

Genetic Analysis of a Metazoan Pathway using Transcriptomic Phenotypes

immediate

This manuscript was compiled on January 24, 2017

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/mprrsq.

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

Genetic 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]. For the most part, the role of transcriptional profiling has been restricted target gene identification. In cell culture, single-cell RNA-seq has seen significant progress towards using transcriptomes as phenotypes with which to test genetic interactions [16, 17]. More recently, we have shown the first identification of a developmental state of *C. elegans* using whole-organism transcriptome profiling [18].

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

Hypoxia Inducible Factors are highly conserved in metazoans [20]. A common mechanism for hypoxia-response induction is heterodimerization between a HIF α and a HIF β subunit. The heterodimer then initiates transcription of target genes [21]. The number and complexity of HIFs varies throughout metazoans, with humans having three HIF α subunits and two HIF β subunits, whereas in the roundworm *Caenorhabditis elegans* (*C. elegans*) there is a single HIF α gene, *hif-1* [22] and a single HIF β gene, *ahr-1* [23]. HIF target genes have been implicated in a wide variety of cellular and extracellular processes such as glycolysis, extracellular matrix modification, autophagy and immunity [19, 24–27].

Levels of HIF α proteins tend to be tightly regulated. Under conditions of normoxia, HIF-1 α exists in the cytoplasm and partakes in a futile cycle of continuous protein production and rapid degradation [28]. HIF-1 α is hydroxylated by three

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.

²To whom correspondence should be addressed. E-mail: pws@caltech.edu

proline hydroxylases in humans (PHD1, PHD2 and PHD3) but is only hydroxylated by one proline hydroxylase (*egl-9*) in *C. elegans* [29]. HIF-1 hydroxylation increases its binding affinity to Von Hippel Lindau Tumor 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 [30].

Here, we show that transcriptomes contain strong, robust signals that can be used to infer relationships between genes in complex metazoans by reconstructing the hypoxia pathway in *C. elegans* using RNA-Seq. 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*. Using a single set of genome-wide measurements, we were able to observe and quantitatively assess significant fraction of the known transcriptional effects of *hif-1* in *C. elegans*. A complete, interactive version of the analysis is also available at www.wormlabcaltech.github.io/mprsq.

Results

The hypoxia pathway controls thousands of genes in *C. elegans*. We performed whole-organism RNA-seq of the hypoxia pathway at a relatively low sequencing depth (7 million mapped reads for each individual replicate). In spite of the low sequencing depth, transcriptome profiling of the hypoxia pathway revealed that this pathway controls thousands of genes in *C. elegans*. The *egl-9* transcriptome showed differential expression of 1,487 genes, similar to the 1,816 genes differentially expressed in *rhy-1* mutants. The *vhl-1* transcriptome showed considerably fewer differentially expressed genes (605), possibly reflecting the known fact that it is a weaker controller of *hif-1* than *egl-9* [31]. The *egl-9;vhl-1* double mutant transcriptome showed 1,989 differentially expressed genes. The *hif-1* mutant also showed a transcriptomic phenotype involving 481 genes. The *egl-9;hif-1* double mutant showed a similar number of genes with altered expression (364).

Clustering visualizes epistatic relationships between genes.

As a first step in our analysis, we analyzed our data using a general linear model (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 β . These coefficients are not identical to the average log-fold change per gene, although they are loosely related to this quantity. In general, larger magnitudes of β correspond to larger perturbations. These coefficients can be used to study the RNA-Seq data in question.

Clustering is a well-known technique in bioinformatics that is used to identify relationships between high dimensional data points [32]. 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

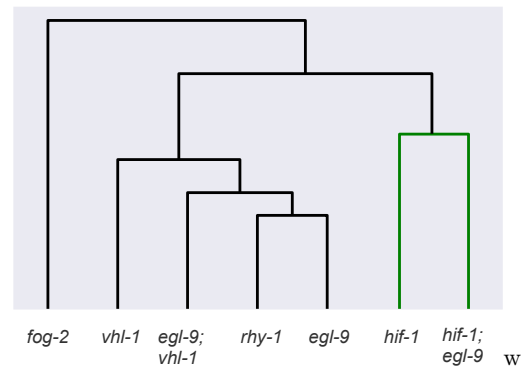


Fig. 1. Unsupervised aggregative clustering of various *C. elegans* mutants. Genes cluster in a manner that is biologically intuitive. Genes that inhibit *hif-1* (i.e. *egl-9*, *vhl-1*, and *rhy-1*) cluster far from *hif-1*. *hif-1* clusters with the suppressed *egl-9*; *hif-1* double mutant. A mutant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

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.

Reconstruction of the hypoxia pathway from first genetic principles.

Having shown that the signal in the mutants we selected was strong enough to cluster mutants using the regression coefficients, we set out to reconstruct the hypoxia pathway from first genetic principles. In general, to reconstruct a pathway, we must assess whether two genes act on the same phenotype (independence); then we must measure whether these genes act additively or epistatically on the measured phenotype; and if there is epistasis we must measure whether it is positive or negative, in order to assess whether the epistatic regulation is a genetic suppression or a synthetic interaction.

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

We observed that all the hypoxia mutants had significant overlap between their transcriptomes (fraction of shared transcriptomes ranged from a minimum of 65 genes shared between *hif-1* and *egl-9;hif-1* to a maximum of 1,249 shared genes between *egl-9* and *egl-9;vhl-1*). For comparison, we also analyzed a previously published *fog-2* transcriptome [18]. *fog-2* is involved in masculinization of the *C. elegans* germline, which enables sperm formation, and has not been described to be involved in the hypoxia pathway. The hypoxia pathway transcriptomes and the *fog-2* transcriptome showed similar overlap as the hypoxia pathway to itself (123–618 genes). Given the similar overlaps between known interactors and an unknown transcriptome, we conclude that the *fog-2* mutant we studied acts on the same phenotype as mutants from the hypoxia pathway.

Although overlapping transcriptomes may be enough to conclude that a set of mutants share a phenotype, we wanted

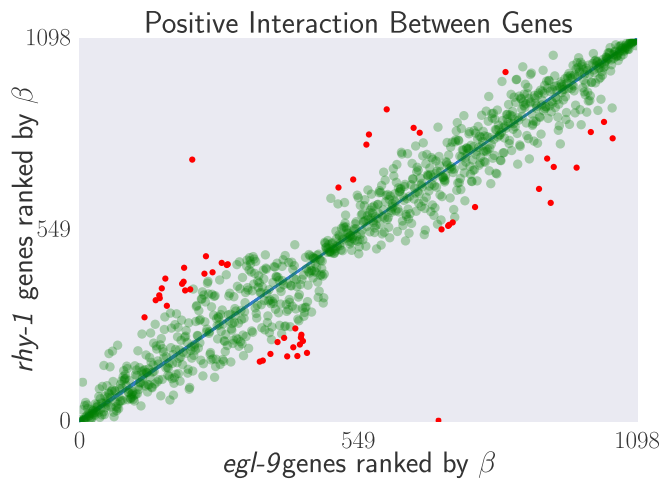


Fig. 2. Strong transcriptional correlations can be identified between genes that share a positive regulatory connection. We took the *egl-9* and the *rhy-1* transcriptomes, identified differentially expressed genes common to both transcriptomes and ranked each gene according to its differential expression coefficient β . 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. Green, transparent large points mark inliers to the regression (blue line); red, opaque, small points mark outliers to the regression. The two furthest outliers are annotated as pseudogenes in WormBase.

to know whether we could draw out more information from looking at quantitative agreement between perturbations. To this end, we rank-transformed the regression coefficients β for each transcriptome, and calculated lines of best fit using Bayesian regression with a Student-T distribution to mitigate noise from outliers (see Fig 2). For transcriptomes associated with the hypoxia pathway, we found that these correlations tended to have values as high as 0.98 with a tight distribution around the line of best fit, whereas the correlations for mutants from the hypoxia pathway with the *fog-2* mutant were considerably weaker, with magnitudes between 0.6–0.85 and a considerably larger spread around the line of best fit. Although *hif-1* is known to be genetically repressed by *egl-9*, *rhy-1* and *vhl-1* [], all the correlations between these genes and *hif-1* were negative. The overlap between *hif-1* and all other genes was 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. After we calculated the pairwise correlation between each transcriptome, we weighted the result of each regression by the number of differentially expressed isoforms shared by two transcriptomes and divided by the total number of differentially expressed isoforms present in the two transcriptomes, $N_{\text{overlap}}/N_{g1 \cup g2}$. The weighted regressions 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 extensive, which makes the weighting factor considerably larger than other pairs. The fine-grained nature of transcriptional phenotypes means that these weighted correlations between transcriptomes of single mutants are predictive of genetic interaction.

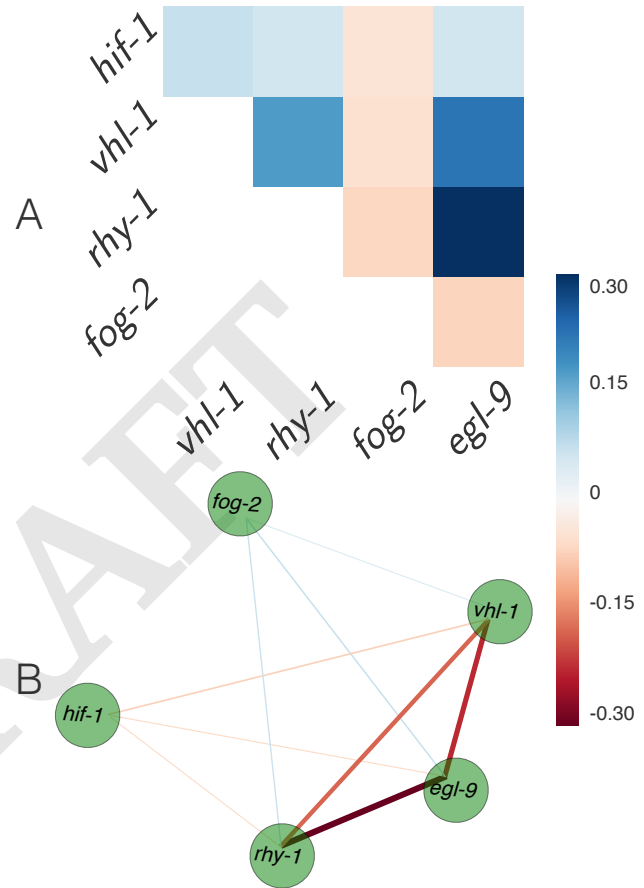


Fig. 3. A: Heatmap showing pairwise regression values between all single mutants. **B:** Correlation network drawn from the diagram. Edge width is proportional to the logarithm of the magnitude of the weighted correlation between two nodes divided by absolute value of the weighted correlation value of smallest magnitude. Edges are also colored according to the heatmap in **A**.

A quality check of the transcriptomic data reveals excellent agreement with the literature. One way to establish whether genes are acting additively or epistatically to each other is to perform qPCR of a reporter gene in the single and double mutants. This approach was used to successfully map the relationships within the hypoxia pathway (see, for example [31, 33]). A commonly used reporter is *nhr-57*, which is known to exhibit large changes in expression upon induction of HIF-1[26, 33–35]. Likewise, *rhy-1* and *egl-9* are both known to be up-regulated when HIF-1 becomes common in the cell [].

Our dataset enables us to perform an equivalent computational experiment to qPCR by selectively looking at expression of a few genes at a time. Therefore, we queried the changes in expression of *rhy-1*, *egl-9*, *nhr-57* and *lam-3* as a negative control. 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.

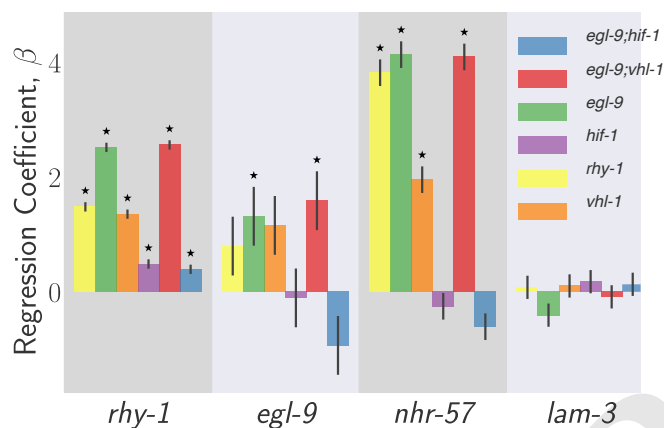


Fig. 4. Top: *In silico* qPCR. We extracted four genes (*rhy-1*, *egl-9*, *nhr-57* and *lam-3*, shown on the x-axis) and plotted their regression coefficients, β , as measured for every genotype (represented by one of six colors) to study the epistatic relationships between each gene. Stars above a bar represent a regression coefficient statistically significantly different from 0, meaning that expression is altered relative to a wild-type control. Error bars show standard error of the mean value of β . *nhr-57* is an expression reporter that has been used previously to identify *hif-1* regulators [31, 33]. 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 [36]. 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.

We also performed *in silico* qPCR of every gene under scrutiny to get a clearer idea of the relationships between them (see Fig. 4). 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 changes in *egl-9* expression were not statistically significant in *rhy-1* and *vhl-1* mutants, the mRNA levels of *egl-9* trended towards increased expression in these genotypes. As with *nhr-57*, the *egl-9* and *rhy-1* expression phenotypes were abrogated in the *egl-9;hif-1* mutant; whereas the *egl-9;vhl-1* mutant showed expression phenotypes identical to the *egl-9* mutant. 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*, which constitutes a novel observation. Taken together, these results indicate that RNA-seq data is at least equivalent to qPCR for purposes of comparing gene expression of a reporter between genotypes. Using a single reporter we would have been able to reconstruct an important fraction of the genetic relationships between the genes in the hypoxia pathway.

Genes in the hypoxia pathway exhibit genome-wide epistasis. As we have shown, it may be sufficient to extract the regression coefficients of a previously known reporter gene and study just that pattern in order rebuild a genetic pathway from RNA-seq data. However, we felt that by relying on a single gene, or even a handful of genes to rebuild the pathway was throwing out all of the valuable information present in our dataset. Therefore, we decided to explore a new epistatic metric—genome-wide epistasis.

Ideally, any measurement of genome-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 behaviour, such that the special values of the statistic (maximum, minimum, zero) should have an unambiguous interpretation.

One way of defining genome-wide epistasis is to use linear regressions to describe the relationship between the change in expression for a set of genes caused by a single mutant and the change in expression in the same set of genes caused by a double mutant containing the single mutant. The set of genes to be studied can be defined as the set of differentially expressed genes common to both genotypes. Once the set is defined, the regression coefficient of each gene in the single mutant can be plotted against the difference between the regression coefficients of the double mutant and the single mutant. We reasoned that under ideal conditions, such a plot would have an intuitive explanation. If two genes are acting entirely independently of each other, the plot will show a line with slope equal to 0. This is because is acting on entirely different sets of genes, so the perturbation caused a gene X^- is unchanged in the double mutant, X^-Y^- . If two genes are acting only additively, then the plot will show a line with slope > 0 (and in fact, the slope should be equal to the slope between a plot of the single mutants X^- and Y^-). If two genes share a negative regulatory interaction, then epistasis will be reflected in the plot as a line with a negative slope that should approach -1 , because the double mutant, X^-Y^- should have regression coefficients near 0, such that the y-axis becomes equal to $-X^-$. If the two genes have a synthetic interaction, we would expect that the slope must be positive and it must be greater than the slope predicted by an additive model.

In our experiment, we studied two double mutants, *egl-9;hif-1* and *egl-9;vhl-1*. We wanted to understand how well the global epistasis agreed with the literature based on qPCR of single reporters. Therefore, we fit weighted linear regressions to each of the four possible combinations (*egl-9* vs. *egl-9;hif-1*; *hif-1* vs. *egl-9;hif-1*; *egl-9* vs. *egl-9;vhl-1*; and *vhl-1* vs. *egl-9;vhl-1*) to measure the slopes of the lines of best fit.

We observe that the *egl-9;vhl-1* mutant has an identical phenotype to the *egl-9* single mutant (slope = 0; see Table. 1).

Table 1. Response Modeling of Double Mutants to Single Mutants

Double Mutant	Single Mutant	Δ	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. Δ 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.

On the other hand, *vhl-1* has a positive slope, indicating that *egl-9* is additive to *vhl-1*. However, this positive slope has to be less than the slope that would be predicted by an additive model because the slope between *egl-9;vhl-1* and *egl-9* is not statistically different from zero. 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, as is well-documented in the literature [31].

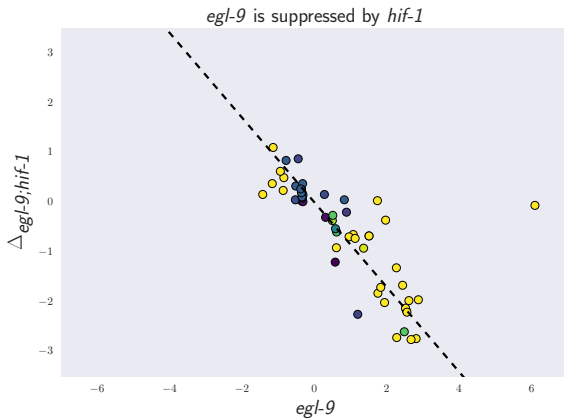


Fig. 5. The mutant *egl-9* transcriptomic phenotype is suppressed by mutations in *hif-1*. The graph shows the β coefficients for *egl-9* in the x-axis, and the change in β coefficient between the *egl-9;hif-1* and *egl-9* mutant. The dotted line is the regression line between the complete *egl-9* and *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*. It is probably downstream of *egl-9*, and is not likely a *hif-1* target.

On the other hand, comparison of the *egl-9;hif-1* double 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. 5). 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.

Because of the feedback between *hif-1* and *egl-9*, we expected a small subset of genes to be differentially expressed in every hypoxia pathway mutant. Therefore, we searched for genes that were differentially expressed in all our hypoxia mutants (except the *hif-1;egl-9* mutant because it has the least number of differentially expressed genes), reasoning that these genes should constitute an extremely high-quality picture of the hypoxia response, and should filter out other pathways. We identified 53 genes that satisfied these conditions, of which 10 genes were up-regulated in every mutant, and 13 genes were down-regulated. These genes constitute a core response around the circuit in question, and their behaviour should reflect the genetic relationships in our system the best. Although we performed the regressions using all the overlapped genes between the single and double mutants, when we plotted only these high-quality genes, we can see that they show beautiful agreement with the global regressions (see www.wormlabcaltech.github.io/mpsq for all interactive graphics).

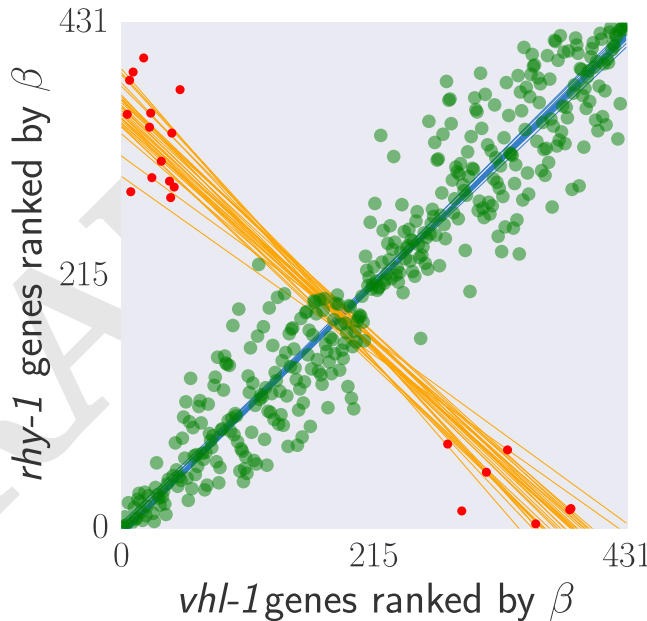


Fig. 6. Top: A feedback loop can generate transcriptomes that are both correlated and anti-correlated. **Bottom:** *hif-1* transcriptome correlated to the *rhy-1* transcriptome. Green large points are inliers to the first regression. Red small points are outliers to the first regression. Only the red 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.

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 pair of genes by $N_{\text{Intersection}}/N_{\text{Union}}$. 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_{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 pathway¹.

Discussion

Reconstructing Hypoxia Circuit in *C. elegans* using Transcriptomic Profiling. Previous work has established a circuit in which *rhy-1* leads to the activation of *egl-9*, 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* [25, 36]. 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* and non-oxygen dependent manner (see Fig. 8 top).

Using only information gathered from transcriptomic profiling, we were able to reconstruct the known regulatory relationships between *egl-9*, *rhy-1* and *vhl-1* using a combination of inter-transcriptome correlations and epistasis measurements. Using clustering as a proxy for phenotype, we were able to infer the relationship between *egl-9* and *hif-1*. Alternatively, we could have used our epistasis measurements to conclude that *egl-9* inhibits *hif-1*. By looking at single gene measurements using *in silico* qPCR, we were able to determine that *rhy-1* transcription is stimulated by HIF-1. Single gene measurements also suggested that *egl-9* is transcriptionally downstream of *hif-1*, although only the measurements from *egl-9* mutants were statistically significantly different from 0 and all other *hif-1*-constitutive mutants only trended towards increased *egl-9* expression. In addition, our dataset also identified *cysl-1* as a *hif-1*-responsive gene. We were unable to detect *vhl-1* transcripts in our RNA-seq measurements, and therefore we cannot rule out effects of *hif-1* on *vhl-1* using

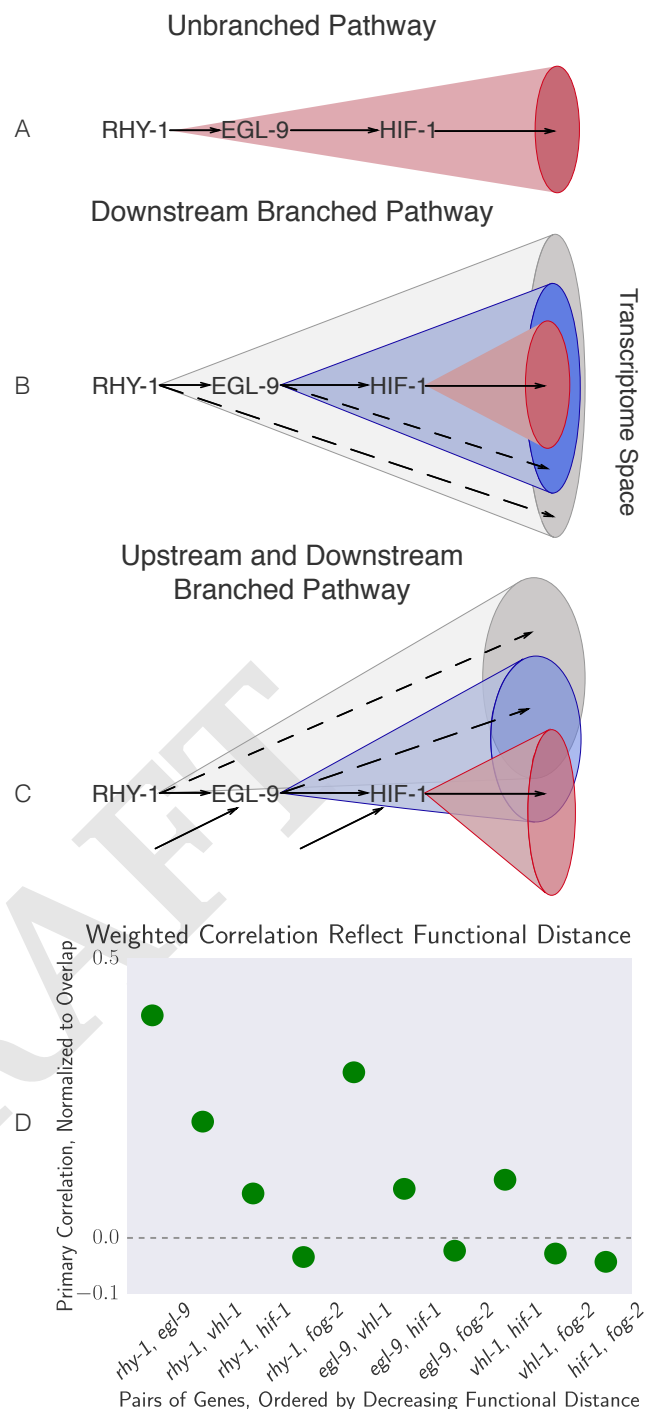


Fig. 7. Theoretically, transcriptomes can be used to order genes in a pathway under certain assumptions. Arrows in the diagrams above are intended to show the direction of flow, and do not indicate valence. **A** A linear pathway in which *rhy-1* is the only gene controlling *egl-9*, which in turn controls *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.

¹ 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.

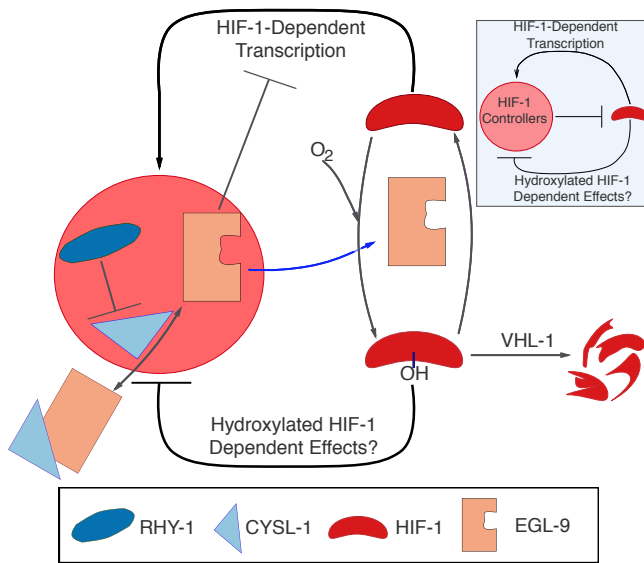


Fig. 8. A schematic of the hypoxia pathway in *C. elegans*. RHY-1 likely inhibits a protein-protein interaction between EGL-9 and CYSL-1 [30]. This interaction inhibits the activity of EGL-9. When active, EGL-9 can inhibit HIF-1 by catalyzing the hydroxylation of HIF-1, which leads to its recognition 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.

this data. All of these measurements are consistent with the previous literature.

Although we reconstructed the reported pathway, some of our data cannot be easily explained under a model in which HIF-1 is transcriptionally active but hydroxylated HIF-1 is not. For example, the *hif-1* transcriptome correlated positively with the *egl-9*, *rhy-1* and *vhl-1* transcriptomes. Moreover, our *in silico* qPCR suggests that *rhy-1* expression levels are increased upon mutation of the *hif-1* gene, which is hard to explain under the standard model, since *hif-1* is an activator of *rhy-1*, most HIF-1 is expected to be hydroxylated under normoxic conditions, and total HIF-1 are not anywhere near their hypoxic levels. The simplest explanation that could explain these results without involving another gene would be to postulate that hydroxylated HIF-1 is an active player with downstream transcriptional consequences, and that one of these consequences is repression of *rhy-1*.

Such a model is appealing because the network that emerges has certain characteristics of homeostatic pathways. First, under oxygen deplete conditions, HIF-1 will no longer be hydroxylated and it will initiate transcription of the hypoxic response. Upregulation of its controller genes generates a negative feedback loop which has positive consequences for the response time of a HIF-1 growth curve [37]. On the other hand, we can envision an environment in which HIF-1 is completely hydroxylated. Presumably, an organism always needs non-hydroxylated HIF-1, or alternatively, needs hydroxylated HIF-1 levels to be below a certain threshold. If hydroxylated HIF-1 can mediate a transcriptional response, then the organism benefits from complete information about the iron and oxygen

balance within it, as both species of HIF-1 signal with strengths proportional to their total levels in the cell.

A model in which the hydroxylation of HIF-1 alters, but does not abrogate, the transcriptional effects of this gene would explain the positive correlations between *hif-1*, *egl-9*, and *rhy-1*. However, we must emphasize that the poor agreement between overlaps (only 5 genes appear in both intersections, and each intersection consists of less than 200 genes) prevents us from making a more complete statement. Moreover, this newly proposed model does not explain why *vhl-1* has a positive correlation with *hif-1*. It may be the case then that another gene is involved in this pathway that positively links *vhl-1* and *hif-1*. An interaction of this sort would explain why the correlation between *vhl-1* and *hif-1* is mediated by genes that are completely different from *egl-9* and *rhy-1*.

Towards A Genetic Theory of Transcriptomics. We have shown that transcriptomes contain sufficient information to be used as semi-quantitative phenotypes in complex 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. As a result of this dynamic range, 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 detectable, reproducible phenotypes at the transcriptomic level, which would facilitate the study of small-effect genes and other quantitative traits.

We have formalized the concept of genetic transcriptome by developing a controlled language and formal notation to study these objects. In our language, perturbations are unary operators that act on a single gene, which has consequences for the function it controls (its transcriptome). This notation makes it possible to think about transcriptome genetics in a manner analogous to classical genetics using scalar phenotypes, and also makes identification of specific transcriptomes easier and more rigorous. For example, we concluded that *rhy-1* does not have a specific transcriptome. All of its transcriptomic consequences appear to emanate from the downstream gene *egl-9*. We arrived at this conclusion through theoretical considerations that suggested that a *rhy-1* specific transcriptome should manifest in a negative correlation when the *egl-9* and *rhy-1* transcriptomes are plotted on a rank-plot. Using simple binary exclusion principles one might arrive at the conclusion that the *rhy-1* specific transcriptome is defined as the genes that are not present in the intersection between *egl-9* and *rhy-1*, which would result in almost 400 genes being assigned to *rhy-1*. Although such subtractive logic is attractive at first sight, it is fundamentally untestable due to the large number of genes involved. On the other hand, using our notation and logic, we can directly query the data for signs of a specific transcriptome. If those signs are not missing, we have some degree of confidence that the specific transcriptome was not measured, could not be identified, was not perturbed or does not exist. Further tests can be envisioned to test each possibility.

Of particular interest to us was the idea that the logic necessary to understand transcriptome genetics is fundamentally

simple. It builds in a logical manner from genetics principles derived from macroscopic observations, and makes predictions that are easily testable through linear regressions. This contrasts starkly with preconceived notions of complicated bioinformatics that are believed to be necessary to extract information from next-generation sequencing data. Even more striking is the fact that transcriptomes seem to exhibit quantitative epistatic behaviour. Finally, we note that the linear structure of the logic we have developed hints at the possibility that an algebra exists that encompasses these rules. Finding such an algebra would constitute an important development which could have important consequences for the kind of algorithms that are developed to deconvolute, dissect and understand transcriptomes in the future.

Materials and Methods

Nematode strains and culture. Strains used were N2 wild-type Bristol, CB5602 *vhl-1*(ok161), CB6088 *egl-9*(sa307) *hif-1*(ia4), CB6116 *egl-9*(sa307) *vhl-1*(ok161), JT307 *egl-9*(sa307), ZG31 *hif-1*(ia4), RB1297 *rhy-1*(ok1402). ZG31 *hif-1*(ia4) is a null mutant of *hif-1* which deletes 1231 bp of the second, third and fourth exons. JT307 contains the null mutant *egl-9*(sa307) which is a 243 bp deletion. RB1297 contains null mutation *rhy-1*(ok1402) with an estimated 700 bp deletion constructed by the OMRF Knockout Group. CB5602 contains the deletion mutation of *vhl-1*(ok161). CB6088 contains *egl-9*(sa307);*hif-1*(ia4). CB6116 contains *egl-9*(sa307) *vhl-1*(ok161). All strains were provided by the CGC, which is funded by NIH Office of Research Infrastructure Programs (P40 OD010440). All lines were grown on standard nematode growth media (NGM) plates with seeded with OP50 *E. coli* at 20°C (Brenner 1974).

RNA Isolation. Unsynchronized lines were grown on NGM plates at 20°C and eggs harvested by sodium hypochlorite treatment. Eggs were plated on 6 to 9 small 5cm NGM plates with ample OP50 *E. coli* at a density chosen to avoid starvation and grown at 20°C. Worms were staged and harvested based on the time after plating, vulva morphology and the absence of eggs. Approximately 30–50 non-gravid young adults (YA) were picked and placed in 100µL of TE pH 8.0 at 4°C in 0.2mL PCR tubes. After settling and a brief spin in microfuge approximately 80µL of TE was removed from the top of the sample and individual replicates were snap frozen in liquid N₂. These replicate samples were then digested with Proteinase K for 15min at 60° in the presence of 1% SDS and 1.25µL RNA Secure (Ambion AM 7005). RNA samples were then taken up in 5 Volumes of Trizol (Tri Reagent Zymo Research) and processed and treated with DNAase I using Zymo MicroPrep RNA Kit (Zymo Research Quick-RNA MicroPrep R1050). RNA was eluted in dH₂O and divided into aliquots and stored at -80°C. One aliquot of each replicate was analyzed by both NanoDrop for impurities, Qubit for concentration and then analyzed on an Agilent 2100 BioAnalyzer. Replicates were selected that had RNA integrity numbers (RIN) equal or greater than 9.0 and showed no evidence of bacterial ribosomal bands, except for the ZG31 mutant where one of three replicates had a RIN of 8.3.

Library Preparation and Sequencing. Forthcoming.

Read Alignment and Differential Expression Analysis. We used Kallisto to perform read pseudo-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 [38–46]. Our analysis is available in a Jupyter Notebook [47]. All code and required data (except the raw reads) are available at <https://github.com/WormLabCaltech/mpsrq> along with version-control information. Our Jupyter Notebook and interactive graphs for this project can be found at <https://wormlabcaltech.github.io/mpsrq/>. Raw reads were deposited at XXXXXXXXXXXX

ACKNOWLEDGMENTS. This article was written with support of the Howard Hughes Medical Institute. We would like to thank Jonathan Liu, Han Wang, and Porfirio Quintero for helpful discussion.

- 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.
- Hughes TR et al. (2000) Functional Discovery via a Compendium of Expression Profiles. *Cell* 102(1):109–126.
- Van Driessche N et al. (2005) Epistasis analysis with global transcriptional phenotypes. *Nature genetics* 37(5):471–477.
- 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.
- Metzker ML (2010) Sequencing technologies - the next generation. *Nature reviews. Genetics* 11(1):31–46.
- 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.
- Bray NL, Pimentel H, Melsted P, Pachter L (2016) Near-optimal probabilistic RNA-seq quantification. *Nature biotechnology* 34(5):525–7.
- Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C (2016) Salmon provides accurate, fast, and bias-aware transcript expression estimates using dual-phase inference. *bioRxiv* p. 021592.
- Pimentel HJ, Bray NL, Puente S, Melsted P, Pachter L (2016) Differential analysis of RNA-Seq incorporating quantification uncertainty. *bioRxiv* p. 058164.
- Trapnell C et al. (2013) Differential analysis of gene regulation at transcript resolution with RNA-seq. *Nature biotechnology* 31(1):46–53.
- 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.
- Shalek AK et al. (2013) Single-cell transcriptomics reveals bimodality in expression and splicing in immune cells. *Nature* 498(7453):236–40.
- 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.
- Van Wolfswinkel JC, Wagner DE, Reddien PW (2014) Single-cell analysis reveals functionally distinct classes within the planarian stem cell compartment. *Cell Stem Cell* 15(3):326–339.
- Scimone ML, Kravarik KM, Lapan SW, Reddien PW (2014) Neoblast specialization in regeneration of the planarian *schmidtea mediterranea*. *Stem Cell Reports* 3(2):339–352.
- Adamson B et al. (2016) A Multiplexed Single-Cell CRISPR Screening Platform Enables Systematic Dissection of the Unfolded Protein Response. *Cell* 167(7):1867–1882.e21.
- Dixit A et al. (2016) Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens. *Cell* 167(7):1853–1866.e17.
- Angeles-Albore D et al. (2016) Transcriptomic Description of an Endogenous Female State in *C. elegans*. *bioRxiv*.
- Semenza GL (2012) Hypoxia-inducible factors in physiology and medicine. *Cell* 148(3):399–408.
- Loenarz C et al. (2011) The hypoxia-inducible transcription factor pathway regulates oxygen sensing in the simplest animal, *Trichoplax adhaerens*. *EMBO reports* 12(1):63–70.
- Jiang BH, Rue E, Wang GL, Roe R, Semenza GL (1996) Dimerization, DNA binding, and transactivation properties of hypoxia-inducible factor 1. *The Journal of biological chemistry* 271(30):17771–17778.
- Jiang H, Guo R, Powell-Coffman JA (2001) The *Caenorhabditis elegans* hif-1 gene encodes a bHLH-PAS protein that is required for adaptation to hypoxia. *Proceedings of the National Academy of Sciences of the United States of America* 98(14):7916–7921.
- Powell-Coffman JA, Bradfield CA, Wood WB (1998) *Caenorhabditis elegans* Orthologs of the Aryl Hydrocarbon Receptor and Its Heterodimerization Partner the Aryl Hydrocarbon Receptor Nuclear Translocator. *Proceedings of the National Academy of Sciences* 95(6):2844–2849.
- Semenza GL, Roth PH, Fang HM, Wang GL (1994) Transcriptional regulation of genes encoding glycolytic enzymes by hypoxia-inducible factor 1. *The Journal of biological chemistry* 269(38):23757–63.
- Bishop T et al. (2004) Genetic Analysis of Pathways Regulated by the von Hippel-Lindau Tumor Suppressor in *Caenorhabditis elegans*. *PLoS Biology* 2(10).
- 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.
- Bellier A, Chen CS, Kao CY, Cinar HN, Aroian RV (2009) Hypoxia and the hypoxic response pathway protect against pore-forming toxins in *C. elegans*. *PLoS Pathogens* 5(12).
- Huang LE, Arany Z, Livingston DM, Franklin Bunn H (1996) Activation of hypoxia-inducible transcription factor depends primarily upon redox-sensitive stabilization of its alpha subunit. *Journal of Biological Chemistry* 271(50):32253–32259.

29. Kaelin WG, Ratcliffe PJ (2008) Oxygen Sensing by Metazoans: The Central Role of the HIF Hydroxylase Pathway. *Molecular Cell* 30(4):393–402.
30. Ma DK, Vozdek R, Bhatia N, Horvitz HR (2012) CYSL-1 Interacts with the O₂-Sensing Hydroxylase EGL-9 to Promote H₂S-Modulated Hypoxia-Induced Behavioral Plasticity in *C. elegans*. *Neuron* 73(5):925–940.
31. 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.
32. Yeung KY, Medvedovic M, Bumgarner RE (2003) Clustering gene-expression data with repeated measurements. *Genome biology* 4(5):R34.
33. 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.
34. Ackerman D, Gems D (2012) Insulin/IGF-1 and hypoxia signaling act in concert to regulate iron homeostasis in *Caenorhabditis elegans*. *PLoS Genetics* 8(3).
35. Park EC et al. (2012) Hypoxia regulates glutamate receptor trafficking through an HIF-independent mechanism. *The EMBO Journal* 31(6):1618–1619.
36. Powell-Coffman JA (2010) Hypoxia signaling and resistance in *C. elegans*. *Trends in Endocrinology and Metabolism* 21(7):435–440.
37. Alon U (2007) Network motifs: theory and experimental approaches. *Nat Rev Genet* 8(6):450–61.
38. Bokeh Development Team (2014) Bokeh: Python library for interactive visualization.
39. McKinney W (2011) pandas: a Foundational Python Library for Data Analysis and Statistics. *Python for High Performance and Scientific Computing* pp. 1–9.
40. Oliphant TE (2007) SciPy: Open source scientific tools for Python. *Computing in Science and Engineering* 9:10–20.
41. Pedregosa F et al. (2012) Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12:2825–2830.
42. Salvatier J, Wiecki T, Fonnesbeck C (2015) Probabilistic Programming in Python using PyMC. *PeerJ Computer Science* 2(e55):1–24.
43. 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.
44. Hunter JD (2007) Matplotlib: A 2D graphics environment. *Computing in Science and Engineering* 9(3):99–104.
45. Angeles-Albores D, N. Lee RY, Chan J, Sternberg PW (2016) Tissue enrichment analysis for *C. elegans* genomics. *BMC Bioinformatics* 17(1):366.
46. Waskom M et al. (2016) seaborn: v0.7.0 (January 2016).
47. 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.

DRAFT