



Gene expression networks in the *Drosophila* Genetic Reference Panel

Logan J Everett, Wen Huang, Shanshan Zhou, et al.

Genome Res. published online March 6, 2020

Access the most recent version at doi:[10.1101/gr.257592.119](https://doi.org/10.1101/gr.257592.119)

P<P	Published online March 6, 2020 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

Gene Expression Networks in the Drosophila Genetic Reference Panel

Logan J. Everett^{1,2,*}, Wen Huang^{1,3,*}, Shanshan Zhou^{1,4}, Mary Anna Carbone¹, Richard F. Lyman^{1,5}, Gunjan H. Arya¹, Matthew S. Geisz^{1,6}, Junwu Ma⁷, Fabio Morgante^{1,8}, Genevieve St. Armour¹, Lavanya Turlapati¹, Robert R. H. Anholt^{1,5}, Trudy F. C. Mackay^{1,5,**}

¹ Program in Genetics, W. M. Keck Center for Behavioral Biology and Department of Biological Sciences, North Carolina State University, Raleigh NC 27695-7614; ² Current Address: Environmental Protection Agency, 109 T. W. Alexander Drive, Durham, NC 27709, USA; ³ Current Address: Department of Animal Science, Michigan State University, 474 S Shaw Lane, East Lansing, MI 48824; ⁴ Current Address: Covance, 100 Perimeter Park, Suite C, Morrisville, NC 27560; ⁵ Current Address: Center for Human Genetics and Department of Genetics and Biochemistry, Clemson University, 114 Gregor Mendel Circle, Greenwood, SC 29646; ⁶ University of North Carolina at Chapel Hill School of Medicine, 321 S Columbia St, Chapel Hill, NC 27516; ⁷ Key Laboratory for Animal Biotechnology of Jiangxi Province and the Ministry of Agriculture of China, JiangXi Agricultural University, JiangXi, China; ⁸ Current Address: Section of Genetic Medicine, Department of Medicine, University of Chicago, Chicago, IL 60637

* Equal contributions, alphabetical order

** Corresponding Author

Tel: 919-604-6531

Email: tmackay@clmson.edu

Running title: Gene expression variation in the DGRP

Keywords: RNA sequencing, eQTLs, systems genetics, novel transcribed regions

1 **ABSTRACT**

2

3 A major challenge in modern biology is to understand how naturally occurring variation in DNA
4 sequences affects complex organismal traits through networks of intermediate molecular
5 phenotypes. This question is best addressed in a genetic mapping population in which all
6 molecular polymorphisms are known and for which molecular endophenotypes and complex
7 traits are assessed on the same genotypes. Here, we performed deep RNA sequencing of 200
8 *Drosophila* Genetic Reference Panel inbred lines with complete genome sequences and for
9 which phenotypes of many quantitative traits have been evaluated. We mapped expression
10 quantitative trait loci for annotated genes, novel transcribed regions, transposable elements and
11 microbial species. We identified host variants that affect expression of transposable elements,
12 independent of their copy number, as well as microbiome composition. We constructed sex-
13 specific expression quantitative trait locus regulatory networks. These networks are enriched for
14 novel transcribed regions and target genes in heterochromatin and euchromatic regions of
15 reduced recombination, and genes regulating transposable element expression. This study
16 provides new insights regarding the role of natural genetic variation in regulating gene
17 expression and generates testable hypotheses for future functional analyses.

1 INTRODUCTION

2

3 Understanding how naturally occurring genetic variation affects variation in organismal
4 quantitative traits by modifying underlying molecular networks is a key challenge in modern
5 biology. Most traits are highly polygenic (Visscher et al. 2012; Mackay et al. 2009; Mackay and
6 huang, 2018) and associated molecular variants have small additive effects on trait variation
7 (Manolio et al. 2009). Most of these variants are in intergenic regions, up- or down-stream of
8 coding regions, or in introns, and presumably play a regulatory role in modulating gene
9 expression.

10 Systems genetics analysis seeks to determine how naturally occurring molecular
11 variation gives rise to genetic variation in organismal phenotypes by examining genetic variation
12 in gene expression (expression quantitative trait loci, or eQTLs) and other intermediate
13 molecular phenotypes (Sieberts et al. 2007; Chen et al. 2008; Emilsson et al. 2008; Rockman
14 2008; Cookson et al. 2009; Mackay et al. 2009; Civelek and Lusis 2014; Albert and Kruglyak
15 2015; Gibson et al. 2015; Ogura and Busch 2016; Schughart and Williams 2017). Polymorphic
16 variants associated with variation in gene expression are classified as *cis*- or *trans*-eQTLs
17 depending on whether they are proximal or distal to the gene encoding the transcript,
18 respectively. Genetic variation in gene expression is pervasive; *cis*-eQTLs can have large
19 effects on gene expression that are detectable in small samples; and variants associated with
20 human diseases and quantitative traits tend to be enriched for *cis*-eQTLs (Sieberts et al. 2007;
21 Chen et al. 2008; Emilsson et al. 2008; Rockman 2008; Cookson et al. 2009; Mackay et al.
22 2009; Nicolai et al. 2010; Civelek and Lusis 2014; Albert and Kruglyak 2015; Gibson et al. 2015;
23 Ogura and Busch 2016; Boyle et al. 2017; Schughart and Williams 2017). eQTLs with both *cis*-
24 and *trans*- effects can be assembled into directed transcriptional networks of regulator and
25 target genes (Liu et al. 2008; Bryois et al. 2014; Fagny et al. 2017). Elucidating such regulatory
26 transcriptional networks will facilitate understanding how the effects of individual variants

1 propagate through the network, how multiple variants together regulate gene expression and
2 affect complex traits (Liu et al. 2008; Nicolae et al. 2010; Bryois et al. 2014; Fagny et al. 2017),
3 and improve genomic prediction (Zhou et al. 2020).

4 Here, we performed deep RNA sequencing of the *Drosophila melanogaster* Genetic
5 Reference Panel (DGRP) of inbred lines with complete DNA sequences (Mackay et al. 2012;
6 Huang et al. 2014). We mapped eQTLs for annotated genes, novel transcribed region (NTRs),
7 transposable elements (TEs) and microbiome composition; constructed *de novo cis-trans* eQTL
8 gene expression networks; and evaluated associations of eQTLs and expression traits with
9 organismal phenotypes.

10

11 **RESULTS**

12

13 We collected and sequenced ribo(-) RNA from replicate pools of young flies from each of 200
14 DGRP lines, separately for males and females. In total we sequenced 1.94 Terabases of RNA,
15 of which on average 13.4 million reads per sample uniquely aligned to the *D. melanogaster*
16 genome (Supplemental Table S1). The sequences were processed through a pipeline
17 (Supplemental Fig S1) that (i) removes adapter and rRNA sequences; (ii) aligns and quantifies
18 expressed TE sequences and microbial transcripts; (iii) verifies the origin of each sample; and
19 (iv) quantifies known and novel *D. melanogaster* transcripts and corrects for potential alignment
20 bias due to line-specific sequence variation. We then analyzed normalized expression values for
21 endogenous genes, TEs and microbial species.

22

23 **Genetic Variation in Gene Expression**

24 We quantified expression levels of all RNA sequences that aligned to the reference genome in
25 each DGRP line. After elimination of sequences with low expression, we found that 12,806 of
26 17,097 known *D. melanogaster* genes (75%) were expressed consistently in young adult males

1 and/or females (Supplemental Table S2A). In addition, we identified 4,282 novel transcribed
2 regions (NTRs) (Supplemental Table S2B) that showed no overlap with exons on the same
3 strand. A total of 3,846 of the NTRs were located in introns; 290 were anti-sense to known
4 genes, and 146 were intergenic. Most (95.6%) of the NTRs are ≥ 200 bp; the majority (4,149 or
5 96.9%) lack protein coding potential (Kang et al. 2017) (Supplemental Table S2C). These NTRs
6 in total represent 5.61 Mb new transcribed mature RNA sequences that eluded prior annotation
7 efforts. This increase is likely due to the multiple genetic backgrounds profiled in this study.
8 Although RNA-seq alignment and assembly alone are not sufficient to prove genuine
9 transcriptional activities, our stringent expression-based filter was able to narrow down the
10 NTRs to a subset that were similar to known genes in terms of mapping ambiguity
11 (Supplemental Fig S2) and expression in at least one *Drosophila* cell line (Supplemental Fig
12 S3).

13 Variation in gene expression among the DGRP lines may be confounded by variation in
14 alignment rate to the reference strain due to variation in DNA sequences between the DGRP
15 lines and the reference. Indeed, 2,735 genes (2,117 known genes and 618 NTRs) were affected
16 by alignment bias (Supplemental Table S2D). We corrected for alignment bias, and partitioned
17 variation in gene expression between males and females, DGRP lines, the sex by line
18 interaction, and residual (environmental) terms (Supplemental Table S2D), using a false
19 discovery rate of $FDR \leq 0.05$. Similar to previous studies (Ayroles et al. 2009; Massouras et al.
20 2012; Huang et al. 2015), we found that gene expression is sexually dimorphic: 98% (96%) of
21 expressed known genes (NTRs) have a significant sex effect (Fig 1A, Supplemental Table
22 S2D). There is genetic variation in the magnitude of sex dimorphism: 69% (10%) of expressed
23 known genes (NTRs) have a significant sex by line interaction (Supplemental Table S2D).
24 Therefore, we assessed genetic variation in gene expression separately for males and females
25 (Supplemental Tables S2D, S2E), and found that 12,151 genes (10,354 known genes and
26 1,797 NTRs) were genetically variable in females (Fig 1B) and 13,819 genes (11,393 known

1 genes and 2,426 NTRs) were genetically variable in males (Fig 1C). These numbers of genes
2 with significant genetic variation are much higher than previously reported studies, which used
3 microarrays (4,308 in females and 5,814 in males) rather than RNA-seq (Huang et al. 2015).
4 Relative to tiling arrays, RNA-seq has a higher dynamic range and greater precision in
5 quantifying gene expression, although the results from both analyses are positively correlated
6 (Supplemental Fig S4).

7 Broad sense heritabilities (proportion of phenotypic variance due to genotype
8 differences) ranged from $H^2 = 0.148 - 0.986$ in females and $H^2 = 0.145 - 0.986$ in males (Figs
9 1B, 1C). Notably, 472 (514) of the genetically variable genes in females (number for males in
10 parenthesis) were located in molecularly defined heterochromatin (*2LHet*, *2RHet*, *3LHet*, *3RHet*,
11 *XHet*, and *YHet*) and Chromosome 4. While there are 6.92× (5.52×) as many annotated genes
12 relative to NTRs in euchromatic regions in females (males); there are 2.21 × (3.18×) as many
13 NTRs in heterochromatin and Chromosome 4 in females (males) (Supplemental Table S2G).
14 Thus, NTRs are highly enriched in heterochromatic regions.

15 We used weighted gene co-expression network analysis (WGCNA, Langfelder and
16 Horvath 2008) to assess the extent to which gene expression levels are genetically correlated in
17 each sex (Figs 1D, 1E, Supplemental Table S3). We found 13 (15) co-expression modules in
18 females (males). We assessed the extent to which each module was significantly ($FDR \leq 0.05$)
19 enriched for Gene Ontology (GO) terms and pathway and protein domain annotations (Lyne et
20 al. 2007) (Supplemental Table S3). For example, female Module 2 (149 genes) is enriched for
21 GO terms involved in ovary function and male Module 6 (365 genes) is enriched for biological
22 process GO terms involved in male reproduction. Female Module 12 (88 genes) and male
23 Modules 13 (35 genes) and 14 (165 genes) are enriched for GO terms affecting small molecule
24 metabolism. Female Modules 3 (26 genes), 6 (27 genes), and 7 (21 genes) and male Modules 9
25 (42 genes) and 12 (44 genes) are enriched for GO terms affecting innate immunity, and female
26 Module 13 (560 genes) is enriched for GO terms affecting chemosensation.

1

2 **Gene Expression QTLs (eQTLs)**

3 We performed genome wide association eQTL analyses for each of the genetically variable
4 genes in each sex. We used ~1,932,427 common (minor allele frequency > 0.05)
5 polymorphisms and accounted for effects of Wolbachia infection, polymorphic inversions and
6 polygenic relatedness on gene expression (Huang et al. 2014; Huang et al. 2015). We mapped
7 90,634 eQTLs in females and 147,412 eQTLs in males (FDR \leq 0.05). A total of 2,053 genes in
8 females (1,818 known genes and 235 NTRs) and 3,178 genes in males (2,790 known genes
9 and 388 NTRs) were associated with at least one significant eQTL. We defined potentially *cis*-
10 and *trans*-regulatory eQTLs as \leq 1 kb and $>$ 1 kb of their respective gene bodies. We mapped
11 putative *cis*-eQTLs to 1,284 (2,154) genes in females (males) and *trans*-eQTLs to 1,653 (2,521)
12 genes in females (males), of which 902 (1,305) were *trans*-eQTLs located on different
13 chromosomal arms (Supplemental Tables S4A, S4B).

14 Due to correlation between genotypes at putative eQTL positions, some genes
15 contained a large number of eQTLs that were not independent from each other. To develop a
16 more parsimonious model, we used forward stepwise model selection to select putative eQTLs
17 from the significant candidates, conservatively requiring that the last eQTL entering the model
18 had a conditional *P* value $<$ 1×10^{-5} . The models contained between one and seven eQTLs, with
19 more than 60% of genes containing only one eQTL (Supplemental Tables S4A, S4B). After
20 model selection, we visualized the significant eQTLs by plotting the polymorphism positions on
21 the *X*- axis and the gene positions on the *Y*-axis such that the diagonal corresponds to *cis*-
22 eQTLs and the off-diagonal to *trans*-eQTLs (Fig 2). We found the majority of eQTLs retained by
23 model selection to be in *cis* with the genes they controlled, though *trans*-eQTLs were not
24 uncommon (Fig 2).

25

26 **eQTL Regulatory Networks**

1 The existence of eQTLs that are *cis*-eQTL for gene X and also *trans*-eQTL for gene Y
2 (Supplemental Tables S5A, S5B) enables us to construct gene regulatory networks based on
3 multifactorial variation in a natural population. Although significant putative eQTLs may not
4 remain in the selected models, we still considered them when constructing regulatory networks
5 because we could not genetically distinguish them and their associations with gene expression
6 when all *P*-values were highly significant. We identified 408 (794) such regulatory interactions
7 supported by at least one *cis-trans* eQTL connecting 257 (471) regulatory genes (*cis* end) to
8 251 (447) target genes (*trans* end) in females (males) (Supplemental Tables S5C, S5D). There
9 are two or three large regulatory networks in each sex, and many smaller networks
10 (Supplemental Figs S5, S6, S7). The regulatory genes are largely distinct between the two
11 sexes, although many target genes are in common between males and females (Fig 3,
12 Supplemental Fig S5, Supplemental Table S5E). Genes from the sex-specific regulatory
13 networks or from the common networks are not enriched for any GO terms. It is not clear from
14 their anatomical gene expression patterns how the sex-specificity could arise, since the majority
15 of these genes are expressed in multiple tissues, including the reproductive tissues of both
16 sexes (Gramates et al. 2017).

17 Examination of *cis* and *trans* eQTLs (Supplemental Table S5) showed that there are
18 more NTRs than expected among genes with *cis-trans* eQTLs based on the total number of
19 NTRs with eQTLs among the target genes ($\chi_1^2 = 29.74$, $P = 4.95 \times 10^{-8}$ in females; $\chi_1^2 = 60.54$, P
20 $= 7.20 \times 10^{-15}$ in males) but not the regulatory genes ($\chi_1^2 = 1.54$, $P = 0.21$ in females; $\chi_1^2 = 1.49$,
21 $P = 0.22$ in males). The regulatory genes tend to be located in pericentromeric regions of
22 reduced recombination (Fiston-Lavier et al. 2010) ($\chi_1^2 = 17.28$, $P = 3.23 \times 10^{-5}$ in females; $\chi_1^2 =$
23 120.28 , $P < 2.2 \times 10^{-16}$ in males) and target gene locations are enriched for heterochromatin and
24 pericentromeric regions of reduced recombination ($\chi_1^2 = 28.53$, $P = 9.21 \times 10^{-8}$ in females; $\chi_1^2 =$
25 147.78 , $P < 2.2 \times 10^{-16}$ in males). Regulatory genes with many target genes thus tend to have

1 multiple *cis*-eQTLs in linkage disequilibrium (LD) near the centromere, and regulate other NTRs
2 both in heterochromatic regions across the genome and euchromatic regions on other
3 chromosomes (Fig 3, Supplemental Figs S5, S6, S7). The smaller networks with fewer
4 regulators and targets tend to consist of genes in euchromatin in regions of normal
5 recombination (Fig 3, Supplemental Figs S5, S6, S7; Supplemental Tables S5C, S5D).
6 Regulatory genes often have many *cis*-eQTLs; a single *cis*-eQTL can regulate multiple target
7 genes; and multiple *cis*-eQTLs (which may or may not be in LD) within a gene can regulate
8 different target genes. It is possible that multiple *cis*-eQTLs in LD can be classified as a *trans*-
9 eQTL for different target genes due to differences in thresholding and ranking of eQTLs among
10 the target genes. Each gene with at least one *cis*-eQTL may itself be regulated in *trans* by *cis*-
11 eQTLs in one or more upstream genes, and the genes regulated by a focal *cis*-eQTL may
12 themselves have *cis*-eQTLs regulating other genes.

13

14 **Genetic Variation in TE Expression**

15 A total of 9% of the *D. melanogaster* genome contains TEs spanning multiple families
16 (Spradling and Rubin 1981). Active retrotransposon sequences are present in our RNA-seq
17 libraries. We aligned reads to the RepBase database of known repetitive elements (Jurka et al.
18 2005), and quantified TE RNA levels based on normalized read counts. Overall, 1.3% of the
19 RNA-seq reads align to RepBase. The most abundant families of TE sequences were *gypsy*,
20 *copia*, *BEL*, *jockey* and *Mariner/Tc1* elements, but all TE families represented in RepBase were
21 detected (Fig 4A, Supplemental Table S6A).

22 Line-specific differences in TE RNA levels can be driven by both differences in
23 underlying copy number (Lee and Langley 2010) and differences in the rate of transcription per
24 genomic copy. We quantified DNA copy variation for each TE sequence (Supplemental Table
25 S6B) and used linear models to estimate the percentage of variation in TE expression that
26 arises from differences in copy number (Supplemental Table S6C). We then partitioned the

1 remaining copy number-independent variation in TE expression between sexes, DGRP lines,
2 the line by sex interaction and residual terms (Supplemental Table S6C), using $FDR \leq 0.05$ as
3 the significance threshold for each term in the analysis. Since the majority (153, 79%) of TEs
4 had a significant sex by line interaction effect, we assessed genetic variation in TE expression
5 for each transposon sequence separately for each sex (Supplemental Tables S6D, S6E). We
6 observed significant genetic variation in expression for 187 (97%) TE sequences in females (Fig
7 4B) and 186 (96%) TE sequences in males (Fig 4C). Broad sense heritabilities of TE expression
8 ranged from $H^2 = 0.15 - 0.99$ in females and $H^2 = 0.15 - 0.98$ in males (Figs 4B, 4C). Thus,
9 there is host genetic control of expression for most *D. melanogaster* TEs.

10 We assessed whether different TE sequences had similar patterns of expression across
11 the DGRP lines (Langfelder and Horvath 2008), separately for males and females (Figs 4D, 4E,
12 Supplemental Tables S6F, S6G). We found minimal correlation structure in the activity scores of
13 different TEs (Supplemental Table S6H), with the strongest correlations between pairs of TE
14 sequences from the same family. This suggests that host genetic factors independently affect
15 variation in expression of each TE family.

16

17 **TE eQTLs**

18 We mapped eQTLs for each of the TEs with genetically variable expression in females and
19 males (Supplemental Table S7). We found 54 TEs with significant eQTLs ($FDR \leq 0.05$), 36 in
20 females and 39 in males. A total of 20 TE sequences were expressed in both males and
21 females; 16 (18) TE sequences were expressed only in females (males). The number of eQTLs
22 per TE sequence ranged from 1-1,020, with on average more eQTL associations for TEs in
23 males than females (Supplemental Tables S7A-C). However, forward model selection retained
24 between one and four eQTLs associated with TE activity, suggesting substantial LD among the
25 eQTLs. Indeed, the large numbers of eQTLs associated with some TEs were located in LD

1 blocks in pericentromeric regions and on the 4th chromosome (Supplemental Fig S8,
2 Supplemental Tables S7D, S7E). Many eQTLs for TEs expressed in both males and females
3 overlapped between the sexes, but typically additional eQTLs were present in males. Although
4 there was little clustering of expression patterns of different TE sequences, 202 (1,032) eQTLs
5 were associated with two or more sequences in females (males) (Supplemental Tables S7F,
6 S7G).

7 Many eQTLs associated with TE expression were within 1 kb of annotated genes and
8 NTRs. Indeed, 19.8% (17.7%) of TE eQTLs were within 1 kb of NTRs in females (males).
9 Known genes near TE eQTLs were enriched (FDR < 0.05) for GO categories related to
10 regulation of gene expression and protein binding (Supplemental Table S7H). We next asked to
11 what extent eQTLs associated with gene expression were also associated with expression of
12 TE sequences. We found 1,206 eQTLs associated with 85 genes (37 known genes and 48
13 NTRs) and 23 TEs in females; and 3,656 eQTLs associated with 166 genes (79 known genes
14 and 87 NTRs) and 30 TEs in males (Supplemental Fig S9, Supplemental Table S8). We could
15 thus incorporate variation in TE expression into the *cis-trans* gene regulatory network via shared
16 eQTLs (Fig 5). These eQTLs are predominantly located in pericentromeric regions, and the
17 genes they regulate are in pericentromeric regions as well as heterochromatin.

18

19 **Genetic Variation in Microbiome Composition**

20 RNA samples extracted from pools of whole flies contain RNA from gut microbial communities,
21 and from microbes on their exoskeleton. We assessed the contribution of microbial sequences
22 to the RNA-seq libraries by aligning reads to a database of candidate microbial genomes
23 (Supplemental Table S9). *Wolbachia pipientis*, a bacterial endosymbiont that infects ~50% of
24 the DGRP lines (Mackay et al. 2012), is the most abundant source of expressed sequence,
25 followed by multiple *Acetobacter* species and genome assemblies (Fig 6A, Supplemental Table
26 S9). We estimated the total gene expression from each microbial species in all samples

1 (Supplemental Table S10A) and partitioned variation in microbial gene expression between
2 sexes, DGRP lines, the sex by line interaction and residual terms, using $FDR \leq 0.05$ as the
3 significance threshold (Supplemental Table S10B). The H^2 of *Wolbachia pipientis* abundance is
4 extremely high ($H^2 = 0.972$), as expected. We next assessed whether the sum of all non-
5 *Wolbachia* microbial species is genetically variable after accounting for any *Wolbachia* effects,
6 and estimated $H^2 = 0.595$ (Fig 6B, Supplemental Table S10B). The sex by line interaction for
7 total microbial gene expression was not significant, indicating that total microbial RNA is highly
8 correlated between males and females. We estimated the heritability of gene expression for the
9 122 non-*Wolbachia* microbial species, and found that 84 microbial species had significant
10 genetic variation in RNA abundance, with broad sense heritabilities ranging from $H^2 = 0.07$ –
11 0.90 (Fig 6C, Supplemental Table S10B). Microbial species that are likely to colonize the
12 *Drosophila* gut (*Acetobacter* and *Lactobacillus* species) were among those with the highest H^2 .

13 We used WGCNA (Langfelder and Horwath 2008) to group species with similar
14 abundance patterns based on the average of male and female line means (Fig 6D,
15 Supplemental Tables S10C, S10D). We found three groups of strongly correlated species,
16 consisting primarily of the gut-related microbes (*Acetobacter* and *Lactobacillus* species), and
17 two additional clusters of microbes primarily consisting of viral and fungal species that are
18 strongly anti-correlated with the abundances of species in the first three clusters. Thus, there is
19 line-specific variation in the microbial communities living in and on DGRP flies. Species which
20 most plausibly colonize the *Drosophila* gut are largely correlated across lines, with some
21 fluctuation in the relative abundance of *Acetobacter* versus *Lactobacillus* species.

22

23 **eQTLs for Microbiome Composition**

24 There was little genetic variation in sexual dimorphism for microbial gene expression; therefore,
25 we performed eQTL mapping using the average expression of males and females for each
26 microbial species. Four microbial species and total microbial sequence expression were

1 associated with significant eQTLs ($FDR \leq 0.05$) (Supplemental Table S11A). The sum of all
2 microbial species is associated with one eQTL that maps to an NTR; the expression of *Borrelia*
3 *coriaceae*, *Acidovorax temperans* and *Podospira anserine* map, respectively, to single eQTLs
4 in *CG2616*, *CG46301*, and to *cic* and an NTR; and *Leuconostoc pseudomesenteroides*
5 expression maps to 39 variants in or near *GC* and *nSyb* (Supplemental Table S11A).

6 We lowered the significance threshold to $P < 10^{-5}$ to explore the extent to which common
7 eQTLs may control the expression of multiple microbial species that cluster together based on
8 the WGCNA analysis (Fig 6D). At this threshold, 1,455 eQTLs are associated with 88 microbial
9 species and the sum of all species (Supplemental Table S11B); 268 variants were associated
10 with expression of more than one microbial species, and five eQTLs were associated with
11 expression of 10 or more microbial species (Supplemental Table S11C). These data suggest
12 that there is genetic variation in host control of microbial gene expression and that some
13 variants have pleiotropic effects on multiple microbial species.

14 We assessed whether the genes to which the eQTLs associated with variation in
15 microbial gene expression were enriched for GO categories ($FDR \leq 0.05$). The most highly
16 enriched Biological Process GO terms were related to development and morphogenesis,
17 including development and function of the nervous system (Supplemental Table S11D).

18

19 **Gene Expression and Complex Traits**

20 To examine the relationship between variation in gene expression and variation in organismal
21 quantitative trait phenotypes, we chose 11 quantitative traits with published phenotypic data
22 (chill coma recovery time and startle response (Mackay et al. 2012); starvation resistance
23 (Huang et al. 2014); day and night sleep bout number, day and night total sleep duration, and
24 total waking activity (Harbison et al. 2013); food consumption (Garlapow et al. 2015); male
25 aggression (Shorptter et al. 2015); phototaxis (Carbone et al. 2016)); and additionally measured

1 five metabolic traits (levels of free glucose, glycogen, free glycerol, triglyceride and protein) and
2 three metrics of body size (body weight, thorax length, thorax width). All traits were quantified in
3 the same laboratory under the same culture conditions used in this study. The line means for all
4 traits are given in Supplemental Table S12; quantitative genetic analyses of the metabolic and
5 body size traits are given in Supplemental Table S13; and the most significant associations ($P <$
6 10^{-5}) from GWA analyses (separately for males and females) for these quantitative traits based
7 on the 200 lines for which we have gene expression data are in Supplemental Table S14.

8 We first assessed whether variants associated with all organismal traits were enriched
9 for eQTLs, as found in human studies (Chen et al. 2008; Emilsson et al. 2008; Cookson et al.
10 2009; Nicolae et al. 2010; Boyle et al. 2017). Of all the eQTLs (prior to model selection) and
11 GWAS hits, only 26 in males and 8 in females were common between eQTLs and GWAS hits,
12 with clear patterns of clustering. We found no enrichment of *cis*-eQTLs ($P = 0.13$ in females and
13 $P = 0.71$ in males), *trans*-eQTLs ($P = 0.98$ in females and $P = 0.28$ in males) or all eQTLs ($P =$
14 0.94 in females and $P = 0.23$ in males) among top GWA hits in either sex. Many top GWA hits
15 as well as eQTLs map to regions greater than 1kb from any gene, and may indicate novel
16 regulatory regions. To exclude the possibility that the lack of overlap was due to using different
17 mapping procedures, we performed QTL mapping for the organismal traits using the same
18 procedure as the eQTL mapping. At an empirical FDR = 0.05, we found four SNPs associated
19 with three traits (chill coma recovery in females, day sleep duration and free glucose level in
20 males) and none was an eQTL.

21 We next performed transcriptome wide association studies (TWAS) for individual
22 genetically variable transcripts for gene expression, TE sequences and microbial species, for
23 each of the 18 (19) genetically variable organismal phenotypes in females (males). We found
24 several significant (Benjamini-Hochberg FDR < 0.05) associations of transcripts with organismal
25 phenotypes (Supplemental Table S15). These associations include a known noncoding RNA
26 (*CR46032*) with male aggression, two NTRs with male waking activity, *Gbs-70E* with free

1 glucose in both sexes, *AkhR* with starvation resistance in males and females, and *Acidovorax*
2 *temperans* with male aggression (Supplemental Table S15).

3

4 **DISCUSSION**

5

6 Deep RNA sequencing gives accurate estimates of gene expression of annotated genes and
7 can implicate novel non-coding RNAs and their regulatory interactions with annotated genes.
8 We have identified 4,282 novel transcribed regions, which are unlikely to be artifacts since the
9 majority are genetically variable, and they are not randomly distributed in the genome but are
10 preferentially located in heterochromatic regions and in pericentromeric euchromatin bordering
11 heterochromatin. Thus, there is genetic variation in heterochromatic gene expression, thought to
12 be largely transcriptionally silent (Riddle et al. 2011). These heterochromatic and
13 pericentromeric NTRs are regulated by pericentromeric *cis*-eQTLs as well as *trans*-eQTLs
14 dispersed throughout the euchromatic genome. Genes associated with eQTLs with both *cis*-
15 and *trans*- effects form sex-specific networks of regulator and target genes, the largest of which
16 is enriched for NTR target genes in heterochromatin and regulator and target genes in
17 pericentromeric euchromatin. The considerable overlap between eQTLs associated with NTRs
18 in the large networks and eQTLs associated with TE expression recruits TEs to the network. We
19 do not know where the TE sequences with genetically variable expression are integrated in the
20 genome; however, heterochromatin is composed of largely silenced TE repeats (Riddle et al.
21 2011), raising the possibility that TEs in heterochromatin are subject to the same regulation as
22 other heterochromatic genes. Further work is needed to confirm the regulatory networks derived
23 from naturally occurring genetic variation and determine the regulatory mechanism(s) through
24 which the NTRs act. We speculate that many of the NTRs may be long noncoding RNAs,
25 operationally defined as encoding transcripts > 200 bp with no significant protein-coding

1 potential, but further work is needed to establish whether this is true (Khalil et al. 2009; Wang et
2 al. 2011; Hacısuleyman et al. 2014; Rogoyski et al. 2017; Ransohoff et al. 2018).

3 The first step in systems genetic analysis is to identify eQTLs associated with both gene
4 expression and organismal quantitative traits, for which variation in gene expression is
5 correlated with variation in the organismal phenotypes (Sieberts and Schadt 2007; Rockman
6 2008; Mackay et al. 2009). We did not find any such trios, although we did find interesting
7 transcript-trait associations. This may be because our sample size is adequate to detect eQTLs
8 but not QTLs affecting organismal traits, which have smaller effects; because eQTLs need to be
9 mapped in tissues relevant to the organismal trait; and because there are non-linear (epistatic)
10 relationships between QTLs for both transcripts and organismal phenotypes. The complex and
11 highly connected *cis-trans* regulatory networks suggest that higher order interactions need to be
12 accommodated in systems genetic modeling, at least at the level of gene expression.

13

14 **METHODS**

15

16 **Drosophila lines**

17 We used 200 inbred, sequenced DGRP lines (Mackay et al. 2012; Huang et al. 2014),
18 established by 20 generations of full sib inbreeding from gravid females collected at the Raleigh,
19 NC USA Farmer's Market. Genome sequences of the lines were obtained previously using the
20 Illumina platform with an average of coverage of 27x. A total of 4,565,215 molecular variants
21 (3,976,011 single/multiple nucleotide polymorphisms (SNPs/MNPs), 169,053 polymorphic
22 insertions (relative to the reference genome), 293,363 polymorphic deletions and 125,788
23 polymorphic microsatellites) segregate in the DGRP.

24

25 **Sample collection**

1 All lines were reared on cornmeal-molasses-agar medium at 25°C, 60–75% relative humidity
2 and a 12-hr light-dark cycle at equal larval densities. We collected two replicates of 25 females
3 and 30 males per line, for a total of 800 samples. We used a strict randomized experimental
4 design for sample collection. We collected mated 3-5 day old flies between 1-3 pm. We
5 transferred the flies into empty culture vials and froze them over ice supplemented with liquid
6 nitrogen, and sexed the frozen flies. The samples were transferred to 2.0 ml nuclease-free
7 microcentrifuge tubes (Ambion) and stored at -80°C until ready to process.

8

9 **RNA sequencing**

10 Total RNA was extracted with QIAzol lysis reagent (Qiagen) and the Quick-RNA MiniPrep Zymo
11 Research Kit (Zymo Research). Ribosomal RNA (rRNA) was depleted from 5 ug of total RNA
12 using the Ribo-Zero™ Gold Kit (Illumina, Inc). Depleted mRNA was fragmented and converted
13 to first-strand cDNA using SuperScript III reverse transcriptase (Invitrogen). During the
14 synthesis of second strand cDNA, dUTP instead of dTTP was incorporated to label the second
15 strand cDNA. cDNA from each RNA sample was used to produce barcoded cDNA libraries
16 using NEXTflex™ DNA Barcodes (Bioo Scientific, Inc.) with an Illumina TruSeq compatible
17 protocol. Libraries were size-selected for 250 bp (insert size ~130 bp) using Agencourt Ampure
18 XP Beads (Beckman Coulter, Inc.). Second strand DNA was digested with Uracil-DNA
19 Glycosylase before amplification to produce directional cDNA libraries. Libraries were quantified
20 using Qubit dsDNA HS Kits (Life Technologies, Inc.) and Bioanalyzer (Agilent Technologies,
21 Inc.) to calculate molarity. Libraries were then diluted to equal molarity and re-quantified. A total
22 of 50 pools of 16 libraries were made, again randomly assigning samples to each pool. Pooled
23 library samples were quantified again to calculate final molarity and then denatured and diluted
24 to 14pM. Pooled library samples were clustered on an Illumina cBot; each pool was sequenced
25 on one lane of Illumina HiSeq 2500 using 125 bp single-read v4 chemistry.

26

1 **RNA sequence analysis**

2 Barcoded sequence reads were demultiplexed using the Illumina pipeline v1.9. Adapter
3 sequences were trimmed using cutadapt v1.6 (Martin 2011) and trimmed sequences shorter
4 than 50bp were discarded from further analysis. Trimmed sequences were then aligned to
5 multiple target sequence databases in the following order, using BWA v0.7.10 (MEM algorithm
6 with parameters '-v 2 -t 4') (Li and Durbin 2010): (1) all trimmed sequences were aligned
7 against a database containing the complete 5S, 18S-5p8S-2S-28S, mt:lrrRNA, and mt:srRNA
8 sequences to filter out residual rRNA that escaped depletion during library preparation; (2)
9 remaining sequences were then aligned against a custom database of potential microbiome
10 component species (see below) using BWA; (3) sequences that did not align to either the rRNA
11 or microbiome databases were aligned to all *D. melanogaster* sequences in Repbase (Jurka et
12 al. 2005). The remaining sequences that did not align to any of the databases above were then
13 aligned to the *D. melanogaster* genome (BDGP5) and known transcriptome (FlyBase v5.57)
14 using STAR v2.4.0e (Dobin et al. 2013). Libraries with fewer than 5 million reads uniquely
15 aligned to the *D. melanogaster* reference genome were re-sequenced to achieve sufficient read
16 depth.

17

18 **Generation of microbiome database**

19 We first performed a preliminary alignment of RNA-seq reads by filtering only rRNA sequences,
20 and then aligning directly to the *D. melanogaster* genome using the tools and parameters
21 described above. Sequences that did not align to the rRNA database or *D. melanogaster*
22 reference genome were then analyzed with Trinity v2.1.1 (Garbherr et al. 2011) to perform *de*
23 *novo* assembly of longer sequences from the short reads. Assembled sequences > 1kb in
24 length were then searched against the RefSeq_genomic database (downloaded from NCBI on
25 1/27/16) using BLAST. We then compiled a list of all RefSeq genomes that were found as a top

1 BLAST hit for at least two assembled sequences. We compiled all FASTA files for each of these
2 RefSeq genomes into a single database for alignment with BWA.

3

4 **Genotype validation**

5 To validate the DGRP line assigned to each RNA-seq sample, we identified single nucleotide
6 polymorphisms (SNPs) from the RNA-seq reads that aligned to the *D. melanogaster* reference
7 genome using STAR as described above. We retained only those SNP calls covered by at least
8 3 reads and at least 75% of all reads supporting the major genotype (note that DGRP lines are
9 inbred and therefore the majority of SNPs are homozygous). This filtering process produced
10 >400k usable SNPs per sample, primarily located in transcribed regions of the genome. We
11 then performed two validation tests of the DGRP line assigned to each sample X by comparing
12 to the previously published genotype calls for each DGRP line
13 (<http://dgrp2.gnets.ncsu.edu/data/website/dgrp2.tgeno>; Huang et al. 2014). First, we computed
14 the “line mean error” (LME) for each line as follows: given the set of homozygous SNPs from
15 line Y that have sufficient coverage (described above) in sample X , $LME(X, Y) = \#$ of
16 mismatching SNPs / total # of comparable SNPs. We confirmed that for each sample X , the
17 DGRP line Y_{lab} labeled for that samples produced the minimum value of $LME(X, Y)$ as compared
18 to all other possible line assignments Y_{alt} , and further confirmed that $LME(X, Y_{lab})$ was below 1%.
19 Second, we performed competitive tests between the labeled line Y_{lab} and each possible
20 alternate line Y_{alt} . Under this test, we considered only the SNPs that are homozygous for
21 different genotypes in Y_{lab} and Y_{alt} (*i.e.*, only the segregating SNPs for the two lines) and which
22 have sufficient coverage in sample X . We then computed the “line error ratio” (LER) = # of
23 SNPs matching Y_{lab} / # of SNPs matching Y_{alt} . We confirmed that for each sample X , the lowest
24 LER for any Y_{alt} was > 1 (*i.e.*, the majority of SNP calls always supported the labeled line
25 compared to any alternative line).

26

1 **Inference of novel transcripts**

2 We constructed a *de novo* transcriptome for each individual sample by inputting the RNA-seq
3 reads aligned to the *D. melanogaster* reference genome into Cufflinks v2.2.1 (Trapnell et al.
4 2012). We also considered the novel transcribed regions (NTRs) identified in a previous study
5 based on unstranded pooled RNA sequencing of the DGRP lines (Huang et al. 2015). However,
6 the previously published data do not provide strand-specific signal, while our current RNA-seq
7 data uses a strand-specific library preparation. Therefore, we reassigned the strand for each of
8 the previously published NTRs that was supported by the greater number of total aligned reads
9 across all samples. We then merged all *de novo* sample transcriptomes and the previously
10 published NTRs using the cuffmerge tool included with Cufflinks v2.2.1, then removed all
11 merged transcript models with any exon overlapping on the same strand any exon in the known
12 *D. melanogaster* transcriptome. We defined the known transcriptome here as all gene models in
13 FlyBase v5.57 plus all subsequently added gene models in FlyBase v6.11 to account for
14 recently discovered lncRNA sequences. Thus, the final output of this analysis was a set of
15 NTRs constructed from both our current RNA-seq data and previously published pooled RNA-
16 seq data that do not overlap any known gene exons on the same strand.

17

18 **Gene expression estimation**

19 Read counts for individual microbial species were computed as all reads aligning to any
20 sequence in any genome for any strain of that species. Reads aligning to multiple species were
21 ignored for individual species read counts. We also aligned microbiome-aligning reads to the *D.*
22 *melanogaster* genome, and removed all reads that aligned to both microbial and *D.*
23 *melanogaster* sequences before computing read counts, to account for several domains which
24 are highly conserved between microbial and metazoan species. Read counts were computed
25 for transposon sequences by computing the number of reads uniquely aligned to each
26 transposon sequence in Repbase. Highly homologous sequences were grouped together for

1 computing transposon read counts. Read counts were computed for known and novel gene
2 models using HTSeq-count (Anders et al. 2015) with the ‘intersection-nonempty’ assignment
3 method for exonic reads only. Tabulated read counts for each expression feature type
4 (microbiome, transposon, endogenous genes) were then normalized across all samples using
5 edgeR (Robinson et al. 2010) as follows. First, genes with low expression overall (<10 aligned
6 reads in >75% of the libraries) were excluded from the analysis. Library sizes were re-computed
7 as the sum of reads assigned to the remaining genes, and further normalized using the
8 Trimmed Mean of M-values (TMM) method (Robinson and Oshlack 2010). At this point, we
9 retained only genes (known or novel) whose expression in both biological replicates was above
10 an empirical threshold in more than 200 line-sex combinations (400 samples total). This criterion
11 retains genes expressed in only one sex. The threshold was determined by fitting all \log_2
12 transformed FPKM expression data points using a 2-component Gaussian mixture model and
13 finding the expression value (FPKM = 0.280263) where the posterior probability of being in the
14 lower expression component is 0.95. Genes on Chr U and Chr Uextra were also removed. We
15 further adjusted transposon expression estimates to account for differences in transposon copy
16 number across lines by fitting a linear model: $\text{RNA} \sim \text{DNA} + \varepsilon$, where RNA = the normalized \log_2
17 (RNA-seq read count); and DNA = normalized \log_2 (DNA read count) derived from the
18 previously published DNA-seq data for each DGRP line (Huang et al. 2014). After fitting the
19 linear model for each transposon sequence, ε estimates the relative transcription rate in each
20 line independent of copy number, and was used as the adjusted transposon expression for all
21 subsequent analysis. We further adjusted endogenous gene expression values by estimating
22 and removing the effect of alignment bias resulting from higher rates of non-reference variants
23 clustering in some lines. We computed the alignment bias score $A(g,L)$ defined as the number
24 of non-reference nucleotides per kb in all exons of gene g in DGRP line L , based on the
25 previous map of genomic variation in the DGRP (Huang et al. 2014). We then fit a linear model
26 for each endogenous gene: $Y = A + \varepsilon$, where Y is the normalized expression profile for gene g

1 after the read counting and edgeR normalization described above. After fitting these linear
2 models, ε represents the alignment bias-corrected expression, and was used as the normalized
3 gene expression in all subsequent analysis. Read mapping ambiguity could affect the confidence
4 in defining NTRs. We assessed this by using blat (<https://genome.ucsc.edu/FAQ/FAQblat.html>)
5 to map all RNA transcripts to the fly genome and identified all possible alignments. We used a
6 metric (Δ Bitscore) to characterize the mapping ambiguity of the full length RNA transcripts for
7 known genes, NTRs filtering for low expression across the DGRP, and NTRs retained after
8 filtering. Δ Bitscore is defined as the difference between the bit score for the best alignment and
9 the second best one. The greater the Δ Bitscore, the less ambiguous is the alignment. In
10 addition, we assessed whether the NTRs identified in the DGRP were present in an
11 independent data set of 41 *Drosophila* cell lines that were either untreated or treated with the
12 hormone ecdysone (Stoiber et al. 2016). We computed the median and maximum expression
13 across all cell lines for each transcript using kallisto (Bray et al. 2016), an alignment-free
14 abundance estimator, and calculated the median and maximum expression RNA transcripts for
15 known genes, NTRs filtered for low expression across the DGRP, and NTRs retained after
16 filtering.

17

18 **Genetics of gene expression**

19

20 For each class of expression features (endogenous genes, transposons, microbiome), we fit
21 mixed-effect models to the gene expression data corresponding to: $Y = S + W + W \times S + L + L \times S$
22 $+ \varepsilon$, where Y is the observed \log_2 (normalized read count), S is sex, W is Wobachia infection
23 status, $W \times S$ is Wolbachia by sex interaction, L is DGRP line, $L \times S$ is the line by sex interaction
24 and ε is the residual error. We also performed reduced analyses ($Y = W + L + \varepsilon$) independently
25 for males and females. We identified genetically variable transcripts as those that passed a 5%

1 FDR threshold (based on Benjamini-Hochberg (1995) corrected P -values) for the L and/or $L \times S$
2 terms. We computed the broad sense heritabilities (H^2) for each gene expression trait
3 separately for males and females as $H^2 = \sigma_L^2 / (\sigma_L^2 + \sigma_\varepsilon^2)$, where σ_L^2 and σ_ε^2 are, respectively, the
4 among line and within line variance components.

6 **Clustering by genetic correlation**

7 For each feature type (microbiome, transposons, endogenous genes), we clustered line means
8 using the WGCNA R package v1.51 (Langfelder and Horvath 2008) as follows. Only genes with
9 sufficient average expression (\log_2 FPKM > 0) and genetic variance (line mean variance > 0.05)
10 were considered in these analyses. First, the Pearson correlation coefficient for every pair of
11 line means, the soft-power threshold was computed using the pickSoftThreshold function, and
12 used to convert the correlation matrix to an adjacency matrix with approximately scale-free
13 connectivity. The adjacency matrix was then converted to a dissimilarity matrix based on the
14 topological overlap map (Langfelder and Horvath 2008). Expression features were then
15 clustered using hierarchical clustering (hclust function) based on the dissimilarity matrix, and
16 split into distinct modules using the cutreeDynamic with deepSplit=4 and minClusterSize=20 (for
17 endogenous gene expression, minClusterSize=4 was used for microbiome and transposon
18 clustering). Module eigengenes were computed for each cluster, and highly similar clusters
19 were combined using the mergeCloseModules function with cutHeight = 0.25. Expression
20 features assigned to module 0 (insufficient similarity) were discarded. Modules consisting of
21 >1,000 features were also discarded as insufficiently split into distinct modules. For each
22 expression feature, the degree was computed as the overage topological overlap with all other
23 features assigned to the same module. The average degree of each module was computed as
24 the mean degree across all features in the module. Modules were sorted by average degree,
25 such that module 1 has the highest average degree in each analysis.

26

1 **Gene set enrichment analyses**

2 Lists of known gene IDs (FlyBase FBgn accessions) were uploaded to FlyMine (Lyne et al.
3 2007) or Panther (Mi et al. 2017) for functional enrichment. For analysis of gene lists from
4 WGCNA clusters, the list of known genes input to WGCNA was used as the background set, to
5 correct for any biases inherent to highly heritable expression patterns in general.

7 **Expression QTL (eQTL) mapping**

8 For each gene expression feature, we performed eQTL analysis as previously described
9 (Huang et al. 2015). Briefly, we adjusted mean expression values in each sex for fixed effects of
10 Wolbachia infection status, five major polymorphic inversions (*In2L(t)*, *In2R(NS)*, *In3R(P)*,
11 *In3R(K)*, *In3R(Mo)*), and the first 10 principal components of the genetic relatedness matrix of all
12 DGRP lines using a linear model. We mapped QTLs for the adjusted line means using PLINK
13 (Purcell et al. 2007) against 1,932,427 SNPs with major allele frequency > 0.05 and missing
14 genotypes in fewer than 25% of the 200 DGRP lines profiled by RNA-seq. Instead of controlling
15 for experiment-wise type I error rate, which can be overly conservative, we controlled for the
16 false discovery rate (FDR, Benjamini and Hochberg 1995). We computed FDR of eQTL calls by
17 comparing observed eQTL *P*-value distributions to those obtained from running PLINK on 100
18 permutations of the observed line means for each expression feature. At any given *P*-value cut-
19 off *X*, the estimated false positive rate of eQTLs for a specific gene expression feature is the
20 average number of eQTLs with *P*-value < *X* across all permutations. The FDR at the same *P*-
21 value is then computed as the estimated false positive rate divided by the number of eQTLs with
22 *P*-value < *X* in the observed data. Using this formulation of FDR, we identified the unadjusted *P*-
23 value cut-off corresponding to 5% FDR for each gene expression feature. No further model
24 selection was performed; however, we classified eQTLs as being either *cis*-eQTLs (within 1kb of
25 the gene body for the associated expression feature) or *trans*-eQTLs (> 1 kb of the gene body).
26 To eliminate eQTLs whose genotypes are correlated with each other and cannot be genetically

1 distinguished, we used forward model selection to iteratively add eQTLs to the model in the
2 order of their conditional association (Huang et al. 2015). The model selection was stopped
3 when none of the remaining putative eQTLs can enter the model with $P < 0.00001$. When two
4 putative eQTLs had equal P -values, the one closer to the transcription start site was added.

5

6 **Construction of eQTL networks**

7 We then constructed regulatory eQTL networks based on individual SNPs which were called as
8 both *cis*- and *trans*-eQTLs for multiple expression features. Specifically, we assign a directed
9 edge $X \rightarrow Y$ if there is at least one variant that is both a *cis*-eQTL for gene X (defined as within
10 1 kb of gene X) and a *trans*-eQTL for gene Y at 5% FDR. We then broke all loops in the
11 regulatory network for each sex by dropping the edge in each loop with the highest minimum P -
12 value from all associated SNPs to create a directed, acyclic network.

13

14 **Quantitative traits**

15 We retrieved phenotypic data documented from previous publications on the same fly lines for
16 male aggression (Shorter et al. 2015); chill coma recovery time and startle response (Mackay et
17 al. 2012); food consumption (Garlapow et al. 2015); phototaxis (Carbone et al. 2016); sleep
18 traits (Harbison et al. 2013) (day and night bout number, day and night total sleep duration, total
19 waking activity); and starvation resistance (Huang et al. 2014).

20 To measure body weight and size, we collected 10 replicates of 10 flies per line and sex
21 into pre-weighed 1.7 ml tubes, and weighed and flash froze them for downstream analyses.
22 Virgin flies were used to avoid body weight variation due to variation in egg production. In
23 addition we measured thorax length and thorax width as metrics for body size.

24 Frozen flies were homogenized in 250 μ L Dulbecco's phosphate-buffered saline, and
25 after gentle centrifugation supernatants were collected for measurements of free glucose,
26 glycogen, free glycerol, triglyceride and total protein (further diluted 10 fold). For free glucose

1 and glycogen, samples were denatured at 95°C for 25 minutes to prevent glycogenolysis.
2 Measurements were performed following protocols provided by the Glycogen
3 Colorimetric/Fluorometric Assay Kit (BioVision Inc.). For free glycerol and triglyceride, we used
4 the Serum Triglyceride Determination Kit (Sigma-Aldrich Inc.), and incubated samples with the
5 Triglyceride Reagent for 1 hour at 37°C. For total protein measurement, we used the Qubit
6 Protein Assay Kit (Thermo Fisher Scientific Inc.).

7

8 **Quantitative trait genetic parameters**

9 We used mixed model, factorial ANOVAs ($Y = S + L + L \times S + Rep(L) + S \times Rep(L) + \varepsilon$, to partition
10 variation of the quantitative traits between sexes (S), DGRP lines (L) and replicate vials within
11 lines (Rep). Broad sense heritabilities were estimated as $H^2 = (\sigma_L^2 + \sigma_{SL}^2) / (\sigma_L^2 + \sigma_{SL}^2 + \sigma_\varepsilon^2)$,
12 where σ_L^2 , σ_{SL}^2 and σ_ε^2 are, respectively, the among line, sex by line and within line variance
13 components.

14

15 **eQTL-GWA enrichment analysis**

16 We performed GWA analyses for all quantitative traits, separately for females and males. All
17 phenotypes (line means) were first adjusted for the effect of Wolbachia infection and major
18 polymorphic inversions using a linear model. The residuals (plus the intercept) from this analysis
19 were then used as phenotype in a linear mixed model to test for the effect of each common
20 variant individually, while adjusting for sample structure using a genomic relationship matrix
21 (GRM), as implemented in GCTA-MLMA (Yang et al. 2011). The GRM was built as $\frac{WW^T}{p}$ where
22 W is a matrix of centered and scaled genotypes for the 200 lines and p is the total number of
23 genetic variants. Similarly, we have also mapped trait QTLs using the same procedure as the
24 eQTL mapping described above, by deriving empirical FDR based on 100 permutations of
25 phenotypes.

1 For each trait and sex, variants with $P < 10^{-5}$ were retained for downstream analysis. We
 2 then combined the lists of variants associated with each trait, separately for females and males,
 3 to obtain a single list of unique variants (*i.e.*, no duplicates) associated with any of the traits of
 4 interest. The enrichment analysis proceeded as described in Nicolae et al. (2010), within each
 5 sex. Briefly, GWAS hits were divided into minor allele frequency bins of width equal to 0.05.
 6 Then, an equal number of common variants (which may or may not have included actual GWAS
 7 hits) per bin were sampled at random and the overlap with eQTLs was calculated. This
 8 procedure was repeated 10,000 times and an empirical P -value for the enrichment was
 9 calculated as the number of replicates where the overlap between randomly sampled variants
 10 and eQTLs was greater than or equal to the observed overlap between GWAS hits and eQTLs
 11 over the total number of replicates.

12

13 **Association of expression and quantitative traits**

14 A transcriptome-wide association study (TWAS), *i.e.*, regressing the phenotype on each gene's
 15 expression level, was performed for each sex separately for each quantitative trait. We
 16 developed a method that accounts for structure present in the transcriptome due correlations
 17 between transcripts. This was achieved by fitting a linear mixed model of the type: $\mathbf{y} = \mathbf{1}\mu +$
 18 $\mathbf{w}\beta + \mathbf{t} + \mathbf{e}$, where \mathbf{y} = n -vector of mean phenotypic values for n lines, μ = fixed population
 19 mean effect, \mathbf{w} = n -vector of the tested gene's centered and scaled expression level, β = fixed
 20 effect of the gene, \mathbf{t} = n -vector of random transcriptomic line effect ($\mathbf{t} \sim N(0, \mathbf{T}\sigma_t^2)$), and \mathbf{e} = n -
 21 vector of random error ($\mathbf{e} \sim N(0, \mathbf{I}\sigma_e^2)$).

22 The key term in the model that accounts for sample structure is \mathbf{T} , the transcriptomic
 23 relationship matrix (TRM). The TRM was computed as $\frac{\mathbf{W}^T\mathbf{W}}{p}$, where \mathbf{W} is a matrix of centered
 24 and scaled gene expression levels for the 200 lines, excluding the gene tested to maximize the

1 power to find an association (Yang et al. 2014), and p is the total number of genes. The TRM in
2 TWAS has similar role to the GRM in GWAS.

3 The effect of each gene's expression level on the phenotype was tested using a Wald
4 test of the form $\frac{\beta^2}{(SE(\beta))^2} \sim \chi_1^2$. Raw P -values and Benjamini-Hochberg (1995) FDR-corrected P -
5 values were computed.

6 The phenotypes were adjusted for the effects of Wolbachia and major polymorphic
7 inversions as described in the previous section. Because the phenotypes were adjusted, we did
8 not adjust gene expression in this analysis to avoid spurious associations due to adjustment on
9 both sides of the equation.

10 We also performed similar associations of quantitative traits with TEs and microbial gene
11 expression, using the same models as for TWAS but substituting TE and microbial expression
12 for gene expression levels. Quantitative trait phenotypes were adjusted for the effects of
13 Wolbachia and major polymorphic inversions but the TE and microbial expression data were
14 not. The TE analysis was performed for males and females separately, while sex-pooled
15 microbe expression data was used with female or male quantitative trait phenotypes since
16 microbial gene expression was not sex-specific.

17

18 **DATA ACCESS**

19 All raw and processed sequencing data generated in this study have been submitted to the
20 NCBI Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession
21 number GSE117850. The DGRP lines are available from the Bloomington Drosophila Stock
22 Center (Bloomington, IN). All analysis codes are available in Supplemental Codes and on
23 GitHub (<https://github.com/qgg-lab/dgrp-rna-seq/>).

24

25 **ACKNOWLEDGEMENTS**

1 This work was supported by National Institutes of Health grants R01 AA016560, R01 AG043490
2 and U01 DA041613 to T. F. C. M and R. R. H. A. and Genomic Selection in Animals and Plants
3 (GenSAP) funded by The Danish Council for Strategic Research to T. F. C. M. and F. M.

4

5 **AUTHOR CONTRIBUTIONS**

6 L. E., W. H. S. Z., F. M., R. R. H. A. and T. F. C. M. wrote the manuscript. L. E., W. H. S. Z., F.
7 M. and T. F. C. M. analyzed the data. M. A. C., R. L., G. A., M. S. G., J. M., G. S. A. and L. T
8 performed the research.

1 REFERENCES

- 2
- 3 Albert FW, Kruglyak L. 2015. The role of regulatory variation in complex traits and disease. *Nat*
- 4 *Rev Genet* **16**: 197-212.
- 5 Anders S, Pyl PT, Huber W. 2015. HTSeq-a Python framework to work with high-throughput
- 6 sequencing data. *Bioinformatics* **31**: 166-169.
- 7 Ayroles JF, Carbone MA, Stone EA, Jordan KW, Lyman RF, Magwire MM, Rollmann SM,
- 8 Duncan LH, Lawrence F, Anholt RR, et al. 2009. Systems genetics of complex traits in
- 9 *Drosophila melanogaster*. *Nat Genet* **41**: 299-307 (2009).
- 10 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful
- 11 approach to multiple testing. *J R Statist Soc B (Methodological)* **57**: 289-300.
- 12 Boyle EA, Li Y I, Pritchard JK. 2017. An expanded view of complex traits: from polygenic to
- 13 omnigenic. *Cell* **169**: 1177-1186.
- 14 Bray NL, Pimentel H, Melsted P, Pachter L. 2016. Near-optimal probabilistic RNA-seq
- 15 quantification. *Nature Biotech* **34**: 525-527.
- 16 Bryois J, Buil A, Evans DM, Kemp JP, Montgomery SB, Conrad DF, Ho KM, Ring S, Hurles M,
- 17 Deloukas P, et al. 2014. Cis and trans effects of human genomic variants on gene
- 18 expression. *PLoS Genet* **10**: e1004461.
- 19 Carbone MA, Yamamoto A, Huang W, Lyman RA, Meadors TB, Yamamoto R, Anholt RR,
- 20 Mackay TFC. 2016. Genetic architecture of natural variation in visual senescence in
- 21 *Drosophila*. *Proc Natl Acad Sci USA* **113**: E6620-E6629.
- 22 Chen Y, Zhu J, Lum PY, Yang X, Pinto S, MacNeil DJ, Zhang C, Lamb J, Edwards S, Sieberts
- 23 SK, et al. 2008. Variations in DNA elucidate molecular networks that cause disease.
- 24 *Nature* **452**: 429-435.
- 25 Civelek M, Lusk AJ. 2014. Systems genetics approaches to understand complex traits. *Nat Rev*
- 26 *Genet* **15**: 34-48.

- 1 Cookson W, Liang L, Abecasis G, Moffatt M, Lathrop M. 2009. Mapping complex disease traits
2 with global gene expression. *Nat Rev Genet* **10**: 184-194.
- 3 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras
4 TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15-21.
- 5 Emilsson V, Thorleifsson G, Zhang B, Leonardson AS, Zink F, Zhu J, Carlson S, Helgason A,
6 Walters GB, Gunnarsdottir S, et al. 2008. Genetics of gene expression and its effect on
7 disease. *Nature* **452**: 423-428.
- 8 Fagny M, Paulson JN, Kuijjer ML, Sonawane AR, Chen CY, Lopes-Ramos CM, Glass K,
9 Quackenbush J, Platig J. 2017. Exploring regulation in tissues with eQTL networks. *Proc*
10 *Natl Acad Sci USA* **114**: E7841-E7850.
- 11 Fiston-Lavier AS, Singh ND, Lipatov M, Petrov DA. 2010. *Drosophila melanogaster*
12 recombination rate calculator. *Gene* **463**: 18-20.
- 13 Garbherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
14 Raychowdhury R, Zeng D, et al. 2011. Full-length transcriptome assembly from RNA-Seq
15 data without a reference genome. *Nature Biotechnol* **29**: 644-652.
- 16 Garlapow ME, Huang W, Yarboro MT, Peterson KR, Mackay TFC. 2015. Quantitative genetics
17 of food intake in *Drosophila melanogaster*. *PLoS One* **10**: e0138129.
- 18 Gibson G, Powell JE, Marigorta UM. 2015. Expression quantitative trait locus analysis for
19 translational medicine. *Genome Medicine* **7**: 60.
- 20 Gramates LS, Marygold SJ, Santos GD, Urbano JM, Antonazzo G, Matthews BB, Rey AJ,
21 Tabone CJ, Crosby MA, Emmert DB, et al. 2017. FlyBase at 25: looking to the future.
22 *Nucleic Acids Res* **45**: D663-D671.
- 23 Hacısuleyman E, Goff LA, Trapnell C, Williams A, Henao-Mejia J, Sun L, McClanahan P,
24 Hendrickson DG, Sauvageau M, Kelley DR, et al. 2014. Topological organization of
25 multichromosomal regions by the long intergenic noncoding RNA Firre. *Nat Struct Mol*
26 *Biol* **21**: 198-206.

- 1 Harbison ST, McCoy LJ, Mackay TFC. 2013. Genome-wide association study of sleep in
2 *Drosophila melanogaster*. *BMC Genomics* **14**: 281.
- 3 Huang W, Massouras A, Inoue Y, Peiffer J, Ràmia M, Tarone AM, Turlapati L, Zichner T, Zhu D,
4 Lyman RF, et al. 2014. Natural variation in genome architecture among 205 *Drosophila*
5 *melanogaster* Genetic Reference Panel lines. *Genome Res* **24**: 1193-1208.
- 6 Huang W, Carbone MA, Magwire MM, Peiffer JA, Lyman RF, Stone EA, Anholt RR, Mackay
7 TFC. 2015. Genetic basis of transcriptome diversity in *Drosophila melanogaster*. *Proc*
8 *Natl Acad Sci USA* **112**: E6010-E6019.
- 9 Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. 2005. Repbase
10 Update, a database of eukaryotic repetitive elements. *Cytogenet Genome Res* **110**: 462-
11 467.
- 12 Kang YJ, Yang DC, Kong L, Hou M, Meng YQ, Wei L, Gao G. 2017. CPC2: a fast and accurate
13 coding potential calculator based on sequence intrinsic features. *Nucleic Acids Res* **3**:
14 45.
- 15 Khalil AM, Guttman M, Huarte M, Garber M, Raj A, Rivea Morales D, Thomas K, Presser A,
16 Bernstein BE, van Oudenaarden A, et al. 2009. Many human large intergenic noncoding
17 RNAs associate with chromatin-modifying complexes and affect gene expression. *Proc*
18 *Natl Acad Sci USA* **106**: 11667-11672.
- 19 Langfelder P, Horvath S. 2008. WGCNA: an R package for weighted correlation network
20 analysis. *BMC Bioinformatics* **9**: 559.
- 21 Lee YC, Langley CH. 2010. Transposable elements in natural populations of *Drosophila*
22 *melanogaster*. *Phil Trans Roy Soc B* **365**: 1219-1228.
- 23 Li H, Durbin R. 2010. Fast and accurate long-read alignment with Burrows-Wheeler transform.
24 *Bioinformatics* **26**: 589-595.
- 25 Liu B, de la Fuente A, Hoeschele I. 2008. Gene network inference via structural equation
26 modeling in genetical genomics experiments. *Genetics* **178**: 1763-1776.

- 1 Lyne R, Smith R, Rutherford K, Wakeling M, Varley A, Guillier F, Janssens H, Ji W, McLaren P,
2 North P, et al. 2007. FlyMine: an integrated database for *Drosophila* and Anopheles
3 genomics. *Genome Biol* **8**: R129.
- 4 Mackay TFC, Huang W. 2018. Charting the genotype-phenotype map: lessons from the
5 *Drosophila melanogaster* Genetic Reference Panel. *Wiley Interdiscip Rev Dev Bio* **7**. doi:
6 10.1002/wdev.289.
- 7 Mackay TFC, Richards S, Stone EA, Barbadilla A, Ayroles JF, Zhu D, Casillas S, Han Y,
8 Magwire MM, Cridland JM, et al. 2012. The *Drosophila melanogaster* Genetic Reference
9 Panel. *Nature* **482**: 173-178.
- 10 Mackay TFC, Stone EA, Ayroles JF. 2009. The genetics of quantitative traits: Challenges and
11 prospects. *Nat Rev Genet* **10**: 565-577.
- 12 Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorf LA, Hunter DJ, McCarthy MI, Ramos
13 EM, Cardon LR, Chakravarti A, et al. 2009. Finding the missing heritability of complex
14 diseases. *Nature* **461**: 747-753.
- 15 Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.
16 *EMBnet Journal* **17**: 10-12.
- 17 Massouras A, Waszak SM, Albarca-Aguilera M, Hens K, Holcombe W, Ayroles JF, Dermitzakis
18 ET, Stone EA, Jensen JD, Mackay TFC, et al. 2012. Genomic variation and its impact on
19 gene expression in *Drosophila melanogaster*. *PLoS Genet* **8**: e1003055.
- 20 Mi H, Huang X, Muruganujan A, Tang H, Mills C, Kang D, Thomas PD. 2017. PANTHER
21 version 11: expanded annotation data from Gene Ontology and Reactome pathways, and
22 data analysis tool enhancements. *Nucleic Acids Res* **45**: D183-D189.
- 23 Nicolae DL, Gamazon E, Zhang W, Duan S, Dolan ME, Cox NJ. 2010. Trait-associated SNPs
24 are more likely to be eQTLs: annotation to enhance discovery from GWAS. *PLoS Genet*
25 **6**: e1000888.

- 1 Ogura T, Busch W. 2016. Genotypes, networks, phenotypes: Moving toward plant systems
2 genetics. *Annu Rev Cell Dev Biol* **32**: 103-126.
- 3 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Maller J, Sklar P, de
4 Bakker PI, Daly MJ, et al. 2007. PLINK: a tool set for whole-genome association and
5 population-based linkage analyses. *Am J Hum Genet* **81**: 559-575.
- 6 Ransohoff JD, Wei Y, Khavari PA. 2018. The functions and unique features of long intergenic
7 non-coding RNA. *Nat Rev Mol Cell Biol* **19**: 143-157.
- 8 Riddle NC, Minoda A, Kharchenko PV, Alekseyenko AA, Schwartz YB, Tolstorukov MY,
9 Gorchakov AA, Jaffe JD, Kennedy C, Linder-Basso D, et al. 2011. Plasticity in patterns of
10 histone modifications and chromosomal proteins in *Drosophila* heterochromatin. *Genome*
11 *Res* **21**: 147-163.
- 12 Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential
13 expression analysis of digital gene expression data. *Bioinformatics* **26**: 139-40.
- 14 Robinson MD, Oshlack A. 2010. A scaling normalization method for differential expression
15 analysis of RNA-seq data. *Genome Biol* **11**: R25.
- 16 Rockman MV. 2008. Reverse engineering the genotype-phenotype map with natural genetic
17 variation. *Nature* **456**: 738-744.
- 18 Rogoyski OM, Pueyo JI, Couso JP, Newbury SF. 2017. Functions of long non-coding RNAs in
19 human disease and their conservation in *Drosophila* development. *Biochem Soc Trans*
20 **45**: 895-904.
- 21 Schughart K, Williams RW. 2017. *Systems Genetics Methods and Protocols*. Humana Press,
22 New York NY.
- 23 Shorter J, Couch C, Huang W, Carbone MA, Peiffer J, Anholt RR, Mackay TFC. 2015. Genetic
24 architecture of natural variation in *Drosophila melanogaster* aggressive behavior. *Proc*
25 *Natl Acad Sci USA* **112**: E3555-E2563.

- 1 Sieberts SK, Schadt EE. 2007. Moving toward a system genetics view of disease. *Mamm*
2 *Genome* **18**: 389-401.
- 3 Spradling AC, Rubin GM. 1981. Drosophila genome organization: conserved and dynamic
4 aspects. *Annu Rev Genet* **15**: 219–264.
- 5 Stoiber M, Celniker S, Cherbas L, Brown B, Cherbas P. 2016. Diverse hormone response
6 networks in 41 independent Drosophila cell lines. *G3* **6**: 683-694.
- 7 Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL, Rinn JL,
8 Pachter L. 2012. Differential gene and transcript expression analysis of RNA-seq
9 experiments with TopHat and Cufflinks. *Nat Protoc* **7**: 562-578.
- 10 Visscher PM, Brown MA, McCarthy MI, Yang J. 2012. Five years of GWAS discovery. *Am J*
11 *Hum Genet* **90**: 7-24.
- 12 Wang J, Samuels DC, Zhao S, Xiang YY, Guo Y. 2017. Current research on non-coding
13 ribonucleic acid (RNA). *Genes* **8**: 366.
- 14 Wang KC, Yang YW, Liu B, Sanyal A, Corces-Zimmerman R, Chen Y, Lajoie BR, Protacio A,
15 Flynn RA, Gupta RA, et al. 2011. A long noncoding RNA maintains active chromatin to
16 coordinate homeotic gene expression. *Nature* **472**: 120-124.
- 17 Yang J, Lee SH, Goddard ME, Visscher PM. 2011. GCTA: a tool for genome-wide complex trait
18 analysis. *Am J Hum Genet* **88**: 76-82.
- 19 Yang J, Zaitlen NA, Goddard ME, Visscher PM, Price AL. 2014. Advantages and pitfalls in the
20 application of mixed-model association methods. *Nat Genet* **46**: 100-106.
- 21 Zhou S, Morgante F, Geisz MS, Ma J, Anholt RRA, Mackay TFC. 2020. Systems genetics of the
22 *Drosophila* metabolome. *Genome Res*, in press.

1 FIGURE LEGENDS

2

3 **Figure 1. Genetic variation of gene expression in the DGRP.** (A) Sexual dimorphism of gene
4 expression. Red (blue) indicates significant up-regulation in females (males). (B) Distribution of
5 H^2 estimates for annotated genes and NTRs in females. (C) Distribution of H^2 estimates for
6 annotated genes and NTRs in males. (D) WGCNA modules for annotated genes and NTRs in
7 females. (E) WGCNA modules for annotated genes and NTRs in males. Heatmaps show the
8 pairwise correlation of all genes in each module, sorted by average connectivity, with the most
9 tightly connected module at the top left.

10

11 **Figure 2. Genomic location of eQTLs for gene expression and genes they regulate.** eQTL
12 chromosome positions (bp) are given on the X-axis, and the genes with which they are
13 associated on the Y-axis. Red points denote female-specific eQTLs, blue indicates male-
14 specific eQTLs, and black shows eQTLs shared by males and females.

15

16 **Figure 3. Large *cis-trans* eQTL genetic network in females and males.** Node interior colors
17 indicate genomic location of genes (yellow: euchromatic regions with normal recombination;
18 gray: euchromatic regions with reduced recombination; blue: heterochromatin). Node border
19 colors denote annotated gene (gray) or NTR (red). Node shape indicates whether a gene is a
20 regulator and/or target (triangles: regulator only; squares: target only; circles: both regulator and
21 target). The node size indicates the number of node connections. Arrows on the edges point to
22 the target. Edges are color coded to show female-specific regulation (red), male-specific
23 regulation (blue) and regulation common to both sexes (black).

24

25 **Figure 4. Genetic variation of TE expression in the DGRP.** (A) Total signal for each TE
26 family, summed over all individual transposon sequences and averaged across all DGRP lines,

1 sex, and replicates. **(B)** Distribution of copy number independent H^2 estimates for TE sequences
2 in females. **(C)** Distribution of copy number independent H^2 estimates for TE sequences in
3 males. **(D)** WGCNA modules of TEs for females. **(E)** WGCNA modules of TEs for males.
4 Heatmaps are depicted as in Figure 1. TE sequences not assigned to any module are included
5 at the bottom right.

6

7 **Figure 5. TE genetic regulatory network.** Symbols and color-coding are as for Figure 3. Black
8 squares denote TE sequences.

9

10 **Figure 6. Genetic variation of microbiome composition.** **(A)** The proportion of microbiome
11 signal in RNA-seq libraries aligned to species in each genus or viral group. **(B)** Line means of
12 total microbial signal (excluding Wolbachia). **(C)** Distribution of H^2 estimates for individual
13 microbe species. **(D)** WGCNA modules for microbial species. Heatmaps are depicted as in
14 Figure 1. Species not assigned to any module are included at the bottom right.











