# PNAS

## www.pnas.org

**Supplementary Information for**

Genomic Regions Influencing Aggressive Behavior in Honey Bees
are Defined by Colony Allele Frequencies

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

**Corresponding Authors:**
Gene E. Robinson; **Email:** generobi@illinois.edu.
Guojie Zhang; **Email:** Guojie.Zhang@bio.ku.dk.
Matthew E. Hudson; **Email:** mhudson@illinois.edu.

**This PDF includes:**

Captions for Dataset S1-2
Figs. S1-7

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

1 of 11

www.pnas.org/cgi/doi/10.1073/pnas.1922927117

20   **Dataset S1 (annotation_data_table_s1.xlsx).**
21           Data table in Excel format with three sheets corresponding to a legend and reference
22   information, colony phenotype and collection information, a gene annotation list for the genes
23   overlapping haplotypes containing a significant SNP, and a list of significant SNPs . The colony
24   information sheet contains specific information including collection and sampling data as well as
25   phenotype details for both measures of colony aggression. The gene annotation list contains
26   information on the 254 genes within genomic regions of significant correlation. Information
27   includes linkage group, number of haplotype blocks with significant SNPs in overlap with the gene,
28   *A. mellifera* NCBI gene IDs and names, and gene ID, symbol and name for the nearest *D.*
29   *melanogaster* homolog. Also highlighted are the subset of 56 genes in overlap with haplotype
30   blocks also containing markers with evidence of selection. The SNP list contains information on
31   the specific markers identified by our colony-level analysis. Included are positional information
32   (linkage group, base pair position), nucleotide information, and resulting p-values for the
33   individual- and colony-level analysis.
34

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao,
Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson
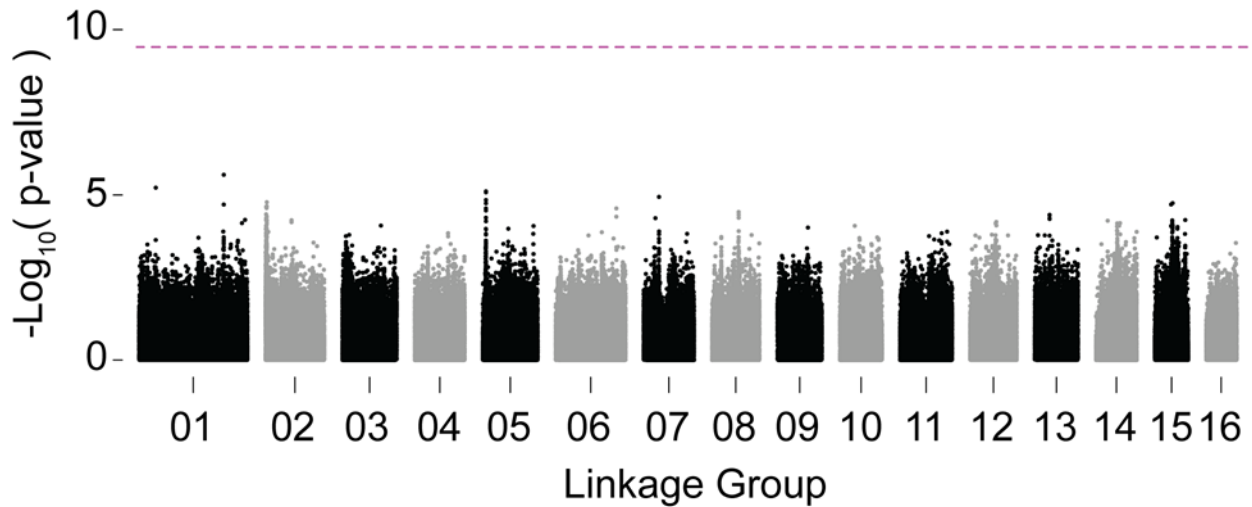
35 **Dataset S2 (simulation_code_s5.r).**
36     Annotated algorithm which conducts a P-value assessment using simulated Principal
37 Components. Annotations summarize conceptual framework and provide the model under
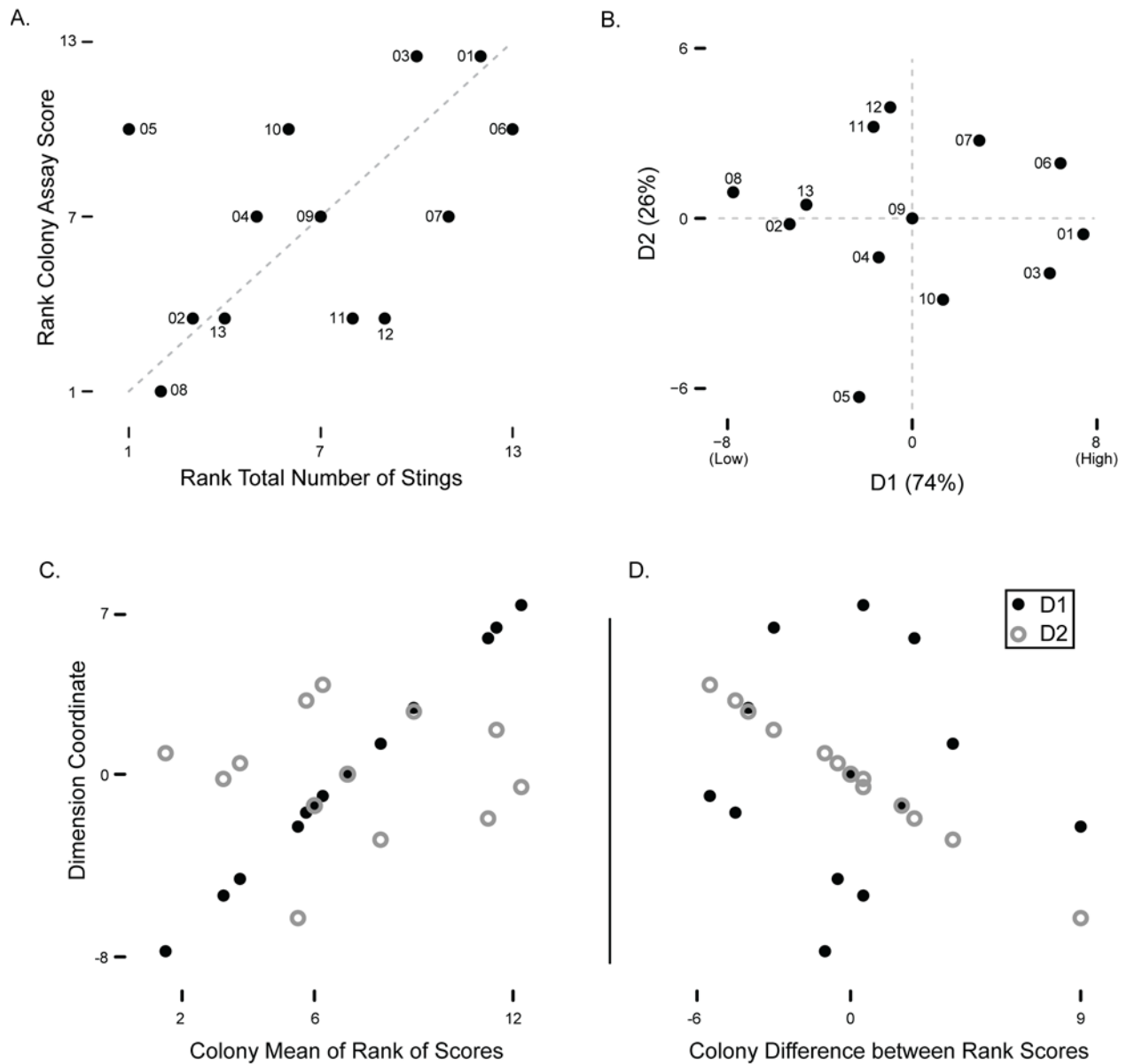38 consideration for SNP x colony phenotype which mirrors the model considered in the manuscript.
39
40
41

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

42

**Fig. S1. Genome-wide associations of aggression at the individual level.** Manhattan plot of P value distributions across the genome for the correlation of individual-level genotype to individual behavioral phenotype (Soldier vs. Forager). The dashed magenta line represents the Bonferroni adjusted threshold ($\alpha$ = 3.35E$^{-10}$) consistent between individual and colony level analyses.
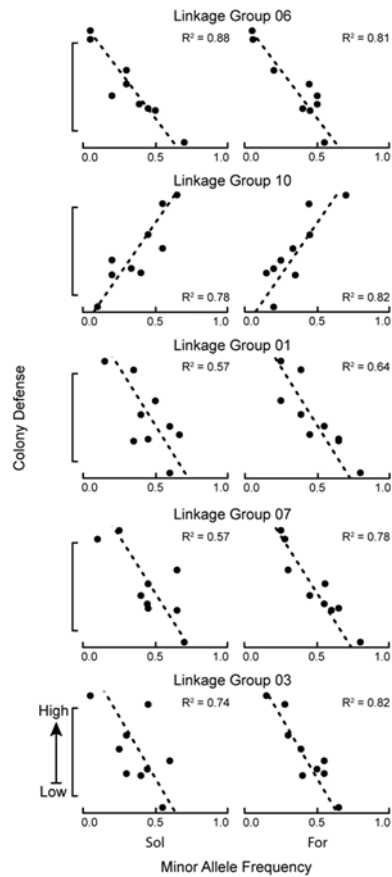
48

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

**Fig. S2. Assessment of colony phenotype.** Panel A: correlation between the rank values of the response for each colony in the two behavioral assays. Panel B: results from Multidimensional Scaling (MDS). In both panels each point corresponds to a colony and is labeled with its colony number. Panels C and D summarize the four correlations between the dimension coordinate value for each of the colonies and the corresponding summary of phenotype. Panel C summarizes the relationship between the dimension value and the per-colony mean of rank scores between the two assays. Panel D summarizes the correlation between the same dimensional position and the per-colony difference between the ranks of the behavioral assays. As in panels A and B, each point is a colony, black points are used when correlating against the first MDS dimension (D1) and open grey circles are used when correlating against the second MDS dimension (D2).

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

62



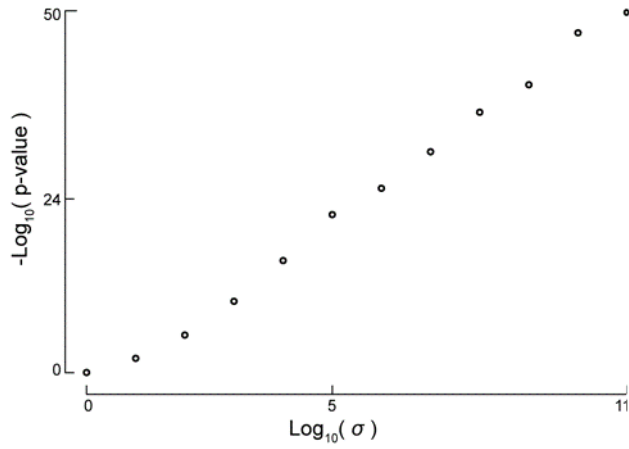Fig. S3. Correlation between minor allele frequencies (MAF) for each behavioral group and the colony phenotype.
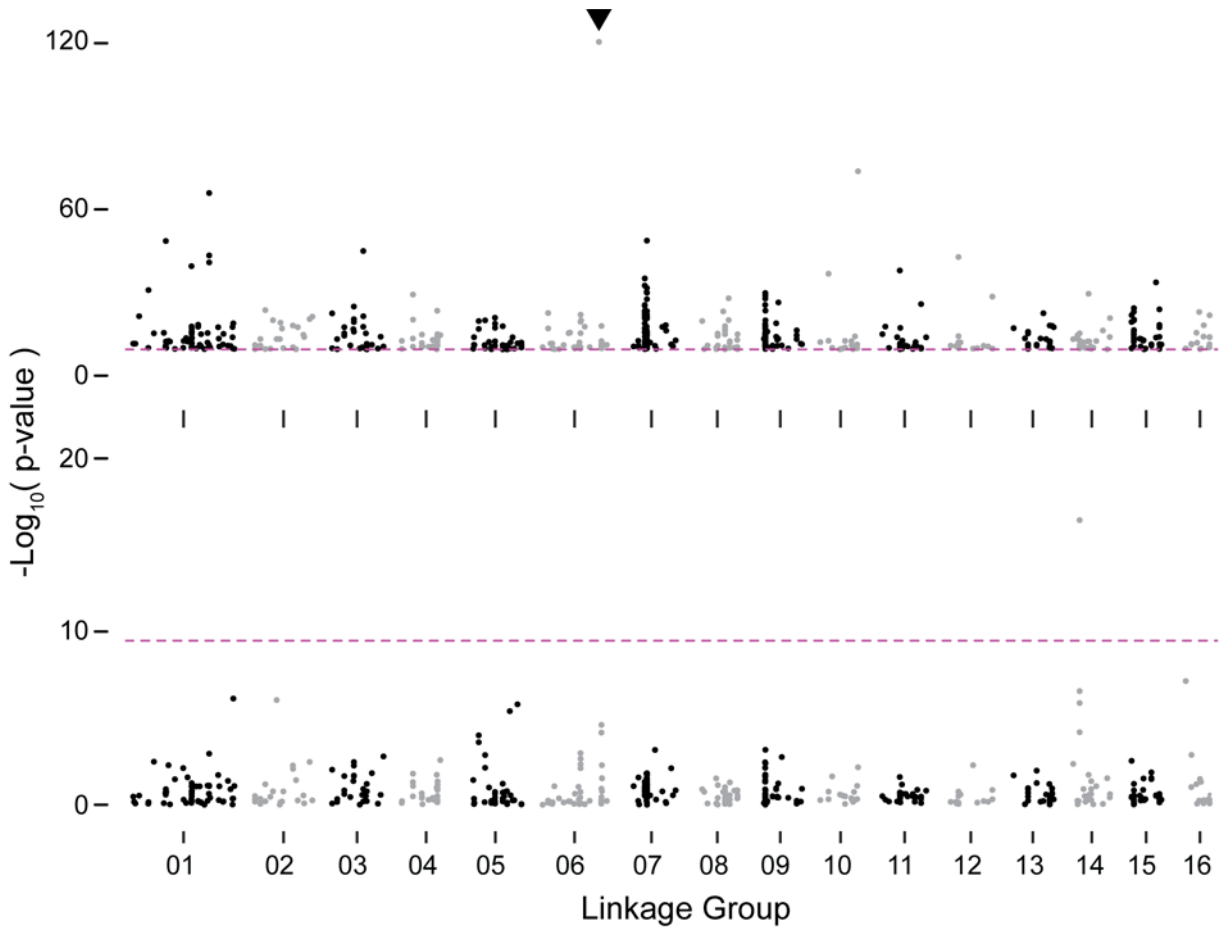
63

64

65 **Fig. S3. Correlation between minor allele frequencies (MAF) for each behavioral group and**
66 **the colony phenotype.** Goodness of fit ($R^2$) was estimated to the model function derived from
67 the colony-level fit (dashed line). A paired set of plots is provided for each of the top 5 focal SNPs
68 of the top 5 peaks of association candidate SNPs (Fig. 1C) one for Soldiers (Sol) and the other
69 for Foragers (For). Y axis corresponds to colony defense (D1; SI Appendix, Fig. S2), X axis to the
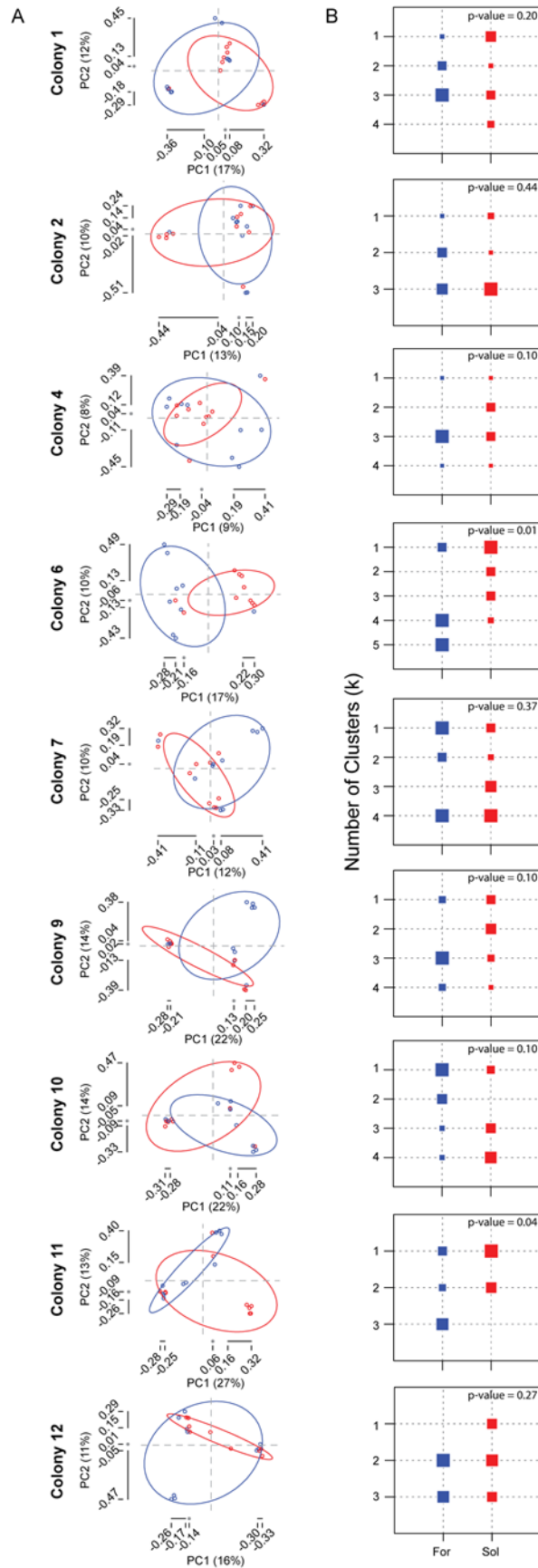70 MAF, and each point represents a colony.

71

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

72
73
74 **Fig. S4. Relationship between a range of simulated residual variances and resulting P**
75 **values.** In our analysis, even with N = 9, as residual variance decreases, p-values reach levels
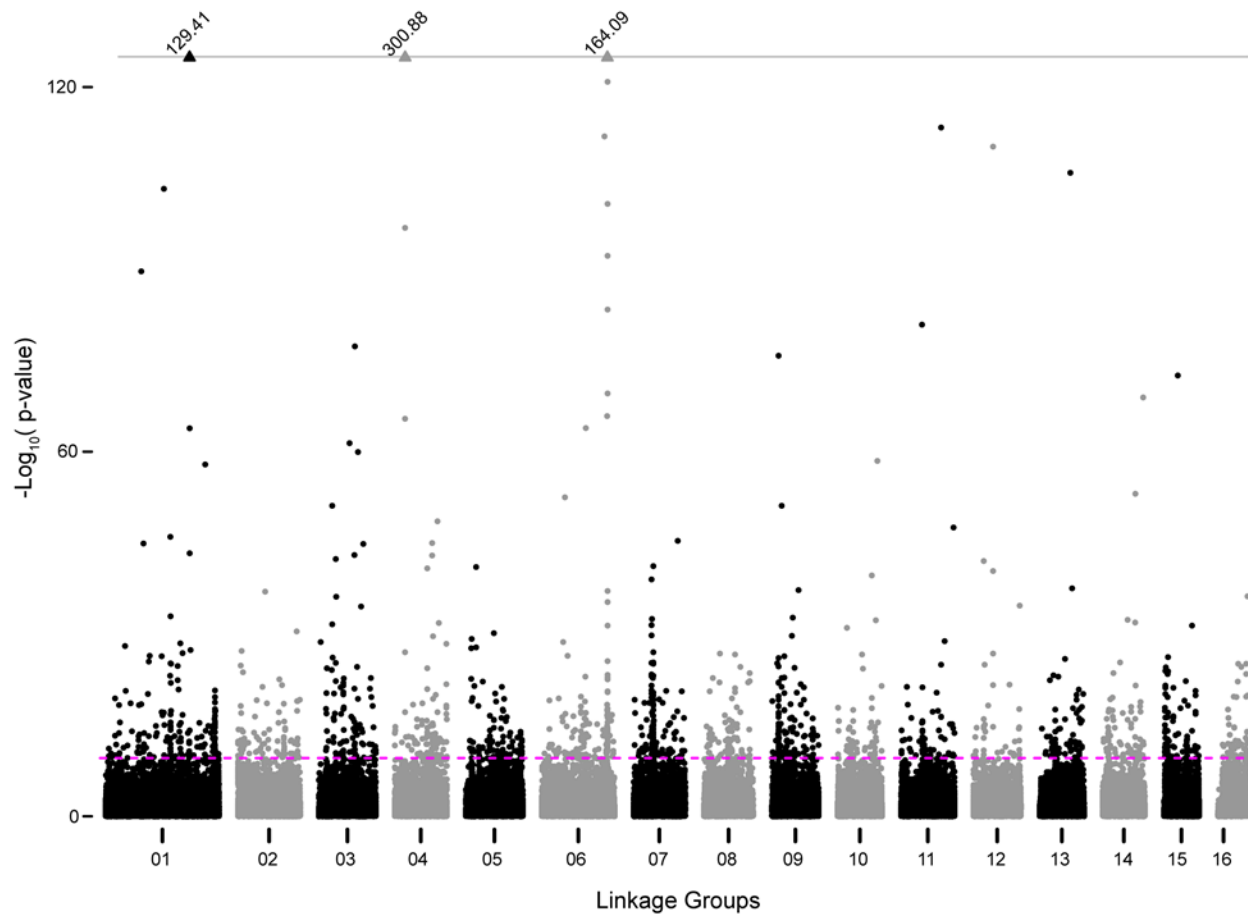76 of 10 - 50 ($1.00E^{-10}$ – $1.00E^{-50}$), as we detected (Fig. 1).
77

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao,
Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

**Fig. S5. Analysis of covariation across top candidate SNPs.** Top panel: P values of SNPs from Fig 1C that pass the significance threshold, plotted against genomic location. Bottom panel: P values of significant SNPs from Fig 1C with the most significant SNP (highlighted by the black triangle) included in the model as a covariate.

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

87 **Fig. S6. Analysis of concordance of genetic diversity with aggressive phenotype within**
88 **each colony.** (A) Principal component analysis of genetic variation is summarized by the first two
89 principal components derived from the genotype matrix for each colony. Each point corresponds
90 to a sample, and each sample is colored by behavioral group: blue = Forager (For), red = Soldier
91 (Sol). An ellipse encapsulating 65% of the samples within a behavioral group (~7 in each group)
92 is provided to further highlight distribution of the behavioral groups across the PC space. (B)
93 Optimal number of clusters was determined via iterative k-means clustering and the elbow method
94 using the within-group total sum of squares. The distribution of behavioral groups between genetic
95 clusters was assessed for each colony using a Fisher's exact test (P value at top right of table).

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson

Arian Avalos, Miaoquan Fang, Hailin Pan, Aixa Ramirez Lluch, Alexander E. Lipka, Sihai Dave Zhao, Tugrul Giray, Gene E. Robinson, Guojie Zhang & Matthew E. Hudson