

# Analysis of allelic series with transcriptomic phenotypes

David Angeles-Albores<sup>1</sup> and Paul W. Sternberg<sup>1,\*</sup>

<sup>1</sup>*Division of Biology and Biological Engineering, Caltech, Pasadena, CA, 91125, USA*

\**Corresponding author. Contact: [pws@caltech.edu](mailto:pws@caltech.edu)*

January 11, 2018

Although transcriptomes have recently been used to perform epistasis analyses, they are not yet used to study intragenic function/structure relationships. We developed a theoretical framework to study allelic series using transcriptomic phenotypes. As a proof-of-concept, we apply our methods to an allelic series of *dpy-22*, a highly pleiotropic *Caenorhabditis elegans* gene orthologous to the human gene *MED12*, which is a subunit of the Mediator complex. Our methods identify functional regions within *dpy-22* that modulate Mediator activity upon various genetic modules.

## 1 Introduction

2 Mutations of a gene can yield a series of alleles with  
3 different phenotypes that reveal multiple functions  
4 encoded by that gene, regardless of the alleles' molec-  
5 ular nature. Homozygous alleles can be ordered by  
6 their phenotypic severity; then, phenotypes of *trans*-  
7 heterozygotes carrying two alleles can reveal which  
8 alleles are dominant for each phenotype. Together,  
9 the severity and dominance hierarchies show intra-  
10 genic functional regions. In *Caenorhabditis elegans*,  
11 these series have helped characterize genes such as  
12 *let-23/EGFR*, *lin-3/EGF* and *lin-12/NOTCH*<sup>1,2,3</sup>.

13 Biology has moved from expression measurements  
14 of single genes towards genome-wide measurements.  
15 Expression profiling via RNA-seq<sup>4</sup> enables simulta-  
16 neous measurement of transcript levels for all genes  
17 in a genome, yielding a transcriptome. These mea-  
18 surements can be made on whole organisms, isolated  
19 tissues, or single cells<sup>5,6</sup>. Transcriptomes have been  
20 successfully used to identify new cell or organismal  
21 states<sup>7,8</sup>. For mutant genes, transcriptomic states  
22 can be used for epistasis analysis<sup>9,10</sup>, but have not  
23 been used to characterize allelic series.

24 We have devised methods for characterizing alle-  
25 lic series with RNA-seq. To test these methods,  
26 we selected three alleles<sup>11,12</sup> of a *C. elegans* Medi-  
27 ator complex subunit gene, *dpy-22*. Mediator is a  
28 macromolecular complex with  $\sim 25$  subunits<sup>13</sup> that  
29 globally regulates RNA polymerase II (Pol II)<sup>14,15</sup>.  
30 The Mediator complex has at least four biochemically  
31 distinct modules: the Head, Middle and Tail mod-  
32 ules and a CDK-8-associated Kinase Module (CKM).

The CKM associates reversibly with other modules, 33  
and appears to inhibit transcription<sup>16,17</sup>. In *C. el-* 34  
*elegans* development, the CKM promotes both male 35  
tail formation<sup>11</sup> (through interactions with the Wnt 36  
pathway), and vulval formation<sup>18</sup> (through inhibi- 37  
tion of the Ras pathway). Homozygotes of allele 38  
*dpy-22(bx93)*, which encodes a premature stop codon 39  
Q2549Amber<sup>11</sup>, appear grossly wild-type. In con- 40  
trast, animals homozygous for a more severe allele, 41  
*dpy-22(sy622)* encoding another premature stop 42  
codon, Q1698Amber<sup>12</sup>, are dumpy (Dpy), have egg- 43  
laying defects (Egl), and have multiple vulvae (Muv). 44  
(see Fig. 1A). In spite of its causative role in a num- 45  
ber of neurodevelopmental disorders<sup>19</sup>, the structural 46  
and functional features of this gene are poorly un- 47  
derstood. In humans, MED12 is known to have a 48  
proline-, glutamine- and leucine-rich domain that in- 49  
teracts with the WNT pathway<sup>20</sup>. However, many 50  
disease-causing variants fall outside of this domain<sup>21</sup>. 51  
To study these variants and how they interfere with 52  
the functionality of *MED12*, quantitative and effi- 53  
cient methods are necessary. 54

RNA-seq phenotypes have the potential to reveal 55  
functional regions within genes, but their phenotypic 56  
complexity makes this difficult. We developed a 57  
method for determining allelic series from transcrip- 58  
tomic phenotypes and used the *C. elegans dpy-22* 59  
gene as a test case. Our analysis revealed functional 60  
regions that act to modulate Mediator activity at 61  
thousands of genetic loci. 62

## Results and Discussion

We adapted the allelic series method, previously used for individual phenotypes, for use with expression profiles as multidimensional phenotypes (see Fig. 1). As a proof of principle, 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 quadruplicates from *trans*-heterozygotes of both alleles. Sequencing was performed at a depth of 20 million reads per sample. Reads were pseudoaligned using Kallisto<sup>22</sup>. We performed a differential expression using a general linear model specified using Sleuth<sup>23</sup> (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 the false discovery rate,  $q$ , was less than or equal to 10%. 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>.

By these criteria, we found 481 genes differentially expressed in *dpy-22(bx93)* homozygotes, and 2,863 differentially expressed genes in *dpy-22(sy622)* homozygotes (see [Basic Statistics Notebook](#)). *Trans*-heterozygotes with the genotype *dpy-6(e14) dpy-22(bx93)/+ dpy-22(sy622)* had 2,214 differentially expressed genes with respect to the wild type.

We used a false hit analysis to identify four non-overlapping phenotypic classes. We use the term genotype-specific to refer to groups of transcripts that were perturbed in one mutant. We use the term genotype-associated to refer to those groups of transcripts whose expression was significantly altered in two or more mutants with respect to the wild type control. The ***dpy-22(sy622)*-associated** phenotypic class consisted of 720 genes differentially expressed in *dpy-22(sy622)* homozygotes and in *trans*-heterozygotes, but which had wild-type expression in *dpy-22(bx93)* homozygotes. The ***dpy-22(bx93)*-associated** phenotypic class contains 403 genes differentially expressed in all genotypes. We also identified a ***dpy-22(sy622)*-specific** phenotypic class (1,841 genes) and a ***trans*-heterozygote-specific** phenotypic class (1,226 genes; see the [Phenotypic Classes Notebook](#)). All genotype-associated phenotypes had Spearman rank correlations  $> 0.8$ , indicating that transcripts within these classes changed in

Phenotypic Class	Dominance
<i>dpy-22(sy622)</i> -specific	$1.00 \pm 0.00$
<i>dpy-22(sy622)</i> -associated	$0.51 \pm 0.01$
<i>dpy-22(bx93)</i> -associated	$0.81 \pm 0.01$

**Table 1.** 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)*.

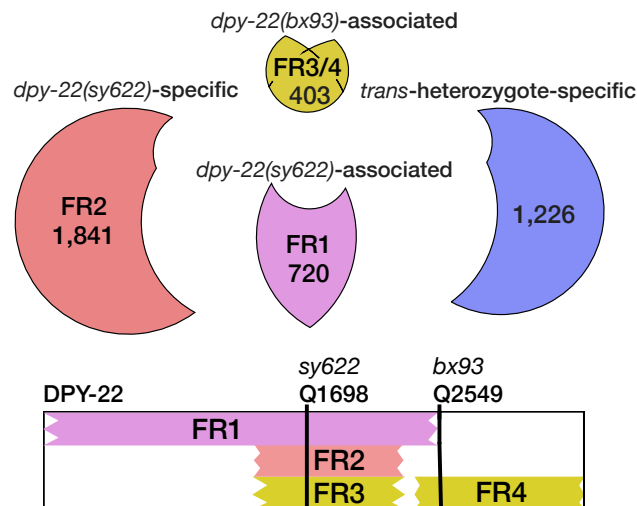
the same direction amongst the genotypes studied.

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. The *dpy-22(sy622)* and *dpy-22(bx93)* alleles are semidominant ( $d_{bx93} = 0.51$ ) to each other for the *dpy-22(sy622)*-associated phenotypic class. The *dpy-22(bx93)* allele is largely dominant over the *dpy-22(sy622)* allele ( $d_{bx93} = 0.81$ ; see Table 1) for the *dpy-22(bx93)*-associated phenotypic class.

Because the mutations we used are truncations, our results suggest the existence of various functional regions in *dpy-22/MDT12* (see Fig. 2). The *dpy-22(sy622)*-specific phenotypic class is likely controlled by a single functional region, functional region 1 (FC1), and the *dpy-22(sy622)*-associated phenotypic class is likely controlled by a second functional region, functional region 2 (FC2). It is unlikely that these regions are identical because their dominance behaviors are very different. The *dpy-22(bx93)* allele was largely dominant over the *dpy-22(sy622)* allele for the *dpy-22(bx93)*-associated class, but gene expression in this class was perturbed in both homozygotes. The perturbations were greater for *dpy-22(sy622)* homozygotes than for *dpy-22(bx93)* homozygotes. This behavior can be explained if the *dpy-22(bx93)*-associated class is controlled jointly by two distinct effectors, functional regions 3 and 4 (FC3, FC4, see Fig. 2). A rigorous examination of this model will require studying alleles that mutate the region between Q1689 and Q2549 using homozygotes and *trans*-heterozygotes.

We also found a class of transcripts that had perturbed levels in *trans*-heterozygotes only; its biological significance is unclear. Phenotypes unique to *trans*-heterozygotes are often the result of physical interactions such as homodimerization, or dosage reduction of a toxic product<sup>24</sup>. In the case of





**Figure 2.** 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. 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 they were merged all classes together (see [Methods](#)).

*dpy-22/MDT12* orthologs, how either mechanism could operate is not obvious, since DPY-22 is expected to assemble in a monomeric manner into the CKM. Massive single-cell RNA-seq of *C. elegans* has recently been reported<sup>25</sup>. When this technique becomes cost-efficient, single-cell profiling of these genotypes may provide information that complements the whole-organism expression phenotypes, perhaps explaining the origin of this phenotype.

Intragenic mapping of functional regions associated with phenotypic classes is important, but their biological meaning remains unclear. To assign biological functionality to phenotypic classes, we extracted transcriptomic signatures associated with a Dumpy (Dpy) phenotype using transcriptomes from *dpy-7* and *dpy-10* mutants (DAA, CPR and PWS unpublished), and a *hif-1*-dependent hypoxia response from a previously published analysis<sup>10</sup> and asked whether any phenotypic class was enriched in either response. The *sy622*-specific and -associated classes were enriched in genes that are transcriptionally associated with a Dpy phenotype (fold-change enrichment = 3,  $p = 2 \cdot 10^{-40}$ , 167 genes observed; fold-change = 1.9,  $p = 9 \cdot 10^{-9}$ , 82 genes observed). The *bx93*-associated class also showed significant enrichment (fold-change = 2.2,  $p = 4 \cdot 10^{-10}$ , 68 genes observed). The class that showed the most extreme deviation from random was the *sy622*-specific class. *dpy-22(sy622)* homozygotes are severely Dpy, whereas *dpy-22(bx93)* homozygotes and *trans*-heterozygotes have a slight Dpy phenotype. Plotting the changes in gene expression for *sy622* homozygotes versus the changes in expression in *dpy-7* mutants revealed that 75% of the transcripts were strongly correlated in both genotypes (see Figure 3). Therefore, the *sy622*-specific phenotypic class contains a transcriptional signature associated with morphological Dpy phenotype (see the [Enrichment Notebook](#)).

*dpy-22* is not known to be upstream of the *hif-1*-dependent hypoxia response in *C. elegans*. Enrichment tests revealed that the hypoxia response was significantly enriched in the *bx93*-associated (fold-change = 2.1,  $p = 10^{-8}$ , 63 genes observed), the *sy622*-associated (fold-change = 1.9,  $p = 4 \cdot 10^{-8}$ , 78 genes observed) and the *sy622*-specific classes (fold-change = 2.4,  $p = 9 \cdot 10^{-55}$ , 186 genes observed). However, there was no correlation between the expression levels of these genes in *dpy-22* genotypes and the expression levels expected from the hypoxia response. Although the hypoxia gene battery can be found in *dpy-22* mutants, these genes are not used to deploy a *hif-1*-dependent hypoxia phenotype. Taken together, our results suggest that transcriptomic signatures can be used to understand the biological func-

212 tionality of phenotypic classes, and they may be use-  
213 ful in associating phenotypic classes with other phe-  
214 notypes. This highlights the importance of gener-  
215 ating an index set of mutants that can be used to  
216 derive a gold standard of transcriptional signatures  
217 with which to test future results.

218 Transcriptomic phenotypes generate large amounts  
219 of differential gene expression data, so false positive  
220 and false negative rates can lead to spurious phe-  
221 notypic classes whose putative biological significance  
222 is badly misleading. Such artifacts are particularly  
223 likely for small phenotypic classes, which should be  
224 viewed with skepticism. Notably, errors of interpreta-  
225 tion cannot be avoided by setting a more stringent  $q$ -  
226 value cut-off: doing so will decrease the false positive  
227 rate, but increase the false negative rate, which will  
228 in turn produce smaller phenotypic classes than ex-  
229 pected. Our method avoids this pitfall by using total  
230 error rate estimates to assess the plausibility of each  
231 class. These conclusions are of broad significance to  
232 research where highly multiplexed measurements are  
233 compared to identify similarities and differences in  
234 the genome-wide behavior of a single variable under  
235 multiple conditions.

236 We have shown that transcriptomes can be used  
237 to study allelic series in the context of a large,  
238 pleiotropic gene. We identified separable phenotypic  
239 classes that would otherwise be obscured by other  
240 methods, correlated each class to a functional region,  
241 and identified sequence requirements for each region.  
242 Given the importance of allelic series for character-  
243 izing gene function and their roles in specific genetic  
244 pathways, we are optimistic that this method will be  
245 a useful addition to the geneticist's arsenal.

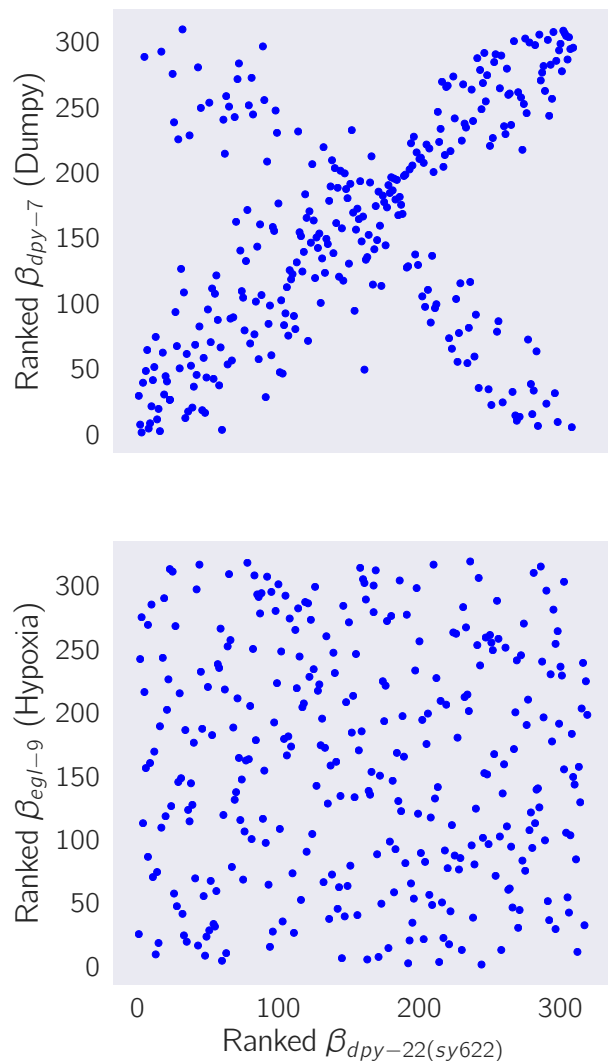
## 246 Methods

### 247 Strains used

248 Strains used were N2 wild-type (Bristol), PS4087  
249 *dpy-22(sy622)*, PS4187 *dpy-22(bx93)*, and PS4176  
250 *dpy-6(e14) dpy-22(bx93) + dpy-22(sy622)*. Lines  
251 were grown on standard nematode growth media  
252 (NGM) Petri plates seeded with OP50 *E. coli* at  
253 20°C<sup>26</sup>.

### 254 Strain synchronization, harvesting and 255 RNA sequencing

256 Strains were synchronized by bleaching P<sub>0</sub>'s into vir-  
257 gin S. basal (no cholesterol or ethanol added) for 8–  
258 12 hours. Arrested L1 larvae were placed in NGM  
259 plates seeded with OP50 at 20°C and grown to



**Figure 3.** *sy622* homozygotes show a transcriptional response associated with the Dpy phenotype. **A** We obtained a set of transcripts associated with the Dpy phenotype from *dpy-7* and *dpy-10* mutants. We identified the transcripts that were differentially expressed in *sy622* homozygotes. We ranked the  $\beta$  values of each transcript in *sy622* homozygotes and plotted them against the ranked  $\beta$  values in *dpy-7* mutants. A significant portion of the genes are correlated between the two genotypes, showing that the signature is largely intact. 25% of the genes are anti-correlated. **B** We performed the same analysis using a set of transcripts associated with the *hif-1*-dependent hypoxia response as a negative control. Although *sy622* is enriched for the transcripts that make up this response, there is no correlation between the  $\beta$  values in *sy622* homozygotes and the  $\beta$  values in *egl-9* homozygotes.



260 the young adult stage (assessed by vulval morphol- 310  
261 ogy and lack of embryos). RNA extraction and se- 311  
262 quencing was performed as previously described by 312  
263 Angeles-Albores *et al*<sup>10,7</sup>. 313

## 264 Read pseudo-alignment and differential 314 265 expression 315

266 Reads were pseudo-aligned to the *C. elegans* genome 316  
267 (WBcel235) using Kallisto<sup>22</sup>, using 200 bootstraps 317  
268 and with the sequence bias (`--seqBias`) flag. The 318  
269 fragment size for all libraries was set to 200 319  
270 and the standard deviation to 40. Quality control was 320  
271 performed on a subset of the reads using FastQC, 321  
272 RNAseQC, BowTie and MultiQC<sup>27,28,29,30</sup>. 322

273 Differential expression analysis was performed using 323  
274 Sleuth<sup>23</sup>. We used a general linear model to iden- 324  
275 tify genes that were differentially expressed between 325  
276 wild-type and mutant libraries. To increase our sta- 326  
277 tistical power, we pooled young adult wild-type repli- 327  
278 cates from other published<sup>10,7</sup> and unpublished anal- 328  
279 yses adjusting for batch effects. 329

## 280 False hit analysis 330

281 To accurately count phenotypes, we developed a false 331  
282 hit algorithm (Algorithm 1). We implemented this 332  
283 algorithm for three-way comparisons in Python. Al- 333  
284 though experimentally restricted, a three-way com- 334  
285 parison can result in  $> 5,000$  possible sets (ignoring 335  
286 size). This large number of models necessitates an 336  
287 algorithmic approach that can at least restrict the 337  
288 possible number of models. Our algorithm uses a 338  
289 noise function that assumes false hit events are non- 339  
290 overlapping (i.e. the same gene cannot be the result 340  
291 of two false positive events in two or more genotypes) 341  
292 to determine the average noise flux between pheno- 342  
293 typic classes. These assumptions break down rapidly 343  
294 if false-positive or negative rates exceed 20%. 344

295 To benchmark our algorithm, we generated one 345  
296 thousand Venn diagrams at random. For each Venn 346  
297 diagram, we calculated the average false positive and 347  
298 false negative flux matrices. Then, we added noise 348  
299 to each phenotypic class in the Venn diagram, as- 349  
300 suming that fluxes were normally distributed with 350  
301 mean and standard deviation equal to the flux co- 351  
302 efficient calculated. We input the noised Venn dia- 352  
303 gram into our false hit analysis and collected clas- 353  
304 sification statistics. For a given signal-to-noise cut- 354  
305 off,  $\lambda$ , classification accuracy varied significantly with 355  
306 changes in the total error rate. In the absence of 356  
307 false negative hits, false hit analysis can accurately 357  
308 identify non-empty genotype-associated phenotypic 358  
309 classes, but identifying genotype-specific classes be-

comes difficult if the experimental false positive rate 310  
is high. On the other hand, even moderate false 311  
negative rates ( $> 10\%$ ) rapidly degrade signal from 312  
genotype-associated classes. For classes that are as- 313  
sociated with three genotypes, an experimental false 314  
negative rate of 30% is enough on average to prevent 315  
this class from being observed. 316

We selected  $\lambda = 3$  because classification using this 317  
threshold was high across a range of false positive and 318  
false negative combinations. A challenge to applying 319  
this algorithm to our data is the fact that the false 320  
negative rate for our experiment is unknown. Al- 321  
though there has been significant progress in control- 322  
ling and estimating false positive rates, we know of no 323  
such attempts for false negative rates. It is unlikely 324  
that the false negative rate for our study is lower than 325  
the false positive rate, because all genotypes except 326  
the controls are likely underpowered. We used false 327  
negative rates between 10–20% for false hit analy- 328  
sis. When the false negative rate was set at 15% 329  
or higher, the algorithm converged on the same five 330  
classes shown above. For false negative rates between 331  
10–15%, the algorithm output the same five classes, 332  
but also accepted the (*dpy-22(sy622),dpy-22(bx93)*)- 333  
associated class. We selected the model correspond- 334  
ing to false negative rates of 15–20% because this 335  
model had lower  $\chi^2$  values than the model selected 336  
with a false negative rate of 10–15% (4,212 versus 337  
100,650). 338

We asked whether re-classification of some classes 339  
into others could improve model fit. We manually 340  
re-classified the (*dpy-22(sy622),dpy-22(bx93)*)- 341  
associated and the (*dpy-22(bx93), trans-* 342  
*heterozygote*)-associated classes into the *bx93*- 343  
associated class (which is associated with all 344  
genotypes), and we compared  $\chi^2$  statistics between 345  
a re-classified reduced model and a reduced model. 346  
The re-classified model had a lower  $\chi^2$  (181). Thus, 347  
we concluded that the re-classified reduced model is 348

349 the most likely model to give rise to our data.

**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 set to initialize the reduced model*

$\mathbf{K} = \{\min_{l \in L} N_l\}$

*Refine the model until the model converges 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**

*Return the reduced model*

$\mathbf{M}_{reduced} = \mathbf{K}$

**return**  $\mathbf{M}_{reduced}$

**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.

## 351 Dominance analysis

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

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

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  $a$ .

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

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

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.

## Code

Code was written in Jupyter notebooks<sup>31</sup> using the Python programming language. The Numpy, pandas and scipy libraries were used for computation<sup>32,33,34</sup> and the matplotlib and seaborn libraries were used for data visualization<sup>35,36</sup>. Enrichment analyses were performed using the WormBase Enrichment Suite<sup>37</sup>. 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.

## Data Availability

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)<sup>38</sup> under the accession code GSE107523 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE107523>).

## Acknowledgements

This work was supported by HHMI with whom PWS was an investigator, by the Millard and Muriel Jacobs Genetics and Genomics Laboratory at California Institute of Technology, and by the NIH grant

392 U41 HG002223. This article would not be possible  
393 without help from Dr. Igor Antoshechkin and Dr.  
394 Vijaya Kumar who performed the library prepara-  
395 tion and sequencing. We would like to thank Carmie  
396 Puckett Robinson for the unpublished Dpy tran-  
397 scriptional signature. Han Wang, Hillel Schwartz,  
398 Erich Schwarz, Porfirio Quintero and Carmie Puck-  
399 ett Robinson provided valuable input throughout the  
400 project.

## 401 References

- 402 1. Aroian, R. V. & Sternberg, P. W. Multiple  
403 functions of let-23, a *Caenorhabditis elegans* re-  
404 ceptor tyrosine kinase gene required for vulval  
405 induction. *Genetics* **128**, 251–67 (1991).
- 406 2. Ferguson, E. & Horvitz, H. R. Identification  
407 and characterization of 22 genes that affect the  
408 vulval cell lineages of *Caenorhabditis elegans*.  
409 *Genetics* **110**, 17–72 (1985).
- 410 3. Greenwald, I. S., Sternberg, P. W. & Robert  
411 Horvitz, H. The lin-12 locus specifies cell fates  
412 in *Caenorhabditis elegans*. *Cell* **34**, 435–444  
413 (1983).
- 414 4. Mortazavi, A., Williams, B. A., McCue, K.,  
415 Schaeffer, L. & Wold, B. Mapping and quanti-  
416 fying mammalian transcriptomes by RNA-Seq.  
417 *Nature Methods* **5**, 621–628 (2008).
- 418 5. Tang, F. *et al.* mRNA-Seq whole-  
419 transcriptome analysis of a single cell.  
420 *Nature Methods* **6**, 377–382 (2009).
- 421 6. Schwarz, E. M., Kato, M. & Sternberg, P. W.  
422 Functional transcriptomics of a migrating cell  
423 in *Caenorhabditis elegans*. *Proceedings of the*  
424 *National Academy of Sciences of the United*  
425 *States of America* **109**, 16246–51 (2012).
- 426 7. Angeles-Albores, D. *et al.* The *Caenorhabdi-*  
427 *tis elegans* Female State: Decoupling the Tran-  
428 scriptional Effects of Aging and Sperm-Status.  
429 *G3: Genes, Genomes, Genetics* (2017).
- 430 8. Villani, A.-C. *et al.* Single-cell RNA-seq re-  
431 veals new types of human blood dendritic  
432 cells, monocytes, and progenitors. *Science* **356**  
433 (2017).
- 434 9. Dixit, A. *et al.* Perturb-Seq: Dissecting Molec-  
435 ular Circuits with Scalable Single-Cell RNA  
436 Profiling of Pooled Genetic Screens. *Cell* **167**,  
437 1853–1866.e17 (2016).
10. Angeles Albores, D., Puckett Robinson, C.,  
Williams, B. A., Wold, B. J. & Sternberg,  
P. W. Reconstructing a metazoan genetic path-  
way with transcriptome-wide epistasis mea-  
surements. *bioRxiv* (2017).
11. Zhang, H. & Emmons, S. W. A *C. elegans*  
mediator protein confers regulatory selectivity  
on lineage-specific expression of a transcription  
factor gene. *Genes and Development* **14**, 2161–  
2172 (2000).
12. Moghal, N. A component of the tran-  
scriptional mediator complex inhibits RAS-  
dependent vulval fate specification in *C. ele-*  
*gans*. *Development* **130**, 57–69 (2003).
13. Jeronimo, C. & Robert, F. The Mediator Com-  
plex: At the Nexus of RNA Polymerase II  
Transcription (2017).
14. Allen, B. L. & Taatjes, D. J. The Mediator  
complex: a central integrator of transcription.  
*Nature reviews. Molecular cell biology* **16**, 155–  
166 (2015).
15. Takagi, Y. & Kornberg, R. D. Mediator as a  
general transcription factor. *The Journal of*  
*biological chemistry* **281**, 80–9 (2006).
16. Knuesel, M. T., Meyer, K. D., Bernecky, C. &  
Taatjes, D. J. The human CDK8 subcomplex  
is a molecular switch that controls Mediator  
coactivator function. *Genes & development* **23**,  
439–51 (2009).
17. Elmlund, H. *et al.* The cyclin-dependent ki-  
nase 8 module sterically blocks Mediator in-  
teractions with RNA polymerase II. *Proceed-*  
*ings of the National Academy of Sciences of*  
*the United States of America* **103**, 15788–93  
(2006).
18. Moghal, N. & Sternberg, P. W. A compo-  
nent of the transcriptional mediator complex  
inhibits RAS-dependent vulval fate specifica-  
tion in *C. elegans*. *Development* **130**, 57–69  
(2003).
19. Graham, J. M. & Schwartz, C. E. MED12 re-  
lated disorders. *American Journal of Medical*  
*Genetics, Part A* **161**, 2734–2740 (2013).
20. Kim, S., Xu, X., Hecht, A. & Boyer, T. G.  
Mediator is a transducer of Wnt/ $\beta$ -catenin sig-  
naling. *Journal of Biological Chemistry* **281**,  
14066–14075 (2006).



- 
- 485 21. Yamamoto, T. & Shimojima, K. A novel  
486 MED12 mutation associated with non-specific  
487 X-linked intellectual disability. *Human*  
488 *Genome Variation* **2**, 15018 (2015).
- 489 22. Bray, N. L., Pimentel, H. J., Melsted, P. &  
490 Pachter, L. Near-optimal probabilistic RNA-  
491 seq quantification. *Nature biotechnology* **34**,  
492 525–7 (2016).
- 493 23. Pimentel, H., Bray, N. L., Puente, S., Mel-  
494 sted, P. & Pachter, L. Differential analysis  
495 of RNA-seq incorporating quantification un-  
496 certainty. *brief communications nature meth-*  
497 *ods* **14** (2017).
- 498 24. Yook, K. Complementation. *WormBook*  
499 (2005).
- 500 25. Cao, J. *et al.* Comprehensive single-cell tran-  
501 scriptional profiling of a multicellular organ-  
502 ism. *Science (New York, N.Y.)* **357**, 661–667  
503 (2017).
- 504 26. Brenner, S. The Genetics of *Caenorhabditis el-*  
505 *egans*. *Genetics* **77**, 71–94 (1974).
- 506 27. Andrews, S. FastQC: A quality control tool for  
507 high throughput sequence data (2010).
- 508 28. Deluca, D. S. *et al.* RNA-SeQC: RNA-seq met-  
509 rics for quality control and process optimiza-  
510 tion. *Bioinformatics* **28**, 1530–1532 (2012).
- 511 29. Langmead, B., Trapnell, C., Pop, M. &  
512 Salzberg, S. L. Bowtie: An ultrafast memory-  
513 efficient short read aligner. *Genome biology*  
514 **R25** (2009).
- 515 30. Ewels, P., Magnusson, M., Lundin, S. & Käller,  
516 M. MultiQC: Summarize analysis results for  
517 multiple tools and samples in a single report.  
518 *Bioinformatics* **32**, 3047–3048 (2016).
- 519 31. Pérez, F. & Granger, B. IPython: A System  
520 for Interactive Scientific Computing Python:  
521 An Open and General- Purpose Environment.  
522 *Computing in Science and Engineering* **9**, 21–  
523 29 (2007).
- 524 32. Van Der Walt, S., Colbert, S. C. & Varoquaux,  
525 G. The NumPy array: A structure for efficient  
526 numerical computation. *Computing in Science*  
527 *and Engineering* **13**, 22–30 (2011).
- 528 33. McKinney, W. pandas: a Foundational  
529 Python Library for Data Analysis and Statis-  
530 tics. *Python for High Performance and Scien-*  
531 *tific Computing* 1–9 (2011).
- 532 34. Oliphant, T. E. SciPy: Open source scientific  
533 tools for Python. *Computing in Science and*  
534 *Engineering* **9**, 10–20 (2007).
- 535 35. Hunter, J. D. Matplotlib: A 2D graphics envi-  
536 ronment. *Computing in Science and Engineer-*  
537 *ing* **9**, 99–104 (2007).
- 538 36. Waskom, M. *et al.* seaborn: v0.7.0 (January  
539 2016) (2016).
- 540 37. Angeles-Albores, D., N. Lee, R. Y., Chan, J.  
541 & Sternberg, P. W. Tissue enrichment analysis  
542 for *C. elegans* genomics. *BMC Bioinformatics*  
543 **17**, 366 (2016).
- 544 38. Edgar, R., Domrachev, M. & Lash, A. E. Gene  
545 Expression Omnibus: NCBI gene expression  
546 and hybridization array data repository. *Nu-*  
547 *cleic acids research* **30**, 207–10 (2002).