



Published in final edited form as:

Nat Genet. 2009 March ; 41(3): 299–307. doi:10.1038/ng.332.

## Systems Genetics of Complex Traits in *Drosophila melanogaster*

Julien F. Ayroles<sup>1,4,6</sup>, Mary Anna Carbone<sup>1,4,6</sup>, Eric A. Stone<sup>3,6</sup>, Katherine W. Jordan<sup>1,4</sup>, Richard F. Lyman<sup>1,4</sup>, Michael M. Magwire<sup>1,4,5</sup>, Stephanie M. Rollmann<sup>1,4,5</sup>, Laura H. Duncan<sup>1,4</sup>, Faye Lawrence<sup>1,4</sup>, Robert R. H. Anholt<sup>1,2,4</sup>, and Trudy F. C. Mackay<sup>1,4</sup>

<sup>1</sup>Department of Genetics, North Carolina State University, Raleigh, NC 27695, USA.

<sup>2</sup>Department of Biology, North Carolina State University, Raleigh, NC 27695, USA.

<sup>3</sup>Department of Statistics, North Carolina State University, Raleigh, NC 27695, USA.

<sup>4</sup>Department of W. M. Keck Center for Behavioral Biology, North Carolina State University, Raleigh, NC 27695, USA.

### SUMMARY

Determining the genetic architecture of complex traits is challenging because phenotypic variation arises from interactions between multiple, environmentally sensitive alleles. We quantified genome-wide transcript abundance and phenotypes for six ecologically relevant traits in *D. melanogaster* wild-derived inbred lines. We observed 10,096 genetically variable transcripts and high heritabilities for all organismal phenotypes. The transcriptome is highly genetically inter-correlated, forming 241 transcriptional modules. Modules are enriched for transcripts in common pathways, gene ontology categories, tissue-specific expression, and transcription factor binding sites. The high transcriptional connectivity allows us to infer genetic networks and the function of predicted genes based on annotations of other genes in the network. Regressions of organismal phenotypes on transcript abundance implicate several hundred candidate genes that form modules of biologically meaningful correlated transcripts affecting each phenotype. Overlapping transcripts in modules associated with different traits provides insight into the molecular basis of pleiotropy between complex traits.

---

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:[http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

Correspondence should be addressed to T. F. C. M. ([trudy\\_mackay@ncsu.edu](mailto:trudy_mackay@ncsu.edu)).

<sup>5</sup>Present address: Institute of Evolutionary Biology, School of Biological Sciences, University of Edinburgh, West Mains Road, Edinburgh EH9 3JT, UK (M. M. M.); Department of Biological Sciences, University of Cincinnati, PO Box 210006, 614 Rieveschl Hall, Cincinnati, OH 45221-0006, USA (S. M. R.).

<sup>6</sup>These authors contributed equally to this work.

**Author Contributions** T. F. C. M., J. F. A., E. A. S. and R. R. H. A. wrote the paper. R. F. L. constructed the *Drosophila* lines. M. A. C. obtained the gene expression data. K. W. J., M. M. M., S. M. R., L. H. D. and F. L. obtained the organismal phenotype data. J. F. A., E. A. S. and K. W. J. performed the statistical analyses.

**Author Information** The raw microarray data are deposited in the ArrayExpress database ([www.ebi.ac.uk/arrayexpress](http://www.ebi.ac.uk/arrayexpress)) under accession number E-MEXP-1594. Fly stocks are available from the Bloomington *Drosophila* Stock Center (Bloomington, Indiana). Additional data are available at [mackaylab.ncsu.edu](http://mackaylab.ncsu.edu). The authors declare no competing financial interests.

## INTRODUCTION

Natural populations harbor a wide range of phenotypic variation for all aspects of morphology, physiology, behaviors and disease susceptibility. Knowledge of the genetic basis of this variation is important for understanding adaptive evolution, deriving elite domestic crop and animal strains, and improving human health. However, determining the genetic architecture of natural phenotypic variation is challenging because most phenotypic variation is attributable to segregating alleles at many interacting genes with environmentally sensitive effects<sup>1, 2</sup>.

Recent genome-wide association studies in large human populations show that most quantitative traits and diseases are associated with loci with small effects that in total account for a few percent of the phenotypic variance<sup>3</sup>. These findings suggest that the bulk of genetic variation for complex traits is due to alleles with small and possibly context-dependent effects. Further, the significant polymorphisms are often in genes with no *a priori* expected relationship to the trait, in computationally predicted genes, or in gene deserts. Studies in model organisms have anticipated these results. In *Drosophila*, quantitative analysis of subtle effects of new mutations have revealed large numbers of novel loci affecting quantitative traits<sup>4</sup>, as have high resolution maps of segregating quantitative trait loci (QTLs) in *Drosophila*<sup>4</sup> and mice<sup>5</sup>. Single alleles often have pleiotropic effects on multiple traits, epistatic interactions among loci affecting the same trait are common, and allelic effects can be conditional on sex and the environment<sup>4</sup>.

Our understanding of the genetic architecture of quantitative traits in model systems, and ultimately humans, will benefit from interrogating a single resource population for variation in DNA sequence, transcript abundance, proteins and metabolites; for multiple organismal phenotypes; and in multiple environments. This ‘systems genetics’ approach will yield a detailed map of genetic variants associated with each organismal phenotype in each environment, provide a functional context for interpreting the phenotypes, elucidate the genetic underpinnings that govern the interdependence of multiple phenotypes, and address the long-standing question of the genetic basis of genotype by environment interaction<sup>6–8</sup>.

Here, we report the first step of a systems genetics analysis of the genetic basis of complex traits in *Drosophila*. We demonstrate ubiquitous variation in transcript abundance among inbred *Drosophila* strains recently derived from the wild, and show that the variable transcripts can be grouped into coherent modules of inter-correlated genes. These lines harbor substantial genetic variation for ecologically relevant complex traits, and variation for hundreds of transcripts is associated with variation for each of the organismal traits. Transcripts associated with each trait form correlated transcriptional modules from which we can construct networks of interacting genes with plausible biological relationships to each other and the traits. These genetic networks provide additional functional annotation of the *Drosophila* genome, and an integrated context in which to frame predictions of the behavior of the system following genetic or environmental perturbations.

## Natural variation in transcript abundance

We derived 40 highly inbred lines from the Raleigh, NC, natural population. These lines are a living library of common polymorphisms affecting complex traits. We assessed whole genome variation in transcript abundance for young males and females of each of these lines using Affymetrix *Drosophila* 2.0 arrays. We standardized the raw array data by median centering and used a statistical approach to identify outlier probes in each perfect match (PM) probe set that contained single feature polymorphisms (SFPs) between the wild-derived lines and the reference strain sequence used to design the array. We used the median log<sub>2</sub> signal intensity of the remaining PM probes in each probe set as the measure of expression. Of the 18,800 transcripts on the array, 14,840 (78.9%) were expressed in young adults (Supplementary Table 1). Many genes that have been characterized for their role in early development are also expressed in adults, and may affect adult phenotypes that cannot be predicted from their prior developmental functions<sup>9, 10</sup>. We used analysis of variance (ANOVA) to partition variation in expression between sexes, among lines, and the sex × line interaction for each expressed transcript.

The sex term was significant at a conservative false discovery rate (FDR<sub>11</sub>) of 0.001 for 13,086 (88.2%) of the expressed transcripts (Fig. 1a, Supplementary Table 1), indicating pervasive sexual dimorphism for gene expression. A total of 3,255 transcripts had two-fold or greater differences in expression between the sexes: 1,690 with female-biased expression, and 1,565 with male-biased expression. Previous studies reported largely male-biased expression in *D. melanogaster*<sup>12, 13</sup>. Our observation of nearly equal numbers of transcripts with strong male and female expression bias is likely attributable to our greater power to detect smaller sex-biased effects, since the absolute magnitude of the difference in expression between the sexes is less for female-biased than for male-biased transcripts (Fig. 1a). Gene ontology analysis<sup>14</sup> showed that both male- and female-biased transcripts were enriched for genes affecting reproduction and gametogenesis. The female-biased transcripts were also enriched for genes affecting basic cellular processes, and the male-biased transcripts for genes involved in reproductive behavior and physiology, mitochondrial energy metabolism and intermediary metabolism (Supplementary Table 2).

The line term was significant at an FDR < 0.001 for 10,096 (68.0%) of the expressed transcripts, indicating considerable genetic variation in gene expression (Supplementary Table 1). Estimates of broad sense heritability ( $H^2$ ) for significant transcripts ranged from  $\approx 0.3 - 1.0$  (Fig 1b). Transcripts with high levels of genetic variation ( $H^2 > 0.8$ ) were enriched for genes involved with responses to the environment, while transcripts with low levels of genetic variation ( $H^2 < 0.2$ ) were enriched for genes affecting processes essential for survival (Supplementary Table 2). The overall correlation ( $r$ ) between  $H^2$  and mean expression level was low, although statistically significant ( $r = 0.078$ ,  $P = 3.13 \times 10^{-21}$ ). The high level of genetic variation in gene expression is partly attributable to the doubling of the additive genetic variation of an outbred population in a population of fully inbred lines, and inflation of broad sense heritability estimates by any non-additive genetic variance<sup>1</sup>. Significant heritability of abundance of a particular transcript does not necessarily mean that *cis*-acting genetic polymorphisms cause the variation; it could be due to *trans*-regulation by another genetically variable transcript<sup>6–8</sup>.

A significant sex  $\times$  line interaction indicates genetic variation in the magnitude of the sex dimorphism among the lines, or equivalently, a significant departure of the cross sex genetic correlation ( $r_{MF}$ ) from unity. The sex  $\times$  line interaction was significant for 4,108 (40.7%) of the expressed transcripts at an FDR  $< 0.001$  (Supplementary Table 1). The average cross-sex genetic correlation of the transcripts exhibiting genetic variation in sexual dimorphism was quite low ( $r_{MF} = 0.234$ ), with a mode at  $r_{MF} = 0$  (Fig. 1c). Variation in expression of many transcripts among the lines was uncorrelated between males and females, arguing for caution when extrapolating inferences about variation in gene expression from expression profiles drawn from one sex to the other sex. These data also reveal great potential for rapid evolution of sex-biased gene expression, as observed when different *Drosophila* species are compared<sup>13</sup>. Low cross-sex genetic correlations were only partly attributable to transcripts that were expressed and variable in only one sex. Many of the transcripts exhibiting variation in sex dimorphism in expression were actually expressed in both sexes, but had much higher heritabilities in one sex compared to the other (Fig. 1d). Transcripts with low cross-sex genetic correlations ( $r_{MF} < 0.2$ ) were enriched for the same gene ontology categories as sex-biased genes, indicating that sex-biased transcripts are also genetically variable in the sex in which they are highly expressed (Supplementary Table 2).

The patterning of genetic variation within a population depends on effective population size, recombination rate, and the selection coefficient of new mutations<sup>1</sup>. These factors vary among chromosomes: *X* chromosomes have smaller effective population sizes than autosomes and are hemizygous in males, while recombination is severely reduced on the *Drosophila* fourth chromosome. Therefore, we asked whether there were differences among chromosomes in mean level and genetic variance of transcript abundance. Consistent with previous studies<sup>12, 15</sup>, we found that male-biased transcripts were strongly underrepresented on the *X* chromosome (and overrepresented on chromosome 2*L*). In contrast, female-biased transcripts were strongly overrepresented on the *X* chromosome (Fig. 1e). Possibly the *X* chromosome is a hospitable location for female- but not male-biased genes because mutations with *X*-linked deleterious effects on male fitness are quickly purged from populations, but mutations with recessive *X*-linked deleterious effects on female fitness can achieve higher frequencies since they are protected from natural selection when they are rare. We found differences in overall transcript abundance among the chromosomes for both sexes ( $P < 0.0001$ ), with chromosome 4 having the highest mean expression level and the *X* chromosome the lowest mean expression level (data not shown). We also found differences in  $H^2$  of transcript abundance between the chromosomes ( $P < 0.0001$ ), largely attributable to reduced genetic variation on the *X* and fourth chromosomes (data not shown).

## Modules of correlated transcripts

We assessed the extent to which the 10,096 variable transcripts were genetically correlated among the lines. We computed pairwise genetic correlations among all the variable transcripts, and computed the mean absolute value of the pairwise genetic correlations of each transcript with all other transcripts as a measure of average connectivity ( $|r|$ , Supplementary Table 1). The distribution of  $|r|$  was strongly skewed towards high values, with a mode at 0.6 (Fig. 2a). Thus, the genome as a whole is highly genetically correlated at

the transcriptional level, which imposes constraints on the evolution of transcriptional genetic networks. There is a strong inverse correlation ( $r = -0.263$ ,  $P = 1.08 \times 10^{-159}$ ) between transcript heritability and average connectivity – the most highly heritable transcripts have low mean values of  $|r|$  (Fig. 2b) and are therefore presumably less evolutionarily constrained. Indeed, there is a significant positive correlation between heritability of transcript abundance and  $\omega$ , the ratio of non-synonymous to synonymous substitutions, among single copy genes with orthologues in six *melanogaster* group species<sup>16</sup> ( $r = 0.132$ ,  $P = 7.56 \times 10^{-25}$ ). Genes encoding transcripts with lower heritabilities experience stronger purifying selection than do genes encoding transcripts with high heritabilities. Although high heritabilities are expected for genes under mutation-drift equilibrium<sup>1</sup>, it is not likely that this mechanism accounts for the observed high heritabilities of transcript abundance, since the estimates of  $\omega$  for these loci are much less than the neutral expectation of unity. In addition, the high heritability transcripts are predominantly for genes affecting responses to the environment, which have been associated with responses to artificial selection for multiple traits<sup>17–20</sup>. Therefore, the high heritabilities for these transcripts could be the result of more complex evolutionary dynamics.

We grouped the genetically variable transcripts into modules using a novel method to identify separable clusters of highly interconnected genes (E.A.S. and J. F. A, personal communication). The 10,096 transcripts fell into 241 modules (Fig. 2c, Supplementary Table 1). The two largest modules (7 and 18) consisted of 1,765 and 4,128 transcripts, with average absolute intra-module correlations of 0.89 and 0.77, respectively. Moreover, the expression of genes in these two modules was strongly negatively correlated among the lines (Fig. 2c). We performed gene ontology enrichment analyses<sup>14</sup> and found that Module 7 was enriched for the same functional categories as male-biased transcripts, and Module 18 was enriched for the same functional categories as female-biased transcripts. We found significant overrepresentations of male-biased genes in Module 7 (1,241 of 1,565 male-biased genes,  $\chi^2_1=4,910$ ;  $P \approx 0$ ) and of female-biased genes in Module 18 (1,381 of 1,690 female-biased genes,  $\chi^2_1=1,400$ ;  $P \approx 0$ ). Thus, the negative correlation between Modules 7 and 18 is attributable to higher levels of expression of genes in Module 7 in males than females, and higher levels of expression of genes in Module 18 in females than males.

There are strong correlations among transcripts in the same functional pathways (for example, amino sugars metabolism and the Notch signaling pathway, Fig. 2d). Genes within a module often show similar tissue-specific expression patterns<sup>21</sup> which suggests a common biological function, genetic variation in organ size, or both (Fig. 3a). Many modules are enriched for known transcription factor binding sites (Supplementary Table 3, Fig. 3b). The high degree of connectivity among transcript variation in a natural population allows us to infer potential interacting partners of focal genes in networks for functional studies. For example, *disco*, *disco-r* and *tsh* interact during embryonic and larval development<sup>22</sup>, and have genetically inter-correlated transcripts in adults in Module 161 (Supplementary Table 1, Fig. 3c). These genes are expressed in the adult nervous system<sup>21</sup>, and are correlated with three other genes in this module (*unc-5*, *drl*, *argos*) that are involved in nervous system development and axon guidance<sup>23</sup> (Fig. 3c). Transcriptional correlation also enables

functional annotation of computationally predicted genes based on known annotations of other genes in the network. *CG15065* is a putative Immune-induced Molecule (*IM*) gene, based on its strong transcriptional correlations with *IM1* ( $r = 0.74$ ), *IM2* ( $r = 0.63$ ) and *IM3* ( $r = 0.67$ ), physical location adjacent to these genes, and remarkable protein sequence similarity to *IM1* and *IM2* (Fig. 3d).

## Associations with organismal phenotypes

We quantified variation among the 40 Raleigh inbred lines for several ecologically relevant traits (resistance to starvation stress, time to recover from a chill-induced coma, life span, a startle-induced locomotor response, and mating speed) as well as a measure of competitive fitness. We found substantial genetic variation for all traits, with estimates of  $H^2$  between 0.25 – 0.58 (Fig. 4a–f, Supplementary Table 4). The range of variation among these lines is comparable to the difference in mean phenotype between lines subjected to long-term divergent artificial selection for these traits<sup>17, 19</sup>. We observed significant sex  $\times$  line interactions for starvation stress resistance, life span and chill coma recovery time (Supplementary Table 4). Estimates of  $r_{MF}$  ( $\pm$  SE) were high for organismal phenotypes ( $0.72 \pm 0.11$ ,  $0.73 \pm 0.11$  and  $0.87 \pm 0.08$  for starvation stress resistance, life span and chill coma recovery, respectively), in contrast to the low cross-sex genetic correlations observed at the level of transcripts.

We asked whether there were significant genetic correlations among these traits, as would occur if segregating alleles have pleiotropic effects on two traits in the same direction<sup>1, 2</sup>. Only five genetic correlations were significantly different from zero (Supplementary Table 4). There is a tendency for lines that recover from chill coma quickly to have high competitive fitness and mate rapidly, but at the expense of surviving starvation stress. Lines that are resistant to starvation stress tend to have longer life spans, but reduced competitive fitness. Thus, there is a trade-off between genetic variants affecting recovery from different environmental stresses. We do not observe high positive correlations between all traits with each other and with fitness, as would be the case if variation among the lines was attributable to the fixation of deleterious alleles.

We observed 3,316 probes containing SFPs, and assessed associations of SFPs with organismal phenotypes. We found 119, 118, 141, 217, 245 and 195 SFPs associated with starvation stress resistance, chill-coma recovery time, life span, locomotor behavior, mating speed and fitness, respectively ( $P < 0.01$ , Supplementary Table 5). Although some SFPs were associated with more than one trait (Supplementary Table 5), the number of SFPs associated with multiple traits did not exceed that expected by chance. Since we can estimate the frequency of SFPs with significant phenotype associations, as well as the homozygous effect of the SFPs, we can evaluate the distribution of allelic effects. The homozygous effects follow an exponential distribution for all traits, with larger effects associated with the rarer SFPs, and smaller effects with the common SFPs (Fig. 5), as previously predicted<sup>24</sup>.

We used regression models to identify transcripts that were associated with each organismal phenotype. At a  $P$ -value of 0.01, we found 355, 1,128, 295, 231, 691 and 414 transcripts

associated with starvation stress resistance, chill-coma recovery time, life span, locomotor behavior, mating speed and fitness, respectively (Supplementary Table 6). There was little overlap between associations of variation in transcripts and SFPs for the same phenotypes, further increasing the number of candidate genes potentially associated with each trait.

Transcripts that are significantly associated with organismal phenotypes are candidate genes affecting the phenotype<sup>25</sup>. We compared phenotypes of *P*-element insertional mutations in or near candidate genes with that of their co-isogenic control lines<sup>9, 10</sup>. Seven of 10 mutations near candidate genes for resistance to starvation stress indeed affected starvation resistance, and 29 of 39 mutations near candidate genes for chill coma recovery time affected this trait (Fig. 6a,b, Supplementary Table 7). Six of nine mutations in candidate genes affecting locomotor reactivity have been shown previously to affect this trait<sup>26</sup> (Fig. 6c).

### Transcriptional networks associated with complex traits

Most transcripts associated with phenotypes were either unexpected based on prior mutational analyses of the traits, or from computationally predicted genes. To gain insight about functional relationships among transcripts associated with each trait, we used the residuals of the significant regressions of organismal phenotypes on gene expression to quantify modules of transcripts with coordinated patterns of expression across the 40 lines (Fig. 7, Supplementary Table 6). Transcripts with spurious association to a phenotype are unlikely to correlate with biologically relevant transcripts after removing the source of the association; conversely, transcripts under coordinated control are likely to exhibit correlated expression patterns even after removing the effect of their common relationship to a phenotype. Each of the correlated transcript modules associated with a trait can be represented as an interaction network, with edges between transcripts in the network determined by genetic correlations in transcript abundance exceeding a threshold value (Fig. 7). We identified 26 modules of correlated transcripts associated with chill coma recovery time, 20 associated with fitness, 11 with starvation stress resistance, 10 with life span, and 9 each with locomotor reactivity and copulation latency (Fig. 7, Supplementary Table 6).

We evaluated the biological significance of these networks by asking whether genes within each module were enriched for gene ontology categories (Supplementary Table 8), expression in particular tissues, known protein-protein interactions or shared domains. As expected, transcripts associated with variation for fitness are enriched for genes that mediate immune response (Modules 6 and 11), visual perception and function of the nervous system (Module 17), chemosensation (Module 20), and for sex-specific transcripts (Modules 7, 8 and 9) (Fig. 7, Supplementary Table 8). Variation for fitness can be maintained if there are negative genetic correlations between fitness components<sup>1</sup>. Transcripts in Modules 7 and 9, which have female-biased expression, are positively genetically correlated with each other but negatively genetically correlated with transcripts in Module 8, which have male-biased expression. The genes of Module 8 encode proteins that are transferred to females on mating, are thought to benefit male fitness<sup>27, 28</sup>, and that evolve rapidly<sup>29, 30</sup>, but which impose a fitness cost on females<sup>31</sup>. The molecular basis of the female response to this sexual conflict is not known, and could plausibly lie in the Module 7 and 9 transcripts.

Transcripts associated with variation in starvation resistance are enriched for genes that mediate antimicrobial response (Module 4), transcription (Module 6) and proteolysis (Module 9) (Supplementary Table 8). One of the most highly connected genes in Module 6, *raptor*, is a member of the TOR (Target Of Rapamycin) pathway, which plays a key role in nutrient-sensitive signaling, regulates cell growth and cellular mass<sup>32</sup> and regulates the use of alternative energy resources under starvation conditions<sup>33</sup>.

Since gene ontology enrichment analysis<sup>14</sup> revealed similarities between modules of correlated transcripts for the six traits (Supplementary Table 8), we tested whether there was more overlap of common genes between modules for the different traits than expected by chance, and uncovered substantial modular pleiotropy (Fig. 8). For example, genes affecting the mitochondrial ribosome are in common between chill coma recovery Module 17 and copulation latency Module 3 ( $P = 4.74 \times 10^{-4}$ ), chill coma recovery Module 17 and starvation resistance Module 8 ( $P = 1.17 \times 10^{-3}$ ), and starvation resistance Module 8 and copulation latency Module 5 ( $P = 7.53 \times 10^{-4}$ ); while genes affecting defense response to bacteria are in common between starvation resistance Module 4 and fitness Module 11 ( $4.57 \times 10^{-9}$ ) (Fig. 8, Supplementary Table 6). These results give insights to the molecular basis of pleiotropy between complex traits.

## DISCUSSION

Our quantitative genetic analysis of whole genome variation in transcript abundance among a wild-derived population of *Drosophila* inbred lines has revealed surprising features of the genetic architecture of transcription. Nearly 80% of the genome is expressed in adult flies, and approximately 90% of the expressed transcripts have sex-biased expression at a stringent false discovery rate. Two-thirds of the expressed transcripts are genetically variable in this sample of lines – a much greater level of genetic variation than indicated in previous studies<sup>34–43</sup>. Over 40% of the genetically variable transcripts also show genetic variation in sex dimorphism. Remarkably, the whole transcriptome is highly genetically intercorrelated, with 60% of the variable transcripts belonging to two large modules with high positive genetic correlations within modules, and high negative correlations between modules. One of the large modules is enriched for male-biased transcripts and the other for female-biased transcripts. The genetically correlated transcriptional modules are biologically plausible, with enrichments for transcripts in common pathways, gene ontology categories, tissue-specific expression, and transcription factor binding sites. The high transcriptional connectivity at the level of genetic correlation of natural variation in gene expression allows us to infer genetic networks from the transcriptional networks, and the function of computationally predicted genes based on known annotations of other genes in the network.

Several hundred transcripts and SFPs are associated with phenotypic variation in each quantitative trait, and 70% of *P*-element insertional mutations tested in these candidate genes indeed significantly affect the traits. The transcripts associated with each trait group into biologically plausible modules of correlated transcripts, which are in turn correlated between traits, providing insight into the molecular basis of genetic correlations. Variation in transcript abundance in young adults reared under standard culture conditions predicts candidate genes and modules of correlated transcripts associated with variation in stress



responses, behaviors, and life span. The lines and transcript data are therefore a valuable resource for the *Drosophila* community, enabling similar analyses for any complex phenotype that can be quantified. Future integration of whole genome DNA sequence variation with variation in transcript abundance and phenotypes will allow us to disentangle causal from consequential associations, and determine the frequency of causal alleles. Further, the lines can be crossed to interrogate heterozygous effects and degrees of dominance of alleles affecting transcriptional and organismal variation<sup>44</sup>. Knowledge of allele frequencies and homozygous and heterozygous effects will yield unprecedented insight into the nature of evolutionary forces maintaining segregating variation for complex traits in natural populations.

## METHODS

### *Drosophila* lines

We derived inbred lines from the Raleigh, USA population by 20 generations of full-sib mating. We used the *C(2L)RM-P1, b<sup>1</sup>*; *C(2R)RM-SKIA, cn<sup>1</sup> bw<sup>1</sup>* compound autosome (CA) stock for fitness assays. *P*-element mutations and co-isogenic control lines were a gift of Dr. Hugo Bellen. We reared flies on cornmeal-molasses-agar-medium, 25°C, 60–75% relative humidity, 12-hr light-dark cycle unless otherwise specified.

### Gene expression

We used Affymetrix *Drosophila* 2.0 arrays to assess transcript profiles of 3- 5-day old flies from the inbred lines. All samples were frozen between 1 – 3 pm. We extracted RNA from two independent pools (25 flies/sex/line), and hybridized 10µg fragmented cRNA to each array. We randomized RNA extraction, labeling, and array hybridization across all samples, and normalized the raw array data across sexes and lines using a median standardization.

Each transcript is represented by 14 Perfect Match (PM) 25bp oligonucleotides. To identify PM probes with single feature polymorphisms (SFPs) between the wild-derived lines and the strain used to design the array, we quantified the maximal degree to which the variation between lines could be reduced by partitioning the lines into two allelic classes. We computed the sum of squared deviations from each class mean and expressed their sum as a fraction of the total sum of squares. The smallest fraction across all bipartitions was used to score each probe. We identified 3,136 candidate SFPs with scores  $\leq 0.1$  (a tenfold reduction in the sum of squares). We validated polymorphisms in 20/21 of these SFPs by designing primers flanking the SFP and sequencing the PCR products (data not shown).

Our measure of expression for each probe set was the median log<sub>2</sub> signal intensity of PM probes without SFPs. We used negative control probes to estimate the background intensity, and removed probes below this threshold.

### Organismal phenotypes

*Starvation stress resistance* We placed 10 same-sex, two day-old flies in vials containing 1.5% agar and 5ml water, and scored survival every eight hours ( $N = 5$  vials/sex/line). *Chill coma recovery*. We placed three-seven day-old flies in empty vials on ice for three hours,

and recorded the time for each individual to right itself after transfer to room temperature ( $N = 20$  flies/sex/line). *Longevity*. We placed five one-two day-old same-sex virgin flies into vials containing 5 ml medium, and recorded survival every two days ( $N = 5$  vials/sex/line). *Locomotor reactivity*. We placed single three-seven day-old flies into vials containing 5ml medium. The following day, between 8am and 12pm, we mechanically disturbed each fly 19, and recorded the total activity in the 45 seconds immediately following the disturbance. We obtained two replicate measurements of 20 flies/sex/replicate/line. *Copulation latency*. We aspirated pairs of three-seven day-old virgin flies into vials containing 5ml medium between 8am and 12pm, and recorded the number of minutes until initiation of copulation, for a maximum of 120 minutes ( $N = 24$  pairs/line). *Reproductive fitness*. We used the competitive index (CI) technique<sup>45, 46</sup>. We reared all wild type and CA parents in constant density (10 pairs) vials. We placed six three-four day-old virgin CA males and females and three three-four day-old wild type males and females in a vial containing 10ml medium, discarding the flies after seven days. The CI was the ratio of the number of wild type to the total number of progeny emerging by day 17 ( $N = 20$  replicate vials/line).

### Quantitative genetic analyses

We used ANOVA to partition phenotypic variation between sexes ( $S$ , fixed), lines ( $L$ , random), the  $S \times L$  interaction (random) and the error variance ( $\epsilon$ ). We also performed reduced ANOVAs by sex. We estimated broad sense heritabilities ( $H^2$ ) as  $H^2 = (\sigma_L^2 + \sigma_{SL}^2) / (\sigma_L^2 + \sigma_{SL}^2 + \sigma_E^2)$ , where  $\sigma_L^2$ ,  $\sigma_{SL}^2$ , and  $\sigma_E^2$  are the among line, sex  $\times$  line and within line variance components, respectively. For the analyses by sex,  $H^2 = \sigma_L^2 / (\sigma_L^2 + \sigma_E^2)$ . We estimated cross-trait (cross-sex) genetic correlations as  $r_G = cov_{ij} / \sigma_i \sigma_j$ , where  $cov_{ij}$  is the covariance of line means between traits  $i$  and  $j$  (males and females), and  $\sigma_i$  and  $\sigma_j$  are the square roots of the among line variance components for the two traits (males and females).

### Transcriptional modules

To identify modules of genetically correlated transcripts, we computed the correlation  $r_{ij}$  between all pairs of significant transcripts  $i$  and  $j$ . The absolute correlations  $|r_{ij}|$  were transformed to define edge weights  $e^{(|r_{ij}|-1)/\sigma^2}$  in a graph of genes indexed by the free parameter  $\sigma$ . We determined the clustering  $P = \{V_1, \dots, V_k\}$  and value of  $\sigma$  that jointly

maximize the modularity function  $Q(P, \sigma) = \sum_{c=1}^k \left[ \frac{A_\sigma(V_c, V_c)}{A_\sigma(V, V)} - \left( \frac{A_\sigma(V_c, V)}{A_\sigma(V, V)} \right)^2 \right]$ , where  $A_\sigma(X, Y)$  denotes the total edge weight in the graph indexed by  $\sigma$  that connects any vertex in set  $X$  to a vertex in set  $Y$ . The optimal partition  $P = \{V_1, \dots, V_k\}$  defines  $k$  transcriptional modules  $V_1, \dots, V_k$ .

### Transcript-phenotype associations

We used regression models ( $Y = \mu + S + T + S \times T + \epsilon$ , where  $S$  denotes sex and  $T$  the trait covariate) to identify transcripts significantly ( $P < 0.01$ ) associated with organismal phenotypic variation in both sexes. We used the residuals from regression models ( $Y = \mu + E + S + S \times E + \epsilon$ , where  $E$  is the covariate median log<sub>2</sub> expression level) to compute genetic

correlations between transcripts significantly associated with each phenotype for module construction.

### Pleiotropic modules

To identify transcriptional modules common to more than one phenotype, we considered pairs of phenotypes, comparing the lists of significant transcripts for each module from the first phenotype to each module from the second. We used Fisher's Exact Test to quantify the extent that the overlap between the two modules exceeded the chance expectation.

### Transcription factor motifs

We scored 5' UTR sequences of each *D. melanogaster* transcript against position weight matrices for 56 transcription factors; 100 permutations of each sequence were used to generate an empirical distribution of scores for each motif. "Present" motifs had scores within the top 5% of the permutation distribution. We determined the genome-wide proportion of genes for which each motif was present, and compared the proportion of genes within a gene list or module for which a motif was present to the genome-wide proportion using a one-sided binomial test.

### Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

### Acknowledgements

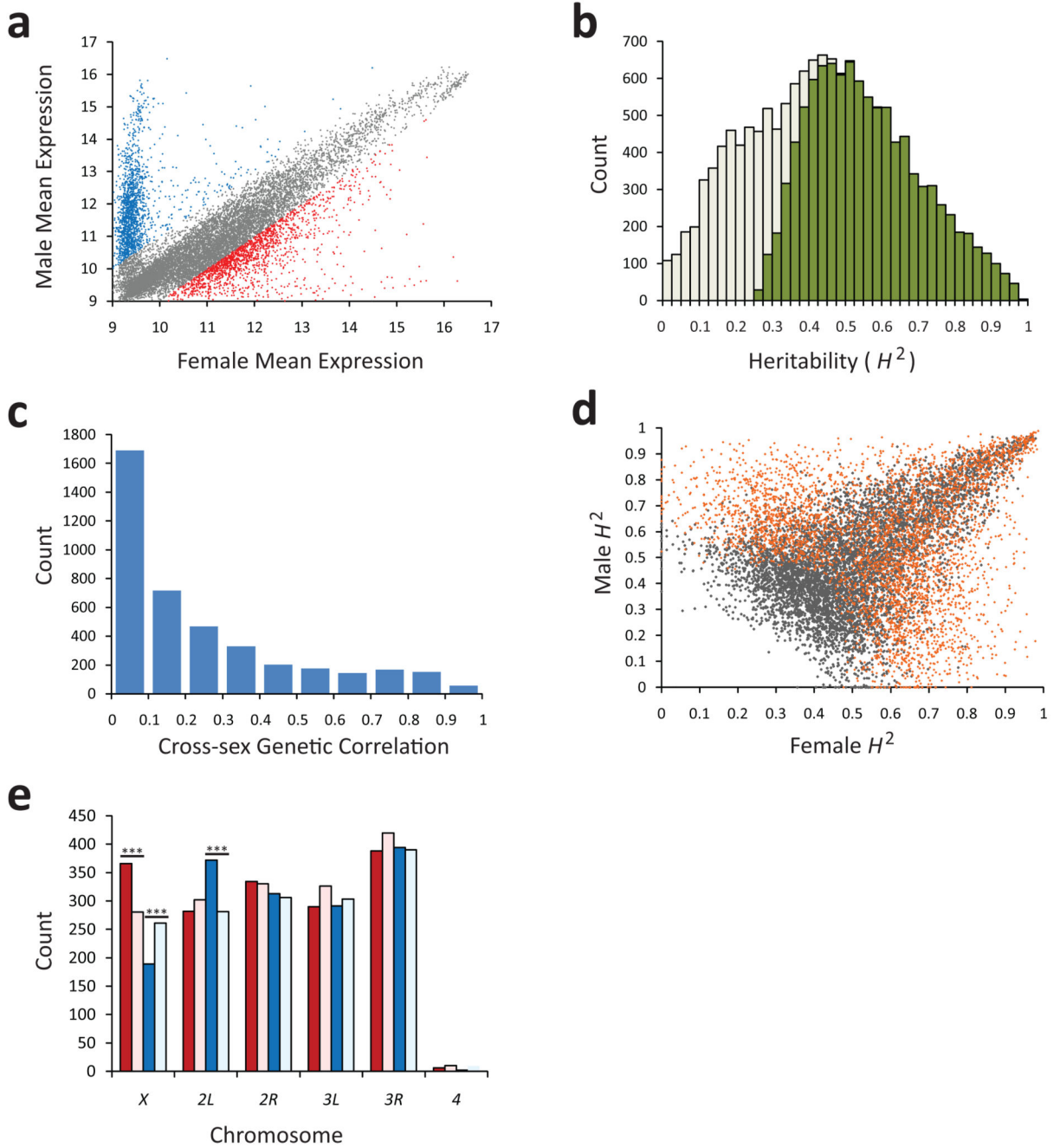
This work was supported by grants from the National Institutes of Health (R01 GM 45146, R01 GM 076083, R01 AA016560 to T. F. C. M. and R01 GM 59469 to R. R. H. A.) The authors thank Stefanie Heinsohn for technical assistance. This is a publication of the W. M. Keck Center for Behavioral Biology.

### REFERENCES

1. Falconer, D.S.; Mackay, TFC. Introduction to Quantitative Genetics. Addison Wesley Longman: Harlow; 1996.
2. Lynch, M.; Walsh, B. Genetics and Analysis of Quantitative Traits. Sunderland, Massachusetts: Sinauer Associates; 1998.
3. Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature*. 2007; 447:661–678. [PubMed: 17554300]
4. Mackay TFC, Anholt RRH. Of flies and man: *Drosophila* as a model for human complex traits. *Ann. Rev. Genomics Hum. Genetics*. 2006; 7:339–367. [PubMed: 16756480]
5. Valdar W, et al. Genetic and environmental effects on complex traits in mice. *Genetics*. 2006; 174:959–984. [PubMed: 16888333]
6. Sieberts SK, Schadt EE. Moving toward a system genetics view of disease. *Mamm. Genome*. 2007; 18:389–401. [PubMed: 17653589]
7. Emilsson V, et al. Genetics of gene expression and its effect on disease. *Nature*. 2008; 452:423–428. [PubMed: 18344981]
8. Chen Y, et al. Variations in DNA elucidate molecular networks that cause disease. *Nature*. 2008; 452:429–435. [PubMed: 18344982]
9. Rollmann SM, et al. Pleiotropic fitness effects of the *Tre1/Gr5a* region in *Drosophila*. *Nat. Genet*. 2006; 38:824–829. [PubMed: 16783380]

10. Sambandan D, Yamamoto A, Fanara JJ, Mackay TFC, Anholt RRH. Dynamic genetic interactions determine odor-guided behavior in *Drosophila melanogaster*. *Genetics*. 2006; 174:1349–1363. [PubMed: 17028343]
11. Storey JD, Tibshirani R. Statistical significance for genomewide studies. *Proc. Natl. Acad. Sci. U.S.A.* 2003; 100:9440–9445. [PubMed: 12883005]
12. Ellegren H, Parsch J. The evolution of sex-biased genes and sex-biased gene expression. *Nat. Rev. Genet.* 2007; 8:689–698. [PubMed: 17680007]
13. Zhang Y, Sturgill D, Parisi M, Kumar S, Oliver B. Constraint and turnover in sex-biased gene expression in the genus *Drosophila*. *Nature*. 2007; 450:233–237. [PubMed: 17994089]
14. Dennis G, et al. DAVID: Database for Annotation, Visualization, and Integrated Discovery. *Genome Biology*. 2003; 4:P3. [PubMed: 12734009]
15. Sturgill D, Zhang Y, Parisi M, Oliver B. Demasculinization of X chromosomes in the *Drosophila* genus. *Nature*. 2007; 450:238–241. [PubMed: 17994090]
16. *Drosophila* 12 Genomes Consortium. Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature*. 2007; 450:203–218. [PubMed: 17994087]
17. Mackay TFC, et al. Genetics and genomics of *Drosophila* mating behavior. *Proc. Natl. Acad. Sci. U. S. A.* 2005; 102:6622–6629. [PubMed: 15851659]
18. Edwards AC, Rollmann SM, Morgan TJ, Mackay TFC. Quantitative genomics of aggressive behavior in *Drosophila melanogaster*. *PLoS Genetics*. 2006; 2:e154. [PubMed: 17044737]
19. Jordan KW, Carbone MA, Yamamoto A, Morgan TJ, Mackay TFC. Quantitative genomics of locomotor behavior in *Drosophila melanogaster*. *Genome Biology*. 2007; 8:R172. [PubMed: 17708775]
20. Morozova TV, Anholt RRH, Mackay TFC. Phenotypic and transcriptional response to selection for alcohol sensitivity in *Drosophila melanogaster*. *Genome Biology*. 2007; 8:R231. [PubMed: 17973985]
21. Chintapalli VR, Wang J, Dow JAT. Using FlyAtlas to identify better *Drosophila* models of human disease. *Nat. Genet.* 2007; 39:715–720. [PubMed: 17534367]
22. Robertson LK, Bowling DB, Mahaffey JP, Imiolyczyk B, Mahaffey JW. An interactive network of zinc-finger proteins contributes to regionalization of the *Drosophila* embryo and establishes the domains of HOM-C protein function. *Development*. 2004; 131:2781–2789. [PubMed: 15142974]
23. Wilson RJ, Goodman JL, Strelets VB. FlyBase: integration and improvements to query tools. *Nucleic Acids Res.* 2008; 36:D588–D593. [PubMed: 18160408]
24. Robertson, A. *Heritage From Mendel*. Brink, A., editor. Madison, Wisconsin; Univ. Wisconsin; 1967. p. 265–280.
25. Passador-Gurgel G, Hsieh WP, Hunt P, Deighton N, Gibson G. Quantitative trait transcripts for nicotine resistance in *Drosophila melanogaster*. *Nat. Genet.* 2007; 39:264–268. [PubMed: 17237783]
26. Yamamoto A, et al. Neurogenetic networks for startle-induced locomotion in *Drosophila melanogaster*. *Proc. Natl. Acad. Sci. U. S. A.* 2008; 105:12393–12398. [PubMed: 18713854]
27. Lung O, Kuo L, Wolfner MF. *Drosophila* males transfer antibacterial proteins from their accessory gland and ejaculatory duct to their mates. *J. Insect. Physiol.* 2001; 47:617–622. [PubMed: 11249950]
28. Wolfner MF. "S.P.E.R.M." (seminal proteins (are) essential reproductive modulators): the view from *Drosophila*. *Soc. Reprod. Fertil. Suppl.* 2007; 65:183–199. [PubMed: 17644962]
29. Date-Ito A, Kasahara K, Sawai H, Chigusa SI. Rapid evolution of the male-specific antibacterial protein andropin gene in *Drosophila*. *J. Mol. Evol.* 2002; 54:665–670. [PubMed: 11965438]
30. Wong A, Turchin MC, Wolfner MF, Aquadro CF. Evidence for positive selection on *Drosophila melanogaster* seminal fluid protease homologs. *Mol. Biol. Evol.* 2008; 25:497–506. [PubMed: 18056920]
31. Wigby S, Chapman T. Sex peptide causes mating costs in female *Drosophila melanogaster*. *Curr Biol.* 2005; 15:316–321. [PubMed: 15723791]
32. Kim D-H, et al. mTOR interacts with Raptor to form a nutrient-sensitive complex that signals to the cell growth machinery. *Cell*. 2002; 110:163–175. [PubMed: 12150925]

33. Kamada Y, et al. Autophagy in yeast: A TOR-mediated response to nutrient starvation. *Curr. Top. Microbiol. Immunol.* 2003; 279:73–84. [PubMed: 14560952]
34. Jin W, et al. The contributions of sex, genotype and age to transcriptional variance in *Drosophila melanogaster*. *Nat. Genet.* 2001; 29:389–395. [PubMed: 11726925]
35. Monks SA, et al. Genetic inheritance of gene expression in human cell lines. *Am. J. Hum. Genet.* 2004; 75:1094–1105. [PubMed: 15514893]
36. Morley M, et al. Genetic analysis of genome-wide variation in human gene expression. *Nature.* 2004; 430:743–747. [PubMed: 15269782]
37. Oleksiak MF, Churchill GA, Crawford DL. Variation in gene expression within and among natural populations. *Nat. Genet.* 2002; 32:261–266. [PubMed: 12219088]
38. Schadt EE, et al. Genetics of gene expression surveyed in maize, mouse and man. *Nature.* 2003; 422:297–302. [PubMed: 12646919]
39. Stranger BE, et al. Genome-wide associations of gene expression variation in humans. *PLoS Genet.* 2005; 1:e78. [PubMed: 16362079]
40. Brem RB, Kruglyak L. The landscape of genetic complexity across 5,700 gene expression traits in yeast. *Proc. Natl. Acad. Sci. U.S.A.* 2005; 102:1572–1577. [PubMed: 15659551]
41. Chesler EJ, et al. Complex trait analysis of gene expression uncovers polygenic and pleiotropic networks that modulate nervous system function. *Nat. Genet.* 2005; 37:233–242. [PubMed: 15711545]
42. Cheung VG, et al. Mapping determinants of human gene expression by regional and genome-wide association. *Nature.* 2005; 437:1365–1369. [PubMed: 16251966]
43. Hubner N, et al. Integrated transcriptional profiling and linkage analysis for identification of genes underlying disease. *Nat. Genet.* 2005; 37:243–253. [PubMed: 15711544]
44. Lemos B, Ararope LO, Fontanillas P, Hartl DL. Dominance and the evolutionary accumulation of *cis*- and *trans*-effects on gene expression. *Proc. Natl. Acad. Sci. U.S.A.* 2008; 105:14471–14476. [PubMed: 18791071]
45. Knight GR, Robertson A. Fitness as a measurable character in *Drosophila*. *Genetics.* 1957; 42:524–530. [PubMed: 17247713]
46. Hartl DL, Jungen H. Estimation of average fitness of populations of *Drosophila melanogaster* and the evolution of fitness in experimental populations. *Evolution.* 1979; 33:371–380.
47. Foronda D, et al. Requirement of *abdominal-A* and *Abdominal-B* in the developing genitalia of *Drosophila* breaks the posterior downregulation rule. *Development.* 2005; 133:117–127. [PubMed: 16319117]
48. DeZazzo J, et al. *nalyot*, a mutation of the *Drosophila myb*-related *Adf1* transcription factor, disrupts synapse formation and olfactory memory. *Neuron.* 2000; 27:145–158. [PubMed: 10939338]



**Figure 1. Variation in transcript abundance among 40 wild-derived inbred lines**

(a) Sex bias for gene expression. Blue and red dots represent genes showing a 2-fold difference in gene expression between males and females, respectively. (b) Distribution of broad-sense heritabilities ( $H^2$ ). Dark green denotes significant  $H^2$  estimates (line FDR < 0.001) and grey indicates non-significant  $H^2$  estimates. (c) Distribution of cross-sex genetic correlations for transcripts exhibiting significant variation in sexual dimorphism (significant sex  $\times$  line interaction variance at FDR < 0.001). (d) Bivariate plot of  $H^2$  estimates in males and females. Orange dots indicate significant line by sex interaction variance. (e)

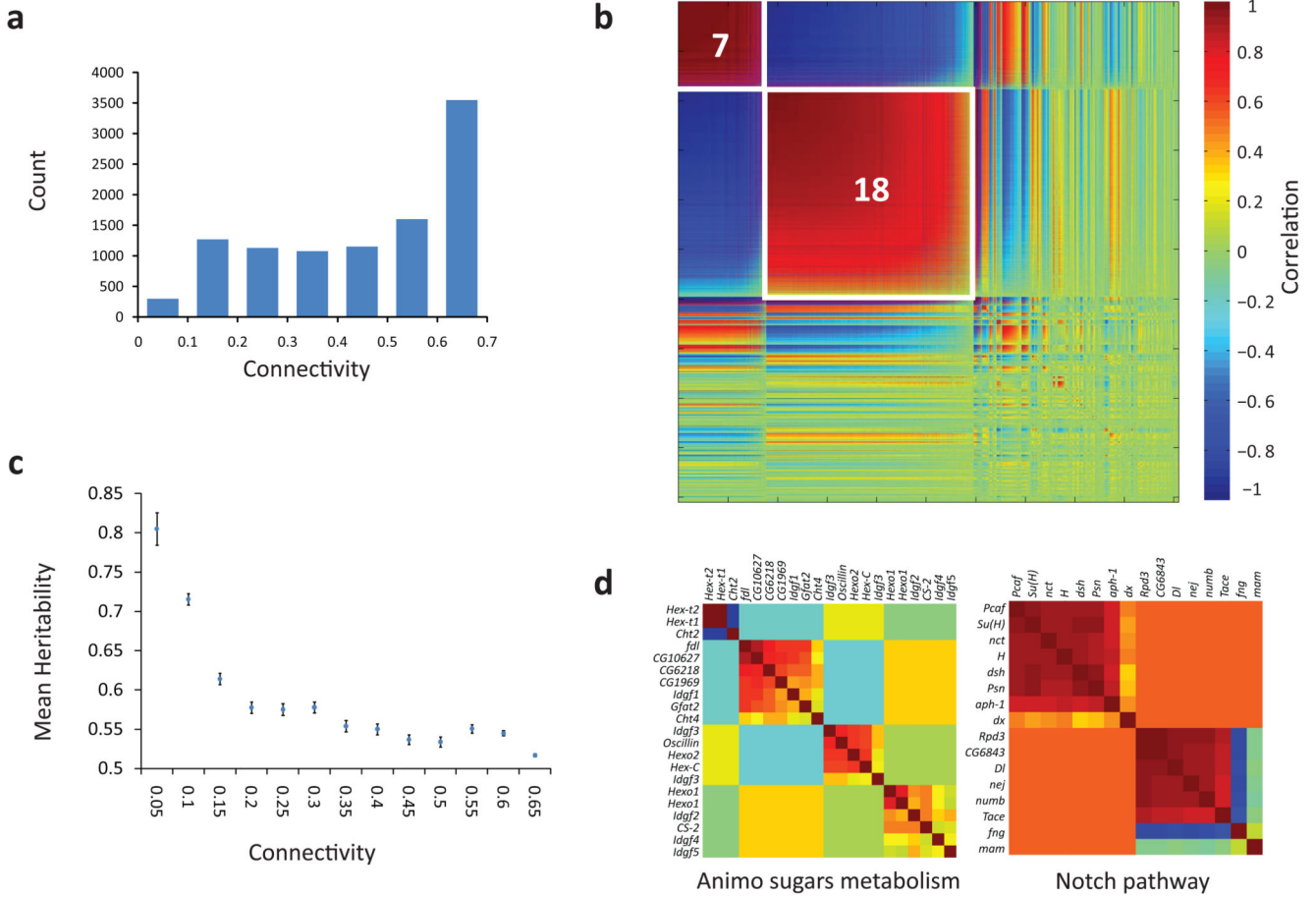
Chromosomal distribution of sex-biased gene expression. The dark blue and red bars are observed male and female counts, respectively, while the light blue and red bars are the expected numbers of male and female transcripts, respectively.

Author Manuscript

Author Manuscript

Author Manuscript

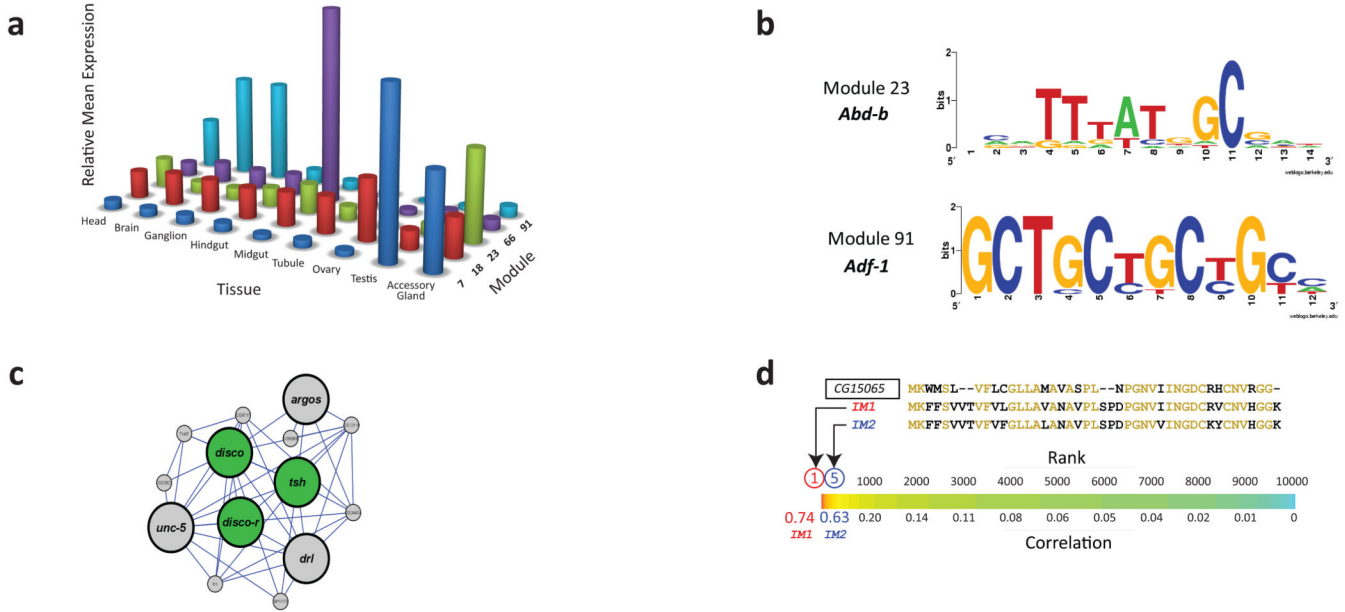
Author Manuscript



**Figure 2. Correlated transcriptional modules**

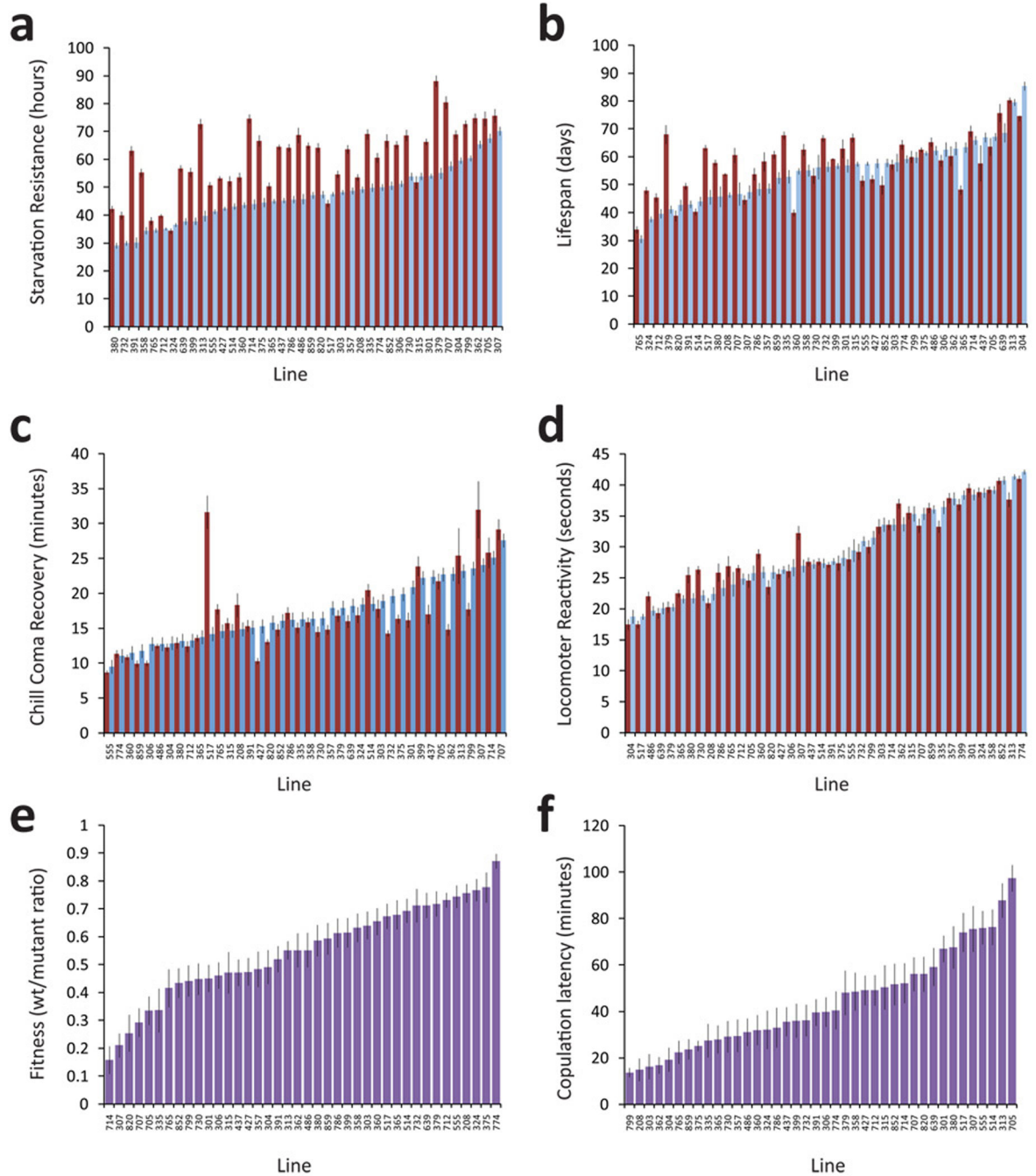
(a) Distribution of connectivity (average  $|r|$ ) for the 10,096 genetically variable transcripts (line FDR < 0.001). (b) Relationship between transcript  $H^2$  and average connectivity. (c) Clustering of the genetically variable transcripts into 241 modules. (d) Correlated transcriptional modules for genes in the amino sugars metabolism and Notch pathway KEGG ontologies. The colors on the off-diagonal represent the average cross-module correlations.





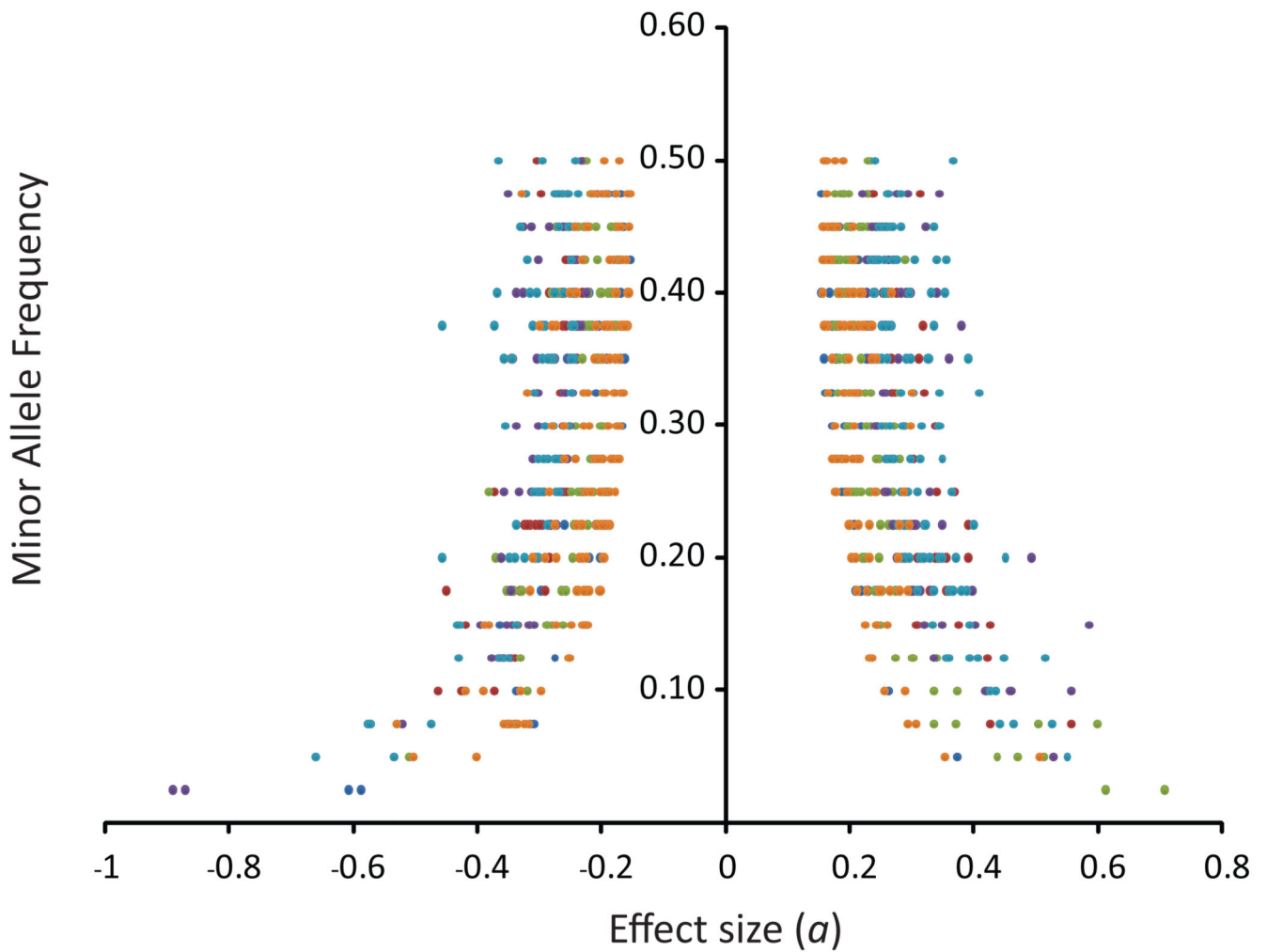
**Figure 3. Biology of transcriptional modules**

(a) Distribution of tissue specific expression in Modules 7, 18, 23, 66, 91. Module 7 is enriched for male-biased transcripts and expression in the testes and accessory glands. Module 18 is enriched for female-biased transcripts and expression in ovaries. Module 23 is enriched for transcripts affecting reproduction and gametogenesis that are highly expressed in ovaries and male accessory glands. Module 66 is enriched for transcripts in the Notch signaling pathway and nervous system development expressed in the midgut. Module 91 is enriched for transcripts affecting the function of the nervous system with high expression in the brain. (b) Modules 23 and 91 are, respectively, enriched for the *Abd-b* ( $P = 0.004$ ) and *Adf-1* ( $P = 0.001$ ) transcription factor binding motifs. *Abd-b* has been implicated in genital disc development<sup>47</sup> and *Adf-1* in memory and synaptogenesis<sup>48</sup>, consistent with the inferred function of genes in these modules. (c) Network representation of module 164, emphasizing the genetic correlations between adult transcripts for three transcription factors that interact during embryonic and larval development. (d) Putative functional annotation of *CG15065* as a gene encoding an Immune induced Molecule (*IM*). Ranking all genetically variable transcripts according to their correlation to *CG15065* shows that *IM1* is the strongest transcriptional correlate ( $r = 0.74$ ) and *IM2* is the fifth strongest ( $r = 0.63$ ). The protein alignments of *CG15065*, *IM1* and *IM2* are highly conserved.



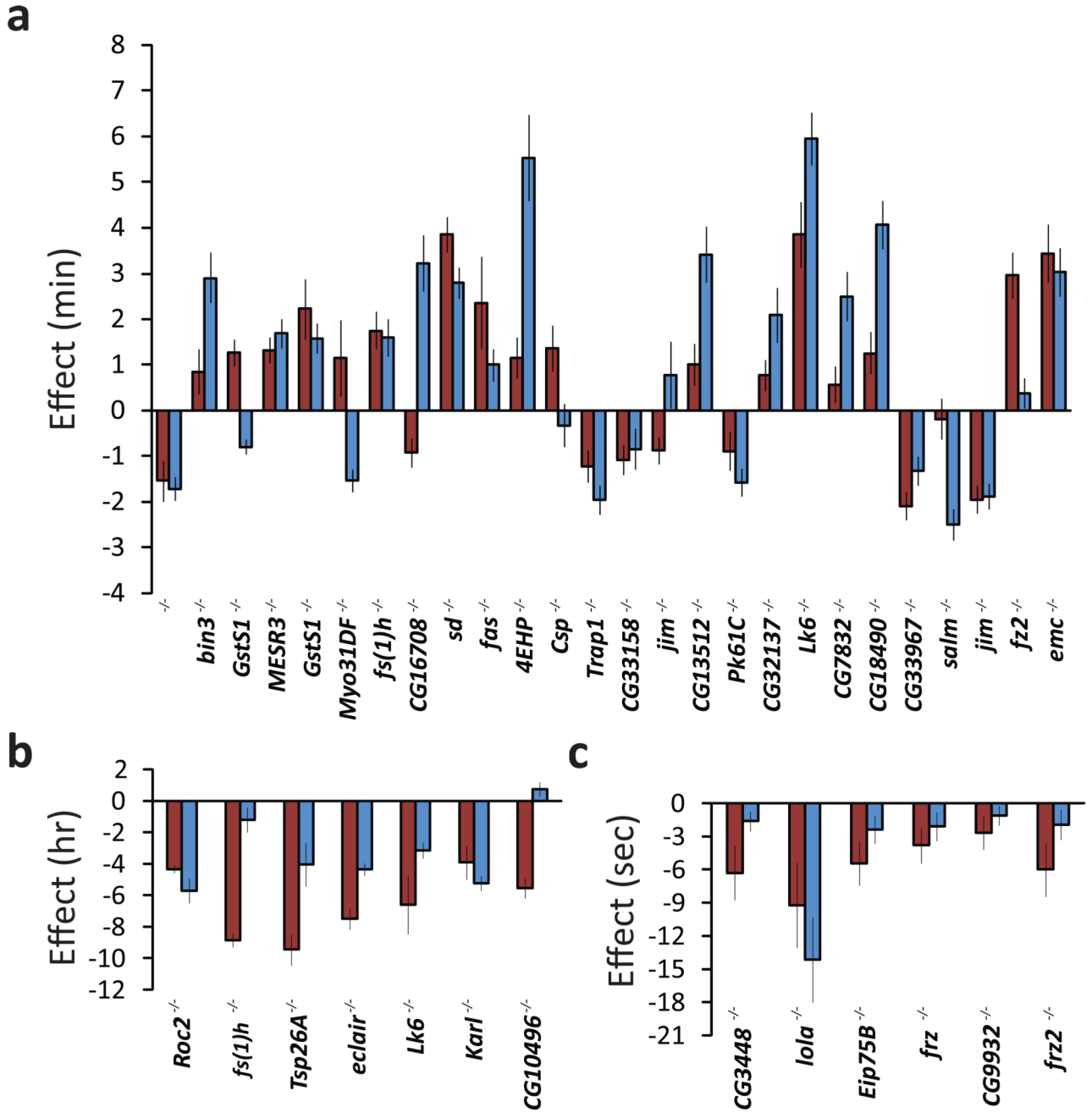
**Figure 4. Variation for organismal phenotypes among 40 wild-derived inbred lines**

Panels **a–f** give the distributions of line means among 40 wild-derived inbred lines. The red and blue bars in panels **a–d** depict females and males, respectively. Sexes were not measured separately in panels **e–f**. Error bars, s. e. (**a**) Starvation stress resistance ( $H^2 = 0.56$ ). (**b**) Chill coma recovery ( $H^2 = 0.23$ ). (**c**) Life span ( $H^2 = 0.54$ ). (**d**) Locomotor reactivity ( $H^2 = 0.58$ ). (**e**) Copulation latency ( $H^2 = 0.25$ ). (**f**) Competitive fitness ( $H^2 = 0.32$ ).

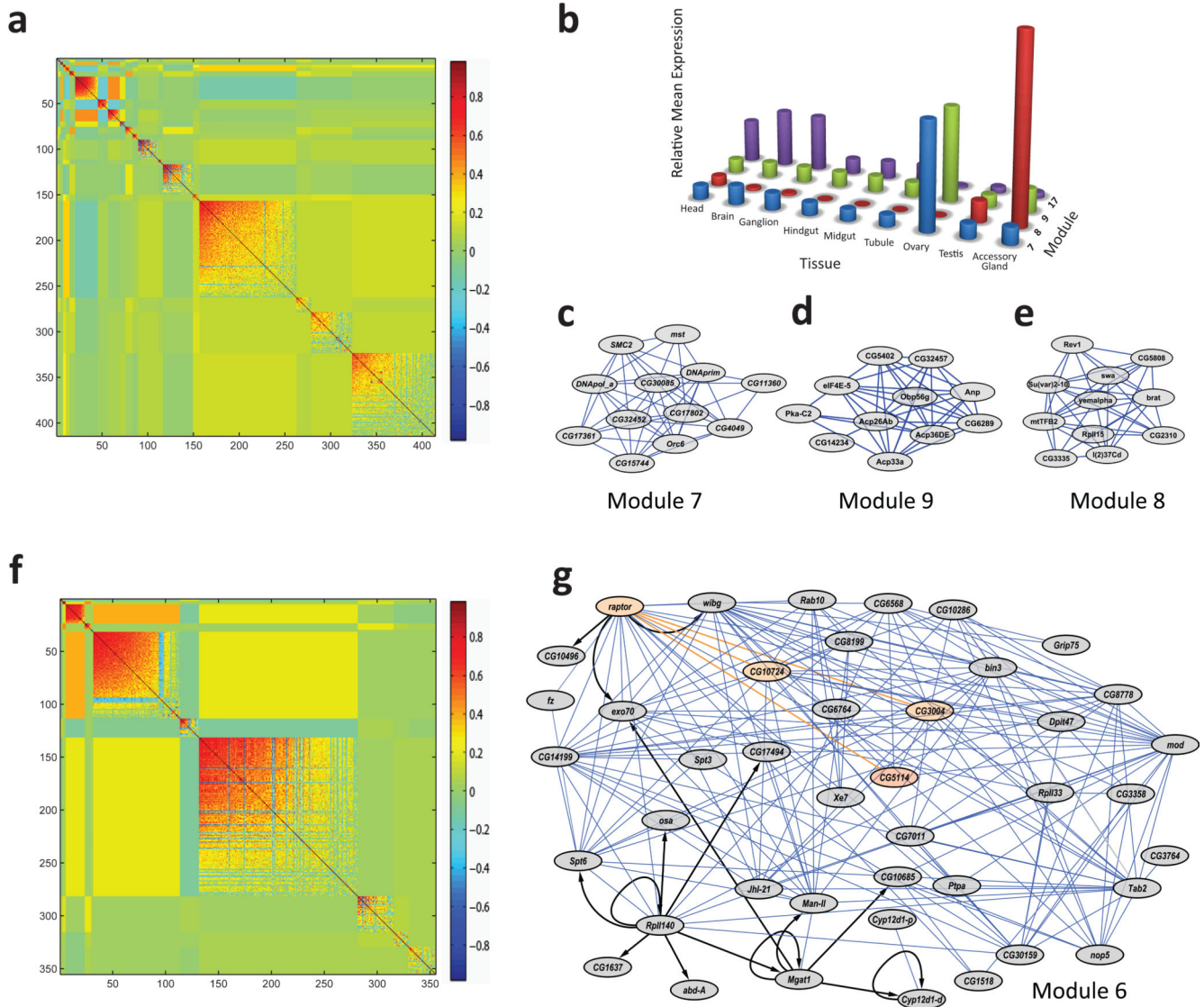


**Figure 5. Distribution of SFP effects**

The  $x$ -axis is the SFP allele effect,  $a/\sigma_G$ , where  $a$  is one half the difference in trait mean between the SFP alleles and  $\sigma_G$  is the genetic standard deviations of each trait. The  $y$ -axis is the minor allele frequency. The traits are color-coded: chill coma recovery (dark blue), starvation resistance (red), fitness (green), lifespan (purple), locomotor reactivity (turquoise), and copulation latency (orange).



**Figure 6. Effects of P-element mutations in candidate genes affecting quantitative traits** Mutational effects are given as deviations from the co-isogenic control line. Red and blue bars represent males and females, respectively. Mutations in all genes shown have significant effects in one or both sexes (Supplementary Table 11). **(a)** Chill coma recovery time. **(b)** Starvation stress resistance. **(c)** Locomotor reactivity (data from Ref. 27).



**Figure 7. Modules of correlated transcripts associated with organismal phenotypes** (a–e) Competitive fitness. (a) Clustering of the 414 transcripts significantly associated with variation in fitness into 20 modules. (b) Tissue-specific expression of transcripts in Modules 7 and 9 (ovaries), Module 8 (accessory glands and testes) and Module 17 (head, brain, and thoracoabdominal ganglion). (c) Interaction network for Module 7. Each node represents a gene and each edge the correlation between a pair of genes. Module 7 is enriched for female-biased transcripts and transcripts affecting DNA replication. (d) Interaction network for Module 9. Module 9 is enriched for female-biased transcripts and transcripts affecting oogenesis and transcriptional regulation. (e) Interaction network for Module 8. Module 8 is dominated by male-biased genes, and is enriched for genes involved in male-induced post-mating behaviors, including three accessory gland proteins (*Acps*). (f–g) Starvation stress resistance. (f) Clustering of the 355 transcripts significantly associated with variation in starvation resistance into 11 modules. (g) Interaction network for Module 6. The black arrows indicate SFP variants in a probe set that are associated with variation in expression of

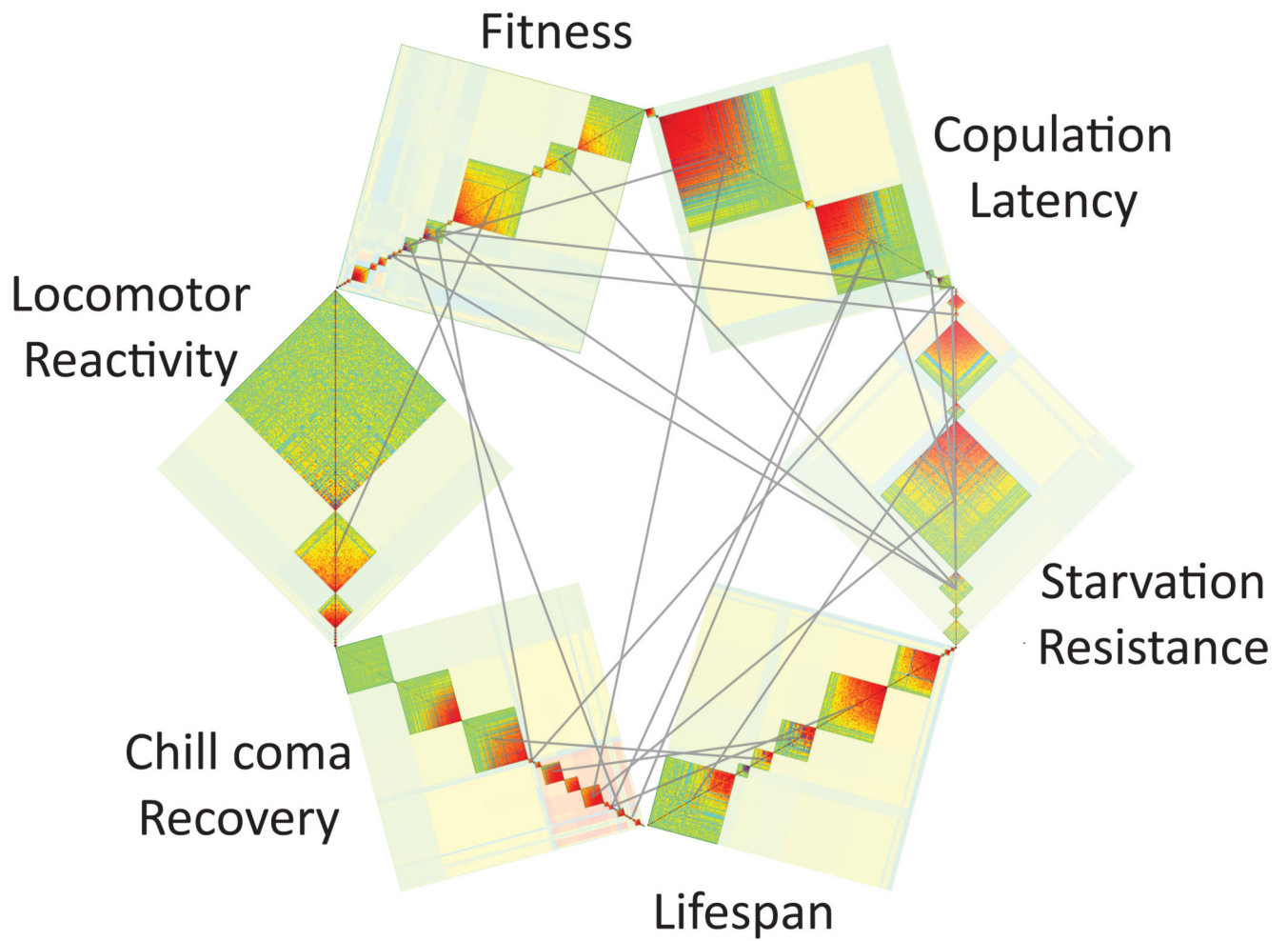
the other probes in that probe set (*cis*-acting variants) and with variation in another transcript (*trans*-acting variants). The orange nodes indicate genes with a WD40 protein domain.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript



**Figure 8. Pleiotropy between phenotypic modules**

Grey lines connect modules with a significant overlap of greater than four genes between gene lists, as determined by Fisher Exact Tests.