

## RESEARCH

# Transcriptomic responses to diet quality and viral infection in *Apis mellifera*

Lindsay Rutter<sup>1</sup>, Bryony C. Bonning<sup>6</sup>, Dianne Cook<sup>2</sup>, Amy L. Toth<sup>3,4</sup> and Adam G. Dolezal<sup>5\*</sup>

\*Correspondence:

adolezal@illinois.edu

<sup>5</sup>Department of Entomology,

University of Illinois at

Urbana-Champaign, Urbana, IL

61801, USA

Full list of author information is available at the end of the article

## Abstract

**Background:** Parts of Europe and the United States have witnessed dramatic losses in commercially managed honey bees over the past decade to what is considered an unsustainable extent. The large-scale loss of bees has considerable implications for the agricultural economy because bees are one of the leading pollinators of numerous crops. Bee declines have been associated with several interactive factors. Recent studies suggest nutritional and pathogen stress can interactively contribute to bee physiological declines, but the molecular mechanisms underlying interactive effects remain unknown. In this study, we provide insight into this question by using RNA-sequencing to examine how monofloral diets and Israeli acute paralysis virus inoculation influence gene expression patterns in bees.

**Results:** We found a considerable nutritional response, with almost 2,000 transcripts changing with diet quality. The majority of these genes were over-represented for nutrient signaling (insulin resistance) and immune response (Notch signaling and JaK-STAT pathways). In our experimental conditions, the transcriptomic response to viral infection was fairly limited. We only found 43 transcripts to be differentially expressed, some with known immune functions (argonaute-2), transcriptional regulation, and muscle contraction. We created contrasts to explore whether protective mechanisms of good diet were due to direct effects on immune function (resistance) or indirect effects on energy availability (tolerance). A similar number of resistance and tolerance candidate differentially expressed genes were found, suggesting both processes may play significant roles in dietary buffering from pathogen infection.

**Conclusions:** Through transcriptional contrasts and functional enrichment analysis, we contribute to our understanding of the mechanisms underlying feedbacks between nutrition and disease in bees. We also show that comparing results derived from combined analyses across multiple RNA-seq studies may allow researchers to identify transcriptomic patterns in bees that are concurrently less artificial and less noisy. This work underlines the merits of using data visualization techniques and multiple datasets to interpret RNA-sequencing studies.

**Keywords:** Honey bee; RNA-sequencing; Israeli acute paralysis virus; Monofloral pollen; Visualization

## 1 Background

- 2 Commercially managed honey bees have undergone unusually large declines in the
- 3 United States and parts of Europe over the past decade [1, 2, 3], with annual

4 mortality rates exceeding what beekeepers consider sustainable [4, 5]. More than 70  
5 percent of major global food crops (including fruits, vegetables, and nuts) at least  
6 benefit from pollination, and yearly insect pollination services are valued worldwide  
7 at \$175 billion [6]. As honey bees are largely considered to be the leading pollinator  
8 of numerous crops, their marked loss has considerable implications for agricultural  
9 sustainability [7].

10 Honey bee declines have been associated with several factors, including pesti-  
11 cide use, parasites, pathogens, habitat loss, and poor nutrition [8, 9]. Researchers  
12 generally agree that these stressors do not act in isolation; instead, they appear  
13 to influence the large-scale loss of honey bees in an interactive fashion as the en-  
14 vironment changes [10]. Nutrition and viral infection are two broad factors that  
15 pose heightened dangers to honey bee health in response to recent environmental  
16 changes.

17 Pollen is a main source of nutrition (including proteins, amino acids, lipids, sterols,  
18 starch, vitamins, and minerals) in honey bees [11, 12]. At the individual level, pollen  
19 supplies most of the nutrients necessary for physiological development [13] and is  
20 believed to have considerable impact on longevity [14]. At the colony level, pollen  
21 enables young workers to produce jelly, which then nourishes larvae, drones, older  
22 workers, and the queen [15, 16]. Various environmental changes (including urban-  
23 ization and monoculture crop production) have significantly altered the nutritional  
24 profile available to honey bees. In particular, honey bees are confronted with a  
25 less diverse selection of pollen, which is of concern because mixed-pollen (polyflo-  
26 ral) diets are generally considered healthier than single-pollen (monofloral) diets  
27 [17, 18, 19]. Indeed, reported colony mortality rates are higher in developed land  
28 areas compared to undeveloped land areas [20], and beekeepers rank poor nutrition  
29 as one of the main reasons for colony losses [21]. Understanding how low diversity

30 diets affect honey bee health will be crucial to resolve problems that may arise as  
31 agriculture continues to intensify throughout the world [22, 23].

32 Viral infection was a comparatively minor problem in honey bees until the last  
33 century when the ectoparasitic varroa mite (*Varroa destructor*) spread worldwide  
34 [24]. This mite feeds on honey bee hemolymph [25], transmits multiple viruses,  
35 and supports replication of some viruses [26, 27, 28, 29]. More than 20 honey bee  
36 viruses have been identified [30]. One of these viruses that has been linked to honey  
37 bee decline is Israeli acute paralysis virus (IAPV), a positive-sense RNA virus of  
38 the family Dicistroviridae [31]. IAPV infection causes shivering wings, decreased  
39 locomotion, muscle spasms, paralysis, and high premature death percentages in  
40 caged infected adult honey bees [32]. IAPV has demonstrated higher infectious  
41 capacities than other honey bee viruses under certain conditions [33] and is more  
42 prevalent in colonies that do not survive the winter [34].

43 Although there is growing interest in how viruses and diet quality affect the health  
44 and sustainability of honey bees, as well as a recognition that such factors might  
45 operate interactively, there are only a small number of experimental studies thus  
46 far directed toward elucidating the interactive effects of these two factors in honey  
47 bees [35, 36, 37, 38, 39]. We recently used laboratory cages and nucleus hive experi-  
48 ments to investigate the health effects of these two factors, and our results show the  
49 importance of the combined effects of both diet quality and virus infection. Specifi-  
50 cally, ingestion by honey bees of high quality pollen is able to mitigate virus-induced  
51 mortality to the level of diverse, polyfloral pollen [40].

52 Following up on these findings, we now aim to understand the corresponding  
53 underlying mechanisms by which high quality diets protect bees from virus-induced  
54 mortality. For example, it is not known whether the protective effect of good diet  
55 is due to direct, specific effects on immune function (resistance), or if it is due  
56 to indirect effects of good nutrition on vigor (tolerance) [41]. Transcriptomics is

one means to better understand the mechanistic underpinnings of dietary and viral effects on honey bee health. Transcriptomic analysis can help us identify 1) the genomic scale of transcriptomic response to diet and virus infection, 2) whether these factors interact in an additive or synergistic way on transcriptome function, and 3) the types of pathways affected by diet quality and viral infection. This information, heretofore lacking in the literature, can help us better understand how good nutrition may be able to serve as a “buffer” against other stressors [42].

There are only a small number of published experiments examining gene expression patterns related to diet effects [43] and virus infection effects [44] in honey bees, but there have been several such studies in model organisms. Model insect studies can inform studies of honey bee transcriptomic responses, using functional inference of as-of-yet uncharacterized honey bee genes based on orthology to *Drosophila* and other model organisms. Previous *Drosophila* studies that examined various diet effects have found gene expression changes related to immunity, metabolism, cell cycle activity, DNA binding, transcription, and insulin signaling [45, 46, 47, 43]. While similar transcriptomic studies have been limited in honey bees, one study found that pollen nutrition upregulates genes involved in macromolecule metabolism, longevity, and the insulin/TOR pathway required for physiological development [43]. Numerous studies on the transcriptomic effects of virus infection in model insect organisms have shown that RNA silencing, transcriptional pausing, Toll pathways, IMD pathways, JAK/STAT pathways, and Toll-7 autophagy pathways play substantial roles in virus-host systems [48, 49]. Virus-bee systems have also revealed key factors of these antiviral conserved defense pathways [50].

As far as we know, there are few to no studies investigating honey bee gene expression patterns specifically related to monofloral diets, and few studies investigating honey bee gene expression patterns related to the combined effects of diet in any broad sense and viral inoculation in any broad sense [38]. In this study, we examine

84 how monofloral diets and viral inoculation influence gene expression patterns in  
85 honey bees by focusing on four treatment groups (low quality diet without IAPV  
86 exposure, high quality diet without IAPV exposure, low quality diet with IAPV  
87 exposure, and high quality diet with IAPV exposure). For our diet factor, we exam-  
88 ined two monofloral pollen diets, rockrose (*Cistus* sp.) and chestnut (*Castanea* sp.).  
89 Rockrose pollen is generally considered less nutritious than chestnut pollen because  
90 it contains smaller amounts of protein, amino acids, antioxidants, calcium, and iron  
91 [40, 51]. We conduct RNA-sequencing analysis on a randomly selected subset of the  
92 honey bees we used in our previous study (as is further described in our methods  
93 section). We then examine pairwise combinations of treatment groups, the main  
94 effect of monofloral diet, the main effect of IAPV exposure, and the combined effect  
95 of the two factors on gene expression patterns.

96 We also compare the main effect of IAPV exposure in our dataset to that obtained  
97 in a previous study conducted by Galbraith and colleagues [44]. While our study  
98 examines honey bees from naturally mated colonies, the Galbraith study examined  
99 honey bees from single-drone colonies. As a consequence, the honey bees in our  
100 study will be on average 25% genetically identical, whereas honey bees from the  
101 Galbraith study will be on average 75% genetically identical [52]. We note that  
102 the difference between these studies may be even greater than this as we used  
103 naturally mated honey bees from 15 different colonies. We should therefore expect  
104 that the Galbraith study may generate data with higher signal:to:noise ratios than  
105 our data due to lower genetic variation between its replicates. At the same time, our  
106 honey bees will be more likely to display the health benefits gained from increased  
107 genotypic variance within colonies, including decreased parasitic load [53], increased  
108 tolerance to environmental changes [54], and increased colony performance [55, 56].  
109 Given that honey bees are naturally very polyandrous [57], our naturally mated  
110 honey bees may also reflect more realistic environmental and genetic simulations.

111 Taken together, each study provides a different point of value: Our study likely  
112 presents less artificial data while the Galbraith data likely presents less messy data.  
113 We wish to explore how the gene expression effects of IAPV inoculation compare  
114 between these two studies that used such different experimental designs. To achieve  
115 this objective, we use visualization techniques to assess the signal:to:noise ratio  
116 between these two datasets, and differential gene expression (DEG) analyses to  
117 determine any significantly overlapping genes of interest between these two datasets.  
118 As RNA-sequencing data can be biased [58, 59, 60], this comparison allowed us  
119 to characterize how repeatable and robust our RNA-sequencing results were in  
120 comparison to previous studies. It also allowed us to shine light on how experimental  
121 designs that control genetic variability to different extents might affect the resulting  
122 gene expression data in honey bees. We suggest that in-depth data visualization  
123 approaches can be useful for cross-study comparisons and validation of noisy RNA-  
124 sequencing data in the future.

## 125 Results

### 126 Mortality and virus titers

127 We reanalyzed our previously published dataset with a subset that focuses on diet  
128 quality and is more relevant to the current study. We show the data subset here to  
129 inform the RNA-sequencing comparison because we reduced the number of treat-  
130 ments from the original published data (from eight to four) [40] as a means to focus  
131 on diet quality effects.

132 As shown in Figure 1, mortality rates of honey bees 72 hours post-inoculation  
133 significantly differed among the treatment groups (mixed model ANOVA across all  
134 treatment groups,  $df = 3, 54$ ;  $F = 10.03$ ;  $p < 2.34e-05$ ). The effect of virus treatment  
135 (mixed model ANOVA,  $df = 1, 54$ ;  $F = 24.73$ ;  $p < 7.04e-06$ ) and diet treatment  
136 (mixed model ANOVA,  $df = 1, 54$ ;  $F = 5.32$ ;  $p < 2.49e-02$ ) were significant, but  
137 the interaction between the two factors (mixed model ANOVA,  $df = 1, 54$ ;  $F =$

138 4.72e-02,  $p = 8.29\text{e-}01$ ) was not significant. We compared mortality levels based  
139 on pairwise comparisons: For a given diet, honey bees exposed to the virus showed  
140 significantly higher mortality rate than honey bees not exposed to the virus. Bees fed  
141 rockrose pollen had significantly elevated mortality with virus infection compared  
142 to uninfected controls (Benjamini-Hochberg,  $p < 1.53\text{e-}03$ ), and bees fed chestnut  
143 pollen similarly had significantly elevated mortality with virus infection compared  
144 to controls (Benjamini-Hochberg,  $p < 3.12\text{e-}03$ ) (Figure 1).

145 As shown in Figure 2, IAPV titers of honey bees 72 hours post-inoculation sig-  
146 nificantly differed among the treatment groups (mixed model ANOVA across all  
147 treatment groups,  $df = 3, 33$ ;  $F = 6.10$ ;  $p < 2.03\text{e-}03$ ). The effect of virus treatment  
148 (mixed model ANOVA,  $df = 1, 33$ ;  $F = 15.04$ ;  $p < 4.75\text{e-}04$ ) was significant, but the  
149 diet treatment (mixed model ANOVA,  $df = 1, 33$ ;  $F = 2.55$ ;  $p = 1.20\text{e-}01$ ) and the  
150 interaction between the two factors (mixed model ANOVA,  $df = 1, 33$ ;  $F = 7.02\text{e-}$   
151  $01$ ,  $p = 4.08\text{e-}01$ ) were not significant. We compared IAPV titers based on pairwise  
152 comparisons: Bees fed rockrose pollen had significantly elevated IAPV titers with  
153 virus infection compared to uninfected controls (Benjamini Hochberg,  $p < 7.56\text{e-}$   
154  $03$ ). However, bees fed chestnut pollen did not have significantly elevated IAPV  
155 titers with virus infection compared to uninfected controls (Benjamini Hochberg,  $p$   
156  $= 6.29\text{e-}02$ ). Overall, we interpreted these findings to mean that high-quality chest-  
157 nut pollen could partially “rescue” high virus titers resulting from the inoculation  
158 treatment, whereas low-quality rockrose pollen could not (Figure 2).

#### 159 Transcriptomic responses to virus infection and diet

160 We observed a substantially larger number of differentially expressed genes (DEGs)  
161 in our diet main effect ( $n = 1,914$ ) than in our virus main effect ( $n = 43$ ) (Sup-  
162 plementary table 1 A and B, Additional file 1). In the diet factor, more DEGs  
163 were upregulated in the more-nutritious chestnut group ( $n = 1,033$ ) than in the  
164 less-nutritious rockrose group ( $n = 881$ ). In the virus factor, there were more virus-

upregulated DEGs ( $n = 38$ ) than control-upregulated DEGs ( $n = 5$ ). While these reported DEG counts are from the DESeq2 package, we saw similar trends for the edgeR and limma package results (Supplementary table 1, Additional file 1 and Additional file 18).

GO analysis of the chestnut-upregulated DEGs revealed the following over-represented biological functions: Wnt signaling, hippo signaling, and dorso-ventral axis formation, as well as pathways related to circadian rhythm, mRNA surveillance, insulin resistance, inositol phosphate metabolism, FoxO signaling, ECM-receptor interaction, phototransduction, Notch signaling, JaK-STAT signaling, MAPK signaling, and carbon metabolism (Supplementary table 2, Additional file 1). GO analysis of the rockrose DEGs revealed pathways related to terpenoid backbone biosynthesis, homologous recombination, SNARE interactions in vesicular transport, aminoacyl-tRNA biosynthesis, Fanconi anemia, and pyrimidine metabolism (Supplementary table 3, Additional file 1).

With so few DEGs ( $n = 43$ ) in our virus main effect comparison, we focused on individual genes and their known functionalities rather than GO over-representation (Table 1). Of the 43 virus-related DEGs, only 10 had GO assignments within the DAVID database. These genes had putative roles in the recognition of pathogen-related lipid products and the cleaving of transcripts from viruses, as well as involvement in ubiquitin and proteasome pathways, transcription pathways, apoptotic pathways, oxidoreductase processes, and several more functions (Table 1).

No interaction DEGs were observed between the diet and virus factors of the study, in any of the pipelines (DESeq2, edgeR, and limma).

The number of DEGs across the six treatment pairings between the diet and virus factor ranged from 0 to 955 (Supplementary table 8, Additional file 1). Again, diet level appeared to have greater influence on the number of DEGs than the virus level. Across every pair comparing the chestnut and rockrose levels, regardless of the



192 virus level, the number of chestnut-upregulated DEGs was higher than the number  
193 of rockrose-upregulated DEGs (Supplementary table 8 C, D, E, F, Additional file 1).  
194 Virus-treated bees showed equal to or more upregulated genes relative to controls,  
195 under both diet treatments (Supplementary table 8 A and B, Additional file 1).  
196 These trends were observed for all three pipelines used (DESeq2, edgeR, and limma).

#### 197 Transcriptomic data visualization and comparison to a previous study

198 We wished to explore the signal:to:noise ratio between the Galbraith dataset and  
199 our dataset. Note that the Galbraith dataset contained three samples for each  
200 virus level, while our dataset contained twelve samples for each virus level. Basic  
201 PCA plots were constructed with the DESeq2 analysis pipeline and showed  
202 that the Galbraith dataset may separate the infected and uninfected honey bees  
203 better than our dataset (Additional file 2). We also noted that the first replicate  
204 of both treatment groups in the Galbraith data did not cluster as cleanly in the  
205 PCA plots. However, through this automatically-generated plot, we can only visu-  
206 alize information at the sample level. Wanting to learn more about the data at the  
207 gene level, we continued with new visualization techniques that are available online  
208 (<https://lrutter.github.io/bigPint>) and are in preparation for publication.

209 We used parallel coordinate lines superimposed onto side-by-side boxplots to visu-  
210 alize the DEGs associated with virus infection in the two studies. The background  
211 side-by-side boxplot represents the distribution of *all* genes in the data, and each  
212 parallel coordinate line represents one DEG. In a parallel coordinate line, connec-  
213 tions between samples with positive correlations should be flat, while connections  
214 between samples with negative correlations should be crossed. We expect DEGs  
215 to show more variability between treatments than between replicates. This means  
216 the parallel coordinate lines should be flat between replicates but crossed between  
217 treatments. However, overplotting problems would obscure our visualization if we  
218 were to plot all DEGs onto the same side-by-side boxplot. As a result, we used

219 hierarchical clustering techniques to separate DEGs into common patterns as is  
220 described in the methods section.

221 We see that the 1,019 DEGs from the Galbraith dataset form relatively clean-  
222 looking visual displays, with consistent replicates and differences between treat-  
223 ments (Figure 3). We do see that the first replicate of the virus group (V.1) appears  
224 somewhat inconsistent with the other virus replicates in Cluster 1, confirming that  
225 the trend we saw in the PCA plot carried through into the DEG results. Cluster  
226 4 reveals somewhat inconsistent replicates in the virus group, although most virus  
227 standardized read counts (group V) remain consistently larger than most control  
228 standardized read counts (group N). In contrast, we see that the 43 virus-related  
229 DEGs from our dataset do not look as clean in their visual displays (Figure 4). The  
230 replicates appear somewhat inconsistent in their estimated expression levels and  
231 there is not always such a large (or even consistent) difference between treatment  
232 groups. We see a similar finding when we also examine a larger subset of 1,914  
233 diet-related DEGs from our study (Additional file 3).

234 We next used repLicate TREatment (“litre”) plots, which we recently developed  
235 and published in our bigPint software package. Litre plots allow users to visualize  
236 one DEG onto the Cartesian coordinates of one scatterplot matrix. In the litre plot,  
237 each gene in the data is plotted once for every combination of replicates between  
238 treatment groups. For example, there are nine ways to pair a replicate from one  
239 treatment group with a replicate from the other treatment group in the Galbraith  
240 dataset (N.1 and V.1, N.1 and V.2, N.1 and V.3, N.2 and V.1, N.2. and V.2, N.2  
241 and V.3, N.3 and V.1, N.3 and V.2, and N.3 and V.3). Hence, each gene in the  
242 Galbraith dataset is plotted as nine points in the litre plot. With 11,825 genes in  
243 the Galbraith data, 106,425 points would need to be plotted. Our dataset is even  
244 more dramatic: There are 144 ways to pair a replicate from one treatment group  
245 with a replicate from the other treatment group, and with 15,314 genes in our data,

we would need to plot 2,205,216 points. For either dataset, plotting all these points would reduce the speed of interactive functionality and cause overplotting problems. As a result, we use hexagon bins to summarize this massive information. Once the background of hexagons has been drawn to reveal the distribution of all between-treatment sample pair combinations for *all* genes, the user can superimpose all between-treatment sample pair combinations for one gene of interest.

Additional file 4 shows nine example litre plots for our dataset. The hexagon background is the same for all nine litre plots because it simply shows the distribution of all between-treatment sample pair combinations for *all* genes in our dataset. In each litre plot, there are 144 magenta points superimposed that show all between-treatment sample pair combinations for one DEG of interest. Additional file 5 and 6 similarly each show nine example litre plots for the Galbraith dataset. We examined individual DEGs from the first cluster (Additional file 5) and second cluster (Additional file 6) of the Galbraith data because the first cluster had previously shown less consistency in the first replicate of the treatment group (Figure 3). Notice that, as previously explained, we now show each DEG as nine points for the Galbraith dataset. We see that indeed the virus DEGs from our data (Additional file 4) show less consistent replications and less differences between the treatment groups compared to the virus DEGs from the Galbraith data (Additional files 5 and 6). We also observe that, in the Galbraith dataset, the DEG points in the first cluster show less tight cluster patterns than the DEG points in the second cluster (Additional files 5 and 6), an observation we saw previously in the parallel coordinate plots (Figure 3).

Finally, we used scatterplot matrices from the bigPint software to further assess the DEGs. A scatterplot matrix is another effective multivariate visualization tool that plots read count distributions across all genes and samples. Specifically, it represents every gene in the dataset as a black point in each scatterplot. DEGs can

273 be superimposed as colored points to assess their patterns against the full dataset.  
274 We expect DEGs to mostly fall along the  $x=y$  line in replicate scatterplots (denot-  
275 ing replicate consistency) but deviate from the  $x=y$  line in treatment scatterplots  
276 (denoting significant treatment changes). The  $x=y$  line is shown in red in our plots.

277 We created standardized scatterplot matrices for each of the four clusters (from  
278 Figure 3) of the Galbraith data (Additional files 7, 8, 9, and 10). We also created  
279 standardized scatterplot matrices for our data. However, as our dataset contained  
280 24 samples, we would need to include 276 scatterplots in our matrix, which would  
281 be too numerous to allow for efficient visual assessment of the data. As a result,  
282 we created four scatterplot matrices of our data, each with subsets of 6 samples  
283 to be more comparable to the Galbraith data (Additional files 11, 12, 13, and 14).

284 We can again confirm through these plots that the DEGs from the Galbraith data  
285 appeared more as expected: They deviated more from the  $x=y$  line in the treatment  
286 scatterplots while staying close to the  $x=y$  line in replicate scatterplots.

287 Despite the virus-related DEGs ( $n = 1,019$ ) from the Galbraith dataset displaying  
288 the expected patterns more than those from our dataset ( $n = 43$ ), there was signif-  
289 icant overlap ( $p\text{-value} < 2.2\text{e-}16$ ) in the DEGs between the two studies, with 26/38  
290 (68%) of virus-upregulated DEGs from our study also showing virus-upregulated  
291 response in the Galbraith study (Figure 6).

## 292 Tolerance versus resistance

293 Using the contrasts specified in Table 2, we discovered 122 “tolerance” candi-  
294 date DEGs and 125 “resistance” candidate DEGs. Within our 122 “tolerance”  
295 gene ontologies, we found functions related to metabolism (such as carbohydrate  
296 metabolism, fructose metabolism, and chitin metabolism). However, we also discov-  
297 ered gene ontologies related to RNA polymerase II transcription, immune response,  
298 and regulation of response to reactive oxygen species (Figure 5A). Within our 125  
299 “resistance” gene ontologies, we found functions related to metabolism (such as car-

300 bohydrate metabolism, chitin metabolism, oligosaccharide biosynthesis, and general  
301 metabolism) (Figure 5B).

302 To visually explore gene expression patterns related to tolerance and resistance,  
303 we used hierarchical clustering to separate candidate DEGs into common patterns,  
304 and then visualized these clusters using parallel coordinate lines superimposed onto  
305 side-by-side boxplots. To reduce overplotting of parallel coordinate lines, we again  
306 used hierarchical clustering techniques to separate DEGs into common patterns.  
307 Perhaps unsurprisingly, we still see a substantial amount of noise (inconsistency  
308 between replicates) in our resulting candidate DEGs (Additional files 15 and 16).  
309 However, the broad patterns we expect to see still emerge: For example, based on  
310 the contrasts we created to obtain the ‘tolerance’ candidate DEGs, we expect them  
311 to display larger count values in the “NC” group compared to the “NR” group and  
312 larger count values in the “VC” group compared to the “VR” group. Indeed, we see  
313 this pattern in the associated parallel coordinate plots (Additional file 15). Likewise,  
314 based on the contrasts we created to obtain the ‘resistance’ candidate DEGs, we  
315 still expect them to display larger count values in the “VC” group compared to  
316 the “VR” group, but we no longer expect to see a difference between the “NC”  
317 and “NR” groups. We do generally see these expected patterns in the associated  
318 parallel coordinate plots: While there are large outliers in the “NC” group, the “NR”  
319 replicates are no longer typically below a standardized count of zero (Additional file  
320 16). The genes in Cluster 3 may follow the expected pattern the most distinctively  
321 (Additional file 16).

## 322 Post hoc analysis

323 To better understand sources of transcriptomic noise, we explored whether pathogen  
324 response measurements (virus titers and mortality), which varied widely across  
325 samples, were correlated with observed patterns in gene expression.

326 The R-squared values between gene read counts and pathogen response measure-  
327 ments were generally low ( $R\text{-squared} < 0.1$ ) across our dataset (Supplementary  
328 table 9, Additional file 1). We further explored whether clusters of DEGs showed  
329 higher correlations with pathogen response measurements than non-DEGs (the lat-  
330 ter serving as a control, where we do not expect a correlation). A Kruskal–Wallis  
331 test was used to determine if R-squared distributions of DEG clusters significantly  
332 differed from those in the rest of the data. The p-values and Bonferroni correction  
333 values for each of the 36 tests (as described in the methods section) is provided  
334 in Supplementary table 9, Additional file 1. An overall trend emerges to suggest  
335 that DEGs may have significantly larger correlation with the pathogen response  
336 measurements compared to non-DEGs.

## 337 Discussion

338 Challenges to honey bee health are a growing concern, in particular the combined,  
339 interactive effects of nutritional stress and pathogens [42]. In this study, we used  
340 RNA-sequencing to probe mechanisms underlying honey bee responses to two ef-  
341 fects, diet quality and infection with the prominent virus of concern, IAPV. In  
342 general, we found a major nutritional transcriptomic response, with nearly 2,000  
343 transcripts changing in response to diet quality (rockrose/poor diet versus chest-  
344 nut/good diet). The majority of these genes were upregulated in response to high  
345 quality diet, and these genes were over-represented for functions such as nutrient  
346 signaling metabolism (insulin resistance), immune response (Notch signaling and  
347 JaK-STAT pathways), and carbon metabolism (Supplementary table 2, Additional  
348 file 1). These data suggest high quality nutrition may allow bees to alter their  
349 metabolism, favoring investment of energy into innate immune responses.

350 One of the few studies that has investigated transcriptomic response to nutrition in  
351 honey bees similarly found that pollen upregulates genes related to macromolecule  
352 metabolism, insulin pathways, and TOR pathways [43]. Diet effects on transcrip-

353 tomics have been more extensively studied in the insect model *Drosophila*. One  
354 recent transcriptomic study in *Drosophila melanogaster* reported an overexpression  
355 of genes related to immunity, metabolism, and hemocyanin in a high-fat diet and  
356 overexpression of genes related to cell cycle activity, DNA binding and transcription,  
357 and CHK kinase-like protein activity in a high-sugar diet [45]. This same study also  
358 discovered an upregulation of genes related to peptide and carbohydrate processing  
359 in both high-fat and high-sugar diets, a finding the authors attributed to a general  
360 increase in caloric intake. Another recent study investigated the transcriptomic ef-  
361 fects of diets high in protein relative to sugar, diets high in sugar relative to protein,  
362 and diets with equal amounts of protein and sugar [46]. *Drosophila mojavensis* and  
363 *Drosophila arizonae* showed substantial differential expression between the dietary  
364 conditions: genes involved in carbohydrate and lipid metabolism were upregulated  
365 in response to high sugar low protein diets and genes involved in juvenile hormone  
366 (JH) and ecdysone were upregulated in response to low sugar high protein diets. In-  
367 terestingly, prior studies have suggested that JH regulates body size by controlling  
368 ecdysone production, which modifies insulin signaling [47]. As we saw in our study,  
369 these studies generally suggest that diet differences may relate to gene expression  
370 changes in metabolism and immune responses in honey bees.

371 While some insect systems have shown relatively low transcriptional responses  
372 to dicistrovirus infection [61, 62], previous work on honey bees has revealed many  
373 hundreds of DEGs [44]. Discrepancies between datasets may be due to noise and  
374 complexity of the honey bee microbiome. The transcriptomic response to virus infec-  
375 tion in our experiment was fairly limited. We found only 43 differentially expressed  
376 transcripts, some with known immune functions such as a gene with similarity to  
377 MD-2 lipid recognition protein and argonaute-2, a protein that plays a central role  
378 in RNA silencing (Table 1). We also found genes related to transcriptional regu-  
379 lation and muscle contraction. The small number of DEGs in this study may be

380 partly explained by the large amount of noise in the data (Figure 4 and Additional  
381 files 2B, 4, 11, 12, 13, and 14).

382 There have been numerous studies on the transcriptomic effects of virus infection  
383 in model organisms like fruit flies and mosquitoes that can provide a useful frame-  
384 work for interpreting virus responses in honey bees. These studies have showed that  
385 RNA silencing is a major antiviral strategy, along with transcriptional pausing, Toll  
386 pathways, IMD pathways, JAK/STAT pathways, and Toll-7-autophagy pathways  
387 [48, 49]. Recent transcriptomic studies in honey bees have shown similar hallmarks  
388 of these same antiviral defense mechanisms, including RNA silencing, Toll path-  
389 ways, IMD pathways, JAK/STAT pathways, autophagy, and endocytosis [50]. It is  
390 important to note that general immune responses to viral infection in insects might  
391 be an indirect result of cellular damage [49]. In fact, every virus-host interaction has  
392 its own particularities derived from the diverse methods of replication and infection  
393 cycle evolved by different viruses. An intricate set of pro- and anti-virus host factors  
394 such as ribosomal proteins and autophagy pathways are involved, but the response  
395 depends on the virus species, as has been elucidated in *Drosophila* [48, 49]. In ad-  
396 dition, a non-sequence-specific antiviral response mediated by unspecific dsRNA  
397 pathway was discovered in honey bees [63, 64]. In the case of dicistroviruses, few  
398 works have studied the impact of IAPV infection at transcriptional level. Chen  
399 et al. 2014 analyzed responses to IAPV infection in larvae and workers using mi-  
400 croarrays [65]. Many of the DEGs found were involved in immune response and  
401 energy-related metabolism, particularly in adults but not in brood. The authors  
402 propose this observed difference could be connected to latent infections in larvae  
403 (where host immunity is not perturbed) versus acute infections in adulthood (in-  
404 duced by stressors faced during development) [65]. IAPV acute infection also alters  
405 the DNA methylation pattern of numerous genes that do not overlap the genes that  
406 are up- or down-regulated at the transcriptional level [44]. These works reiterate the



407 conclusion that viruses trigger particular antiviral mechanisms by different means  
408 and depending on several factors. The honey bee antiviral pathways induced by  
409 specific viruses were recently reviewed [50]; it is noteworthy that many honey bee  
410 factors discovered by transcriptomics need further characterization to uncover their  
411 role in controlling (or promoting) viral infection in honey bees.

412 Given the noisy nature of our data, and our desire to hone in on genes with real  
413 expression differences, we compared our data to the Galbraith study [44], which  
414 also examined bees response to IAPV infection. In contrast to our study, Galbraith  
415 et al. identified a large number of virus responsive transcripts, and generally had  
416 less noise in their data (Figure 3 and Additional files 2A, 5, 6, 7, 8, 9, and 10). To  
417 identify the most consistent virus-responsive genes from our study, we looked for  
418 overlap in the DEGs associated with virus infection on both experiments. We found  
419 a large, statistically significant ( $p\text{-value} < 2.2\text{e-}16$ ) overlap, with 26/38 (68%) of  
420 virus-responsive DEGs from our study also showing response to virus infection in  
421 Galbraith et al. (Figure 6). This result gives us confidence that, although noisy, we  
422 were able to uncover reliable, replicable gene expression responses to virus infection  
423 with our data.

424 Data visualization is a useful method to identify noise and robustness in RNA-  
425 sequencing data [66]. In this study, we used extensive data visualization to improve  
426 the interpretation of our RNA-sequencing results. For example, the DESeq2 pack-  
427 age comes with certain visualization options that are popular in RNA-sequencing  
428 analysis. One of these visualization is the principal component analysis (PCA) plot,  
429 which allows users to visualize the similarity between samples within a dataset. We  
430 could determine from this plot that indeed the Galbraith data may show more simi-  
431 larity between its replicates and differences between its treatments compared to our  
432 data (Additional file 2). However, the PCA plot only shows us information at the  
433 sample level. We wanted to investigate how these differences in the signal:to:noise

ratios of the datasets would affect the structure of any resulting DEGs. As a result, we also used three plotting techniques from the *bigPint* package: We investigated the 1,019 virus-related DEGs from the Galbraith dataset and the 43 virus-related DEGs from our dataset using parallel coordinate lines, scatterplot matrices, and litre plots. To prevent overplotting issues in our graphics, we used a hierarchical clustering technique for the parallel coordinate lines to separate the set of DEGs into smaller groups. We also needed to examine four subsets of samples from our dataset to make effective use of the scatterplot matrices. After these tailorizations, we determined that the same patterns we saw in the PCA plots regarding the entire dataset extended down the pipeline analysis into the DEG calls: Even the DEGs from the Galbraith dataset showed more similarity between their replicates and differences between their treatments compared to those from our data. However, the 365 DEGs from the Galbraith data in Cluster 1 of Figure 3 showed an inconsistent first replicate in the treatment group (“V.1”), which was something we observed in the PCA plot. This indicates that this feature also extended down the analysis pipeline into DEG calls. Despite the differences in signal between these two datasets, there was substantial overlap in the resulting DEGs. We believe these visualization applications can be useful for future researchers analyzing RNA-sequencing data to quickly and effectively ensure that the DEG calls look reliable or at least overlap with DEG calls from similar studies that look reliable. We also expect this type of visualization exploration can be especially crucial when studying wild populations with high levels of genetic and environmental variation between replicates and/or when using experiments that may lack rigid design control.

One of the goals of this study was to use our RNA-sequencing data to assess whether transcriptomic responses to diet quality and virus infection provide insight into whether high quality diet can buffer bees from pathogen stress via mechanisms of “resistance” or “tolerance”. Recent evidence has suggested that overall immu-

nity is determined by more than just “resistance” (the reduction of pathogen fitness within the host by mechanisms of avoidance and control) [67]. Instead, overall immunity is related to “resistance” in conjunction with “tolerance” (the reduction of adverse effects and disease resulting from pathogens by mechanisms of healing) [41, 67]. Immune-mediated resistance and diet-driven tolerance mechanisms are costly and may compete with each other [41, 68]. Data and models have suggested that selection can favor an optimum combination of both resistance and tolerance [69, 70, 71, 72]. We attempted to address this topic through specific gene expression contrasts (Table 2), accompanied by GO analysis of the associated gene lists. We found an approximately equal number of resistance ( $n = 125$ ) and tolerance ( $n = 122$ ) related candidate DEGs, suggesting both processes may be playing significant roles in dietary buffering from pathogen induced mortality. Resistance candidate DEGs had functions related to several forms of metabolism (chitin and carbohydrate), regulation of transcription, and cell adhesion (Figure 5B). Tolerance candidate DEGs had functions related to carbohydrate metabolism and chitin metabolism; however, they also showed functions related to immune response, including RNA polymerase II transcription (Figure 5A). Previous studies have shown that transcriptional pausing of RNA polymerase II may be an innate immune response in *D. melanogaster* that allows for a more rapid response by increasing the accessibility of promoter regions of virally induced genes [73]. These possible immunological defense mechanisms within our “tolerance” candidate DEGs and metabolic processes within our “resistance” candidate DEGs may provide additional evidence of feedbacks between diet and disease in honey bees [42].

There were several limitations in this study that could be improved upon in future studies. For instance, our comparison between the Galbraith data (single-drone colonies) and our data (naturally mated colonies) was limited by numerous extraneous variables between these studies. In addition to different molecular pipelines

488 and bioinformatic preprocessing pipelines used between these studies, the Galbraith  
489 study focused on one-day old worker honey bees that were fed sugar and artificial  
490 pollen diet, whereas our study focused on adult worker honey bees that were fed  
491 bee-collected monofloral diets. Furthermore, the Galbraith data used eviscerated  
492 abdomens with attached fat bodies and only considered symptomatic honey bees  
493 for their infected treatment group, whereas we used whole bodies and considered  
494 both asymptomatic and symptomatic honey bees for our infected treatment group.  
495 There are also differences in the hours post inoculation and possible differences in  
496 the inoculation amount between the studies. Further differences between the studies  
497 can be found in their corresponding published methods sections [40, 44]. The dif-  
498 ferent factors between these two studies may be critical because particular antiviral  
499 factors in honey bees are linked to specific viruses, specific developmental stages,  
500 the analyzed tissue, the route of inoculation, and the time (post-inoculation) dur-  
501 ing which the study was performed. This was clearly demonstrated when comparing  
502 honey bee responses to two related iflaviruses with very different infection dynamics,  
503 sacbrood bee virus (SBV) vs. deformed wing virus (DWV) [74]. Authors observed  
504 differences in induction of defensin and hymenoptaecin immune-related genes, and  
505 suggested the results reflect adaptations to the different routes of transmission [74].

506 Moreover, our comparative visualization assessment between these two datasets  
507 was also somewhat limited because the virus effect in the Galbraith study used  
508 three replicates for each level, whereas the virus effect in our study used twelve  
509 replicates for each level that were actually further subdivided into six replicates for  
510 each diet level. Hence the apparent reduction in noise observed in the Galbraith  
511 data compared to our data in the PCA plots, parallel coordinate plots, scatterplot  
512 matrices, and litre plots may be an inadvertent product of the smaller number of  
513 replicates used and the lack of a secondary treatment group rather than solely the  
514 reduction in genetic variability through the single-drone colony design itself. With

515 this in mind, while our current efforts may be a starting point, future studies can  
516 shed more light on signal:to:noise and differential expression differences between  
517 naturally mated colony designs and single-drone colony designs by controlling for  
518 extraneous factors more strictly than what we were able to do in the current line  
519 of work.

520 In addition, this study used a whole body RNA-sequencing approach. In future  
521 related studies, it may be informative to use tissue-specific methods. Previous work  
522 has shown that even though IAPV replication occurs in all honey bee tissues, it  
523 localizes more in gut and nerve tissues and in the hypopharyngeal glands. Likewise,  
524 the highest IAPV titers have been observed in gut tissues [34]. Recent evidence has  
525 suggested that RNA-sequencing approaches toward composite structures in honey  
526 bees leads to false negatives, implying that genes strongly differentially expressed  
527 in particular structures may not reach significance within the composite structure  
528 [75]. These studies have also found that within a composite extraction, structures  
529 therein may contain opposite patterns of differential expression. We can provide  
530 more detailed answers to our original transcriptomic questions if we were to repeat  
531 this same experimental design only now at a more refined tissue level. Another  
532 future direction related to this work would be to integrate multiple omics datasets  
533 to investigate monofloral diet quality and IAPV infection in honey bees. Indeed,  
534 previous studies in honey bees have found that multiple omics datasets do not  
535 always align in a clear-cut manner, and hence may broaden our understanding of  
536 the molecular mechanisms being explored [44].

## 537 **Conclusions**

538 To the best of our knowledge, there are few to no studies investigating honey bee  
539 gene expression specifically related to monofloral diets, and few to no studies ex-  
540 amining honey bee gene expression related to the combined effects of diet in any  
541 general sense and viral inoculation in any general sense. It also remains unknown

whether the protective effects of good diet in honey bees is due to direct effects on immune function (resistance) or indirect effects of energy availability on vigor and health (tolerance). We attempted to address these unresolved areas by conducting a two-factor RNA-sequencing study that examined how monofloral diets and IAPV inoculation influence gene expression patterns in honey bees. Overall, our data suggest complex transcriptomic responses to multiple stressors in honey bees. Diet has the capacity for large and profound effects on gene expression and may set up the potential for both resistance and tolerance to viral infection, adding to previous evidence of possible feedbacks between diet and disease in honey bees [42].

Moreover, this study also demonstrated the benefits of using data visualizations and multiple datasets to address inherently messy biological data. For instance, by verifying the substantial overlap in our DEG lists to those obtained in another study that addressed a similar question using specimens with less genetic variability, we were able to place much higher confidence in the differential gene expression results from our otherwise noisy data. We also suggested that comparing results derived from multiple studies varying in level of genetic and environmental variability may allow researchers to identify transcriptomic patterns that are concurrently more realistic and less noisy. Altogether, we hope our results underline the merits of using data visualization techniques and multiple datasets to understand and interpret RNA-sequencing datasets.

## Methods

### Mortality and virus titers

Details of the procedures we used to prepare virus inoculum, infect and feed caged honey bees, and quantify IAPV can be reviewed in our previous work [40, 33]. A linear mixed effects model was used to relate the mortality rates and IAPV titers to the main and interaction effects of the diet and virus factors. The model was fitted to the data by restricted maximum likelihood (REML) using the “lme” function

in the R package “nlme”. A random (intercept) effect for experimental setup was included in the model. Post-hoc pairwise comparisons of the four (diet and virus combination) treatment groups were performed and Benjamini-Hochberg adjusted p-values were calculated to limit familywise Type I error rates [76].

### Design of two-factor experiment

For our nutrition factor, we examined two monofloral pollen diets, rockrose (*Cistus* sp.) and chestnut (*Castanea* sp.). Rockrose pollen is generally considered less nutritious than chestnut pollen due to its lower levels of protein, amino acids, antioxidants, calcium, and iron [40, 51]. For our virus factor, one level contained bees that were infected with IAPV and another level contained bees that were not infected with IAPV. This experimental design resulted in four treatment groups (rockrose pollen without IAPV exposure, chestnut pollen without IAPV exposure, rockrose pollen with IAPV exposure, and chestnut pollen with IAPV exposure) that allowed us to assess main effects and interactive effects between diet quality and IAPV infection in honey bees.

There are several reasons why our design focused only on diet quality (monofloral diets) as opposed to diet diversity (monofloral diets versus polyfloral diets). First, when assessing diet diversity, a sugar diet is often used as a control. However, such an experimental design does not reflect real-world conditions for honey bees as they rarely face a total lack of pollen [51]. Second, in studies that compared honey bee health using monofloral and polyfloral diets at the same time, if the polyfloral diet and one of the high-quality monofloral diets both exhibited similarly beneficial effects, then it was difficult for the authors to assess if the polyfloral diet was better than most of the monofloral diets because of its diversity or because it contained as a subset the high-quality monofloral diet [51]. Third, as was previously mentioned, honey bees are now confronted with less diverse sources of pollen. As a result, there is a need to better understand how monofloral diets affect honey bee health.

## 596 RNA extraction

597 Fifteen cages per treatment were originally produced for monitoring of mortality.  
598 From these, six live honey bees were randomly selected from each cage 36 hours  
599 post inoculation and placed into tubes [33]. Tubes were kept on dry ice and then  
600 transferred into a -80C freezer until processing. From the fifteen possible cages,  
601 eight were randomly selected for RNA-sequencing. From these eight cages, two of  
602 the honey bees per cage were randomly selected from the original six live honey  
603 bees per cage. These two bees were combined to form a pooled sample representing  
604 the cage. Whole body RNA from each pool was extracted using Qiagen RNeasy  
605 MiniKit followed by Qiagen DNase treatment. Samples were suspended in water to  
606 200-400 ng/ $\mu$ l. All samples were then tested on a Bioanalyzer at the Iowa State  
607 University DNA Facility to ensure quality (RIN > 8).

## 608 Gene expression

609 Samples were sequenced starting on January 14, 2016 at the Iowa State University  
610 DNA Facility (Platform: Illumina HiSeq Sequencing 2500 in rapid run mode; Cat-  
611 egory: Single End 100 cycle sequencing). A standard Illumina mRNA library was  
612 prepared by the DNA facility. Reads were aligned to the BeeBase Version 3.2 genome  
613 [77] from the Hymenoptera Genome Database [78] using the programs GMAP and  
614 GSNAP [79]. There were four lanes of sequencing with 24 samples per lane. Each  
615 sample was run twice. Approximately 75-90% of reads were mapped to the honey  
616 bee genome. Each lane produced around 13 million single-end 100 basepair reads.

617 We tested all six pairwise combinations of treatments for DEGs (pairwise DEGs).  
618 We also tested the diet main effect (diet DEGs), virus main effect (virus DEGs), and  
619 interaction term for DEGs (interaction DEGs). We then also tested for virus main  
620 effect DEGs (virus DEGs) in public data derived from a previous study exploring  
621 the gene expression of IAPV virus infection in honey bees [44]. We tested each  
622 DEG analysis using recommended parameters with DESeq2 [80], edgeR [66], and



623 LimmaVoom [81]. In all cases, we used a false discovery rate (FDR) threshold of 0.05  
624 [82]. Fisher's exact test was used to determine significant overlaps between DEG  
625 sets (whether from the same dataset but across different analysis pipelines or from  
626 different datasets across the same analysis pipelines). The eulerr shiny application  
627 was used to construct Venn diagram overlap images [83]. In the end, we focused on  
628 the DEG results from DESeq2 [80] as this pipeline was also used in the Galbraith  
629 study [44]. We used the independent filtering process built into the DESeq2 software  
630 that mitigates multiple comparison corrections on genes with no power rather than  
631 defining one filtering threshold.

#### 632 Comparison to prior studies on transcriptomic response to viral infection

633 We compare the main effect of IAPV exposure in our dataset to that obtained in a  
634 previous study conducted by Galbraith and colleagues [44] who also addressed honey  
635 bee transcriptomic responses to virus infection. We applied the same downstream  
636 bioinformatics analyses between our count table and the count table provided in  
637 the Galbraith study. When we applied our bioinformatics pipeline to the Galbraith  
638 count table, we obtained different differential expression counts compared to the  
639 results published in the Galbraith study. However, there was substantial overlap and  
640 we considered this justification to use the differential expression list we obtained in  
641 order to keep the downstream bioinformatics analyses as similar as possible between  
642 the two datasets (Additional file 17).

643 We used honey bees from naturally mated colonies, whereas Galbraith et al. [44]  
644 used honey bees from single-drone colonies. In light of this, we should expect the  
645 Galbraith et al. dataset to contain lower genetic variation between its replicates  
646 and higher signal:to:noise ratios than our dataset. We use visualization techniques  
647 to assess the signal:to:noise ratio between these two datasets, and differential gene  
648 expression (DEG) analyses to determine any significantly overlapping genes of in-  
649 terest between these two datasets.

## 650 Visualization

651 We used an array of visualization tools as part of our analysis. We used the PCA plot  
652 [84] from the DESeq2 package, a well-known and established tool. Along with that,  
653 we used lesser-known multivariate visualization tools from our work-in-progress R  
654 package called bigPint. Specifically, we used parallel coordinate plots [85], scatter-  
655 plot matrices [86], and litre plots (which we recently developed based on “replicate  
656 line plots” [87]) to assess the variability between the replicates and the treatments  
657 in our data. We also used these plotting techniques to assess for normalization  
658 problems and other common problems in RNA-sequencing analysis pipelines [87].

659 Furthermore, we used statistical graphics to better understand patterns in our  
660 DEGs. However, in cases of large DEG lists, these visualization tools had overplot-  
661 ting problems (where multiple objects are drawn on top of one another, making  
662 it impossible to detect individual values). To remedy this problem, we first stan-  
663 dardized each DEG to have a mean of zero and standard deviation of unity [88, 89].  
664 Then, we performed hierarchical clustering on the standardized DEGs using Ward’s  
665 linkage. This process divided large DEG lists into smaller clusters of similar pat-  
666 terns, which allowed us to more efficiently visualize the different types of patterns  
667 within large DEG lists (see Figures 3 and 4 for examples).

## 668 Gene ontology

669 DEGs were uploaded as a background list to DAVID Bioinformatics Resources 6.7  
670 [90, 91]. The overrepresented gene ontology (GO) terms of DEGs were determined  
671 using the BEEBASE.ID identifier option (honey bee gene model) in the DAVID  
672 software. To fine-tune the GO term list, only terms correlating to Biological Pro-  
673 cesses were considered. The refined GO term list was then imported into REVIGO  
674 [92], which uses semantic similarity measures to cluster long lists of GO terms.

## 675 Probing tolerance versus resistance

676 To investigate whether the protective effect of good diet is due to direct, specific  
677 effects on immune function (resistance), or if it is due to indirect effects of good nu-  
678 trition on energy availability and vigor (tolerance), we created contrasts of interest  
679 (Table 2). In particular, we assigned “resistance candidate DEGs” to be the ones  
680 that were upregulated in the chestnut group within the virus infected bees but not  
681 upregulated in the chestnut group within the non-infected bees. Our interpretation  
682 of these genes is that they represent those that are only activated in infected bees  
683 that are fed a high quality diet. We also assigned “tolerance candidate DEGs” to  
684 be the ones that were upregulated in the chestnut group for both the virus infected  
685 bees and non-infected bees. Our interpretation of these genes is that they represent  
686 those that are constitutively activated in bees fed a high quality diet, regardless  
687 of whether they are experiencing infection or not. We then determined how many  
688 genes fell into these two categories and analyzed their GO terminologies.

## 689 Post hoc analysis

690 We found considerable noisiness in our data and saw, through gene-level visual-  
691 izations, that our DEGs contained outliers and inconsistent replicates. Hence, we  
692 wanted to explore whether our DEG read counts correlated with pathogen response  
693 metrics, including IAPV titers, sacbrood bee virus (SBV) titers, and mortality rates.  
694 For this process, we considered virus main effect DEGs (Figure 4), “tolerance can-  
695 didate” DEGs (Additional file 15), and “resistance candidate” DEGs (Additional  
696 file 16). For each DEG in each cluster, we calculated a coefficient of determination  
697 (R-squared) value to estimate the correlation between its raw read counts and the  
698 pathogen response metrics across its 24 samples. We then used the Kruskal–Wallis  
699 test to determine if the distribution of the R-squared values in any of the DEG clus-  
700 ters significantly differed from those in the non-DEG genes (the rest of the data).  
701 As there were four clusters for each of the nine combinations of DEG lists (“tol-

erance” candidate DEGs, “resistance” candidate DEGs, and virus-related DEGs) and pathogen response measurements (IAPV titer, SBV titer, and mortality rate), this process resulted in 36 statistical tests.

#### **Ethics approval and consent to participate**

All honey bees used in this work were sampled in the United States, and no ethical use approval is required for this species in this country.

#### **Consent for publication**

Not applicable.

#### **Availability of data and materials**

The data discussed in this publication have been deposited in NCBI's Gene Expression Omnibus [93] and are accessible through GEO Series accession number GSE121885 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE121885>). The scripts to reproduce analyses and figures in this publication are available online (<https://github.com/lrutter/HoneyBeePaper>).

#### **Competing interests**

The authors declare that they have no competing interests.

#### **Funding**

This work was supported by the United States Department of Agriculture, Agriculture and Food Research Initiative (USDA-AFRI) 2011-04894.

#### **Author's contributions**

LR performed the bioinformatic and statistical analyses, produced the figures and tables, and drafted the manuscript. BB conceptualized the study and critically revised the manuscript. AD contributed to experimental design, carried out the laboratory experiments, and processed samples for virus titers and RNA-seq. DC advised on statistical analyses and visualization.

#### **Acknowledgements**

We would like to thank Giselle Narvaez for assisting with cage experiments.

#### **Author details**

<sup>1</sup>Bioinformatics and Computational Biology Program, Iowa State University, Ames, IA 50011, USA. <sup>2</sup>Econometrics and Business Statistics, Monash University, Clayton, VIC 3800, Australia. <sup>3</sup>Department of Entomology, Iowa State University, Ames, IA 50011, USA. <sup>4</sup>Department of Ecology, Evolution, and Organismal Biology, Iowa State University, Ames, IA 50011, USA. <sup>5</sup>Department of Entomology, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA. <sup>6</sup>Department of Entomology and Nematology, University of Florida, Gainesville, FL 32611, USA.

#### **References**

- van Engelsdorp, D., Evans, J.D., Saegerman, C., Mullin, C., Haubruge, E., Nguyen, B.K., Frazier, M., Frazier, J., Cox-Foster, D., Chen, Y., Underwood, R., Tarry, D.R., Pettis, J.S.: Colony collapse disorder: A descriptive study. *PLoS ONE* **4**, 6481 (2009)
- Kulhanek, K., Steinhauer, N., Rennich, K., Caron, D.M., Sagili, R.R., Pettis, J.S., Ellis, J.D., Wilson, M.E., Wilkes, J.T., Tarry, D.R., Rose, R., Lee, K., Rangel, J., vanEngelsdorp, D.: A national survey of managed honey bee 2014–2015 annual colony losses in the USA. *Journal of Apicultural Research* **56**, 328–340 (2017)
- Laurent, M., Hendrikx, P., Ribiere-Chabert, M., Chauzat, M.-P.: A pan-European epidemiological study on honeybee colony losses 2012–2014. *Epilobee* **2013**, 44 (2016)
- Caron, D., Sagili, R.: Honey bee colony mortality in the Pacific Northwest: Winter 2009/2010. *Am Bee J* **151**, 73–76 (2011)
- Bond, J., Plattner, K., Hunt, K.: Fruit and Tree Nuts Outlook: Economic Insight U.S. Pollination- Services Market. Economic Research Service Situation and Outlook FTS-357SA, USDA (2014)
- Gallai, N., Salles, J.-M., Settele, J., Vaissière, B.B.: Economic valuation of the vulnerability of world agriculture confronted with pollinator decline. *Ecol. Econ.* **68**, 810–821 (2009)
- Klein, A.-M., Vaissière, B.E., Cane, J.H., Steffan-Dewenter, I., Cunningham, S.A., Kremen, C., Tscharntke, T.: Importance of pollinators in changing landscapes for world crops. *Proc Biol Sci* **274**, 303–313 (2007)
- Potts, S.G., Biesmeijer, J.C., Kremen, C., Neumann, P., Schweiger, O., Kunin, W.E.: Global pollinator declines: trends, impacts and drivers **25**, 345–353 (2010)
- Spivak, M., Mader, E., Vaughan, M., Euliss, N.H.: The Plight of the Bees. *Environ Sci Technol* **45**, 34–38 (2011)
- Goulson, D., Nicholls, E., Botías, C., Rotheray, E.L.: Bee declines driven by combined stress from parasites, pesticides, and lack of flowers. *Science* **347**, 1255957 (2015)
- Roulston, T.H., Buchmann, S.L.: A phylogenetic reconsideration of the pollen starch-pollination correlation. *Evol Ecol Res* **2**, 627–643 (2000)
- Stanley, R.G., Linsens, H.F.: Pollen: Biology, Biochemistry, Management
- Brodtschneider, R., Crailsheim, K.: Nutrition and health in honey bees. *Apidologie* **41**, 278–294 (2010)
- Haydak, M.H.: Honey bee nutrition. *Annu Rev Entomol* **15**, 143–156 (1970)

15. Crailsheim, K., Schneider, L.H.W., Hrassnigg, N., Bühlmann, G., Brosch, U., Gmeinbauer, R., Schöffmann, B.: Pollen consumption and utilization in worker honeybees (*Apis mellifera carnica*): dependence on individual age and function. *J Insect Physiol* **38**, 409–419 (1992)
16. Crailsheim, K.: The flow of jelly within a honeybee colony. *J Comp Physiol B* **162**, 681–689 (1992)
17. Schmidt, J.O.: Feeding preference of *Apis mellifera* L. (Hymenoptera: Apidae): Individual versus mixed pollen species. *J. Kans. Entomol. Soc.* **57**, 323–327 (1984)
18. Schmidt, J.O., Thoenes, S.C., Levin, M.D.: Survival of honey bees, *Apis mellifera* (Hymenoptera: Apidae), fed various pollen sources. *J. Econ. Entomol.* **80**, 176–183 (1987)
19. Alaux, C., Ducloz, F., Conte, D.C.Y.L.: Diet effects on honeybee immunocompetence. *Biol. Lett.* **6**, 562–565 (2010)
20. Naug, D.: Nutritional stress due to habitat loss may explain recent honeybee colony collapses. *Biol Conserv* **142**, 2369–2372 (2009)
21. Engelsdorp, D.V., Hayes, J.J., Underwood, R.M., Pettis, J.: A survey of honey bee colony losses in the U.S., fall 2007 to spring 2008. *PLoS ONE* **3**, 4071 (2008)
22. Neumann, P., Carreck, N.L.: Honey bee colony losses. *J Apicult Res* **49**, 1–6 (2010)
23. Engelsdorp, D.V., Meixner, M.D.: A historical review of managed honey bee populations in Europe and the United States and the factors that may affect them. *J Invertebr Pathol* **103**, 80–95 (2010)
24. Rosenkranz, P., Aumeier, P., Ziegelmann, B.: Biology and control of *Varroa destructor*. *J Invertebr Pathol* **103**, 96–119 (2010)
25. Weinberg, K.P., Madel, G.: The influence of the mite *Varroa Jacobsoni* Oud. on the protein concentration and the haemolymph volume of the brood of worker bees and drones of the honey bee *Apis Mellifera* L. *Apidologie* **16**, 421–436 (1985)
26. Shen, M.Q., Cui, L.W., Ostiguy, N., Cox-Foster, D.: Intricate transmission routes and interactions between picorna-like viruses (Kashmir bee virus and sacbrood virus) with the honeybee host and the parasitic varroa mite. *J Gen Virol* **86**, 2281–2289 (2005)
27. Yang, X., Cox-Foster, D.: Effects of parasitization by *Varroa destructor* on survivorship and physiological traits of *Apis mellifera* in correlation with viral incidence and microbial challenge. *Parasitology* **134**, 405–412 (2007)
28. Yang, X.L., Cox-Foster, D.L.: Impact of an ectoparasite on the immunity and pathology of an invertebrate: Evidence for host immunosuppression and viral amplification. *P Natl Acad Sci USA* **102**, 7470–7475 (2005)
29. Emsen, B., Hamiduzzaman, M.M., Goodwin, P.H., Guzman-Novoa, E.: Lower virus infections in *Varroa destructor*-infested and uninfested brood and adult honey bees (*Apis mellifera*) of a low mite population growth colony compared to a high mite population growth colony. *PLoS ONE* **10**, 0118885 (2015)
30. Chen, Y.P., Siede, R.: Honey bee viruses. *Adv Virus Res* **70**, 33–80 (2007)
31. Bonning, B.C., Miller, W.A.: Dicistroviruses. *Annu Rev Entomol* **55**, 129–150 (2010)
32. Maori, E., Paldi, N., Shafir, S., Kalev, H., Tsur, E., Glick, E., Sela, I.: IAPV, a bee-affecting virus associated with Colony Collapse Disorder can be silenced by dsRNA ingestion. *Insect Mol Biol* **18**, 55–60 (2009)
33. Carrillo-Tripp, J., Dolezal, A.G., Goblirsch, M.J., Miller, W.A., Toth, A.L., Bonning, B.C.: In vivo and in vitro infection dynamics of honey bee viruses. *Sci Rep* **6**, 22265 (2016)
34. Chen, Y.P., Pettis, J.S., Corona, M., Chen, W.P., Li, C.J., Spivak, M., Visscher, P.K., DeGrandi-Hoffman, G., Boncristiani, H., Zhao, Y., van Engelsdorp, D., Delaplane, K., Solter, L., Drummond, F., Kramer, M., Lipkin, W.I., Palacios, G., Hamilton, M.C., Smith, B., Huang, S.K., Zheng, H.Q., Li, J.L., Zhang, X., Zhou, X.F., Wu, L.Y., Zhou, J.Z., Lee, M.-L., Teixeira, E.W., Li, Z.G., Evans, J.D.: Israeli acute paralysis virus: Epidemiology, pathogenesis and implications for honey bee health. *PLoS Pathog* **10**, 1004261 (2014)
35. DeGrandi-Hoffman, G., Chen, Y.: Nutrition, immunity and viral infections in honey bees. *Current Opinion in Insect Science* **10**, 170–176 (2015)
36. DeGrandi-Hoffman, G., Chen, Y., Huang, E., Huang, M.H.: The effect of diet on protein concentration, hypopharyngeal gland development and virus load in worker honey bees (*Apis mellifera* L.). *J Insect Physiol* **56**, 1184–1191 (2010)
37. Le Conte, Y., BRUNET, J.-L., McDonnell, C., Dussaubat, C., Alaux, C.: Interactions Between Risk Factors in Honey Bees
38. Annoscia, D., Zanni, V., Galbraith, D., Quirici, A., Grozinger, C., Bortolomeazzi, R., Nazzi, F.: Elucidating the mechanisms underlying the beneficial health effects of dietary pollen on honey bees (*Apis mellifera*) infested by *Varroa* mite ectoparasites. *Scientific Reports* **7**, 6258 (2017)
39. Nazzi, F., Pennacchio, F.: Honey bee antiviral immune barriers as affected by multiple stress factors: A novel paradigm to interpret colony health decline and collapse. *Viruses* **10**, 159 (2018)
40. Dolezal, A.G., Carrillo-Tripp, J., Judd, T., Miller, A., Bonning, B., Toth, A.: Interacting stressors matter: Diet quality and virus infection in honey bee health. In prep (2018)
41. Miller, C.V.L., Cotter, S.C.: Resistance and tolerance: The role of nutrients on pathogen dynamics and infection outcomes in an insect host. *Journal of Animal Ecology* **87**, 500–510 (2017)
42. Dolezal, A.G., Toth, A.L.: Feedbacks between nutrition and disease in honey bee health. *Current Opinion in Insect Science* **26**, 114–119 (2018)
43. Alaux, C., Dantec, C., Parrinello, H., Conte, Y.L.: Nutrigenomics in honey bees: digital gene expression analysis of pollen's nutritive effects on healthy and varroa-parasitized bees. *BMC Genomics* **12**, 496 (2011)
44. Galbraith, D.A., Yang, X., Niño, E.L., Yi, S., Grozinger, C.: Parallel epigenomic and transcriptomic responses to viral infection in honey bees (*Apis mellifera*). *PLoS Pathogens* **11**, 1004713 (2015)
45. Hemphill, W., Rivera, O., Talbert, M.: RNA-Sequencing of *Drosophila melanogaster* head tissue on high-sugar and high-fat diets. *G3: Genes, Genomes, Genetics* **8**, 279–290 (2018)
46. Nazario-Yepiz, N.O., Loustalot-Laclette, M.R., Carpinteyro-Ponce, J., Abreu-Goodger, C., Markow, T.A.: Transcriptional responses of ecologically diverse *Drosophila* species to larval diets differing in relative sugar and protein ratios. *PLoS ONE* **12**, 0183007 (2017)
47. Mirth, C.K., Tang, H.Y., Makohon-Moore, S.C., Salhadar, S., Gokhale, R.H., Warner, R.D., Koyama, T., Riddiford, L.M., Shingleton, A.W.: Juvenile hormone regulates body size and perturbs insulin signaling in

- Drosophila*. Proceedings of the National Academy of Sciences **25**, 201313058 (2014)
48. Xu, J., Cherry, S.: Viruses and antiviral immunity in *Drosophila*. *Dev Comp Immunol* **42**, 67–84 (2014)
  49. Swevers, L., Liu, J., Smagghe, G.: Defense Mechanisms against Viral Infection in *Drosophila*: RNAi and Non-RNAi. *Viruses* **10**, 230 (2018)
  50. McMenamin, A.J., Daughenbaugh, K.F., Parekh, F., Pizzorno, M.C., Flenniken, M.L.: Honey Bee and Bumble Bee Antiviral Defense. *Viruses* **10**, 395 (2018)
  51. Pasquale, G.D., Salignon, M., Conte, Y.L., Belzunces, L.P., Decourtye, A., Kretzschmar, A., Suchail, S., Brunet, J.-L., Alaux, C.: Influence of pollen nutrition on honey bee health: Do pollen quality and diversity matter? *PLoS ONE* **8**, 72016 (2013)
  52. Page, R.E., Laidlaw, H.H.: Full sisters and supersisters: A terminological paradigm. *Anim. Behav.* **36**, 944–945 (1988)
  53. Sherman, P.W., Seeley, T.D., Reeve, H.K.: Parasites, pathogens, and polyandry in social Hymenoptera. *Am. Nat* **131**, 602–610 (1988)
  54. Crozier, R.H., Page, R.E.: On being the right size: Male contributions and multiple mating in social Hymenoptera. *Behav. Ecol. Sociobiol.* **18**, 105–115 (1985)
  55. Mattila, H.R., Seeley, T.D.: Genetic diversity in honey bee colonies enhances productivity and fitness. *Science* **317**, 362–364 (2007)
  56. Tarpy, D.R.: Genetic diversity within honeybee colonies prevents severe infections and promotes colony growth. *Proc. R. Soc. Lond. B* **270**, 99–103 (2003)
  57. Brodschneider, R., Arnold, G., Hrassnigg, N., Crailsheim, K.: Does patriline composition change over a honey bee queen's lifetime? *Insects* **3**, 857–869 (2012)
  58. Hansen, K.D., Brenner, S.E., Dudoit, S.: Biases in Illumina transcriptome sequencing caused by random hexamer priming. *Nucleic Acids Research* **38**, 131 (2010)
  59. Oshlack, A., Robinson, M.D., Young, M.D.: From RNA-seq reads to differential expression results. *Genome Biology* **11**, 220 (2010)
  60. McIntyre, L.M., Lopiano, K.K., Morse, A.M., Amin, V., Oberg, A.L., Young, L.J., Nuzhdin, S.V.: RNAseq: Technical variability and sampling. *BMC Genomics* **12**, 293 (2011)
  61. Merkling, S.H., Overheul, G.J., van Mierlo, J.T., Arends, D., Gilissen, C., van Rij, R.P.: The heat shock response restricts virus infection in *Drosophila*. *Scientific Reports* **5**, 12758 (2015)
  62. Dostert, C., Jouanguy, E., Irving, P., Troxler, L., Galiana, D., Hetru, C., Hoffmann, J.A., Imler, J.-L.: The JAK-STAT signaling pathway is required but not sufficient for the antiviral response of *Drosophila*. *Nature Immunology* **6**, 946–953 (2005)
  63. Flenniken, M.L., Andino, R.: Non-specific dsRNA-mediated antiviral response in the honey bee. *PLoS ONE* **8**, 77263 (2013)
  64. Brutscher, L.M., Daughenbaugh, K.F., Flenniken, M.L.: Virus and dsRNA-triggered transcriptional responses reveal key components of honey bee antiviral defense. *Scientific Reports* **7**, 6448 (2017)
  65. Chen, Y.P., Pettis, J.S., Corona, M., Chen, W.P., Li, C.J., Spivak, M., Visscher, P.K., DeGrandi-Hoffman, G., Boncristiani, H., Zhao, Y., vanEngelsdorp, D., Delaplane, K., Solter, L., Drummond, F., Kramer, M., Lipkin, W.I., Palacios, G., Hamilton, M.C., Smith, B., Huang, S.K., Zheng, H.Q., Li, J.L., Zhang, X., Zhou, A.F., Wu, L.Y., Zhou, J.Z., Lee, M.-L., Teixeira, E.W., Li, Z.G., Evans, J.D.: Israeli Acute Paralysis Virus: Epidemiology, pathogenesis and implications for honey bee health. *PLoS Pathogens* **10**, 1004261 (2014)
  66. Robinson, M.D., McCarthy, D.J., Smyth, G.K.: edgeR: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140 (2010)
  67. Carval, D., Ferriere, R.: A unified model for the coevolution of resistance, tolerance, and virulence. *Evolution* **64**, 2988–3009 (2010)
  68. Moret, Y.: Trans-generational immune priming: Specific enhancement of the antimicrobial immune response in the mealworm beetle, *Tenebrio molitor*. *Proceedings of the Royal Society B: Biological Sciences* **273**, 1399–1405 (2006)
  69. Mauricio, R., Rausher, M.D., Burdick, D.S.: Variation in the defense strategies of plants: are resistance and tolerance mutually exclusive? *Ecology* **78**, 1301–1310 (1997)
  70. Fornoni, J., Nunez-Farfan, J., Valverde, P.L., Rausher, M.D.: Evolution of mixed plant defense allocation against natural enemies. *Evolution* **58**, 1685–1695 (2004)
  71. Restif, O., Koella, J.C.: Shared control of epidemiological traits in a coevolutionary model of host-parasite interactions. *The American Naturalist* **161**, 827–836 (2003)
  72. Chambers, M.C., Schneider, D.S.: Balancing resistance and infection tolerance through metabolic means. *PNAS* **109**, 13886–13887 (2012)
  73. Xu, J., Grant, G., Sabin, L.R., Gordesky-Gold, B., Yasunaga, A., Tudor, M., Cherry, S.: Transcriptional pausing controls a rapid antiviral innate immune response in *Drosophila*. *Cell Host Microbe* **12**, 531–543 (2012)
  74. Ryabov, E.V., Fannon, J.M., Moore, J.D., Wood, G.R., Evans, D.J.: The Iflaviruses Sacbrood virus and Deformed wing virus evoke different transcriptional responses in the honeybee which may facilitate their horizontal or vertical transmission. *PeerJ* **4**, 1591 (2016)
  75. Johnson, B.R., Atallah, J., Plachetzki, D.C.: The importance of tissue specificity for RNA-seq: highlighting the errors of composite structure extractions. *BMC Genomics* **14**, 586 (2013)
  76. Thissen, D., Steinberg, L., Kuang, D.: Quick and easy implementation of the Benjamini-Hochberg procedure for controlling the false positive rate in multiple comparisons. *J Educ Behav Stat* **27**, 77–83 (2002)
  77. Consortium, H.B.G.S.: Finding the missing honey bee genes: lessons learned from a genome upgrade. *BMC Genomics* **15**, 86 (2014)
  78. Elsik, C.G., Tayal, A., Diesh, C.M., Unni, D.R., Emery, M.L., Nguyen, H.N., Hagen, D.E.: Hymenoptera Genome Database: integrating genome annotations in HymenopteraMine. *Nucleic Acids Research* **4**, 793–800 (2016)
  79. Wu, T.D., Reeder, J., Lawrence, M., Becker, G., Brauer, M.J.: GMAP and GSNAP for genomic sequence alignment: Enhancements to speed, accuracy, and functionality. *Methods Mol Biol* **1418**, 283–334 (2016)
  80. Love, M.I., Huber, W., Anders, S.: Moderated estimation of fold change and dispersion for RNA-seq data with

- 905 DESeq2. *Genome Biology* **15**, 550 (2014)
- 906 81. Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W., Smyth, G.K.: limma powers differential  
907 expression analyses for rna-sequencing and microarray studies. *Nucleic Acids Research* **43**(7), 47 (2015)
- 908 82. Benjamini, Y., Hochberg, Y.: Controlling the false discovery rate: A practical and powerful approach to multiple  
909 testing. *Journal of the Royal Statistical Society. Series B (Methodological)* **57**, 289–300 (1995)
- 910 83. Larsson, J.: eulerr: Area-Proportional Euler and Venn Diagrams with Ellipses. (2018). R package version 4.0.0.  
911 <https://cran.r-project.org/package=eulerr>
- 912 84. I.T. Jolliffe: *Principal Component Analysis*. Springer, Berlin, Heidelberg (2002)
- 913 85. Inselberg, A.: The plane with parallel coordinates. *The Visual Computer* **1**, 69–91 (1985)
- 914 86. W.S. Cleveland: *Visualizing Data*. Hobart Press, Summit, New Jersey (1993)
- 915 87. Cook, D., Hofmann, H., Lee, E., Yang, H., Nikolau, B., Wurtele, E.: Exploring gene expression data, using  
916 plots. *Journal of Data Science* **5**, 151–182 (2007)
- 917 88. Chandrasekhar, T., Thangavel, K., Elayaraja, E.: Effective Clustering Algorithms for Gene Expression Data.  
918 *International Journal of Computer Applications* **32**, 4 (2011)
- 919 89. de Souto D. de Araujo, M., Costa, I., Soares, R., Luderemir, T., Schliep, A.: Comparative Study on  
920 Normalization Procedures for Cluster Analysis of Gene Expression Datasets. *International Joint Conference on*  
921 *Neural Networks*, 2793–2799 (2008)
- 922 90. Huang, D.W., Sherman, B.T., Lempicki, R.: Systematic and integrative analysis of large gene lists using DAVID  
923 bioinformatics resources. *Nat Protoc* **4**, 44–57 (2009)
- 924 91. Huang, D.W., Sherman, B.T., Lempicki, R.A.: Bioinformatics enrichment tools: paths toward the  
925 comprehensive functional analysis of large gene lists. *Nucleic Acids Res* **37**, 1–13 (2009)
- 926 92. Supek, F., Bošnjak, M., Škunca, N., Šmuc, T.: REVIGO summarizes and visualizes long lists of Gene Ontology  
927 terms. *PLoS ONE* **6**, 21800 (2011)
- 928 93. Edgar, R., Domrachev, M., Lash, A.E.: Gene Expression Omnibus: NCBI gene expression and hybridization  
929 array data repository. *Nucleic Acids Res* **30**, 207–210 (2002)
- 930 94. Schlicker, A., Domingues, F.S., Rahnenfuhrer, J., Lengauer, T.: A new measure for functional similarity of gene  
931 products based on Gene Ontology. *BMC Bioinformatics* **7**, 302 (2006)

932 **Figures**

**Figure 1 Mortality rates for the four treatment groups, two virus groups, and two diet groups.** Left to right: Mortality rates for the four treatment groups, two virus groups, and two diet groups. “N” represents non-inoculation, “V” represents viral inoculation, “C” represents chestnut pollen, and “R” represents rockrose pollen. The mortality rate data included 59 samples with 15 replicates per treatment group, except for the “NC” group having 14 replicates. ANOVA values and p-values for the statistical tests are listed in the text of the paper. The letters above the bars represent significant differences with a confidence level of 95%.

**Figure 2 IAPV titers for the four treatment groups, two virus groups, and two diet groups.** Left to right: IAPV titers for the four treatment groups, two virus groups, and two diet groups. “N” represents non-inoculation, “V” represents viral inoculation, “C” represents chestnut pollen, and “R” represents rockrose pollen. The IAPV titer data included 38 samples with 10 replicates per treatment group, except for the “NR” group having 8 replicates. ANOVA values and p-values for the statistical tests are listed in the text of the paper. The letters above the bars represent significant differences with a confidence level of 95%.

**Figure 3 Parallel coordinate plots of the 1,019 DEGs after hierarchical clustering of size four between the virus-infected and control groups of the Galbraith data [44].** Parallel coordinate plots of the 1,019 DEGs after hierarchical clustering of size four between the virus-infected and control groups of the Galbraith study. “N” represents non-inoculation, “V” represents viral inoculation. Clusters 1, 2, and 4 seem to represent DEGs that were overexpressed in the virus inoculated group, and Cluster 3 seems to represent DEGs that were overexpressed in the non-inoculated control group. In general, the DEGs appeared as expected, but there is rather noticeable deviation of the first replicate from the virus-treated sample (“V.1”) from the other virus-treated replicates in Cluster 1.

**Figure 4 Parallel coordinate plots of the 43 DEGs after hierarchical clustering of size four between the virus-infected and control groups of our study.** Parallel coordinate plots of the 43 DEGs after hierarchical clustering of size four between the virus-infected and control groups of our study. “N” represents non-infected control group, and “V” represents treatment of virus. The vertical red line indicates the distinction between treatment groups. We see from this plot that the DEG designations for this dataset do not appear as clean compared to what we saw in the Galbraith dataset in Figure 3.

**Figure 5 Gene ontology analysis results for the 122 DEGs related to our “tolerance” hypothesis and for the 125 DEGs related to our “resistance” hypothesis.** GO analysis results for the 122 DEGs related to our “tolerance” hypothesis (A) and for the 125 DEGs related to our “resistance” hypothesis (B). The color and size of the circles both represent the number of genes in that ontology. The x-axis and y-axis are organized by SimRel, a semantic similarity metric [94].

**Figure 6 Venn diagrams comparing the virus-related DEG overlaps between our dataset and the Galbraith dataset.** Venn diagrams comparing the virus-related DEG overlaps between the Galbraith study (labeled as “G”) and our study (labeled as “R”). From left to right: Total virus-related DEGs (subplot A), virus-upregulated DEGs (subplot B), control-upregulated DEGs (subplot C). Both the total virus-related and virus-upregulated DEGs showed significant overlap between the studies ( $p\text{-value} < 2.2\text{e-}16$ ) as per Fisher’s Exact Test for Count Data. There was one gene that was virus-upregulated in the Galbraith study but control-upregulated in our study.



933 **Tables**

BeeBase ID	Gene Name	Known functions	Us	Galbraith
GB41545	MD-2-related lipid-recognition protein-like	Implicated in lipid recognition, particularly in the recognition of pathogen related products	N	-
GB50955	Protein argonaute-2	Interacts with small interfering RNAs to form RNA-induced silencing complexes which target and cleave transcripts that are mostly from viruses and transposons	V	V
GB48755	UBA-like domain-containing protein 2	Found in diverse proteins involved in ubiquitin/proteasome pathways	V	V
GB47407	Histone H4	Capable of affecting transcription, DNA repair, and DNA replication when post-transcriptionally modified	V	V
GB42313	Leishmanolysin-like peptidase	Encodes a protein involved in cell migration and invasion; implicated in mitotic progression in <i>D. melanogaster</i>	V	V
GB50813	Rho guanine nucleotide exchange factor 11	Implicated in regulation of apoptotic processes, cell growth, signal transduction, and transcription	V	V
GB54503	Thioredoxin domain-containing protein	Serves as a general protein disulphide oxidoreductase	N	-
GB53500	Transcriptional regulator Myc-B	Regulator gene that codes for a transcription factor	V	V
GB51305	Tropomyosin-like	Related to protein involved in muscle contraction	N	N
GB50178	Cilia and flagella-associated protein 61-like	Induces components required for wild-type motility and stable assembly of motile cilia	V	V

**Table 1** Known functions of the mapped subset of 43 DEGs in the virus main effect of our study. Whether the gene was overrepresented in the virus or non-virus group is also indicated for both our study and the Galbraith study. Functionalities were extracted from Flybase, National Center for Biotechnology Information and The European Bioinformatics Institute databases.

Contrast	DEGs	Interpretation	Results
V (all) vs N (all)	43	Genes that change expression due to virus effect regardless of diet status in bees	Table 1
NC vs NR	941	Genes that change expression due to diet effect in uninfected bees	Supplementary tables 4 and 5, Additional file 1
VC vs VR	376	Genes that change expression due to diet effect in infected bees	Supplementary tables 6 and 7, Additional file 1
VC upregulated in VC vs VR, and NC upregulated in NC vs NR	122	“Tolerance” genes that turn on by good diet regardless of virus infection status in bees	Figure 5A
VC upregulated in VC vs VR, but NC not upregulated in NC vs NR	125	“Resistance” genes that turn on by good diet only in infected bees	Figure 5B

**Table 2** Contrasts in our study for assessing GO and pathways analysis.

#### Additional Files

Additional file 1 — Supplementary tables.

**Table 1:** Number of DEGs across three analysis pipelines for (A) the diet main effect in our study, (B) the virus main effect in our study, and (C) the virus main effect in the Galbraith study. For the diet effects, “C” represents chestnut diet and “R” represents rockrose diet. For the virus effects, “N” represents control non-inoculated and “V” represents virus-inoculated. **Table 2:** Pathways related to the 1,033 DEGs that were upregulated in the chestnut treatment from the diet main effect. **Table 3:** Pathways related to the 881 DEGs that were upregulated in the rockrose treatment from the diet main effect. **Table 4:** GO analysis results for the 601 DEGs that were upregulated in the NC treatment from the NC versus NR treatment pair analysis. These DEGs represent genes that are upregulated when non-infected honey bees are given high quality chestnut pollen compared to being given low quality rockrose pollen. **Table 5:** GO analysis results for the 340 DEGs that were upregulated in the NR treatment from the NC versus NR treatment pair analysis. These DEGs represent genes that are upregulated when non-infected honey bees are given low quality rockrose pollen compared to being given high quality chestnut pollen. **Table 6:** GO analysis results for the 247 DEGs that were upregulated in the VC treatment from the VC versus VR treatment pair analysis. These DEGs represent genes that are upregulated when infected honey bees are given high quality chestnut pollen compared to being given low quality rockrose pollen. **Table 7:** GO analysis results for the 129 DEGs that were upregulated in the VR treatment from the VC versus VR treatment pair analysis. These DEGs represent genes that are upregulated when infected honey bees are given low quality rockrose pollen compared to being given high quality chestnut pollen. **Table 8:** Number of DEGs across three analysis pipelines for all six treatment pair combinations between the diet and virus factor. “C” represents chestnut diet, “R” represents rockrose diet, “V” represents virus-inoculated, and “N” represents control non-inoculated. **Table 9:** Kruskal-Wallis p-value and Bonferroni corrections for the 36 combinations of DEG lists, pathogen response metrics, and cluster number. (XLS).

Additional file 2 — PCA plots for the Galbraith dataset and for our dataset.

PCA plots for the Galbraith dataset (A) and for our dataset (B). “V” represents virus-inoculated, and “N” represents control non-inoculated. The x-axis represents the principal component with the most variation and the y-axis represents the principal component with the second-most variation (PNG).

Additional file 3 — Parallel coordinate lines of the diet-related DEGs of our dataset.

Parallel coordinate plots of the 1,914 DEGs after hierarchical clustering of size six between the chestnut and rockrose groups of our study. Here “C” represents chestnut samples, and “R” represents rockrose samples. The vertical red line indicates the distinction between treatment groups. We see from this plot that the DEG designations for this dataset do not appear as clean compared to what we saw in the Galbraith dataset in Figure 3 (PNG).

Additional file 4 — Example litre plots from the virus-related DEGs of our dataset.

Example litre plots of the nine DEGs with the lowest FDR values from the 43 virus-related DEGs of our dataset. “N” represents non-infected control samples and “V” represents virus-treated samples. Most of the magenta points (representing the 144 combinations of samples between treatment groups for a given DEG) do not reflect the expected pattern as clearly compared to what we saw in the litre plots of the Galbraith data. They are not as clustered together (representing replicate inconsistency) and they sometimes cross the  $x=y$  line (representing lack of difference between treatment groups). This finding reflects what we saw in the messy looking parallel coordinate lines of Figure 4 (PNG).

973 Additional file 5 — Example litre plots of DEGs from Cluster 1 of the Galbraith dataset.

974 Example litre plots of the nine DEGs with the lowest FDR values from the 365 DEGs in Cluster 1 (originally shown  
975 in Figure 3) of the Galbraith dataset. "N" represents non-infected control samples and "V" represents virus-treated  
976 samples. Most of the light orange points (representing the nine combinations of samples between treatment groups  
977 for a given DEG) deviate from the  $x=y$  line in a tight bundle as expected (PNG).

978 Additional file 6 — Example litre plots of DEGs from Cluster 2 of the Galbraith dataset.

979 Example litre plots of the nine DEGs with the lowest FDR values from the 327 DEGs in Cluster 2 (originally shown  
980 in Figure 3) of the Galbraith dataset. "N" represents non-infected control samples and "V" represents virus-treated  
981 samples. Most of the dark orange points (representing the nine combinations of samples between treatment groups  
982 for a given DEG) deviate from the  $x=y$  line in a compact clump as expected. However, they are not as tightly  
983 bunched together compared to what we saw in the example litre plots of Cluster 1 (shown in Additional file 5). As a  
984 result, what we see in these litre plots reflects what we saw in the parallel coordinate lines of Figure 3: The replicate  
985 consistency in the Cluster 1 DEGs is not as clean as that in the Cluster 2 DEGs, but is still relatively clean (PNG).

986 Additional file 7 — Scatterplot matrix of DEGs from Cluster 1 of the Galbraith dataset.

987 The 365 DEGs from the first cluster of the Galbraith dataset (originally shown in Figure 3) superimposed as light  
988 orange dots onto all genes as black dots in the form of a scatterplot matrix. The data has been standardized. "N"  
989 represents non-infected control samples and "V" represents virus-treated samples. We confirm that the DEGs  
990 mostly follow the expected structure, with their placement deviating from the  $x=y$  line in the treatment  
991 scatterplots, but adhering to the  $x=y$  line in the replicate scatterplots. However, we do see that sample "V.1" may  
992 be somewhat inconsistent in these DEGs, as its presence in the replicate scatterplots shows DEGs deviating from  
993 the  $x=y$  line more than expected and its presence in the treatment scatterplots shows DEGs adhering to the  $x=y$   
994 line more than expected. This inconsistent sample was something we observed in Figure 3 (PNG).

995 Additional file 8 — Scatterplot matrix of DEGs from Cluster 2 of the Galbraith dataset.

996 The 327 DEGs from the second cluster of the Galbraith dataset (originally shown in Figure 3) superimposed as dark  
997 orange dots onto all genes as black dots in the form of a scatterplot matrix. The data has been standardized. "N"  
998 represents non-infected control samples and "V" represents virus-treated samples. We confirm that the DEGs  
999 mostly follow the expected structure, with their placement deviating from the  $x=y$  line in the treatment  
1000 scatterplots, but adhering to the  $x=y$  line in the replicate scatterplots (PNG).

1001 Additional file 9 — Scatterplot matrix of DEGs from Cluster 3 of the Galbraith dataset.

1002 The 224 DEGs from the third cluster of the Galbraith dataset (originally shown in Figure 3) superimposed as  
1003 turquoise dots onto all genes as black dots in the form of a scatterplot matrix. The data has been standardized. "N"  
1004 represents non-infected control samples and "V" represents virus-treated samples. We confirm that the DEGs  
1005 mostly follow the expected structure, with their placement deviating from the  $x=y$  line in the treatment  
1006 scatterplots, but adhering to the  $x=y$  line in the replicate scatterplots (PNG).

1007 Additional file 10 — Scatterplot matrix of DEGs from Cluster 4 of the Galbraith dataset.

1008 The 103 DEGs from the fourth cluster of the Galbraith dataset (originally shown in Figure 3) superimposed as pink  
1009 dots onto all genes as black dots in the form of a scatterplot matrix. The data has been standardized. "N"  
1010 represents non-infected control samples and "V" represents virus-treated samples. We confirm that the DEGs  
1011 mostly follow the expected structure, with their placement deviating from the  $x=y$  line in the treatment  
1012 scatterplots, but adhering to the  $x=y$  line in the replicate scatterplots. We also see that the second replicate from  
1013 the virus-treated sample ("V.2") may be somewhat inconsistent in these DEGs, as its presence in the replicate  
1014 scatterplots results in the DEGs unexpectedly deviating from the  $x=y$  line and its presence in the treatment  
1015 scatterplots results in the DEGs unexpectedly adhering to the  $x=y$  line (PNG).

1016 Additional file 11 — Scatterplot matrix of virus-related DEGs from our dataset, showing only replicates 1, 2, and 3.

1017 The 43 virus-related DEGs from our dataset superimposed as magenta dots onto all genes in the form of a  
1018 scatterplot matrix. Only replicates 1, 2, and 3 are shown from both treatment groups. The data has been  
1019 standardized. "N" represents non-infected control samples and "V" represents virus-treated samples. We see that,  
1020 compared to the scatterplot matrices from certain clusters of the Galbraith data, the 43 DEGs from this subset of  
1021 six samples from our data do not paint as clear of a picture, sometimes unexpectedly deviating from the  $x=y$  line in  
1022 the replicate plots and sometimes unexpectedly adhering to the  $x=y$  line in the treatment plots (PNG).

1023 Additional file 12 — Scatterplot matrix of virus-related DEGs from our dataset, showing only replicates 4, 5, and 6.

1024 The 43 virus-related DEGs from our dataset superimposed as magenta dots onto all genes in the form of a  
1025 scatterplot matrix. Only replicates 4, 5, and 6 are shown from both treatment groups. The data has been  
1026 standardized. "N" represents non-infected control samples and "V" represents virus-treated samples. We see that,  
1027 compared to the scatterplot matrices from certain clusters of the Galbraith data, the 43 DEGs from this subset of  
1028 six samples from our data do not paint as clear of a picture, and most of them unexpectedly adhere to the  $x=y$  line  
1029 in the treatment plots (PNG).

1030 Additional file 13 — Scatterplot matrix of virus-related DEGs from our dataset, showing only replicates 7, 8, and 9.  
 1031 The 43 virus-related DEGs from our dataset superimposed as magenta dots onto all genes in the form of a  
 1032 scatterplot matrix. Only replicates 7, 8, and 9 are shown from both treatment groups. The data has been  
 1033 standardized. "N" represents non-infected control samples and "V" represents virus-treated samples. We see that,  
 1034 compared to the scatterplot matrices from certain clusters of the Galbraith data, the 43 DEGs from this subset of  
 1035 six samples from our data do not paint as clear of a picture, sometimes unexpectedly deviating from the  $x=y$  line in  
 1036 the replicate plots and sometimes unexpectedly adhering to the  $x=y$  line in the treatment plots (PNG).

1037 Additional file 14 — Scatterplot matrix of virus-related DEGs from our dataset, showing only replicates 10, 11, and  
 1038 12.  
 1039 The 43 virus-related DEGs from our dataset superimposed onto all genes in the form of a scatterplot matrix. Only  
 1040 replicates 10, 11, and 12 are shown from both treatment groups. The data has been standardized. "N" represents  
 1041 non-infected control samples and "V" represents virus-treated samples. We see that, compared to the scatterplot  
 1042 matrices from certain clusters of the Galbraith data, the 43 DEGs from this subset of six samples from our data do  
 1043 not paint as clear of a picture, and most of them unexpectedly deviate from the  $x=y$  line in the virus-related  
 1044 replicate plots (PNG).

1045 Additional file 15 — Parallel coordinate plots of the "tolerance" candidate DEGs.  
 1046 Parallel coordinate plots of the 122 DEGs after hierarchical clustering of size four between the "tolerance" candidate  
 1047 DEGs. Here "N" represents non-infected control group, "V" represents treatment of virus, "C" represents  
 1048 high-quality chestnut diet, and "R" represents low-quality rockrose diet. The vertical red line indicates the  
 1049 distinction between treatment groups. We see there is considerable noise in the data (non-consistent replicate  
 1050 values), but that the general patterns of the DEGs follow what we expect based on our "tolerance" contrast (PNG).

1051 Additional file 16 — Parallel coordinate plots of the "resistance" candidate DEGs.  
 1052 Parallel coordinate plots of the 125 DEGs after hierarchical clustering of size four between the "resistance"  
 1053 candidate DEGs. Here "N" represents non-infected control group, "V" represents treatment of virus, "C" represents  
 1054 high-quality chestnut diet, and "R" represents low-quality rockrose diet. The vertical red line indicates the distinction  
 1055 between treatment groups. We see there is considerable noise in the data (non-consistent replicate values), but that  
 1056 the general patterns of the DEGs follow what we expect based on our "resistance" contrasts (PNG).

1057 Additional file 17 — Venn diagrams comparing the virus-related DEG overlaps in the Galbraith data using our  
 1058 pipeline and the pipeline used by Galbraith *et al.*  
 1059 Venn diagrams comparing the virus-related DEG overlaps of the Galbraith data from the DESeq2 bioinformatics  
 1060 pipelines used in the Galbraith study (labeled as "G.O.") and the DESeq2 bioinformatics pipelines used in our study  
 1061 (labeled as "G.R"). While we were not able to fully replicate the DEG list published in the Galbraith study, our DEG  
 1062 list maintained significant overlaps with their DEG list. From left to right: Total virus-related DEGs (subplot A),  
 1063 virus-upregulated DEGs (subplot B), control-upregulated DEGs (subplot C) (PNG).

1064 Additional file 18 — Venn diagrams of main effect DEG overlaps across DESeq2, edgeR, and limma  
 1065 Venn diagrams comparing DEG overlaps across DESeq2, edgeR, and limma for our diet main effect (top row), our  
 1066 virus main effect (middle row), and the Galbraith virus main effect (bottom row). Within a given subplot, "D"  
 1067 represents DESeq2, "E" represents edgeR, and "L" represents limma. From left to right on top row: Total  
 1068 diet-related DEGs (subplot A), chestnut-upregulated DEGs (subplot B), rockrose-upregulated DEGs (subplot C).  
 1069 From left to right on middle row: Total virus-related DEGs (subplot D), virus-upregulated DEGs (subplot E),  
 1070 control-upregulated DEGs in our data (subplot F). From left to right on bottom row: Total virus-related DEGs  
 1071 (subplot G), virus-upregulated DEGs (subplot H), control-upregulated DEGs in the Galbraith data (subplot I)  
 1072 (PNG). With the exception of the limma pipeline resulting in zero DEGs in our virus main effect analysis, we found  
 1073 significant overlaps between DEG lists across the different pipelines (DESeq2, edgeR, and limma). In general,  
 1074 DESeq2 resulted in the largest number of DEGs and limma resulted in the least number of DEGs (PNG).

1075 Additional file 19 — Analysis of correlation between DEG read counts and pathogen response metrics  
 1076 Distribution of R-squared values for DEG cluster read counts and pathogen response metrics. Columns left to right:  
 1077 SBV titers, mortality rates, and IAPV titers. Rows top to bottom: Tolerance candidate DEGs, resistance candidate  
 1078 DEGs, and virus-related DEGs. Each subplot includes five boxplots which represent the R-squared value distributions  
 1079 for four DEG clusters and all remaining non-DEGs in the data. The top number above each boxplot represents the  
 1080 number of genes included. The first four boxplots also include a bottom number, which represents the  
 1081 Kruskal-Wallis p-value of the comparison of the R-squared distribution of the cluster and the R-squared distribution  
 1082 of the non-DEG data (PNG).