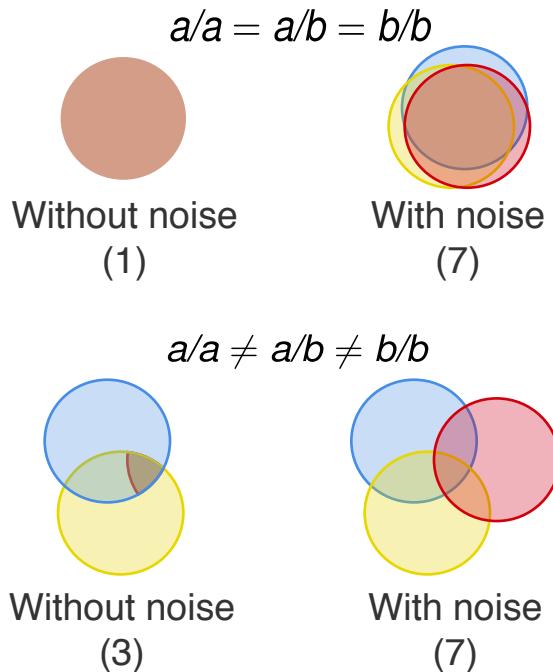


**Figure 13** Genes that are differentially expressed in genotypes containing mutant ( $a, b$ ) alleles relative to a wild type homozygote can be categorized into phenotypic classes. Each phenotypic class can in turn be associated with a dominance behavior. The Venn diagram represents differentially expressed transcripts in each genotype relative to the wild-type control. Each of the possible 7 intersections is labelled with its dominance interpretation if the intersection is real. In this context, semi-recessiveness means that one allele is partially or completely dominant to the other along a continuous spectrum between 0 and 1. The dominance sign between an allele and the heterozygote genotype indicates heterosis or over-dominance.

arising from allelic series analyses can be readily interpreted in terms of inter-allelic complementation, a phenomenon that is extremely well characterized in genetics. Allelic series provide an excellent testing ground in which to explore algorithms to partition transcriptomes into gene batteries that have sufficient statistical support, since it is possible to have an intuition for artifactual classes.

### Open problems and opportunities

RNA-sequencing is becoming increasingly easier and cheaper. RNA-seq offers a powerful, unbiased approach to genetics that can be multiplexed in many systems relatively easily. We expect that genetics using expression profiles will be an



**Figure 14** RNA-seq artifacts can greatly exaggerate apparent biological complexity. We considered the case where we have two phenotypically identical alleles that can be used to generate the genotypes  $a/a$ ,  $b/b$  and  $a/b$ . In the absence of artifacts, the set of differentially expressed transcripts relative to a wild-type control should be the same amongst all three genotypes. However, if measurement error occurs, then instead of observing a single Venn intersection, we will observe seven intersections. If these intersections are not identified as false, we would wrongly conclude that allele  $a$  and  $b$  are not phenotypically equivalent, incurring in an error rate of 600%. Even in the case where the three genotypes are not equivalent, statistical noise will tend to significantly increase the apparent biological complexity present in the system (from 3 to 7 in this example). In general, statistical artifacts are so common in genomic assays that they will tend to generate all the possible intersections in a comparison. This highlights the need to apply simplifying constraints on transcriptome data before interpreting the results.

excellent first-pass assay because of the speed and sheer amount of information 805 associated with the generation of expression profiles. These properties make RNA- 806 seq particularly advantageous for groups that are studying relatively unknown genes 807 or genes with subtle phenotypes. RNA-seq may also be a powerful method to 808 complement genetics in emerging model organisms where conventional genetics 809 may be laborious and where researchers may wish to minimize the number of 810 experiments performed while maximizing the amount they can learn. 811

A major challenge moving forward will be mixed epistasis analyses with allelic 812 series. Such mixed analyses try to identify the sequence requirements of one gene 813 to participate in an epistatic interaction, and to test whether the observed epistatic 814 interaction between two genes reflects a single biochemical function or the joint 815 activity of distinct molecular properties. For example, in *C. elegans* the inhibition 816 of *hif-1* by *egl-9* is mediated partially by the hydroxylation of HIF-1 by EGL- 817 9, and partially through a hydroxylation-independent mechanism that is not well 818 understood (Shao, Zhang, and Powell-Coffman, 2009). The high false positive and 819 false negative rates inherent to RNA-seq means that all interactions amongst all 820 genes will appear to be the compounded result of many independent activities. The 821 solution to this problem will require methods that can incorporate information not 822 just between single and double mutants, or homozygotes and heterozygotes, but 823 amongst epistatic modules and dominance modules while searching for the most 824 parsimonious structure that can explain all the expression profiles. 825

A second challenge will be the association of gene batteries with other observable 826 phenotypes to develop signatures that allow us to read and interpret a transcriptome 827 in terms of biological covariates. In other words, we would like signatures that 828 allowed us to infer what the organism was doing when the RNA was extracted, 829 what pathways had been disrupted or activated, what cellular or morphological 830 phenotypes it exhibited. Such signatures could be derived by allowing organisms 831

to undergo a specific life history, then extracting the transcriptome and associating 832  
the differentially expressed genes in response to this life history relative to a control 833  
history to derive a signature. Alternatively, single-cell or single-organism methods 834  
may be able to track organisms, recording their behavior, before extracting their 835  
RNA (Lane et al., 2017). These signatures, although useful, should not be treated as 836  
causal, because the derivation of these signatures is through correlation. Deriving 837  
causal signatures would be very interesting and potentially useful as well, since this 838  
would make the discovery and association of novel pathways considerably easier. A 839  
significant weakness of expression signatures is that they only make sense relative 840  
to a baseline control, and therefore signatures can only be associated with events 841  
that have a sufficient dynamic range relative to the baseline. Another problem with 842  
signatures is the arbitrary definition, since they will inevitably be defined according 843  
to a *q*-value cut-off. It seems reasonable to postulate that eventually we must 844  
abandon the concept of differential expression: It is too brittle, too relativistic and 845  
prevents us from thinking about the transcriptome as a complete object. 846

Without transcriptional signatures of some sort, understanding modules will be all 847  
but impossible. Even with signatures, modules will be explained only phenomeno- 848  
logically: We know this signature is correlated to this phenotype, therefore this 849  
module is correlated to the same phenotype. With time, we may be able to under- 850  
stand mechanistically why specific phenotypes are correlated with the expression of 851  
specific genes. For the moment, such understanding seems far from our reach. 852

In the end, the major challenge for transcriptome genetics is likely to be our own 853  
creativity. New phenotypes always have their difficulties and drawbacks, and expres- 854  
sion profiles are no exception. Expression profiles will not, on their own, reconstruct 855  
every network or solve all of biology. However, expression profiles are an object 856  
of a new kind, with behaviors that we do not fully understand hiding novel biolog- 857  
ical phenomena. It has become evident that genetics is applicable at an enormous 858

range of phenotypes, from population phenotypes to organismal to macroscopic and mesoscopic phenotypes. Transcriptomes represent a new phenotype at the microscopic and genomic level. Perhaps surprisingly, these microscopic phenotypes, in spite of all their complexity, seem to obey the genetic properties that bind all other phenotypes. The challenge, then, is how to use transcriptomes to discover biological principles that help us understand how the hierarchy of cells, organs, organisms and populations emerges from the collective actions of a string of atoms.

## References

- Adamson, Britt et al. (2016). “A Multiplexed Single-Cell CRISPR Screening Platform Enables Systematic Dissection of the Unfolded Protein Response”. In: *Cell* 167.7, 1867–1882.e21. ISSN: 00928674. doi: [10.1016/j.cell.2016.11.048](https://doi.org/10.1016/j.cell.2016.11.048).
- Angeles-Albores, David, Daniel H. W. Leighton, et al. (2017). “The *Caenorhabditis elegans* Female-Like State: Decoupling the Transcriptomic Effects of Aging and Sperm Status”. In: *G3: Genes, Genomes, Genetics* 7.9.
- Angeles-Albores, David, Raymond Y. N. Lee, et al. (2016). “Tissue enrichment analysis for *C. elegans* genomics”. In: *BMC Bioinformatics* 17.1, p. 366. ISSN: 1471-2105. doi: [10.1186/s12859-016-1229-9](https://doi.org/10.1186/s12859-016-1229-9).
- Angeles-Albores, David, Carmie Puckett Robinson, et al. (2018). “Reconstructing a metazoan genetic pathway with transcriptome-wide epistasis measurements”. In: *Proceedings of the National Academy of Sciences*, p. 201712387. ISSN: 0027-8424. doi: [10.1073/pnas.1712387115](https://doi.org/10.1073/pnas.1712387115).
- Angeles-Albores, David and Paul W Sternberg (2018). “Using Transcriptomes as Mutant Phenotypes Reveals Functional Regions of a Mediator Subunit in *Caenorhabditis elegans*.” In: *Genetics*, genetics.301133.2018. ISSN: 1943-2631. doi: [10.1534/genetics.118.301133](https://doi.org/10.1534/genetics.118.301133).
- Avery, Leon and Steven Wasserman (1992). *Ordering gene function: the interpretation of epistasis in regulatory hierarchies*. doi: [10.1016/0168-9525\(92\)90263-4](https://doi.org/10.1016/0168-9525(92)90263-4). arXiv: [NIHMS150003](https://arxiv.org/abs/NIHMS150003).
- Bateson, William and Gregor Mendel (2009). *Mendel's principles of heredity: A defence, with a translation of mendel's original papers on hybridisation*, pp. 1–212. ISBN: 9780511694462. doi: [10.1017/CBO9780511694462](https://doi.org/10.1017/CBO9780511694462).
- Bray, Nicolas L et al. (2016). “Near-optimal probabilistic RNA-seq quantification.” In: *Nature biotechnology* 34.5, pp. 525–7. ISSN: 1546-1696. doi: [10.1038/nbt.3519](https://doi.org/10.1038/nbt.3519). arXiv: [1505.02710](https://arxiv.org/abs/1505.02710).

- Brem, Rachel B. et al. (2002). "Genetic Dissection of Transcriptional Regulation in Budding Yeast". In: *Science* 296.5568. 893  
894
- Britten, Roy J. and Eric H. Davidson (1969). "Gene regulation for higher cells: A theory". In: *Science* 165.3891, pp. 349–357. issn: 00368075. doi: [10.1126/science.165.3891.349](https://doi.org/10.1126/science.165.3891.349). 895  
896  
897
- Capaldi, Andrew P et al. (Nov. 2008). "Structure and function of a transcriptional network activated by the MAPK Hog1". In: *Nature Genetics* 40.11, pp. 1300–1306. issn: 1061-4036. doi: [10.1038/ng.235](https://doi.org/10.1038/ng.235). 898  
899  
900
- Chan, Io Long, Oliver J Rando, and Colin C Conine (2018). "Effects of Larval Density on Gene Regulation in *Caenorhabditis elegans* During Routine L1 Synchronization". In: *G3: Genes|Genomes|Genetics* 8.5, 1787 LP –1793. issn: 2160-1836. doi: [10.1534/g3.118.200056](https://doi.org/10.1534/g3.118.200056). 901  
902  
903  
904
- Cordell, Heather J (2002). "Epistasis: what it means, what it doesn't mean, and statistical methods to detect it in humans". In: *Human Molecular Genetics* 11.20, pp. 2463–2468. doi: [10.1093/hmg/11.20.2463](https://doi.org/10.1093/hmg/11.20.2463). 905  
906  
907
- Costanzo, Michael et al. (2016). "A global genetic interaction network maps a wiring diagram of cellular function". In: *Science* 353.6306. issn: 10959203. doi: [10.1126/science.aaf1420](https://doi.org/10.1126/science.aaf1420). arXiv: [arXiv:1011.1669v3](https://arxiv.org/abs/1011.1669v3). 908  
909  
910
- DeCook, Rhonda et al. (2006). "Genetic regulation of gene expression during shoot development in *Arabidopsis*". In: *Genetics* 172.2, pp. 1155–1164. issn: 00166731. 911  
912  
913
- Dixit, Atray et al. (2016). "Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens". In: *Cell* 167.7, 1853–1866.e17. issn: 00928674. doi: [10.1016/j.cell.2016.11.038](https://doi.org/10.1016/j.cell.2016.11.038). 914  
915  
916
- Fisher, R. A. (1919). "XV.—The Correlation between Relatives on the Supposition of Mendelian Inheritance". In: *Transactions of the Royal Society of Edinburgh* 52.2, pp. 399–433. issn: 00804568. doi: [10.1017/S0080456800012163](https://doi.org/10.1017/S0080456800012163). 917  
918  
919
- Huang, Linda S and Paul W Sternberg (2006). "Genetic dissection of developmental pathways." In: *WormBook: the online review of C. elegans biology* 1995, pp. 1–19. issn: 1551-8507. doi: [10.1895/wormbook.1.88.2](https://doi.org/10.1895/wormbook.1.88.2). 920  
921  
922
- Kirst, Matias et al. (2004). "Coordinated genetic regulation of growth and lignin revealed by quantitative trait locus analysis of cDNA microarray data in an interspecific backcross of eucalyptus." In: *Plant physiology* 135.4, pp. 2368–78. issn: 0032-0889. doi: [10.1104/pp.103.037960](https://doi.org/10.1104/pp.103.037960). 923  
924  
925  
926
- Lane, Keara et al. (Apr. 2017). "Measuring Signaling and RNA-Seq in the Same Cell Links Gene Expression to Dynamic Patterns of NF- $\kappa$ B Activation". In: *Cell Systems* 4.4, 458–469.e5. issn: 24054712. doi: [10.1016/j.cels.2017.03.010](https://doi.org/10.1016/j.cels.2017.03.010). 927  
928  
929  
930

- Lehner, Ben (Aug. 2011). “Molecular mechanisms of epistasis within and between genes”. In: *Trends in Genetics* 27.8, pp. 323–331. ISSN: 01689525. DOI: [10.1016/j.tig.2011.05.007](https://doi.org/10.1016/j.tig.2011.05.007). 931  
932  
933
- Lott, Susan E. et al. (2011). “Noncanonical compensation of zygotic X transcription in early *Drosophila melanogaster* development revealed through single-embryo RNA-Seq”. In: *PLoS Biology* 9.2. ISSN: 15449173. DOI: [10.1371/journal.pbio.1000590](https://doi.org/10.1371/journal.pbio.1000590). 934  
935  
936  
937
- Mi, Huaiyu et al. (2009). “PANTHER version 7: Improved phylogenetic trees, orthologs and collaboration with the Gene Ontology Consortium”. In: *Nucleic Acids Research* 38.SUPPL.1. ISSN: 03051048. DOI: [10.1093/nar/gkp1019](https://doi.org/10.1093/nar/gkp1019). 938  
939  
940
- Mortazavi, Ali et al. (2008). “Mapping and quantifying mammalian transcriptomes by RNA-Seq”. In: *Nature Methods* 5.7, pp. 621–628. ISSN: 1548-7091. DOI: [10.1038/nmeth.1226](https://doi.org/10.1038/nmeth.1226). arXiv: [1111.6189v1](https://arxiv.org/abs/1111.6189v1). 941  
942  
943
- Patro, Rob, Geet Duggal, et al. (2016). “Salmon provides accurate, fast, and bias-aware transcript expression estimates using dual-phase inference”. In: *bioRxiv*, p. 021592. DOI: [10.1101/021592](https://doi.org/10.1101/021592). arXiv: [1505.02710](https://arxiv.org/abs/1505.02710). 944  
945  
946
- Patro, Rob, Stephen M. Mount, and Carl Kingsford (2014). “Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms”. In: *Nature biotechnology* 32.5, pp. 462–464. ISSN: 1546-1696. DOI: [10.1038/nbt.2862](https://doi.org/10.1038/nbt.2862). arXiv: [1308.3700](https://arxiv.org/abs/1308.3700). 947  
948  
949  
950
- Phillips, Patrick C (2008). “Epistasis—the essential role of gene interactions in the structure and evolution of genetic systems”. In: *Nat Rev Genet* 9.11, pp. 855–867. ISSN: 1471-0056. DOI: [10.1038/nrg2452](https://doi.org/10.1038/nrg2452). 951  
952  
953
- Picelli, Simone et al. (2014). “Full-length RNA-seq from single cells using Smart-seq2.” In: *Nature protocols* 9.1, pp. 171–81. ISSN: 1750-2799. DOI: [10.1038/nprot.2014.006](https://doi.org/10.1038/nprot.2014.006). 954  
955  
956
- Pimentel, Harold et al. (2017). “Differential analysis of RNA-seq incorporating quantification uncertainty”. In: *Nature Methods* 14.7, pp. 687–690. ISSN: 15487105. DOI: [10.1038/nmeth.4324](https://doi.org/10.1038/nmeth.4324). 957  
958  
959
- Sameith, Katrin et al. (2015). “A high-resolution gene expression atlas of epistasis between gene-specific transcription factors exposes potential mechanisms for genetic interactions”. In: *BMC Biology* 13.1. ISSN: 17417007. DOI: [10.1186/s12915-015-0222-5](https://doi.org/10.1186/s12915-015-0222-5). 960  
961  
962  
963
- Schadt, Eric E. et al. (Mar. 2003). “Genetics of gene expression surveyed in maize, mouse and man”. In: *Nature* 422.6929, pp. 297–302. ISSN: 00280836. DOI: [10.1038/nature01434](https://doi.org/10.1038/nature01434). 964  
965  
966
- Serra, Lorryne et al. (2018). “Adapting the Smart-seq2 Protocol for Robust Single Worm RNA-seq”. In: *BIO-PROTOCOL* 8.4. ISSN: 2331-8325. DOI: [10.21769/BioProtoc.2729](https://doi.org/10.21769/BioProtoc.2729). 967  
968  
969

- Shao, Zhiyong, Yi Zhang, and Jo Anne Powell-Coffman (2009). “Two Distinct Roles for EGL-9 in the Regulation of HIF-1-mediated gene expression in *Caenorhabditis elegans*”. In: *Genetics* 183.3, pp. 821–829. ISSN: 00166731. DOI: [10.1534/genetics.109.107284](https://doi.org/10.1534/genetics.109.107284). 970  
971  
972  
973
- Tang, Fuchou et al. (May 2009). “mRNA-Seq whole-transcriptome analysis of a single cell”. In: *Nature Methods* 6.5, pp. 377–382. ISSN: 1548-7091. DOI: [10.1038/nmeth.1315](https://doi.org/10.1038/nmeth.1315). 974  
975  
976
- Van De Peppel, Jeroen et al. (2005). “Mediator expression profiling epistasis reveals a signal transduction pathway with antagonistic submodules and highly specific downstream targets”. In: *Molecular Cell* 19.4, pp. 511–522. ISSN: 10972765. DOI: [10.1016/j.molcel.2005.06.033](https://doi.org/10.1016/j.molcel.2005.06.033). 977  
978  
979  
980
- Van Driessche, Nancy et al. (2005). “Epistasis analysis with global transcriptional phenotypes”. In: *Nature genetics* 37.5, pp. 471–477. ISSN: 1061-4036. DOI: [10.1038/ng1545](https://doi.org/10.1038/ng1545). 981  
982  
983
- Yook, Karen (2005). “Complementation”. In: *WormBook*. ISSN: 15518507. DOI: [10.1895/wormbook.1.24.1](https://doi.org/10.1895/wormbook.1.24.1). 984  
985

Chapter 2

986

# A STATISTICAL MECHANICAL THEORY OF GENETICS USING GENE EXPRESSION PHENOTYPES

987

988

## Abstract

989

Genetics has been a cornerstone of biology since its inception almost 100 years ago. However, the rules of genetic inference remain axiomatic, and as a result, confusion regarding these rules abounds. Here, I show that Batesonian epistasis, which has been used to construct genetic networks, emerges naturally from a statistical mechanical framework. I will show how the inferences of genetic interactions on the basis of Batesonian epistasis emerge from this framework, thus establishing that the methods geneticists use to identify interactions amongst genes can be viewed as a non-analytical variational method to probe an unknown partition function. By establishing Batesonian epistasis, this method can constrain the properties of this unknown partition function and reveal how variables within this partition function are functionally related to each other.

1001

## Introduction

1002

Biological processes can be broadly subdivided into two large categories: Processes where the players and interactions are known, and processes where the players or the interactions are unknown. Necessarily, the study of processes in either category varies substantially, as do the available methodologies. Processes where the components are unknown cannot be studied in a predictive fashion and mathematical methods cannot be used to identify the components that participate in the process (though they can make statements about the properties these components must exhibit). On the other hand, when the major components and interactions of a

1009

system are known, detailed models can be built to study emergent properties present 1011  
in the system or to identify the regimes in which the model breaks down, thus 1012  
indicating missing information about the system. 1013

A particularly effective method for modeling a large number of physical phenom- 1014  
ena is statistical mechanics. Statistical mechanics dictates that if all the possible 1015  
configurations of a system are known (all the states and their weights), then the 1016  
probability that an event happens is equal to the weighted sum of the states that 1017  
can lead to that event divided by the weighted sum of all the possible states. The 1018  
weighted sum of all the possible states is particularly important and is called the 1019  
Partition Function. Though this theory was originally applied to ideal gases, it 1020  
has become widespread, finding applications in all areas of physics, but also in 1021  
biology (Garcia et al., 2007), neurobiology (Schneidman et al., 2006) and even the 1022  
social sciences (Lee, Broedersz, and Bialek, 2015). Statistical mechanical models 1023  
demand a detailed and quantitative knowledge of the system to be tractable. 1024

In contrast to the quantitative description of a physical phenomenon obtained from 1025  
statistical mechanical models, genetics is a favored approach when the system is 1026  
completely unknown, because the genetical framework is explicitly designed to 1027  
make as few assumptions as possible. Saturating genetic screens can be used to find 1028  
all of the genes that cause a specific phenotype when mutated (Sulston and Brenner, 1029  
1974; Jürgens et al., 1984). Once these genes are known, epistatic analyses using 1030  
null mutants of these genes can establish genetic interactions between these genes. 1031

Genetic interactions describe the broadest possible mode of interaction in a bio- 1032  
logical system. Formally, two genes genetically interact if they affect the same 1033  
phenotype. However, not all gene interactions are considered equally informative; 1034  
the set of uninformative genetic interactions are defined by the choice of null hy- 1035  
pothesis, which is often selected to be a multiplicative null hypothesis (also known 1036

as log-additive null hypothesis). A more strict and informative definition of genetic interactions requires:

1038

1. The two genes control the same phenotype; 1039
2. The phenotype of the double null mutant deviates from that expected under a 1040  
null model of interaction 1041

The deviation between the expected phenotype and the observed phenotype of the double mutant is an important quantity that we refer to as generalized epistasis (Fisher, 1919). Although generalized epistasis can take on any value in theory, in the context of developmental genetics, generalized epistasis measured between two genes using null mutants often reduces to a limited set of specific values. These values correspond to a phenomenon we refer to as complete epistasis, which is important for classical developmental genetics. Complete epistasis refers to a situation where the double mutant exhibits exactly the same phenotype as one of the single mutants, originally defined by Bateson (Bateson and Mendel, 2009). Measuring this equality is equivalent to constraining the value of the generalized epistasis term to a unique value. The relationship between classical and generalized epistasis has been a source of great confusion, not least because multiple fields use conflicting definitions for the term and use it to measure different things (Cordell, 2002; Phillips, 2008). A major problem is that the presence of generalized epistasis depends strongly on the choice of null model. In this text, we only use the term epistasis, either complete or generalized, to refer to the situation where the double mutant exhibits the same phenotype as one of the corresponding single mutants. We will use the word ‘mutant’ to refer to a null allele, unless otherwise specified. Finally, we will use the term ‘epistasis coefficient’ to refer to the statistical quantity of generalized epistasis, which depends on a null model.

1061

Although genetics is used to study physical systems, there is not a systematic map- 1062  
 ping of the language of genetics into the language of statistical mechanics. A 1063  
 consequence of this is that genetic interactions are often considered vague to the 1064  
 point of being uninformative. Here, we show that genetics is a method for studying 1065  
 properties of an unknown partition function. Identifying genetic interactions be- 1066  
 tween two components helps constrain the functional form of the partition function, 1067  
 which in turn constrains the states that are accessible to the system. 1068

### Statistical Mechanics of Complete Epistasis

We have previously observed complete epistasis using gene expression profiles as 1070  
 phenotypes (Angeles-Albores, Leighton, et al., 2017; Angeles-Albores, Puckett 1071  
 Robinson, et al., 2018). To study how this can occur, we use a model of gene expres- 1072  
 sion derived from statistical mechanics that has been used to accurately describe 1073  
 various promoter architectures in multiple organisms (Bintu et al., 2005; Garcia 1074  
 et al., 2007; Bothma et al., 2015; Raveh-Sadka, Levo, and Segal, 2009). Briefly, the 1075  
 model can be expressed in the following functional form (for a detailed and friendly 1076  
 derivation, see Bintu et al. (2005) or Garcia et al. (2007)): 1077

$$p_{\text{bound}}(A, B, \dots) = \frac{1}{1 + \frac{1}{\rho F_{\text{reg}}(A, B, \dots)}}. \quad (2.1)$$

Here,  $p_{\text{bound}}(A, B, \dots)$  is the probability that RNA Polymerase is bound to the locus 1078  
 of interest (this probability is in turn proportional to the expression level of this 1079  
 gene ()), and which depends on several factors  $\rho, A, B, \dots$ ;  $\rho$  describes the concen- 1080  
 tration of RNA polymerase bound to the promoter;  $F_{\text{reg}}$  is a factor that modulates the 1081  
 effective number of RNA polymerases at the promoter of interest. This factor can 1082  
 be used to model a variety of transcription factors, such as activators, or inhibitors, 1083  
 and can be used to model interactions between these factors. Its exact value will de- 1084  
 pend on factors  $A, B, \dots$ . These factors represent the activities of the gene products 1085

of genes  $a, b, \dots$ . For simplicity, capital letters will always represent compound activities of the final gene products (proteins), whereas lower-case italicized letters will always represent the genes encoding said products, in accord with *C. elegans* genetic nomenclature rules. Equation 2.1, although superficially simple, is able to accomodate an enormous variety of promoter architectures through its dependence on  $F_{reg}$ .

For example, consider a system where two proteins,  $A$  and  $B$  can independently bind different sequences on a promoter, and can independently bind RNA Polymerase. Suppose that  $A$  and  $B$  do not interact with each other. Then,  $F_{reg}$  takes the following functional form:

$$F_{reg}(A, B) = \frac{1 + Ae^{-\varepsilon_{AP}} + Be^{-\varepsilon_{BP}}}{1 + A + B}, \quad (2.2)$$

where  $\varepsilon_{XP}$  is the interaction energy between protein  $X$  and RNA Polymerase. If  $A$  and  $B$  can bind to each other as well as to RNA polymerase, then  $F_{reg}$  becomes:

$$F_{reg}(A, B) = \frac{\frac{1+Be^{-\varepsilon_{BP}}}{1+B} + Ae^{-\varepsilon_{AP}} \cdot \frac{1+Be^{-\varepsilon_{AB}-\varepsilon_{BP}}}{1+B}}{1 + A \frac{1+Be^{-\varepsilon_{AB}}}{1+B}}, \quad (2.3)$$

where  $\varepsilon_{XY}$  denotes the energy of interaction between protein  $X$  and protein  $Y$ . Finally, in the case where protein  $B$  binds DNA but does not interact with RNA polymerase ( $\varepsilon_{BP} = 0$ ), this equation simplifies to:

$$F_{reg}(A, B) = \frac{1 + Ae^{-\varepsilon_{AP}} \cdot \frac{1+Be^{-\varepsilon_{AB}}}{1+B}}{1 + A \frac{1+Be^{-\varepsilon_{AB}}}{1+B}}, \quad (2.4)$$

We have shown how  $F_{reg}$  accomodates a variety of activator architectures.  $F_{reg}$  can similarly accomodate a variety of repressive architectures as well as mixed cases where proteins are both repressors and activators.

Next, we suppose that the two genetic factors we are studying can be expressed via  $F_{reg}$ . Using our transcriptional reporter as a phenotype, we would like to know what

the conditions are for complete epistasis to occur. Let the two genes under study 1106  
 be  $a$  and  $b$ , with protein products  $A$  and  $B$  respectively. Since epistasis analyses 1107  
 can only be safely carried out using null alleles, we can specify epistasis in our toy 1108  
 system as an equation (we let  $a$  be the epistatic factor without loss of generality)<sup>1</sup>: 1109

$$p_{bound}(0, 0) = p_{bound}(0, B). \quad (2.5)$$

Since  $p_{bound}$  depends on  $A$  and  $B$  only through the factor,  $F_{reg}(A, B)$ , Eq. 2.5 can 1110  
 be re-written to more explicitly show the requirement for epistasis in terms of the 1111  
 regulatory function,  $F_{reg}$ : 1112

$$F_{reg}(0, 0) = F_{reg}(0, B). \quad (2.6)$$

We are interested in identifying a mathematical condition that will satisfy Eq. 2.5 and 1113  
 Eq. 2.6 and which is biologically relevant. We noticed that Eq. 2.6 can be guaranteed 1114  
 if the variables  $A$  and  $B$  can be combined into a single compound variable: 1115

$$F_{reg}(A, B) = F_{reg}[A \cdot G_{reg}(B)]. \quad (2.7)$$

In other words, Eq. 2.5 is satisfied if the only function of  $b$  is to genetically alter the 1116  
 effective genetic activity of  $a$  through a regulatory function,  $G_{reg}$ . This condition 1117  
 is relatively lax, and in fact is frequently a property of signaling pathways() and of 1118  
 many promoter architectures (Bintu et al., 2005). No conditions are imposed on 1119  
 how  $a$  interacts with the promoter. 1120

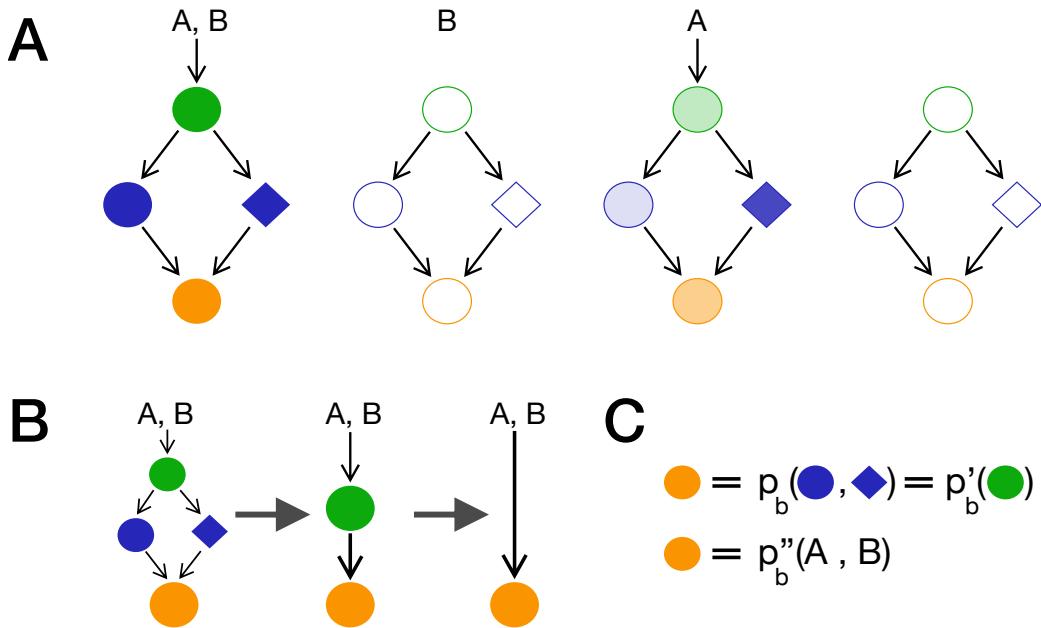
---

<sup>1</sup>For simplicity, we represent the effect of a null mutation as deleting all protein product. Another way to generate null alleles is by eliminating all biochemical activities, which would correspond to setting all the interaction energies for a given protein to zero. This alternative representation does not affect our argument.

**Genetic suppression does not require the two factors interact physically on the promoter sequence** 1121  
**1122**

In the above section, we modeled a situation where the protein products of our genes interacted directly with the promoter that drives the transcriptional reporter. 1123  
 However, the situation does not change at all if the protein products never interact 1125 with the promoter, but instead drive activity of another factor, *alpha*, that is the physical interactor. 1126 The numeric results may change, but the equality holds, as long as *alpha* accepts the inputs *A* and *B* in compound form. 1127 In other words, if  $\alpha(A, B)$  can be written in the form  $\alpha[h(A) \cdot j(B)]$ , then the equality in Eq. 2.5 holds. 1128  
 In other words, complete epistasis can occur between factors that are physically 1130 and/or temporally separated. Conversely, if two factors show complete epistasis, 1131 the only constraint this imposes on their interactions is that one gene must alter 1132 the effective gene activity of the other, and that the hypostatic gene (the gene for 1133 which the phenotype is masked) not interact with the phenotype through another 1134 pathway that is independent of the epistatic gene. Thus, genetic interaction can 1135 abbreviate the pathway components between the interacting genes, *a* and *b*, and 1136 represent them through single arrows. The above arguments strongly suggest that the 1137 genetic ‘distance’ between the genes under study and the transcriptional phenotype 1138 is irrelevant. 1139

Another important question is whether there is an ‘optimal’ choice of expression phenotype. 1140 In other words, for a given graph containing nodes (genes) that can 1141 interact (arrows) with other nodes, is there a ‘best’ node to select as a readout of the 1142 pathway (see Fig. 21)? For simplicity, but without loss of generality, we envision 1143 a graph where a single, primary node accepts as an input our two factors *A* and *B*. 1144 This node can in turn drive expression of a second layer of genes through arbitrary 1145 activation functions, and this second layer can drive expression of a third layer of 1146 (completely different) genes. Genes in this second layer are allowed to have any 1147



**Figure 21** Signaling pathways can be abbreviated in genetic pathway representations. **A.** Two genes,  $a$  and  $b$ , encoding protein products A and B respectively, interact to drive expression of a green gene, which in turn drives expression of two blue genes (through potentially distinct mechanisms). The blue genes then interact to drive expression of an orange gene.  $a$  is epistatic over  $b$ . Removing  $a$  sends the green node to a specific expression state. This state drives the blue genes into a specific state as well; then the blue genes drive the orange gene into a well-defined state. Removing  $b$  sends the network into a completely different state. Since  $a$  is epistatic over  $b$ , mutating both genes sends the network into the same state as mutating  $a$  alone. **B.** Since the network can only exist in three states (a wild type state, an  $a^-$  state and a  $b^-$  state), the network can be abbreviated to reflect that the state of the orange node can be completely specified by knowing the state (present or missing) of the upstream factors,  $a$  and  $b$ . **C.** This corresponds to stating that the probability that the orange node is transcribed can be expressed solely in terms of the protein products A and B when the products are expressed at wild-type levels or when one or the other product is completely missing.

arbitrary interaction with other genes in the second layer to drive expression of 1148 genes in the third layer. This layering can repeat arbitrarily. We seek to answer 1149 the question: For an arbitrary, acyclical graph, is there an ‘optimal’ node that will 1150 reveal the epistatic interactions between the two genetic factors  $a$  and  $b$ , assuming 1151 perfect measurements and zero biological noise? 1152

Let us represent the expression level of the  $j$ th node in the  $i$ th layer as  $p_{b,ij}$  following 1153 our previous notation. Clearly, the expression of the primary node,  $p_{b,00}$  depends 1154 on  $A$  and  $B$ . We can write the dependence of the state on these factors explicitly: 1155  $p_{b,00}(A, B)$ . Then, it follows that loss of  $A$  or  $B$  will send the state of this node to 1156 the specific states,  $p_{b,00}(0, B)$  and  $p_{b,00}(A, 0)$  respectively. Since  $a$  is epistatic over 1157  $b$ , and  $A$  and  $B$  interact to drive this node directly, then it follows that  $p_{b,00}(0, B) = 1158 p_{b,00}(0, 0)$ . The value  $p_{b,00}(0, B)$  is an important quantity, therefore we assign it the 1159 name  $P_{00}$ . This proves node 00 can function as a suitable phenotype for genetic 1160 analysis (by definition). 1161

The second set of nodes  $p_{b,1j}$  are regulated directly by node 00. We take the activity 1162 of node 00 to be directly proportional to  $p_{b,00}(A, B)$ . Thus, for any node  $1j$ , we 1163 can represent its expression level as  $p_{b,1j}(T[p_{b,00}(A, B)])$ , where  $T[\cdot]$  represents an 1164 arbitrary time-independent function that transforms the expression level of 00,  $p_{b,00}$ , 1165 into its product activity. Deletion of  $a$  or  $b$  transmits directly into the expression states 1166  $p_{b,1j}(T[P_{00}])$  and  $p_{b,1j}(T[p_{b,00}(A, 0)])$  respectively. Deleting  $a$  and  $b$  simultaneously 1167 will cause the any node in layer 1 to enter the state  $p_{b,1j}(T[P_{00}])$ . Since  $T[\cdot]$  is 1168 a time-independent function, it follows that any node in layer 1 can be used as a 1169 gene expression phenotype, regardless of any well-defined, non-trivial dependence 1170 it may have on node 00. We can iterate this logic throughout any arbitrary number of 1171 layers. A layer can drive expression of the next layer in any way without altering the 1172 ability of the next set of nodes to function as indicators of the epistasis relationship 1173 between  $a$  and  $b$ . 1174

Taken together, these results explain why genetic diagrams need not represent all factors between a pathway, and why genetic pathways are broadly insensitive to spatial or temporal separation between components: Biological and experimental noise would be the only limitations to detecting epistasis deep into a directed acyclic graph. These results also highlight the fact that the *genetic diagrams* that emerge from time-independent epistatic analysis are not necessarily informative about the dynamical behaviors of the *molecular interactions* that they result from because our genetic diagrams are constructed in a digital fashion (presence or absence of a factor) and the phenotypes do not contain temporal information that could inform us about these dynamics.

1184

### The requirements for epistatic analysis from a statistical mechanical perspective

1186

#### **Null alleles.**

1187

Geneticists have known for a long time that epistatic analyses must be carried out using verified null alleles (alleles that do not generate product, or which product is devoid of all biochemical activity). The reason for this requirement becomes abundantly clear from Eq. 2.6. Since the genetic activity of  $a$  is modulated by  $b$ , failure to use a null allele of  $a$  means that the loss of the activity provided by  $b$  in the double mutant will result in a measurable decrease of the activity of the gene product,  $A$ . Thus, the  $F_{reg}$  factor in the double mutant will not be the same as the  $F_{reg}$  factor in the mutant of  $a$ , and the equality requirement for epistasis will be violated. Notably, a null allele of  $b$  is not essential, strictly speaking, to detect epistasis. However, failure to obtain a null allele of  $b$  will inhibit our ability to test whether the pathway via which these genes interact is branched or not (see below).

1198

**Sensitization.**

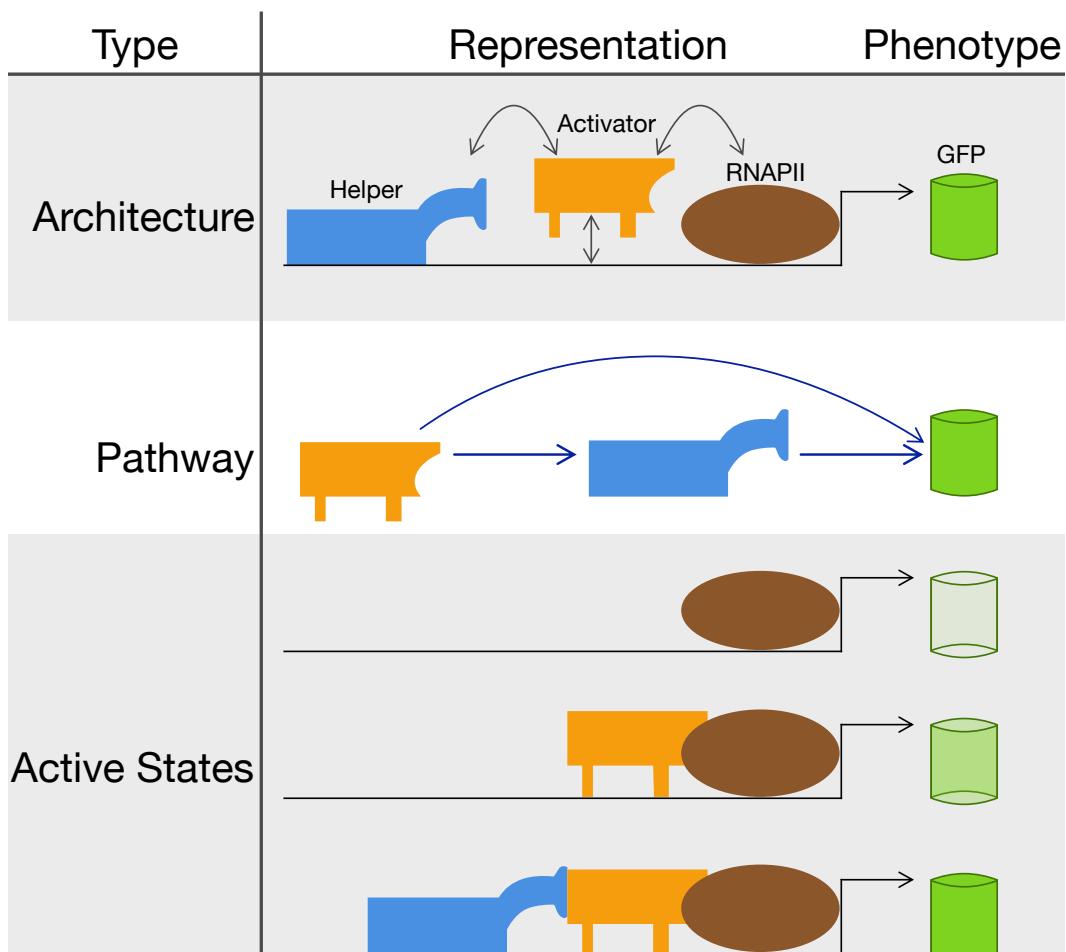
1199

Another difficulty with genetic analysis of a pathway can occur if the genes under 1200 study are not major determinants of the phenotypic outcome. This is equivalent 1201 to stating that  $F_{reg}$  depends on  $A$ ,  $B$  and at least one other factor,  $C$ ; and that 1202  $p_{bound}(A, B, C) \sim p_{bound}(0, 0, C)$ . The only case when Eq. 2.1 suggests this can 1203 happen (in the absence of measurement error) is when the effective polymerase 1204 activity,  $\rho \cdot F_{reg}(0, 0, C)$  is very large. To study the functions of  $a$  and  $b$  productively, 1205 it is necessary to re-scale  $F_{reg}$  to an appropriate dynamic range such that  $p_{bound}$  is 1206 elastic to changes in the gene activities of  $a$  and  $b$ . Since  $F_{reg}$  depends on the gene 1207 activity of  $c$ , the correct approach is to use an allele of  $c$  that decreases  $F_{reg}$  towards 1208 the desired direction. It is imperative that a null allele of  $c$  **not** be used, since  $c$  may 1209 be epistatic to  $a$ ,  $b$  or both, in which case we would be unable to study interactions 1210 between  $a$  and  $b$  in this genetic background. An alternative to this strategy would 1211 be to decrease  $\rho$ , but since RNA polymerase is a crucial component of all cells, this 1212 is not advisable. 1213

**Genetic pathways in statistical mechanics**

1214

Genetic diagrams are representations that satisfy observed epistasis relationships 1215 between genes. Since we have found a mapping between epistasis relationships 1216 and statistical mechanics, we can find the genetic pathway equivalent to the states 1217 oriented picture of promoter architectures from statistical mechanics (see Fig. 22). 1218 As a representative example, envision a promoter architecture where there is an 1219 activator,  $A$ , which recruits RNA polymerase to the promoter and which can bind 1220 DNA with a specific affinity; a helper or tethering protein,  $H$ , that can bind the 1221 activator as well as DNA, but cannot bind RNA polymerase; and RNA polymerase 1222 present at some concentration. In this system, the helper molecule is not necessary 1223 for activator function. The  $F_{reg}$  function for this system takes the following analytical 1224



**Figure 22** Promoter architectures can be represented in three different ways. A promoter architecture can be represented by showing all the interactions between components simultaneously. An alternative representation is via genetic diagrams. The genetic diagram obviates the epistatic element. Finally, the statistical mechanical view is an enumeration of all the possible states in the promoter. In the above figure, only the active states are shown for brevity. A state is considered active if RNAPII is bound to the promoter.

form:

1225

$$F_{reg}(A, H) = \frac{1 + Ae^{-\varepsilon_{AP} \frac{1+He^{-\varepsilon_{AH}}}{1+H}}}{1 + A \frac{1+He^{-\varepsilon_{AH}}}{1+H}}. \quad (2.8)$$

Here,  $A$  and  $H$  represent the amount of activator and helper;  $\varepsilon_{AH}$  is the energy 1226  
of binding of activator to helper; and  $\varepsilon_{AP}$  is the energy of binding of activator to 1227  
polymerase. For simplicity, energies are expressed in units of  $k_B T$ . 1228

As a first step, we confirm that this function satisfies the nested function requirement 1229  
for epistasis. We could re-write the above equation into a nested function as follows: 1230

$$F_{reg}(A \cdot G_{reg}(H)) = \frac{1 + Ae^{-\varepsilon_{AP}} \cdot G_{reg}(H)}{1 + A \cdot G_{reg}(H)}. \quad (2.9)$$

The fact that we can write this nested function means that the activator gene is 1231  
epistatic over the helper gene. Next, we will order the mutants according to the 1232  
magnitude of their  $F_{reg}$  factor. Since this is a system of activators, the wild-type 1233  
 $F_{reg}(A, B)$  must be greater than the value in any of the mutants. Since  $a$  is epistatic 1234  
over  $b$ , we know that  $F_{reg}$  for the single mutant of  $a$  and the double mutant have to 1235  
be the same. We also know that  $F_{reg}$  always increases with increasing effective gene 1236  
activity of  $a$ ,  $A \cdot G_{reg}(H)$ . Finally, from inspection of Eq. 2.8,  $G_{reg}(0) = 1$ . Putting 1237  
all of this together, we can order the factors: 1238

$$F_{reg}[0] \leq F_{reg}[A \cdot G_{reg}(0)] \leq F_{reg}[A \cdot G(H)]. \quad (2.10)$$

This means that the phenotype of the double mutant  $ab$  and  $a$  are of equal severity, 1239  
and show the greatest perturbation relative to the wild-type, while the single mutant 1240  
of  $h$  has a phenotype intermediate to the phenotype of mutants of  $a$  and the wild-type 1241  
control. Since the epistasis coefficient depends only on the  $F_{reg}$  associated with each 1242  
mutant, it follows that if we study how flexible this hierarchy can be, we can study 1243

how this hierarchy can be modified by tuning the parameters in Eq. 2.8, and how 1244  
this in turn can affect the value of the epistasis coefficient. 1245

We would like to know when  $F_{reg}(A \cdot G_{reg}(0)) \sim F_{reg}[A \cdot G(H)]$ . This is equivalent 1246  
to stating that loss of  $h$  does not measurably affect gene expression. This will happen 1247  
in the limit of saturating gene activity of  $a$ , which results from the combination of 1248  
protein copy number, DNA binding affinity, and affinity for RNA polymerase. An 1249  
epistatic analysis will conclude that  $a$  and  $h$  do not interact along this phenotype and 1250  
would identify a single pathway, involving the activator gene,  $a$ . 1251

Another relevant limit is  $F_{reg}(0) \sim F_{reg}(A \cdot G_{reg}(0))$ . For this to be true, two 1252  
conditions must hold. First, the gene activity of  $a$  must be low. Either the protein 1253  
product is present at very low copy number or the DNA binding affinity of the 1254  
product is very low by itself. Second, the gene activity of  $h$  must not be low: It must 1255  
have a reasonable combination of protein copy number, DNA binding affinity and 1256  
binding affinity to the product of  $a$  such that its regulatory function,  $G_{reg}(H)$  has a 1257  
value much greater than unity. In this regime, the single mutants of  $a$  and  $h$ , as well 1258  
as the double mutant, would have the same fold-change relative to the wild-type. 1259  
This reflects the fact that the only complex making a measurable contribution to 1260  
polymerase binding is the complex where the activator is bound to the helper. In 1261  
other words, these two genes act within a single, unbranched pathway that involves 1262  
both  $a$  and  $h$ . 1263

In between these two regimes, epistasis indicates the existence of two pathways, one 1264  
of which involves both  $a$  and  $h$ , and one of which involves only  $a$ . These pathways 1265  
reflect the importance of an activator-polymerase complex that does not include the 1266  
helper protein; genetically, this secondary complex constitutes a secondary, helper- 1267  
independent pathway. The relative importance of each pathway will depend strongly 1268  
on the wild-type dosage levels. This shows that statistical mechanical microstates 1269

can be mapped to genetic pathways involving multiple genes. This mapping is not 1270  
guaranteed to be one-to-one (the mappings will usually be many-to-one), nor is it 1271  
guaranteed to describe all the states available to the pathway. 1272

Although at no point is the equality in Eq. 2.5 violated, the interpretation of the 1273  
epistasis relationship is highly dependent on gene dosage. Epistasis relationships are 1274  
not immutable with dosage, and a given relationship may change from independent 1275  
to additive to suppressive along a given dosage curve. This highlights the importance 1276  
of varying the reference gene activity levels using partial loss-of-function alleles to 1277  
establish that Eq. 2.5 holds along an entire curve. 1278

## Genetic Morphs

Genetics can generate alleles with vastly different properties. Our formalism allows 1280  
us to immediately derive the most commonly encountered allele classes (see Sup- 1281  
plementary Information Section XXX: Common allelic classes and their statistical 1282  
mechanical definitions). In the following text, we assume the gene in question 1283  
promotes transcription; for an inhibitor, the definitions are the same except that the 1284  
effects on active and inactive states are flipped. An active state is a state that has 1285  
RNA polymerase bound. 1286

- *Hypomorph.* **Genetic definition:** an allele with reduced gene activity, either 1287  
by reduced product copy number or decreased biochemical activity, causing 1288  
a loss-of-function mutant phenotype. **Statistical mechanical definition:** an 1289  
allele which decreases the relative proportion of active states compared with 1290  
the wild type homozygote. 1291
- *Hypermorph.* **Genetic definition:** an allele with increased gene activity, 1292  
either by increased copy number or improved biochemical activity, causing 1293  
an increased-function mutant phenotype; the hypermorphic phenotype can be 1294

phenocopied by overexpression of the wild type allele. **Statistical mechanical definition:** an allele that increases the relative proportion of active states compared with the wild type homozygote.

- *Neomorph.* **Genetic definition:** an allele which has a novel functionality causing a mutant phenotype which cannot be phenocopied simply by over-expressing wild type product. The neomorphic allele generates product at a similar rate as the wild type allele. **Statistical mechanical definition:** an allele encoding a modified  $F_{reg}$  function, generating novel states (active or inactive) not accessible to the wild type product at any concentration.

Grouping alleles into the correct change-of-function classes is an important part of genetic analysis because hypermorphic alleles from one gene can be used in combination with null alleles at a second locus to order genes in a pathway. However, neomorphic and hypermorphic alleles can be extremely difficult to discern and can confound genetic analysis. From a statistical mechanical perspective, this confounding arises because the neomorphic allele changes the functional form of the interactions, which qualitatively alters the system by adding neofunctionalized states, whereas the hypermorphic allele exclusively changes the gene dosage relative to the wild-type levels.

## Discussion

We have shown that genetics is a method to study partition functions that are not known analytically. The method relies on systematically setting component values to zero, individually and pairwise, and identifying circumstances when the partition function exhibits epistasis. Although equality between two sets of perturbations can occur as a coincidence, the stringent requirement of equality makes this an unlikely occurrence. Therefore, if there is epistasis between two components, it is reasonable to infer that these perturbed variables function as a single variable equal to a product

of two functions that each take as input one of the two component. This methodology 1321  
is known in genetics as epistatic analysis, and it provides a way to study, in an analog 1322  
fashion, a molecular system where the components and interactions may be entirely 1323  
unknown. 1324

The results from epistatic analyses can be represented as genetic diagrams, where 1325  
the genes represent a component in the system and arrows indicate interactions. 1326  
Because of the analog nature of the analysis (components are either ‘ON’ or ‘OFF’), 1327  
genetic diagrams do not have to yield information about the dynamical capacities 1328  
of the system, nor do they represent direct biochemical interactions. Rather, genetic 1329  
pathways are a graphical method of representing epistatic interactions between 1330  
components in a circuit, allowing geneticists to rapidly identify the epistatic and 1331  
hypostatic components in any given comparison, and the weighting given to different 1332  
arrows depends strongly on the dosage of all the components connected by them. 1333  
When these pathways are drawn to represent elements of a promoter, pathways can 1334  
be mapped onto subsets of active states involving the elements connected by arrows. 1335

Gene expression levels are increasingly being used as a phenotype for genetic anal- 1336  
ysis. A particularly attractive feature of these phenotypes is our ability to measure 1337  
expression levels in a highly multiplexed fashion in a single experiment, but the 1338  
resulting datasets have proven hard to analyze due to their high-dimensional na- 1339  
ture. Dimensionality reduction methods such as principal component analysis or 1340  
non-negative matrix factorization have become popular ways to analyze transcrip- 1341  
tomic data. However, a drawback of these methods is that they are derived from a 1342  
statistical framework without connection to the underlying biology, and the trans- 1343  
formations that they apply to the data make it difficult to interpret signals in terms of 1344  
biological components. Another powerful approach is to fit general linear models 1345  
with interactions terms to transcriptomic dataset. In this framework, the expression 1346  
level of an individual transcript is modeled with two first order coefficients that 1347

quantify the independent effect of each component on the expression level of that transcript, and an interaction term that is non-zero when the double mutant shows non-additive expression levels. Here, we argue that this approach does not extract all the available information encoded in transcriptomes because it does not, on its own, identify cases where classical epistasis is occurring.

1352

We have shown that, for systems at equilibrium, epistasis arises naturally. Moreover, we have shown that epistasis percolates through a network at equilibrium regardless of its depth, and we speculate whether this result may explain the unreasonable effectiveness of genetics at enormous differences in scale, since macroscopic (organismal or population) phenotypes to microscopic phenotypes all show epistasis. However, a major caveat to this work is the fact that biological systems are not at thermal equilibrium. The commonality of classical epistasis in biological systems suggests that many of our results will have a non-equilibrium analog with some changes. Certainly, networks that are not at equilibrium will not reveal epistasis at arbitrary depth. For non-equilibrium networks, it will be particularly interesting to study how epistasis is affected by relaxing the acyclic requirement we imposed. Feedback might be expected to hide epistasis; however, biological experiments suggest that certain topologies, even with feedback, show epistasis robustly. The study of the physical basis of epistasis in biological systems seems poised to reveal fascinating new insights into the principles of how these networks are built.

1367

## References

1368

- Angeles-Albores, David, Daniel H. W. Leighton, et al. (2017). “The *Caenorhabditis elegans* Female-Like State: Decoupling the Transcriptomic Effects of Aging and Sperm Status”. In: *G3: Genes, Genomes, Genetics* 7.9.
- Angeles-Albores, David, Carmie Puckett Robinson, et al. (2018). “Reconstructing a metazoan genetic pathway with transcriptome-wide epistasis measurements”. In: *Proceedings of the National Academy of Sciences*, p. 201712387. ISSN: 0027-8424. doi: [10.1073/pnas.1712387115](https://doi.org/10.1073/pnas.1712387115).

1375

- Bateson, William and Gregor Mendel (2009). *Mendel's principles of heredity: A defence, with a translation of mendel's original papers on hybridisation*, pp. 1– 1376  
212. ISBN: 9780511694462. doi: [10.1017/CBO9780511694462](https://doi.org/10.1017/CBO9780511694462). 1377  
1378
- Bintu, Lacramioara et al. (Apr. 2005). “Transcriptional regulation by the numbers: Models”. In: *Current Opinion in Genetics and Development* 15.2, pp. 116– 1379  
124. ISSN: 0959437X. doi: [10.1016/j.gde.2005.02.007](https://doi.org/10.1016/j.gde.2005.02.007). arXiv: [0412011](https://arxiv.org/abs/0412011) 1380  
[q-bio]. 1381  
1382
- Bothma, Jacques P et al. (Aug. 2015). “Enhancer additivity and non-additivity are determined by enhancer strength in the *Drosophila* embryo”. In: *eLife* 4. ISSN: 1383  
2050-084X. doi: [10.7554/eLife.07956](https://doi.org/10.7554/eLife.07956). 1384  
1385
- Cordell, Heather J (2002). “Epistasis: what it means, what it doesn't mean, and statistical methods to detect it in humans”. In: *Human Molecular Genetics* 11.20, 1386  
pp. 2463–2468. doi: [10.1093/hmg/11.20.2463](https://doi.org/10.1093/hmg/11.20.2463). 1387  
1388
- Fisher, R. A. (1919). “XV.—The Correlation between Relatives on the Supposition of Mendelian Inheritance”. In: *Transactions of the Royal Society of Edinburgh* 1389  
52.2, pp. 399–433. ISSN: 00804568. doi: [10.1017/S0080456800012163](https://doi.org/10.1017/S0080456800012163). 1390  
1391
- Garcia, Hernan G. et al. (2007). “A First Exposure to Statistical Mechanics for Life Scientists”. In: p. 27. ISSN: 0036-8075. arXiv: [0708.1899](https://arxiv.org/abs/0708.1899). 1392  
1393
- Jürgens, G. et al. (1984). “Mutations affecting the pattern of the larval cuticle in *Drosophila melanogaster* - II. Zygotic loci on the third chromosome”. In: *Wilhelm Roux's Archives of Developmental Biology* 193.5, pp. 283–295. ISSN: 03400794. 1394  
DOI: [10.1007/BF00848157](https://doi.org/10.1007/BF00848157). 1395  
1396  
1397
- Lee, Edward D., Chase P. Broedersz, and William Bialek (July 2015). “Statistical Mechanics of the US Supreme Court”. In: *Journal of Statistical Physics* 160.2, 1398  
pp. 275–301. ISSN: 00224715. doi: [10.1007/s10955-015-1253-6](https://doi.org/10.1007/s10955-015-1253-6). arXiv: [1306.5004](https://arxiv.org/abs/1306.5004). 1399  
1400  
1401
- Phillips, Patrick C (2008). “Epistasis—the essential role of gene interactions in the structure and evolution of genetic systems”. In: *Nat Rev Genet* 9.11, pp. 855–867. 1402  
ISSN: 1471-0056. doi: [10.1038/nrg2452](https://doi.org/10.1038/nrg2452). 1403  
1404
- Raveh-Sadka, Tali, Michal Levo, and Eran Segal (Aug. 2009). “Incorporating nucleosomes into thermodynamic models of transcription regulation.” In: *Genome research* 19.8, pp. 1480–96. ISSN: 1088-9051. doi: [10.1101/gr.088260.108](https://doi.org/10.1101/gr.088260.108). 1405  
1406  
1407
- Schneidman, Elad et al. (Apr. 2006). “Weak pairwise correlations imply strongly correlated network states in a neural population”. In: *Nature* 440.7087, pp. 1007– 1408  
1012. ISSN: 0028-0836. doi: [10.1038/nature04701](https://doi.org/10.1038/nature04701). 1409  
1410
- Sulston, J. E. and S. Brenner (1974). “The DNA of *Caenorhabditis elegans*.” In: *Genetics* 77.1, pp. 95–104. ISSN: 00166731. 1411  
1412

RECONSTRUCTING A METAZOAN GENETIC PATHWAY 1414  
WITH TRANSCRIPTOME-WIDE EPISTASIS MEASUREMENTS 1415

**Abstract** 1416

**RNA-seq is commonly used to identify genetic modules that respond to pertur-** 1417  
**bations. In single cells, transcriptomes have been used as phenotypes, but this** 1418  
**concept has not been applied to whole-organism RNA-seq. Also, quantifying** 1419  
**and interpreting epistatic effects using expression profiles remains a challenge.** 1420  
**We developed a single coefficient to quantify transcriptome-wide epistasis that** 1421  
**reflects the underlying interactions and which can be interpreted intuitively.** 1422  
**To demonstrate our approach, we sequenced four single and two double mu-** 1423  
**tants of *Caenorhabditis elegans*. From these mutants, we reconstructed the** 1424  
**known hypoxia pathway. In addition, we uncovered a class of 56 genes with** 1425  
***hif-1*-dependent expression that have opposite changes in expression in mu-** 1426  
**tants of two genes which cooperate to negatively regulate HIF-1 abundance;** 1427  
**however, the double mutant of these genes exhibits suppression epistasis. This** 1428  
**class violates the classical model of HIF-1 regulation, but can be explained by** 1429  
**postulating a role of hydroxylated HIF-1 in transcriptional control.** 1430

**Introduction** 1431

Genetic analysis of molecular pathways has traditionally been performed through 1432  
epistatic analysis. If the mutants of two distinct genes have a quantifiable phenotype, 1433  
and the double mutant has a phenotype that is not the sum of the phenotypes of 1434  
the single mutants, this non-additivity is referred to as generalized epistasis, and 1435

indicates that these genes interact functionally. Such interactions can occur at the 1436 biochemical level between their products or as a consequence of their functions (L. S. 1437 Huang and Paul W Sternberg, 2006). Epistasis analysis remains a cornerstone of 1438 genetics today (Phillips, 2008). 1439

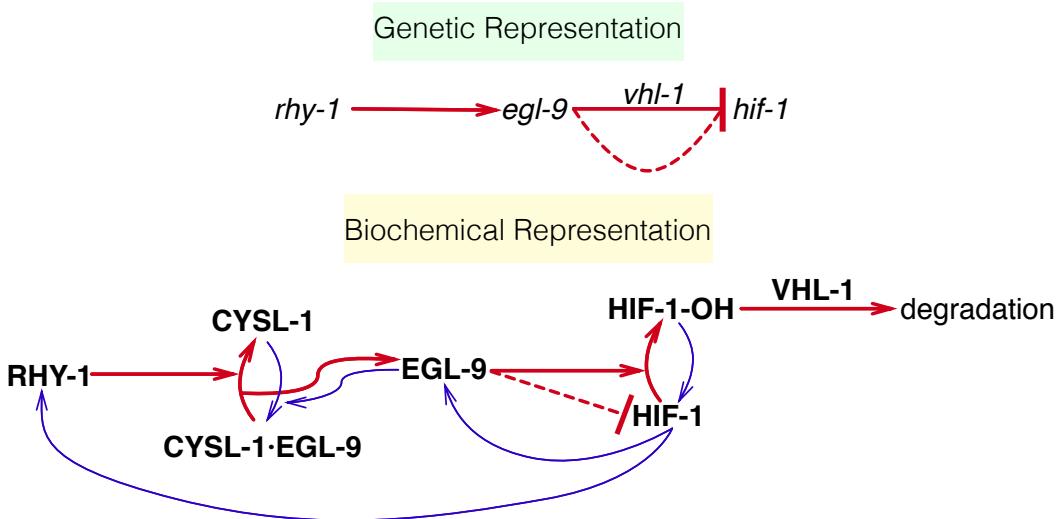
Recently, biological studies have shifted in focus from studying single genes to study- 1440 ing all genes in parallel. In particular, RNA-seq (Mortazavi et al., 2008) enables 1441 biologists to identify genes that change expression in response to a perturbation. 1442 RNA-seq has been used to identify genetic modules involved in a variety of pro- 1443 cesses, such as in the *Caenorhabditis elegans* linker cell migration (Schwarz, Kato, 1444 and Paul W. Sternberg, 2012), planarian stem cell maintenance (Van Wolfswinkel, 1445 Wagner, and Reddien, 2014; Scimone et al., 2014). The role of transcriptional 1446 profiling has been restricted to target gene identification, and so far there are only a 1447 few examples where transcriptomes have been used to generate quantitative genetic 1448 models of any kind. In quantitative genetics, eQTL studies have established the 1449 power of transcriptomes for genetic mapping (Brem et al., 2002; Schadt et al., 2003; 1450 Li et al., 2006; King et al., 2014). Genetic pathway analysis via epistasis has been 1451 performed in *Saccharomyces cerevisiae* (Hughes et al., 2000; Capaldi et al., 2008) 1452 and in *Dictyostelium discoideum* (Van Driessche et al., 2005). Recently, Dixit *et al* 1453 described a protocol for epistasis analysis in dendritic and K562 cells using single- 1454 cell RNA-seq (Dixit et al., 2016). Epistasis analysis of single cells or single-celled 1455 organisms is popular because of the concern that whole-organism sequencing will 1456 mix information from multiple cell types, preventing the accurate reconstruction 1457 of genetic interactions. Using whole-organism transcriptome profiling, we have 1458 recently identified a new developmental state of *C. elegans* caused by loss of a 1459 single cell type (sperm cells) (Angeles-Albores, Leighton, et al., 2017), which sug- 1460 gests that whole-organism transcriptome profiling contains sufficient information 1461 for epistatic analysis. To investigate the ability of whole-organism transcriptomes to 1462

serve as quantitative phenotypes for epistatic analysis in metazoans, we sequenced 1463  
the transcriptomes of four well-characterized loss-of-function mutants in the *C. ele- 1464  
gans* hypoxia pathway (Epstein et al., 2001; Shen, Shao, and Powell-Coffman, 2006; 1465  
Shao, Zhang, and Powell-Coffman, 2009; H. Jiang, Guo, and Powell-Coffman, 1466  
2001). 1467

Metazoans depend on the presence of oxygen in sufficient concentrations to support 1468  
aerobic metabolism. Hypoxia inducible factors (HIFs) are an important group 1469  
of oxygen-responsive genes that are highly conserved in metazoans (Loenarz et al., 1470  
2011). A common mechanism for hypoxia-response induction is heterodimerization 1471  
between a HIF $\alpha$  and a HIF $\beta$  subunit; the heterodimer then initiates transcription 1472  
of target genes (B. H. Jiang et al., 1996). The number and complexity of HIFs 1473  
varies throughout metazoans. In the roundworm *C. elegans* there is a single HIF $\alpha$  1474  
gene, *hif-1* (H. Jiang, Guo, and Powell-Coffman, 2001), and a single HIF $\beta$  gene, 1475  
*ahr-1* (Powell-Coffman, Bradfield, and Wood, 1998). 1476

Levels of HIF $\alpha$  proteins are tightly regulated. Under conditions of normoxia, HIF- 1477  
1 $\alpha$  exists in the cytoplasm and partakes in a futile cycle of protein production 1478  
and rapid degradation (L. E. Huang et al., 1996). In *C. elegans*, HIF-1 $\alpha$  is hy- 1479  
droxylated by a proline hydroxylase (EGL-9) (Kaelin and Ratcliffe, 2008). HIF-1 1480  
hydroxylation increases its binding affinity to Von Hippel-Lindau tumor suppressor 1481  
1 (VHL-1), which in turn allows ubiquitination of HIF-1 leading to its degradation. 1482  
In *C. elegans*, EGL-9 activity is inhibited by binding of CYSL-1, a homolog of 1483  
sulphydrylases/cysteine synthases; and CYSL-1 activity is in turn inhibited by the 1484  
putative transmembrane O-acyltransferase RHY-1, possibly by post-translational 1485  
modifications to CYSL-1 (Ma et al., 2012) (see Fig. 31). 1486

Our reconstruction of the hypoxia pathway in *C. elegans* shows that whole-animal 1487  
transcriptome profiles can be used as phenotypes for genetic analysis and that epis- 1488



**Figure 31** Genetic and biochemical representation of the hypoxia pathway in *C. elegans*. Red arrows are arrows that lead to inhibition of HIF-1, and blue arrows are arrows that increase HIF-1 activity or are the result of HIF-1 activity. EGL-9 is known to exert VHL-1-dependent and independent repression on HIF-1 by EGL-9 is denoted by a dashed line and is not dependent on the hydroxylating activity of EGL-9. RHY-1 inhibits CYSL-1, which in turn inhibits EGL-9, but this interaction was abbreviated in the genetic diagram for clarity.

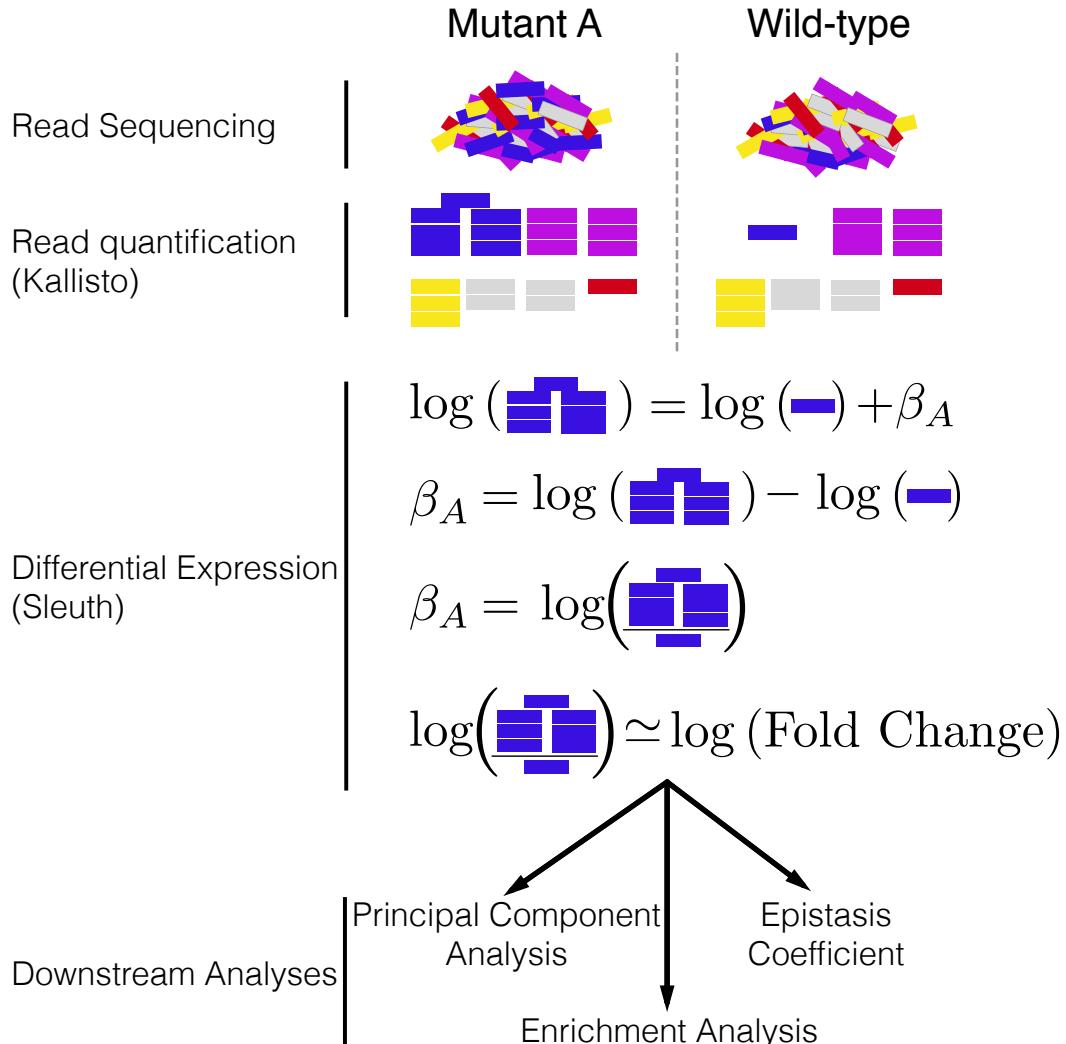
tasis, a hallmark of genetic interaction observed in double mutants, holds at the 1489 molecular systems level. We demonstrate that transcriptomes can aid in ordering 1490 genes in a pathway using only single mutants. We were able to identify genes that 1491 appear to be downstream of *vhl-1*, but not downstream of *hif-1*. Using a single set of 1492 transcriptome-wide measurements, we observed most of the known transcriptional 1493 effects of *hif-1* as well as novel effects not described before in *C. elegans*. Taken 1494 together, this analysis demonstrates that whole-animal RNA-seq is a fast and power- 1495 ful method for genetic analyses in an area where phenotypic measurements are now 1496 the rate-limiting step. 1497

**Results** 1498

**The hypoxia pathway controls thousands of genes in *C. elegans*** 1499

We selected four null single mutants within the hypoxia pathway for expression pro- 1500 filing: *egl-9(sa307)*, *rhy-1(ok1402)*, *vhl-1(ok161)*, *hif-1(ia4)*. We also sequenced 1501 the transcriptomes of two double mutants, *egl-9; vhl-1* and *egl-9 hif-1* as well as 1502 wild type (N2). Each genotype was sequenced in triplicate using mRNA extracted 1503 from 30 worms at a depth of 15 million reads per sample. Of these 15 million 1504 reads, 50% of the reads mapped to the *C. elegans* genome on average. All samples 1505 were analyzed under normoxic conditions. We measured differential expression of 1506 19,676 isoforms across all replicates and genotypes (~70% of the protein coding 1507 isoforms in *C. elegans*; see [Basic Statistics Notebook](#)). We included in our analysis 1508 a *fog-2(q71)* mutant we have previously studied (Angeles-Albores, Leighton, et al., 1509 2017), because *fog-2* is not reported to interact with the hypoxia pathway. We 1510 analyzed our data using a general linear model on logarithm-transformed counts. 1511 Changes in gene expression are reflected in the regression coefficient  $\beta$ , which is 1512 specific to each isoform within a genotype (excluding wild type, which is used 1513 as baseline). Statistical significance is achieved when the q-value of a  $\beta$  coeffi- 1514 cient (*p*-values adjusted for multiple testing) are less than 0.1. Transcripts that are 1515 differentially expressed between the wild type and a given mutant have  $\beta$  values 1516 that are statistically significantly different from 0 (i.e. greater than 0 or less than 1517 0).  $\beta$  coefficients are analogous to the logarithm of the fold-change between the 1518 mutant and the wild type. Larger magnitudes of  $\beta$  correspond to larger perturba- 1519 tions (see Fig. 32). When we refer to  $\beta$  coefficients and *q*-values, it will always 1520 be in reference to isoforms. However, we report the sizes of each gene set by 1521 the number of differentially expressed genes (DEGs), not isoforms, they contain. 1522 For the case of *C. elegans*, this difference is negligible since the great majority 1523 of protein-coding genes have a single isoform. We have opted for this method of 1524

referring to gene sets because it simplifies the language considerably. A complete 1525 version of the code used for this analysis with ample documentation, is available at 1526 <https://wormlabcaltech.github.io/mprsq>. 1527



**Figure 32** Analysis workflow. After sequencing, reads are quantified using Kallisto. Bars show estimated counts for each isoform. Differential expression is calculated using Sleuth, which outputs one  $\beta$  coefficient per isoform per genotype.  $\beta$  coefficients are analogous to the natural logarithm of the fold-change relative to a wild type control. Downstream analyses are performed with  $\beta$  coefficients that are statistically significantly different from 0.  $q$ -values less than 0.1 are considered statistically different from 0.

Transcriptome profiling of the hypoxia pathway revealed that this pathway controls 1528 thousands of genes in *C. elegans* (see Table 31, see SI File 1 for a complete list 1529

of differentially expressed genes). The *egl-9(lf)* transcriptome showed differential expression of 2,549 genes. 3,005 genes were differentially expressed in *rhy-1(lf)* mutants. The *vhl-1(lf)* transcriptome showed considerably fewer DEGs (1,275), possibly because *vhl-1* is a weaker inhibitor of *hif-1* than *egl-9* (Shao, Zhang, and Powell-Coffman, 2009). The *egl-9(lf);vhl-1(lf)* double mutant transcriptome showed 3,654 DEGs. The *hif-1(lf)* mutant showed a transcriptomic phenotype involving 1,075 genes. The *egl-9(lf) hif-1(lf)* double mutant showed a similar number of genes with altered expression (744 genes). We do not think that this transcriptional response is due to transiently induced hypoxia during harvesting. If the wild type strain had become hypoxic, then the *hif-1(lf)* genotype should show significantly lower levels of *nhr-57*, a marker that increases during hypoxia. We do not observe altered levels of *nhr-57* when comparing the wild type and *hif-1(lf)* mutant, nor between the wild type and *egl-9(lf) hif-1(lf)* double mutant. Finally, the *egl-9(lf)*, *vhl-1(lf)*, *rhy-1(lf)* and *egl-9(lf); vhl-1(lf)* mutants did show altered *nhr-57* transcript levels (see Quality Control Notebook, SI Figure 1). Of the differentially expressed genes in *hif-1(lf)* mutants, 161/1,075 were also differentially expressed in *egl-9(lf) hif-1(lf)* mutants, which suggests these transcripts are *hif-1*-dependent under normoxia. For the remaining genes, we cannot rule out cumulative effects from loss of *hif-1*, strain-specific eQTLs present in the strain background or that loss of *egl-9* suppresses the mutant phenotype. We designed our experiments to probe the constitutive hypoxia response, and not the effects of *hif-1* under normoxia, which we did not foresee. As a result, we have limited resolving power to explain the transcriptome of *hif-1(lf)* mutants.

1552

Genotype	Differentially Expressed Genes
<i>egl-9(lf)</i>	2,549
<i>rhy-1(lf)</i>	3,005
<i>vhl-1(lf)</i>	1,275
<i>hif-1(lf)</i>	1,075
<i>egl-9(lf); vhl-1(lf)</i>	3,654
<i>egl-9(lf) hif-1(lf)</i>	744
<i>fog-2(lf)</i>	2,840

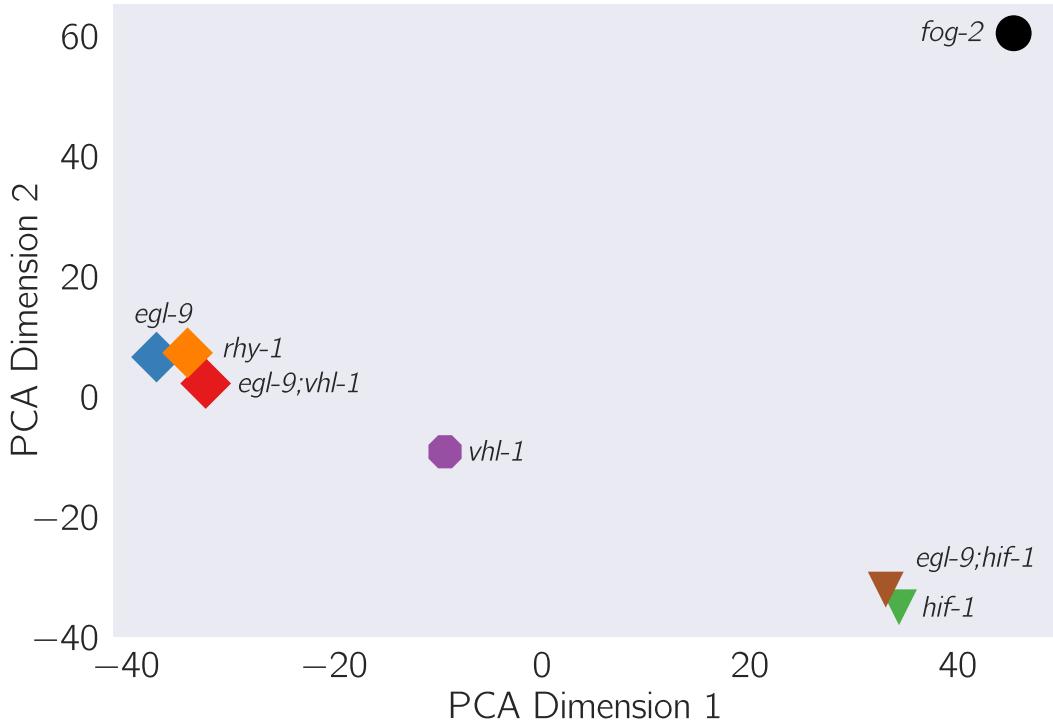
**Table 31** Number of differentially expressed genes in each mutant strain with respect to the wild type (N2).

### Principal Component Analysis visualizes epistatic relationships between geno- 1553 types 1554

Principal component analysis (PCA) is used to identify relationships between high- 1555 dimensional data points (Yeung and Ruzzo, 2001). We used PCA examine whether 1556 each genotype clustered in a biologically relevant manner. PCA identifies the vec- 1557 tor that explains most of the variation in the data; this is called the first principal 1558 component. PCA can identify the first  $n$  components that explain more than 95% of 1559 the data variance. Clustering in these  $n$  dimensions can indicate biological relation- 1560 ships, although interpreting principal components can be difficult. In our analysis, 1561 the first principal component discriminated mutants that have constitutive high lev- 1562 els of HIF-1 from mutants that have no HIF-1, whereas the second component was 1563 able to discriminate between mutants within the hypoxia pathway and outside the 1564 hypoxia pathway (see Fig. 33; *fog-2* is not reported to act in the hypoxia pathway 1565 and acts as a negative control; see [Genetic Interactions Notebook](#)). 1566

### Reconstruction of the hypoxia pathway from first genetic principles 1567

To reconstruct a genetic pathway, we must assess whether two genes act on the same 1568 phenotype. If they do not act on the same phenotype (two mutations do not cause the 1569 same genes to become differentially expressed relative to wild type), these mutants 1570



**Figure 33** Principal component analysis of various *C. elegans* mutants. Genotypes that have an constitutive hypoxia response (i.e. *egl-9(lf)*) cluster far from genotypes that do not have a hypoxic response (i.e. *hif-1(lf)*) along the first principal component. The second principal component separates genotypes that do not participate hypoxic response pathway.

are independent. Otherwise, we must measure whether these genes act additively 1571 or epistatically on the phenotype of interest; if there is epistasis we must measure 1572 whether it is positive or negative, in order to assess whether the epistatic relationship 1573 is a genetic suppression or a synthetic interaction. To allow coherent comparisons 1574 of different mutant transcriptomes (the phenotype we are studying here), we define 1575 the shared transcriptomic phenotype (STP) between two mutants as the shared set 1576 of genes or isoforms whose expression in both mutants are different from wild-type, 1577 regardless of the direction of change. 1578

**Genes in the hypoxia mutant act on the same transcriptional phenotype**

1579

All the hypoxia mutants had a significant STP: the fraction of differentially expressed genes that was shared between mutants ranged from a minimum of 10% between *hif-1(lf)* and *egl-9(lf)*; *vhl-1(lf)* to a maximum of 32% between *egl-9(lf)* and *egl-9(lf); vhl-1(lf)* (see SI Table 1). For comparison, we also analyzed a previously published *fog-2(lf)* transcriptome (Angeles-Albores, Leighton, et al., 2017). The *fog-2* gene is involved in masculinization of the *C. elegans* germline, which enables sperm formation, and is not known to be involved in the hypoxia pathway. The hypoxia pathway mutants and the *fog-2(lf)* mutant also had STPs (8.8%–14%).

Next, we analyzed pairwise correlations between all mutant pairs. We rank-transformed the  $\beta$  coefficients of each isoform between the STP of two mutants, and plotted the transcript ranks between genotypes (see Fig 34). Although *hif-1* is known to be genetically repressed by *egl-9*, *rhy-1* and *vhl-1* (Epstein et al., 2001; Shen, Shao, and Powell-Coffman, 2006), all the correlations between mutants of these genes and *hif-1(lf)* were positive (see [Genetic Interactions Notebook](#)). We reasoned that this apparent contradiction could be due to either strain-specific effects in our N2 background (an artifactual signal) or that it could reflect a previously unrecognized aspect of HIF-1 biology. This motivated us to look for genes that exhibited verifiable extreme patterns of anomalous behavior and led us to propose a new model of the hypoxia pathway (see Identification of non-classical epistatic interactions).

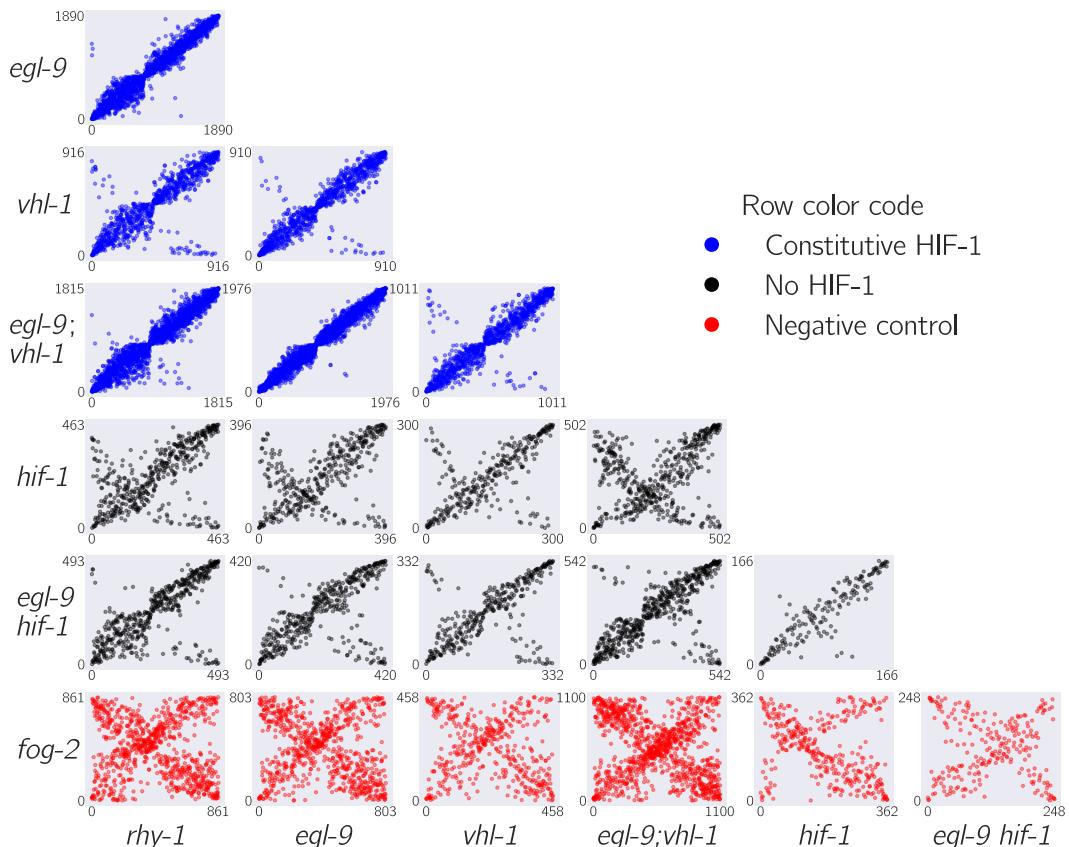
1599

**Transcriptome-wide epistasis**

1600

Ideally, any measurement of transcriptome-wide epistasis should conform to certain expectations. First, it should make use of the regression coefficients of as many genes as possible. Second, it should be summarizable in a single, well-defined number. Third, it should have an intuitive behavior, such that special values of the

1604



**Figure 34** Interacting genes have correlated transcriptional signatures. The rank order of transcripts contained in the shared transcriptional phenotype is plotted for each pairwise combination of genotypes. Correlations between in-pathway genotypes are strong whereas comparisons with a *fog-2(lf)* genotype are dominated by noise. Comparisons between some genotypes show populations of transcripts that are anticorrelated, possibly as a result of feedback loops. Plots are color-coded by row. Comparisons with genotypes with a constitutive hypoxia response are in blue; comparisons with genotypes negative for *hif-1(lf)* are black; and comparisons involving *fog-2(lf)* are red. X- and y-axes show the rank of each transcript within each genotype.

statistic have an unambiguous interpretation.

1605

We found an approach that satisfies all of the above conditions and which can be graphed in an epistasis plot (see Fig 35) In an epistasis plot, the X-axis represents the expected  $\beta$  coefficient for given gene in a double mutant  $a^-b^-$  if  $a$  and  $b$  interact log-additively. In other words, each individual isoform's x-coordinate is the sum of the regression coefficients from the single mutants  $a^-$  and  $b^-$ . The Y-axis represents the deviations from the log-additive (null) model, and can be calculated as the difference between the predicted and the observed  $\beta$  coefficients. Only isoforms that are differentially expressed in all three genotypes are plotted. This attempts to ensure that the isoforms to be examined are regulated by both genes. These plots will generate specific patterns that can be described through linear regressions. The slope of these lines, to which we assign the mathematical notation  $s(a, b)$ , is the transcriptome-wide epistasis coefficient. Importantly, the transcriptome-wide epistasis coefficient is fundamentally distinct from Pearson or Spearman correlation coefficients and need not have a simple linear mapping. In other words, negative correlation coefficients do not imply a specific sign of the epistasis coefficient, and vice versa. For suppression to occur, for example, the only requirement is that the phenotype of the double mutant should match one, and only one, of the two single mutants. The value of the correlation coefficient is not relevant.

1623

Transcriptome-wide epistasis coefficients can be understood intuitively for simple cases of genetic interactions if complete genetic nulls are used. If two genes act additively on the same set of differentially expressed isoforms then all the plotted points will fall along the line  $y = 0$ . If two genes act positively in an unbranched pathway, then all the mutants should have the same phenotype. It follows that data from this pathway will form line with slope equal to  $-\frac{1}{2}$ . On the other hand, in the limit of complete genetic inhibition of  $b$  by  $a$  in an unbranched pathway (i.e.,  $a$  is in great excess over  $b$ , such that under the conditions measured  $b$  has no activity),

1631

the plots should show a line of best fit with slope equal to  $-1$ . Genes that interact synthetically (*i.e.*, through an OR-gate) will fall along lines with slopes  $> 0$ . When there is epistasis of one gene over another, the points will fall along one of two possible slopes that must be determined empirically from the single mutant data. We can use both single mutant data to predict the distribution of slopes that results for the cases stated above. Thus, the transcriptome-wide epistasis coefficient integrates information from many different isoforms into a single number (see Fig. 35). 1638

In our experiment, we studied two double mutants, *egl-9(lf)* *hif-1(lf)* and *egl-9(lf); vhl-1(lf)*. We wanted to understand how well an epistatic analysis based on transcriptome-wide coefficients agreed with the epistasis results reported in the literature, which were based on qPCR of single genes. Therefore, we determined the epistasis coefficient of the two gene combinations we studied (*egl-9* and *vhl-1*, and *egl-9* and *hif-1*). In addition to computing an epistasis coefficient from these factors, we would like to know which gene is suppressed in the double mutant. Suppression means that the double mutant should have exactly the phenotype of one and only one mutant, we can simulate the double mutant by replacing the double mutant data with either of the two single mutants and matching the simulated result to the observed result. The result that most closely matches the real data will reveal which gene is being suppressed, which in turn allows us to order the genes along a pathway. 1650

We measured the epistasis coefficient between *egl-9* and *vhl-1*,  $s(\text{egl-9 } \text{vhl-1}) = -0.41 \pm 0.01$  (see [Epistasis Notebook](#)). Simulations using just the single mutant data showed that the double mutant exhibited the *egl-9(lf)* phenotype (see Fig. 35). We used Bayesian model selection to reject a linear pathway (odds ratio (OR)  $> 10^{92}$ ), which leads us to conclude *egl-9* is upstream of *vhl-1* acting on a phenotype in a branched manner. We also measured epistasis between *egl-9* and *hif-1*,  $s(\text{egl-9, hif-1}) = -0.80 \pm 0.01$  (see SI Figure 2), and we found that this behavior could be predicted by modeling *hif-1* downstream of *egl-9*. We also rejected the null

hypothesis that these two genes act in a positive linear pathway ( $OR > 10^{93}$ ). Taken 1659  
together, this leads us to conclude that *egl-9* strongly inhibits *hif-1*. 1660

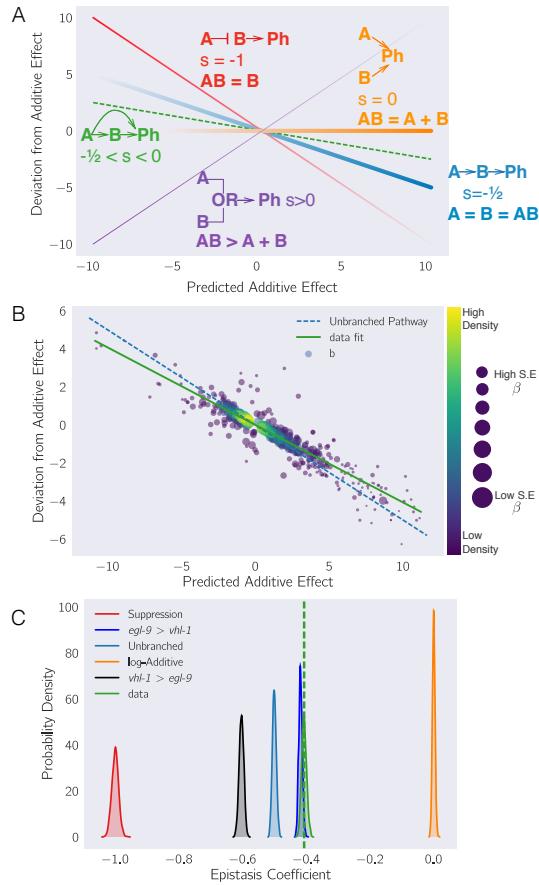
### **Epistasis between two genes can be predicted using an upstream component** 1661

Given our success in measuring epistasis coefficients, we wanted to know whether 1662  
it would be possible to predict the epistasis coefficient between *egl-9* and *vhl-1* 1663  
in the absence of the *egl-9(lf)* genotype. Since RHY-1 indirectly activates EGL- 1664  
*9*, we reasoned that the *rhy-1(lf)* transcriptome should contain almost equivalent 1665  
information to the *egl-9(lf)* transcriptome. Therefore, we generated predictions of 1666  
the epistasis coefficient between *egl-9* and *vhl-1* by substituting in the *rhy-1(lf)* data, 1667  
predicting  $s(rhy - 1, vhl - 1) = -0.45$ . Similarly, we used the *egl-9(lf); vhl-1(lf)* 1668  
double mutant to measure the epistasis coefficient while replacing the *egl-9(lf)* 1669  
dataset with the *rhy-1(lf)* dataset. We found that the epistasis coefficient using 1670  
this substitution was  $-0.38 \pm 0.01$ . This coefficient was different from  $-0.50$  (OR 1671  
 $> 10^{102}$ ), reflecting the same qualitative conclusion that *vhl-1* represents a branch 1672  
in the hypoxia pathway. We were able to obtain a close prediction of the epistasis 1673  
coefficient for two mutants using the transcriptome of a related, upstream mutant. 1674

### **Transcriptomic decorrelation can be used to infer functional distance** 1675

So far, we have shown that RNA-seq can accurately measure genetic interactions. 1676  
However, genetic interactions do not require two gene products to interact biochem- 1677  
ically, nor even to be physically close to each other. RNA-seq cannot measure 1678  
physical interactions between genes, but we wondered whether expression profiling 1679  
contains sufficient information to order genes along a pathway. 1680

Single genes are often regulated by multiple independent sources. The connection 1681  
between two nodes can in theory be characterized by the strength of the edges 1682  
connecting them (the thickness of the edge); the sources that regulate both nodes 1683



**Figure 35** (A) Schematic diagram of an epistasis plot. The X-axis on an epistasis plot is the expected coefficient for a double mutant under an log-additive model (null model). The Y-axis plots deviations from this model. Double mutants that deviate in a systematic manner from the null model exhibit transcriptome-wide epistasis ( $s$ ). To measure  $s$ , we find the line of best fit and determine its slope. Genes that act log-additively on a phenotype (**Ph**) will have  $s = 0$  (null hypothesis, orange line); whereas genes that act along an unbranched pathway will have  $s = -1/2$  (blue line). Strong repression is reflected by  $s = -1$  (red line), whereas  $s > 0$  correspond to synthetic interactions (purple line). (B) Epistasis plot showing that the *egl-9(lf)*; *vhl-1(lf)* transcriptome deviates significantly from a null additive. Points are colored qualitatively according to density (purple—low, yellow—high) and size is inversely proportional to the standard error (S.E.) of the y-axis. The green line is the line of best fit from an orthogonal distance regression. (C) Comparison of simulated epistatic coefficients against the observed coefficient. Green curve shows the bootstrapped observed transcriptome-wide epistasis coefficient for *egl-9* and *vhl-1*. Dashed green line shows the mean value of the data. Simulations use only the single mutant data to idealize what expression of the double mutant should look like.  $a > b$  means that the phenotype of  $a$  is observed in a double mutant  $a^-b^-$ .

(the fraction of inputs common to both nodes); and the genes that are regulated by both nodes (the fraction of outputs that are common to both nodes). In other words, we expected that expression profiles associated with a pathway would respond quantitatively to quantitative changes in activity of the pathway. Targeting a pathway at multiple points would lead to expression profile divergence as we compare nodes that are separated by more degrees of freedom, reflecting the flux in information between them.

1690

We investigated this possibility by weighting the robust Bayesian regression between each pair of genotypes by the size of the shared transcriptomic phenotype of each pair divided by the total number of isoforms differentially expressed in either mutant ( $N_{\text{Intersection}}/N_{\text{Union}}$ ). We plotted the weighted correlation of each gene pair, ordered by increasing functional distance (see Fig. 36). In every case, we see that the weighted correlation decreases monotonically due mainly, but not exclusively, to a smaller STP (see [Decorrelation Notebook](#)).

1697

We believe that this result is not due to random noise or insufficiently deep sequencing. Instead, we propose a framework in which every gene is regulated by multiple different molecular species, which induces progressive decorrelation. This decorrelation in turn has two consequences. First, decorrelation within a pathway implies that two nodes may be almost independent of each other if the functional distance between them is large. Second, it may be possible to use decorrelation dynamics to infer gene order in a branching pathway, as we have done with the hypoxia pathway.

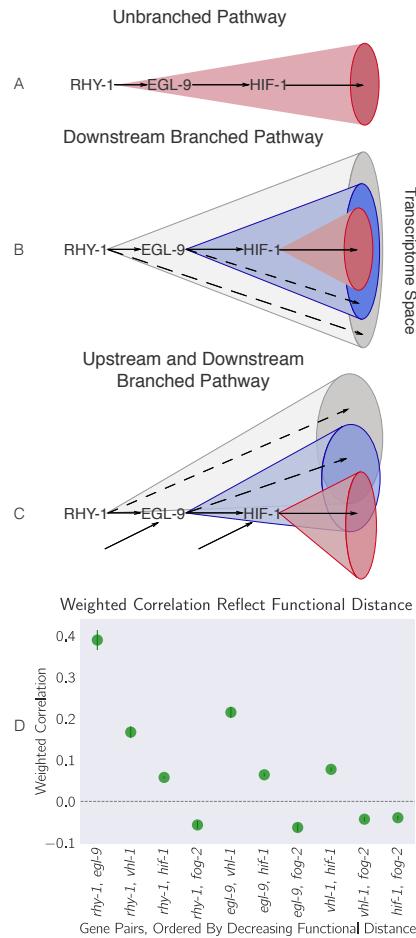
1704

### **Classical epistasis identifies a core hypoxic response**

1705

We searched for genes whose expression obeyed the two epistatic equality relationships,  $hif-1(lf)=egl-9(lf)$   $hif-1(lf)$  and  $egl-9(lf)=egl-9(lf)$ ;  $vhl-1(lf)$ , since these equalities define the hypoxia pathway. We excluded genes whose expression deviated from this relationship by more than 2 standard deviations or that had opposite changes in

1709



**Figure 36** Transcriptomes can be used to order genes in a pathway under certain assumptions. Arrows in the diagrams above are intended to show the direction of flow, and do not indicate valence. **A.** A linear pathway in which *rhy-1* is the only gene controlling *egl-9*, which in turn controls *hif-1* does not contain information to infer the order between genes. **B.** If *rhy-1* and *egl-9* have transcriptomic effects that are separable from *hif-1*, then the *rhy-1* transcriptome should contain contributions from *egl-9*, *hif-1* and *egl-9*- and *hif-1*-independent pathways. This pathway contains enough information to infer order. **C.** If a pathway is branched both upstream and downstream, transcriptomes will show even faster decorrelation. Nodes that are separated by many edges may begin to behave almost independently of each other with marginal transcriptomic overlap or correlation. **D.** The hypoxia pathway can be ordered. We hypothesize the rapid decay in correlation is due to a mixture of upstream and downstream branching that happens along this pathway. Bars show the standard error of the weighted coefficient from the Monte Carlo Markov Chain computations.

direction. Using these criteria, we identified 1,258 genes in the hypoxia response. 1710  
Tissue Enrichment Analysis showed that the intestine and epithelial system were 1711  
enriched in this response ( $q < 10^{-10}$  for both terms), consistent with previous re- 1712  
ports (Budde and Roth, 2010). Gene Enrichment Analysis (Angeles-Albores, N. 1713  
Lee, et al., 2018) showed enrichment in the mitochondrion and in collagen trimers 1714  
( $q < 10^{-10}$ ) (see [Enrichment Analysis Notebook](#) and SI Figures 3 and 4). This re- 1715  
sponse included 15 transcription factors. Even though HIF-1 is an activator, not all 1716  
of these genes were up-regulated. We reasoned that only genes that are up-regulated 1717  
in HIF-1-inhibitor mutants are candidates for direct regulation by HIF-1. We found 1718  
264 such genes. 1719

### Feedback can be inferred

While some of the rank plots contained a clear positive correlation, others showed 1721  
a discernible cross-pattern (see Fig. 34). In particular, this cross-pattern emerged 1722  
between *vhl-1(lf)* and *rhy-1(lf)* or between *vhl-1(lf)* and *egl-9(lf)*, even though *vhl-1*, 1723  
*rhy-1* and *egl-9* are all inhibitors of *hif-1(lf)*. Such cross-patterns could be indicative 1724  
of feedback loops or other complex interaction patterns. If the above is correct, 1725  
then it should be possible to identify genes that are regulated by *rhy-1* in a logically 1726  
consistent way: Since loss of *egl-9* causes *rhy-1* mRNA levels to increase, if this 1727  
increase leads to a significant change in RHY-1 activity, then it follows that the 1728  
*egl-9(lf)* and *rhy-1(lf)* should show anti-correlation in a subset of genes. Since we 1729  
do not observe many genes that are anti-correlated, we conclude that is unlikely that 1730  
the change in *rhy-1* mRNA expression causes a significant change in RHY-1 activity 1731  
under normoxic conditions. We also searched for genes with *hif-1*-independent, 1732  
*vhl-1*-dependent gene expression and found 71 genes (SI File 1). 1733

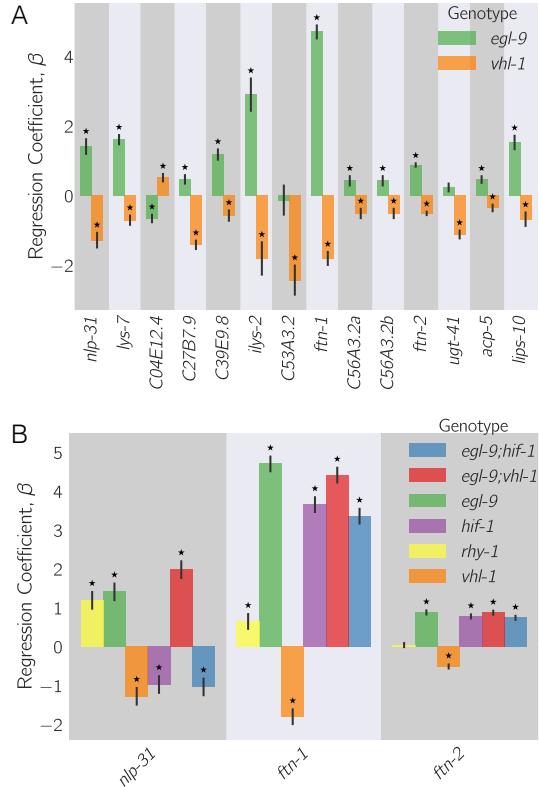
## Identification of non-classical epistatic interactions

1734

*hif-1(lf)* has traditionally been viewed as existing in a genetic OFF state under 1735 normoxic conditions. However, our dataset indicates that 1,075 genes show altered 1736 expression when *hif-1* function is removed in normoxic conditions. Moreover, we 1737 observed positive correlations between *hif-1(lf)*  $\beta$  coefficients and *egl-9(lf)*, *vhl-1(lf)* 1738 and *rhy-1(lf)*  $\beta$  coefficients in spite of the negative regulatory relationships between 1739 these genes and *hif-1*. Such positive correlations could indicate a relationship 1740 between these genes that has not been reported previously. 1741

We identified genes that exhibited violations of the canonical genetic model of the 1742 hypoxia pathway (see Fig. 37; also [Non-canonical epistasis notebook](#)). We searched 1743 for genes that changed in different directions between *egl-9(lf)* and *vhl-1(lf)*, or, 1744 equivalently, between *rhy-1(lf)* and *vhl-1(lf)* (we assume that all results from the 1745 *rhy-1(lf)* transcriptome reflect a complete loss of *egl-9* activity) without specifying 1746 any further conditions. We found 56 that satisfied this condition (see Fig. 37, SI 1747 File 1). When we checked expression of these genes in the double mutant, we found 1748 that *egl-9* remained epistatic over *vhl-1* for this class of genes. This class of genes 1749 may in fact be larger because it overlooks genes that have wild-type expression in an 1750 *egl-9(lf)* background, altered expression in a *vhl-1(lf)* background, and suppressed 1751 (wild-type) expression in an *egl-9(lf); vhl-1(lf)* background. As a result, it could 1752 help explain why the *hif-1(lf)* mutant transcriptome is positively correlated with its 1753 inhibitors. 1754

Although this entire class had similar behavior, we focused on two genes, *nlp-31* 1755 and *ftn-1* which have representative expression patterns. *ftn-1* is described to be 1756 responsive to mutations in the hypoxia pathway and has been reported to have 1757 aberrant behaviors; specifically, loss of function of *egl-9* and *vhl-1* have opposing 1758 effects on *ftn-1* expression (Ackerman and Gems, [2012](#); Romney et al., [2011](#)). 1759 These studies showed the same *ftn-1* expression phenotypes using RNAi and alleles, 1760



**Figure 37 A.** 56 genes in *C. elegans* exhibit non-classical epistasis in the hypoxia pathway, characterized by opposite effects on gene expression, relative to the wild type, of the *vhl-1(lf)* compared to *egl-9(lf)* (or *rhy-1(lf)*) mutants. Shown are a random selection of 15 out of 56 genes for illustrative purposes. **B.** Genes that behave non-canonically have a consistent pattern. *vhl-1(lf)* mutants have an opposite effect to *egl-9(lf)*, but *egl-9* remains epistatic to *vhl-1* and loss-of-function mutations in *hif-1* suppress the *egl-9(lf)* phenotype. Asterisks show  $\beta$  values significantly different from 0 relative to wild type ( $q < 10^{-1}$ ).

allaying concerns of strain-specific interference. We observed that *hif-1* was epistatic 1761 to *egl-9*, and that *egl-9* and *hif-1* both promoted *ftn-1* expression. 1762

Analysis of *ftn-1* expression reveals that *egl-9* is epistatic to *hif-1*; that *vhl-1* has 1763 opposite effects to *egl-9*, and that *vhl-1* is epistatic to *egl-9*. Analysis of *nlp-31* reveals 1764 similar relationships. *nlp-31* expression is decreased in *hif-1(lf)*, and increased in 1765 *egl-9(lf)*. However, *egl-9* is epistatic to *hif-1*. Like *ftn-1*, *vhl-1* has the opposite 1766 effect to *egl-9*, yet is epistatic to *egl-9*. We propose in the Discussion a novel model 1767 for how HIF-1 might regulate these targets. 1768

**Discussion**

1769

**The *C. elegans* hypoxia pathway can be reconstructed *de novo* from RNA-seq data**

1770

We have shown that whole-organism transcriptomic phenotypes can be used to 1772 reconstruct genetic pathways and to discern previously uncharacterized genetic in- 1773 teractions. We successfully reconstructed the hypoxia pathway including the order of 1774 action of the genetic components and its branching pattern. These results highlight 1775 the potential of whole-animal expression profiles for dissecting molecular pathways 1776 that are expressed in a large number of cells within an organism. While our results 1777 are promising, it remains to be seen whether our approach will also work for path- 1778 ways that act in a few cells. We selected a previously characterized pathway because 1779 *C. elegans* is less amenable to high-throughput screens compared to cultured cells. 1780 That said, the striking nature of our results makes us optimistic that this technique 1781 could be successfully used to reconstruct unknown pathways. 1782

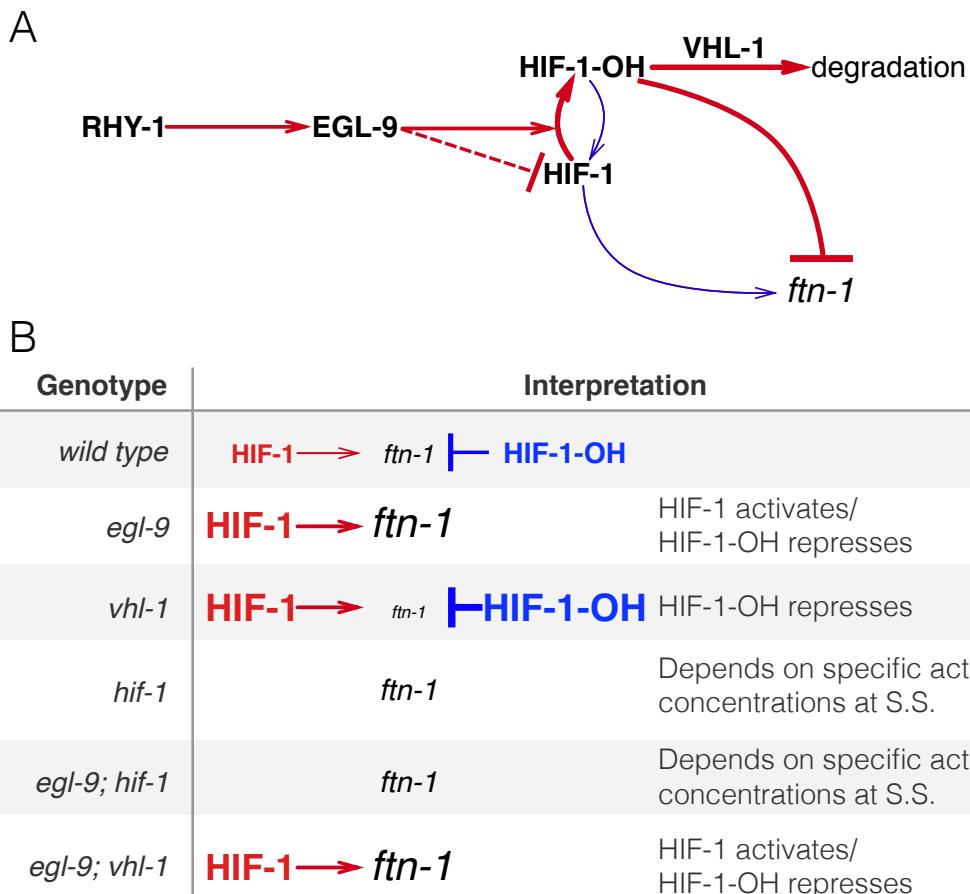
**Interpretation of the non-classical epistasis in the hypoxia pathway**

1783

The 56 genes that exhibit a striking pattern of non-classical epistasis suggest the 1784 existence of previously undescribed aspects of the hypoxia pathway. Some of these 1785 non-classical behaviors had been observed previously (Ackerman and Gems, 2012; 1786 Romney et al., 2011; Luhachack et al., 2012), but no satisfactory mechanism has 1787 been proposed to explain them. Previous studies (Romney et al., 2011; Ackerman 1788 and Gems, 2012) suggested that HIF-1 integrates information on iron concentration 1789 in the cell to determine its binding affinity to the *ftn-1* promoter, but could not 1790 definitively establish a mechanism. It is unclear why deletion of *hif-1* and deletion 1791 of *egl-9* both cause induction of *ftn-1* expression, but deletion of *vhl-1* abolishes this 1792 induction. Moreover, Luchachack et al (Luhachack et al., 2012) have previously 1793 reported that certain genes important for the *C. elegans* immune response against 1794

pathogens reflect similar non-canonical expression patterns. Their interpretation 1795 was that *swan-1*, which encodes a binding partner to EGL-9 (Shao, Zhang, Ye, 1796 et al., 2010), is important for modulating HIF-1 activity in some manner. The lack 1797 of a conclusive double mutant analysis in this work means the role of SWAN-1 in 1798 modulation of HIF-1 activity remains to be demonstrated. Other mechanisms, such 1799 as tissue-specific differences in the pathway (Budde and Roth, 2010) could also 1800 modulate expression, though it is worth pointing out that *ftn-1* expression appears 1801 restricted to a single tissue, the intestine (Kim et al., 2004). Another possibility is 1802 that *egl-9* controls *hif-1* mRNA stability via other *vhl-1*-independent pathways, but 1803 we did not see a decreases in *hif-1* level in *egl-9(lf)*, *rhy-1(lf)* or *vhl-1(lf)* mutants. 1804 Another possibility, such as control of protein stability via *egl-9* independently of 1805 *vhl-1* (Chintala et al., 2012) will not lead to splitting unless it happens in a tissue- 1806 specific manner. 1807

One parsimonious solution is to consider HIF-1 as a protein with both activating and 1808 inhibiting states. In fact, HIF-1 already exists in two states in *C. elegans*: unmodi- 1809 fied HIF-1 and HIF-1-hydroxyl (HIF-1-OH). Under this model, the effects of HIF-1 1810 for certain genes like *ftn-1* or *nlp-31* are antagonized by HIF-1-hydroxyl, which is 1811 present at only a low level in the cell in normoxia because it is degraded in a *vhl-1*- 1812 dependent fashion. This means that loss of *vhl-1* stabilizes HIF-1-hydroxyl. If *vhl-1* 1813 is inactivated, genes that are sensitive to HIF-1-hydroxyl will be inhibited as a result 1814 of the increase in HIF-1-hydroxyl, despite the increased levels of non-hydroxylated 1815 HIF-1. On the other hand, *egl-9(lf)* abrogates the generation of HIF-1-hydroxyl, 1816 stimulating accumulation of non-hydroxylated HIF-1 and promoting gene expres- 1817 sion. Whether deletion of *hif-1(lf)* is overall activating or inhibiting will depend 1818 on the relative activity of each protein state under normoxia (see Fig. 38). HIF-1- 1819 hydroxyl is challenging to study genetically, and if it does have the activity suggested 1820 by our genetic evidence this may have prevented such a role from being detected. 1821



**Figure 38** A hypothetical model showing a mechanism where HIF-1-hydroxyl antagonizes HIF-1 in normoxia. **A.** Diagram showing that RHY-1 activates EGL-9. EGL-9 hydroxylates HIF-1 in an oxygen-dependent manner. HIF-1 is rapidly hydroxylated and the product, HIF-1-OH is rapidly degraded in a VHL-1-dependent fashion. EGL-9 can also inhibit HIF-1 in an oxygen-independent fashion. In our model, HIF-1 and HIF-1-OH have opposing effects on transcription. The width of the arrows represents rates in normoxic conditions. **B.** Table showing the effects of loss-of-function mutations on HIF-1 and HIF-1-OH activity, showing how this can potentially explain the *ftn-1* expression levels in each case. S.S = Steady-state.

No mimetic mutations are known with which to study the pure hydroxylated HIF-1 species, and mutations in the Von Hippel-Lindau gene that stabilize the hydroxyl species also increase the quantity of non-hydroxylated HIF-1 by mass action.

Because HIF-1 is detected at low levels in cells under normoxic conditions (Wang and Semenza, 1993), total HIF-1 protein levels are assumed to be so low as to be biologically inactive. However, our data show 1,075 genes change expression in response to loss of *hif-1* under normoxic conditions, which establishes that there is sufficient total HIF-1 protein to be biologically active. Our analyses also revealed that *hif-1(lf)* shares positive correlations with *egl-9(lf)*, *rhy-1(lf)* and *vhl-1(lf)*, and that each of these genotypes also shows a secondary negative rank-ordered expression correlation with each other.

A homeostatic argument can be made in favor of the activity of HIF-1-hydroxyl. The cell must continuously monitor multiple metabolite levels. The *hif-1*-dependent hypoxia response integrates information from O<sub>2</sub>, α-ketoglutarate and iron concentrations in the cell. One way to integrate this information is by encoding it within the effective hydroxylation rate of HIF-1 by EGL-9. Then the dynamics in this system will evolve exclusively as a result of the total amount of HIF-1 in the cell. Such a system can be sensitive to fluctuations in the absolute concentration of HIF-1 (Goen- toro et al., 2009). Since the absolute levels of HIF-1 are low in normoxic conditions, small fluctuations in protein copy-number can represent a large fold-change in HIF-1 levels. These fluctuations might not be problematic for genes that must be turned on only under conditions of severe hypoxia—presumably, these genes would be activated only when HIF-1 levels increase far beyond random fluctuations.

For yet other sets of genes that must change expression in response to the hypoxia pathway, it may not be sufficient to integrate metabolite information exclusively via EGL-9-dependent hydroxylation of HIF-1. In particular, genes that may function

to increase survival in mild hypoxia may benefit from regulatory mechanisms that 1848 can sense minor changes in environmental conditions and which therefore benefit 1849 from robustness to transient changes in protein copy number. Likewise, genes 1850 that are involved in iron or  $\alpha$ -ketoglutarate metabolism (such as *ftn-1*) may benefit 1851 from being able to sense, accurately, small and consistent deviations from basal 1852 concentrations of these metabolites. For these genes, the information may be better 1853 encoded by using HIF-1 and HIF-1-hydroxyl as an activator/repressor pair. Such 1854 circuits are known to possess distinct advantages for controlling output robustly to 1855 transient fluctuations in the levels of their components (Hart, Antebi, et al., 2012; 1856 Hart and Alon, 2013). 1857

Our RNA-seq data suggests that one of these atypical targets of HIF-1 may be RHY- 1858 1. Although *rhy-1* does not exhibit non-classical epistasis, all genotypes containing 1859 a *hif-1(lf)* mutation had increased expression levels of *rhy-1*. We speculate that if 1860 *rhy-1* is controlled by both HIF-1 and HIF-1-hydroxyl, then this might imply that 1861 HIF-1 auto-regulates both positively and negatively. 1862

### Strengths and weaknesses of the methodology

We have described a set of methods that can in principle be applied to any multi- 1864 dimensional phenotype. Although we have not applied these methods to *de novo* 1865 pathway discovery, we believe that they will be broadly applicable to a wide variety 1866 of genetic problems. One aspect of our methodology is the use of whole-organism 1867 expression data. Data collection from whole-organisms can be rapid with low 1868 technical barriers. On the other hand, a concern is that whole-organism data will 1869 average signals across tissues, which would limit the scope of this technology to the 1870 study of genetic pathways that are systemic or expressed in large tissues. In real- 1871 ity, our method may be applicable for pathways that are expressed even in a small 1872 number of cells in an organism. If a pathway is active in a single cell, this does 1873

not mean that it does not have cell-non-autonomous effects that could be detected 1874 on an organism-wide level. Thus, pathways that act in single cells could still be 1875 characterized via whole-organism transcriptome profiling. If the non-autonomous 1876 effects are long-lasting, then the profiling could take place after the time-of-action 1877 of this pathway. In fact, this is how the female-like state in *C. elegans* was recently 1878 identified (Angeles-Albores, Leighton, et al., 2017): *fog-2* is involved in translation 1879 repression of *tra-2* in the somatic gonad, thereby promoting sperm formation in 1880 late larvae (Clifford et al., 2000). Loss of this gene causes non-cell-autonomous 1881 effects that can be detected well after the time-of-action of *fog-2* in the somatic 1882 gonad has ended. Therefore, we believe that our methodology will be applicable to 1883 many genetic cases, with the exception of pathways that acts in complex, antago- 1884 nistic manners depending on the cell type, or if the pathway minimally affects gene 1885 expression. 1886

Genetic analysis of transcriptomic data has proved challenging as a result of its com- 1887 plexity. Although dimensionality reduction techniques such as PCA have emerged 1888 as powerful methods with which to understand these data, these methods generate 1889 reduced coordinates which are difficult or impossible to interpret. As an example, 1890 the first principal component in this paper (see Fig. 33) could be interpreted as HIF-1 1891 pseudo-abundance (Lönnberg et al., 2017). However, another equally reasonable, 1892 yet potentially completely different interpretation, is as a pseudo-HIF-1/HIF-1-OH 1893 ratio. Another way to analyze genetic interactions is via general linear models 1894 (GLMs) that include interaction terms between two or more genes. GLMs can 1895 quantify the genetic interactions on single transcripts. We and others (Dixit et al., 1896 2016; Angeles-Albores, Leighton, et al., 2017) have used GLMs to perform epis- 1897 tasis analyses of pathways using transcriptomic phenotypes. GLMs are powerful, 1898 but they generate a different interaction coefficient for each gene measured. The 1899 large number of coefficients makes interpretation of the genetic interaction between 1900

two mutants difficult. Previous approaches (Dixit et al., 2016) visualize these co-  
efficients via clustered heatmaps. However, two clusters cannot be assumed to be  
evidence that two genes interact via entirely distinct pathways. Indeed, the non-  
classical epistasis examples we described here might cluster separately even though  
a reasonable model can be invoked that does not require any new molecular players.

The epistasis plots shown here are a useful way to visualize epistasis in vectorial  
phenotypes. We have shown how an epistasis plot can be used to identify inter-  
actions between two genes by examining the transcriptional phenotypes of single  
and double mutants. Epistasis plots can accumulate an arbitrary number of points  
within them, possess a rich structure that can be visualized and have straightfor-  
ward interpretations for special slope values. Epistasis plots and GLMs are not  
mutually exclusive. A GLM could be used to quantify epistasis interactions at  
single-transcript resolution, and the results then analyzed using an epistasis plot (for  
a non-genetic example, see Angeles-Albores, Leighton, et al. (2017)). A benefit of  
epistasis plots is that they enable the computation of a single, aggregate statistic that  
describes the ensemble behavior of a set of genes. This aggregate statistic is not  
enough to describe all possible behaviors in a system, but it can be used to establish  
whether the genes under study are part of a single pathway. In the case of the  
hypoxia pathway, phenotypes that are downstream of the hypoxia pathway should  
conform to the genetic equalities,  $egl-9(lf)$   $hif-1(lf) = hif-1(lf)$  AND  $egl-9(lf); vhl-1(lf) = egl-9(lf)$ . Genes whose expression levels behave strangely, yet satisfy these equal-  
ities are downstream of the hypoxia pathway. These anomalous genes cannot be  
identified via the epistasis coefficient but the epistasis coefficient does provide a uni-  
fying framework with which to analyze them by constraining the space of plausible  
hypotheses.

Until relatively recently, the rapid generation and molecular characterization of  
null mutants was a major bottleneck for genetic analyses. Advances in genomic

1925

1926

1927

engineering mean that, for a number of organisms, production of mutants is now 1928  
rapid and efficient. As mutants become easier to produce, biologists are realizing 1929  
that phenotyping and characterizing the biological functions of individual genes is 1930  
challenging. This is particularly true for whole organisms, where subtle phenotypes 1931  
can go undetected for long periods of time. We have shown that whole-animal 1932  
RNA-sequencing is a sensitive method that can be seamlessly incorporated with 1933  
genetic analyses of epistasis. 1934

## Methods

1935

### Nematode strains and culture

1936

Strains used were N2 (Bristol), JT307 *egl-9(sa307)*, CB5602 *vhl-1(ok161)*, ZG31 1937  
*hif-1(ia4)*, RB1297 *rhy-1(ok1402)*, CB6088 *egl-9(sa307) hif-1(ia4)*, CB6116 1938  
*egl-9(sa307);vhl-1(ok161)*, Lines were grown on standard nematode growth media 1939  
Petri plates seeded with OP50 *E. coli* at 20°C (Sulston and Brenner, 1974). 1940

### RNA isolation

1941

Lines were synchronized by harvesting eggs via sodium hypochlorite treatment and 1942  
subsequently plating eggs on food. Worms were staged and based on the time after 1943  
plating, vulva morphology and the absence of eggs. 30–50 non-gravid young adults 1944  
were picked and placed in 100 µL of TE pH 8.0 (Ambion AM9849) in 0.2 mL PCR 1945  
tubes on ice. Worms were allowed to settle or spun down by centrifugation and 1946  
~ 80 µL of supernatant removed before flash-freezing in liquid  $N_2$ . These samples 1947  
were digested with Recombinant Proteinase K PCR Grade (Roche Lot No. 03115 1948  
838001) for 15 min at 60° in the presence of 1% SDS and 1.25 µL RNA Secure 1949  
(Ambion AM7005). 5 volumes of Trizol (Tri-Reagent Zymo Research) were added 1950  
to the RNA samples and treated with DNase I using Zymo Research Quick-RNA 1951  
MicroPrep R1050. Samples were analyzed run on an Agilent 2100 BioAnalyzer 1952  
(Agilent Technologies). Replicates were selected that had RNA integrity numbers 1953

equal to or greater than 9.0 and without bacterial ribosomal bands, except for the 1954  
ZG31 mutant where one of three replicates had a RIN of 8.3. 1955

### Library preparation and sequencing 1956

10 ng of total RNA from each sample was reverse-transcribed into cDNA using the 1957  
Clontech SMARTer Ultra Low Input RNA for Sequencing v3 kit (catalog #634848) 1958  
in the SMARTSeq2 protocol (Picelli et al., 2014). RNA was denatured at 70°C 1959  
for 3 min in the presence of dNTPs, oligo dT primer and spiked-in quantitation 1960  
standards (NIST/ERCC from Ambion, catalog #4456740). After chilling to 4°C, 1961  
the first-strand reaction was assembled using a LNA TSO primer (Picelli et al., 1962  
2014), and run at 42°C for 90 minutes, followed by denaturation at 70°C for 10 1963  
min. The first strand reaction was used as template for 13 cycles of PCR using the 1964  
Clontech v3 kit. Reactions were purified with Ampure XP SPRI beads (catalog 1965  
#A63880). After quantification using the Qubit High Sensitivity DNA assay, a 3 ng 1966  
aliquot of the cDNA was run on the Agilent HS DNA chip to confirm the length 1967  
distribution of the amplified fragments. The median value for the average cDNA 1968  
lengths from all length distributions was 1,076 bp. Tagmentation of the full length 1969  
cDNA was performed using the Illumina/Nextera DNA library prep kit (catalog 1970  
#FC-121-1030). Following Qubit quantitation and Agilent BioAnalyzer profiling, 1971  
the tagmented libraries were sequenced on an Illumina HiSeq2500 machine in single 1972  
read mode with a read length of 50 nt to a depth of 15 million reads per sample. 1973  
Base calls were performed with RTA 1.13.48.0 followed by conversion to FASTQ 1974  
with bcl2fastq 1.8.4. 1975

### Read alignment and differential expression analysis 1976

We used Kallisto (Bray et al., 2016) to perform read pseudo-alignment and per- 1977  
formed differential analysis using Sleuth (Pimentel et al., 2016). We fit a general 1978

linear model for an isoform  $t$  in sample  $i$ :

1979

$$y_{t,i} = \beta_{t,0} + \beta_{t,genotype} \cdot X_{t,i} + \beta_{t,batch} \cdot Y_{t,i} + \epsilon_{t,i} \quad (3.1)$$

where  $y_{t,i}$  was the logarithm transformed counts of isoform  $t$  in sample  $i$ ;  $\beta_{t,genotype}$  and  $\beta_{t,batch}$  were parameters of the model for the isoform  $t$ , and which could be interpreted as biased estimators of the log-fold change;  $X_{t,i}, Y_{t,i}$  were indicator variables describing the experimental conditions of the isoform  $t$  in sample  $i$ ; and  $\epsilon_{t,i}$  was the noise associated with a particular measurement. After fitting the general linear model, we tested isoforms for differential expression using the built-in Wald-test in Sleuth (Pimentel et al., 2016), which outputs a  $q$ -value that has been corrected for multiple hypothesis testing.

1987

## Genetic Analysis, Overview

1988

The processed data were analyzed using Python 3.5. We used the Pandas, Matplotlib, Scipy, Seaborn, Sklearn, Networkx, PyMC3, and TEA libraries (McKinney, 2011; Oliphant, 2007; Pedregosa et al., 2012; Salvatier, Wiecki, and Fonnesbeck, 2015; Van Der Walt, Colbert, and Varoquaux, 2011; Hunter, 2007; Angeles-Albores, N. Lee, et al., 2016; Waskom et al., 2016). Our analysis is available in Jupyter Notebooks (Pérez and Granger, 2007). All code and processed data are available at <https://github.com/WormLabCaltech/mprsq> along with version-control information. Our Jupyter Notebook and interactive graphs for this project can be found at <https://wormlabcaltech.github.io/mprsq/> in html format. Raw reads were deposited in the Short Read Archive under the study accession number SRP100886 and in the GEO under the accession number GSE97355.

1999

## Weighted correlations

2000

Correlations between mutants were calculated by identifying their STP. Transcripts were rank-ordered according to their regression coefficient,  $\beta$ . Regressions were

2002

performed using a Student-T distribution with the PyMC3 library (Salvatier, Wiecki, 2003  
and Fonnesbeck, 2015) (`pm.glm.families.StudenT` in Python). If the correlations had an average value  $> 1$ , the average correlation coefficient was set to 1. Weights were calculated as the number of genes that were inliers divided by the number of DEGs present in either mutant.

2007

## Epistatic analysis

2008

The epistasis coefficient between two null mutants  $a$  and  $b$  was calculated as:

2009

$$s(a, b) = \frac{\beta_{a,b} - \beta_a - \beta_b}{\beta_a + \beta_b} \quad (3.2)$$

Null models for various epistatic relationships were generated by sampling the single mutants in an appropriate fashion. For example, to generate the distribution for two mutants that obey the epistatic relationship  $a^- = a^-b^-$ , we substituted  $\beta_{a,b}$  with  $\beta_a$  and bootstrapped the result.

2013

To select between theoretical models, we implemented an approximate Bayesian Odds Ratio. We defined a free-fit model,  $M_1$ , that found the line of best fit for the data:

2016

$$P(\alpha | M_1, D) \propto \prod_{(x_i, y_i, \sigma_i) \in D} \exp \left[ \frac{(y_i - \alpha \cdot x_i)^2}{2\sigma_i^2} \right] \cdot (1 + \alpha^2)^{-3/2}, \quad (3.3)$$

where  $\alpha$  was the slope to be determined,  $x_i, y_i$  are the of each point, and  $\sigma_i$  was the standard error associated with the y-value. We used equation 3.3 to obtain the most likely slope given the data,  $D$ , via minimization (`scipy.optimize.minimize` in Python). Finally, we approximated the odds ratio as:

2020

$$OR = \frac{P(D | \alpha^*, M_1) \cdot (2\pi)^{1/2} \sigma_{\alpha^*}}{P(D | M_i)}, \quad (3.4)$$

where  $\alpha^*$  was the slope found after minimization,  $\sigma_{\alpha^*}$  was the standard deviation of the parameter at the point  $\alpha^*$  and  $P(D | M_i)$  was the probability of the data given the parameter-free model,  $M_i$ .

2023

**Enrichment analysis**

2024

Tissue, Phenotype and Gene Ontology Enrichment Analysis were carried out using the WormBase Enrichment Suite for Python (Angeles-Albores, N. Lee, et al., 2018; Angeles-Albores, N. Lee, et al., 2016).

2027

**References**

2028

- Ackerman, Daniel and David Gems (2012). “Insulin/IGF-1 and hypoxia signaling act in concert to regulate iron homeostasis in *Caenorhabditis elegans*”. In: *PLoS Genetics* 8.3. ISSN: 15537390. doi: [10.1371/journal.pgen.1002498](https://doi.org/10.1371/journal.pgen.1002498). 2029  
2030  
2031
- Angeles-Albores, David, Daniel H. W. Leighton, et al. (2017). “The *Caenorhabditis elegans* Female-Like State: Decoupling the Transcriptomic Effects of Aging and Sperm Status”. In: *G3: Genes, Genomes, Genetics* 7.9. 2032  
2033  
2034
- Angeles-Albores, David, Raymond Y. N. Lee, et al. (2016). “Tissue enrichment analysis for *C. elegans* genomics”. In: *BMC Bioinformatics* 17.1, p. 366. ISSN: 1471-2105. doi: [10.1186/s12859-016-1229-9](https://doi.org/10.1186/s12859-016-1229-9). 2035  
2036  
2037
- (2018). “Two new functions in the WormBase Enrichment Suite”. In: *Micropublication: biology. Dataset*. doi: <https://doi.org/10.17912/W25Q2N>. 2038  
2039
- Bray, Nicolas L et al. (2016). “Near-optimal probabilistic RNA-seq quantification.” In: *Nature biotechnology* 34.5, pp. 525–7. ISSN: 1546-1696. doi: [10.1038/nbt.3519](https://doi.org/10.1038/nbt.3519). arXiv: [1505.02710](https://arxiv.org/abs/1505.02710). 2040  
2041  
2042
- Brem, Rachel B. et al. (2002). “Genetic Dissection of Transcriptional Regulation in Budding Yeast”. In: *Science* 296.5568. 2043  
2044
- Budde, Mark W. and Mark B. Roth (2010). “Hydrogen Sulfide Increases Hypoxia-inducible Factor-1 Activity Independently of von Hippel–Lindau Tumor Suppressor 1 in *C. elegans*”. In: *Molecular biology of the cell* 21, pp. 212–217. ISSN: 1939-4586. doi: [10.1091/mbc.E09](https://doi.org/10.1091/mbc.E09). 2045  
2046  
2047  
2048
- Capaldi, Andrew P et al. (Nov. 2008). “Structure and function of a transcriptional network activated by the MAPK Hog1”. In: *Nature Genetics* 40.11, pp. 1300–1306. ISSN: 1061-4036. doi: [10.1038/ng.235](https://doi.org/10.1038/ng.235). 2049  
2050  
2051
- Chintala, Sreenivasulu et al. (2012). “Prolyl hydroxylase 2 dependent and Von-Hippel-Lindau independent degradation of Hypoxia-inducible factor 1 and 2 alpha by selenium in clear cell renal cell carcinoma leads to tumor growth inhibition”. In: *BMC Cancer* 12.1, p. 293. ISSN: 1471-2407. doi: [10.1186/1471-2407-12-293](https://doi.org/10.1186/1471-2407-12-293). 2052  
2053  
2054  
2055  
2056
- Clifford, Robert et al. (2000). “FOG-2, a novel F-box containing protein, associates with the GLD-1 RNA binding protein and directs male sex determination in the *C. elegans* hermaphrodite germline.” In: *Development (Cambridge, England)* 127.24, pp. 5265–5276. ISSN: 0950-1991. 2057  
2058  
2059  
2060

- Dixit, Atry et al. (2016). "Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens". In: *Cell* 167.7, pp. 1853–1866.e17. ISSN: 00928674. doi: [10.1016/j.cell.2016.11.038](https://doi.org/10.1016/j.cell.2016.11.038). 2061  
2062  
2063
- Epstein, Andrew C. R. et al. (2001). "*C. elegans* EGL-9 and mammalian homologs define a family of dioxygenases that regulate HIF by prolyl hydroxylation". In: *Cell* 107.1, pp. 43–54. ISSN: 00928674. doi: [10.1016/S0092-8674\(01\)00507-4](https://doi.org/10.1016/S0092-8674(01)00507-4). 2064  
2065  
2066
- Goentoro, Lea et al. (2009). "The Incoherent Feedforward Loop Can Provide Fold-Change Detection in Gene Regulation". In: *Molecular Cell* 36.5, pp. 894–899. ISSN: 10972765. doi: [10.1016/j.molcel.2009.11.018](https://doi.org/10.1016/j.molcel.2009.11.018). arXiv: [NIHMS150003](https://arxiv.org/abs/150003). 2067  
2068  
2069  
2070
- Hart, Yuval and Uri Alon (2013). "The Utility of Paradoxical Components in Biological Circuits". In: *Molecular Cell* 49.2, pp. 213–221. ISSN: 10972765. doi: [10.1016/j.molcel.2013.01.004](https://doi.org/10.1016/j.molcel.2013.01.004). 2071  
2072  
2073
- Hart, Yuval, Yaron E Antebi, et al. (2012). "Design principles of cell circuits with paradoxical components". In: *Proceedings of the National Academy of Sciences of the United States of America* 109.21, pp. 8346–8351. ISSN: 0027-8424. doi: [10.1073/pnas.1117475109](https://doi.org/10.1073/pnas.1117475109). 2074  
2075  
2076  
2077
- Huang, L. Eric et al. (1996). "Activation of hypoxia-inducible transcription factor depends primarily upon redox-sensitive stabilization of its alpha subunit". In: *Journal of Biological Chemistry* 271.50, pp. 32253–32259. ISSN: 00219258. doi: [10.1074/jbc.271.50.32253](https://doi.org/10.1074/jbc.271.50.32253). 2078  
2079  
2080  
2081
- Huang, Linda S and Paul W Sternberg (2006). "Genetic dissection of developmental pathways." In: *WormBook: the online review of C. elegans biology* 1995, pp. 1–19. ISSN: 1551-8507. doi: [10.1895/wormbook.1.88.2](https://doi.org/10.1895/wormbook.1.88.2). 2082  
2083  
2084
- Hughes, Timothy R. et al. (2000). "Functional Discovery via a Compendium of Expression Profiles". In: *Cell* 102.1, pp. 109–126. ISSN: 00928674. doi: [10.1016/S0092-8674\(00\)00015-5](https://doi.org/10.1016/S0092-8674(00)00015-5). 2085  
2086  
2087
- Hunter, John D. (2007). "Matplotlib: A 2D graphics environment". In: *Computing in Science and Engineering* 9.3, pp. 99–104. ISSN: 15219615. doi: [10.1109/MCSE.2007.55](https://doi.org/10.1109/MCSE.2007.55). arXiv: [0402594v3 \[arXiv:cond-mat\]](https://arxiv.org/abs/0402594v3). 2088  
2089  
2090
- Jiang, B H et al. (1996). "Dimerization, DNA binding, and transactivation properties of hypoxia-inducible factor 1." In: *The Journal of biological chemistry* 271.30, pp. 17771–17778. ISSN: 00219258. doi: [10.1074/jbc.271.30.17771](https://doi.org/10.1074/jbc.271.30.17771). 2091  
2092  
2093
- Jiang, Huaqi, Rong Guo, and Jo Anne Powell-Coffman (2001). "The *Caenorhabditis elegans hif-1* gene encodes a bHLH-PAS protein that is required for adaptation to hypoxia." In: *Proceedings of the National Academy of Sciences of the United States of America* 98.14, pp. 7916–7921. ISSN: 0027-8424. doi: [10.1073/pnas.141234698](https://doi.org/10.1073/pnas.141234698). 2094  
2095  
2096  
2097  
2098

- Kaelin, William G. and Peter J. Ratcliffe (2008). "Oxygen Sensing by Metazoans: 2099  
The Central Role of the HIF Hydroxylase Pathway". In: *Molecular Cell* 30.4, 2100  
pp. 393–402. ISSN: 10972765. doi: [10.1016/j.molcel.2008.04.009](https://doi.org/10.1016/j.molcel.2008.04.009). 2101
- Kim, Young-Il et al. (2004). "Transcriptional Regulation and Life-span Modulation 2102  
of Cytosolic Aconitase and Ferritin Genes in *C.elegans*". In: *Journal of Molecular 2103  
Biology* 342.2, pp. 421–433. ISSN: 00222836. doi: [10.1016/j.jmb.2004.07.036](https://doi.org/10.1016/j.jmb.2004.07.036). 2104  
2105
- King, Elizabeth G. et al. (May 2014). "Genetic Dissection of the *Drosophila* 2106  
*melanogaster* Female Head Transcriptome Reveals Widespread Allelic Hetero- 2107  
geneity". In: *PLoS Genetics* 10.5. Ed. by Greg Gibson, e1004322. ISSN: 1553- 2108  
7404. doi: [10.1371/journal.pgen.1004322](https://doi.org/10.1371/journal.pgen.1004322). 2109
- Li, Yang et al. (2006). "Mapping Determinants of Gene Expression Plasticity by 2110  
Genetical Genomics in *C. elegans*". In: *PLoS Genetics* 2.12, e222. ISSN: 1553- 2111  
7390. doi: [10.1371/journal.pgen.0020222](https://doi.org/10.1371/journal.pgen.0020222). 2112
- Loenarz, Christoph et al. (2011). "The hypoxia-inducible transcription factor path- 2113  
way regulates oxygen sensing in the simplest animal, *Trichoplax adhaerens*". In: 2114  
*EMBO reports* 12.1, pp. 63–70. ISSN: 1469-221X. doi: [10.1038/embor.2010.170](https://doi.org/10.1038/embor.2010.170). 2115  
2116
- Lönnberg, Tapiro et al. (2017). "Single-cell RNA-seq and computational analysis 2117  
using temporal mixture modeling resolves TH1/TFH fate bifurcation in malaria". 2118  
In: *Science Immunology* 2.9. 2119
- Luhachack, Lyly G. et al. (2012). "EGL-9 Controls *C. elegans* Host Defense Speci- 2120  
ficity through Prolyl Hydroxylation-Dependent and -Independent HIF-1 Path- 2121  
ways". In: *PLoS Pathogens* 8.7, p. 48. ISSN: 15537366. doi: [10.1371/journal.ppat.1002798](https://doi.org/10.1371/journal.ppat.1002798). 2122  
2123
- Ma, Dengke K. et al. (2012). "CYSL-1 Interacts with the O 2-Sensing Hydroxylase 2124  
EGL-9 to Promote H 2S-Modulated Hypoxia-Induced Behavioral Plasticity in 2125  
*C. elegans*". In: *Neuron* 73.5, pp. 925–940. ISSN: 08966273. doi: [10.1016/j.neuron.2011.12.037](https://doi.org/10.1016/j.neuron.2011.12.037). arXiv: [NIHMS150003](https://arxiv.org/abs/1503.05003). 2126  
2127
- McKinney, Wes (2011). "pandas: a Foundational Python Library for Data Analysis 2128  
and Statistics". In: *Python for High Performance and Scientific Computing*, pp. 1– 2129  
9. 2130
- Mortazavi, Ali et al. (2008). "Mapping and quantifying mammalian transcriptomes 2131  
by RNA-Seq". In: *Nature Methods* 5.7, pp. 621–628. ISSN: 1548-7091. doi: [10.1038/nmeth.1226](https://doi.org/10.1038/nmeth.1226). arXiv: [1111.6189v1](https://arxiv.org/abs/1111.6189v1). 2132  
2133
- Oliphant, Travis E (2007). "SciPy: Open source scientific tools for Python". In: 2134  
*Computing in Science and Engineering* 9, pp. 10–20. ISSN: 1521-9615. doi: [10.1109/MCSE.2007.58](https://doi.org/10.1109/MCSE.2007.58). 2135  
2136

- Pedregosa, Fabian et al. (2012). “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12, pp. 2825–2830. issn: 15324435. doi: 2138 [10.1007/s13398-014-0173-7](https://doi.org/10.1007/s13398-014-0173-7). arXiv: [1201.0490](https://arxiv.org/abs/1201.0490). 2139
- Pérez, F. and B.E. Granger (2007). “IPython: A System for Interactive Scientific Computing Python: An Open and General- Purpose Environment”. In: *Computing in Science and Engineering* 9.3, pp. 21–29. issn: 15219615. doi: [doi:10.1109/MCSE.2007.53](https://doi.org/10.1109/MCSE.2007.53). 2140 2141 2142 2143
- Phillips, Patrick C (2008). “Epistasis—the essential role of gene interactions in the structure and evolution of genetic systems”. In: *Nat Rev Genet* 9.11, pp. 855–867. 2144 2145 issn: 1471-0056. doi: [10.1038/nrg2452](https://doi.org/10.1038/nrg2452). 2146
- Picelli, Simone et al. (2014). “Full-length RNA-seq from single cells using Smart- seq2.” In: *Nature protocols* 9.1, pp. 171–81. issn: 1750-2799. doi: [10.1038/nprot.2014.006](https://doi.org/10.1038/nprot.2014.006). 2147 2148 2149
- Pimentel, Harold J et al. (2016). “Differential analysis of RNA-Seq incorporating quantification uncertainty”. In: *bioRxiv*, p. 058164. doi: [10.1101/058164](https://doi.org/10.1101/058164). 2150 2151
- Powell-Coffman, Jo Anne, Christopher A. Bradfield, and William B. Wood (1998). “*Caenorhabditis elegans* Orthologs of the Aryl Hydrocarbon Receptor and Its Heterodimerization Partner the Aryl Hydrocarbon Receptor Nuclear Translocator”. In: *Proceedings of the National Academy of Sciences* 95.6, pp. 2844–2849. 2152 2153 2154 2155 issn: 0027-8424. doi: [10.1073/pnas.95.6.2844](https://doi.org/10.1073/pnas.95.6.2844). 2156
- Romney, Steven Joshua et al. (2011). “HIF-1 regulates iron homeostasis in *Caenorhabditis elegans* by activation and inhibition of genes involved in iron uptake and storage”. In: *PLoS Genetics* 7.12. issn: 15537390. doi: [10.1371/journal.pgen.1002394](https://doi.org/10.1371/journal.pgen.1002394). 2157 2158 2159 2160
- Salvatier, John, Thomas Wiecki, and Christopher Fonnesbeck (2015). “Probabilistic Programming in Python using PyMC”. In: *PeerJ Computer Science* 2.e55, pp. 1–24. issn: 2376-5992. doi: [10.7717/peerj-cs.55](https://doi.org/10.7717/peerj-cs.55). arXiv: [1507.08050](https://arxiv.org/abs/1507.08050). 2161 2162 2163
- Schadt, Eric E. et al. (Mar. 2003). “Genetics of gene expression surveyed in maize, mouse and man”. In: *Nature* 422.6929, pp. 297–302. issn: 00280836. doi: [10.1038/nature01434](https://doi.org/10.1038/nature01434). 2164 2165 2166
- Schwarz, Erich M., Mihoko Kato, and Paul W. Sternberg (Oct. 2012). “Functional transcriptomics of a migrating cell in *Caenorhabditis elegans*.” In: *Proceedings of the National Academy of Sciences of the United States of America* 109.40, pp. 16246–51. issn: 1091-6490. doi: [10.1073/pnas.1203045109](https://doi.org/10.1073/pnas.1203045109). 2167 2168 2169 2170
- Scimone, M. Lucila et al. (2014). “Neoblast specialization in regeneration of the planarian *Schmidtea mediterranea*”. In: *Stem Cell Reports* 3.2, pp. 339–352. issn: 22136711. doi: [10.1016/j.stemcr.2014.06.001](https://doi.org/10.1016/j.stemcr.2014.06.001). 2171 2172 2173

- Shao, Zhiyong, Yi Zhang, and Jo Anne Powell-Coffman (2009). “Two Distinct Roles for EGL-9 in the Regulation of HIF-1-mediated gene expression in *Caenorhabditis elegans*”. In: *Genetics* 183.3, pp. 821–829. ISSN: 00166731. doi: [10.1534/genetics.109.107284](https://doi.org/10.1534/genetics.109.107284).
- Shao, Zhiyong, Yi Zhang, Qi Ye, et al. (2010). “*C. elegans swan-1* binds to *egl-9* and regulates *hif-1*-mediated resistance to the bacterial pathogen *Pseudomonas aeruginosa* PAO1”. In: *PLoS Pathogens* 6.8, pp. 91–92. ISSN: 15537366. doi: [10.1371/journal.ppat.1001075](https://doi.org/10.1371/journal.ppat.1001075).
- Shen, Chuan, Zhiyong Shao, and Jo Anne Powell-Coffman (2006). “The *Caenorhabditis elegans rhy-1* Gene Inhibits HIF-1 Hypoxia-Inducible Factor Activity in a Negative Feedback Loop That Does Not Include *vhl-1*”. In: *Genetics* 174.3, pp. 1205–1214. ISSN: 00166731. doi: [10.1534/genetics.106.063594](https://doi.org/10.1534/genetics.106.063594).
- Sulston, J. E. and S. Brenner (1974). “The DNA of *Caenorhabditis elegans*.” In: *Genetics* 77.1, pp. 95–104. ISSN: 00166731.
- Van Der Walt, Stéfan, S. Chris Colbert, and Gaël Varoquaux (2011). “The NumPy array: A structure for efficient numerical computation”. In: *Computing in Science and Engineering* 13.2, pp. 22–30. ISSN: 15219615. doi: [10.1109/MCSE.2011.37](https://doi.org/10.1109/MCSE.2011.37). arXiv: [1102.1523](https://arxiv.org/abs/1102.1523).
- Van Driessche, Nancy et al. (2005). “Epistasis analysis with global transcriptional phenotypes”. In: *Nature genetics* 37.5, pp. 471–477. ISSN: 1061-4036. doi: [10.1038/ng1545](https://doi.org/10.1038/ng1545).
- Van Wolfswinkel, Josien C., Daniel E. Wagner, and Peter W. Reddien (2014). “Single-cell analysis reveals functionally distinct classes within the planarian stem cell compartment”. In: *Cell Stem Cell* 15.3, pp. 326–339. ISSN: 18759777. doi: [10.1016/j.stem.2014.06.007](https://doi.org/10.1016/j.stem.2014.06.007).
- Wang, G L and G L Semenza (1993). “Characterization of hypoxia-inducible factor 1 and regulation of DNA binding activity by hypoxia.” In: *The Journal of biological chemistry* 268.29, pp. 21513–8. ISSN: 0021-9258.
- Waskom, Michael et al. (2016). “seaborn: v0.7.0 (January 2016)”. In: doi: [10.5281/zenodo.45133](https://doi.org/10.5281/zenodo.45133).
- Yeung, K. Y. and W. L. Ruzzo (Sept. 2001). “Principal component analysis for clustering gene expression data.” In: *Bioinformatics (Oxford, England)* 17.9, pp. 763–774. ISSN: 1367-4803. doi: [10.1093/bioinformatics/17.9.763](https://doi.org/10.1093/bioinformatics/17.9.763).

THE CAENORHABDITIS ELEGANS FEMALE-LIKE STATE:	2208
DECOUPLING THE TRANSCRIPTOMIC EFFECTS OF AGING	2209
AND SPERM STATUS	2210

Abstract	2211
----------	------

Understanding genome and gene function in a whole organism requires us to fully comprehend the life cycle and the physiology of the organism in question. *Caenorhabditis elegans* XX animals are hermaphrodites that exhaust their sperm after 3 days of egg-laying. Even though *C. elegans* can live for many days after cessation of egg-laying, the molecular physiology of this state has not been as intensely studied as other parts of the life cycle, despite documented changes in behavior and metabolism. To study the effects of sperm depletion and aging of *C. elegans* during the first 6 days of adulthood, we measured the transcriptomes of 1st day adult hermaphrodites; 6th day sperm-depleted adults; and at the same time points, mutant *fog-2(lf)* worms that have a feminized germline phenotype. We found that we could separate the effects of biological aging from sperm depletion. For a large subset of genes, young adult *fog-2(lf)* animals had the same gene expression changes as sperm-depleted 6th day wild-type hermaphrodites, and these genes did not change expression when *fog-2(lf)* females reached the 6th day of adulthood. Taken together, this indicates that changing sperm status causes a change in the internal state of the worm, which we call the female-like state. Our data provide a high-quality picture of the changes that happen in global gene expression throughout the period of early aging in the worm.

Transcriptome analysis by RNA-seq (Mortazavi et al., 2008) has allowed for in-depth analysis of gene expression changes between life stages and environmental conditions in many species (Gerstein et al., 2014; Blaxter et al., 2012). *Caenorhabditis elegans*, a genetic model nematode with extremely well defined and largely invariant development (Sulston and Horvitz, 1977; Sulston, Schierenberg, et al., 1983), has been subjected to extensive transcriptomic analysis across all stages of larval development (Hillier et al., 2009; Boeck et al., 2016; Murray et al., 2012) and many stages of embryonic development (Boeck et al., 2016). Although RNA-seq was used to develop transcriptional profiles of the mammalian aging process soon after its invention (Magalhães, Finch, and Janssens, 2010), few such studies have been conducted in *C. elegans* past the entrance into adulthood.

A distinct challenge to the study of aging transcriptomes in *C. elegans* is the hermaphroditic lifestyle of wild-type individuals of this species. Young adult hermaphrodites are capable of self-fertilization (Sulston and Brenner, 1974; Corsi, Wightman, and Chalfie, 2015), and the resulting embryos will contribute RNA to whole-organism RNA extractions. Most previous attempts to study the *C. elegans* aging transcriptome have addressed the aging process only indirectly, or relied on the use of genetically or chemically sterilized animals to avoid this problem (Murphy et al., 2003; Halaschek-Wiener et al., 2005; Lund et al., 2002; McCormick et al., 2012; Eckley et al., 2013; Boeck et al., 2016; Rangaraju et al., 2015). In addition, most of these studies obtained transcriptomes using microarrays, which are less accurate than RNA-seq, especially for genes expressed at low levels (Wang et al., 2014).

Here, we investigate what we argue is a distinct state in the *C. elegans* life cycle. Although *C. elegans* hermaphrodites emerge into adulthood replete with sperm, after about 3 days of egg-laying the animals become sperm-depleted and can only reproduce by mating. This marks a transition into what we define as the endogenous

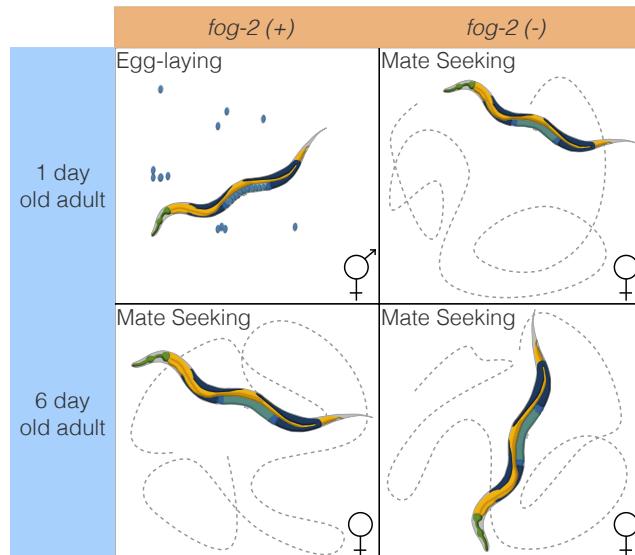
female-like state. This state is behaviorally distinguished by increased male-mating 2258 success (Garcia et al., 2007), which may be due to an increased attractiveness 2259 to males (Morsci, Haas, and Barr, 2011). This increased attractiveness acts at 2260 least partially through production of volatile chemical cues (Leighton et al., 2014). 2261 These behavioral changes are also coincident with functional deterioration of the 2262 germline (Andux and Ellis, 2008), muscle (Herndon et al., 2002), intestine (McGee 2263 et al., 2011) and nervous system (J. Liu et al., 2013), changes traditionally attributed 2264 to the aging process (T. R. Golden and Melov, 2007). 2265

To decouple the effects of aging and sperm-loss, we devised a two factor experiment. 2266 We examined wild-type XX animals at the beginning of adulthood (before worms 2267 contained embryos, referred to as 1st day adults) and after sperm depletion (6 days 2268 after the last molt, which we term 6th day adults). Second, we examined feminized 2269 XX animals that fail to produce sperm but are fully fertile if supplied sperm by mating 2270 with males (see Fig. 41). We used *fog-2* null mutants to obtain feminized animals. 2271 *fog-2* is involved in germ-cell sex determination in the hermaphrodite worm and 2272 is required for sperm production (Schedl and Kimble, 1988; Clifford et al., 2000). 2273 *C. elegans* defective in sperm formation will emerge from the larval stage as female 2274 adults. As time moves forward, these spermless worms only exhibit changes related 2275 to biological aging. As a result, *fog-2(lf)* mutants should show fewer gene changes 2276 during the first 6 days of adulthood compared to their egg-laying counterparts that 2277 age and also transition from egg-laying into a sperm depleted stage. 2278

Here, we show that we can detect a transcriptional signature associated with loss of 2279 hermaphroditic sperm marking entrance into the endogenous female-like state. We 2280 can also detect changes associated specifically with biological aging. Biological 2281 aging causes transcriptomic changes consisting of 5,592 genes in *C. elegans*. 4,552 2282 of these changes occur in both genotypes we studied, indicating they do not depend on 2283 sperm status. To facilitate exploration of the data, we have generated a website where 2284

we have deposited additional graphics, as well as all of the code used to generate these 2285 analyses: [https://wormlabcaltech.github.io/Angeles\\_Leighton\\_2016/](https://wormlabcaltech.github.io/Angeles_Leighton_2016/) 2286

2287



**Figure 41** Experimental design to identify genes associated with sperm loss and with aging. Studying the wild-type worm alone would measure time- and sperm-related changes at the same time, without allowing us to separate these changes. Studying the wild-type worm and a *fog-2(lf)* mutant would enable us to measure sperm-related changes but not time-related changes. By mixing both designs, we can measure and separate both modules.

## Materials and Methods

2288

### Strains

2289

Strains were grown at 20°C on NGM plates containing *E. coli* OP50. We used the 2290 laboratory *C. elegans* strain N2 as our wild-type strain (Sulston and Brenner, 1974). 2291 We also used the N2 mutant strain JK574, which contains the *fog-2(q71)* allele, for 2292 our experiments. 2293

### RNA extraction

2294

Synchronized worms were grown to either young adulthood or the 6th day of 2295 adulthood prior to RNA extraction. Synchronization and aging were carried out 2296

according to protocols described previously (Leighton et al., 2014). 1,000–5,000 worms from each replicate were rinsed into a microcentrifuge tube in S basal (5.85 g/L NaCl, 1 g/L K<sub>2</sub>HPO<sub>4</sub>, 6 g/L KH<sub>2</sub>PO<sub>4</sub>), and then spun down at 14,000 rpm for 30 s. The supernatant was removed and 1mL of TRIzol was added. Worms were lysed by vortexing for 30 s at room temperature and then 20 min at 4°. The TRIzol lysate was then spun down at 14,000 rpm for 10 min at 4°C to allow removal of insoluble materials. Thereafter the Ambion TRIzol protocol was followed to finish the RNA extraction (MAN0001271 Rev. Date: 13 Dec 2012). 3 biological replicates were obtained for each genotype and each time point.

2305

### RNA-Seq

2306

RNA integrity was assessed using RNA 6000 Pico Kit for Bioanalyzer (Agilent Technologies #5067–1513) and mRNA was isolated using NEBNext Poly(A) mRNA Magnetic Isolation Module (New England Biolabs, NEB, #E7490). RNA-Seq libraries were constructed using NEBNext Ultra RNA Library Prep Kit for Illumina (NEB #E7530) following manufacturer's instructions. Briefly, mRNA isolated from ~ 1 µg of total RNA was fragmented to the average size of 200 nt by incubating at 94°C for 15 min in first strand buffer, cDNA was synthesized using random primers and ProtoScript II Reverse Transcriptase followed by second strand synthesis using Second Strand Synthesis Enzyme Mix (NEB). Resulting DNA fragments were end-repaired, dA tailed and ligated to NEBNext hairpin adaptors (NEB #E7335). After ligation, adaptors were converted to the 'Y' shape by treating with USER enzyme and DNA fragments were size selected using Agencourt AMPure XP beads (Beckman Coulter #A63880) to generate fragment sizes between 250 and 350 bp. Adaptor-ligated DNA was PCR amplified followed by AMPure XP bead clean up. Libraries were quantified with Qubit dsDNA HS Kit (ThermoFisher Scientific #Q32854) and the size distribution was confirmed with High Sensitivity DNA Kit

2321

2322

for Bioanalyzer (Agilent Technologies #5067–4626). Libraries were sequenced on 2323 Illumina HiSeq2500 in single read mode with the read length of 50nt following man- 2324 ufacturer's instructions. Base calls were performed with RTA 1.13.48.0 followed by 2325 conversion to FASTQ with bcl2fastq 1.8.4. 2326

## Statistical Analysis

### RNA-Seq Analysis.

RNA-Seq alignment was performed using Kallisto (Bray et al., 2016) with 200 bootstraps. Differential expression analysis was performed using Sleuth (Pimentel et al., 2016). The following General Linear Model (GLM) was fit:

$$\log(y_i) = \beta_{0,i} + \beta_{G,i} \cdot G + \\ \beta_{A,i} \cdot A + \beta_{A::G,i} \cdot A \cdot G,$$

where  $y_i$  are the TPM counts for the  $i$ th gene;  $\beta_{0,i}$  is the intercept for the  $i$ th gene;  $\beta_{X,i}$  2329 is the regression coefficient for variable  $X$  for the  $i$ th gene;  $A$  is a binary age variable 2330 indicating 1st day adult (0) or 6th day adult (1);  $G$  is the genotype variable indicating 2331 wild-type (0) or *fog-2(lf)* (1);  $\beta_{A::G,i}$  refers to the regression coefficient accounting 2332 for the interaction between the age and genotype variables in the  $i$ th gene. Genes 2333 were called significant if the FDR-adjusted q-value for any regression coefficient 2334 was less than 0.1. Our script for differential analysis is available on GitHub. 2335

Regression coefficients and TPM counts were processed using Python 3.5 in a 2336 Jupyter Notebook (Pérez and Granger, 2007). Data analysis was performed using 2337 the Pandas, NumPy and SciPy libraries (McKinney, 2011; Van Der Walt, Colbert, 2338 and Varoquaux, 2011; Oliphant, 2007). Graphics were created using the Matplotlib 2339 and Seaborn libraries (Waskom et al., 2016; Hunter, 2007). Interactive graphics 2340 were generated using Bokeh (Bokeh Development Team, 2014). 2341

Tissue, Phenotype and Gene Ontology Enrichment Analyses (TEA, PEA and 2342

GEA, respectively) were performed using the WormBase Enrichment Suite for 2343  
 Python (Angeles-Albores, N. Lee, et al., 2016; Angeles-Albores, Lee, et al., 2018). 2344  
 Briefly, the WormBase Enrichment Suite accepts a list of genes and identifies the 2345  
 terms to which these genes are annotated. Terms are annotated by frequency of 2346  
 occurrence, and the probability that a term appears at this frequency under random 2347  
 sampling is calculated using a hypergeometric probability distribution. The hyper- 2348  
 geometric probability distribution is extremely sensitive to deviations from the null 2349  
 distribution, which allows it to identify even small deviations from the null. 2350

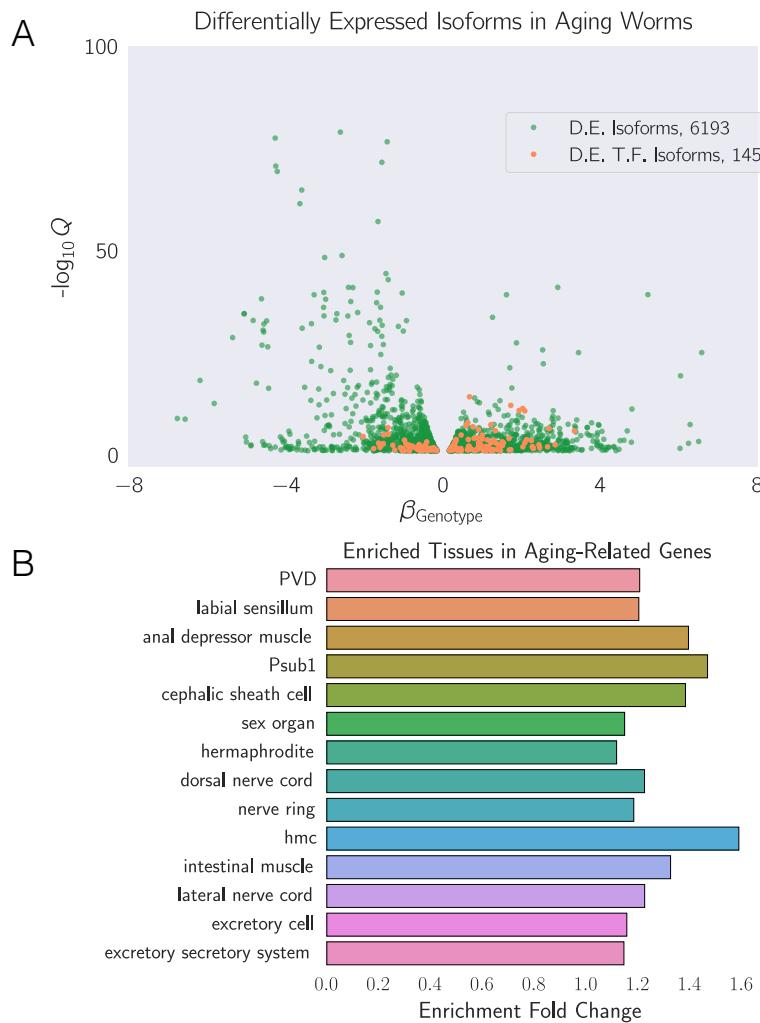
## Data Availability

Strains are available from the *Caenorhabditis* Genetics Center. All of the data and 2352  
 scripts pertinent for this project except the raw reads can be found on our Github 2353  
 repository [https://github.com/WormLabCaltech/Angeles\\_Leighton\\_2016](https://github.com/WormLabCaltech/Angeles_Leighton_2016). 2354  
 File S1 contains the list of genes that were altered in aging regardless of genotype. 2355  
 File S2 contains the list of genes and their associations with the *fog-2(lf)* pheno- 2356  
 type. File S3 contains genes associated with the female-like state. Raw reads were 2357  
 deposited to the Sequence Read Archive under the accession code SUB2457229. 2358

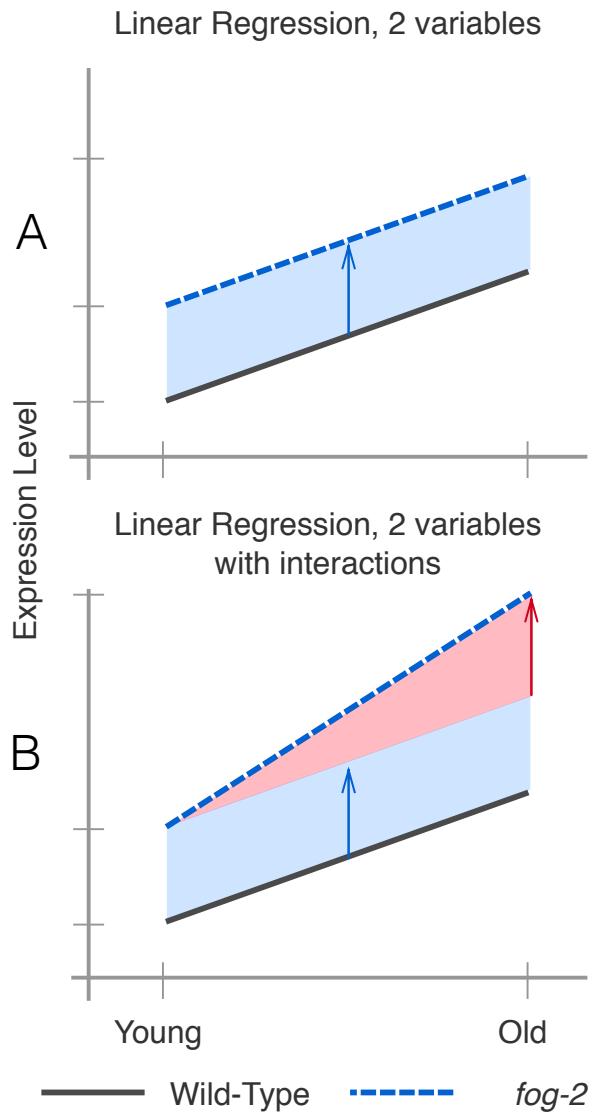
## Results and Discussion

**Decoupling time-dependent effects from sperm-status via general linear models** 2359  
 In order to decouple time-dependent effects from changes associated with loss of 2361  
 hermaphroditic sperm, we measured wild-type and *fog-2(lf)* adults at the 1st day 2362  
 adult stage (before visible embryos were present) and 6th day adult stage, when all 2363  
 wild-type hermaphrodites have laid all their eggs (see Fig 41), but mortality is still 2364  
 low (< 10%) (Stroustrup et al., 2013). We obtained 16–19 million reads mappable 2365  
 to the *C. elegans* genome per biological replicate, which enabled us to identify 2366  
 14,702 individual genes totalling 21,143 isoforms (see Figure 42a). 2367

One way to analyze the data from this two-factor design is by pairwise comparison 2368



**Figure 42 A.** Differentially expressed isoforms in the aging category. We identified a common aging expression signature between N2 and *fog-2(lf)* animals, consisting of 6,193 differentially expressed isoforms totaling 5,592 genes. The volcano plot is randomly down-sampled 30% for ease of viewing. Each point represents an individual isoform.  $\beta_{\text{Aging}}$  is the regression coefficient. Larger magnitudes of  $\beta$  indicate a larger log-fold change. The y-axis shows the negative logarithm of the q-values for each point. Green points are differentially expressed isoforms; orange points are differentially expressed isoforms of predicted transcription factor genes (Reece-Hoyes et al., 2005). An interactive version of this graph can be found on our [website](#). **B.** Enriched tissues in aging-associated genes. Tissue Enrichment Analysis (Angeles-Albores, N. Lee, et al., 2016) showed that genes associated with muscle tissues and the nervous system are enriched in aging-related genes. Only statistically significantly enriched tissues are shown. Enrichment Fold Change is defined as *Observed/Expected*. hmc stands for head mesodermal cell.



**Figure 43** Explanation of linear regressions with and without interactions. **A.** A linear regression with two variables, age and genotype. The expression level of a hypothetical gene increases by the same amount as worms age regardless of genotype. However, *fog-2(lf)* has higher expression of this gene than the wild-type at all stages (blue arrow). **B.** A linear regression with two variables and an interaction term. In this example, the expression level of this hypothetical gene is different between wild-type worms and *fog-2(lf)* (blue arrow). Although the expression level of this gene increases with age, the slope is different between wild-type and *fog-2(lf)*. The difference in the slope can be accounted for through an interaction coefficient (red arrow).

of the distinct states. However, such an analysis would not make full use of all the 2369 statistical power afforded by this experiment. Another method that makes full use of 2370 the information in our experiment is to perform a linear regression in 3 dimensions 2371 (2 independent variables, age and genotype, and 1 output). A linear regression with 2372 1 parameter (age, for example) would fit a line between expression data for young 2373 and old animals. When a second parameter is added to the linear regression, said 2374 parameter can be visualized as altering the y-intercept, but not the slope, of the first 2375 line in question (see Fig. 43a). 2376

Although a simple linear model is oftentimes useful, sometimes it is not appropriate 2377 to assume that the two variables under study are entirely independent. For example, 2378 in our case, three out of the four timepoint-and-genotype combinations we studied 2379 did not have sperm, and sperm-status is associated with both the *fog-2(lf)* self-sterile 2380 phenotype and with biological age of the wild-type animal. One way to statistically 2381 model such correlation between variables is to add an interaction term to the linear 2382 regression. This interaction term allows extra flexibility in describing how changes 2383 occur between conditions. For example, suppose a given theoretical gene *X* has 2384 expression levels that increase in a *fog-2*-dependent manner, but also increases in an 2385 age-dependent manner. However, aged *fog-2(lf)* animals do not have the expression 2386 levels of *X* that would be expected from adding the effect of the two perturbations; 2387 instead, the expression levels of *X* in this animal are considerably above what is 2388 expected. In this case, we could add a positive interaction coefficient to the model to 2389 explain the effect of genotype on the y-intercept as well as the slope (see Fig. 43b). 2390 When the two perturbations affect a single genetic pathway, these interactions can 2391 be interpreted as epistatic interactions. 2392

For these reasons, we used a general linear model with interactions to identify a 2393 transcriptomic profile associated with the *fog-2(lf)* genotype independently of age, 2394 as well as a transcriptomic profile of *C. elegans* aging common to both genotypes. 2395

The change associated with each variable is referred as  $\beta$ ; this number, although related to the natural logarithm of the fold change, is not equal to it. However, it is true that larger magnitudes of  $\beta$  indicate greater change. Thus, for each gene we performed a linear regression, and we evaluated the whether the  $\beta$  values associated with each coefficient were significantly different from 0 via a Wald test corrected for multiple hypothesis testing. A coefficient was considered to be significantly different from 0 if the q-value associated with it was less than 0.1.

2402

### **A quarter of all genes change expression between the 1st day of adulthood and the 6th day of adulthood in *C. elegans***

2404

We identified a transcriptomic signature consisting of 5,592 genes that were differentially expressed in 6th day adult animals of either genotype relative to 1st day adult animals (see S1). This constitutes more than one quarter of the genes in *C. elegans*. Tissue Enrichment Analysis (TEA) (Angeles-Albores, N. Lee, et al., 2016) showed that nervous tissues including the ‘nerve ring’, ‘dorsal nerve cord’, ‘PVD’ and ‘labial sensillum’ were enriched in genes that become differentially expressed through aging. Likewise, certain muscle groups (‘anal depressor muscle’, ‘intestinal muscle’) were enriched. (see Figure 42b). Gene Enrichment Analysis (GEA) (Angeles-Albores, Lee, et al., 2018) revealed that genes that were differentially expressed during the course of aging were enriched in terms involving respiration (‘respiratory chain’, ‘oxoacid metabolic process’); translation (‘cytosolic large ribosomal subunit’); and nucleotide metabolism (‘purine nucleotide’, ‘nucleoside phosphate’ and ‘ribose phosphate’ metabolic process). Phenotype Enrichment Analysis (PEA) (Angeles-Albores, Lee, et al., 2018) showed this gene list was associated with phenotypes that affect the *C. elegans* gonad, including ‘gonad vesiculated’, ‘gonad small’, ‘oocytes lack nucleus’ and ‘rachis narrow’.

2420

To verify the quality of our dataset, we generated a list of 1,056 golden standard

2421

genes expected to be altered in 6th day adult worms using previous literature reports 2422  
 including downstream genes of *daf-12*, *daf-16*, and aging and lifespan extension 2423  
 datasets (Murphy et al., 2003; Halaschek-Wiener et al., 2005; Lund et al., 2002; 2424  
 McCormick et al., 2012; Eckley et al., 2013). Of 1,056 standard genes, we found 2425  
 506 genes in our time-responsive dataset. This result was statistically significant 2426  
 with a p-value <  $10^{-38}$ . 2427

Next, we used a published compendium (Reece-Hoyes et al., 2005) to search for 2428  
 known or predicted transcription factors. We found 145 transcription factors in the 2429  
 set of genes with differential expression in aging nematodes. We subjected this list 2430  
 of transcription factors to TEA to understand their expression patterns. 6 of these 2431  
 transcription factors were expressed in the ‘hermaphrodite specific neuron’ (HSN), 2432  
 a neuron physiologically relevant for egg-laying (*hlh-14*, *sem-4*, *ceh-20*, *egl-46*, 2433  
*ceh-13*, *hlh-3*), which represented a statistically significant 2-fold enrichment of this 2434  
 tissue ( $q < 10^{-1}$ ). The term ‘head muscle’ was also overrepresented at twice the 2435  
 expected level ( $q < 10^{-1}$ , 13 genes). 2436

### The whole-organism *fog-2(lf)* differential expression signature 2437

We identified 1,881 genes associated with the *fog-2(lf)* genotype, including 60 tran- 2438  
 scription factors (see S2). TEA showed that the terms ‘AB’, ‘somatic gonad’, ‘uterine 2439  
 muscle’, ‘cephalic sheath cell’, ‘spermathecal-uterine junction’, and ‘PWD’ were en- 2440  
 riched in this gene set. The ‘somatic gonad’ and ‘spermathecal-uterine junction’ 2441  
 are both near the site of action of *fog-2(lf)* (the germline) and possibly reflect phys- 2442  
 iological changes from a lack of sperm. Phenotype ontology enrichment analysis 2443  
 showed that only a single phenotype term, ‘spindle orientation variant’ was enriched 2444  
 in the *fog-2(lf)* transcriptional signature ( $q < 10^{-1}$ , 38 genes, 2-fold enrichment). 2445  
 Most genes annotated as ‘spindle orientation variant’ were slightly upregulated, 2446  
 and therefore are unlikely to uniquely reflect reduced germline proliferation. GO 2447

term enrichment was very similar to the aging gene set and reflected enrichment in annotations pertaining to translation and respiration. Unlike the aging gene set, the *fog-2(lf)* signature was significantly enriched in ‘myofibril’ and ‘G-protein coupled receptor binding’ ( $q < 10^{-1}$ ). Enrichment of the term ‘G-protein coupled receptor binding’ was due to 14 genes: *cam-1*, *mom-2*, *dsh-1*, *spp-10*, *fip-6*, *fip-7*, *fip-9*, *fip-13*, *fip-14*, *fip-18*, *K02A11.4*, *nlp-12*, *nlp-13*, and *nlp-40*. *dsh-1*, *mom-2* and *cam-1* are members of the Wnt signaling pathway. Most of these genes’ expression levels were up-regulated, suggesting increased G-protein binding activity in *fog-2(lf)* mutants.

2456

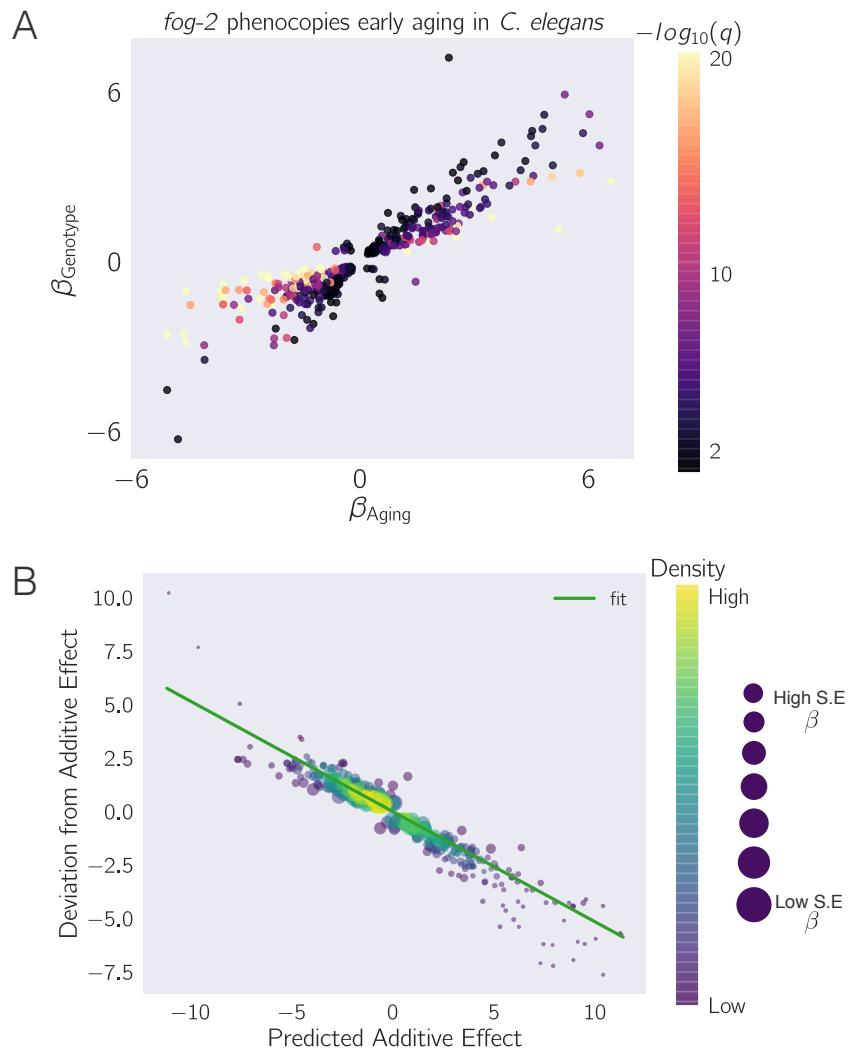
**The *fog-2(lf)* expression signature overlaps significantly with the aging signature**

Of the 1,881 genes that we identified in the *fog-2(lf)* signature, 1,040 genes were also identified in our aging set. Moreover, of these 1,040 genes, 905 genes changed in the same direction in response to either aging or germline feminization. The overlap between these signatures suggests an interplay between sperm-status and age. The nature of the interplay should be captured by the interaction coefficients in our model. There are four possibilities. First, the *fog-2(lf)* worms may have a fast-aging phenotype, in which case the interaction coefficients should match the sign of the aging coefficient. Second, the *fog-2(lf)* worms may have a slow-aging phenotype, in which case the interaction coefficients should have an interaction coefficient that is of opposite sign, but not greater in magnitude than the aging coefficient (if a gene increases in aging in a wild-type worm, it should still increase in a *fog-2(lf)* worm, albeit less). Third, the *fog-2(lf)* worms exhibit a rejuvenation phenotype. If this is the case, then these genes should have an interaction coefficient that is of opposite sign and greater magnitude than their aging coefficient, such that the change of these genes in *fog-2(lf)* mutant worms is reversed relative to the wild-type. Finally, if these genes are indicative of a female-like state, then these genes should not change with

age in *fog-2(lf)* animals, since these animals do not exit this state during the course 2474  
of the experiment. Moreover, because wild-type worms become female as they age, 2475  
a further requirement for a transcriptomic signature of the female-like state is that 2476  
aging coefficients for genes in this signature should have genotype coefficients of 2477  
equal sign and magnitude. In other words, entrance into the female-like state should 2478  
be not be path-dependent. 2479

To evaluate which of these possibilities was most likely, we selected the 1,040 genes 2480  
that had aging, genotype and interaction coefficients significantly different from zero 2481  
and we plotted their temporal coefficients against their genotype coefficients (see 2482  
Fig. 44a). We observed that the aging coefficients were strongly predictive of the 2483  
genotype coefficients. Most of these genes fell near the line  $y = x$ , suggesting that 2484  
these genes define a female-like state. 2485

We considered how loss-of-function of *fog-2* and aging could both interact to cause 2486  
entry into this state. We reasoned that a plausible mechanism is that *fog-2* pro- 2487  
motes sperm-production, and aging promotes sperm-depletion. This simple path- 2488  
way model suggests that a double perturbation consisting of aging and loss of 2489  
function of *fog-2* should show non-additivity of phenotypes (generalized epistasis). 2490  
To test whether these two perturbations deviate from additivity, we generated an 2491  
epistasis plot using this gene set. We have previously used epistasis plots to measure 2492  
transcriptome-wide epistasis between genes in a pathway (Angeles-Albores, Puckett 2493  
Robinson, et al., 2018). Briefly, an epistasis plot shows the expected expression of 2494  
a double perturbation under an additive model (null model) on the x-axis, and the 2495  
deviation from this null model in the y-axis. In other words, we calculated the 2496  
x-coordinates for each point by adding  $\beta_{\text{Genotype}} + \beta_{\text{Aging}}$ , and the y-coordinates are 2497  
equal to  $\beta_{\text{Interaction}}$  for each isoform. Previously we have shown that if two genes or 2498  
perturbations act within a linear pathway, an epistasis plot will generate a line with 2499  
slope equal to  $-0.5$ . When we generated an epistasis plot and found the line of best 2500



**Figure 44** *fog-2(lf)* partially phenocopies early aging in *C. elegans*. The  $\beta$  in each axes is the regression coefficient from the GLM, and can be loosely interpreted as an estimator of the log-fold change. Loss of *fog-2* is associated with a transcriptomic phenotype involving 1,881 genes. 1,040/1,881 of these genes are also altered in wild-type worms as they progress from young adulthood to old adulthood, and 905 change in the same direction. However, progression from young to old adulthood in a *fog-2(lf)* background results in no change in the expression level of these genes. **A.** We identified genes that change similarly during feminization and aging. The correlation between feminization and aging is almost 1:1. **B.** Epistasis plot of aging versus feminization. Epistasis plots indicate whether two genes (or perturbations) act on the same pathway. When two effects act on the same pathway, this is reflected by a slope of  $-0.5$ . The measured slope was  $-0.51 \pm 0.01$ .

fit, we observed a slope of  $-0.51 \pm 0.01$ , which suggests that the *fog-2* gene and time 2501 are acting to generate a single transcriptomic phenotype along a single pathway. 2502 Overall, we identified 405 genes that changed in the same direction through age or 2503 mutation of the *fog-2(lf)* gene and that had an interaction coefficient of opposite sign 2504 to the aging or genotype coefficient (see S3). Taken together, these observations 2505 suggests that these 405 genes define a female-like state in *C. elegans*. 2506

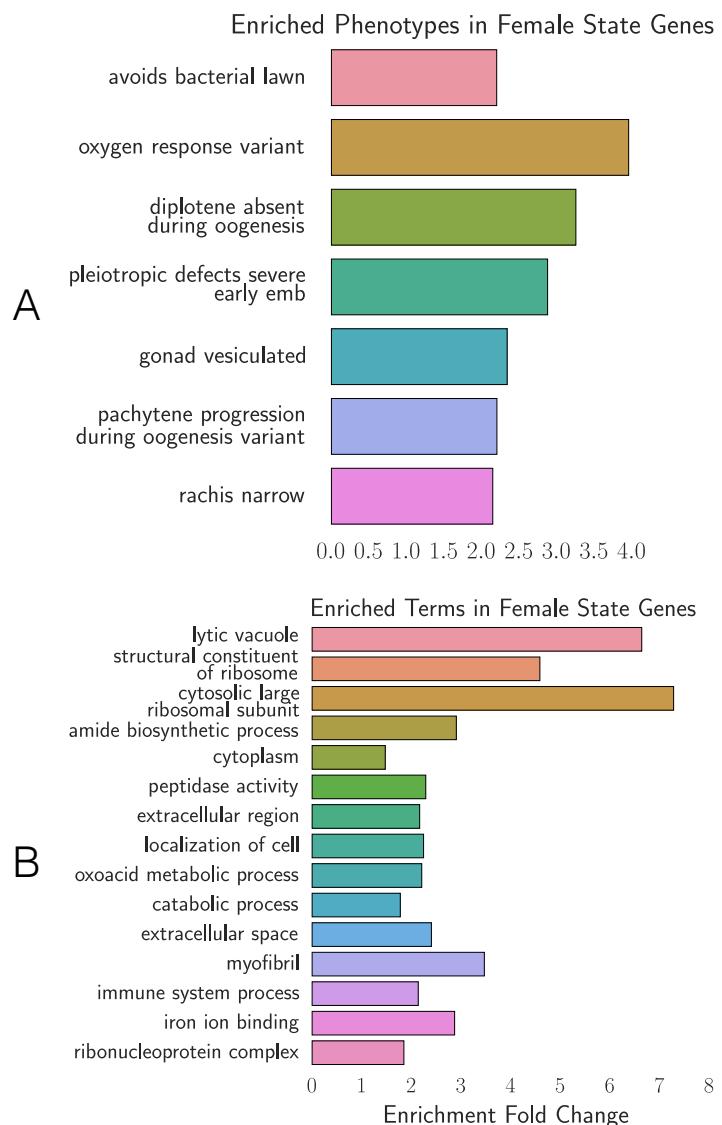
### **Analysis of the female-like state expression signature**

To better understand the changes that happen after sperm loss, we performed tissue 2508 enrichment, phenotype enrichment and gene ontology enrichment analyses on the 2509 set of 405 genes that we associated with the female-like state (see Fig. 45). TEA 2510 showed no tissue enrichment using this gene-set. GEA showed that this gene list 2511 was enriched in constituents of the ribosomal subunits almost four times above 2512 background ( $q < 10^{-5}$ , 17 genes). The enrichment of ribosomal constituents 2513 in this gene set in turn drives the enriched phenotypes: ‘avoids bacterial lawn’, 2514 ‘diplotene absent during oogenesis’, ‘gonad vesiculated’, ‘pachytene progression 2515 during oogenesis variant’, and ‘rachis narrow’. The expression of most of these 2516 ribosomal subunits is down-regulated in aged animals or in *fog-2(lf)* mutants. 2517

## **Discussion**

### **Defining an Early Aging Phenotype**

Our experimental design enables us to decouple the effects of egg-laying from aging. 2520 As a result, we identified a set of almost 4,000 genes that are altered similarly be- 2521 tween wild-type and *fog-2(lf)* mutants. Due to the read depth of our transcriptomic 2522 data (20 million reads) and the number of samples measured (3 biological replicates 2523 for 4 different life stages/genotypes), this dataset constitutes a high-quality descrip- 2524 tion of the transcriptomic changes that occur in aging populations of *C. elegans*. 2525 Although our data only capture ~ 50% of the expression changes reported in earlier 2526



**Figure 45** Phenotype and GO enrichment of genes involved in the female-like state. **A.** Phenotype Enrichment Analysis. **B.** Gene Ontology Enrichment Analysis. Most of the terms enriched in PEA reflect the abundance of ribosomal subunits present in this gene set.

aging transcriptome literature, this disagreement can be explained by a difference in methodology; earlier publications typically addressed the aging of fertile wild-type hermaphrodites only indirectly, or queried aging animals at a much later stage of their life cycle.

2530

### General linear models enable epistasis measurements

2531

We set out to study the self-fertilizing (hermaphroditic) to self-sterile (female-like) transition by comparing wild-type animals with *fog-2(lf)* mutants as they aged. Our computational approach enabled us to separate between two biological processes that are correlated within samples. Because of this intra-sample correlation, identifying this state via pairwise comparisons would not have been straightforward. Although it is a favored method amongst biologists, such pairwise comparisons suffer from a number of drawbacks. First, pairwise comparisons are unable to draw on the full statistical power available to an experiment because they discard almost all information except the samples being compared. Second, pairwise comparisons require a researcher to define *a priori* which comparisons are informative. For experiments with many variables, the number of pairwise combinations is explosively large. Indeed, even for this two-factor experiment, there are 6 possible pairwise comparisons. On the other hand, by specifying a linear regression model, each gene can be summarized with three variables, each of which can be analyzed and understood without the need to resort to further pairwise combinations.

2546

### The *C. elegans* female-like state

2547

Our explorations have shown that the loss of *fog-2(lf)* partially phenocopies the transcriptional events that occur naturally as *C. elegans* ages from the 1st day of adulthood to the 6th day of adulthood. Moreover, epistasis analysis of these perturbations suggests that they act on the same pathway, namely sperm generation and depletion (see Fig. 46). Self-sperm generation promotes the hermaphrodite

2552

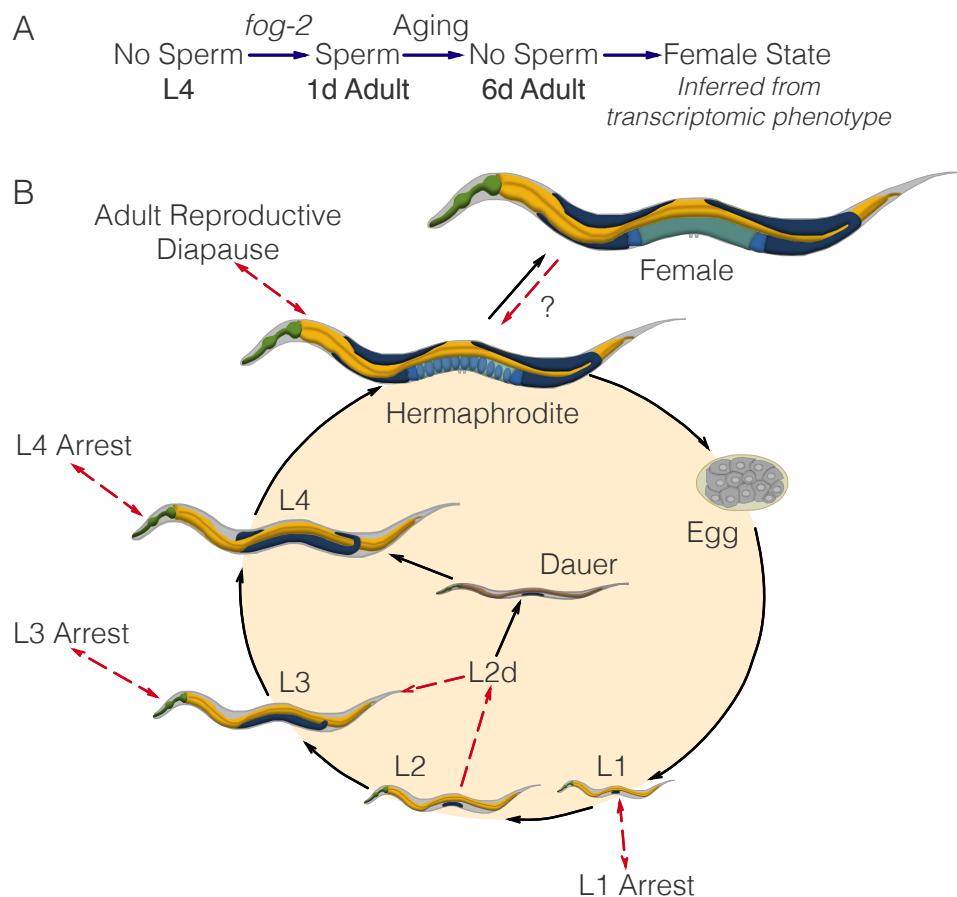
state, whereas sperm depletion marks entry into the female-like state. Given the 2553 enrichment of neuronal transcription factors that are associated with sperm loss 2554 in our dataset, we believe this dataset should contain some of the transcriptomic 2555 modules that are involved in these pheromone production and behavioral pathways, 2556 although we have been unable to find these genes. 2557

Behavioral and physiological changes upon mating are not unknown in other species. 2558 In particular, in the fruit fly *Drosophila melanogaster*, sex peptide present in the 2559 male seminal fluid is known to drive changes in gene expression (H. Liu and Kubli, 2560 2003; Xue and Noll, 2000; Avila et al., 2011; Heifetz et al., 2014; Rezával et al., 2561 2014; Mack et al., 2006) as well as behavior. More recently, sperm was found to be 2562 necessary to drive changes in aggression in the fruit fly (Bath et al., 2017). These 2563 changes are often reversible upon the disappearance of seminal fluid or sperm. In 2564 the case of *C. elegans*, we have observed that sperm loss is associated with gene 2565 expression changes that probably reflect physiological changes in the worm. Our 2566 experimental design did not include a test for reversibility of these changes. The 2567 possibility of a rescue experiment with males raises interesting possibilities: What 2568 fraction of the changes observed upon loss of self-sperm are reversible? Do male 2569 seminal fluid or male sperm cause changes beyond rescue? 2570

### **The *C. elegans* life cycle, life stages and life states**

2571

*C. elegans* has a complicated life cycle, with two alternative developmental pathways 2572 that have multiple stages (larval development and dauer development), followed by 2573 reproductive adulthood. In addition to its developmental stages, researchers have 2574 recognized that *C. elegans* has numerous life states that it can enter into when 2575 given instructive environmental cues. One such state is the L1 arrest state, where 2576 development ceases entirely upon starvation (Johnson et al., 1984; Baugh and 2577 Sternberg, 2006). More recently, researchers have described additional diapause 2578



**Figure 46 A.** A substrate-dependent model showing how *fog-2* promotes sperm generation, whereas aging promotes sperm depletion, leading to entry to the female-like state. Such a model can explain why *fog-2* and aging appear epistatic to each other. **B.** The complete *C. elegans* life cycle. Recognized stages of *C. elegans* are marked by black arrows. States are marked by red arrows to emphasize that at the end of a state, the worm returns to the developmental timepoint it was at before entering the state. The L2d state is an exception. It is the only stage that does not return to the same developmental timepoint; rather, the L2d state is a permissive state that allows entry into either dauer or the L3 stage. We have presented evidence of a female-like state in *C. elegans*. At this point, it is unclear whether the difference between hermaphrodites and females is reversible by males. Therefore, it remains unclear whether it is a stage or a true state.

states that the worm can access at the L3, L4 and young adult stages under conditions 2579 of low food (Angelo and Gilst, 2009; Seidel and Kimble, 2011; Schindler, Baugh, 2580 and Sherwood, 2014). Not all states of *C. elegans* are arrested, however (see Fig. 46). 2581 For example, the L2d state is induced by crowded and nutrient poor conditions (J. W. 2582 Golden and Riddle, 1984). While within this state, the worm is capable of entry 2583 into either dauer or the L3 larval stage, depending on environmental conditions. 2584 Thus, the L2d state is a permissive state, and marks the point at which the nematode 2585 development is committed to a single developmental pathway. 2586

Identification of the *C. elegans* life states has often been performed by morphological 2587 studies (as in the course of L4 arrest or L2d) or via timecourses (L1 arrest). However, 2588 not all states may be visually identifiable, or even if they are, the morphological 2589 changes may be very subtle, making positive identification difficult. However, the 2590 detailed information afforded by a transcriptome should in theory provide suffi- 2591 cient information to definitively identify a state, since transcriptomic information 2592 underlies morphology. Moreover, transcriptomics can provide an insight into the 2593 physiology of complex metazoan life states. By identifying differentially expressed 2594 genes and using ontology enrichment analyses to identify gene functions, sites of 2595 expression or phenotypes that are enriched in a given gene set, we can obtain a 2596 clear picture of the changes that occur in the worm analogous to identifying gross 2597 morphological changes. 2598

RNA-seq is a powerful technology that has been used successfully in the past as 2599 a qualitative tool for target acquisition, though recent work has successfully used 2600 RNA-seq to measure genetic interactions via epistasis (Dixit et al., 2016; Angeles- 2601 Albores, Puckett Robinson, et al., 2018). Here, we have shown that whole-organism 2602 RNA-seq data can also be used to successfully identify internal states in a multi- 2603 cellular organism. 2604

- References** 2605
- Andux, Sara and Ronald E. Ellis (2008). “Apoptosis maintains oocyte quality in aging *Caenorhabditis elegans* females”. In: *PLoS Genetics* 4.12. ISSN: 15537390. doi: [10.1371/journal.pgen.1000295](https://doi.org/10.1371/journal.pgen.1000295). 2606  
2607
- Angeles-Albores, David, Raymond YN Lee, et al. (2018). “Two new functions in the WormBase Enrichment Suite”. In: *Micropublication: biology. Dataset*. DOI: <https://doi.org/10.17912/W25Q2N>. 2609  
2610  
2611
- Angeles-Albores, David, Raymond Y. N. Lee, et al. (2016). “Tissue enrichment analysis for *C. elegans* genomics”. In: *BMC Bioinformatics* 17.1, p. 366. ISSN: 1471-2105. doi: [10.1186/s12859-016-1229-9](https://doi.org/10.1186/s12859-016-1229-9). 2612  
2613  
2614
- Angeles-Albores, David, Carmie Puckett Robinson, et al. (Mar. 2018). “Reconstructing a metazoan genetic pathway with transcriptome-wide epistasis measurements.” In: *Proceedings of the National Academy of Sciences of the United States of America* 115.13, E2930–E2939. ISSN: 1091-6490. doi: [10.1073/pnas.1712387115](https://doi.org/10.1073/pnas.1712387115). 2615  
2616  
2617  
2618  
2619
- Angelo, Giana and Marc R Van Gilst (2009). “Cells and Extends Reproductive”. In: *Science* 326.November, pp. 954–958. 2620  
2621
- Avila, Frank W. et al. (Jan. 2011). “Insect Seminal Fluid Proteins: Identification and Function”. In: *Annual Review of Entomology* 56.1, pp. 21–40. ISSN: 0066-4170. doi: [10.1146/annurev-ento-120709-144823](https://doi.org/10.1146/annurev-ento-120709-144823). 2622  
2623  
2624
- Bath, Eleanor et al. (May 2017). “Sperm and sex peptide stimulate aggression in female *Drosophila*”. In: *Nature Ecology & Evolution* 1.6, p. 0154. ISSN: 2397-334X. doi: [10.1038/s41559-017-0154](https://doi.org/10.1038/s41559-017-0154). 2625  
2626  
2627
- Baugh, L. Ryan and Paul W. Sternberg (2006). “DAF-16/FOXO Regulates Transcription of *cki-1/Cip/Kip* and Repression of *lin-4* during *C. elegans* L1 Arrest”. 2628  
2629
- Blaxter, M. et al. (2012). “Genomics and transcriptomics across the diversity of the Nematoda”. In: *Parasite Immunology* 34.2-3, pp. 108–120. ISSN: 01419838. doi: [10.1111/j.1365-3024.2011.01342.x](https://doi.org/10.1111/j.1365-3024.2011.01342.x). 2630  
2631  
2632
- Boeck, Max E et al. (2016). “The time-resolved transcriptome of *C. elegans*”. In: *Genome Research*, pp. 1–10. ISSN: 15495469. doi: [10.1101/gr.202663.115](https://doi.org/10.1101/gr.202663.115). 2633  
2634  
2635
- Bokeh Development Team (2014). “Bokeh: Python library for interactive visualization”. In: 2636  
2637
- Bray, Nicolas L et al. (2016). “Near-optimal probabilistic RNA-seq quantification.” In: *Nature biotechnology* 34.5, pp. 525–7. ISSN: 1546-1696. doi: [10.1038/nbt.3519](https://doi.org/10.1038/nbt.3519). arXiv: [1505.02710](https://arxiv.org/abs/1505.02710). 2638  
2639  
2640

- Clifford, Robert et al. (2000). “FOG-2, a novel F-box containing protein, associates with the GLD-1 RNA binding protein and directs male sex determination in the *C. elegans* hermaphrodite germline.” In: *Development (Cambridge, England)* 127.24, pp. 5265–5276. ISSN: 0950-1991.
- Corsi, Ann K., Bruce Wightman, and Martin Chalfie (2015). “A transparent window into biology: A primer on *Caenorhabditis elegans*”. In: *Genetics* 200.2, pp. 387–407. ISSN: 19432631. DOI: [10.1534/genetics.115.176099](https://doi.org/10.1534/genetics.115.176099).
- Dixit, Atray et al. (2016). “Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens”. In: *Cell* 167.7, 1853–1866.e17. ISSN: 00928674. DOI: [10.1016/j.cell.2016.11.038](https://doi.org/10.1016/j.cell.2016.11.038).
- Eckley, D. Mark et al. (2013). “Molecular characterization of the transition to mid-life in *Caenorhabditis elegans*”. In: *Age* 35.3, pp. 689–703. ISSN: 01619152. DOI: [10.1007/s11357-012-9401-2](https://doi.org/10.1007/s11357-012-9401-2).
- Garcia, Hernan G. et al. (2007). “A First Exposure to Statistical Mechanics for Life Scientists”. In: p. 27. ISSN: 0036-8075. arXiv: [0708.1899](https://arxiv.org/abs/0708.1899).
- Gerstein, Mark B. et al. (2014). “Comparative analysis of the transcriptome across distant species”. In: *Nature* 512, pp. 445–448. ISSN: 0028-0836. DOI: [10.1038/nature13424](https://doi.org/10.1038/nature13424). arXiv: [NIHMS150003](https://arxiv.org/abs/NIHMS150003).
- Golden, James W. and Donald L. Riddle (1984). “The *Caenorhabditis elegans* dauer larva: Developmental effects of pheromone, food, and temperature”. In: *Developmental Biology* 102.2, pp. 368–378. ISSN: 00121606. DOI: [10.1016/0012-1606\(84\)90201-X](https://doi.org/10.1016/0012-1606(84)90201-X).
- Golden, Tamara R and Simon Melov (2007). “Gene expression changes associated with aging in *C. elegans*.” In: *WormBook : the online review of C. elegans biology*, pp. 1–12. ISSN: 1551-8507. DOI: [10.1895/wormbook.1.127.2](https://doi.org/10.1895/wormbook.1.127.2).
- Halaschek-Wiener, Julius et al. (2005). “Analysis of long-lived *C. elegans* daf-2 mutants using serial analysis of gene expression”. In: *Genome Research*, pp. 603–615. DOI: [10.1101/gr.3274805..](https://doi.org/10.1101/gr.3274805..)
- Heifetz, Yael et al. (Mar. 2014). “Mating Regulates Neuromodulator Ensembles at Nerve Termini Innervating the *Drosophila* Reproductive Tract”. In: *Current Biology* 24.7, pp. 731–737. ISSN: 09609822. DOI: [10.1016/j.cub.2014.02.042](https://doi.org/10.1016/j.cub.2014.02.042).
- Herndon, Laura a et al. (2002). “Stochastic and genetic factors influence tissue-specific decline in ageing *C. elegans*.” In: *Nature* 419.6909, pp. 808–814. ISSN: 0028-0836. DOI: [10.1038/nature01135](https://doi.org/10.1038/nature01135).
- Hillier, Ladeana W. et al. (2009). “Massively parallel sequencing of the polyadenylated transcriptome of *C. elegans*”. In: *Genome Research* 19.4, pp. 657–666. ISSN: 10889051. DOI: [10.1101/gr.088112.108](https://doi.org/10.1101/gr.088112.108).

- Hunter, John D. (2007). "Matplotlib: A 2D graphics environment". In: *Computing in Science and Engineering* 9.3, pp. 99–104. ISSN: 15219615. doi: [10.1109/MCSE.2007.55](https://doi.org/10.1109/MCSE.2007.55). arXiv: [0402594v3 \[arXiv:cond-mat\]](https://arxiv.org/abs/0402594v3). 2679  
2680  
2681
- Johnson, Thomas E. et al. (1984). "Arresting development arrests aging in the nematode *Caenorhabditis elegans*". In: *Mechanisms of Ageing and Development* 28.1, pp. 23–40. ISSN: 00476374. doi: [10.1016/0047-6374\(84\)90150-7](https://doi.org/10.1016/0047-6374(84)90150-7). 2682  
2683  
2684
- Leighton, Daniel H. W. et al. (2014). "Communication between oocytes and somatic cells regulates volatile pheromone production in *Caenorhabditis elegans*". In: *Proceedings of the National Academy of Sciences* 111.50, pp. 17905–17910. ISSN: 1091-6490. doi: [10.1073/pnas.1420439111](https://doi.org/10.1073/pnas.1420439111). 2685  
2686  
2687  
2688
- Liu, Huanfa and Eric Kubli (Aug. 2003). "Sex-peptide is the molecular basis of the sperm effect in *Drosophila melanogaster*." In: *Proceedings of the National Academy of Sciences of the United States of America* 100.17, pp. 9929–33. ISSN: 0027-8424. doi: [10.1073/pnas.1631700100](https://doi.org/10.1073/pnas.1631700100). 2689  
2690  
2691  
2692
- Liu, Jie et al. (2013). "Functional aging in the nervous system contributes to age-dependent motor activity decline in *C. elegans*". In: *Cell Metabolism* 18.3, pp. 392–402. ISSN: 15504131. doi: [10.1016/j.cmet.2013.08.007](https://doi.org/10.1016/j.cmet.2013.08.007). arXiv: [NIHMS150003](https://arxiv.org/abs/NIHMS150003). 2693  
2694  
2695  
2696
- Lund, James et al. (2002). "Transcriptional profile of aging in *C. elegans*". In: *Current Biology* 12.18, pp. 1566–1573. ISSN: 09609822. doi: [10.1016/S0960-9822\(02\)01146-6](https://doi.org/10.1016/S0960-9822(02)01146-6). 2697  
2698  
2699
- Mack, Paul D et al. (July 2006). "Mating-responsive genes in reproductive tissues of female *Drosophila melanogaster*." In: *Proceedings of the National Academy of Sciences of the United States of America* 103.27, pp. 10358–63. ISSN: 0027-8424. doi: [10.1073/pnas.0604046103](https://doi.org/10.1073/pnas.0604046103). 2700  
2701  
2702  
2703
- Magalhães, Jp De, Ce Finch, and G Janssens (2010). "Next-generation sequencing in aging research: emerging applications, problems, pitfalls and possible solutions". In: *Ageing research reviews* 9.3, pp. 315–323. ISSN: 1872-9649. doi: [10.1016/j.arr.2009.10.006](https://doi.org/10.1016/j.arr.2009.10.006). Next-generation. 2704  
2705  
2706  
2707
- McCormick, Mark et al. (2012). "New genes that extend *Caenorhabditis elegans*' lifespan in response to reproductive signals". In: *Aging Cell* 11.2, pp. 192–202. ISSN: 14749718. doi: [10.1111/j.1474-9726.2011.00768.x](https://doi.org/10.1111/j.1474-9726.2011.00768.x). arXiv: [NIHMS150003](https://arxiv.org/abs/NIHMS150003). 2708  
2709  
2710  
2711
- McGee, Matthew D. et al. (2011). "Loss of intestinal nuclei and intestinal integrity in aging *C. elegans*". In: *Aging Cell* 10.4, pp. 699–710. ISSN: 14749718. doi: [10.1111/j.1474-9726.2011.00713.x](https://doi.org/10.1111/j.1474-9726.2011.00713.x). 2712  
2713  
2714
- McKinney, Wes (2011). "pandas: a Foundational Python Library for Data Analysis and Statistics". In: *Python for High Performance and Scientific Computing*, pp. 1–9. 2715  
2716  
2717

- Morsci, Natalia S., Leonard A. Haas, and Maureen M. Barr (2011). “Sperm status regulates sexual attraction in *Caenorhabditis elegans*”. In: *Genetics* 189.4, pp. 1341–1346. ISSN: 00166731. DOI: [10.1534/genetics.111.133603](https://doi.org/10.1534/genetics.111.133603). 2718  
2719  
2720
- Mortazavi, Ali et al. (2008). “Mapping and quantifying mammalian transcriptomes by RNA-Seq”. In: *Nature Methods* 5.7, pp. 621–628. ISSN: 1548-7091. DOI: [10.1038/nmeth.1226](https://doi.org/10.1038/nmeth.1226). arXiv: [1111.6189v1](https://arxiv.org/abs/1111.6189v1). 2721  
2722  
2723
- Murphy, Coleen T. et al. (2003). “Genes that act downstream of DAF-16 to influence the lifespan of *Caenorhabditis elegans*.” In: *Nature* 424.6946, pp. 277–283. ISSN: 00280836. DOI: [10.1038/nature01789](https://doi.org/10.1038/nature01789). 2724  
2725  
2726
- Murray, John Isaac et al. (2012). “Multidimensional regulation of gene expression in the *C. elegans* embryo”. In: pp. 1282–1294. ISSN: 1088-9051. DOI: [10.1101/gr.131920.111](https://doi.org/10.1101/gr.131920.111). 2727  
2728  
2729
- Oliphant, Travis E (2007). “SciPy: Open source scientific tools for Python”. In: *Computing in Science and Engineering* 9, pp. 10–20. ISSN: 1521-9615. DOI: [10.1109/MCSE.2007.58](https://doi.org/10.1109/MCSE.2007.58). 2730  
2731  
2732
- Pérez, F. and B.E. Granger (2007). “IPython: A System for Interactive Scientific Computing Python: An Open and General-Purpose Environment”. In: *Computing in Science and Engineering* 9.3, pp. 21–29. ISSN: 15219615. DOI: [doi:10.1109/MCSE.2007.53..](https://doi.org/10.1109/MCSE.2007.53..). 2733  
2734  
2735  
2736
- Pimentel, Harold J et al. (2016). “Differential analysis of RNA-Seq incorporating quantification uncertainty”. In: *bioRxiv*, p. 058164. DOI: [10.1101/058164](https://doi.org/10.1101/058164). 2737  
2738
- Rangaraju, Sunitha et al. (2015). “Suppression of transcriptional drift extends *C. elegans* lifespan by postponing the onset of mortality”. In: *eLife* 4.December2015, pp. 1–39. ISSN: 2050084X. DOI: [10.7554/eLife.08833](https://doi.org/10.7554/eLife.08833). 2739  
2740  
2741
- Reece-Hoyes, John S. et al. (2005). “A compendium of *Caenorhabditis elegans* regulatory transcription factors: a resource for mapping transcription regulatory networks.” In: *Genome biology* 6.13, R110. ISSN: 1474-760X. DOI: [10.1186/gb-2005-6-13-r110](https://doi.org/10.1186/gb-2005-6-13-r110). 2742  
2743  
2744  
2745
- Rezával, Carolina et al. (Mar. 2014). “Sexually Dimorphic Octopaminergic Neurons Modulate Female Postmating Behaviors in *Drosophila*”. In: *Current Biology* 24.7, pp. 725–730. ISSN: 09609822. DOI: [10.1016/j.cub.2013.12.051](https://doi.org/10.1016/j.cub.2013.12.051). 2746  
2747  
2748
- Schedl, Tim and Judith Kimble (1988). “fog-2, a germ-line-specific sex determination gene required for hermaphrodite spermatogenesis in *Caenorhabditis elegans*.” In: *Genetics* 119.1, pp. 43–61. ISSN: 00166731. 2749  
2750  
2751
- Schindler, Adam J., L. Ryan Baugh, and David R. Sherwood (2014). “Identification of Late Larval Stage Developmental Checkpoints in *Caenorhabditis elegans* Regulated by Insulin/IGF and Steroid Hormone Signaling Pathways”. In: *PLoS Genetics* 10.6, pp. 13–16. ISSN: 15537404. DOI: [10.1371/journal.pgen.1004426](https://doi.org/10.1371/journal.pgen.1004426). 2752  
2753  
2754  
2755

- Seidel, Hannah S. and Judith Kimble (2011). “The oogenic germline starvation response in *C. elegans*”. In: *PLoS ONE* 6.12. issn: 19326203. doi: [10.1371/journal.pone.0028074](https://doi.org/10.1371/journal.pone.0028074). 2756  
2757  
2758
- Stroustrup, Nicholas et al. (2013). “The *Caenorhabditis elegans* Lifespan Machine.” In: *Nature methods* 10.7, pp. 665–70. issn: 1548-7105. doi: [10.1038/nmeth.2475](https://doi.org/10.1038/nmeth.2475). arXiv: [NIHMS150003](https://arxiv.org/abs/NIHMS150003). 2759  
2760  
2761
- Sulston, J. E. and S. Brenner (1974). “The DNA of *Caenorhabditis elegans*.” In: *Genetics* 77.1, pp. 95–104. issn: 00166731. 2762  
2763
- Sulston, J. E. and H. R. Horvitz (Mar. 1977). “Post-embryonic cell lineages of the nematode, *Caenorhabditis elegans*”. In: *Developmental Biology* 56.1, pp. 110–156. issn: 00121606. doi: [10.1016/0012-1606\(77\)90158-0](https://doi.org/10.1016/0012-1606(77)90158-0). 2764  
2765  
2766
- Sulston, J. E., E. Schierenberg, et al. (Nov. 1983). “The embryonic cell lineage of the nematode *Caenorhabditis elegans*”. In: *Developmental Biology* 100.1, pp. 64–119. issn: 00121606. doi: [10.1016/0012-1606\(83\)90201-4](https://doi.org/10.1016/0012-1606(83)90201-4). 2767  
2768  
2769
- Van Der Walt, Stéfan, S. Chris Colbert, and Gaël Varoquaux (2011). “The NumPy array: A structure for efficient numerical computation”. In: *Computing in Science and Engineering* 13.2, pp. 22–30. issn: 15219615. doi: [10.1109/MCSE.2011.37](https://doi.org/10.1109/MCSE.2011.37). arXiv: [1102.1523](https://arxiv.org/abs/1102.1523). 2770  
2771  
2772  
2773
- Wang, Charles et al. (2014). “The concordance between RNA-seq and microarray data depends on chemical treatment and transcript abundance.” In: *Nature biotechnology* 32.9, pp. 926–32. issn: 1546-1696. doi: [10.1038/nbt.3001](https://doi.org/10.1038/nbt.3001). arXiv: [NIHMS150003](https://arxiv.org/abs/NIHMS150003). 2774  
2775  
2776  
2777
- Waskom, Michael et al. (2016). “seaborn: v0.7.0 (January 2016)”. In: doi: [10.5281/zenodo.45133](https://doi.org/10.5281/zenodo.45133). 2778  
2779
- Xue, L. and M. Noll (Mar. 2000). “*Drosophila* female sexual behavior induced by sterile males showing copulation complementation”. In: *Proceedings of the National Academy of Sciences* 97.7, pp. 3272–3275. issn: 0027-8424. doi: [10.1073/pnas.97.7.3272](https://doi.org/10.1073/pnas.97.7.3272). 2780  
2781  
2782  
2783

USING TRANSCRIPTOMES AS MUTANT PHENOTYPES REVEALS FUNCTIONAL REGIONS OF A MEDIATOR SUBUNIT IN <i>C. ELEGANS</i>	2785
	2786
	2787

<b>Abstract</b>	2788
-----------------	------

Although transcriptomes have recently been used as phenotypes with which 2789 to perform epistasis analyses, they are not yet used to study intragenic func- 2790 tion/structure relationships. We developed a theoretical framework to study 2791 allelic series using transcriptomic phenotypes. As a proof-of-concept, we apply 2792 our methods to an allelic series of *dpy-22*, a highly pleiotropic *Caenorhabditis el-* 2793 *elegans* gene orthologous to the human gene *MED12*, which encodes a subunit 2794 of the Mediator complex. Our methods identify functional units within *dpy-22* 2795 that modulate Mediator activity upon various genetic programs, including the 2796 Wnt and Ras modules. 2797

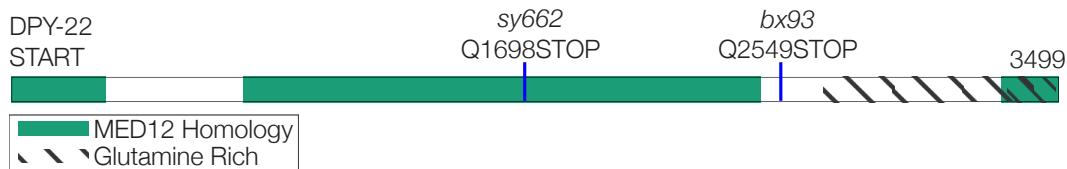
<b>Introduction</b>	2798
---------------------	------

Mutations of a gene can yield a series of alleles with different phenotypes that 2799 reveal multiple functions encoded by that gene, regardless of the alleles' molecular 2800 nature. In *Caenorhabditis elegans*, allelic series have characterized genes such as 2801 *let-23/EGFR*, *lin-3/EGF* and *lin-12/NOTCH* (Aroian and Paul W Sternberg, 1991; 2802 Ferguson and Horvitz, 1985; Greenwald, Paul W. Sternberg, and Robert Horvitz, 2803 1983). Allelic series provide a way to probe genes where biochemical approaches 2804 would be difficult, slow or uninformative with regards to the biological phenomenon 2805 of interest. Their power derives from the ability to draw broad conclusions about 2806 the gene of interest in terms of gene dosage and functional units, to the extent 2807 that these two factors are separable, without regard to the molecular identity of the 2808

mutations that created these alleles. Here, gene dosage is defined as the combined 2809 effects of transcriptional and translational expression, gene product localization, 2810 and biochemical kinetics of the final gene product *in situ*. To study allelic series, 2811 we must first enumerate the phenotypes each allele affects, and subsequently order 2812 the alleles into severity and dominance hierarchies per phenotype. The resulting 2813 hierarchies enable us to better understand how a given gene, which may be highly 2814 pleiotropic, can give rise to highly specific mutant phenotypes when mutated in just 2815 the right way. 2816

Biology has moved from expression measurements of single genes towards genome- 2817 wide measurements. Expression profiling via RNA-seq (Mortazavi et al., 2008) 2818 enables simultaneous measurement of transcript levels for all genes in a genome, 2819 yielding a transcriptome. These measurements can be made on whole organisms, 2820 isolated tissues, or single cells (Tang et al., 2009; Schwarz, Kato, and Paul W. 2821 Sternberg, 2012). Transcriptomes have been successfully used to identify new cell 2822 or organismal states (Angeles-Albores, Leighton, et al., 2017; Villani et al., 2017). 2823 Transcriptomic states can be used to perform epistatic analyses (Dixit et al., 2016; 2824 Angeles-Albores, Puckett Robinson, Brian A Williams, et al., 2018), but have not 2825 been used to characterize allelic series. 2826

We have devised methods for characterizing allelic series using RNA-seq. To test 2827 these methods, we selected three alleles (Zhang and Emmons, 2000; Moghal and 2828 Paul W. Sternberg, 2003) of a *C. elegans* Mediator complex subunit gene, *dpy-22*. 2829 Mediator is a macromolecular complex with ~ 25 subunits (Jeronimo and Robert, 2830 2017) that globally regulates RNA polymerase II (Pol II) (Allen and Taatjes, 2015; 2831 Takagi and Kornberg, 2006). The Mediator complex has at least four biochemi- 2832 cally distinct modules: the Head, Middle and Tail modules and a CDK-8-associated 2833 Kinase Module (CKM). The CKM associates reversibly with other modules, and 2834 appears to inhibit transcription (Knuesel et al., 2009; Elmlund et al., 2006). In *C. el-* 2835



**Figure 51** Protein sequence schematic for DPY-22. The positions of the nonsense mutations used are shown.

*egans* development, the CKM promotes the formation of the male tail (Zhang and Emmons, 2000) (through interactions with the Wnt pathway), as well as formation of the hermaphrodite vulva (Moghal and Paul W. Sternberg, 2003) (through inhibition of the Ras pathway). Null alleles of *dpy-22* are likely to be lethal, based on embryonic lethal phenotypes observed after RNAi (Wang et al., 2004; Lehner et al., 2006) and the severe phenotypes of a strong *dpy-22* hypomorphic allele, *dpy-22(e652)* (homozygous hermaphrodites are very sick) (Riddle et al., 1997). Homozygotes of allele *dpy-22(bx93)*, which encodes a premature stop codon Q2549Amber (Zhang and Emmons, 2000), appear grossly wild-type, though this allele does not have complete wild-type functionality, since it fails to fully complement the Muv phenotype of another allele, *sy622*, in a sensitized *let-23* background. In contrast, animals homozygous for a more severe allele, *dpy-22(sy622)* encoding another premature stop codon, Q1698Amber (Moghal and Paul W. Sternberg, 2003), are dumpy (Dpy), have egg-laying defects (Egl), and have multiple vulvae (Muv) (Fig. 51). In humans, MED12 is known to have a proline-, glutamine- and leucine-rich domain that interacts with the WNT pathway (Kim et al., 2006). However, many disease-causing variants fall outside of this domain (Yamamoto and Shimojima, 2015). In spite of its causative role in a number of neurodevelopmental disorders (Graham and Schwartz, 2013), the structural and functional features of this gene are poorly understood, partially because genetic approaches towards studying pleiotropic genes have proved difficult in the past, highlighting the need for new methods.

2856

<b>Methods</b>	2857
<b>Strains used</b>	2858
Strains used were N2 wild-type (Bristol) (Sulston and Brenner, 1974), PS4087 <i>dpy-22(sy622)</i> (Moghal and Paul W. Sternberg, 2003), PS4187 <i>dpy-22(bx93)</i> (Zhang and Emmons, 2000), PS4176 <i>dpy-6(e14) dpy-22(bx93)/+ dpy-22(sy622)</i> (Moghal and Paul W. Sternberg, 2003), MT4866 <i>let-60(n2021)</i> (Beitel, Clark, and Horvitz, 1990), MT2124 <i>let-60(n1046gf)</i> (Beitel, Clark, and Horvitz, 1990) and EW15 <i>bar-1(ga80)</i> (Eisenmann et al., 1998). Lines were grown on standard nematode growth media (NGM) Petri plates seeded with OP50 <i>E. coli</i> at 20°C (Sulston and Brenner, 1974).	2859 2860 2861 2862 2863 2864 2865 2866
<b>Strain synchronization, harvesting and RNA sequencing</b>	2867
With the exception of strain MT4866, strains were synchronized by bleaching P <sub>0</sub> 's into virgin S. basal (no cholesterol or ethanol added) for 16–18 hours. Arrested L1 larvae were placed in NGM plates seeded with OP50 at 20°C and grown to the young adult stage (assessed by vulval morphology and lack of embryos). We discovered that MT4866 dies upon L1 starvation for this period of time. As a result, we syn- chronized this strain by double bleaching. Animals were picked if they were young adults, regardless of whether any vulval or morphological phenotypes were present. RNA extraction and sequencing was performed as previously described by Angeles- Albores, Puckett Robinson, Brian A Williams, et al. (2018) and Angeles-Albores, Leighton, et al. (2017). Briefly, young adults were placed in 10 μL of TE buffer, and digested using Recombinant Proteinase K PCR Grade (Roche Lot 656 No. 03115 838001) incubated with 1% SDS 657 and 1.25 μL RNA Secure (Ambion AM7005). Total RNA was extracted using the Zymo Research Directzol RNA MicroPrep Kit (Zymo Research, SKU R2061). mRNA was subsequently purified using a NEBNext Poly(A) mRNA Magnetic Isolation Module (New England Biolabs, NEB, #E7490).	2868 2869 2870 2871 2872 2873 2874 2875 2876 2877 2878 2879 2880 2881 2882

Sequencing libraries were generated using the NEBNext Ultra RNA Library Prep 2883  
Kit for Illumina (NEB #E7530). These libraries were sequenced using an Illumina 2884  
HiSeq2500 machine in single-read mode with a read length of 50 nucleotides. 2885

### Read pseudo-alignment and differential expression 2886

Reads were pseudo-aligned to the *C. elegans* genome (WBcel235) using Kallisto (Bray 2887  
et al., 2016), using 200 bootstraps and with the sequence bias (-seqBias) flag. The 2888  
fragment size for all libraries was set to 200 and the standard deviation to 40. Quality 2889  
control was performed on a subset of the reads using FastQC, RNAseQC, BowTie 2890  
and MultiQC (Andrews, 2010; Deluca et al., 2012; Langmead et al., 2009; Ewels 2891  
et al., 2016). 2892

Differential expression analysis was performed using Sleuth (Pimentel et al., 2017). 2893  
We used a general linear model to identify genes that were differentially expressed 2894  
between wild-type and mutant libraries. To increase our statistical power, we pooled 2895  
young adult wild-type replicates from other published (Angeles-Albores, Puckett 2896  
Robinson, Brian A Williams, et al., 2018; Angeles-Albores, Leighton, et al., 2017) 2897  
and unpublished analyses adjusting for batch effects. Briefly, batch effects were 2898  
controlled by including the identity of the person who collected the worms and the 2899  
method by which the libraries were generated as covariates. 2900

### False hit analysis 2901

To accurately count phenotypes, we developed a false hit algorithm (Algorithm 1). 2902  
We implemented this algorithm for comparisons of three genotypes using Python. 2903  
Such an experiment can result in 128 possible combinations of phenotypic classes 2904  
(ignoring size). This large number of models necessitates an algorithmic approach 2905  
that can restrict the number of models. Our algorithm uses a noise function that 2906  
assumes false hit events are non-overlapping (i.e. the same gene cannot be the result 2907  
of two false positive events in two or more genotypes) to determine the average noise 2908

flux between phenotypic classes. These assumptions break down if false-positive or 2909 negative rates are large (>25%). 2910

To benchmark our algorithm, we generated one thousand Venn diagrams at ran- 2911 dom. For each Venn diagram, we calculated the average false positive and false 2912 negative flux matrices. Then, we added noise to each phenotypic class in the Venn 2913 diagram, assuming that fluxes were normally distributed with mean and standard 2914 deviation equal to the flux coefficient calculated. We input the noised Venn diagram 2915 into our false hit analysis and collected classification statistics. For a given signal- 2916 to-noise cutoff,  $\lambda$ , classification accuracy varied significantly with changes in the 2917 total error rate. In the absence of false negative hits, false hit analysis can accu- 2918 rately identify non-empty genotype-associated phenotypic classes, but identifying 2919 genotype-specific classes becomes difficult if the experimental false positive rate is 2920 high. On the other hand, even moderate false negative rates (> 10%) rapidly de- 2921 grade signal from genotype-associated classes. For classes that are associated with 2922 three genotypes, an experimental false negative rate of 30% is enough on average to 2923 prevents this class from being observed. 2924

We selected  $\lambda = 3$  because classification using this threshold was high across a 2925 range of false positive and false negative combinations. A challenge to applying 2926 this algorithm to our data is the fact that the false negative rate for our experiment is 2927 unknown. Although there has been significant progress in controlling and estimating 2928 false positive rates, we know of no such attempts for false negative rates. It is unlikely 2929 that the false negative rate for our study is lower than the false positive rate, because 2930 all genotypes except the controls are likely underpowered. We used false negative 2931 rates between 10–20% for false hit analysis. All analyses returned the same final 2932 model. 2933

We asked whether re-classification of some classes into others could improve model 2934

fit. We manually re-classified the (*dpy-22(sy622)*,*dpy-22(bx93)*)-associated and the 2935 (*dpy-22(bx93)*, *trans-heterozygote*)-associated classes into the *bx93*-associated class 2936 (which is associated with all genotypes), and compared  $\chi^2$  statistics between a re- 2937 classified reduced model ( $\chi^2 = 72$ ) and a reduced model ( $\chi^2 = 130$ ). Based on the 2938 lower  $\chi^2$  of the re-classified reduced model, we concluded that it is the most likely 2939 model given our data. 2940

**Algorithm 1** False Hit Algorithm. Briefly, the algorithm initializes a reduced model with the phenotypic class or classes labelled by the largest number of genotypes. This reduced model is used to estimate noise fluxes, which in turn can be used to estimate a signal-to-noise metric between observed and modelled classes. Classes that exhibit a high signal-to-noise are incorporated into the reduced model.

**Data:**  $\mathbf{M}_{obs} = \{N_l\}$ , an observed set of classes, where each class is labelled by  $l \in L$  and is of size  $N_l$ .  $f_p, f_n$ , the false positive and negative rates respectively.  $\alpha$ , the signal-to-noise threshold for acceptance of a class.

**Result:**  $\mathbf{M}_{reduced}$ , a reduced model that fits the data.

```

begin
  Define a minimal model,  $\mathbf{K}$ 
  Refine the model until convergence or iterations max out
   $i \leftarrow 0$ 
   $\mathbf{K}_{prev} \leftarrow \emptyset$ 
  while ( $i < i_{max}$ ) | ( $\mathbf{K}_{prev} \neq \mathbf{K}$ ) do
     $\mathbf{K}_{prev} \leftarrow \mathbf{K}$ 
    Define a noise function to estimate error flows in  $\mathbf{K}$   $\mathbf{F} \leftarrow \text{noise}(\mathbf{K}, f_p, f_n)$ 
    for  $l \in L$  do
      Calculate signal to noise for each labelled class False negatives can
      result in  $\lambda < 0$   $\lambda_l \leftarrow \mathbf{M}_{obs,l}/F_l$  if ( $\lambda > \alpha$ ) | ( $\lambda < 0$ ) then
        |  $\mathbf{K}_l \leftarrow \mathbf{M}_{obs,l}$ 
      end
    end
     $i++$ 
  end
end
 $\mathbf{M}_{reduced} = \mathbf{K}$ 
return  $\mathbf{M}_{reduced}$ 
```

---

**Dominance analysis**

2941

We modeled allelic dominance as a weighted average of allelic activity:

2942

$$\beta_{a/b,i,\text{Pred}}(d_a) = d_a \cdot \beta_{a/a,i} + (1 - d_a) \cdot \beta_{b/b,i}, \quad (5.1)$$

where  $\beta_{k/k,i}$  refers to the  $\beta$  value of the  $i$ th isoform in a genotype  $k/k$ , and  $d_a$  is the dominance coefficient for allele  $a$ .

2943

To find the parameters  $d_a$  that maximized the probability of observing the data, we found the parameter,  $d_a$ , that maximized the equation:

2945

$$P(d_a | D, H, I) \propto \prod_{i \in S} \exp -\frac{(\beta_{a/b,i,\text{Obs}} - \beta_{a/b,i,\text{Pred}}(d_a))^2}{2\sigma_i^2} \quad (5.2)$$

where  $\beta_{a/b,i,\text{Obs}}$  was the coefficient associated with the  $i$ th isoform in the *trans*-het  $a/b$  and  $\sigma_i$  was the standard error of the  $i$ th isoform in the *trans*-heterozygote samples as output by Kallisto.  $S$  is the set of isoforms that participate in the regression (see main text). This equation describes a linear regression which was solved numerically.

2950

**Code**

2951

Code was written in Jupyter notebooks (Pérez and Granger, 2007) using the Python programming language. The Numpy, pandas and scipy libraries were used for computation (Van Der Walt, Colbert, and Varoquaux, 2011; McKinney, 2011; Oliphant, 2007) and the matplotlib and seaborn libraries were used for data visualization (Hunter, 2007; Waskom et al., 2016). Enrichment analyses were performed using the WormBase Enrichment Suite (Angeles-Albores, N. Lee, et al., 2016; Angeles-Albores, Puckett Robinson, Brian A. Williams, et al., 2018). For all enrichment analyses, a  $q$ -value of less than  $10^{-3}$  was considered statistically significant. For gene ontology enrichment analysis, terms were considered statistically significant only if they also showed an enrichment fold-change greater than 2.

2961

**Data Availability**

2962

Raw and processed reads were deposited in the Gene Expression Omnibus. Scripts for the entire analysis can be found with version control in our Github repository, <https://github.com/WormLabCaltech/med-cafe>. A user-friendly, commented website containing the complete analyses can be found at <https://wormlabcaltech.github.io/med-cafe/>. Raw reads and quantified abundances for each sample were deposited at the NCBI Gene Expression Omnibus (GEO) (Edgar, Domrachev, and Lash, 2002) under the accession code GSE107523 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE107523>). 2963 2964 2965 2966 2967 2968 2969 2970

**Results**

2971

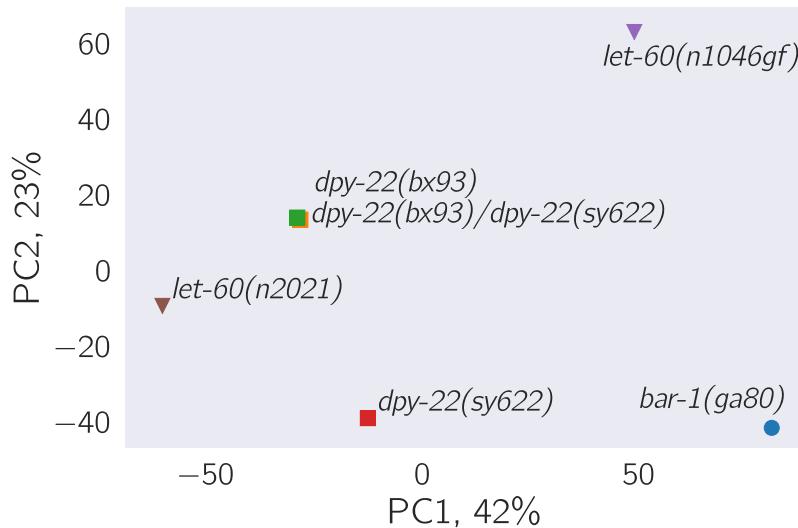
**RNA-sequencing of three *dpy-22* alleles and two known interactor genes**

2972

We carried out RNA-seq on biological triplicates of mRNA extracted from *dpy-22(sy622)* homozygotes, *dpy-22(bx93)* homozygotes, and wild type controls, along with quadruplicate from *trans*-heterozygotes of both alleles with the genotype *dpy-6(e14) dpy-22(bx93)/+ dpy-22(sy622)*. We also sequenced mRNA extracted from *bar-1(ga80)* (the  $\beta$ -catenin ortholog in *C. elegans*), *let-60(n2021)* and *let-60(n1046gf)* (the Ras ortholog in *C. elegans*) mutants in triplicate because these genes have been previously described to interact with *dpy-22* to form the vulva (Moghal and Paul W. Sternberg, 2003) and the male tail (Zhang and Emmons, 2000). Sequencing was performed at a depth of 20 million reads per sample. Reads were pseudoaligned using Kallisto (Bray et al., 2016). We performed a differential expression using a general linear model specified using Sleuth (Pimentel et al., 2017) (see Methods). Differential expression with respect to the wild type control for each transcript  $i$  in a genotype  $g$  is measured via a coefficient  $\beta_{g,i}$ , which can be loosely interpreted as the natural logarithm of the fold-change. Transcripts were considered to have differential expression between wild-type and a mutant if their false discovery rate,  $q$ , was less than 0.05.

Genotype	Differentially Expressed Genes
<i>dpy-22(bx93)</i>	266
<i>dpy-6(e14) dpy-22(bx93) / + dpy-22(sy622)</i>	2,128
<i>dpy-22(sy622)</i>	2,036
<i>bar-1(ga80)</i>	4613
<i>let-60(n2021)</i>	509
<i>let-60(n1046gf)</i>	2526

**Table 51** The number of differentially expressed genes relative to the wild-type control for each genotype with a significance threshold of 0.1.



**Figure 52** Principal component analysis of the analyzed genotypes. The analysis was performed using only those transcripts that were differentially expressed in at least one genotype. The plot shows that the *trans*-heterozygotes phenocopy the *dpy-22(bx93)* homozygotes along the first two principal dimensions.

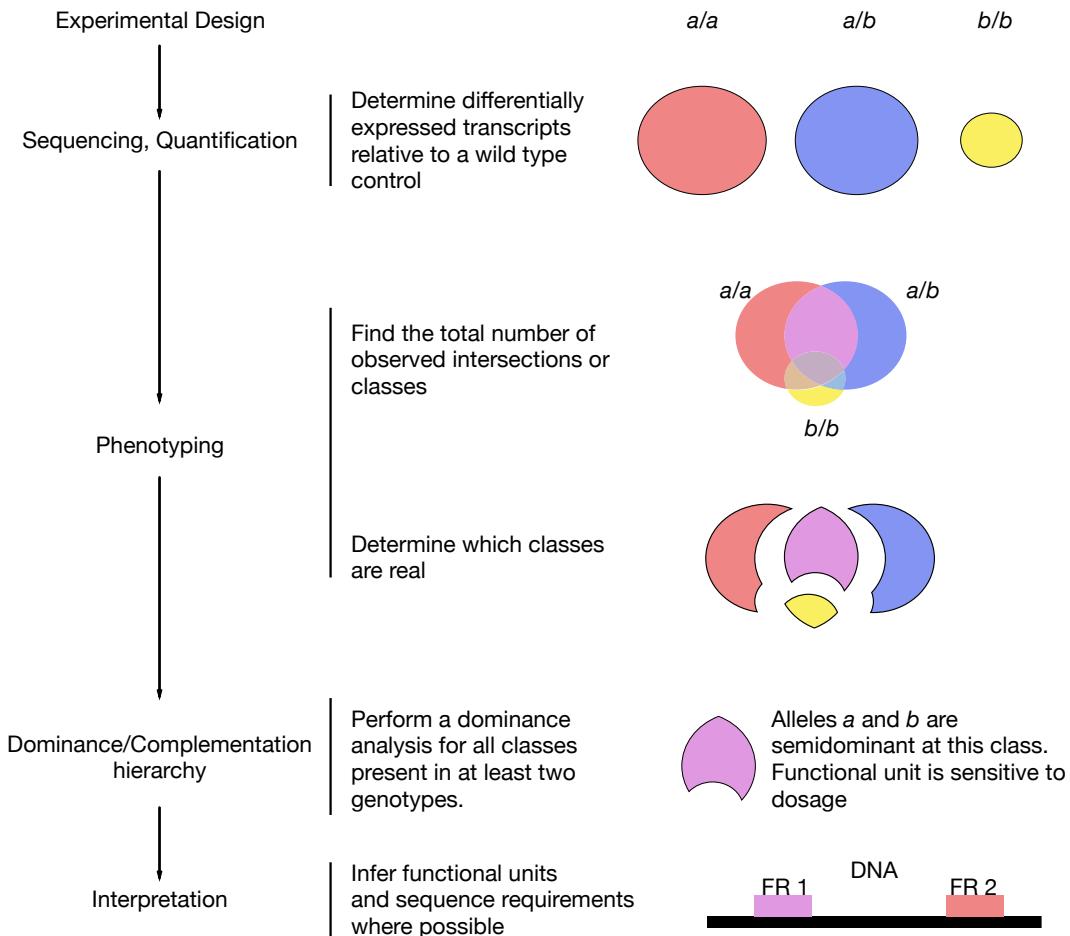
than or equal to 10%. We used this method to identify the differentially expressed genes associated with each mutant (Table 51; [Basic Statistics Notebook](#)) Supplementary File 1 contains all the beta values associated with this project. We have also generated a website containing complete details of all the analyses available at the following URL: <https://wormlabcaltech.github.io/med-cafe/analysis>.

**Principal component analysis visualizes the allelic dominance of the *dpy-22(bx93)*<sub>2993</sub> allele over *dpy-22(sy622)*** 2994

As a first step in our analysis, we performed dimensionality reduction on the transcriptomes we sequenced using Principal Component Analysis (PCA). Briefly, PCA identifies the vectors along which there is most variation in the data. These vectors can be used to project the data into lower dimensions to assess whether samples cluster, though interpreting the biological reasons for this clustering can be challenging. To perform PCA, we selected only those transcripts that were differentially expressed in at least one genotype, and used the  $\beta$  coefficients associated with these genes to perform PCA. Projecting the data into two dimensions maintains 65% of the variation. The first dimension separates the gain and loss of function *let-60* mutants. The second dimension separates the *dpy-22* mutants (Fig. 52). On the PCA plot, the *trans*-heterozygote mutants appear to phenocopy the *dpy-22(bx93)* mutants, recapitulating previous experiments that showed the *dpy-22(bx93)* allele to be dominant over the *dpy-22(sy622)* allele. 3007

**Three *dpy-22* genotypes have shared transcriptomic phenotypes** 3008

We would like to understand the degree and nature of the dominance between these *dpy-22* alleles. To construct a severity and dominance hierarchy, we must establish how many transcriptomic phenotypes are represented among the three *dpy-22* genotypes, and of those phenotypes, how many of them are shared transcriptomic phenotypes (STPs). Shared transcriptomic phenotypes are defined as the set of genes that are commonly differentially expressed in two mutant genotypes relative to a wild-type control, regardless of the direction of change, as defined previously (Angeles-Albores, Puckett Robinson, Brian A Williams, et al., 2018). We use the term in the plural version, because the shared genes may represent multiple independent modules that formally constitute different phenotypic classes. 3018



**Figure 53** Flowchart for an analysis of arbitrary allelic series. A set of alleles is selected, and the corresponding genotypes are sequenced. Independent phenotypic classes are identified, and classes that are the result of noise are discarded via a false hit analysis. For each phenotypic class, the alleles are ordered in a dominance/complementation hierarchy, which can then be used to infer functional regions (FR) within the genes in question.

We identified significant pairwise STPs between all *dpy-22* mutants. The transcripts that were differentially expressed in *dpy-22(bx93)* homozygotes were almost all differentially expressed in *dpy-22(sy622)* homozygotes (189/266) and in *trans-* heterozygotes (192/266). On the other hand, although *dpy-22(sy622)* homozygotes and *trans*-heterozygotes exhibited a similar number of differentially expressed genes, less than half of these were shared between the two genotypes.

3024

**False hit analysis identifies four non-overlapping phenotypic classes**

3025

Severity and dominance hierarchies must be calculated with respect to each independent phenotype associated with the alleles under study. A challenge with expression profiles is to identify these independent phenotypes. We reasoned that comparing the expression profiles of the two *dpy-22* homozygotes and the *trans*-heterozygote would naturally partition the expression profiles into groups that would constitute phenotypic classes. However, a three-way comparison can give rise to  $7 (2^3 - 1)$  possible groupings: transcripts perturbed in only a single genotype (3), transcripts perturbed in two genotypes (3) and transcripts perturbed in all three genotypes (1). A shortcoming of RNA-seq is that it is prone to false positive and false negative artifacts, and these artifacts could be numerous enough to cause the appearance of certain groups that would not be there otherwise. In other words, we might find a subset of genes that are differentially expressed in a single genotype, but if this subset is small enough, we ought to be concerned that this subset is caused by false positive hits within this genotype or false negative hits in the other genotypes. This thought experiment highlights the need to assess which groups have sufficient statistical support to consider as phenotypic classes.

3041

We developed a method to assess whether groups in a Venn diagram are likely to be the result of statistical artifacts. Briefly, the algorithm works by first assuming all of the data is the result of false positive and false negative hits except for the group of transcripts that is differentially expressed in most genotypes. Then, using estimates for the false positive and negative response, we calculate the expected sizes of all the groups after adding noise under this model. If an observed group is much larger than expected by noise, we refine the data model to accept the group. This process is iterated until the data model converges. We called this method a false hit analysis.

We used false hit analysis to identify four non-overlapping phenotypic classes (Fig. 53). We use the term genotype-specific to refer to groups of transcripts

3050

3051

that were perturbed in one mutant genotype. We use the term genotype-associated 3052 to refer to those groups of transcripts whose expression was significantly altered 3053 in two or more mutants genotypes with respect to the wild type control. The 3054 ***dpy-22(sy622)*-associated** phenotypic class consisted of 665 genes differentially ex- 3055 pressed in *dpy-22(sy622)* homozygotes and in *trans*-heterozygotes, but which had 3056 wild-type expression in *dpy-22(bx93)* homozygotes. The ***dpy-22(bx93)*-associated** 3057 phenotypic class contains 229 genes differentially expressed in all genotypes. The 3058 *dpy-22(bx93)*-associated class included re-classified transcripts that had been found 3059 to be differentially expressed in the *dpy-22(bx93)* homozygote and one other geno- 3060 type, because these were very likely to be the result of false negative hits in the 3061 missing genotype, and re-classifying these transcripts improved our model substan- 3062 tially. We also identified a ***dpy-22(sy622)*-specific** phenotypic class (1,213 genes) 3063 and a ***trans*-heterozygote-specific** phenotypic class (1,302 genes; see the [Phenotypic 3064 Classes Notebook](#)). 3065

### Severity hierarchy of a *dpy-22* allelic series

3066

Having separated the expression profiles into phenotypic classes, we can ask what 3067 the severity hierarchy is between the *dpy-22(bx93)* allele and the *dpy-22(sy622)* 3068 allele. Broadly speaking, there are two ways to assess severity. First, we can ask 3069 which allele causes more mutant phenotypes or phenotypic groups as a homozygote 3070 (**allelic pleiotropy**). Alternatively, we can identify the allele which causes the 3071 greatest change in expression in a homozygote at each shared phenotype among the 3072 homozygotes of both alleles, which we refer to as **allelic magnitude**. An important 3073 caveat is that magnitude only makes sense if the homozygotes of each allele are well 3074 correlated (i.e., they have a linear relationship with small spread). If the phenotypes 3075 have zero or negative correlation between two homozygotes, then the two alleles 3076 under inspection are not of the same kind, i.e., they cannot both be loss-of-function 3077

alleles or gain-of-function alleles for this phenotype, though the converse is not necessarily true.

3079

The *dpy-22(sy622)* homozygote shows more differentially expressed genes that participate in a greater number of phenotypic classes relative to the *dpy-22(bx93)* homozygote. Thus, the *dpy-22(sy622)* allele is a more pleiotropic mutation than the *dpy-22(bx93)* allele. Since the homozygotes of each allele only share a single phenotypic class in common, we need only assess magnitude along this single phenotype. To calculate a magnitude coefficient, for genes in the *dpy-22(bx93)*-associated phenotypic class, we plotted the  $\beta$  coefficients from the *dpy-22(sy622)* homozygote against the  $\beta$  coefficients from the *dpy-22(bx93)* homozygote (see Fig. 54) and performed a linear regression to find the slope of this line. Using this method, we found that the *dpy-22(bx93)* homozygote has a magnitude that is  $62\% \pm 2\%$  of the *dpy-22(sy622)* homozygote. Taken together, these results suggest that the *dpy-22(sy622)* allele represents a more severe alteration-of-function mutation than the mutation within the *dpy-22(bx93)* allele.

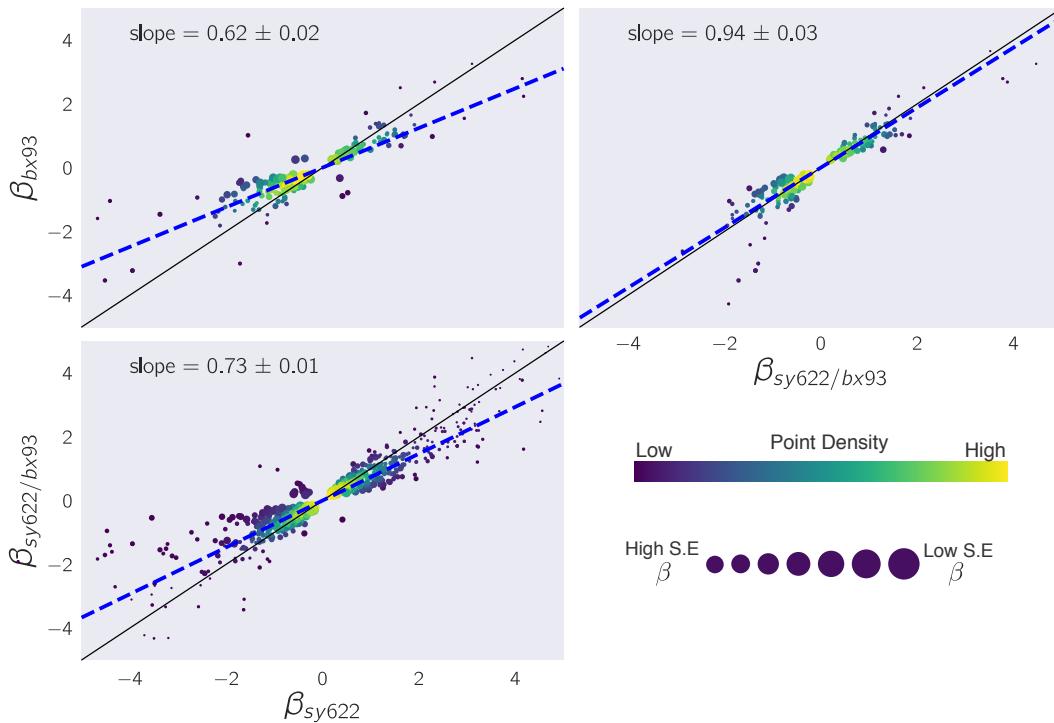
3092

### Dominance hierarchy of a *dpy-22* allelic series

3093

We measured allelic dominance for each class using a dominance coefficient (see [Methods](#)). The dominance coefficient is a measure of the contribution of each allele to the total expression level in *trans*-heterozygotes. By definition, the *dpy-22(sy622)* allele is completely recessive to *dpy-22(bx93)* for the *dpy-22(sy622)*-specific phenotypic class. To determine the dominance coefficient for the remaining phenotypic classes, we first selected the transcripts within those classes, and asked what linear combination of the homozygotic  $\beta$  coefficients best approximated the  $\beta$  coefficients of the *trans*-heterozygote, subject to the constraint that the sum of the weights for the two homozygotes should be equal to unity. We solved this problem by finding the maximum likelihood estimate for these weights. Using this method, we found

3103



**Figure 54** Shared Transcriptomic Phenotypes amongst the *dpy-22* genotypes are regulated in the same direction. For each pairwise comparison, we found those transcripts that were commonly differentially expressed in both genotypes relative to the wild-type control and plotted the  $\beta$  coefficients for each. We performed a linear regression on each plot to find the line of best fit (broken blue line). Only the comparison between *dpy-22(sy622)* and *dpy-22(bx93)* homozygotes was used to establish that the magnitude of the *dpy-22(sy622)* allele is greater than the magnitude of the *dpy-22(bx93)* allele. The other comparisons are shown for completeness.

that the *dpy-22(sy622)* and *dpy-22(bx93)* alleles are semidominant ( $d_{bx93} = 0.48$ )  
3104 to each other for the *dpy-22(sy622)*-associated phenotypic class. The *dpy-22(bx93)*  
3105 allele is largely dominant over the *dpy-22(sy622)* allele ( $d_{bx93} = 0.82$ ; see Table 52)  
3106 for the *dpy-22(bx93)*-associated phenotypic class.  
3107

### Phenotypic classes reflect morphological phenotypes

We performed enrichment analysis of anatomical, phenotypic and gene ontology  
3109 terms using the WormBase Enrichment Suite (Angeles-Albores, N. Lee, et al.,  
3110 2016; Angeles-Albores, Puckett Robinson, Brian A. Williams, et al., 2018). The  
3111

Phenotypic Class	Dominance
<i>dpy-22(sy622)</i> -specific	$1.00 \pm 0.00$
<i>dpy-22(sy622)</i> -associated	$0.48 \pm 0.01$
<i>dpy-22(bx93)</i> -associated	$0.82 \pm 0.01$

**Table 52** Dominance analysis for the *dpy-22/MDT12* allelic series. Dominance values closer to 1 indicate *dpy-22(bx93)* is dominant over *dpy-22(sy622)*, whereas 0 indicates *dpy-22(sy622)* is dominant over *dpy-22(bx93)*.

*dpy-22(bx93)*-associated phenotypic class was enriched in genes involved in ‘immune system processes’ ( $q < 10^{-5}$ ), and was enriched in genes expressed in the ‘intestine’ ( $q < 10^{-4}$ ). The *dpy-22(sy622)*-associated class was enriched in genes expressed in the ‘cephalic sheath cell’ ( $q < 10^{-4}$ ). Using ontology enrichment analysis from the WormBase Enrichment Suite, we found that the *dpy-22(sy622)*-associated class is enriched in histones and histone-like proteins (‘DNA packaging complex’  $q < 10^{-3}$ ) as well as genes involved in ‘immune system processes’ ( $q < 10^{-5}$ ). The *dpy-22(sy622)*-specific class was enriched in genes that have expression in the ‘intestine’ ( $q < 10^{-7}$ ), ‘muscular system’ ( $q < 10^{-3}$ ) and ‘epithelial system’ ( $q < 10^{-2}$ ). The genes in this class are known to cause bacterial lawn avoidance when knocked down or knocked out ( $q < 10^{-2}$ ). Finally, GO enrichment showed that the *dpy-22(sy622)*-specific class is specifically enriched in ‘structural constituents of cuticle’ ( $q < 10^{-12}$ ), and in genes involved in respiration ( $q < 10^{-6}$ ). This last result recapitulates the fact that *dpy-22(sy622)* homozygotes show a severe Dumpy phenotype. The *trans*-heterozygote specific class was enriched in genes expressed in ‘male’ animals ( $q < 10^{-63}$ ) and genes expressed in the ‘reproductive system’ ( $q < 10^{-21}$ ). GO enrichment of genes in the *trans*-heterozygote specific class showed enrichment of the genes involved in the ‘regulation of cell shape’ ( $q < 10^{-6}$ ) and in a variety of terms involving phosphate metabolism, such as ‘nucleoside phosphate binding’ ( $q < 10^{-5}$ ), ‘dephosphorylation’ ( $q < 10^{-3}$ ) or ‘phosphorylation’ ( $q < 10^{-2}$ ), suggesting that this class may be enriched in genes involved in

volved in signal transduction though the reason for this enrichment remains unclear. 3133  
The *dpy-22(bx93)*-specific class did not show enrichment on any test, consistent 3134  
with our interpretation that this class is the result of random false positive hits. 3135

### Predicted interactions of Mediator with Wnt and Ras pathways in *C. elegans* 3136

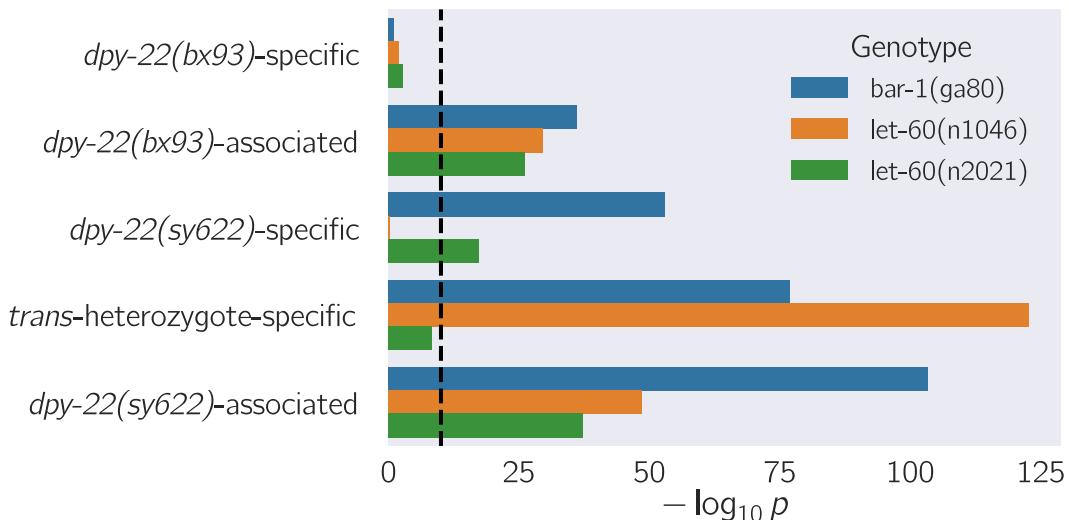
Previous work in *C. elegans* (Moghal and Paul W. Sternberg, 2003; Zhang and 3137  
Emmons, 2000) has implicated *dpy-22* as an inhibitor of the Wnt and Ras path- 3138  
ways during the formation of the vulva and the male tail. We obtained expression 3139  
profiles for *bar-1(ga80)* mutants as well as loss-of-function and gain-of-function 3140  
Ras mutants, *let-60(n2021)* and *let-60(n1046gf)* respectively. We predicted that the 3141  
*dpy-22(sy622)*-specific phenotypic class would exhibit the most significant overlap 3142  
(assessed by a hypergeometric enrichment test) with differentially expressed genes 3143  
in *let-60(n1046gf)* mutants, whereas the *dpy-22(bx93)*-associated phenotypic class 3144  
would exhibit the most significant overlap with *bar-1(ga80)* mutants. 3145

The *dpy-22(bx93)*-specific class did not show a transcriptomic signature associated 3146  
with either the Wnt or the Ras pathway, consistent with our interpretation of this 3147  
class as false positive (Fig. 55). All other classes showed significant enrichment 3148  
with genes perturbed in *bar-1(ga80)*. Similarly, *let-60(n2021)* showed enrichment 3149  
in all real phenotypic classes, with the exception of the *trans*-heterozygote specific 3150  
class. Contrary to our hypotheses, differentially expressed genes in *let-60(n1046gf)* 3151  
did not show significant overlap with the *dpy-22(sy622)*-specific phenotype, but they 3152  
did show significant overlap with all remaining real phenotypic classes. 3153

## Discussion 3154

### A conceptual framework for analyses of allelic series using transcriptomic 3155 phenotypes 3156

Although transcriptomic phenotypes have been used for epistatic analyses (Dixit 3157  
et al., 2016; Angeles-Albores, Puckett Robinson, Brian A Williams, et al., 2018; 3158



**Figure 55** *dpy-22* phenotypic classes are statistically significantly enriched for signatures of *let-60* (ras) and *bar-1* (wnt) signaling. We tested whether the overlap between the differentially expressed genes in *bar-1(ga80)*, *let-60(n1046)* or *let-60(n2021)* and the *dpy-22* phenotypic classes was statistically significant using a hypergeometric enrichment test. Since the hypergeometric enrichment test is very sensitive to deviations from random, and since we suspect that there may be a broad genotoxic response to all mutants, we used a statistical significance threshold of  $p < 10^{-10}$  (dashed black line).

Angeles-Albores, Leighton, et al., 2017), they have not been used to study gene 3159 function in the context of an allelic series. Outstanding challenges for transcriptomes 3160 in allelic series were how to count or identify distinct phenotypes within the different 3161 transcriptomes, how to order alleles in a severity hierarchy and how to order alleles 3162 in a dominance hierarchy. In this work, we present solutions to these problems, 3163 and propose a set of unifying concepts that we believe will be useful for future 3164 analyses. We re-analyzed an allelic series of the Mediator subunit gene *dpy-22* that 3165 had been studied previously (Moghal and Paul W. Sternberg, 2003), recapitulating 3166 and extending previous results as a proof of principle for our methodology. In our 3167 results, we derived a set of methods that do not rely on the nature of the mutations. 3168 In the subsequent discussion, we use the fact that the mutations we used were 3169 truncations to derive further insights into the functional units present in this gene. 3170

To interpret our phenotypic classes in a biological context, we investigated whether 3171

these phenotypic classes contained Ras and Wnt expression signatures. Our attempts 3172 were partially successful, but a more rigorous analysis awaits the availability of a 3173 larger mutant set to establish empirically the overlap that is biologically significant. 3174 In part, we reason that some genes may form part of a broad stress response. If that 3175 were the case, many mutants may share similar transcriptomic signatures. 3176

### Phenotypic classes and their sequence requirements

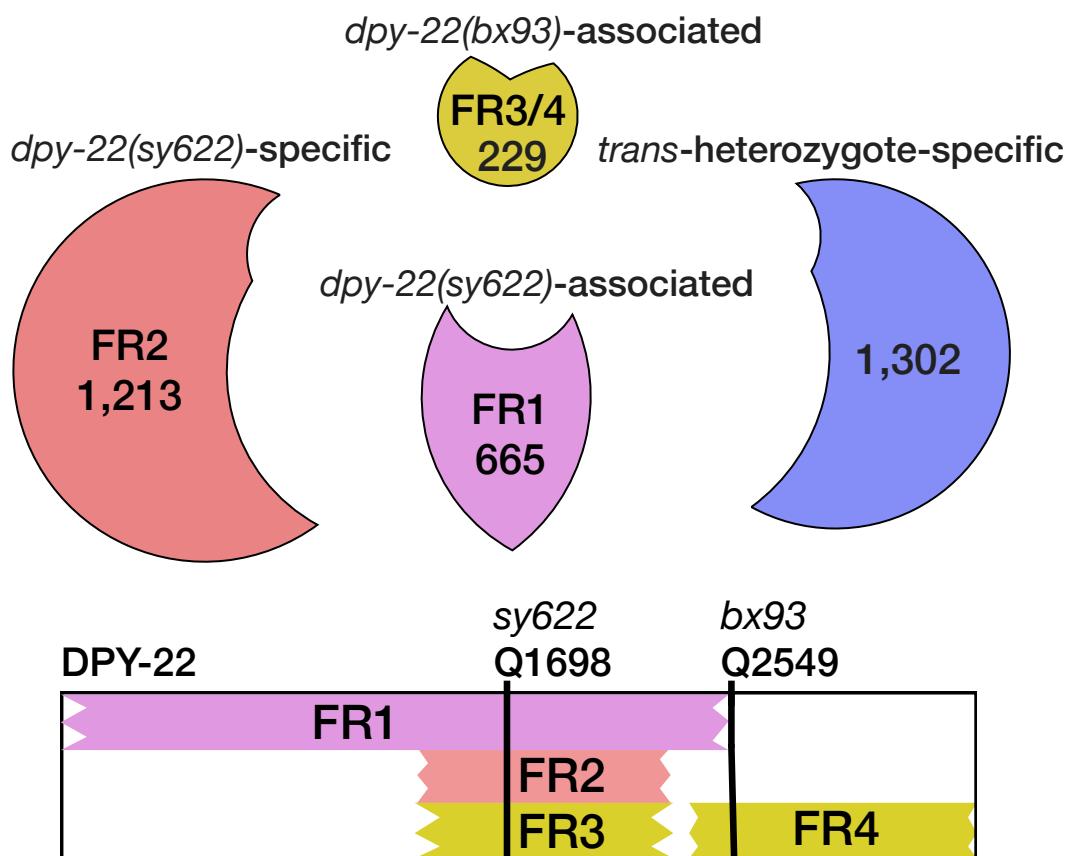
3177

Because the mutations we used are truncations, our results suggest the existence of 3178 various functional regions in *dpy-22/MDT12* (Fig. 56). These functional regions 3179 could encode protein domains with biochemical activity, or they could encode 3180 biochemically active amino acid motifs, such as nuclear localization sequences or 3181 protein binding sites. These functional regions could confer stability to the protein, 3182 thereby regulating its levels. As a caveat, we note that we have interpreted the effects 3183 these mutations have in terms of their putative effects at the protein level. In the 3184 case of our alleles, the relevant homozygotes had wild-type *dpy-22* mRNA levels, 3185 suggesting that these mutations do not affect the stability of the mRNA. 3186

The *dpy-22(sy622)*-specific phenotypic class is likely controlled by a single func- 3187 tional region, functional region 1 (FR1). Sequence necessary for wild-type FR1 3188 functionality is encoded between amino acid positions 1 and 2,549, since this is the 3189 sequence that is intact in the *bx93* allele. We speculate that this functional region 3190 may be the reason that *bx93* is unable to complement the Muv phenotype of *sy622* in 3191 a sensitized *let-23* background, since *trans*-heterozygotes in this background exhibit 3192 a semidominant Muv phenotype. The *dpy-22(sy622)*-associated phenotypic class 3193 is likely controlled by a second functional region, functional region 2 (FR2), and 3194 some necessary sequences for wild-type function are encoded between amino acid 3195 positions 1,698 and 2,549, but additional sequence could lie between amino acids 3196 1 and 1,698. It is unlikely that FR1 and FR2 are identical because their dominance 3197

behaviors are very different. The *dpy-22(bx93)* allele was largely dominant over 3198  
 the *dpy-22(sy622)* allele for the *dpy-22(bx93)*-associated class, but gene expression 3199  
 in this class was perturbed in both homozygotes. The perturbations were greater 3200  
 for *dpy-22(sy622)* homozygotes than for *dpy-22(bx93)* homozygotes. This behavior 3201  
 can be explained if the *dpy-22(bx93)*-associated class is controlled jointly by two 3202  
 distinct effectors, functional regions 3 and 4 (FR3, FR4, see Fig. 56). Such a model 3203  
 would propose that the sequences necessary for FR3 functionality are within the in- 3204  
 terval 1 and 2,549, and some sequences necessary for FR4 functionality are encoded 3205  
 between positions 2549 and 3499. This model explains how expression levels of 3206  
 the *bx93*-associated phenotypic class in the *trans*-heterozygote are complemented 3207  
 to the levels of the *bx93* homozygote, because FR3 is complemented in *trans*, but 3208  
 FR4 is defective. Thus, FR3 encodes a functionality that is not dosage-dependent. 3209  
 One possibility is that FR3 is equivalent to FR1 or FR2, and FR4 modifies activity 3210  
 of either of these regions at a subset of loci. A rigorous examination of this model 3211  
 will require studying many alleles that mutate the region between Q1689 and Q2549 3212  
 using homozygotes and *trans*-heterozygotes. 3213

We also found a class of transcripts that had perturbed levels in *trans*-heterozygotes 3214  
 only; its biological significance is unclear. Phenotypes unique to *trans*-heterozygotes 3215  
 are often the result of physical interactions such as homodimerization, or dosage 3216  
 reduction of a toxic product (Yook, 2005). In the case of *dpy-22/MDT12* orthologs, 3217  
 these explanations seem unlikely since DPY-22 is a monomeric subunit of the 3218  
 CKM. Another possibility is that the *trans*-heterozygote-specific class is the result 3219  
 of complex tissue cross-talk. Massive single-cell RNA-seq of *C. elegans* has recently 3220  
 been reported (Cao et al., 2017), and this tool could provide valuable information 3221  
 regarding this hypothesis. Another possibility is that the *cis*-marker we used for 3222  
 the *bx93* allele, *dpy-6(e14)*, which we assumed to be recessive in all phenotypes, 3223  
 actually has dominant transcriptomic phenotype. 3224



**Figure 56** The functional regions associated with each phenotypic class can be mapped intragenically. The number of genes associated with each class is shown. The *dpy-22(bx93)*-associated class may be controlled by two functional regions. FR1 is a dosage-sensitive unit. FR2 and FR3 could be redundant if FR4 is a modifier of FR2 functionality at *dpy-22(bx93)*-associated loci. Note that the *dpy-22(bx93)*-associated phenotypic class is actually three classes merged together. Two of these classes are DE in *dpy-22(bx93)* homozygotes and one other genotype. Our analyses suggested that these two classes are likely the result of false negative hits and genes in these classes should be differentially expressed in all three genotypes, so we merged these three classes together (see [Methods](#)).

**Occam's razor** 3225

Transcriptomic phenotypes generate large amounts of differential gene expression 3226 data, so false positive and false negative rates can lead to spurious phenotypic classes 3227 whose putative biological significance is misleading. Such artifacts are particularly 3228 likely when a phenotypic class is small. Notably, errors of interpretation cannot be 3229 avoided by setting a more stringent  $q$ -value cut-off: doing so will decrease the false 3230 positive rate, but increase the false negative rate, which will in turn produce smaller 3231 phenotypic classes than expected. Our method tries to avoid this pitfall by using total 3232 error rate estimates to assess the plausibility of each class, though a major drawback 3233 is that it relies on a subjective estimation of the false negative rate. These conclusions 3234 are of broad significance to research where highly multiplexed measurements are 3235 compared to identify similarities and differences in the genome-wide behavior of a 3236 single variable under multiple conditions. 3237

We have shown that transcriptomes can be used to study allelic series in the context 3238 of a large, pleiotropic gene. We identified separable phenotypic classes that would 3239 otherwise be obscured by other methods, correlated each class to a functional region, 3240 and identified sequence requirements for each region. Given the importance of allelic 3241 series for characterizing gene function and their roles in specific genetic pathways, 3242 we are optimistic that this method will be a useful addition to the geneticist's arsenal. 3243

**References** 3244

- Allen, Benjamin L and Dylan J Taatjes (2015). "The Mediator complex: a central 3245 integrator of transcription." In: *Nature reviews. Molecular cell biology* 16.3, 3246 pp. 155–166. ISSN: 1471-0080. doi: [10.1038/nrm3951](https://doi.org/10.1038/nrm3951). 3247
- Andrews, S. (2010). *FastQC: A quality control tool for high throughput sequence 3248 data.* doi: [citeulike-article-id:11583827](https://doi.org/10.11583/11583827). 3249
- Angeles-Albores, David, Daniel H. W. Leighton, et al. (2017). "The *Caenorhabditis 3250 elegans* Female-Like State: Decoupling the Transcriptomic Effects of Aging 3251 and Sperm Status". In: *G3: Genes, Genomes, Genetics* 7.9. 3252

- Angeles-Albores, David, Raymond Y. N. Lee, et al. (2016). “Tissue enrichment analysis for *C. elegans* genomics”. In: *BMC Bioinformatics* 17.1, p. 366. ISSN: 1471-2105. doi: [10.1186/s12859-016-1229-9](https://doi.org/10.1186/s12859-016-1229-9). 3253  
3254  
3255
- Angeles-Albores, David, Carmie Puckett Robinson, Brian A. Williams, et al. (2018). “Reconstructing a metazoan genetic pathway with transcriptome-wide epistasis measurements”. In: *Proceedings of the National Academy of Sciences*, p. 201712387. ISSN: 0027-8424. doi: [10.1073/pnas.1712387115](https://doi.org/10.1073/pnas.1712387115). 3256  
3257  
3258  
3259
- Angeles-Albores, David, Carmie Puckett Robinson, Brian A Williams, et al. (Mar. 2018). “Reconstructing a metazoan genetic pathway with transcriptome-wide epistasis measurements.” In: *Proceedings of the National Academy of Sciences of the United States of America* 115.13, E2930–E2939. ISSN: 1091-6490. doi: [10.1073/pnas.1712387115](https://doi.org/10.1073/pnas.1712387115). 3260  
3261  
3262  
3263  
3264
- Aroian, Raffi V and Paul W Sternberg (1991). “Multiple functions of let-23, a *Caenorhabditis elegans* receptor tyrosine kinase gene required for vulval induction.” In: *Genetics* 128.2, pp. 251–67. ISSN: 0016-6731. 3265  
3266  
3267
- Beitel, Greg J., Scott G. Clark, and H. Robert Horvitz (Dec. 1990). “*Caenorhabditis elegans* ras gene *let-60* acts as a switch in the pathway of vulval induction”. In: *Nature* 348.6301, pp. 503–509. ISSN: 0028-0836. doi: [10.1038/348503a0](https://doi.org/10.1038/348503a0). 3268  
3269  
3270
- Bray, Nicolas L et al. (2016). “Near-optimal probabilistic RNA-seq quantification.” In: *Nature biotechnology* 34.5, pp. 525–7. ISSN: 1546-1696. doi: [10.1038/nbt.3519](https://doi.org/10.1038/nbt.3519). arXiv: [1505.02710](https://arxiv.org/abs/1505.02710). 3271  
3272  
3273
- Cao, Junyue et al. (Aug. 2017). “Comprehensive single-cell transcriptional profiling of a multicellular organism.” In: *Science (New York, N.Y.)* 357.6352, pp. 661–667. ISSN: 1095-9203. doi: [10.1126/science.aam8940](https://doi.org/10.1126/science.aam8940). 3274  
3275  
3276
- Deluca, David S. et al. (2012). “RNA-SeQC: RNA-seq metrics for quality control and process optimization”. In: *Bioinformatics* 28.11, pp. 1530–1532. ISSN: 13674803. doi: [10.1093/bioinformatics/bts196](https://doi.org/10.1093/bioinformatics/bts196). 3277  
3278  
3279
- Dixit, Atray et al. (2016). “Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens”. In: *Cell* 167.7, 1853–1866.e17. ISSN: 00928674. doi: [10.1016/j.cell.2016.11.038](https://doi.org/10.1016/j.cell.2016.11.038). 3280  
3281  
3282
- Edgar, Ron, Michael Domrachev, and Alex E Lash (Jan. 2002). “Gene Expression Omnibus: NCBI gene expression and hybridization array data repository.” In: *Nucleic acids research* 30.1, pp. 207–10. ISSN: 1362-4962. 3283  
3284  
3285
- Eisenmann, D M et al. (1998). “The beta-catenin homolog BAR-1 and LET-60 Ras coordinately regulate the Hox gene *lin-39* during *Caenorhabditis elegans* vulval development.” In: *Development (Cambridge, England)* 125, pp. 3667–3680. ISSN: 0950-1991. 3286  
3287  
3288  
3289

- Elmlund, Hans et al. (Oct. 2006). “The cyclin-dependent kinase 8 module sterically blocks Mediator interactions with RNA polymerase II.” In: *Proceedings of the National Academy of Sciences of the United States of America* 103.43, pp. 15788– 93. ISSN: 0027-8424. doi: [10.1073/pnas.0607483103](https://doi.org/10.1073/pnas.0607483103). 3290  
3291  
3292  
3293
- Ewels, Philip et al. (2016). “MultiQC: Summarize analysis results for multiple tools and samples in a single report”. In: *Bioinformatics* 32.19, pp. 3047–3048. ISSN: 14602059. doi: [10.1093/bioinformatics/btw354](https://doi.org/10.1093/bioinformatics/btw354). 3294  
3295  
3296
- Ferguson, E and H. Robert Horvitz (1985). “Identification and characterization of 22 genes that affect the vulval cell lineages of *Caenorhabditis elegans*”. In: *Genetics* 110.1, pp. 17–72. 3297  
3298  
3299
- Graham, John M. and Charles E. Schwartz (Nov. 2013). “MED12 related disorders”. In: *American Journal of Medical Genetics, Part A* 161.11, pp. 2734–2740. ISSN: 15524825. doi: [10.1002/ajmg.a.36183](https://doi.org/10.1002/ajmg.a.36183). 3300  
3301  
3302
- Greenwald, Iva S., Paul W. Sternberg, and H. Robert Horvitz (Sept. 1983). “The *lin-12* locus specifies cell fates in *Caenorhabditis elegans*”. In: *Cell* 34.2, pp. 435– 444. ISSN: 00928674. doi: [10.1016/0092-8674\(83\)90377-X](https://doi.org/10.1016/0092-8674(83)90377-X). 3303  
3304  
3305
- Hunter, John D. (2007). “Matplotlib: A 2D graphics environment”. In: *Computing in Science and Engineering* 9.3, pp. 99–104. ISSN: 15219615. doi: [10.1109/MCSE.2007.55](https://doi.org/10.1109/MCSE.2007.55). arXiv: [0402594v3 \[arXiv:cond-mat\]](https://arxiv.org/abs/0402594v3). 3306  
3307  
3308
- Jeronimo, Célia and François Robert (Oct. 2017). *The Mediator Complex: At the Nexus of RNA Polymerase II Transcription*. doi: [10.1016/j.tcb.2017.07.001](https://doi.org/10.1016/j.tcb.2017.07.001). 3309  
3310
- Kim, Seokjoong et al. (May 2006). “Mediator is a transducer of Wnt/β-catenin signaling”. In: *Journal of Biological Chemistry* 281.20, pp. 14066–14075. ISSN: 00219258. doi: [10.1074/jbc.M602696200](https://doi.org/10.1074/jbc.M602696200). 3311  
3312  
3313
- Knuesel, Matthew T et al. (Feb. 2009). “The human CDK8 subcomplex is a molecular switch that controls Mediator coactivator function.” In: *Genes & development* 23.4, pp. 439–51. ISSN: 1549-5477. doi: [10.1101/gad.1767009](https://doi.org/10.1101/gad.1767009). 3314  
3315  
3316
- Langmead, Ben et al. (2009). “Bowtie: An ultrafast memory-efficient short read aligner.” In: *Genome biology* 10, R25. doi: [10.1186/gb-2009-10-3-r25](https://doi.org/10.1186/gb-2009-10-3-r25). 3317  
3318
- Lehner, Ben et al. (Jan. 2006). “Loss of LIN-35, the *Caenorhabditis elegans* ortholog of the tumor suppressor p105Rb, results in enhanced RNA interference”. In: *Genome Biology* 7.1, R4. ISSN: 14656906. doi: [10.1186/gb-2006-7-1-r4](https://doi.org/10.1186/gb-2006-7-1-r4). 3319  
3320  
3321
- McKinney, Wes (2011). “pandas: a Foundational Python Library for Data Analysis and Statistics”. In: *Python for High Performance and Scientific Computing*, pp. 1– 9. 3322  
3323  
3324
- Moghal, N. and Paul W. Sternberg (2003). “A component of the transcriptional mediator complex inhibits RAS-dependent vulval fate specification in *C. elegans*.” In: *Development* 130.1, pp. 57–69. ISSN: 09501991. doi: [10.1242/dev.00189](https://doi.org/10.1242/dev.00189). 3325  
3326  
3327

- Mortazavi, Ali et al. (2008). “Mapping and quantifying mammalian transcriptomes by RNA-Seq”. In: *Nature Methods* 5.7, pp. 621–628. ISSN: 1548-7091. doi: [10.1038/nmeth.1226](https://doi.org/10.1038/nmeth.1226). arXiv: [1111.6189v1](https://arxiv.org/abs/1111.6189v1). 3328  
3329  
3330
- Oliphant, Travis E (2007). “SciPy: Open source scientific tools for Python”. In: *Computing in Science and Engineering* 9, pp. 10–20. ISSN: 1521-9615. doi: [10.1109/MCSE.2007.58](https://doi.org/10.1109/MCSE.2007.58). 3331  
3332  
3333
- Pérez, F. and B.E. Granger (2007). “IPython: A System for Interactive Scientific Computing Python: An Open and General-Purpose Environment”. In: *Computing in Science and Engineering* 9.3, pp. 21–29. ISSN: 15219615. doi: [doi:10.1109/MCSE.2007.53](https://doi.org/10.1109/MCSE.2007.53). 3334  
3335  
3336  
3337
- Pimentel, Harold et al. (2017). “Differential analysis of RNA-seq incorporating quantification uncertainty”. In: *Nature Methods* 14.7, pp. 687–690. ISSN: 15487105. doi: [10.1038/nmeth.4324](https://doi.org/10.1038/nmeth.4324). 3338  
3339  
3340
- Riddle, Donald L et al. (1997). *C. elegans II*. ISBN: 0879695323. doi: [NBK20183](https://doi.org/10.1101/1091-6490). 3341
- Schwarz, Erich M., Mihoko Kato, and Paul W. Sternberg (Oct. 2012). “Functional transcriptomics of a migrating cell in *Caenorhabditis elegans*.” In: *Proceedings of the National Academy of Sciences of the United States of America* 109.40, pp. 16246–51. ISSN: 1091-6490. doi: [10.1073/pnas.1203045109](https://doi.org/10.1073/pnas.1203045109). 3342  
3343  
3344  
3345
- Sulston, J. E. and S. Brenner (1974). “The DNA of *Caenorhabditis elegans*.” In: *Genetics* 77.1, pp. 95–104. ISSN: 00166731. 3346  
3347
- Takagi, Yuichiro and Roger D Kornberg (Jan. 2006). “Mediator as a general transcription factor.” In: *The Journal of biological chemistry* 281.1, pp. 80–9. ISSN: 0021-9258. doi: [10.1074/jbc.M508253200](https://doi.org/10.1074/jbc.M508253200). 3348  
3349  
3350
- Tang, Fuchou et al. (May 2009). “mRNA-Seq whole-transcriptome analysis of a single cell”. In: *Nature Methods* 6.5, pp. 377–382. ISSN: 1548-7091. doi: [10.1038/nmeth.1315](https://doi.org/10.1038/nmeth.1315). 3351  
3352  
3353
- Van Der Walt, Stéfan, S. Chris Colbert, and Gaël Varoquaux (2011). “The NumPy array: A structure for efficient numerical computation”. In: *Computing in Science and Engineering* 13.2, pp. 22–30. ISSN: 15219615. doi: [doi:10.1109/MCSE.2011.37](https://doi.org/10.1109/MCSE.2011.37). arXiv: [1102.1523](https://arxiv.org/abs/1102.1523). 3354  
3355  
3356  
3357
- Villani, Alexandra-Chloé et al. (2017). “Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors”. In: *Science* 356.6335. ISSN: 0036-8075. doi: [10.1126/science.aah4573](https://doi.org/10.1126/science.aah4573). 3358  
3359  
3360
- Wang, Jen-Chywan et al. (July 2004). “The *Caenorhabditis elegans* ortholog of TRAP240, CeTRAP240/let-19, selectively modulates gene expression and is essential for embryogenesis.” In: *The Journal of biological chemistry* 279.28, pp. 29270–7. ISSN: 0021-9258. doi: [10.1074/jbc.M401242200](https://doi.org/10.1074/jbc.M401242200). 3361  
3362  
3363  
3364
- Waskom, Michael et al. (2016). “seaborn: v0.7.0 (January 2016)”. In: doi: [10.5281/zenodo.45133](https://doi.org/10.5281/zenodo.45133). 3365  
3366

- Yamamoto, Toshiyuki and Keiko Shimojima (June 2015). “A novel MED12 mutation associated with non-specific X-linked intellectual disability”. In: *Human Genome Variation* 2, p. 15018. ISSN: 2054-345X. DOI: [10.1038/hgv.2015.18](https://doi.org/10.1038/hgv.2015.18). 3367  
3368  
3369
- Yook, Karen (2005). “Complementation”. In: *WormBook*. ISSN: 15518507. DOI: [10.1895/wormbook.1.24.1](https://doi.org/10.1895/wormbook.1.24.1). 3370  
3371
- Zhang, H. and S. W. Emmons (2000). “A *C. elegans* mediator protein confers regulatory selectivity on lineage-specific expression of a transcription factor gene”. In: *Genes and Development* 14.17, pp. 2161–2172. ISSN: 08909369. DOI: [10.1101/gad.814700](https://doi.org/10.1101/gad.814700). 3372  
3373  
3374  
3375

**TISSUE ENRICHMENT ANALYSIS FOR *C. ELEGANS* GENOMICS**

3377

3378

**Abstract**

3379

**Background** Over the last ten years, there has been explosive development in 3380 methods for measuring gene expression. These methods can identify thousands 3381 of genes altered between conditions, but understanding these datasets and forming 3382 hypotheses based on them remains challenging. One way to analyze these datasets 3383 is to associate ontologies (hierarchical, descriptive vocabularies with controlled 3384 relations between terms) with genes and to look for enrichment of specific terms. 3385 Although Gene Ontology (GO) is available for *Caenorhabditis elegans*, it does not 3386 include anatomical information. 3387

**Results** We have developed a tool for identifying enrichment of *C. elegans* tissues 3388 among gene sets and generated a website GUI where users can access this tool. 3389 Since a common drawback to ontology enrichment analyses is its verbosity, we 3390 developed a very simple filtering algorithm to reduce the ontology size by an order 3391 of magnitude. We adjusted these filters and validated our tool using a set of 30 gold 3392 standards from Expression Cluster data in WormBase. We show our tool can even 3393 discriminate between embryonic and larval tissues and can even identify tissues 3394 down to the single-cell level. We used our tool to identify multiple neuronal tissues 3395 that are down-regulated due to pathogen infection in *C. elegans*. 3396

**Conclusions** Our Tissue Enrichment Analysis (TEA) can be found within Worm- 3397 Base, and can be downloaded using Python's standard pip installer. It tests a 3398 slimmed-down *C. elegans* tissue ontology for enrichment of specific terms and 3399

provides users with a text and graphic representation of the results. 3400

## Background 3401

RNA-seq and other high-throughput methods in biology have the ability to iden- 3402  
tify thousands of genes that are altered between conditions. These genes are often 3403  
correlated in their biological characteristics or functions, but identifying these func- 3404  
tions remains challenging. To interpret these long lists of genes, biologists need 3405  
to abstract genes into concepts that are biologically relevant to form hypotheses 3406  
about what is happening in the system. One such abstraction method relies on Gene 3407  
Ontology (GO). GO provides a controlled set of hierarchically ordered terms (Ash- 3408  
burner et al., 2000; The Gene Ontology Consortium, 2015) that provide detailed 3409  
descriptions about the molecular, cellular or biochemical functions of any gene. For 3410  
a given gene list, existing software programs can query whether a particular term 3411  
is enriched (Mi, Dong, et al., 2009; McLean et al., 2010; Huang, Lempicki, and 3412  
Brad T Sherman, 2009; Pathan et al., 2015). One area of biological significance 3413  
that GO does not include is anatomy. One way to address this shortcoming is to use 3414  
a ‘tissue ontology’ that provides a complete anatomical description for an organism 3415  
(e.g. ‘tissue’, ‘organ’ or ‘specific cell’), in this case for *C. elegans*. Such an ontol- 3416  
ogy has been described previously for this organism (R. Y. N. Lee and Sternberg, 3417  
2003). Cells and tissues are physiologically relevant units with broad, relatively 3418  
well-understood functionalities amenable to hypothesis formation. The *C. elegans* 3419  
database, WormBase (Howe et al., 2016), maintains a curated list of gene expression 3420  
data from the literature. Here we provide a new framework that analyzes a user-input 3421  
list for enrichment of specific cells and tissues. 3422

Another problem frequently associated with GO enrichment analysis is that it is often 3423  
difficult to interpret due to the large number of terms associated with a given gene 3424  
(which we refer to as ‘result verbosity’). DAVID, a common tool for GO enrichment 3425

analysis, clusters enriched terms into broad categories (Huang, Brad T. Sherman, et al., 2007), whereas PANTHER (Mi, Dong, et al., 2009; Mi, Muruganujan, and Thomas, 2013) attempts to solve this issue by employing a manually reduced ontology, GOslim (pers.\_ comm., H. Yu and P. Thomas). To reduce verbosity, we have filtered our ontology using a small set of well-defined criteria to remove terms that do not contribute additional information. To our knowledge, such filtering has not been performed in an algorithmic fashion for a biological ontology before; indeed, DAVID does not employ term trimming *a priori* of testing, but rather fuzzy clustering *post* testing to reduce the number of ontology terms. Other pruning methods do exist (see for example (J. W. Kim, Caralt, and Hilliard, 2007; Garrido and Requena, 2012)), but the pruning is query-dependent or generates a brand new ‘brief ontology’ which satisfies a set of logic relationships and has certain connectivity requirements. We do not propose to regenerate a new ‘brief ontology’, but instead we use our approach to select those nodes that have sufficient annotated evidence for statistical testing. We believe our trimming methodology strikes a good balance between detailed tissue calling and conservative testing.

3441

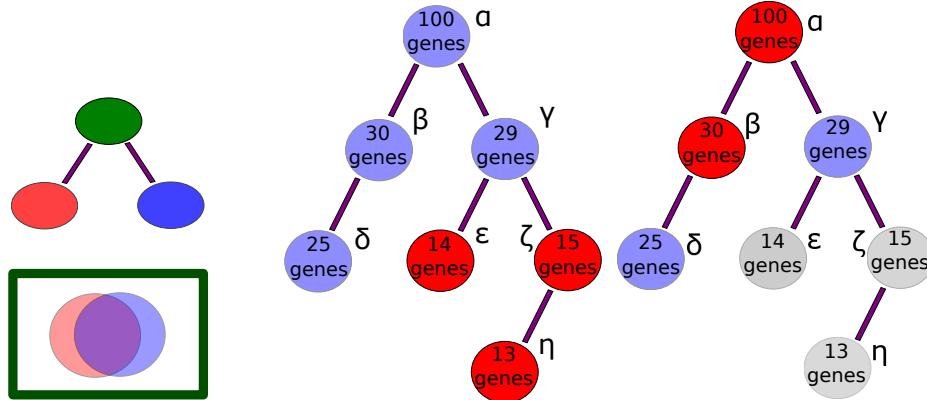
We have developed a tool that tests a user-provided list of genes for term enrichment using a nematode-specific tissue ontology. This ontology, which is not a module of Gene Ontology, is verbose. We select nodes from the ontology for statistical testing using an algorithmic approach, outlined below, that reduces multiple hypothesis testing issues by limiting testing to terms that are well-annotated. The results are provided to the user in a GUI that includes a table of results and an automatically generated bar-chart. This software addresses a previously unmet need in the *C. elegans* community for a tool that reliably and specifically links gene expression with changes in specific cells, organs or tissues in the worm.

3450

<b>Results</b>	3451
<b>Generating a Gene-Tissue Dictionary by Specific Node Selection</b>	3452
<b>Reducing term redundancy through a similarity metric</b>	3453

For our tool, we employ a previously generated cell and tissue ontology for *C. elegans* (R. Y. N. Lee and Sternberg, 2003), which is maintained and curated by WormBase. This ontology contains thousands of anatomiy terms, but not every term is equally well-annotated. As a first step to generate our tissue enrichment software, we wished to select tissue terms that were reasonably well-annotated, yet specific enough to provide insight and not redundant with other terms. For example, nematodes have a number of neurons that are placed symmetrically along the left/right body axis, and are functionally similar. These left/right neuronal pairs (which are sisters in the ontology) have almost identical annotations, with at most one or two gene differences between them, and therefore we cannot have statistical confidence in differentiating between them. As a result, testing these sister terms provides no additional information compared with testing only the parent node to these sisters. To identify redundancy, we defined two possible similarity metrics (see *Methods* section and Figure 61) that can be used to identify ontology sisters that have very high similarity between them. Intuitively, a set of sisters can be considered very similar if they share most gene annotations. Within a given set of sisters, we can calculate a similarity score for a single node by counting the number of unique annotations it contains and dividing by the total number of unique annotations in the sister set. Having assigned to each sister a similarity score, we can identify the **average** similarity score for this set of sisters, and if this average value exceeds a threshold, these sisters are not considered testable candidates. An alternative method is check whether **any** of the scores exceeds a predetermined threshold, and if so remove this sister set from the ontology. We referred to these two scoring criteria as ‘**avg**’ and ‘**any**’ respectively.

3477



**Figure 61** Schematic representation of trimming filters for an acyclical ontology. **a.** The parent node (green) contains at least as many annotations as the union of the two sisters. These two sisters share annotations extensively, as expressed by the overlap in the Venn diagram, so they qualify for removal. **b.** Nodes with less than a threshold number of genes are trimmed (red) and discarded from the dictionary. Here, the example threshold is 25 genes. Nodes  $\epsilon, \zeta, \eta$ , shown in red are removed. **c.** Parent nodes are removed recursively, starting from the root, if all their daughter nodes have more than the threshold number of annotations. Nodes in grey ( $\epsilon, \zeta, \eta$ ) were removed in the previous step. Nodes  $\alpha, \beta$  shown in red are trimmed because each one has a complete daughter set. Only nodes  $\gamma$  and  $\delta$  will be used to generate the static dictionary.

### Terminal branch terms and parent terms can be safely removed in an algorithmic fashion

3478

3479

Another problem arises from the ontology being scarcely populated. Many nodes 3480 have 0–10 annotations, which we consider too few to accurately test. To solve 3481 this issue, we implemented another straightforward node selection strategy. For a 3482 given terminal node, we test whether the node has more than a threshold number of 3483 annotations. If it does not, the node is not used for statistical testing. The next higher 3484 node in the branch is tested and removed recursively until a node that satisfies the 3485 condition is found. At that point, no more nodes can be removed from that branch. 3486 This completion is guaranteed by the structure of the ontology: parent nodes inherit 3487 all of the annotations of all of their descendants, so the number of annotated terms 3488 monotonically increases with increasing term hierarchy (see Figure 61). In this way, 3489

we ensure that our term dictionary includes only those tissues that are considered 3490 sufficiently well annotated for statistical purposes. 3491

Additionally, we reasoned that for any parent node if all its daughters were selected 3492 for testing, there was no additional benefit to test the parent. We removed parent 3493 nodes from the analysis if all their daughter nodes passed the annotation threshold 3494 (see Figure 61). We called this a ceiling filter. Applying these three filters reduced 3495 the number of ontology terms by an order of magnitude. 3496

### **Filtering greatly reduces the number of nodes used for analysis** 3497

By itself, each of these filters can reduce the number of nodes employed for analysis, 3498 but applying the filters in different orders removes different numbers of nodes (not all 3499 the filters are commutative). We chose to always execute annotation and similarity 3500 thresholding first, followed by the ceiling filter. For validation (see below) we made 3501 a number of different dictionaries. The original ontology has almost 6,000 terms 3502 of which 1675 have at least 5 gene annotations. After filtering, dictionary sizes 3503 ranged from 21 to a maximum of 460 terms, which shows the number of terms in a 3504 scarcely annotated ontology can be reduced by an order of magnitude through the 3505 application of a few simple filters (see Table 61). These filters were used to compile 3506 a static dictionary that we employ for all analyses (see *Validation of the algorithm* 3507 and *parameter selection* section for details). Our trimming pipeline is applied as 3508 part of each new WormBase release. This ensures that the ontology database we are 3509 using remains up-to-date with regards to both addition or removal of specific terms 3510 as well as with regard to gene expression annotations. 3511

### **Tissue enrichment testing via a hypergeometric model** 3512

Having built a static dictionary, we generated a Python script that implements a 3513 significance testing algorithm based on the hypergeometric model. Briefly, the 3514

**Table 61** Parameter specifications and number of tissues for all dictionaries. The ‘Method’ column refers to the trimming criterion for the similarity metric. We used two such criteria, ‘any’ and ‘avg’. ‘any’: For a given sister set, if any sister had a similarity exceeding the corresponding threshold, all sisters were removed from the final dictionary. ‘avg’: For a given sister set, if the average similarity across all the sisters in the set was greater than the corresponding threshold, all sisters were removed from the final dictionary.

Annotation Cutoff	Similarity Threshold	Method	No. Of Terms in Dictionary
25	0.9	any	460
25	0.9	avg	461
25	0.95	any	466
25	0.95	avg	468
25	1.0	any	476
25	1.0	avg	476
33	0.9	any	261
33	0.9	avg	255
33	0.95	any	261
33	0.95	avg	262
33	1.0	any	247
33	1.0	avg	247
50	0.9	any	83
50	0.9	avg	77
50	0.95	any	82
50	0.95	avg	81
50	1.0	any	70
50	1.0	avg	70
100	0.9	any	45
100	0.9	avg	35
100	0.95	any	42
100	0.95	avg	36
100	1.0	any	21
100	1.0	avg	21

hypergeometric model tests the probability of observing  $n_i$  occurrences of a tissue  $i$  3515  
in a list of size  $M$  if there are  $m_i$  labels for that tissue in a dictionary of total size  $N$  3516  
that are drawn without replacement. Mathematically, this is expressed as: 3517

$$P(n_i|N, m_i, M) = \frac{\binom{m_i}{n_i} \binom{M - m_i}{N - n_i}}{\binom{N}{n_i}}. \quad (6.1)$$

Although a user will input gene IDs, we test the number of occurrences of a term 3518  
within the gene list, so a single gene can contribute to multiple terms. Due to 3519  
the discrete nature of the hypergeometric distribution, this algorithm can generate 3520  
artifacts when the list is small. To avoid spurious results, a tissue is never considered 3521  
significant if there are no annotations for it in the user-provided list. 3522

Once the p-values for each term have been calculated, we apply a standard FDR cor- 3523  
rection using a Benjamini-Hochberg step-up algorithm (Benjamini and Hochberg, 3524  
[1995](#)). FDR corrected p-values are called q-values. Genes that have a q-value less 3525  
than a given alpha are considered significant. Our default setting is an alpha of 3526  
0.1, which is a standard threshold broadly agreed upon by the scientific community 3527  
(see for example (Love, Huber, and Anders, [2014](#); Pawitan et al., [2005](#); Storey and 3528  
Tibshirani, [2003](#))). This threshold cannot be altered in the web GUI, but is user 3529  
tunable through our command-line implementation. 3530

Users input a gene list using any valid gene name for *C. elegans*. These names are 3531  
processed into standard WormBase gene IDs (WBGene IDs). The program returns 3532  
a table containing all the enriched terms and associated information such as number 3533  
of terms in gene list and expected number of terms. Finally, the program can also 3534  
return a bar chart of the enrichment fold change for the fifteen tissues with the 3535  
lowest measured q-values. The bars in the graph are sorted in ascending order of 3536

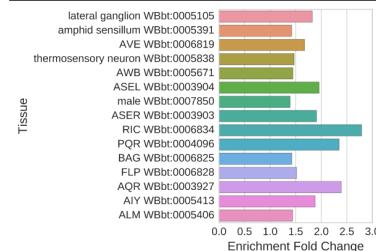
q-value and then in descending order of fold-change. Bars are colored for ease of viewing, and color does not convey information. Our software is implemented in an easy to use GUI (see Figure 62). Anatomy terms are displayed in human-readable format followed by their unique ontology ID (WBbt ID). In summary, each time the ontology annotations are updated, a new trimmed ontology is generated using our filters; in parallel, users can submit their gene lists through WormBase for testing, with results output in a number of formats (see Figure 63).

3543

### Tissue Enrichment Analysis Results •

Return up to 15 most significant anatomy terms.

	Tissue	Expected	Observed	Enrichment Fold Change	P value	Q value
lateral ganglion	WBbt:000510577	141	1.8	<10 <sup>-6</sup>	<10 <sup>-6</sup>	
amphid sensillum	WBbt:00053912e+02	291	1.4	<10 <sup>-6</sup>	<10 <sup>-6</sup>	
AVE	WBbt:000681968	114	1.7	<10 <sup>-6</sup>	<10 <sup>-6</sup>	
thermosensory neuron	WBbt:00583813e+02	187	1.5	<10 <sup>-6</sup>	<10 <sup>-6</sup>	
AWB	WBbt:00567113e+02	187	1.4	<10 <sup>-6</sup>	1.3e-06	
ASEL	WBbt:00390420	39	2	6.5e-06	0.00029	
male	WBbt:0078501e+02	140	1.4	1.1e-05	0.00043	
ASER	WBbt:00390320	38	1.9	1.6e-05	0.00054	
PQR	WBbt:00040968.5	20	2.3	3.7e-05	0.001	
RIC	WBbt:00068345	14	2.8	3.5e-05	0.001	
BAG	WBbt:000682573	104	1.4	5e-05	0.0012	
FLP	WBbt:000682647	72	1.5	7.8e-05	0.0017	
AQR	WBbt:0039275	12	2.4	0.00064	0.013	
AIY	WBbt:00541312	22	1.9	0.0008	0.015	
ALM	WBbt:00540642	60	1.4	0.0011	0.02	
SMDDL	WBbt:0049725.3	11	2.1	0.0041	0.057	



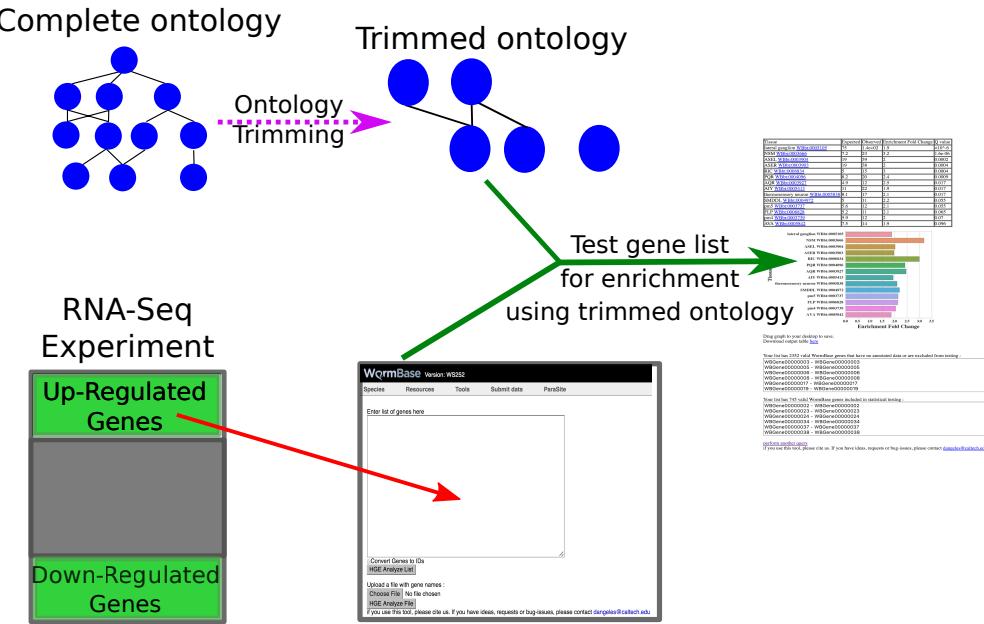
Drag graph to your desktop to save.  
Download output table here

Your list has 966 valid WormBase genes that have no annotated data or are excluded from testing :  
 WBGene000000005 - WBGene000000005  
 WBGene000000009 - WBGene000000009  
 WBGene000000017 - WBGene000000017  
 WBGene000000021 - WBGene000000021  
 WBGene000000025 - WBGene000000025  
 WBGene000000133 - WBGene000000133

Your list has 2131 valid WormBase genes included in statistical testing :  
 WBGene000000002 - WBGene000000002  
 WBGene000000003 - WBGene000000003  
 WBGene000000006 - WBGene000000006  
 WBGene000000019 - WBGene000000019  
 WBGene000000023 - WBGene000000023  
 WBGene000000024 - WBGene000000024

perform another query

**Figure 62** Screenshot of results from the web GUI. After inputting a gene-list, the user is provided with the results. An HTML table is output with hyperlinks to the ontology terms. A publication-ready graph is provided below, which can be saved by dragging to the desktop. The graph is colored for better visualization; color is not intended to convey information. The graph and the table show anatomy terms in human-readable format, followed by their unique WBbt ID. Finally, lists of the genes used and discarded for the analysis are also presented.



**Figure 63** TEA Workflow. The complete ontology is annotated continuously by WormBase curators. After each update, the ontology is processed to remove uninformative terms, and the remaining terms are used for statistical testing. Users can select a gene list and input it into our tool using our WormBase portal. The gene list is tested for enrichment using the trimmed ontology, and results are output in tabular and graphic formats for analysis.

### Validation of the algorithm and optimizing parameter selection

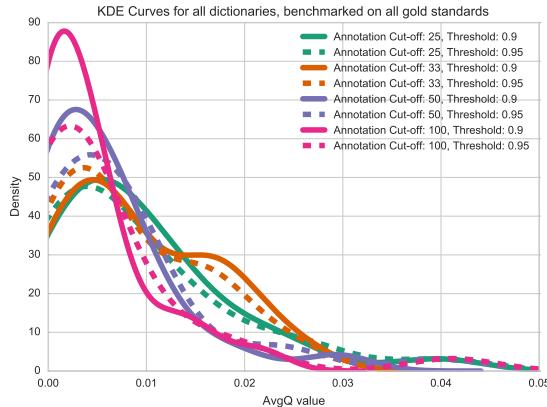
3544

We wanted to select a dictionary that included enough terms to be specific beyond 3545 the most basic *C. elegans* tissues, yet would minimize the number of spurious 3546 results and which had a good dynamic range in terms of enrichment fold-change. 3547 Larger tissues are correlated with better annotation, so increasing term specificity is 3548 associated with losses in statistical power. To help us select an appropriate dictionary 3549 and validate our tool, we used a set of 30 gold standards based on microarray and 3550 RNA-seq literature which are believed to be enriched in specific tissues (Gaudet 3551 et al., 2004; Spencer et al., 2011; Cinar, Keles, and Jin, 2005; Watson et al., 2008; 3552 Pauli et al., 2006; Portman and Emmons, 2004; Fox et al., 2007; Smith et al., 2010). 3553 These data sets are annotated gene lists derived from the corresponding Expression 3554 Cluster data in WormBase. Some of these studies have been used to annotate gene 3555

expression, and so they did not constitute an independent testing set. To correct this 3556 flaw, we built a clean dictionary that specifically excluded all annotation evidence 3557 that came from these studies. 3558

As a first attempt to select a dictionary, we generated all possible combinations of 3559 dictionaries with minimal annotations of 10, 25, 33, 50 and 100 genes and similarity 3560 cutoffs of 0.9, 0.95 and 1, using ‘avg’ or ‘any’ similarity thresholding methods 3561 (see Table 61). The number of remaining ontology terms was inversely correlated 3562 to the minimum annotation cutoff, and was largely insensitive to the similarity 3563 threshold in the range we explored. Next, we analyzed all 30 datasets using each 3564 dictionary. Because of the large number of results, instead of analyzing each set of 3565 terms individually, we measured the average q-value for significantly enriched terms 3566 in each dataset without regard for the perceived accuracy of the terms that tested 3567 significant. We found that the similarity threshold mattered relatively little for any 3568 dictionary. We also noticed that the ‘any’ thresholding method resulted in tighter 3569 histograms with a mode closer to 0. For this reason, we chose the ‘any’ method for 3570 dictionary generation. The average q-value increased with decreasing annotation 3571 cut-off (see Figure 64), which reflects the decreasing statistical power associated 3572 with fewer annotations per term, but we remained agnostic as to how significant is 3573 the trade-off between power and term specificity. Based on these observations, we 3574 ruled out the dictionary with the 100 gene annotation cut-off: it had the fewest terms 3575 and its q-values were not low enough in our opinion to compensate for the trade-off 3576 in specificity. 3577

To select between dictionaries generated between 50, 33 and 25 annotation cut-offs, 3578 and also to ensure the terms that are selected as enriched by our algorithm are 3579 reasonable, we looked in detail at the enrichment analysis results. Most results 3580 were comparable and expected. For some sets, all dictionaries performed well. 3581 For example, in our ‘all neuron enriched sets’ (Spencer et al., 2011; Watson et al., 3582



**Figure 64** Kernel density estimates (KDE) for 30 gold standard datasets. We ran TEA on 30 datasets we believed to be enriched in particular tissues and pooled all the results to observe the distribution of q-values. The mode of the distribution for dictionaries with annotation cut-offs of 100 and 50 genes are very similar; however, when the cut-off is lowered to 25 genes, the mode of the distribution shifts to the left, potentially signalling a decrease in measurement power.

2008) all terms were neuron-related regardless of the dictionary used (see Table 62). 3583 On the other hand, for a set enriched for germline precursor expression in the 3584 embryo (Spencer et al., 2011), the 50 cutoff dictionary was only able to identify 3585 ‘oocyte WBbt:006797’, which is not a germline precursor although it is germline 3586 related; whereas the two smaller dictionaries singled out actual germline precursor 3587 cells—at the 33 cutoff, our tool identified the larval germline precursor cells ‘Z2’ and 3588 ‘Z3’ as enriched, and at the 25 gene cutoff the embryonic germline precursor terms 3589 ‘P<sub>4</sub>’, ‘P<sub>3</sub>’ and ‘P<sub>2</sub>’ were identified in addition to ‘Z2’ and ‘Z3’. We also queried 3590 an intestine precursor set (Spencer et al., 2011). Notably, this gene set yielded 3591 no enrichment when using the 25 cutoff dictionary, nor when using the 50 cutoff 3592 dictionary. However, the 33 cutoff dictionary identified the E lineage, which is the 3593 intestinal precursor lineage in *C. elegans*, as enriched. Both of these results capture 3594 specific aspects of *C. elegans* that are well known to developmental biologists. 3595

Not all queries worked equally well. For example, a number of intestinal sets (Spencer 3596 et al., 2011; Pauli et al., 2006) were not enriched in intestine-related terms in any 3597

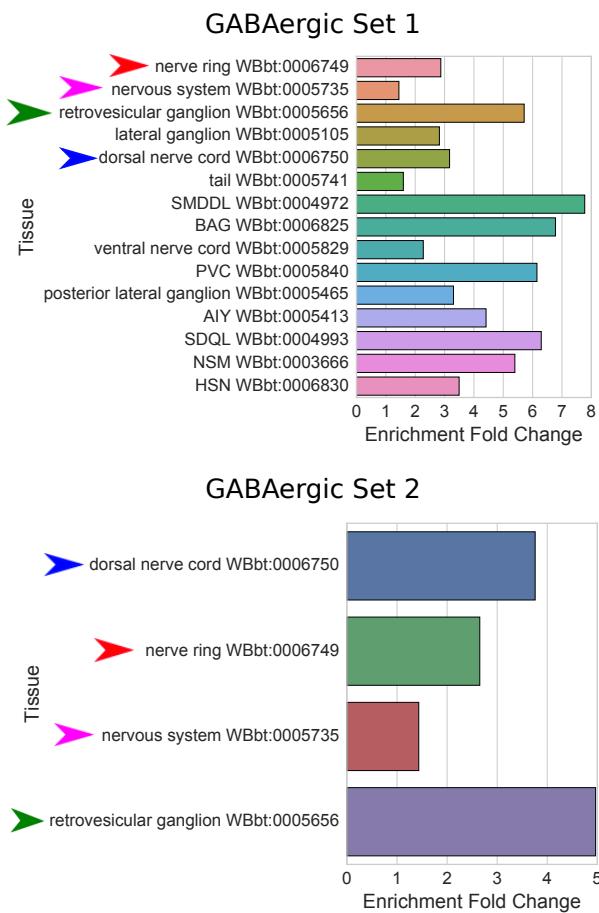
dictionary, but were enriched for pharynx and hypodermis. We were surprised that 3598  
intestinal gene sets performed poorly, since the intestine is a relatively well-annotated 3599  
tissue. 3600

We assessed the internal agreement of our tool by using independent gene sets 3601  
that we expected to be enriched in the same tissues. We used two pan-neuronal 3602  
sets (Spencer et al., 2011; Watson et al., 2008); two PVD sets (Spencer et al., 2011; 3603  
Smith et al., 2010); and two GABAergic sets (Spencer et al., 2011; Cinar, Keles, 3604  
and Jin, 2005). Overall, the tool has good internal agreement. On most sets, the 3605  
same terms were enriched, although order was somewhat variable (see Table 65), and 3606  
most high-scoring terms were preserved between sets. All comparisons can be found 3607  
online in our Github repository (see Availability of data and materials). Overall, the 3608  
dictionary generated by a 33 gene annotation cutoff with 0.95 redundancy threshold 3609  
using the ‘any’ criterion performed best, with a good balance between specificity, 3610  
verbosity and accuracy, so we selected this parameter set to generate our static 3611  
dictionary. As of this publication, the testable dictionary contains 261 terms. 3612

## Applying the tool

We applied our tool to the RNA-seq datasets developed by Engelmann et al. (En- 3614  
gelmann et al., 2011) to gain further understanding of their underlying biology. 3615  
Engelmann et al.\_ exposed young adult worms to 5 different pathogenic bacteria or 3616  
fungi for 24 hours, after which mRNA was extracted from the worms for sequenc- 3617  
ing. We ran TEA on the genes Engelmann *et al* identified as up- or down-regulated. 3618  
Initially we noticed that genes that are down-regulated tend to be twice as better 3619  
annotated on average than genes that were up-regulated, suggesting that our under- 3620  
standing of the worm immune system is scarce, in spite of important advances made 3621  
over the last decade. Up-regulated tissues, when detected, almost always included 3622  
the hypodermis and excretory duct. Three of the five samples showed enrichment of 3623

**Table 62** Comparison of results for a GABAergic neuronal-enriched gene set from Watson (Watson et al., 2008) showing that results are similar regardless of annotation cutoff. We ran the same gene list on a dictionary with a minimum annotation cutoff of 50, similarity threshold of 0.95 and similarity method ‘any’ versus another with a minimum annotation cutoff of 33, similarity threshold of 0.95 and similarity method ‘any’. In the table, columns are labeled with their significance value (Q-value) or enrichment fold change followed by a hyphen and a number which indicates which the cutoff for the dictionary that was used for testing. Not all tissues are present in either dictionary. Hyphens denote not-applicable values, which occurs when a particular tissue is not present in both dictionaries.



**Figure 65** Independently derived gene sets show similar results when tested with the same dictionary. **Set 1.** GABAergic gene set from Watson (Watson et al., 2008). **Set 2.** GABAergic gene set from Spencer (Spencer et al., 2011). Arrowheads highlight identical terms between both analyses. All terms refer to neurons or neuronal tissues and are GABA-associated. Dictionary with cutoff: 33; threshold: 0.95; method: ‘any’.

neuronal tissues or neuronal precursor tissues among the down-regulated genes. As 3624 an independent verification, we also performed GO analysis using PANTHER on 3625 the down-regulated genes for *D. coniospora*. These results also showed enrichment 3626 in terms associated with neurons (see Figure 66). A possible explanation for this 3627 neuronal association might be that the infected worms are sick and the neurons are 3628 beginning to shut down; an alternative hypothesis would be that the worm is down- 3629 regulating specific neuronal pathways as a behavioral response against the pathogen. 3630 Indeed, several studies (Meisel and D. H. Kim, 2014; Zhang, Lu, and Bargmann, 3631 2005) have provided evidence that *C. elegans* uses chemosensory neurons to identify 3632 pathogens. Our results highlight the involvement of various *C. elegans* neuronal 3633 tissues in pathogen defense. 3634

## Discussion

3635

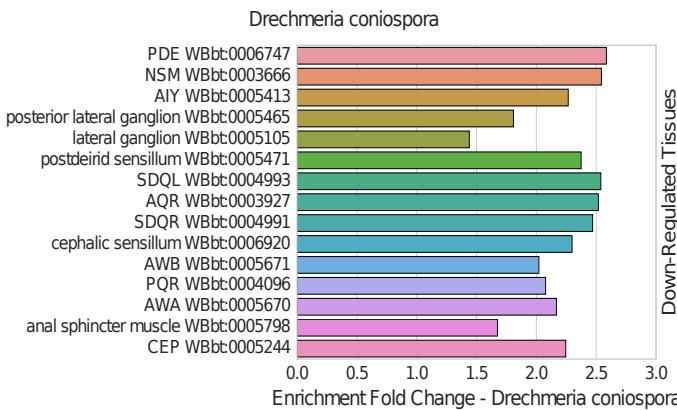
We have presented a tissue enrichment analysis tool that employs a standard hyper- 3636 geometric model to test the *C. elegans* tissue ontology. We use a hypergeometric 3637 function to test a user-provided gene list for enrichment of anatomical terms in 3638 *C. elegans*. Our hope is that the physiological relevance of anatomical terms will 3639 enable researchers to make hypotheses about high-dimensionality data. Specifically, 3640 we believe an enriched term may broadly suggest one of two hypotheses: if a list 3641 is enriched in a particular anatomical region, that anatomical region is affected by 3642 the experimental treatment; alternatively, the anatomical regions that are enriched 3643 reflect biologically relevant interactions between tissues. We believe the first hy- 3644 pothesis is a reasonable one to make in the case of whole-worm RNA-seq data for 3645 example, whereas the second hypothesis may be more plausible in cases where a 3646 researcher already knows what tissues a particular gene list came from, as may be 3647 the case in single-cell RNA-seq. 3648

Our tool relies on an annotation dictionary that is continuously updated primarily 3649

### a) GO Enrichment Analysis

PANTHER GO-Slim Biological Process		#	# expected	Fold Enrichment	+/-	P value
Unclassified		12877	745	925.08	.81	+ 0.00E00
translation		362	46	26.01	1.77	+ 4.19E-02
↳protein metabolic process		1879	246	134.99	1.82	+ 1.16E-17
↳primary metabolic process		4498	603	323.14	1.87	+ 6.74E-58
↳metabolic process		5383	697	386.71	1.80	+ 6.88E-65
sensory perception		454	59	32.62	1.81	+ 3.14E-03
↳neurological system process		631	116	45.33	2.56	+ 4.50E-17
↳system process		887	183	63.72	2.87	+ 4.45E-34
↳single-multicellular organism process		956	198	68.68	2.85	+ 2.39E-36
↳multicellular organismal process		956	198	68.68	2.85	+ 2.39E-36
cellular protein modification process		904	123	64.94	1.89	+ 5.36E-09
regulation of transcription from RNA polymerase II promoter		695	95	49.93	1.90	+ 8.54E-07
↳transcription from RNA polymerase II promoter		893	115	64.15	1.79	+ 5.34E-07
↳transcription, DNA-dependent		928	123	66.67	1.84	+ 2.64E-08
↳RNA metabolic process		1258	160	90.37	1.77	+ 8.08E-10
↳nucleobase-containing compound metabolic process		1936	247	139.08	1.78	+ 2.25E-16
↳regulation of nucleobase-containing compound metabolic process		826	116	59.34	1.95	+ 3.13E-09
↳regulation of biological process		1619	213	116.31	1.83	+ 3.53E-15
↳biological regulation		2130	326	153.02	2.13	+ 6.25E-37
response to stress		370	53	26.58	1.99	+ 6.24E-04

### b) TEA



**Figure 66** *D. coniospora* Gene Enrichment Analysis and Tissue Enrichment Analysis results. We compared and contrasted the results from a gene enrichment analysis program, pantherDB, with TEA by analyzing genes that were significantly down-regulated when *C. elegans* was exposed to *D. coniospora* in a previously published dataset by Engelmann *et al* (Engelmann et al., 2011) with both tools. **a.** pantherDB screenshot of results, sorted by p-value. Only top hits shown. **b.** TEA results, sorted by q-value (lowest on top) and fold-change. Both pantherDB and TEA identify terms associated with neurons (red square). The two analyses provide complementary, not redundant, information.

with data from single gene qualitative analyses, does not require retraining and 3650  
does not require ranked genes. To our knowledge, this is the first tool that tests 3651  
tissue enrichment in *C. elegans* via the hypergeometric method, but similar projects 3652  
exist for humans and zebrafish (Y. S. Lee et al., 2013; Prykhozhij, Marsico, and 3653  
Meijsing, 2013), highlighting the relevance of our tool for high-dimensionality 3654  
biology. Chikina *et al* (Chikina et al., 2009) have previously reported a tissue 3655  
enrichment model for *C. elegans* based on a Support Vector Machine classifier that 3656  
has been trained on microarray studies. SVMs are powerful tools, but they require 3657  
continuous retraining as more tissue expression data becomes available. Moreover, 3658  
classifiers require that data be rank-ordered by some metric, something which is not 3659  
possible for certain studies. Furthermore, this tissue enrichment tool provides users 3660  
with enrichment results for only 6 large tissues. In contrast, our tool routinely tests 3661  
a much larger number of terms, and we have shown it can even accurately identify 3662  
enrichment of embryonic precursor lineages for select data sets. 3663

We have also presented the first, to our knowledge, ontology term filtering algorithm 3664  
applied to biomedical ontologies. This algorithm, which is very easy to execute, 3665  
identifies terms that have specificity and statistical power for hypothesis testing. Due 3666  
to the nature of all ontologies as hierarchical, acyclical graphs with term inheritance, 3667  
term annotations are correlated along any given branch. This correlation reduces 3668  
the benefits of including all terms for statistical analysis: for any given term along 3669  
a branch, if that term passes significance, there is a high probability that many 3670  
other terms along that branch will also pass significance. If the branch is enriched 3671  
by random chance, error propagation along a branch means that many more false 3672  
positives will follow. Thus, a researcher might be misled by the number of terms of 3673  
correlated function and assign importance to this finding; the fact that the branching 3674  
structure of GO amplifies false positive signals is a powerful argument for either 3675  
reducing branch length or branch intracorrelation, or both. On the other hand, if a 3676

term is actually enriched, we argue that there is little benefit to presenting the user 3677 with additional terms along that branch. Instead, a user will benefit most from testing 3678 sparsely along the tree at a suitable specificity for hypothesis formation. Related 3679 terms of the same level should only be tested when there is sufficient annotation 3680 to differentiate, with statistical confidence, whether one term is enriched above the 3681 other. Our algorithm reduces branch length by identifying and removing nodes that 3682 are insufficiently annotated and parents that are likely to include sparse information. 3683

We endeavoured to benchmark our tool well, but our analysis cannot address prob- 3684 lems related to spurious term enrichment. Although we were unable to determine 3685 false-positive and false-negative rates, we do not believe this should deter scientists 3686 from using our tool. Rather, we encourage researchers to use our tool as a guide, 3687 integrating evidence from multiple sources to inform the most likely hypotheses. 3688 As with any other tool based on statistical sampling, our analysis is most vulnerable 3689 to bias in the data set. For example, expression reports are negatively biased against 3690 germline expression because of the difficulties associated with expressing transgenes 3691 in this tissue (Kelly et al., 1997). As time passes, we are certain the accuracy and 3692 power of this tool will improve thanks to the continuing efforts of the worm research 3693 community; indeed, without the community reports of tissue expression in the first 3694 place, this tool would not be possible. 3695

## Conclusions

3696

We have built a tissue enrichment tool that employs a tissue ontology previously 3697 developed by WormBase. We use a simple algorithm to identify the best ontology 3698 terms for statistical testing and in this way minimize multiple testing problems. Our 3699 tool is available within WormBase or can be downloaded for offline use via ‘pip 3700 install’. 3701

**Methods** 3702

**Fetching annotation terms** 3703

We used WormBase-curated gene expression data, which includes annotated descriptions of spatial-temporal expression patterns of genes, to build our dictionary. Gene lists per anatomy term were extracted from a Solr document store of gene expression data from the WS252 database provided by WormBase (Howe et al., 2016). We used the Solr document store because it provided a convenient access to expression data that included inferred annotations. That is, for each anatomy term, the expression gene list includes genes that were directly annotated to the term, as well as those that were annotated to the term's descendant terms (if there were any). Descendant terms were those connected with the focus term by is\_a/part\_of relationship chains defined in the anatomy term ontology hierarchy. 3713

**Filtering nodes** 3714

**Defining a Similarity Metric** 3715

To identify redundant sisters, we defined the following similarity metric: 3716

$$s_i = \frac{|g_i|}{|\bigcup_{i=0}^k g_i|} \quad (6.2)$$

Where  $s_i$  is the similarity for a tissue  $i$  with  $k$  sisters;  $g_i$  refers to the set of tissues associated with tissue  $i$  and  $|g|$  refers to the cardinality of set  $g$ . For a given set of sisters, we called them redundant if they exceeded a given similarity threshold. We envisioned two possible criteria and built different dictionaries using each one. Under a threshold criterion 'any' with parameter  $S$  between  $(0, 1)$ , a given set of sisters  $j$  was considered redundant if the condition 3722

$$s_{i,j} > S \quad (6.3)$$

was true for any sister  $i$  in set  $j$ . Under a threshold criterion ‘avg’ with parameter  $S$ ,  
3723  
a given set of sisters  $j$  was considered redundant if the condition  
3724

$$\text{E}[s_i]_j > S \quad (6.4)$$

was true for the set of sisters  $j$  (see Figure 61).  
3725

Since nodes can have multiple parents (and therefore multiple sister sets), a complete  
3726  
set of similarity scores was calculated before trimming the ontology, and nodes were  
3727  
removed from the ontology if they exceeded the similarity threshold at least once in  
3728  
any comparison.  
3729

## Implementation

3730

All scripts were written in Python 3.5. Our software relies on the pandas, NumPy,  
3731  
Seaborn and SciPy modules to perform all statistical testing and data handling (McK-  
3732  
inney, 2011; Van Der Walt, Colbert, and Varoquaux, 2011; Oliphant, 2007).  
3733

## Availability of data and materials

3734

Our web implementation is available at <http://www.wormbase.org/tools/enrichment/tea/tea.cgi>. Our software can also be downloaded using Python’s  
3735  
3736  
3737  
pip installer via the command

`pip install tissue_enrichment_analysis` 3738

Alternatively, our software is available for download at: <http://dangeles.github.io/TissueEnrichmentAnalysis> 3739  
3740

All benchmark gene sets, benchmarking code and Figures can also be found at the  
3741  
same address, under the ‘tests’ folder.  
3742

<b>Abbreviations</b>	3743
• TEA—Tissue Enrichment Analysis	3744
• GO—Gene Ontology	3745
• WBbt ID—A unique ID assigned to reference ontology terms	3746
• WBgene ID—A unique ID assigned to reference nematode genes	3747
<b>Additional Files</b>	3748
<b>Additional file 1 — TEA Tutorial</b>	3749
Tutorial for users interested in using our software within a python script	3750
<b>Additional file 2 — Folder Structure for SI files 3 and 4</b>	3751
A file detailing the folder structure of the zipped folders 3 and 4.	3752
<b>Additional file 3 — Golden Gene Sets</b>	3753
A list of all the genes used for our benchmarking process.	3754
<b>Additional file 4 — Results</b>	3755
A folder containing a complete version of the results we generated for this paper.	3756
<b>References</b>	3757
Ashburner, M et al. (2000). “Gene Ontology: tool for the unification of biology”. In: <i>Nature Genetics</i> 25.1, pp. 25–29. ISSN: 1061-4036. doi: <a href="https://doi.org/10.1038/75556">10.1038/75556</a> . arXiv: <a href="https://arxiv.org/abs/10614036">10614036</a> .	3758 3759 3760
Benjamini, Yoav and Yosef Hochberg (1995). <i>Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing</i> . doi: <a href="https://doi.org/10.2307/2346101">10.2307/2346101</a> .	3761 3762
Chikina, Maria D. et al. (2009). “Global prediction of tissue-specific gene expression and context-dependent gene networks in <i>Caenorhabditis elegans</i> ”. In: <i>PLoS Computational Biology</i> 5.6. ISSN: 1553734X. doi: <a href="https://doi.org/10.1371/journal.pcbi.1000417">10.1371/journal.pcbi.1000417</a> .	3763 3764 3765 3766

- Cinar, Hulusi, Sunduz Keles, and Yishi Jin (2005). “Expression profiling of GABAergic motor neurons in *Caenorhabditis elegans*”. In: *Current Biology* 15.4, pp. 340– 346. ISSN: 09609822. doi: [10.1016/j.cub.2005.02.025](https://doi.org/10.1016/j.cub.2005.02.025). 3769
- Engelmann, Ilka et al. (2011). “A comprehensive analysis of gene expression changes provoked by bacterial and fungal infection in *C. elegans*”. In: *PLoS ONE* 6.5. ISSN: 19326203. doi: [10.1371/journal.pone.0019055](https://doi.org/10.1371/journal.pone.0019055). 3770 3771 3772
- Fox, Rebecca M et al. (2007). “The embryonic muscle transcriptome of *Caenorhabditis elegans*”. In: *Genome Biol* 8.9, R188. ISSN: 14656906. doi: [10.1186/gb-2007-8-9-r188](https://doi.org/10.1186/gb-2007-8-9-r188). 3773 3774 3775
- Garrido, Julián and Ignacio Requena (2012). “Towards summarizing knowledge: Brief ontologies”. In: *Expert Systems with Applications* 39.3, pp. 3213–3222. ISSN: 09574174. doi: [10.1016/j.eswa.2011.09.008](https://doi.org/10.1016/j.eswa.2011.09.008). 3776 3777 3778
- Gaudet, Jeb et al. (2004). “Whole-genome analysis of temporal gene expression during foregut development”. In: *PLoS Biology* 2.11. ISSN: 15449173. doi: [10.1371/journal.pbio.0020352](https://doi.org/10.1371/journal.pbio.0020352). 3779 3780 3781
- Howe, Kevin L et al. (2016). “WormBase 2016: expanding to enable helminth genomic research.” In: *Nucleic acids research* 44.November 2015, pp. D774– D780. ISSN: 1362-4962. doi: [10.1093/nar/gkv1217](https://doi.org/10.1093/nar/gkv1217). 3782 3783 3784
- Huang, Da Wei, Richard a Lempicki, and Brad T Sherman (2009). “Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources.” In: *Nature Protocols* 4.1, pp. 44–57. ISSN: 1750-2799. doi: [10.1038/nprot.2008.211](https://doi.org/10.1038/nprot.2008.211). 3785 3786 3787 3788
- Huang, Da Wei, Brad T. Sherman, et al. (2007). “DAVID Bioinformatics Resources: Expanded annotation database and novel algorithms to better extract biology from large gene lists”. In: *Nucleic Acids Research* 35.SUPPL.2. ISSN: 03051048. doi: [10.1093/nar/gkm415](https://doi.org/10.1093/nar/gkm415). 3789 3790 3791 3792
- Kelly, William G. et al. (1997). “Distinct requirements for somatic and germline expression of a generally expressed *Caenorhabditis elegans* gene”. In: *Genetics* 146.1, pp. 227–238. ISSN: 00166731. 3793 3794 3795
- Kim, Jong Woo, Jordi Conesa Caralt, and Julia K. Hilliard (2007). “Pruning bio- ontologies”. In: *Proceedings of the Annual Hawaii International Conference on System Sciences*, pp. 1–10. ISSN: 15301605. doi: [10.1109/HICSS.2007.455](https://doi.org/10.1109/HICSS.2007.455). 3796 3797 3798
- Lee, R. Y N and Paul W. Sternberg (2003). *Building a cell and anatomy ontology of Caenorhabditis elegans*. doi: [10.1002/cfg.248](https://doi.org/10.1002/cfg.248). 3799 3800
- Lee, Young Suk et al. (2013). “Ontology-aware classification of tissue and cell-type signals in gene expression profiles across platforms and technologies”. In: *Bioinformatics* 29.23, pp. 3036–3044. ISSN: 13674803. doi: [10.1093/bioinformatics/btt529](https://doi.org/10.1093/bioinformatics/btt529). 3801 3802 3803 3804

- Love, Michael I., Wolfgang Huber, and Simon Anders (2014). “Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2.” In: *Genome biology* 15.12, p. 550. ISSN: 1474-760X. doi: [10.1186/s13059-014-0550-8](https://doi.org/10.1186/s13059-014-0550-8).
- McKinney, Wes (2011). “pandas: a Foundational Python Library for Data Analysis and Statistics”. In: *Python for High Performance and Scientific Computing*, pp. 1–9.
- McLean, Cory Y et al. (2010). “GREAT improves functional interpretation of cis-regulatory regions.” In: *Nature biotechnology* 28.5, pp. 495–501. ISSN: 1087-0156. doi: [10.1038/nbt.1630](https://doi.org/10.1038/nbt.1630).
- Meisel, Joshua D. and Dennis H. Kim (2014). “Behavioral avoidance of pathogenic bacteria by *Caenorhabditis elegans*”. In: *Trends in Immunology* 35.10, pp. 465–470. ISSN: 14714981. doi: [10.1016/j.it.2014.08.008](https://doi.org/10.1016/j.it.2014.08.008).
- Mi, Huaiyu, Qing Dong, et al. (2009). “PANTHER version 7: Improved phylogenetic trees, orthologs and collaboration with the Gene Ontology Consortium”. In: *Nucleic Acids Research* 38.SUPPL.1. ISSN: 03051048. doi: [10.1093/nar/gkp1019](https://doi.org/10.1093/nar/gkp1019).
- Mi, Huaiyu, Anushya Muruganujan, and Paul D. Thomas (2013). “PANTHER in 2013: Modeling the evolution of gene function, and other gene attributes, in the context of phylogenetic trees”. In: *Nucleic Acids Research* 41.D1. ISSN: 03051048. doi: [10.1093/nar/gks1118](https://doi.org/10.1093/nar/gks1118).
- Oliphant, Travis E (2007). “SciPy: Open source scientific tools for Python”. In: *Computing in Science and Engineering* 9, pp. 10–20. ISSN: 1521-9615. doi: [10.1109/MCSE.2007.58](https://doi.org/10.1109/MCSE.2007.58).
- Pathan, Mohashin et al. (2015). “FunRich: An open access standalone functional enrichment and interaction network analysis tool”. In: *Proteomics* 15.15, pp. 2597–2601. ISSN: 16159861. doi: [10.1002/pmic.201400515](https://doi.org/10.1002/pmic.201400515).
- Pauli, Florencia et al. (2006). “Chromosomal clustering and GATA transcriptional regulation of intestine-expressed genes in *C. elegans*.” In: *Development (Cambridge, England)* 133.2, pp. 287–295. ISSN: 0950-1991. doi: [10.1242/dev.02185](https://doi.org/10.1242/dev.02185).
- Pawitan, Yudi et al. (2005). “False discovery rate, sensitivity and sample size for microarray studies”. In: *Bioinformatics* 21.13, pp. 3017–3024. ISSN: 13674803. doi: [10.1093/bioinformatics/bti448](https://doi.org/10.1093/bioinformatics/bti448).
- Portman, Douglas S. and Scott W. Emmons (2004). “Identification of *C. elegans* sensory ray genes using whole-genome expression profiling”. In: *Developmental Biology* 270.2, pp. 499–512. ISSN: 00121606. doi: [10.1016/j.ydbio.2004.02.020](https://doi.org/10.1016/j.ydbio.2004.02.020).

- Prykhozhij, Sergey V, Annalisa Marsico, and Sebastiaan H Meijssing (2013). “Zebrafish Expression Ontology of Gene Sets (ZEOGS): a tool to analyze enrichment of zebrafish anatomical terms in large gene sets.” In: *Zebrafish* 10.3, pp. 303–15. ISSN: 1557-8542. doi: [10.1089/zeb.2012.0865](https://doi.org/10.1089/zeb.2012.0865). 3842  
3843  
3844  
3845
- Smith, Cody J. et al. (2010). “Time-lapse imaging and cell-specific expression profiling reveal dynamic branching and molecular determinants of a multi-dendritic nociceptor in *C. elegans*”. In: *Developmental Biology* 345.1, pp. 18–33. ISSN: 00121606. doi: [10.1016/j.ydbio.2010.05.502](https://doi.org/10.1016/j.ydbio.2010.05.502). 3846  
3847  
3848  
3849
- Spencer, W. Clay et al. (2011). “A spatial and temporal map of *C. elegans* gene expression”. In: *Genome Research* 21.2, pp. 325–341. ISSN: 10889051. doi: [10.1101/gr.114595.110](https://doi.org/10.1101/gr.114595.110). 3850  
3851  
3852
- Storey, John D and Robert Tibshirani (2003). “Statistical significance for genomewide studies.” In: *Proceedings of the National Academy of Sciences of the United States of America* 100.16, pp. 9440–5. ISSN: 0027-8424. doi: [10.1073/pnas.1530509100](https://doi.org/10.1073/pnas.1530509100). 3853  
3854  
3855  
3856
- The Gene Ontology Consortium (2015). “Gene Ontology Consortium: going forward”. In: *Nucleic Acids Research* 43.D1, pp. D1049–D1056. ISSN: 0305-1048. doi: [10.1093/nar/gku1179](https://doi.org/10.1093/nar/gku1179). 3857  
3858  
3859
- Van Der Walt, Stéfan, S. Chris Colbert, and Gaël Varoquaux (2011). “The NumPy array: A structure for efficient numerical computation”. In: *Computing in Science and Engineering* 13.2, pp. 22–30. ISSN: 15219615. doi: [10.1109/MCSE.2011.37](https://doi.org/10.1109/MCSE.2011.37). arXiv: [1102.1523](https://arxiv.org/abs/1102.1523). 3860  
3861  
3862  
3863
- Watson, Joseph D et al. (2008). “Complementary RNA amplification methods enhance microarray identification of transcripts expressed in the *C. elegans* nervous system”. In: *BMC Genomics* 9, p. 84. ISSN: 1471-2164. doi: [10.1186/1471-2164-9-84](https://doi.org/10.1186/1471-2164-9-84). 3864  
3865  
3866  
3867
- Zhang, Y, H Lu, and C I Bargmann (2005). “Pathogenic bacteria induce aversive olfactory learning in *Caenorhabditis elegans*”. In: *Nature* 438.7065, pp. 179–184. ISSN: 0028-0836. doi: [10.1038/nature04216](https://doi.org/10.1038/nature04216). 3868  
3869  
3870  
3871

TWO NEW FUNCTIONS IN THE WORMBASE ENRICHMENT 3872  
SUITE 3873

7.1 Description 3874

Genome-wide experiments routinely generate large amounts of data that can be 3875 hard to interpret biologically. A common approach to interpreting these results 3876 is to employ enrichment analyses of controlled languages, known as ontologies, 3877 that describe various biological parameters such as gene molecular or biological 3878 function. In *C. elegans*, three distinct ontologies, the Gene Ontology (GO), Anatomy 3879 Ontology (AO), and the Worm Phenotype Ontology (WPO) are used to annotate 3880 gene function, expression and phenotype, respectively (Ashburner et al., 2000; 3881 Lee and Sternberg, 2003; Schindelman et al., 2011). Previously, we developed 3882 software to test datasets for enrichment of anatomical terms, called the Tissue 3883 Enrichment Analysis (TEA) tool (Angeles-Albores et al., 2016). Using the same 3884 hypergeometric statistical method, we extend enrichment testing to include WPO 3885 and GO, offering a unified approach to enrichment testing in *C. elegans*. The 3886 WormBase Enrichment Suite can be accessed via a user-friendly interface at <http://www.wormbase.org/tools/enrichment/tea/tea.cgi>. 3888

To validate the tools, we analyzed a previously published extracellular vesicle 3889 (EV)-releasing neuron (EVN) signature gene set derived from dissociated ciliated 3890 EV neurons(Wang et al., 2015) using the WormBase Enrichment Suite based on 3891 the WS262 WormBase release. TEA correctly identified the CEM, hook sensil- 3892 lum and IL2 neuron as enriched tissues. The top phenotype associated with the 3893 EVN signature was chemosensory behavior. Gene Ontology enrichment analysis 3894 showed that cell projection and cell body were the most enriched cellular com- 3895

ponents in this gene set, followed by the biological processes neuropeptide sig- 3896  
naling pathway and vesicle localization further down. The tutorial script used to 3897  
generate the figure above can be viewed at: <https://github.com/dangeles/TissueEnrichmentAnalysis/blob/master/tutorial/Tutorial.ipynb> 3898

The addition of Gene Enrichment Analysis (GEA) and Phenotype Enrichment Anal- 3900  
ysis (PEA) to WormBase marks an important step towards a unified set of analyses 3901  
that can help researchers to understand genomic datasets. These enrichment anal- 3902  
yses will allow the community to fully benefit from the data curation ongoing at 3903  
WormBase. 3904

## Methods

3905

Using the methods described in Angeles-Albores et al. (2016), we generated on- 3906  
tology dictionaries using the Anatomy, Phenotype and Gene Ontology annotations 3907  
for *C. elegans*. The dictionary similarity parameter was set to 95 for all ontolo- 3908  
gies. The annotation per term minimum was set to 33 annotations for the AO, 3909  
a 50 annotations for the WPO, and 33 annotations for GO. Terms within the 3910  
dictionary are tested using a hypergeometric probability test and corrected using 3911  
the Benjamini-Hochberg step-up algorithm. In WS262, there are 1320 anatomy 3912  
terms, 1117 phenotypes, and 3025 GO terms that have at least 11 genes annotated 3913  
to them. The dictionaries are freely accessible using the Python version of the 3914  
Suite, which can be installed using the pip tool for Python libraries: `pip install 3915  
tissue_enrichment_analysis`. The dictionary can then be automatically down- 3916  
loaded by importing the enrichment analysis library in a Python script by writing 3917  
`import tissue_enrichment_analysis as ea`. The dictionaries can then be 3918  
downloaded by typing `ea.fetch_dictionary(dict)` into Python, where ‘dict’ 3919  
is one of the strings ‘tissue’, ‘phenotype’ or ‘go’ to specify which dictionary to 3920  
download. If the function does not receive an argument, the dictionary correspond- 3921

ing to the AO is downloaded by default. See the tutorial above for an example 3922  
implementation. 3923

## References 3924

- Angeles-Albores, David et al. (2016). “Tissue enrichment analysis for *C. elegans* 3925 genomics”. In: *BMC Bioinformatics* 17.1, p. 366. ISSN: 1471-2105. doi: [10.1186/s12859-016-1229-9](https://doi.org/10.1186/s12859-016-1229-9). 3926  
3927
- Ashburner, M et al. (2000). “Gene Ontology: tool for the unification of biology”. In: 3928 *Nature Genetics* 25.1, pp. 25–29. ISSN: 1061-4036. doi: [10.1038/75556](https://doi.org/10.1038/75556). arXiv: 3929  
[10614036](https://arxiv.org/abs/10614036). 3930
- Lee, R. Y N and Paul W. Sternberg (2003). *Building a cell and anatomy ontology* 3931  
of *Caenorhabditis elegans*. doi: [10.1002/cfg.248](https://doi.org/10.1002/cfg.248). 3932
- Schindelman, Gary et al. (2011). “Worm Phenotype Ontology: integrating pheno- 3933 type data within and beyond the *C. elegans* community.” In: *BMC bioinformatics* 3934 12, p. 32. ISSN: 1471-2105. doi: [10.1186/1471-2105-12-32](https://doi.org/10.1186/1471-2105-12-32). 3935
- Wang, Juan et al. (2015). “Cell-Specific Transcriptional Profiling of Ciliated Sensory 3936 Neurons Reveals Regulators of Behavior and Extracellular Vesicle Biogenesis”. 3937 In: *Current Biology* 25.24, pp. 3232–3238. ISSN: 09609822. doi: [10.1016/j.cub.2015.10.057](https://doi.org/10.1016/j.cub.2015.10.057). 3938  
3939

## CONCLUSION

3940

I have tried to demonstrate that transcriptomes are phenotypes, and I have tried to 3941 show examples of how these phenotypes can be rigorously manipulated. The work 3942 is not without flaws, and some of it may even be wrong. However, I believe the 3943 principles of the work are correct: Batesonian epistasis exists in transcriptomes and 3944 should be a quantity of immense interest to us because it does not require a null 3945 hypothesis; allelic series must be explored using transcriptome-wide dominance; 3946 and we must continue to develop tools that make use of the enormous amount of 3947 information that the scientific community is actively generating. It is my hope that 3948 these concepts prove useful to the greater scientific community. 3949