1 **Untangling an insect virome from its endogenous viral elements**

2

3 Paula Rozo-Lopez,[1,*] William Brewer,[1] Simon Käfer,[2] McKayla M. Martin,[1] & Benjamin J. Parker
4 [1,*]

5

6 [1] Department of Microbiology, University of Tennessee, Knoxville, TN 37916, USA

7

8 [2] Institut für Biologie und Umweltwissenschaften, Carl von Ossietzky Universität Oldenburg,
9 26129 Oldenburg, Germany

10

11 * authors for correspondence: PRL: plopez2@utk.edu; BJP: bjp@utk.edu

12

13 ORCIDs: PRL - 0000-0001-9207-6579; SK - 0000-0003-3270-8348; MMM - 0000-0002-7471-
14 6023; BJP - 0000-0002-0679-4732

15 ABSTRACT
16
17 Insects are an important reservoir of viral biodiversity, but the vast majority of viruses associated
18 with insects have not been discovered. Recent studies have therefore employed high-
19 throughput sequencing of RNA, which has led to rapid advances in our understanding of insect
20 viral diversity. However, insect genomes frequently contain transcribed endogenous viral
21 elements with significant homology to exogenous viruses, complicating the use of RNAseq for
22 viral discovery. In this study, we use a multi-pronged sequencing approach to study the virome
23 of an important agricultural pest and prolific vector of plant pathogens, the potato aphid
24 *Macrosiphum euphorbiae*. We first used rRNA-reduced RNAseq to characterize the bacteria
25 and viruses found in individual insects. We then characterized the frequency of a heritable
26 Flavivirus and an Ambidensovirus in our population. We next generated a quality draft genome
27 assembly for *M. euphorbiae* using Illumina-corrected Nanopore sequencing. This analysis
28 showed that the Ambidensovirus, previously described from an RNAseq viral screen, is not a
29 exogenous virus and instead is a transcribed endogenous viral element in the *M. euphorbiae*
30 genome. Our study generates key insight into an important agricultural pest and highlights a
31 widespread challenge for the study of viral diversity using RNAseq.
32
33 KEYWORDS
34
35 Viral discovery; RNAseq; insects; aphids; endogenous viral elements

36 INTRODUCTION
37
38 The last decade has transformed our understanding of the viral communities associated with
39 insects, the most abundant and diversified animal group [1-4]. Insect viruses have been
40 primarily studied in the context of vector-borne pathogens, which are transmitted horizontally
41 between insect vectors and amplifying hosts and often have medical or agricultural relevance.
42 Other viruses, however, only replicate within the insect and are maintained in natural
43 populations through horizontal and/or vertical transmission. These insect-specific viruses have
44 been shown to have important impacts on host biology [5-7], but much work remains to be done
45 to describe insect-specific viral diversity and uncover the hidden role they play in insect
46 phenotypes and evolution [8-10].
47
48 To address this gap, researchers have employed high-throughput approaches to viral discovery,
49 including next-generation sequencing and analysis of RNA. Recent studies have used this
50 approach to characterize and discover an enormous diversity of viruses [2, 11-15]. However,
51 there are several serious limitations to this approach. For example, it is unclear from RNAseq
52 data whether viral reads come from microbes infecting insect cells or if they are present from an
53 organism ingested by the insect. Another potential challenge with using RNAseq for viral
54 discovery is that insects often harbor fragments of viral sequences in their genomes. The
55 endogenous viral elements (EVEs) described to date have homology with multiple clades of
56 single- and double-stranded DNA and RNA viral families [16]. We have a limited understanding
57 of the role EVEs are playing in insect biology, but transcriptionally active EVEs have been
58 shown to play functional roles in regulating host genome stability and as an antiviral defense
59 against exogenous viruses [17-19]. EVEs are remarkably common across insects [20], and thus
60 EVEs could represent a widespread challenge facing the field. As such, studies are needed to
61 uncover the contribution of EVEs to insect 'viromes'.
62
63 Aphids (Hemiptera: Aphidoidea) are hosts to diverse viruses, including plant pathogens with
64 agricultural significance and insect-specific viruses [21, 22]. Recent studies have used
65 metatranscriptome sequencing to describe viral diversity in aphids [23-27], and have described
66 insect-specific DNA viruses in the family Parvoviridae (Ambidensovirus) and RNA viruses in the
67 Bunyaviruses, Dicistroviruses, Flaviviruses, Iflaviruses, and Mesoviruses families [21].
68 The potato aphid *Macrosiphum euphorbiae* (Thomas, 1878) is an important cosmopolitan
69 agricultural pest that infests tomatoes, potatoes, and other economically important crops [28].
70 *M. euphorbiae* is also an important vector of plant viruses (Families Bromoviridae,
71 Closteroviridae, Geminiviridae, Potyviridae, and Solemoviridae) and was recently shown to host
72 several insect-specific viruses belonging to the families Flaviviridae (Flavivirus) and
73 Parvoviridae [24, 29, 30]. Despite *M. euphorbiae's* economic importance, no genomic resources
74 are available outside the body and salivary gland transcriptomes [24, 31, 32].
75
76 The genomes of multiple aphid species have been shown to harbor EVEs that mediate growth,
77 development, and wing plasticity [33-37]. In this study, we use next-generation sequencing and
78 molecular techniques to show that aphid EVEs have led to the misidentification of aphid viruses
79 from RNAseq data. First, we used RNAseq to characterize the microbial diversity of field-
80 collected *M. euphorbiae* adults, and we found evidence of two insect-specific viruses infecting
81 aphids collected from the field, including a Flavivirus and Ambidensovirus. Then, we generated
82 a high-quality draft genome sequence of this species. Our genome showed that insect-specific
83 Ambidensoviral hits correspond to transcriptionally active EVEs, indicating that a previously
84 described virus is actually an endogenous viral element in the *M. euphorbiae* genome. Our
85 results illustrate how careful analysis using multiple methods is needed to untangle insect

86  viromes from EVEs, and this study furthers our understanding of the surprisingly widespread
87  presence of densoviral EVEs in aphid genomes.
88
89
90  METHODS
91
92  **Aphid collection:** We collected asexual winged and wingless female *M. euphorbiae* adults from
93  cultivated tomato plants (var Husky Cherry Red) in Knoxville, TN, USA, between April and June
94  2021 and 2022. We stored individual aphids in 1.5 mL Eppendorf tubes (Eppendorf, Hamburg,
95  Germany) at -80°C until processing. To validate our ability to identify *M. euphorbiae* (NCBI
96  TaxID: 13131), we used COI barcoding (LCO1490 5'-GGTCAACAAATCATAAAGATATTGG-3'
97  and HCO2198 5'-TAAACTTCAGGGTGACCAAAAAATCA-3'), sanger sequencing, and
98  comparisons of our COI sequences to the Barcode of Life Data System
99  (https://www.boldsystems.org/) [38]. Our partial COI barcode sequence was uploaded to NCBI
100  with accession number OQ588703.
101
102  **Cultivation of *M. euphorbiae* strain Me57:** To establish a colony of *M. euphorbiae* in the
103  laboratory, we used a single asexual female collected in 2021. After colonization, we maintained
104  this line on tomato plants (Husky Cherry Red) at 20°C 16L:8D. We screened the line for the
105  seven species of facultative symbionts found in aphids using established PCR protocols [39,
106  40]. For this screen, we extracted DNA using 'Bender buffer' and ethanol precipitation as in
107  previous studies [41, 42]. We then used PCR with species-specific primers [39, 43] to screen for
108  *Hamiltonella defensa, Fukatsuia symbiotica* (X-type), *Regiella insecticola*, *Rickettsia* sp.,
109  *Ricketsiella* sp., *Serratia symbiotica*, and *Spiroplasma* sp. following the recommended thermal
110  profiles (94°C for 2 min, 11 cycles of 94°C for 20 sec, 56°C (declining 1°C each cycle) for 50
111  sec, 72°C for 30 sec, 25 cycles of 94°C for 2 min, 45°C for 50 sec, 72°C for 2 min, and a final
112  extension of 72°C for 5 min).
113
114  **RNA extraction and sequencing**: We homogenized individual aphids with a pestle in 500 μL of
115  TRIzol (Invitrogen; Thermo Fisher Scientific, Inc., Waltham, MA, USA) and extracted total RNA
116  using BCP (1-bromo-3-chloropropane; Life Technologies, Thermo Fisher Scientific, Inc.,
117  Waltham, MA, USA) with isopropanol precipitation. We used the Zymo RNA Clean &
118  Concentrator kit (Zymo Genetics Inc., Seattle, WA, USA) to improve the purity and to remove
119  gDNA using DNAse I. We then performed metatranscriptome Sequencing at Novogene
120  (Novogene Corporation Inc., Sacramento, CA, USA). Library preparation was conducted using
121  ribosomal RNA (rRNA) depletion by Illumina TruSeg Stranded Total RNA with Ribo-Zero Plus
122  and NEBNext rRNA Depletion Kit (Zymo Genetics, Inc., Seattle, WA, USA). The libraries were
123  sequenced to approximately 9 billion base pairs (bp) per sample with 150 bp paired-end reads
124  on an Illumina NovaSeq platform. Raw reads were deposited into the NCBI Sequence Read
125  Archive under BioProject ID PRJNA942253 with BioSample accessions SAMN33770905-
126  SAMN33770908, and data accessions SRR23870213-SRR23870216.
127
128  **Microbial analysis using CZID:** We assessed the success of ribosomal reduction in the
129  metatranscriptome libraries using riboPicker [44] and the reference database SILVA_138 [45]
130  (supplementary file reads_report.csv). We then used the CZ ID platform pipeline V7.1
131  (https://czid.org) [46], a cloud-based, open-source bioinformatics platform designed to detect
132  microbes from metagenomic data. We removed host-specific reads (STAR host subtraction)
133  using the *Acyrthosiphon pisum* genome [47], trimmed adapters using Trimmomatic [48],
134  removed low-quality reads with PriceSeqFilter [49], and aligned the remaining reads to the NCBI
135  NT and NR databases using Minimap2 [50] and Diamond [51]. In parallel, short reads were *de*

136 *novo* assembled using SPADES [52] and mapped back to the resulting contigs using bowtie2
137 [53] to identify the contig to which each raw read belongs. We used the CZ ID water background
138 model, which evaluates the significance (z-scores) of relative abundance estimates for microbial
139 taxa in each sample. Potential bacterial reads were distinguished from contaminating
140 environmental sequences by establishing z-score metrics ≥10, alignment length over 50
141 matching nucleotides (NT L ≥50), and a minimum of five reads per million aligning to the
142 reference protein database (NR rPM ≥ 5). Potential viruses were established by z-score metrics
143 of ≥1, NT L ≥50, and NR rPM ≥ 5 [46, 54, 55]. Bacterial and viral hits were confirmed with
144 BLASTX and BLASTN manual searches. Only annotated non-host hits with revised Taxonomy
145 IDs and BLAST-based match refinement were used for further analysis. The "Macrosiphum
146 euphorbiae" project is viewable and searchable to anyone in CZ ID.

147

148 **Densoviral analysis using *de novo* assembly and Travis:** We conducted an additional
149 screening and viral genome assembly of potential Ambidensoviruses using *de novo*
150 transcriptome assemblies obtained as follows. We used Trimmomatic v.0.39 [48] to trim the
151 sequence adapters and filtered low-quality/complexity reads, and assessed for post-trimming
152 quality using FastQC v.0.11.9 [56]. Then, we used Trinity v.2.14 [57] to *de novo* assemble the
153 remaining reads. We used TRAVIS (v.20221029, https://github.com/kaefers/travis) to scan the
154 assembled transcriptomes for Densovirus-like sequences. We built the reference database
155 according to the currently accepted Densovirinae (ICTV, 29. Oct 2022, see supplementary file
156 parvoviridae_reference_library.csv), extracted open reading frames between 100 and 2000
157 amino acids from the assembled transcriptomes, and screened using HMMER v3.3.1 [58],
158 MMSeqs2 [59], BLASTP v2.12.0 [60], and Diamond v2.0.15 [51]. We set the e-value cutoff at 1
159 ×$10^{-6}$, where applicable. All hits were again searched with Diamond against the non-redundant
160 protein database (NCBI, downloaded on 29 Oct 2022).

161

162 **MeV-1 genome analysis:** We used the CZ ID viral consensus genomes pipeline to build a
163 consensus genome from the sample with MeV-1 present at high levels. In short, contigs were
164 aligned to the reference MeV-1 genome (NCBI Entry KT309079.1) using minimap2 [50] and
165 then trimmed using TrimGalore (Phred score <20) [61]. The consensus genome was generated
166 with iVar consensus using a depth of five or more reads [62].

167

168 **MeV-1, MeV-2, and *Hamiltonella defensa* screening**: Like all aphids, *M. euphorbiae* hosts an
169 obligate heritable bacterial symbiont called *Buchnera aphidicola* that synthesizes amino acids
170 missing from the aphid's diet of plant phloem, and can also harbor several other facultative
171 symbiotic bacteria (listed above) [43]. To screen for these microbes, we used 1 µg of total RNA
172 extracted (as above) from each of the 23 adults collected during 2022 for cDNA synthesis with
173 iScript cDNA synthesis kit (Bio-Rad Laboratories, Inc., Hercules, CA, USA). To screen for the
174 Flavivirus Macrosiphum euphorbiae virus 1 (MeV-1), we used 100 ng of cDNA, the primers
175 MevirF1 (5'-GTACACTTGCCTTACCTTACTGT-3') and MevirR1b (5'-
176 AACACGGGTCACGACCTTAG-3'), and the PCR conditions previously described [30]. To
177 screen for the Ambidensovirus Macrosiphum euphorbiae virus 2 (MeV-2), we used 100 ng of
178 cDNA, the MeV2-F (5'-CCGGATGACAAATCCCACGA-3') and MeV2-R (5'-
179 AATAGGCGCAGAGATGGACG-3') primers, and the recommended PCR conditions [24]. In
180 addition, we extracted DNA from colonized Me57 aphids (as above) and used 40 ng of genomic
181 DNA to screen for MeV-2. The aphid Glyceraldehyde 3-phosphate dehydrogenase (G3PDH)
182 was used as internal control (primers G3PDH_F (5'-CGGGAATTTCATTGAACGAC-3') and
183 G3PDH_R (5'- TCCACAACACGGTTGGAGTA-3') [35]).

184

185 We used 200 ng of the cDNA previously synthesized for MeV-1 and MeV-2 screening and the
186 protocols for *Hamiltonella defensa* PCR screening (as described above) to evaluate the

187 proportion of field-collected aphids harboring this bacterial symbiont (supplementary file
188 samples_metadata.csv). Furthermore, we used a non-parametric (Spearman) correlation to
189 investigate the potential interaction between *Hamiltonella* and MeV-1.
190
191 **DNA extraction and sequencing**: We pooled seven genetically identical adult unwinged
192 aphids from cultivated lab line Me57 and isolated genomic DNA (gDNA) using a
193 phenol/chloroform extraction. We then sheared the gDNA to approximately 20kb fragments
194 using Covaris G-tubes (Covaris LLC., Woburn, MA, USA) at 4200 RMP for 1 minute, followed
195 by tube inversion. For library preparation, we used the NEB Next PPFE repair kit with Ultra II
196 end prep reaction (New England Biolabs, Ipswich, MA, USA) under recommended conditions
197 and Nanopore ligation sequencing kit SQK-LSK110. For sequencing, we used a Nanopore
198 R9.4.1 (FLO-MIN106D) flow cell and a MinION MIN-101B sequencing device (Oxford Nanopore
199 Technologies, Oxford, UK). We ran the flow cell for 24 hours, followed by a wash with Flow Cell
200 Wash Kit (EXP-WSH004); we then reloaded the flow cell with a second library prep and ran the
201 sequencer for an additional 48 hours. We stopped the second sequencing run at 72 hours (~22
202 Gbps of sequencing). In addition, we performed an additional 5.3 Gb of 150 bp paired-end
203 sequencing to polish the assembly on an Illumina NovaSeq platform. DNA was extracted as
204 above, and library prep and sequencing were performed by Novogene Inc. Raw reads were
205 filtered for low quality and adapter contamination by Novogene Inc.
206
207 *M. euphorbiae* **whole genome assembly:** We used Guppy (Oxford Nanopore Technologies)
208 for base-calling and quality trimming raw reads. For the removal of *Buchnera* reads, we used
209 minimap2 v.2.24 [50] in conjunction with SAMtools v.1.15.1 [63] to map our reads against the
210 *Buchnera aphidicola* (strain *Macrosiphum euphorbiae*) genome (NCBI accession
211 NZ_CP029205) and the corresponding plasmids (NCBI accession number NZ_CP029203 and
212 NZ_CP029204). We only kept unmapped reads for aphid genome assembly. We assembled
213 Nanopore reads using CANU v.2.0 [64] with an estimated genome size of 541 Mbp. We
214 removed allelic variants from the assembly using the purge_haplotigs v.1.1.2 [65], first by
215 mapping reads to the assembly using minimap2 v2.24-r1122 with Samtools v.1.15.1 and
216 manually choosing cutoffs for haploid vs. diploid coverage based on a histogram plot (v -l 5  -m
217 27  -h 60), and then by purging duplicated contigs based on coverage level (-j 80 -s 50). For
218 assembly polishing, we used the Illumina reads after quality assessment using FastQC V0.11.9
219 [56]. Then we used these reads to polish the purged assembly using Pilon v.1.24 with default
220 parameters [66]. We used BlobTools2 [67] to identify remaining contaminating contigs. For this,
221 we used blast results obtained from the BLASTN function against the NR database using blast
222 plus v.2.12.0 [68], read coverage obtained by mapping the Illumina reads to the assembly using
223 minimap2 v.2.24 [50], and GC content in this analysis. Based on these results, we removed all
224 the short contigs with strong homology to the plant genus *Solanum* (which includes the tomato
225 host plant species of *M. euphorbiae*) as we suspect these contigs were assembled from host
226 plant contamination in the guts of sequenced aphids. We also removed two short contigs with
227 homology to other bacterial contaminants such as *Escherichia coli* and *Pseudomonas* sp.
228 Lastly, we removed a contig nearly identical to the pLeu plasmid found in *Buchnera aphidicola*
229 and a small portion of two large contigs matching the *Buchnera* genome. The final annotation
230 was assessed using BUSCO v.5.3.2 [69] with the MetaEuk gene predictor [70] implemented in
231 galaxy.org, using the hemiptera_odb10 (2020-08-05) lineage dataset. The *M. euphorbiae*
232 genome is available in NCBI with BioProject ID PRJNA942253 and BioSample
233 SAMN33681650. The raw Nanopore (SRR23851809) and Illumina reads (SRR23919025)
234 associated with the genome are available through the Sequence Read Archive, and the finished
235 assembly is available with accession number JARHUA000000000.
236

237 **Characterizing endogenous viral elements in the *M. euphorbiae* genome:** DNA Illumina
238 raw reads were used as input to the CZ ID platform pipeline V7.1 (https://czid.org) and a z-score
239 metrics of ≥1 and NT L ≥50 as described above [46, 54]. Additionally, to screen for actively
240 transcribed Ambidensovirus-like EVES in the *M. euphorbiae* genome, we used BLASTN
241 searches using the seven viral hits provided in individual Trinity contigs flagged by TRAVIS
242 (supplementary file contigs_TRAVIS.fasta) against the genome scaffolds. All non-redundant hits
243 from these searches with E-values < $1.10^{-3}$ were extracted and used in further analyses [33].
244
245
246 RESULTS
247
248 **Analysis of non-host sequences detected in single aphid:** We used the pea aphid (*A.*
249 *pisum*) genome to subtract host reads from our transcriptome data set. On average, 81.8% of
250 the reads mapped to this host and were subtracted from further analysis (see supplementary file
251 reads_report.csv). We then analyzed the remaining distribution of non-pea aphid reads, within a
252 single *M. euphorbiae* aphid, as the overall proportion of reads assembled into contigs that could
253 be assigned to bacterial, eukaryotic, and viral taxa (public project Macrosiphum euphorbiae at
254 https://czid.org). Bacterial taxa dominated the microbial signature (Figure 1A), and as expected,
255 the highest number of reads assembled into contigs matched the aphid obligate symbiont
256 *Buchnera aphidicola* with over 45,000 reads per million aligning to the nucleotide database (NT
257 rPM>45,000). Reads from an aphid facultative symbiont *Hamiltonella defensa,* were found in
258 two samples (NT rPM>8,700). One sample (Me152) showed a strong signature of bacterial
259 contaminants (*E. coli, Pseudomonas, Halomonas,* and *Agrobacterium*) commonly present in soil
260 and plant surfaces.
261
262 In terms of eukaryotes (Figure 1B), we found hits to Solanaceae, which includes the host plant
263 species of *M. euphorbiae,* and Brachonidae parasitoid wasps (Insecta: Hymenoptera) in two
264 samples (NT rPM>18,000). *M. euphorbiae* is known to be parasitized by hymenopterous wasps
265 belonging to the superfamilies Ichneumonoidea (Braconidae) and Chalcidoidea [71]. In addition,
266 there were some *M. euphorbiae* species-specific reads remaining, which did not map to the pea
267 aphid reference genome but showed some homology to other aphid species (Insecta:
268 Hemiptera).
269
270 We detected the presence of two insect-specific viruses in our metatranscriptome data (Figure
271 1C). The highest number of hits matched a previously described insect-specific Flavivirus,
272 called Macrosiphum euphorbiae virus 1 (MeV-1), which we detected in two samples (NT rPM =
273 234 and 4055 for Me112 and Me202, respectively). We also detected viral hits to an insect-
274 specific Ambidensovirus (Me202 and Me152; NT rPM>60). Other viral reads in our samples
275 included a Bracovirus in one of the samples that was parasitized with the Brachonidae wasp
276 (Me202; NT rPM=1) and a Tombusvirus (Me152; NT rPM=2.9), a family of plant pathogenic
277 viruses with a single-stranded positive-sense RNA genome. Lastly, we detected two phage
278 genera, the *Hamitonella*-specific phage APSE (NT rPM>310) in the same samples found
279 positive for this symbiont (Me112 and Me202) and Acinetobacter phage (NT rPM 0.5-18), a
280 bacteriophage largely prevalent in the environment [72].
281
282 **Figure 1.** Details of the per aphid breakdown of non-host reads aligning to specific bacteria
283 (**Figure 1a**), eukaryotic (**Figure 1b**), and viral (**Figure 1c**) taxa. Reads per million aligning to the
284 nucleotide database (NT rPM) used as the quantitative metric in the heatmaps (see
285 supplementary files heatmap_metrics.csv for metric details).
286

287 **Comprehensive and quantitative analysis of insect-specific viruses:** Using the CZ ID
288 platform, we aligned five assembled contigs to the MeV-1 reference genome (NCBI accession
289 KT309079) and found that they ranged between 85.8-97.2% nucleotide identity to the reference
290 genome (Figure 2A). Our transcriptome retrieved 17,397 informative nucleotides allowing the
291 assembly of a nearly complete genome for MeV-1. Our MeV-1 consensus genome has a
292 coverage breadth of 79% and a coverage depth of 673.2x (NCBI accession OQ504571)
293 (supplementary figure MeV1_coverage.tif). This single-stranded positive-sense RNA genome
294 contains a single large ORF encoding a polyprotein of 7,333 amino acids, which is subsequently
295 processed to generate structural and non-structural proteins [73]. Previous analysis indicated
296 that the polyprotein motifs of MeV-1 helicase, methyltransferase, and RdRp are similar to
297 domains in other *Flavivirus* (family Flaviviridae) [21, 30]. The characteristic secondary structures
298 (RNA stem-loop) in *Flavivirus* genomes most likely contributed to the 5,283 missing bases in our
299 MeV-1 consensus genome assembly [74].

300

301 **Figure 2.** Assembled *M. euphorbiae* transcriptome contigs aligning to previously described
302 insect Flavivirus (**Figure 2a**) and Ambidensovirus (**Figure 2b**) (see supplementary files
303 contigs_CZID.fasta and contigs_TRAVIS.fasta for sequence details).

304

305 In addition, using the CZ ID platform, we detected two contigs with 80% similarity to the non-
306 structural protein 1 (NS1) of Dysaphis plantaginea Densovirus (DplDNV), single-stranded DNA
307 insect-specific *Ambidensovirus* (family Parvoviridae) (supplementary file contigs_CZID.fasta).
308 Due to the lack of a publicly available genome or partial viral sequences of Macrosiphum
309 euphorbiae virus 2 (MeV-2), an Ambidensovirus previously described in the same aphid species
310 [24], we were not able to explore the homology between both viruses. Therefore, we conducted
311 a more extensive analysis of our RNAseq data using TRAVIS, a consistency-based virus
312 detection pipeline for sensitive mass screening of transcriptomic data directed toward
313 Parvoviridae proteins. While degrees of sequence identity between Densovirinae (a subfamily of
314 viral species exclusively infecting arthropods) is very low, viral species often express NS1 and
315 VP proteins, which are useful for parvovirus phylogenetic inferences [75]. We used the seven
316 viral hits provided in individual Trinity contigs flagged by TRAVIS (supplementary file
317 contigs_TRAVIS.fasta) to identify the ORF orientation and similarity and to construct a
318 hypothetical genome assembly using DplDNV as the closest reference available (Figure 2B).
319 We found three contigs with 68.8% to 81.3% similarity to the non-structural ORF1 (encoding for
320 the NS1 protein) and two contigs with 68.8% to 86.2% similarity to the structural ORF (encoding
321 for the VP protein). None of the assembled contigs showed similarity to DplDNV ORF2
322 (encoding for the NS2 protein). We only detected 70% similarity with the ORF2 of a distantly
323 related Ambidensovirus (NCBI accession AMG693112), which genomic organization differs
324 from previously reported aphid densoviruses [21]. Importantly, all densoviral NS1-like
325 sequences also showed a high nucleotide similarity (72-85%) to the pea aphid APNS-2 (NCBI
326 accession NC_042493.1 and NC_042494.1), an EVE that contributes to wing phenotypic
327 plasticity in this species [35].

328

329 **Insect-specific virus frequency in natural populations:** To further investigate the infection
330 frequency of MeV-1 and MeV-2 infections in natural populations, we used a PCR approach to
331 screen 23 individual adult aphids collected during 2022 as well as aphids from our colonized
332 *Macrosiphum* line (Me57). We found only 13 field-collected aphids positive for MeV-1 (54.2%)
333 and 21 aphids (87.5%) positive for MeV-2, including the colonized individuals (Figure 3). We
334 also tested the cDNA of field-collected aphids (previously screened for MeV-1) for the presence
335 of *Hamiltonella defensa* and found that 54.2% of the aphids (n=13) were harboring this bacterial
336 symbiont. We found that 41.7% of individuals (n=10) shared a co-infection between this
337 Flavivirus and *Hamiltonella* (Figure 3), but this

338     association was not statistically significant (p-value= 0.078; r= 0.375).

339

340     **Figure 3.** Frequency of Macrosiphum euphorbiae virus 1 (MeV-1), Macrosiphum euphorbiae
341     virus 2 (MeV-2), and *Hamiltonella denfesa* infections in wild-collected (n=23) and colonized
342     (n=1) aphids. All samples tested using cDNA for PCR screenings.

343

344     **Genome sequencing for analysis of endogenous viral elements (EVEs):** Our laboratory line
345     (Me57) was found to be PCR positive for MeV-2, and we then used DNA sequences from a
346     pooled sample of Me57 aphids to look for viral reads. We used the CZ ID platform as above to
347     identify viral taxa using the Illumina DNA reads from our colonized Me57 aphid line.
348     Surprisingly, we detected only a single contig with a low number of Ambidensoviral hits (NT
349     rPM>0.329), which also showed 79.0% similarity to the DpIDNV NS1 viral protein and 84.34%
350     similarity to an uncharacterized genomic transcript in pea aphids (NCBI accession
351     XM_029492170.1). Since both of our transcriptome and genomic data were unable to recover a
352     complete or near-to-complete Ambidensovirus genome, we then suspected that these viral
353     reads could correspond instead to actively transcribed EVEs, as previously reported in other
354     closely related aphid species [33, 35].

355

356     To determine with certainty whether the Ambidensovirus hits found in our transcriptome data
357     corresponded to an actively transcribed EVE, we assembled the first *M. euphorbiae* genome
358     publicly available. We obtained a total of 4,223,264 nanopore reads (at an average of 5.21kb)
359     and 35,578,886 Illumina reads (PE 150bp) from sequencing. After assembly, haplotig purging,
360     polishing, and manual removal of plant and bacterial contigs, our assembly contained 2,176
361     contigs with an N50 length of 665kb and a total length of 545.7 Mb (Figure 4A). *M. euphorbiae*
362     has a similar GC content (29.96%; Figure 4B) to other sequenced aphids (e.g., *A. pisum* at
363     29.6%, *M. persicae* at 30.1%, and *A. glycines* at 27.8%) [76, 77]. The size of our assembly is
364     close to a recent estimation of the *M. euphorbia* genome size based on flow cytometry which
365     was estimated at 531.7 Mb [76]. Similarly, an analysis of single-copy orthologs showed our
366     assembly contains 98.5% complete BUSCOs, with 94% present in single copies and 4.5%
367     duplicated (Figure 4C). An additional 1.2% of BUSCOs are fragmented, and 0.3% are missing.
368     Together these results suggest that this draft of the genome is highly complete.

369

370     **Figure 4.** M. *euphorbiae* genome assembly metrics (**Figure 4a**), GC content and coverage
371     (**Figure 4b**), and BUSCO metrics (**Figure 4c**).

372

373     We used the genome as a reference to screen for the seven individual Trinity contigs flagged by
374     TRAVIS as potential Ambidensovirus in our previous analysis. Initially, we selected hits with E-
375     values < $1.10^{-3}$ [33]; however, most of the 3,044 hits represent shorter sequences rather than
376     the actual transcript length (see supplementary file expressed_densoviral_EVEs.csv); therefore,
377     we restricted the search to matches consistently to the entire length of each transcript and E-
378     values=0 (Table 1). No full-length hits in the genome were found for the two largest viral contigs
379     assembled from transcriptome data (contig3 and contig4); instead, the best hits for these two
380     contigs corresponded to 16-17% of the total length. In insects, the EVE repertoire varies
381     between distinct populations of a given species and, in some cases, even between individuals
382     within the same population [78]. This phenomenon explains why all the field aphid samples
383     (n=3) that tested negative for MeV-2 by PCR amplified a product of approximately 500 bp,
384     which is about half of the expected size reported for the primers used. Given that the genome
385     assemblies and RNAseq data sets were derived from different aphid strains, it is not surprising
386     the wide range of partial-length Ambidensovirus hits obtained in our analysis. However, we are
387     confident that five full-length viral transcripts are constitutively expressed from three regions of
388     the *M. euphorbiae* genome (tig00030708_pilon, tig00029914_pilon, and tig00027226_pilon).

389

390 **Table 1.** List of Ambidensoviral transcripts and the corresponding integrations in *M. euphorbiae*
391 genome.

| Transcriptome contig | Transcript length | Percentage of identical sites | Hit end | Hit start | Genome contig | Query end | Query start |
|---|---|---|---|---|---|---|---|
| Travis_contig1 | 783 | 96.70% | 783 | 1 | tig00030708_pilon | 198038 | 197258 |
| Travis_contig1 | 783 | 98.90% | 1 | 783 | tig00029914_pilon | 60345 | 59559 |
| Travis_contig2 | 466 | 96.20% | 416 | 1 | tig00030708_pilon | 198433 | 198018 |
| Travis_contig2 | 466 | 99.80% | 1 | 466 | tig00029914_pilon | 59579 | 59114 |
| Travis_contig3 | 2155 | - | - | - | - | - | - |
| Travis_contig4 | 2878 | - | - | - | - | - | - |
| Travis_contig5 | 1174 | 99.90% | 1 | 1174 | tig00030708_pilon | 92758 | 91585 |
| Travis_contig6 | 1040 | 99.80% | 1 | 1040 | tig00030708_pilon | 85562 | 84525 |
| Travis_contig7 | 635 | 100.00% | 635 | 1 | tig00030708_pilon | 86191 | 85557 |
| Travis_contig7 | 635 | 84.90% | 635 | 1 | tig00027226_pilon | 138266 | 137632 |

392
393
394 DISCUSSION
395
396 RNAseq is becoming an essential tool for virus discovery. Our study illustrates how endogenous
397 viral elements in insect genomes can be an obstacle to using RNAseq for characterizing viral
398 diversity in arthropods. We used rRNA-depleted RNAseq along with bioinformatic tools to
399 characterize the virome of an important insect pest species, the potato aphid *Macrosiphum*
400 *euphorbiae.* Our analysis found two insect-specific viruses from the families Flavivirus and
401 Ambidensovirus described in previous RNAseq studies. However, by sequencing and
402 assembling the genome of this insect using long-read sequencing, we found that the
403 Ambidensovirus is a transcriptionally active EVE rather than an exogenous virus. Endogenous
404 viral elements encoded in the host genome are abundant in arthropod genomes, and thus EVE
405 sequences in RNAseq studies are an important consideration for future studies of viral diversity
406 in arthropods.
407
408 Densoviral EVEs have been shown to be transcriptionally active in two other aphid species:
409 *Myzus persicae* and *A. pisum*. In pea aphids, two copies of a transcribed densoviral non-
410 structural protein (termed the "APNS" genes) were found to be upregulated in response to
411 crowded conditions and to be functionally linked to the plastic production of wings [35]. These
412 genes had close homology with the non-structural genes of Dysaphis plantaginea densovirus
413 (DplDNV), which, when infecting rosy apple aphids, causes them to be winged [79], suggesting
414 the function of these viral genes had been conserved after endogenization. The transcribed
415 EVEs we found in *M. euphorbiae* have significant homology to the pea aphid APNS genes, and
416 it seems likely that these genes may also be playing a role in wing plasticity in *M. euphorbiae*
417 though additional data is needed.
418
419 Our study contributes to the growing list of sequenced aphid genomes, which together show
420 that transcribed densoviral EVEs are common in this insect group [33, 35, 37, 80]. Most
421 identified EVEs in insect genomes correspond to unclassified single-stranded RNA viruses and
422 viruses belonging to the families Rhabdoviridae and Parvoviridae [78]. Unlike RNA viruses,

423 which may produce abundant short mRNAs that favor virus endogenization [20], Parvoviruses
424 undergo a double-stranded DNA intermediate during nuclear replication, which along with the
425 endonuclease activity of NS1 and the eukaryote double-stranded break repair mechanism may
426 largely favor endogenization of this virus family [81, 82]. Previous studies have estimated that
427 around 10% of the parvoviruses described in animals are likely integrated into host genomes,
428 but in most cases, the EVE status remains uncertain due to unavailable or incomplete genomes
429 for those species in which transcriptome data is available [75]. Multiple recent studies have
430 described the presence of "new" densoviruses in aphid's transcriptome [23, 26, 83]; however,
431 our combined transcriptomic and genomic analyses suggest that some of those viral transcripts
432 may likely correspond to actively transcribed EVEs instead of heritable exogenous viruses
433 infecting aphids at very high rates.

434

435 Last, our study shed light on the biology of MeV-1, an insect-specific Flavivirus (family
436 Flaviviridae), previously characterized by RNAseq studies of *M. euphorbiae* populations
437 collected in France [30]. We found that this virus, contrary to previous reports, is present in a
438 North American population of *M. euphorbiae*, and we found that it is highly prevalent. By
439 assembling the genome of MeV-1 from our RNAseq data, we found that our local population is
440 infected with a potentially distinct viral strain from previous studies. No obvious infection
441 symptoms or abnormal phenotypes were observed in MeV-1-infected *M. euphorbiae* adults, and
442 future studies are needed to determine what phenotypic effects this virus has on its host. Other
443 heritable viruses have been found to interact with the secondary symbiotic bacteria found in
444 aphids [84-86] but we did not find significant patterns of co-infections with the bacterial symbiont
445 *Hamiltonella defensa*.

446

447 EVEs are common in insect genomes, and our results highlight this widespread challenge in
448 studying insect viromes. Our study further emphasizes how combining sequencing
449 methodologies is necessary to overcome the potential pitfalls of only RNAseq-based viral
450 discovery. Careful consideration of the biological characteristics and genome structure of
451 viruses discovered through RNAseq is essential [87]. In aphids and other widely study systems,
452 the development of cultured cell lines is also imperative to isolate viral species described by
453 sequence-based methods, to characterize viral replication, and for use in large-scale virus
454 production that will facilitate future investigation of the complex interaction of aphid viruses and
455 their hosts [21].

456

457 The relatively high transcription level of some EVEs suggests that viral integration may have
458 important biological implications for the fitness of aphids. Likewise, uncovering the phenotypic
459 effects of accurately described insect-specific viruses may also show promising targets for
460 alternative control strategies of agriculturally destructive organisms while providing important
461 foundational resources in the study of host-virus dynamics. Research efforts need to be done on
462 the evolutionary dynamics of heritable viruses to better understand how they are acting as
463 hidden drivers of host phenotypes.

464

473

485    REFERENCES
486
487    [1] Koonin, E.V. & Dolja, V.V. 2018 Metaviromics: a tectonic shift in understanding virus
488    evolution. *Virus Research* **246**, A1-A3. (doi:https://doi.org/10.1016/j.virusres.2018.02.001).
489    [2] Zhang, Y.-Z., Shi, M. & Holmes, E.C. 2018 Using Metagenomics to Characterize an
490    Expanding Virosphere. *Cell* **172**, 1168-1172. (doi:https://doi.org/10.1016/j.cell.2018.02.043).
491    [3] Greninger, A.L. 2018 A decade of RNA virus metagenomics is (not) enough. *Virus Research*
492    **244**, 218-229. (doi:https://doi.org/10.1016/j.virusres.2017.10.014).
493    [4] Stork, N.E. 2018 How Many Species of Insects and Other Terrestrial Arthropods Are There
494    on Earth? *Annual Review of Entomology* **63**, 31-45. (doi:10.1146/annurev-ento-020117-
495    043348).
496    [5] Coatsworth, H., Bozic, J., Carrillo, J., Buckner, E.A., Rivers, A.R., Dinglasan, R.R. & Mathias,
497    D.K. 2022 Intrinsic variation in the vertically transmitted core virome of the mosquito Aedes
498    aegypti. *Mol Ecol* **31**, 2545-2561. (doi:10.1111/mec.16412).
499    [6] Longdon, B. & Jiggins, F.M. 2012 Vertically transmitted viral endosymbionts of insects: do
500    sigma viruses walk alone? *Proc Biol Sci* **279**, 3889-3898. (doi:10.1098/rspb.2012.1208).
501    [7] Longdon, B., Day, J.P., Schulz, N., Leftwich, P.T., de Jong, M.A., Breuker, C.J., Gibbs, M.,
502    Obbard, D.J., Wilfert, L., Smith, S.C., et al. 2017 Vertically transmitted rhabdoviruses are found
503    across three insect families and have dynamic interactions with their hosts. *Proc Biol Sci* **284**.
504    (doi:10.1098/rspb.2016.2381).
505    [8] González, R., Butković, A. & Elena, S.F. 2020 Chapter Three - From foes to friends: Viral
506    infections expand the limits of host phenotypic plasticity. In *Advances in Virus Research* (eds.
507    M. Kielian, T.C. Mettenleiter & M.J. Roossinck), pp. 85-121, Academic Press.
508    [9] Simmonds, P., Aiewsakun, P. & Katzourakis, A. 2019 Prisoners of war — host adaptation
509    and its constraints on virus evolution. *Nature Reviews Microbiology* **17**, 321-328.
510    (doi:10.1038/s41579-018-0120-2).
511    [10] Mauck, K.E. 2016 Variation in virus effects on host plant phenotypes and insect vector
512    behavior: what can it teach us about virus evolution? *Current Opinion in Virology* **21**, 114-123.
513    (doi:https://doi.org/10.1016/j.coviro.2016.09.002).
514    [11] Bolling, B.G., Weaver, S.C., Tesh, R.B. & Vasilakis, N. 2015 Insect-Specific Virus
515    Discovery: Significance for the Arbovirus Community. In *Viruses* (pp. 4911-4928.
516    [12] Li, C.-X., Shi, M., Tian, J.-H., Lin, X.-D., Kang, Y.-J., Chen, L.-J., Qin, X.-C., Xu, J., Holmes,
517    E.C. & Zhang, Y.-Z. 2015 Unprecedented genomic diversity of RNA viruses in arthropods
518    reveals the ancestry of negative-sense RNA viruses. *eLife* **4**, e05378.
519    (doi:10.7554/eLife.05378).
520    [13] Käfer, S., Paraskevopoulou, S., Zirkel, F., Wieseke, N., Donath, A., Petersen, M., Jones,
521    T.C., Liu, S., Zhou, X., Middendorf, M., et al. 2019 Re-assessing the diversity of negative strand
522    RNA viruses in insects. *PLOS Pathogens* **15**, e1008224. (doi:10.1371/journal.ppat.1008224).
523    [14] Liu, S., Chen, Y. & Bonning, B.C. 2015 RNA virus discovery in insects. *Current Opinion in
524    Insect Science* **8**, 54-61. (doi:https://doi.org/10.1016/j.cois.2014.12.005).
525    [15] Wu, H., Pang, R., Cheng, T., Xue, L., Zeng, H., Lei, T., Chen, M., Wu, S., Ding, Y., Zhang,
526    J., et al. 2020 Abundant and Diverse RNA Viruses in Insects Revealed by RNA-Seq Analysis:
527    Ecological and Evolutionary Implications. *mSystems* **5**, e00039-00020.
528    (doi:10.1128/mSystems.00039-20).
529    [16] Gilbert, C. & Belliardo, C. 2022 The diversity of endogenous viral elements in insects. *Curr
530    Opin Insect Sci* **49**, 48-55. (doi:10.1016/j.cois.2021.11.007).
531    [17] Veglia, A.J., Bistolas, K.S.I., Voolstra, C.R., Hume, B.C.C., Planes, S., Allemand, D.,
532    Boissin, E., Wincker, P., Poulain, J., Moulin, C., et al. 2022 Endogenous viral elements reveal
533    associations between a non-retroviral RNA virus and symbiotic dinoflagellate genomes. *bioRxiv*,
534    2022.2004.2011.487905. (doi:10.1101/2022.04.11.487905).

[18] Blair, C.D., Olson, K.E. & Bonizzoni, M. 2020 The widespread occurrence and potential biological roles of endogenous viral elements in insect genomes. *Curr Issues Mol Biol* **34**, 13-30. (doi:10.21775/cimb.034.013).

[19] Frank, J.A. & Feschotte, C. 2017 Co-option of endogenous viral sequences for host cell function. *Curr Opin Virol* **25**, 81-89. (doi:10.1016/j.coviro.2017.07.021).

[20] Holmes, E.C. 2011 The evolution of endogenous viral elements. *Cell Host Microbe* **10**, 368-377. (doi:10.1016/j.chom.2011.09.002).

[21] Guo, Y., Ji, N., Bai, L., Ma, J. & Li, Z. 2022 Aphid viruses: A brief view of a long history. *Frontiers in Insect Science* **2**. (doi:10.3389/finsc.2022.846716).

[22] Ray, S. & Casteel, C.L. 2022 Effector-mediated plant–virus–vector interactions. *The Plant Cell* **34**, 1514-1531. (doi:10.1093/plcell/koac058).

[23] Feng, Y., Krueger, E.N., Liu, S., Dorman, K., Bonning, B.C. & Miller, W.A. 2017 Discovery of Known and Novel Viral Genomes in Soybean Aphid by Deep Sequencing. *Phytobiomes Journal* **1**, 36-45. (doi:10.1094/PBIOMES-11-16-0013-R).

[24] Teixeira, M.A., Sela, N., Atamian, H.S., Bao, E., Chaudhary, R., MacWilliams, J., He, J., Mantelin, S., Girke, T. & Kaloshian, I. 2018 Sequence analysis of the potato aphid *Macrosiphum euphorbiae* transcriptome identified two new viruses. *PLOS ONE* **13**, e0193239. (doi:10.1371/journal.pone.0193239).

[25] Kondo, H., Fujita, M., Hisano, H., Hyodo, K., Andika, I.B. & Suzuki, N. 2020 Virome Analysis of Aphid Populations That Infest the Barley Field: The Discovery of Two Novel Groups of Nege/Kita-Like Viruses and Other Novel RNA Viruses. *Frontiers in Microbiology* **11**. (doi:10.3389/fmicb.2020.00509).

[26] Li, T., Li, H., Wu, Y., Li, S., Yuan, G. & Xu, P. 2022 Identification of a Novel Densovirus in Aphid, and Uncovering the Possible Antiviral Process During Its Infection. *Frontiers in Immunology* **13**. (doi:10.3389/fimmu.2022.905628).

[27] Wamonje, F.O., Michuki, G.N., Braidwood, L.A., Njuguna, J.N., Musembi Mutuku, J., Djikeng, A., Harvey, J.J.W. & Carr, J.P. 2017 Viral metagenomics of aphids present in bean and maize plots on mixed-use farms in Kenya reveals the presence of three dicistroviruses including a novel Big Sioux River virus-like dicistrovirus. *Virology Journal* **14**, 188. (doi:10.1186/s12985-017-0854-x).

[28] Blackman, R.L., Eastop, V.F. & Museum, N.H. 2000 *Aphids on the World's Crops: An Identification and Information Guide*, Wiley.

[29] Xu, Y. & Gray, S.M. 2020 Aphids and their transmitted potato viruses: A continuous challenges in potato crops. *Journal of Integrative Agriculture* **19**, 367-375. (doi:https://doi.org/10.1016/S2095-3119(19)62842-X).

[30] Teixeira, M., Sela, N., Ng, J., Casteel, C.L., Peng, H.-C., Bekal, S., Girke, T., Ghanim, M. & Kaloshian, I. 2016 A novel virus from Macrosiphum euphorbiae with similarities to members of the family Flaviviridae. *Journal of General Virology* **97**, 1261-1271. (doi:https://doi.org/10.1099/jgv.0.000414).

[31] Planelló, R., Llorente, L., Herrero, Ó., Novo, M., Blanco-Sánchez, L., Díaz-Pendón, J.A., Fernández-Muñoz, R., Ferrero, V. & de la Peña, E. 2022 Transcriptome analysis of aphids exposed to glandular trichomes in tomato reveals stress and starvation related responses. *Sci Rep* **12**, 20154. (doi:10.1038/s41598-022-24490-1).

[32] Atamian, H.S., Chaudhary, R., Cin, V.D., Bao, E., Girke, T. & Kaloshian, I. 2013 In planta expression or delivery of potato aphid Macrosiphum euphorbiae effectors Me10 and Me23 enhances aphid fecundity. *Mol Plant Microbe Interact* **26**, 67-74. (doi:10.1094/mpmi-06-12-0144-fi).

[33] Clavijo, G., van Munster, M., Monsion, B., Bochet, N. & Brault, V. 2016 Transcription of densovirus endogenous sequences in the Myzus persicae genome. *Journal of General Virology* **97**, 1000-1009. (doi:https://doi.org/10.1099/jgv.0.000396).
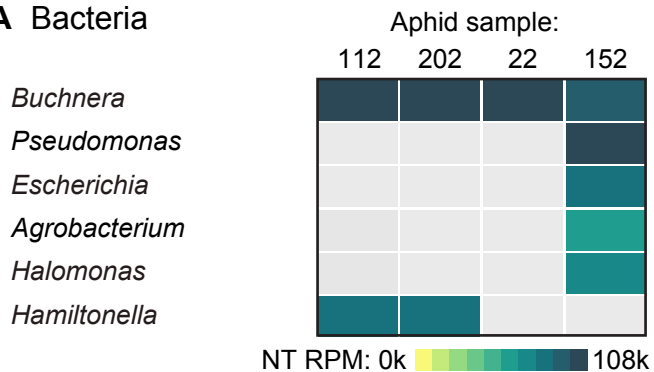
585  [34] Jayasinghe, W.H., Kim, H., Nakada, Y. & Masuta, C. 2021 A plant virus satellite RNA
586  directly accelerates wing formation in its insect vector for spread. *Nature Communications* **12**,
587  7087. (doi:10.1038/s41467-021-27330-4).
588  [35] Parker, B.J. & Brisson, J.A. 2019 A laterally transferred viral gene modifies aphid wing
589  plasticity. *Current Biology* **29**, 2098-2103.e2095. (doi:https://doi.org/10.1016/j.cub.2019.05.041).
590  [36] Shang, F., Niu, J., Ding, B.-Y., Zhang, W., Wei, D.-D., Wei, D., Jiang, H.-B. & Wang, J.-J.
591  2020 The miR-9b microRNA mediates dimorphism and development of wing in aphids.
592  *Proceedings of the National Academy of Sciences* **117**, 8404-8409.
593  (doi:10.1073/pnas.1919204117).
594  [37] Liu, S., Coates, B.S. & Bonning, B.C. 2020 Endogenous viral elements integrated into the
595  genome of the soybean aphid, Aphis glycines. *Insect Biochem Mol Biol* **123**, 103405.
596  (doi:10.1016/j.ibmb.2020.103405).
597  [38] Foottit, R.G., Maw, H.E., CD, V.O.N.D. & Hebert, P.D. 2008 Species identification of aphids
598  (Insecta: Hemiptera: Aphididae) through DNA barcodes. *Mol Ecol Resour* **8**, 1189-1201.
599  (doi:10.1111/j.1755-0998.2008.02297.x).
600  [39] Henry, L.M., Peccoud, J., Simon, J.C., Hadfield, J.D., Maiden, M.J., Ferrari, J. & Godfray,
601  H.C. 2013 Horizontally transmitted symbionts and host colonization of ecological niches. *Curr*
602  *Biol* **23**, 1713-1717. (doi:10.1016/j.cub.2013.07.029).
603  [40] McLean, A.H.C., Hrcek, J., Parker, B.J., Mathe-Hubert, H., Kaech, H., Paine, C. & Godfray,
604  H.C.J. 2020 Multiple phenotypes conferred by a single insect symbiont are independent. *Proc*
605  *Biol Sci* **287**, 20200562. (doi:10.1098/rspb.2020.0562).
606  [41] Goldstein, E.B., de Anda Acosta, Y., Henry, L.M. & Parker, B.J. 2022 Variation in density,
607  immune gene suppression, and co-infection outcomes among strains of the aphid endosymbiont
608  <em>Regiella insecticola</em>. *bioRxiv*, 2022.2008.2028.505589.
609  (doi:10.1101/2022.08.28.505589).
610  [42] Bender, W., Spierer, P., Hogness, D.S. & Chambon, P. 1983 Chromosomal walking and
611  jumping to isolate DNA from the Ace and rosy loci and the bithorax complex in Drosophila
612  melanogaster. *Journal of Molecular Biology* **168**, 17-33. (doi:https://doi.org/10.1016/S0022-
613  2836(83)80320-9).
614  [43] Henry, L.M., Maiden, M.C.J., Ferrari, J. & Godfray, H.C.J. 2015 Insect life history and the
615  evolution of bacterial mutualism. *Ecology Letters* **18**, 516-525.
616  (doi:https://doi.org/10.1111/ele.12425).
617  [44] Schmieder, R., Lim, Y.W. & Edwards, R. 2012 Identification and removal of ribosomal RNA
618  sequences from metatranscriptomes. *Bioinformatics* **28**, 433-435.
619  (doi:10.1093/bioinformatics/btr669).
620  [45] Pruesse, E., Quast, C., Knittel, K., Fuchs, B.M., Ludwig, W., Peplies, J. & Glöckner, F.O.
621  2007 SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA
622  sequence data compatible with ARB. *Nucleic Acids Research* **35**, 7188-7196.
623  (doi:10.1093/nar/gkm864).
624  [46] Kalantar, K.L., Carvalho, T., de Bourcy, C.F.A., Dimitrov, B., Dingle, G., Egger, R., Han, J.,
625  Holmes, O.B., Juan, Y.-F., King, R., et al. 2020 IDseq—An open source cloud-based pipeline
626  and analysis service for metagenomic pathogen detection and monitoring. *GigaScience* **9**.
627  (doi:10.1093/gigascience/giaa111).
628  [47] Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P.,
629  Chaisson, M. & Gingeras, T.R. 2012 STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*
630  **29**, 15-21. (doi:10.1093/bioinformatics/bts635).
631  [48] Bolger, A.M., Lohse, M. & Usadel, B. 2014 Trimmomatic: a flexible trimmer for Illumina
632  sequence data. *Bioinformatics* **30**, 2114-2120. (doi:10.1093/bioinformatics/btu170).
633  [49] Ruby, J.G., Bellare, P. & Derisi, J.L. 2013 PRICE: software for the targeted assembly of
634  components of (Meta) genomic sequence data. *G3 (Bethesda)* **3**, 865-880.
635  (doi:10.1534/g3.113.005967).

[50] Li, H. 2018 Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094-3100. (doi:10.1093/bioinformatics/bty191).

[51] Buchfink, B., Reuter, K. & Drost, H.G. 2021 Sensitive protein alignments at tree-of-life scale using DIAMOND. *Nat Methods* **18**, 366-368. (doi:10.1038/s41592-021-01101-x).

[52] Bankevich, A., Nurk, S., Antipov, D., Gurevich, A.A., Dvorkin, M., Kulikov, A.S., Lesin, V.M., Nikolenko, S.I., Pham, S., Prjibelski, A.D., et al. 2012 SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol* **19**, 455-477. (doi:10.1089/cmb.2012.0021).

[53] Langmead, B. & Salzberg, S.L. 2012 Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357-359. (doi:10.1038/nmeth.1923).

[54] Bohl, J.A., Lay, S., Chea, S., Ahyong, V., Parker, D.M., Gallagher, S., Fintzi, J., Man, S., Ponce, A., Sreng, S., et al. 2022 Discovering disease-causing pathogens in resource-scarce Southeast Asia using a global metagenomic pathogen monitoring system. *Proc Natl Acad Sci U S A* **119**, e2115285119. (doi:10.1073/pnas.2115285119).

[55] Batson, J., Dudas, G., Haas-Stapleton, E., Kistler, A.L., Li, L.M., Logan, P., Ratnasiri, K. & Retallack, H. 2021 Single mosquito metatranscriptomics identifies vectors, emerging pathogens and reservoirs in one assay. *eLife* **10**, e68353. (doi:10.7554/eLife.68353).

[56] Andrews, S. 2010 FastQC: A Quality Control Tool for High Throughput Sequence Data [Online]. (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/.

[57] Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X., Fan, L., Raychowdhury, R., Zeng, Q., et al. 2011 Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nature Biotechnology* **29**, 644-652. (doi:10.1038/nbt.1883).

[58] Wheeler, T.J. & Eddy, S.R. 2013 nhmmer: DNA homology search with profile HMMs. *Bioinformatics* **29**, 2487-2489. (doi:10.1093/bioinformatics/btt403).

[59] Steinegger, M. & Söding, J. 2017 MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat Biotechnol* **35**, 1026-1028. (doi:10.1038/nbt.3988).

[60] Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W. & Lipman, D.J. 1997 Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* **25**, 3389-3402. (doi:10.1093/nar/25.17.3389).

[61] Martin, M. 2011 Cutadapt removes adapter sequences from high-throughput sequencing reads. *2011* **17**, 3. (doi:10.14806/ej.17.1.200).

[62] Grubaugh, N.D., Gangavarapu, K., Quick, J., Matteson, N.L., De Jesus, J.G., Main, B.J., Tan, A.L., Paul, L.M., Brackney, D.E., Grewal, S., et al. 2019 An amplicon-based sequencing framework for accurately measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biology* **20**, 8. (doi:10.1186/s13059-018-1618-7).

[63] Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G. & Durbin, R. 2009 The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079. (doi:10.1093/bioinformatics/btp352).

[64] Koren, S., Walenz, B.P., Berlin, K., Miller, J.R., Bergman, N.H. & Phillippy, A.M. 2017 Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res* **27**, 722-736. (doi:10.1101/gr.215087.116).

[65] Roach, M.J., Schmidt, S.A. & Borneman, A.R. 2018 Purge Haplotigs: allelic contig reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics* **19**, 460. (doi:10.1186/s12859-018-2485-7).

[66] Walker, B.J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C.A., Zeng, Q., Wortman, J., Young, S.K., et al. 2014 Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement. *PLOS ONE* **9**, e112963. (doi:10.1371/journal.pone.0112963).
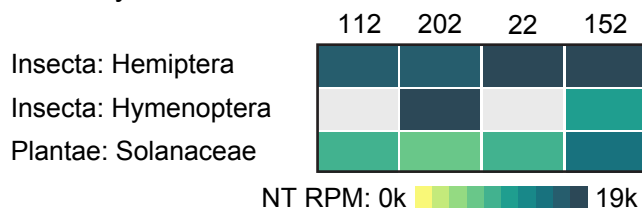
685 [67] Challis, R., Richards, E., Rajan, J., Cochrane, G. & Blaxter, M. 2020 BlobToolKit –
686 Interactive Quality Assessment of Genome Assemblies. *G3 Genes|Genomes|Genetics* **10**,
687 1361-1374. (doi:10.1534/g3.119.400908).
688 [68] Sayers, E.W., Bolton, E.E., Brister, J.R., Canese, K., Chan, J., Comeau, D.C., Connor, R.,
689 Funk, K., Kelly, C., Kim, S., et al. 2022 Database resources of the national center for
690 biotechnology information. *Nucleic Acids Res* **50**, D20-d26. (doi:10.1093/nar/gkab1112).
691 [69] Manni, M., Berkeley, M.R., Seppey, M., Simão, F.A. & Zdobnov, E.M. 2021 BUSCO
692 Update: Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic
693 Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes. *Molecular Biology and*
694 *Evolution* **38**, 4647-4654. (doi:10.1093/molbev/msab199).
695 [70] Levy Karin, E., Mirdita, M. & Söding, J. 2020 MetaEuk—sensitive, high-throughput gene
696 discovery, and annotation for large-scale eukaryotic metagenomics. *Microbiome* **8**, 48.
697 (doi:10.1186/s40168-020-00808-x).
698 [71] Clarke, H.V. 2013 Genotypic and endosymbiont-mediated variation in parasitoid
699 susceptibility and other fitness traits of the potato aphid, Macrosiphum euphorbiae.  (
700 [72] Turner, D., Ackermann, H.W., Kropinski, A.M., Lavigne, R., Sutton, J.M. & Reynolds, D.M.
701 2017 Comparative Analysis of 37 Acinetobacter Bacteriophages. *Viruses* **10**.
702 (doi:10.3390/v10010005).
703 [73] Blitvich, B.J. & Firth, A.E. 2015 Insect-specific flaviviruses: a systematic review of their
704 discovery, host range, mode of transmission, superinfection exclusion potential and genomic
705 organization. *Viruses* **7**, 1927-1959. (doi:10.3390/v7041927).
706 [74] Mazeaud, C., Freppel, W. & Chatel-Chaix, L. 2018 The Multiples Fates of the Flavivirus
707 RNA Genome During Pathogenesis. *Frontiers in Genetics* **9**. (doi:10.3389/fgene.2018.00595).
708 [75] François, S., Filloux, D., Roumagnac, P., Bigot, D., Gayral, P., Martin, D.P., Froissart, R. &
709 Ogliastro, M. 2016 Discovery of parvovirus-related sequences in an unexpected broad range of
710 animals. *Scientific Reports* **6**, 30880. (doi:10.1038/srep30880).
711 [76] Wenger, J.A., Cassone, B.J., Legeai, F., Johnston, J.S., Bansal, R., Yates, A.D., Coates,
712 B.S., Pavinato, V.A. & Michel, A. 2017 Whole genome sequence of the soybean aphid, Aphis
713 glycines. *Insect Biochem Mol Biol.* (doi:10.1016/j.ibmb.2017.01.005).
714 [77] International Aphid Genomics, C. 2010 Genome sequence of the pea aphid Acyrthosiphon
715 pisum. *PLoS Biol* **8**, e1000313. (doi:10.1371/journal.pbio.1000313).
716 [78] Horst, A.M.t., Nigg, J.C., Dekker, F.M. & Falk, B.W. 2019 Endogenous Viral Elements Are
717 Widespread in Arthropod Genomes and Commonly Give Rise to PIWI-Interacting RNAs.
718 *Journal of Virology* **93**, e02124-02118. (doi:doi:10.1128/JVI.02124-18).
719 [79] Ryabov, E.V., Keane, G., Naish, N., Evered, C. & Winstanley, D. 2009 Densovirus induces
720 winged morphs in asexual clones of the rosy apple aphid, *Dysaphis plantaginea. Proceedings of*
721 *the National Academy of Sciences* **106**, 8465-8470. (doi:doi:10.1073/pnas.0901389106).
722 [80] Nigg, J.C., Kuo, Y.W. & Falk, B.W. 2020 Endogenous viral element-derived piwi-interacting
723 RNAs (piRNAs) are not required for production of ping-pong-dependent piRNAs from
724 Diaphorina citri Densovirus. *mBio* **11**. (doi:10.1128/mBio.02209-20).
725 [81] Kapoor, A., Simmonds, P. & Lipkin, W.I. 2010 Discovery and characterization of
726 mammalian endogenous parvoviruses. *J Virol* **84**, 12628-12635. (doi:10.1128/jvi.01732-10).
727 [82] Liu, H., Fu, Y., Xie, J., Cheng, J., Ghabrial, S.A., Li, G., Peng, Y., Yi, X. & Jiang, D. 2011
728 Widespread endogenization of densoviruses and parvoviruses in animal and human genomes.
729 *J Virol* **85**, 9863-9876. (doi:10.1128/jvi.00828-11).
730 [83] Pinheiro, P.V., Wilson, J.R., Xu, Y., Zheng, Y., Rebelo, A.R., Fattah-Hosseini, S., Kruse, A.,
731 Dos Silva, R.S., Xu, Y., Kramer, M., et al. 2019 Plant viruses transmitted in two different modes
732 produce differing effects on small RNA-mediated processes in their aphid vector. *Phytobiomes*
733 *Journal* **3**, 71-81. (doi:10.1094/pbiomes-10-18-0045-r).
734 [84] Altinli, M., Schnettler, E. & Sicard, M. 2021 Symbiotic Interactions Between Mosquitoes and
735 Mosquito Viruses. *Front Cell Infect Microbiol* **11**, 694020. (doi:10.3389/fcimb.2021.694020).

[85] Wu, W., Shan, H.W., Li, J.M., Zhang, C.X., Chen, J.P. & Mao, Q. 2022 Roles of Bacterial Symbionts in Transmission of Plant Virus by Hemipteran Vectors. *Front Microbiol* **13**, 805352. (doi:10.3389/fmicb.2022.805352).

[86] Jia, D., Mao, Q., Chen, Y., Liu, Y., Chen, Q., Wu, W., Zhang, X., Chen, H., Li, Y. & Wei, T. 2017 Insect symbiotic bacteria harbour viral pathogens for transovarial transmission. *Nature Microbiology* **2**, 17025. (doi:10.1038/nmicrobiol.2017.25).

[87] Depledge, D.P., Mohr, I. & Wilson, A.C. 2019 Going the Distance: Optimizing RNA-Seq Strategies for Transcriptomic Analysis of Complex Viral Genomes. *J Virol* **93**. (doi:10.1128/jvi.01342-18).

**A: Macrosiphum euphorbiae virus 1 (NCBI Reference Sequence NC_028137.1)**

5'    1000bp   3'

Genomic polyprotein

C | PR | M | E | NS1/NS2A | NS2B/NS3 | NS4A/B | NS5

Contig 1    92.3%
Contig 2    85.8%
Contig 3    97.6%
Contig 4    96.6%
Contig 5    97.2%

**B: Dysaphis plantaginea Densovirus (NCBI Reference Sequence NC_034532.1)**

1000bp

NS1    VP
NS2    VP

Unaligned

Contig 1
Contig 2    66.9%
Contig 3    80.7%
Contig 4    80.7%
Contig 5
Contig 6    69.2%
Contig 7    81%

MeV-1    *Hamiltonella*

1

10    3

2

6

MeV-2

No infection

2

**A**

550M
490M
90
440M 80
380M
70
60
330M
270M 50
40 220M
30 160M
20
110M
10
55M
0%

10k
100k
1.0M

GC (30.0%)
AT (70.0%)

Log10 scaffold count (total 2.2k)
Scaffold length (total 550M)
Longest scaffold (4.3M)
N50 length (660k)
N90 length(120k)

**B**

Scaffold Size:
1Mbp
100kb
10kb

1000

100

10

1

Map Coverage (Illumina sequencing)

25    30    35    40    45    50    55

Percent GC Content

**C**

BUSCO (hemiptera_odb10)

Complete (94%)
Duplicated (4.5%)
Fragmented (1.2%)
Missing (0.3%)