

# Genetic Analysis of a Metazoan Pathway using Transcriptomic Phenotypes

immediate

This manuscript was compiled on November 14, 2016

**RNA-Seq is a technology that is commonly used to identify genetic modules that are responsive to a perturbation. In theory, global gene expression could also be used as a phenotype in complex metazoans, with all the implications that has for genetic analysis. To that end, we sequenced the transcriptome of four single mutants and two double mutants of the hypoxia pathway in *C. elegans*. We successfully analyzed the single mutants in a blinded fashion to predict the genetic relationships between the genes, and used the double mutants as a test of our predictions and to infer the directionality of the relationship. We show that genes along a pathway tend to decorrelate as a result of alternative regulatory modes and crosstalk with other pathways; and that this decorrelation accurately reflects functional distance between genes. As a by-product of our analysis, we predict 133 genes under the regulation of *hif-1*, and 36 genes under the regulation of *vhl-1*. Transcriptomic perturbations suggest an important role of *hif-1*-dependent response in chromatin remodelling in *C. elegans*. Interactive graphics for this paper can be found at [www.wormlabcaltech.github.io/mprsq](http://www.wormlabcaltech.github.io/mprsq).**

genetics | RNA-Seq | epistasis | hypoxia | transcriptomics | systems biology

**G**enetic analysis of molecular pathways has traditionally been performed through epistatic analysis. Epistasis occurs when two genes interact, either directly (biochemical interaction) or through a molecular pathway or physical interaction (genetic interaction). If two genes interact, and the mutants of these genes have a quantifiable phenotype, the double mutant will have a phenotype that is not the sum of the phenotypes of the single mutants that make up its genotype. Epistatic analysis remains a cornerstone of genetics today [1].

Previous work in *S. cerevisiae* and *D. discoideum* using microarrays has shown that transcriptomes can be used to infer genetic relationships in simple eukaryotes [2, 3]. Developments in the area of transcriptomics have brought forward new protocols, such as RNA-Seq [4], and have also made important progress towards cheaper sequencing [5], better and faster abundance quantification [6–8] and improved differential analysis of gene expression [9, 10]. As a result, RNA-Seq has been successfully used to identify genetic modules involved in a variety of processes, including T-cell regulation [11, 12], the *C. elegans* linker cell migration [13], or planarian stem cell maintenance [14, 15]. However, even in these novel applications, transcriptional profiling largely serves a descriptive role in target gene identification.

To investigate the ability of transcriptomes to serve as quantitative phenotypes, we selected mutants in the *C. elegans* hypoxia pathway for transcriptome sequencing. Metazoans depend on the presence of oxygen in sufficient concentrations to support aerobic metabolism. Genetic pathways evolved to rapidly respond to any acute or chronic changes in oxygen levels at the cellular or organismal level. These oxygen sensitive pathways are involved in a broad range of human pathologies and they have been subject to investigation biochemical

and genetic approaches [16]. These approaches identified the Hypoxia Inducible Factors (HIFs) as an important group of oxygen responsive genes.

Hypoxia Inducible Factors are highly conserved in metazoans [1]. A common mechanism for hypoxia-response induction is heterodimerization between a HIF $\alpha$  and a HIF $\beta$  subunit. The heterodimer then initiates transcription of target genes [1]. The number and complexity of HIFs varies throughout metazoans, with humans having three HIF $\alpha$  subunits and two HIF $\beta$  subunits, whereas in the roundworm *Caenorhabditis elegans* (*C. elegans*) there is a single HIF $\alpha$  gene, *hif-1* and a single HIF $\beta$  gene, *ahr-1*. HIF target genes have been implicated in a wide variety of cellular and extracellular processes such as glycolysis, extracellular matrix modification, autophagy and immunity [1].

Levels of HIF $\alpha$  proteins tend to be tightly regulated. Under conditions of normoxia, HIF-1 $\alpha$  exists in the cytoplasm and partakes in a futile cycle of continuous protein production and rapid degradation with a half-life of minutes [1]. HIF-1 $\alpha$  is hydroxylated by three proline hydroxylases in humans (PHD1, PHD2 and PHD3) but is only hydroxylated by one proline hydroxylase (*egl-9*) in *C. elegans* [1]. HIF-1 hydroxylation increases its binding affinity to Von Hippel Lindau Tumor

## Significance Statement

Measurements of global gene expression are often used as descriptive tools capable of identifying genes that are downstream a perturbation. In theory, there is no reason why measurements of global transcriptomes could not be used as a quantitative phenotype for genetic analysis in multicellular organisms. In fact, qPCR measurements of single or a few reporter genes are already used to perform genetic network analysis. Here, we show that transcriptomes can be used for epistasis analysis in a metazoan, and that transcriptomes afford far more information per experiment than classic genetic analysis. By using transcriptomes as quantitative phenotypes, we can accurately predict interactions between genes, while at the same time identifying genes common to a pathway. When pathways branch, it is also possible to identify gene batteries that are associated with each end of the branch point. Finally, genes that would result in invisible visible phenotypes in an animal are not likely to be invisible at the transcriptome phenotype due to the exquisite granularity present in these structures, which represents an important advance towards studying small effect genes that make up the majority of animals' genetic repertoire.

DA, CPR and PWS designed the experiments. CPR selected the genes and extracted mRNA from all mutants. BW made the libraries. IA performed all sequencing. DA developed the mathematical theory. DA wrote all computer code and performed all analyses. DA made all the reporter strains and performed all microscopy. DA, CPR and PWS wrote the manuscript.

The authors declare no conflict of interest.

<sup>2</sup>To whom correspondence should be addressed. E-mail: [pws@caltech.edu](mailto:pws@caltech.edu)

Suppressor 1 (*vhl-1*), which allows ubiquitination of HIF-1 leading to its subsequent degradation. In *C. elegans*, EGL-9 activity is inhibited by binding of CYSL-1, and CYSL-1 activity is in turn inhibited at the protein level by RHY-1, possibly by post-translational modifications to CYSL-1 [17].

Here, we show that transcriptomes contain strong, robust signals that can be used to infer relationships between genes in complex metazoans by reconstructing a the hypoxia pathway in *C. elegans* using RNA-Seq. To reconstruct this pathway, we developed new analytic mathematical tools to query genetic interactions at the transcriptome level. We used these analytic tools to generate a method for analyzing transcriptome genetics data, which we termed Metazoan RNA-based Gene Analysis (MoRGAn). We show that MoRGAn provides a sound theoretical and experimental framework with which to dissect genetic pathways. Using MoRGAn, we were able to reconstruct interactions between genes in the hypoxia pathway. Furthermore, we show that the phenomenon of phenotypic epistasis, a hallmark of genetic interaction, holds at the molecular systems level. We also demonstrate that transcriptomes contain sufficient information, under certain circumstances, to order genes in a pathway using only single mutants. Finally, we were able to identify genes that appear to be downstream of *egl-9* and *vhl-1*, but are almost certainly not targets of *hif-1*. A complete, interactive version of the analysis is also available at [www.wormlabcaltech.github.io/mprsq](http://www.wormlabcaltech.github.io/mprsq).

## Results

**Development of a New Genetic Logic for Vectorial Phenotypes.** Transcriptomes can be understood as information-rich phenotypes that are fine-grained representations of an organism's internal state (citation?) []. As transcriptomes have become more prominent, the single-cell community in particular has begun to use them to define and understand cellular identities []. However, most approaches rely on data-driven, computational methods to define the important or relevant aspects of a transcriptomes and the community is actively developing new algorithms for increasingly complex data. On the other hand, whole-organism transcriptomes have not been analyzed quite so exhaustively, partially because these transcriptomes contain the identities of multiple cell-types in one convoluted measurement.

Although computational advances for RNA-seq analysis are important, an equally important facet of transcriptome analysis is the development of analytical mathematical tools with which to understand these objects. In an effort to understand transcriptome genetics at a profound level, we developed extensions to standard genetic theory borrowing from concepts in group and operator theory that enabled us to think logically about genetics using a vectorial phenotype.

Briefly, we envisioned that a single gene is always associated with the same specific transcriptome under a set of specified conditions (age of the organism, ambient temperature, food status, etc...). We refer to the set of genes that are differentially expressed in response to a perturbation of a gene as a specific transcriptome. A specific transcriptome can be formally thought of as a group of genes with an associated metric for each gene (such as abundance, fold change or log-fold change), or equivalently as a function that takes as input a gene ID and outputs a scalar. However, although a global transcriptome is fundamentally an observable quantity, the

specific transcriptome associated with a gene is not, since we do not know the identities of the genes that incorporate it. In order to observe a specific transcriptome, it is necessary to perform at least two experiments: first, we must measure a basal transcriptome; second, we must measure a transcriptome in a set of organisms that have a perturbed activity of the gene in question.

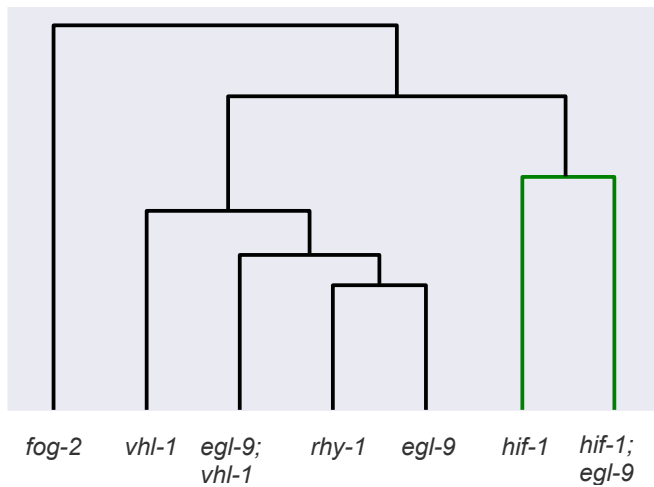
Having established the concept of a specific transcriptome as a group of ordered tuples of the form (gene, measurement), we reasoned that transcriptomes should obey the laws of genetic networks. Namely, we reasoned that if two hypothetical genes, *A* and *B*, act in a linear unbranched pathway, then *A* and *B* should have equivalent specific transcriptomes. We also reasoned that if *A* acts upon two different pathways, one which involves *B* and one which doesn't, the specific transcriptome of *A* should itself be separable into two independent parts: the specific transcriptome associated with the pathway  $A \rightarrow B$ , and the specific transcriptome associated with *A* but not *B*.

In order to solve genetic pathways, it is not sufficient to measure two transcriptomes and compare the overlap between the genes they contain. To solve complex genetic pathways, we must understand the valence of a genetic interaction between genes. In order to do this, we created linear operators to model perturbations that increase or decrease the activity of a gene relative to the basal measured state. These linear operators allow us to propagate a perturbation in one gene through a network and predict the expected changes in specific transcriptomes associated with the other genes in a network subject to the genetic model that has been specified by the investigator.

Using this formalism, we are, in theory, able to infer genetic relationships between genes. However, this formalism also provides a rational framework with which to identify specific transcriptomes associated with a pathway, and it also provides a theoretical manner with which to order genes acting along a pathway if the pathway branches at multiple points. Finally, our methods can allow us to identify regulatory relationships that are acting on different levels (either at the transcriptional level or at the protein level). For example, assuming that the mode of interaction between *rhy-1* and *egl-9* is the same as the mode of interaction between *hif-1* and *rhy-1* (transcriptional control) leads to aberrant behavior that is not observed empirically. Only by assuming that the regulatory mode between *rhy-1* and *egl-9* is different than between *hif-1* and *rhy-1* does the model generate predictions that parallel reality. Our new genetic notation, and the concepts contained therein, represent an important advance towards harnessing the power of transcriptomes as quantitative phenotypes. For a full description of this notation, see the S.I..

## Clustering visualizes epistatic relationships between genes.

As a first step in our analysis, we analyzed our data using a generalized linear model with a genotype term (see 1) on logarithm-transformed counts. Genes that are significantly altered between wild-type and a given mutant have a genotype coefficient that is statistically significantly different from 0. We refer to these coefficients through the greek letter  $\beta$ . These coefficients are not identical to the average log-fold change per gene, although they are loosely related to this quantity. In general, larger  $\beta$  magnitudes correspond to larger perturbations. These coefficients can be used to study the RNA-Seq data in question.



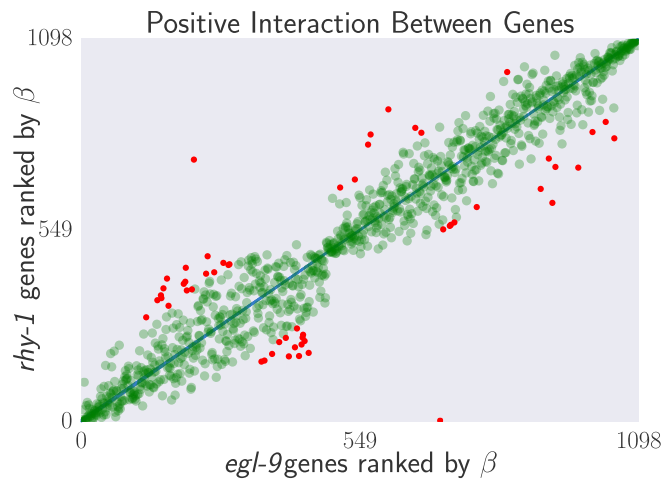
**Fig. 1.** Blind unsupervised clustering of various *C. elegans* mutants. Genes cluster in a manner that is biologically intuitive. Genes that inhibit *hif-1* (i.e. *egl-9*, *vhl-1*, and *rhy-1*) cluster far from *hif-1*. *hif-1* clusters with the suppressed *egl-9*; *hif-1* double mutant. A control gene, *fog-2*, clusters farthest away.

Clustering is a well-known technique in bioinformatics that is used to identify relationships between high dimensional data points [18]. We wanted to make sure that clustering by differential expression yielded genetically relevant information. *hif-1* exhibits no obvious phenotypes under normoxic conditions, in contrast to *egl-9*, which exhibits an egg-laying (*egl*) phenotype in the same environment. *egl-9*; *hif-1* mutants suppress the *egl* phenotype. If transcriptomic phenotypes behave similarly to their macroscopic counterparts, *hif-1* should cluster with the *egl-9*; *hif-1* double mutant, whereas *egl-9* should cluster away from the *hif-1* mutant. Indeed, when blind, unsupervised clustering was performed on the data, three clusters emerged. *hif-1* and *egl-9*; *hif-1* clustered together, indicating suppression of the *egl-9* phenotype; whereas *egl-9*, *egl-9*; *vhl-1*, *vhl-1* and *rhy-1* all clustered separately. Finally, our negative control *fog-2* was in its own cluster (see Fig. 1). We conclude that expression data contains enough signal to cluster genes in a meaningful manner in complex metazoans.

### Transcriptomic correlations can predict genetic regulation.

Theoretically, two genes that share linear positive regulation should be positively correlated in their overlapping transcriptomes, whereas two genes that share linear negative regulation should be negatively correlated in their transcriptomes. Conversely, it follows that if two mutants have overlapping transcriptomes that are strongly positively correlated, it is likely that these two genes share a positive regulatory association. In other words, transcriptomic correlation is a good predictor of genetic regulation. For a formal introduction to the genetic logic, see S.I..

Although transcriptomic correlations could theoretically be used for the purposes of identifying genetic regulation, noise can cause serious interference with any inferences. Additionally, genes sometimes experience multiple modes of regulation, including positive and negative regulation, from the same gene or pathway. Because we are measuring the system at steady state, both modes of regulation will be measured simultaneously. If a positive and a negative signal of equal strength are both present in a transcriptome, running a naive regression



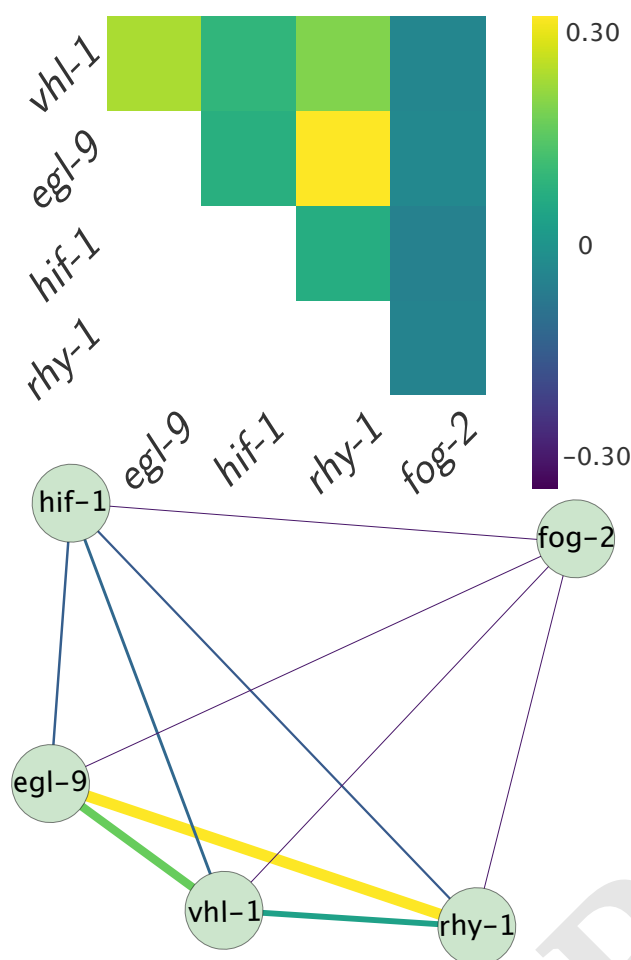
**Fig. 2.** Strong transcriptional correlations can be identified between genes that share a positive regulatory connection. We took the *egl-9* and the *rhy-1* transcriptomes, identified differentially expressed genes common to both transcriptomes and ranked each gene according to its differential expression coefficient  $\beta$ . We then plotted the rank of each gene in *rhy-1* versus the rank of the same gene in the *egl-9* transcriptome. The result is an almost perfect correlation. For unknown genes, such a correlation would be predictive of an interaction.

may result in a value close to zero. Therefore, we took steps to mitigate noise emanating from frequent outliers and to identify multiple regulatory signals.

To mitigate noise, we rank-transformed the  $\beta$  coefficients for each mutant. This has the effect of mitigating outliers by resetting the difference between adjacent coefficients to unity. Next, we performed robust Bayesian regressions using a Student-T distribution as a prior. A Student-T distribution decays less quickly than a normal distribution, which causes the model to consider outliers to be less informative than traditional regressions.

We saw that certain gene pairs correlated very well when genes were ranked by their expression changes (see Fig. 2). We generated all pairwise correlations between transcriptomes and we weighted the correlations by the number of genes that participated in the correlation (that were not outliers) divided by the total number of genes detected in all samples. We were able to identify a strong positive interaction between *egl-9* and *rhy-1*. The transcriptomes for these genes consisted of 1,487 and 1,816 significantly altered genes respectively and the overlap between both transcriptomes was extensive. On the other hand, none of the primary correlations between *hif-1* and its controlling genes were negative. We were unable to definitively determine the reason for behind this. The overlap between *hif-1* and all other genes was relatively small, and each overlap involved different sets of genes, which suggests that we did not sequence deeply enough to identify the nature of these positive interactions. See SI A for an exhaustive analysis of the expected and observed correlation between each gene pair in this circuit respectively. The regression slopes recapitulated a network with three ‘modules’: A control module, a responder module and an uncorrelated module (see Fig. 3). We were able to identify a strong positive interaction between *egl-9* and *rhy-1*. The magnitude of this weighted correlation is derived from the fact that the transcriptomes for these genes consisted of 1,487 and 1,816 significantly altered genes respectively and the overlap between both genes was

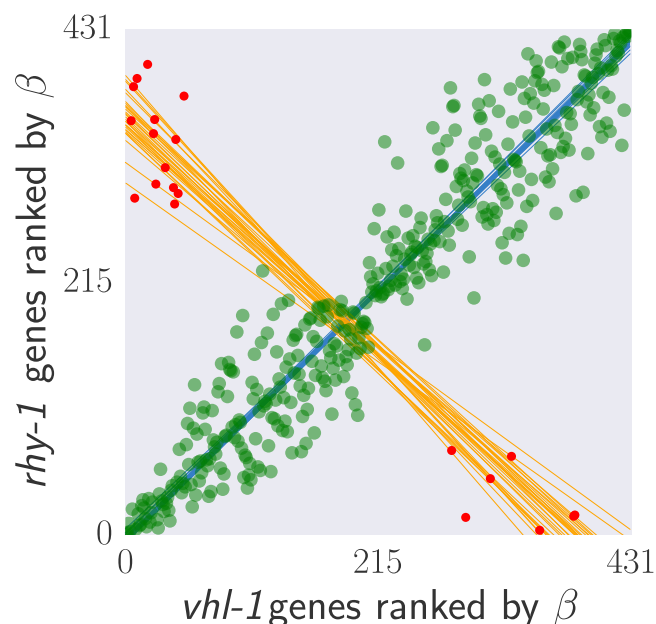
quite extensive, which makes the weighting factor considerably larger than other pairs.



**Fig. 3. Top:** Heatmap showing pairwise regression values between all single mutants. **Bottom:** Correlation network drawn from the diagram. Edge width is proportional to the weighted correlation between two nodes raised to the power of 1.15. Making the edge width proportional to a power of the regression makes differences in the weighted correlations easier to see.

Previous work in the hypoxia pathway has found extensive feedback loops in this pathway. Using the genetic formalism we developed, we realized that due to the fine-grained nature of interactomes we can use them to measure two regulatory interactions of opposite sign simultaneously in a single gene pair. This should lead to a characteristic *X* pattern in the ranked data. Such cross patterns can be indicative of feedback loops or of incoherent regulation at two distinct levels of gene expression (see Fig 4). For any gene that has a *X*-patterned diagram, we refer to the correlation that contains the largest number of participating isoforms as the primary correlation, and the correlation that contains the lesser number of points is referred to as a secondary correlation. In our dataset, all primary correlations were positive.

We investigated whether any pairwise comparisons between our single mutants generated this cross pattern. Indeed, we found that comparing *hif-1* with *rhy-1*, and *hif-1* with *egl-9* yielded negative correlations, as did *rhy-1* and *vhl-1*. While the number of genes that lead these negative correlations is



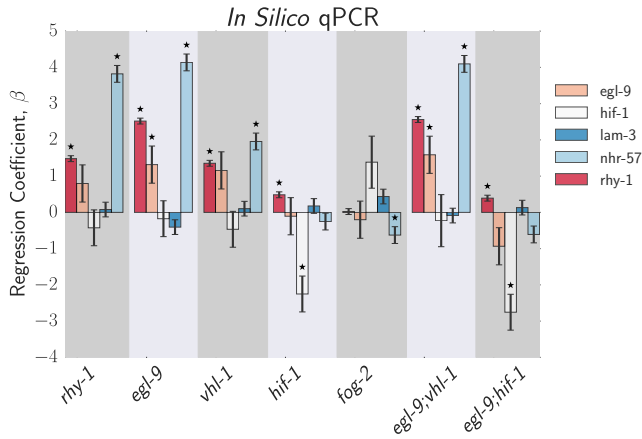
**Fig. 4. Top:** A feedback loop can generate transcriptomes that are both correlated and anti-correlated. **Bottom:** *hif-1* transcriptome correlated to the *rhy-1* transcriptome. Green large points are inliers to the first regression. Red small points are outliers to the first regression. Only the red small points were used for the secondary regression. Blue lines are representative samples of the primary bootstrapped regression lines. Orange lines are representative samples of the secondary bootstrapped regression lines.

small (10–30 genes in any comparison) and is not significantly different from random expectation as assessed by a hypergeometric test, these outliers are expected for this circuit (see SI). On the other hand, unweighted secondary correlations had coefficients near unity for most comparisons, and they are predicted by theory. Statistical information should be integrated holistically with genetic models to assess whether outliers are meaningful or not. In our case, we believe that these outliers reflect meaningful interactions between genes in the hypoxia circuit, not random noise.

**in silico qPCR reveals extensive feedback in the hypoxia pathway.** Our dataset enables us to perform a sort of *in silico* qPCR by selectively looking at expression of a few genes at a time. To verify the quality of our data, we queried the changes in expression of *nhr-57*. This reporter has been shown to have *hif-1* dependent expression [19–22]. In our dataset, this gene be upregulated in *egl-9*, *rhy-1* and *vhl-1*, but remains unchanged in *hif-1*. The *egl-9*; *vhl-1* had an expression level similar to *egl-9*; whereas the *egl-9*; *hif-1* mutant showed suppression of the reporter expression. All of these interactions reflect the literature.

We also performed *in silico* qPCR of every gene under scrutiny to get a clearer idea of the relationships between them (see Fig. 5). We observed changes in *rhy-1* expression consistent with previous literature [ ] when *hif-1* is activated. We also observed changes in *egl-9* expression when *egl-9* was mutated, and previous literature has identified *egl-9* as a hypoxia responsive gene [ ]. Although *egl-9* was not increased in *rhy-1* and *vhl-1* mutants, the mRNA levels of *egl-9* trended towards increased expression. As with *nhr-57*, the *egl-9* and *rhy-1* expression phenotypes were abrogated in the *egl-9*; *hif-1*





**Fig. 5. Top:** *In silico* qPCR results. *nhr-57* is an expression reporter that has been used previously to identify *hif-1* regulators [19, 23]. The *nhr-57* mRNA levels replicate what is observed in the literature. *lam-3* is shown here as a negative control that should not be altered by mutations in this pathway. The increases in the levels of *egl-9* and *rhy-1* when repressors of *hif-1* are knocked out are in agreement with previous literature [24]. We measured modest increases in the levels of *rhy-1* mRNA when *hif-1* is knocked out. The mechanism behind this is unclear. Negative and positive feedback loops from *hif-1* into its inhibiting genes could be a homeostatic mechanism.

mutant; whereas the *egl-9;vhl-1* mutant showed expression phenotypes identical to the *egl-9* mutant, in support of *egl-9* and *vhl-1* acting in an AND-gated manner. Our dataset also shows that knockout of *hif-1* resulted in a modest increase in the levels of *rhy-1*. This suggests that *hif-1* is also a negative regulator of *rhy-1*.

In summary, the *in silico* qPCR results recapitulate previous results that show *egl-9* and *vhl-1* act in concert to inhibit *hif-1*. Moreover, these results taken together with the transcriptome-wide cross-patterns that emerge from pairwise comparisons between genes in the hypoxia pathway suggest that there are both positive and negative feedback loops feeding into *rhy-1* and possibly *egl-9*. These feedback loops could explain why *hif-1* is positively transcriptomically correlated with *egl-9*.

**Epistasis effects can be detected and quantified..** To quantify any epistasis between *egl-9* and *vhl-1* in our dataset, we identified the genes that were shared between each single mutant and the double mutant *egl-9;vhl-1*. If two genes act, for example, in a linear manner, then the double mutant should exhibit an identical phenotype to each single mutant. To test such a relationship, we can plot the difference in a gene *i* between the log change in expression between the two mutants,  $\Delta_i = \beta_{\text{Double Mutant},i} - \beta_{\text{Single Mutant},i}$ , against the log change in the single mutant,  $\beta_{\text{Single Mutant},i}$ . We can then fit a weighted linear regression to measure the slope of best fit. Genes that act in a linear pathway should yield lines with a slope of 0. Genes that have some additive flavor should have slopes greater than 0. Suppression, a hallmark of inhibition, should yield a slope less than 0.

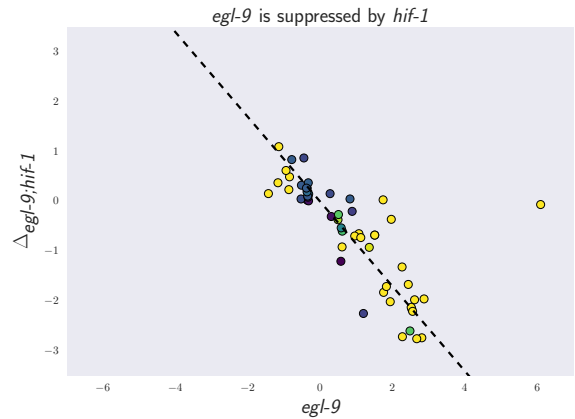
We observe that the *egl-9;vhl-1* mutant has an identical phenotype to the *egl-9* single mutant (slope = 0; see Table 1). On the other hand, *vhl-1* has a positive slope, indicating that *egl-9* is additive to *vhl-1*. Such partial additivity can be explained if *egl-9* is inhibiting *hif-1* in a *vhl-1*-dependent as well as a *vhl-1*-independent manner [23].

On the other hand, comparison of the *egl-9;hif-1* double

**Table 1. Response Modeling of Double Mutants to Single Mutants**

Double Mutant	Single Mutant	$\Delta$	SE	p-value
1. <i>egl-9;vhl-1</i>	<i>egl-9</i>	0.00	0.01	0.81
2. <i>egl-9;vhl-1</i>	<i>vhl-1</i>	0.28	0.033	$10^{-15}$
3. <i>egl-9;hif-1</i>	<i>egl-9</i>	-0.85	0.074	$10^{-13}$
4. <i>egl-9;hif-1</i>	<i>hif-1</i>	-0.18	0.10	0.10

Table showing changes between single and double mutants.  $\Delta$  is the result of a weighted-linear regression (WLS) between  $\beta_{\text{Single Mutant}}$  and  $\Delta = \beta_{\text{Double Mutant}} - \beta_{\text{Single Mutant}}$ .  $\Delta > 0$  represents a more severe phenotype than the single mutant.  $\Delta < 0$  represents a suppressed phenotype relative to the single mutant.  $\Delta = 0$  is expected for linear pathways or genes that are acting in linear or AND-gated fashion.  $\Delta > 0$  is expected for genes that are acting additively on a pathway. WLS were performed only on genes that were significantly altered in both single mutants and the double mutant.  $1 + \Delta$  is a very close approximation to the line of best fit between single mutant and double mutant.



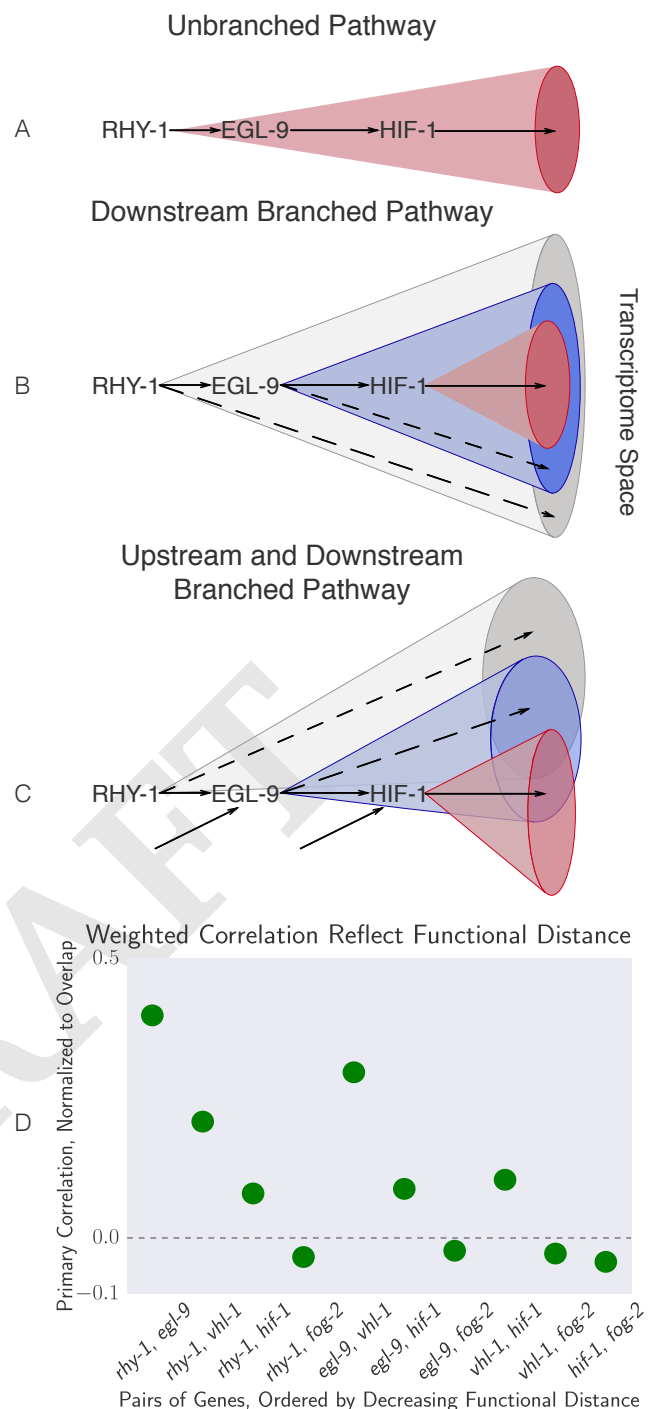
**Fig. 6.** The mutant *egl-9* transcriptomic phenotype is suppressed by mutations in *hif-1*. The graph shows the  $\beta$  coefficients for *egl-9* in the x-axis, and the change in  $\beta$  coefficient between the *egl-9;hif-1* and *egl-9* mutant. The dotted line is the regression line between the complete *egl-9* and *egl-9;hif-1* shared transcriptome. For clarity, only genes that were differentially expressed in the *egl-9*, *rhy-1*, *vhl-1*, *hif-1* and *egl-9;vhl-1* datasets are shown. These points constitute a very high-quality subset of the measured hypoxia response, as each isoform was identified as differentially expressed in 5 independent genotypes. The single outlier near (6, 0) is *nog-1*, and it is probably downstream of *egl-9*, and is not likely a *hif-1* target.

mutant showed suppression of the *egl-9* transcriptomic phenotype. This suppression is expressed in various ways. First, the double mutant shows less statistically significantly differentially expressed genes than either single mutant. Secondly, the genes that are common to the *egl-9* and *egl-9;hif-1* transcriptomes show decreased expression in the *egl-9;hif-1* mutant than they do in *egl-9* on average (see Fig. 6). Likewise, the genes that are common to *hif-1* and *egl-9;hif-1* show no change in expression on average between these two mutants.

**Transcriptomic decorrelation can be used to infer functional distance.** We were interested in figuring out whether RNA-Seq could be used to identify functional interactions within a genetic pathway. Although there is no *a priori* reason why global gene expression should reflect functional interactions, the strength of the unweighted correlations between genes in the hypoxia pathway made us wonder how much information can be extracted from this dataset. Single genes are often regulated by multiple independent sources. The connection between two nodes can in theory be characterized by the strength of the edges connecting them (the thickness of the edge); the fraction of sources that regulate both nodes (the fraction of common inputs); and the fraction of genes that are regulated by both nodes (the fraction of common outputs). 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.

We investigated the possibility that transcriptomic signals do in fact contain relevant information about the degrees of separation by weighting the robust bayesian regression of each pairwise analysis by  $N_{\text{Overlap}}/N_{\text{detected}}$ . We plotted the weighted correlation of each gene pair, ordered by increasing functional distance (see Fig. 7). In every case, we see that the weighted correlation decreases monotonically due mainly, but not exclusively, to decreasing  $N_{\text{Overlap}}$ . 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 pathway, as we have done with the hypoxia pathway<sup>1</sup>.

**Identification of novel targets and biological processes in the hypoxia response.** So far, our analysis has focused mainly on extracting genetic relationships between the set of mutants we sequenced. Our dataset also provides us with a unique view of the *hif-1*-dependent response in *C. elegans*. In total, we identified 3,211 differentially expressed genes that are altered in any of the hypoxia pathway mutants. Of these 3,211 genes, 53 genes were differentially expressed in all the hypoxia mutants. Because of the extensive feedback between *hif-1* and *egl-9*,



**Fig. 7.** Theoretically, 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 transcriptomes with enough information to infer the order between genes. **B** On the other hand, 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 in both upstream and downstream directions, observed 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, reflecting the weak control distant nodes exert on each other. **D** The hypoxia pathway can be ordered according to functional distance. The rapid decay in correlation is probably due to a mixture of upstream and downstream branching that happens along this pathway.

<sup>1</sup> An important question is whether a looped circuit like the hypoxia pathway can be ordered in the way we have ordered it in Fig. 7 since a loop does not technically have a beginning. One explanation is that we studied the hypoxia pathway under normoxic conditions, and therefore the control of *hif-1* over *rhy-1* and *egl-9* is weak, effectively turning the looped pathway into a linear one. Probably, under hypoxic conditions the pathway would effectively be reversed.

we expected to identify a small subset of genes that were up-regulated or down-regulated consistently in every hypoxia mutant except the *egl-9;hif-1* double mutant. We identified 10 genes that were up-regulated in this manner, and 13 genes that were down-regulated (see SI for gene identities). These genes likely constitute a core response around the circuit in question, and their behaviour should reflect the genetic relationships the best. Indeed, graphing these genes shows beautiful agreement with predictions (see [www.wormlabcaltech.github.io/mpsq](http://www.wormlabcaltech.github.io/mpsq) for interactive graphics).

In order to identify affected biological processes, we performed an in-house gene ontology enrichment analysis using annotations provided by WormBase, following the procedure shown in TEA [25]. Top enriched terms included ‘hydrolase activity’ (869 observed hits; 7.8 fold change; p-value <  $10^{-10}$ ); ‘organic anion transport’ (803 hits; 7.5; p-value <  $10^{-10}$ ); ‘spliceosomal complex’ (647 hits; 8.2 p-value <  $10^{-10}$ ); ‘SAM-dependent methyltransferase activity’ (1215 hits; 6.6; p-value <  $10^{-10}$ ); and ‘cell division’ (1251 hits; 7.9; p-value <  $10^{-10}$ ). In mammals, the mammalian target of rapamycin pathway (mTOR), which is intimately associated with the hypoxia pathway, has been previously linked to osmotic stress responses [26]. Our findings also suggest that the *hif-1*-dependent response causes important changes in chromatin structure via activation or recruitment of chromatin remodelling factors, as well as changes in isoform processing.

We identified downstream targets of the genes we studied that were not associated with other genes in the *hif-1* circuit. *vhl-1* targets were particularly easy to isolate because *vhl-1* does not seem to participate in the *rhy-1*, *egl-9*, *hif-1* feedback circuit, and as a result it is easy to isolate targets for that are *hif-1*-independent. We found 36 genes downstream of *vhl-1*. These 36 genes include *pole-1*, an ortholog of human polymerase  $\epsilon$  catalytic subunit; *F33H2.6*, an ortholog of the human regulator of microtubule dynamics 1 (RMDN1) [1]; and many solute carriers. Reflecting this, enriched GO terms were ‘ion binding’, ‘growth’, ‘cell division’, ‘cell projection assembly’ as well as ‘ion binding’ and ‘divalent metal ion transport’. *vhl-1* has been previously implicated as a controller of mitotic fidelity in renal cell carcinoma [27]. Our findings support a role of *vhl-1* in chromosomal integrity and mitotic fidelity. Furthermore, recent reports suggest that solute carriers may be associated with poor prognosis in clear-cell renal carcinoma [28], which highlights the biological relevance of our predictions.

We identified 133 genes that are activated by HIF-1. We verified that the genes we identified are actually *hif-1* targets by searching for a set of 20 gold-standard genes from the literature [19, 20] in our gene set, and found that *hif-1* targets were significantly enriched in these genes ( $p < 10^{-7}$ ). GO term enrichment indicated that this list was associated with ‘cell division’, ‘SAM methyltransferase activity’ and ‘cellular modified amino acid metabolic processes’. A full list of *hif-1* targets can be found in S.I..

**Discovery of New Pathways: *nog-1*.** Enzymes rarely have a single substrate. Although the role of *egl-9* in *hif-1* hydroxylation has been clearly elucidated, the roles of *hif-1*-independent *egl-9* effects remain largely unknown and had not been explored until relatively recently [1]. Using our experimental design, we were able to measure transcriptomic epistasis in our double mutant. A corollary of measuring epistasis is that we can now identify genes that violate the epistasis relationships

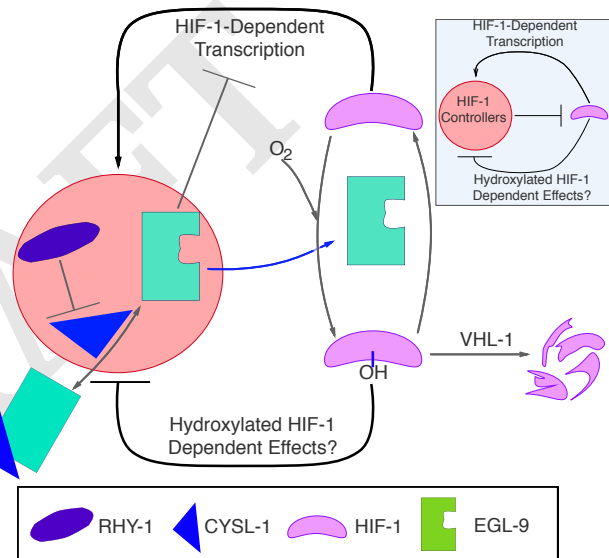
**Table 2. Using *nog-1* as a genetic reporter reconstructs a new pathway and identifies *nog-1* downstream to *egl-9*.**

Single Mutant	Double Mutant	$\beta_{\text{single}}$	$\beta_{\text{double}}$
1. <i>egl-9</i>	<i>egl-9;vhl-1</i>	6.10	6.23
2. <i>vhl-1</i>	<i>egl-9;vhl-1</i>	6.45	6.23
3. <i>egl-9</i>	<i>egl-9;hif-1</i>	6.08	6.02
4. <i>hif-1</i>	<i>egl-9;hif-1</i>	3.56	6.02

Our genetics design allows us to identify genes that don’t agree with the model we have built. *nog-1* appears to be under strong control by *egl-9* and *vhl-1*, probably through hydroxylation and ubiquitination of an unknown factor that drives transcription to *hif-1*. *nog-1* is not a target of *hif-1*, since *hif-1* loss of function has a weaker effect than *egl-9* and the *egl-9;hif-1* double mutant shows the same expression level as *egl-9*.

between the genes we studied. As a particular example, we focused on a gene called *nog-1*.

## Discussion



**Fig. 8.** A schematic of the hypoxia pathway in *C. elegans*. RHY-1 likely inhibits a protein-protein interaction between EGL-9 and CYSL-1 [17]. This interaction inhibits the activity of EGL-9. When active, EGL-9 can inhibit HIF-1 by catalyzing the hydroxylation of HIF-1, which leads to its recognition and ubiquitination by VHL-1 and ultimately leads to rapid degradation of HIF-1 protein. Independently of its enzymatic function, EGL-9 can also inhibit HIF-1 transcriptional activity. Our results identified *rhy-1*, *egl-9* and *cysl-1* as downstream targets of the *hif-1*-dependent transcriptional response. However, we also identified increased in the mRNA levels of *rhy-1* and *cysl-1* when *hif-1* was knocked out. Moreover, the *hif-1* knock-out transcriptome correlated positively with *rhy-1*, *egl-9* and *vhl-1*. One plausible explanation for these observations is that hydroxylated HIF-1 has transcriptional consequences. Inset shows a simplified diagram of the hypoxia pathway.

**The Hypoxia Circuit in *C. elegans*.** Previous work has established a circuit in which *rhy-1* activates *egl-9* in a linear pathway, and *egl-9* inhibits *hif-1* in an oxygen-dependent manner. Hydroxylated HIF-1 can then be degraded in a *vhl-1*-dependent manner. There is also evidence that *egl-9* and *rhy-1* are in turn activated by *hif-1* [24, 29]. Finally, there is evidence that although the interaction between *egl-9* and *vhl-1* is important for *hif-1* repression, *egl-9* can also act in a

non-*vhl-1* dependent manner (see Fig. 8 top).

We were able to impute the positive regulatory relationship between *egl-9* and *rhy-1*. We would not have been able to infer the order of the regulation without additional information. Using clustering as a proxy for phenotype, we were able to infer the relationship between *egl-9* and *hif-1*. We were also able to infer a positive (linear or AND) relationship between *egl-9* and *vhl-1* using clustering. Alternatively, we gained the same information by performing *in silico* qPCR on the genes under study. *In silico* qPCR also revealed that *hif-1* has two states with different activities: Non-hydroxylated HIF-1 increases levels of *rhy-1*, and hydroxylated HIF-1 inhibits *rhy-1* and possibly *egl-9* as well, although the double mutant did not recapitulate that interaction. We also revealed that *hif-1* is an autoregulator.

These discoveries are consistent with a homeostatic circuit. By autoregulating itself, *hif-1* can maintain appropriate protein levels both in normoxic and hypoxic conditions. Inhibition of *rhy-1*, and possibly of *egl-9*, ensures that an appropriate equilibrium is maintained between hydroxylated and non-hydroxylated protein, which may have functional consequences for the cell if both forms are active.

In addition to these biological findings, our dataset allows us to generate predictions of genes that may be under direct *hif-1* regulation. Assuming that non-hydroxylated HIF-1 has different activities from hydroxylated HIF-1, we identified 5 genes that are candidates for activation by hydroxylated HIF-1. These genes have been implicated in the *C. elegans* immune response, or have behavioural phenotypes, underscoring the importance of *hif-1* in neurobiology and immunology [30–33].

**Towards A Genetic Theory of Transcriptomics.** We have shown that transcriptomes contain sufficient information to be used as semi-quantitative phenotypes in metazoans. These phenotypes can be interpreted globally via correlation tests, clustering or other probabilistic methods; alternatively, they can be used to query single reporter genes in a manner similar to qPCR today. Transcriptomic phenotypes have distinct advantages over physical traits. Firstly, due to their increased complexity, the genotype-phenotype mapping degeneracy ought to be greatly reduced, which facilitates predictions of genetic interaction. Secondly, genes that result in subtle or no visible traits when mutated may have strong (detectable), reproducible phenotypes at the transcriptomic level, which would make the study of small-effect genes significantly easier.

RNA-Seq and microarray datasets have been used previously by bioinformaticians to generate high-throughput predictions of genetic interactions and consortiums such as the The Cancer Genome Atlas have sequenced RNA from many different cancers in the hope of identifying clinically or biologically relevant interactions []. By correlating many different datasets in many different conditions, it is possible in theory to predict genetic interaction. Our approach differs from these high-throughput methods in that we are not attempting to generate large scale networks. Rather, the strength in our analysis derives from our experimental design, which allows us to ask and answer a large number of questions about the functional interactions between genes. As a by-product, we are also able to identify genes related to the core circuit studied in question, but our main goal is not to generate databases or predict large numbers of interactions between a large number of genes. We have shown that transcriptomic phenotypes

can capture distinct interaction modes in a single experiment, making it possible to infer complex regulatory relationships between genes. By measuring these transcriptomes under a rigorous experimental design, it is possible to identify many relationships simultaneously. With the advent of fast pseudo-alignment tools and ever cheaper sequencing techniques, biologists should consider using global transcriptomes as a tool beyond hypothesis generation or target acquisition. We have developed a genetic framework to deal with sequence-based phenotypes. Here, we have shown that global transcriptomes can be readily dealt with under this framework. In the future, we hope that other sequence-based phenotypes will shed new insights into genetic relationships between genes.

## Materials and Methods

### RNA-Seq. Tagmentation etc

We used Kallisto to perform pseudo-read alignment and performed differential analysis using Sleuth. We fit a generalized linear model for a transcript  $t$  in sample  $i$ :

$$y_{t,i} = \beta_{t,0} + \beta_{t,genotype} \cdot X_{t,i} + \beta_{t,batch} \cdot Y_{t,i} + \epsilon_{t,i} \quad [1]$$

where  $y_{t,i}$  are the logarithm transformed counts;  $\beta_{t,genotype}$  and  $\beta_{t,batch}$  are parameters of the model, and which can be interpreted as biased estimators of the log-fold change;  $X_{t,i}$ ,  $Y_{t,i}$  are indicator variables describing the conditions of the sample; and  $\epsilon_{t,i}$  is the noise associated with a particular measurement.

**Genetic Analysis.** Genetic analysis of the processed data was performed in Python 3.5. Our scripts made extensive use of the Pandas, Matplotlib, Scipy, Seaborn, Sklearn, Networkx, Bokeh, PyMC3, and TEA libraries [25, 34–41]. Our analysis is available in a Jupyter Notebook [42]. All code and required data (except the raw reads) 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/>. Raw reads were deposited at XXXXXXXXXXXX

**ACKNOWLEDGMENTS.** This article was written with support of the Howard Hughes Medical Institute.

1. Phillips PC (2008) Epistasis — the essential role of gene interactions in the structure and evolution of genetic systems. *Nat Rev Genet* 9(11):855–867.
2. Hughes TR et al. (2000) Functional Discovery via a Compendium of Expression Profiles. *Cell* 102(1):109–126.
3. Van Driessche N et al. (2005) Epistasis analysis with global transcriptional phenotypes. *TL - 37. Nature genetics* 37 VN - r(5):471–477.
4. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B (2008) Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nature Methods* 5(7):621–628.
5. Metzker ML (2010) Sequencing technologies - the next generation. *Nature reviews. Genetics* 11(1):31–46.
6. Patro R, Mount SM, Kingsford C (2014) Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms. *Nature biotechnology* 32(5):462–464.
7. Bray NL, Pimentel H, Melsted P, Pachter L (2015) Near-optimal RNA-Seq quantification. *arXiv*.
8. Patro R, Duggal G, Kingsford C (2015) Salmon: Accurate, Versatile and Ultrafast Quantification from RNA-seq Data using Lightweight-Alignment. *bioRxiv* p. 021592.
9. Pimentel HJ, Bray N, Puente S, Melsted P, Pachter L (2016) Differential analysis of RNA-Seq incorporating quantification uncertainty. *bioRxiv* p. 058164.
10. Trapnell C et al. (2013) Differential analysis of gene regulation at transcript resolution with RNA-seq. *Nature biotechnology* 31(1):46–53.
11. Singer M et al. (2016) A Distinct Gene Module for Dysfunction Uncoupled from Activation in Tumor-Infiltrating T Cells. *Cell* 166(6):1500–1511.e9.
12. Shalek AK et al. (2013) Single-cell transcriptomics reveals bimodality in expression and splicing in immune cells. *Nature* 498(7453):236–40.
13. Schwarz EM, Kato M, Sternberg PW (2012) Functional transcriptomics of a migrating cell in *Caenorhabditis elegans*. *Proceedings of the National Academy of Sciences of the United States of America* 109(40):16246–51.
14. Van Wolfswinkel JC, Wagner DE, Reddini PW (2014) Single-cell analysis reveals functionally distinct classes within the planarian stem cell compartment. *Cell Stem Cell* 15(3):326–339.
15. Scimone ML, Kravarik KM, Lapan SW, Reddini PW (2014) Neoblast specialization in regeneration of the planarian *schmidtea mediterranea*. *Stem Cell Reports* 3(2):339–352.



16. Semenza GL (2012) Hypoxia-inducible factors in physiology and medicine. *Cell* 148(3):399–408.
17. Ma DK, Vozdek R, Bhatla N, Horvitz HR (2012) CYSL-1 Interacts with the O<sub>2</sub>-Sensing Hydroxylase EGL-9 to Promote H<sub>2</sub>S-Modulated Hypoxia-Induced Behavioral Plasticity in *C. elegans*. *Neuron* 73(5):925–940.
18. Yeung KY, Medvedovic M, Bumgarner RE (2003) Clustering gene-expression data with repeated measurements. *Genome biology* 4(5):R34.
19. Shen C, Shao Z, Powell-Coffman JA (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. *Genetics* 174(3):1205–1214.
20. Shen C, Nettleton D, Jiang M, Kim SK, Powell-Coffman JA (2005) Roles of the HIF-1 hypoxia-inducible factor during hypoxia response in *Caenorhabditis elegans*. *Journal of Biological Chemistry* 280(21):20580–20588.
21. Ackerman D, Gens D (2012) Insulin/IGF-1 and hypoxia signaling act in concert to regulate iron homeostasis in *Caenorhabditis elegans*. *PLoS Genetics* 8(3).
22. Park EC et al. (2012) Hypoxia regulates glutamate receptor trafficking through an HIF-independent mechanism. *The EMBO Journal* 31(6):1618–1619.
23. Shao Z, Zhang Y, Powell-Coffman JA (2009) Two distinct roles for EGL-9 in the regulation of HIF-1-mediated gene expression in *Caenorhabditis elegans*. *Genetics* 183(3):821–829.
24. Powell-Coffman JA (2010) Hypoxia signaling and resistance in *C. elegans*. *Trends in Endocrinology and Metabolism* 21(7):435–440.
25. Angeles-Albores D, N. Lee RY, Chan J, Sternberg PW (2016) Tissue enrichment analysis for *C. elegans* genomics. *BMC Bioinformatics* 17(1):366.
26. Zhou B et al. (2007) Hypertonic induction of aquaporin-5: novel role of hypoxia-inducible factor-1alpha. *Am J Physiol Cell Physiol* 292(4):C1280–90.
27. Hell MP, Duda M, Weber TC, Moch H, Krek W (2014) Tumor suppressor vhl functions in the control of mitotic fidelity. *Cancer Research* 74(9):2422–2431.
28. Liu Y et al. (2015) High expression of Solute Carrier Family 1, member 5 (SLC1A5) is associated with poor prognosis in clear-cell renal cell carcinoma. *Scientific reports* 5(October):16954.
29. Bishop T et al. (2004) Genetic analysis of pathways regulated by the von Hippel-Lindau tumor suppressor in *Caenorhabditis elegans*. *PLoS Biology* 2(10).
30. Gray JM et al. (2004) Oxygen sensation and social feeding mediated by a *C. elegans* guanylate cyclase homologue. *Nature* 430(6997):317–322.
31. Cheung BHH, Cohen M, Rogers C, Albayram O, De Bono M (2005) Experience-dependent modulation of *C. elegans* behavior by ambient oxygen. *Current Biology* 15(10):905–917.
32. Chang AJ, Bargmann CI (2008) Hypoxia and the HIF-1 transcriptional pathway reorganize a neuronal circuit for oxygen-dependent behavior in *Caenorhabditis elegans*. *Proceedings of the National Academy of Sciences of the United States of America* 105(20):7321–7326.
33. Ma DK et al. (2013) Cytochrome P450 drives a HIF-regulated behavioral response to reoxygenation by *C. elegans*. *Science (New York, N.Y.)* 341(6145):554–8.
34. Team BD (2014) Bokeh: Python library for interactive visualization.
35. McKinney W (2011) pandas: a Foundational Python Library for Data Analysis and Statistics. *Python for High Performance and Scientific Computing* pp. 1–9.
36. Oliphant TE (2007) SciPy: Open source scientific tools for Python. *Computing in Science and Engineering* 9:10–20.
37. Pedregosa F et al. (2012) Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12:2825–2830.
38. Salvatier J, Wiecki T, Fonnesbeck C (2015) Probabilistic Programming in Python using PyMC. *Arxiv* pp. 1–24.
39. Van Der Walt S, Colbert SC, Varoquaux G (2011) The NumPy array: A structure for efficient numerical computation. *Computing in Science and Engineering* 13(2):22–30.
40. Developers M (year?) matplotlib: v1.5.3.
41. Waskom M et al. (year?) seaborn: v0.7.0 (January 2016).
42. Pérez F, Granger B (2007) IPython: A System for Interactive Scientific Computing Python: An Open and General- Purpose Environment. *Computing in Science and Engineering* 9(3):21–29.