

1 **Immune suppressive extracellular vesicle proteins of *Leptopilina heterotoma* are**  
2 **encoded in the wasp genome**

3 Brian Wey<sup>\*‡</sup>, Mary Ellen Heavner<sup>\*†\*\*</sup>, Kameron T. Wittmeyer<sup>††</sup>, Thomas Briese<sup>§</sup>,  
4 Keith R. Hopper<sup>††</sup>, Shubha Govind<sup>\*†‡</sup>

5 <sup>\*</sup>Biology Department, The City College of New York, 160 Convent Avenue, New York,  
6 NY 10031, USA

7 <sup>†</sup> PhD Program in Biochemistry, The Graduate Center of the City University of New  
8 York, 365 Fifth Avenue, New York, NY 10016, USA

9 <sup>‡</sup> PhD Program in Biology, The Graduate Center of the City University of New York, 365  
10 Fifth Avenue, New York, NY 10016, USA

11 <sup>§</sup> Center of Infection and Immunity, and Department of Epidemiology, Mailman School  
12 of Public Health, Columbia University, NY, NY 10032

13 <sup>\*\*</sup>Laboratory of Host-Pathogen Biology, Rockefeller University, 1230 York Ave, New  
14 York, NY 10065

15 <sup>††</sup> USDA-ARS, Beneficial Insect Introductions Research Unit, Newark, DE 19713

16

17 **Short Title:** *L. heterotoma* MSEV proteins are encoded in the genome

18 **Key words:** Extracellular vesicle, whole genome sequencing, *Leptopilina heterotoma*,

19 endoparasitoid wasp, *Drosophila*, VLP, host-parasite, organelle, immune suppression

20 **Corresponding Author:**

21 Shubha Govind

22 160 Convent Avenue

23 Biology Department

24 The City College of New York

25 New York, New York, 10031

26 sgovind@ccny.cuny.edu

27 Phone number: 212-650-8571

28 BW ORCID: <https://orcid.org/0000-0003-0895-5584>

29 SG ORCID: <https://orcid.org/0000-0002-6436-639X>

## **ABSTRACT**

30  
31 *Leptopilina heterotoma* are obligate parasitoid wasps that develop in the body of their  
32 *Drosophila* hosts. During oviposition, female wasps introduce venom into the larval  
33 hosts' body cavity. The venom contains discrete, 300 nm-wide, mixed-strategy  
34 extracellular vesicles (MSEVs), until recently referred to as virus-like particles. While the  
35 crucial immune suppressive functions of *L. heterotoma* MSEVs have remained  
36 undisputed, their biotic nature and origin still remain controversial. In recent proteomics  
37 analyses of *L. heterotoma* MSEVs, we identified 161 proteins in three classes:  
38 conserved eukaryotic proteins, infection and immunity related proteins, and proteins  
39 without clear annotation. Here we report 246 additional proteins from the *L. heterotoma*  
40 MSEV proteome. An enrichment analysis of the entire proteome supports vesicular  
41 nature of these structures. Sequences for more than 90% of these proteins are present  
42 in the whole-body transcriptome. Sequencing and *de novo* assembly of the 460 Mb-  
43 sized *L. heterotoma* genome revealed 90% of MSEV proteins have coding regions  
44 within the genomic scaffolds. Altogether, these results explain the stable association of  
45 MSEVs with their wasps, and like other wasp structures, their vertical inheritance. While  
46 our results do not rule out a viral origin of MSEVs, they suggest that a similar strategy  
47 for co-opting cellular machinery for immune suppression may be shared by other wasps  
48 to gain advantage over their hosts. These results are relevant to our understanding of  
49 the evolution of figitid and related wasp species.

50

## INTRODUCTION

51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73

Parasitic wasps are amongst the most abundant insects; they are vital to biodiversity and contribute to biological control of agricultural pests (Narendran 2001; Rodriguez et al. 2013). A common strategy for reproductive success of parasitic wasps is suppression of immunity in their larval hosts. Parasitic wasps produce viruses or virus-like particles in tissues associated with the ovary. Wasps of the Ichneumonoidea superfamily produce symbiotic polydnviruses (PDVs), which package circular dsDNA. PDV (Bracovirus (BV) in braconid wasps; Ichnovirus (IV) in ichneumonid wasps) genomes are integrated within the wasp genome as islands of viral genes. Upon oviposition, PDVs suppress host immunity. BVs and IVs derive from nudivirus and large DNA cytoplasmic viruses, respectively (reviewed in (Strand and Burke 2015; Drezen et al. 2017; Gauthier et al. 2018) and references therein).

Immune-suppressive virus-like particles (VLPs) (e.g., VcVLPs in the ichneumonid *Venturia canescens* and FaENVs in the braconid *Fopius arisanus*) lack proviral DNA segments, but are of viral origin and transfer virulence proteins into host cells (Pichon et al. 2015; Burke et al. 2018). Viral genes encoding VLP proteins are either dispersed in the wasp genome (as in VcVLP) or present in discrete genomic areas (as in FaENV). Thus, various independent viral endogenization events have been important for successful parasitism by these wasps (Strand and Burke 2015; Gauthier et al. 2018).

Here, we focus on immune-suppressive particles of figitid wasps of the genus *Leptopilina*, that infect *Drosophila* spp. and are gaining importance as models for natural host-parasite interactions (Keebaugh and Schlenke 2014). *Leptopilina heterotoma* (*Lh*), *L. victoriae* (*Lv*), and *L. boulardi* (*Lb*) produce VLPs in their venom glands. The VLPs of

74 *Leptopilina* spp. and their proteins have been linked to parasite success (Rizki et al.  
75 1990; Dupas et al. 1996; Morales et al. 2005; Labrosse et al. 2005; Chiu et al. 2006;  
76 Heavner et al. 2014). Evidence for DNA in *Leptopilina* VLPs is lacking, and because of  
77 the absence of a published wasp genome, the chromosomal versus extrachromosomal  
78 location of MSEV protein genes is not known. Our goals here are (a) to describe  
79 additional proteins in the MSEV proteome and examine their relationship with PDV and  
80 other viral proteins, and (b) determine whether MSEV genes are encoded within the  
81 wasp genome.

82         We recently described 161 proteins in the VLPs from two *Lh* strains in three  
83 classes: conserved eukaryotic with cellular function (Class 1), infection- and immunity-  
84 related (Class 2), and unannotated (novel) proteins without similarity to known proteins  
85 (Class 3) (Heavner et al. 2017). Class 1 proteins include several vesicular transport and  
86 endomembrane system proteins. Class 2 proteins include predicted modulators of  
87 immune response, e.g., metalloendopeptidases, RhoGAPs, a knottin-like protein, and a  
88 new family of prokaryotic-like GTPases whose genes lack introns. A striking example of  
89 Class 3 proteins is p40, with three-dimensional structural similarity to Type 3 secretion  
90 system (T3SS) needle tip proteins, IpaD/SipD/BipD from Gram-negative bacteria,  
91 *Salmonella*, *Shigella* and *Burkholderia*. Earlier results have indicated that the *p40* gene  
92 (unlike the *GTPase* genes) is expected to have introns. These results suggested that *Lh*  
93 VLPs have novel properties with elements of the prokaryotic and eukaryotic secretion  
94 systems and possess a functionally diverse array of immune-suppressive proteins. We  
95 therefore renamed VLPs as Mixed Strategy Extracellular Vesicles (MSEVs). Their  
96 variable morphologies distinguish them from ordered PDV morphologies. Additionally,

97 genes encoding abundant MSEV proteins p40 and GTPase are present even in  
98 antibiotic-treated wasps. These results favored a non-microbial nature of MSEVs  
99 (Heavner et al. 2017).

100 Here, we present an analysis of an additional 246 proteins from the *Lh* 14 MSEV  
101 proteome to obtain a more comprehensive description. A combined analysis of these  
102 and previous results reinforce the idea that the MSEV proteome is enriched in exosomal  
103 proteins and that Class 3 proteins are not shared with either *Lb* or an unrelated  
104 *Gnaspis* spp. Whole-body transcriptome of adult *Lh* wasps validated the expression of  
105 the MSEV genes. *De novo* genomic assembly and analyses revealed 90% of conserved  
106 Insecta Benchmarking Universal Single-Copy Orthologs (BUSCOs), as well as a  
107 majority (375/407; ~90%) of the MSEV proteins are encoded in the wasp genome.  
108 While we cannot rule out a viral origin of MSEVs, in aggregate, our results provide a  
109 clearer understanding of the extant nature of these complex structures and strengthen  
110 the idea that specialized extracellular vesicles transfer wasp virulence factors and other  
111 parasite proteins into *Drosophila* host cells.

112

## **MATERIALS AND METHODS**

113 **Insects:** Isogenized *Lh* strains New York (NY; (Chiu and Govind 2002; Chiu et al.  
114 2006)) and *Lh* 14 (Schlenke et al. 2007), were raised on the *y w* strain of *D.*  
115 *melanogaster* that were reared on standard cornmeal, yeast, and agar fly food at 25°C  
116 as described (Small et al. 2012). Adult wasps were collected from parasitized hosts, 25  
117 days after infection at 25°C. Male and female wasps were stored on fly food with 70%  
118 honey on “buzz” plugs.

119

120 **Analysis of MSEV super-set ORFs:** Previously undescribed open reading frames  
121 (ORFs) from the *Lh* 14 MSEV proteome and sequenced as part of (Heavner et al. 2017)  
122 (PXD005632) are analyzed in the context of the published female abdominal (Goecks et  
123 al. 2013) and whole body (this study) *Lh* 14 transcriptomes. We have not observed any  
124 difference in venom activities of *Lh* 14 and *Lh* NY (Morales et al. 2005; Schlenke et al.  
125 2007), or in wasp success under laboratory conditions. The *Lh* 14 ORFs were aligned  
126 against transcripts from BioProject: PRJNA202370, Accession number GAJJC0000000  
127 (Goecks et al. 2013) as previously described in (Heavner et al. 2017). Proteins with an  
128 ORF and a transcript were run through the BLAST2GO (v 5.2; downloaded June 2018)  
129 annotation pipeline with an E-value threshold of  $1 \times 10^{-7}$  (Conesa et al. 2005; Götz et al.  
130 2008). Results were organized and classified based on Gene Ontology (GO) terms from  
131 UniProt and InterPro (Ashburner et al. 2000; Jones et al. 2014; Consortium 2015;  
132 Carbon et al. 2019). Proteins were considered “virulence-related” based on GO terms  
133 indicating involvement with infection, host evasion, inflammation, and immune  
134 response. ORFs that did not return results via BLAST or InterProScan (Class 3

135 proteins) were run through Conserved Domain Search (CDD) on NCBI (version 3.16)  
136 (Marchler-Bauer et al. 2017). The E-value cut off for CDD search was  $1 \times 10^{-2}$ . Proteins  
137 were considered to have a signal peptide if one was predicted using Phobius and  
138 Signal P (Kall et al. 2004; Kall et al. 2007; Nielsen 2017; Almagro Armenteros et al.  
139 2019). Transmembrane domains were considered to be present if they were predicted  
140 using Phobius and TMHMM (Sonnhammer et al. 1998; Krogh et al. 2001; Kall et al.  
141 2004; Kall et al. 2007).

142 The GhostKOALA algorithm (Kanehisa et al. 2016) was used to assign KEGG  
143 ortholog (KO) numbers for the MSEV superset protein sequences. If a primary KO  
144 number failed to be assigned by GhostKOALA, a secondary number assignment with a  
145 score  $\geq 50$  was used. Redundant KO numbers were excluded.

146 MSEV proteins were included in the enrichment analyses only if a human  
147 ortholog exists; the gene identifiers for human orthologs were obtained from the MSEV  
148 KO and the UniProt mapping utility (Li et al. 2015; Consortium 2015). (Human orthologs  
149 were chosen because a robust proportion of Vesiclepedia's data is derived from human  
150 vesicle proteomes.) The orthologs of human genes were analyzed for enrichment with  
151 the FunRich algorithm (Pathan et al. 2015; Pathan et al. 2017) against the Vesiclepedia  
152 database (Kalra et al. 2012; Pathan et al. 2019).

153 Finally, the MSEV proteome was used as a query using BLASTp for the following  
154 databases: "non-redundant" (nr), nr restricted to Taxid: Viridae (10239), nr restricted to  
155 Taxid: Polydnaviridae (10482), and nr restricted to Taxid: Unclassified Polydnaviridae  
156 (40273) (E-value threshold:  $1.0 \times 10^{-3}$ , %ID minimum: 20%, performed 04/16/2019).  
157 tBLASTn of *L. boulardi* and *G. hookeri* (previously called *Ganaspis spp. 1*) (Goecks et



158 al. 2013) transcriptomes was performed on 03/10/2019; the threshold for homologs in  
159 *Lb* and *G. hookeri* were 25% ID and an E-value of  $1.0 \times 10^{-10}$ .

160  
161 **Genomes sequencing and assembly:** Library preparations, sequencing reactions,  
162 and associated validations were conducted by GENEWIZ, Inc. (South Plainfield, NJ,  
163 USA). Genomic DNA was extracted from ~50 mg of tissue (~100 wasps) of *Lh* males  
164 and females separately using mixed bead beating and PureLink Genomic DNA  
165 extraction kits following manufacturer's protocol. Quantification of extracted DNA was  
166 performed using Nanodrop and Qubit2.0 Fluorometer (Live Technologies, Carlsbad,  
167 CA, USA). Integrity of genomic DNA was verified by gel electrophoresis (0.6% agarose).  
168 DNA libraries were prepared for each wasp gender by acoustic shearing fragmentation  
169 using a Covaris S220. Fragments were end repaired and adenylated prior to adapter  
170 ligation on 3' ends (NEB NextUltra DNA Library Preparation kit, Illumina, San Diego,  
171 CA, USA). Enrichment and indexing of adapter-ligated DNA was done through limited  
172 cycle PCR. DNA library validation was performed using TapeStation (Agilent  
173 Technologies, Palo Alto, CA, USA). Libraries were quantified using Qubit 2.0  
174 Fluorometer.

175 Real time PCR (Applied Biosystems, Carlsbad, CA, USA) was used to quantify  
176 DNA molar mass for each library before multiplexing in equal molar mass. DNA libraries  
177 were sequenced using a 2x150 paired-end (PE) configuration on one lane on an  
178 Illumina HiSeq 4000. Image analysis and base calling were performed using the HiSeq  
179 Control Software (HCS) on the HiSeq instrument.

180 The average size of inserts (without adaptors) in the Illumina library was ~300-  
181 350 bp. *De novo* assembly of reads and scaffolding of contigs was performed using  
182 ABySS 2.2 (Jackman et al. 2017) by the New York Genome Center. *De novo* assembly  
183 of combined male/female genome was performed using Platanus-allee (Kajitani et al.  
184 2019) and scaffolding was improved using AGOUTI (Zhang et al. 2016) on the  
185 University of Delaware's BIOMIX cluster.

186 Sequences from *Drosophila*-associated bacteria such as *Wolbachia* spp.,  
187 *Acetobacter pasteurianus*, and *Lactobacillus plantarum* were identified in both  
188 assemblies. *Wolbachia* are endosymbionts of many insects including *Leptopilina* spp.  
189 (Pannebakker et al. 2004; Werren et al. 2008; Gueguen et al. 2012). *Lactobacilli* and  
190 *Acetobacter* are symbionts and commensals of sugar-consuming insects (Crotti et al.  
191 2010; Engel and Moran 2013). Among the three bacterial species, *Wolbachia*  
192 sequences were the most abundant. BLASTx analysis showed that predicted genes  
193 from *Wolbachia* scaffolds were associated with *Wolbachia* proteins in GenBank. These  
194 bacterial and mitochondrial sequence-containing scaffolds were identified during the  
195 NCBI submission process and were manually removed from the submission.

196  
197 **Evaluation of genome assemblies:** Assemblies made with ABySS and Platanus-allee  
198 with AGOUTI were run through QUAST v4.0 (Mikheenko et al. 2016) to determine  
199 scaffold number, N50, and GC%. All assemblies were examined for conserved genes  
200 and orthologs with BUSCO v9 (Simão et al. 2015; Waterhouse et al. 2017) using the  
201 Insecta set and training parameters set to "Nasonia". NCBI BLAST+ (v 2.7.1) was used  
202 to compare select scaffolds produced from male and female genome assemblies

203 (Johnson et al. 2008; Camacho et al. 2009). E-value threshold was set at  $1 \times 10^{-7}$ . E-  
204 values of alignments were considered acceptable if within the range of 0 to  $1 \times 10^{-10}$ .

205 K-mer analysis was performed using the K-mer Analysis Toolkit (KAT) (Mapleson  
206 et al. 2016) and heat maps were used to compare multiplicity (coverage plus repeats) of  
207 K-mers to GC content of the reads, coloring bins according to the number of distinct K-  
208 mers in each. This analysis was used to determine whether there were separate  
209 clusters of multiplicity/GC content that might arise from different sources, such as  
210 contamination. BLAST (Altschul et al. 1990; Camacho et al. 2009; Johnson et al. 2008)  
211 was used to search for homologs of a random sample of genomic scaffolds to which  
212 reads from each cluster mapped. The joint assembly of the *Lh 14* genome was  
213 compared to the published *L. clavipes* genome (Bioproject: PRJNA84205 (Kraaijeveld  
214 et al. 2016)) through maps of 27-mer multiplicity versus GC content. Finally, 27-mer  
215 multiplicity/GC content of the scaffolds (9.6 Mb) containing MSEV genes was compared  
216 to a random subset of scaffolds (9.6 Mb) without MSEV genes. Statistical differences  
217 between *Lh 14* and *L. clavipes* genomes and between MSEV-gene containing scaffolds  
218 and non-MSEV-gene containing scaffolds were calculated using a multivariate Cramér  
219 test (Ihaka and Gentleman 1996; Baringhaus and Franz 2004; Franz 2019).

220

221 **Gene predictions, gene annotation, and viral gene searches:** Gene predictions were  
222 performed on parallel and anti-parallel strands using AUGUSTUS (v3.3.1; August 2018)  
223 (Stanke et al. 2004; Stanke and Morgenstern 2005; Keller et al. 2011) with the *Nasonia*  
224 training set. The AUGUSTUS readout was separated into mRNA, coding DNA  
225 sequence (CDS), and translations by gffread (Trapnell et al. 2012).

226 Gene predictions were annotated by performing a BLASTx of all gene predictions  
227 against the entire nr database (Downloaded on January, 2019) and InterProScan on the  
228 University of Delaware BIOMIX Cluster before using BLAST2GO (Conesa et al. 2005;  
229 Götz et al. 2008) to finish annotation based on BLASTx and InterProScan results.

230 NCBI BLAST+ (v 2.7.1) was used on a local machine to search predictions and  
231 scaffolds, cutoff was %ID >70%, E-value < 1E-50, and query coverage > 70%. MSEV  
232 genes and  $1 \times 10^{-2}$  for Polydnavirus and Nudivirus proteins. Family *Polydnaviridae* and  
233 *Nudiviridae* protein sequences for the 11 species available on OrthoDB v9 were  
234 downloaded on February 2019 (Johnson et al. 2008; Camacho et al. 2009).

235

236 **Whole-body transcriptome sequencing and assembly:** Total RNA extraction, library  
237 preparations, sequencing reactions, and bioinformatics analysis were conducted at  
238 GENEWIZ, INC (South Plainfield, NJ, USA). RNA was extracted from frozen tissue with  
239 the Qiagen RNeasy Plus Universal mini kit using manufacturer's instructions (Qiagen,  
240 Hilden, Germany). The extracted RNA was quantified using a Qubit 2.0 Fluorometer  
241 and its integrity was checked with the 4200 TapeStation (Agilent Technologies, Palo  
242 Alto, CA, USA).

243 RNA samples were enriched for mRNA using Oligo d(T) beads. RNA sequencing  
244 libraries were prepared using the NEBNext Ultra RNA Library Prep Kit for Illumina  
245 following manufacturer's instructions (NEB, Ipswich, MA, USA). The sequencing  
246 libraries were validated by using the Agilent TapeStation. Quantification was performed  
247 using the Qubit 2.0 Fluorometer and quantitative PCR (KAPA Biosystems, Wilimington,  
248 MA, USA).

249 Sequencing libraries were clustered on a single lane of a flow cell and  
250 sequenced on the Illumina HiSeq 4000 instrument using a 2x150 PE configuration.  
251 Image analysis and base calling were conducted by the HCS. Raw sequence data (.bcl  
252 files) was converted into fastq files and de-multiplexed using Illumina's bcl2fastq 2.17  
253 software. One mismatch was allowed for index sequence identification.

254 The Trinity v2.5 (Grabherr et al. 2011), *de novo* assembler was used to assemble  
255 the *Lh 14* transcripts. The *de novo* assembled transcriptome was created with a  
256 minimum contig length of 200 bp per sample. Transrate v1.0.3 (Smith-Unna et al. 2016)  
257 was used to generate statistics for the *de novo* assembled transcriptome. EMBOSS  
258 tools getorf were then used to find the ORFs within the *de novo* assembled  
259 transcriptome. The *de novo* transcriptome assembly was then annotated using Diamond  
260 BLASTx (Buchfink et al. 2015).

261 The transcriptome reads were mapped to the genomic scaffolds for downstream  
262 analyses using HISAT2 or BWA (Li and Durbin 2009; Kim et al. 2015).

263

264 **Preparation of template DNA and PCR:** Male and female wasps (n=12, for each sex),  
265 were separated and washed in 70% ethanol, and then rinsed twice in deionized water.  
266 Genomic DNA (gDNA) was extracted using a Qiagen DNeasy Blood and Tissue kit  
267 following provided protocols. gDNA was eluted in Tris-EDTA buffer, pH 8.0, and stored  
268 at 4°C. The concentration of gDNA was determined by NanoDrop (Thermo Fisher).

269 For cDNA preparation, male and female wasps (n=12 for each sex), were  
270 separated and washed in 70% ethanol and rinsed twice in deionized water. Total body  
271 RNA was extracted using 100 µL of Trizol (Invitrogen) following manufacturer's

272 protocols. RNA was resuspended in 0.1% DEPC treated water and treated with DNase I  
273 to remove contaminating DNA (Thermo Fisher Scientific). The RNA concentration was  
274 determined by NanoDrop (Thermo Fisher). cDNA was synthesized using Proto-Script  
275 First Strand cDNA Synthesis Kit (New England Biolabs).

276

277 **Analysis of select genes:** Primers for *p40* and *SmGTPase01* are as follows:

278 *p40* forward: GAATCATTGTTTCGTTTGCTTGAAGAAAGAATTGG

279 *p40* reverse: CATTATTAATGGGCCTTTACAATAATTTTAGCC

280 *SmGTPase01* forward: CGTTGCACTACCTTGTTTGTCA

281 *SmGTPase01* reverse: TTGTCTTTGCCCTGAGCGTT

282 PCR products were performed with Taq polymerase (gift of C. Li lab, CCNY), PCR buffer  
283 (300 mM Tris HCl pH 9.5, 75 mM (NH<sub>4</sub>)<sub>2</sub>SO<sub>4</sub>, 10 mM MgCl<sub>2</sub>) and deoxyribonucleotides  
284 (0.2 mM; Thermo Fisher Scientific). The PCR products were resolved on a 1% agarose  
285 gel in Tris acetic acid EDTA buffer (40 mM Tris HCl pH 7.6, 20 mM acetic acid, 1mM  
286 EDTA pH 8.0). Ethidium bromide (Sigma Aldrich)-stained gels were visualized on an  
287 ultra violet Trans-Illuminator (UVP) and gel images were taken using the DigiDoct  
288 Imaging System (UVP). Gel images were processed in Adobe Photoshop for clarity  
289 only.

290 gDNA or cDNA-containing PCR products were cloned into pCR TOPO II  
291 plasmids (Invitrogen) and transformed into DH10β competent cells (New England  
292 Biolabs). For plasmid preparation, colonies were screened via PCR and positive

293 colonies were cultured in Luria Broth with ampicillin (100 µg/mL) at 37°C overnight.  
294 Plasmids were extracted using Plasmid Miniprep kit (Qiagen) and sequenced  
295 (GENEWIZ, INC. South Plainfield, NJ, USA). Sequences were aligned using NCBI  
296 BLAST+ (Johnson et al. 2008; Camacho et al. 2009) and Clustal Omega (Li et al. 2015).  
297 Expected PCR band sizes were determined using SerialCloner (v2.6.1).

298

299 **Data availability:** *L. heterotoma* strains (Chiu et al. 2006; Schlenke et al. 2007) are  
300 available upon request. All supplemental files and figures can be downloaded from  
301 figshare (<https://figshare.com/s/3a4598308909307c2ae0>). File S1 contains details of  
302 supplemental files and tables. File S2 contains listing of accession numbers for all  
303 sequences reported in this work. Figure S1 contains the 27-mer vs GC count  
304 comparison of MSEV containing scaffolds to non-MSEV containing scaffolds. Table S1  
305 contains annotations and related data for proteins. Table S2 contains BLAST search  
306 results of the MSEV proteome against the nr database. Table S3 contains all BUSCOs  
307 found in male, female, and joint genome assemblies. MSEV protein sequences are  
308 available upon request. Accession numbers for datasets are as follows: *Leptopilina*  
309 *heterotoma* strain *Lh* 14, genome assembly: Male genome: QYUB0000000, Female  
310 genome: QYUC0000000, Joint genome: VOOK000000000. *Leptopilina heterotoma* strain  
311 *Lh* 14, whole-body transcriptome: GHUQ000000000. *Leptopilina heterotoma* abdominal  
312 transcriptome by Goecks et al.: GAJC000000000. *Leptopilina clavipes* genome Bioproject:  
313 PRJNA84205. *Leptopilina heterotoma* strain *Lh* 14 proteome: PRIDE: PXD005632

314

## RESULTS AND DISCUSSION

315 **The MSEV proteome superset:** A comparative study of the proteomes of the MSEVs  
316 from *Lh* 14 and *Lh* NY strains previously generated a list of 161 “common” MSEV  
317 proteins (Heavner et al. 2017). More than 90% of the 161 proteins are part of the *Lh* 14  
318 MSEV proteome. To describe MSEVs more completely, we characterized a larger set of  
319 407 MSEV proteins from *Lh* 14 (161 common and 246 *Lh* 14) and define this set as the  
320 *Lh* MSEV “super-set” (Fig. 1A). Key results from annotation-based classification,  
321 analysis for signal peptide and/or transmembrane domain, and presence/absence of  
322 proteins in related wasps are summarized below and in Table S1.

323         The presence/absence of signal peptide (SP) alone, or SP with/without the  
324 transmembrane (TM) domain(s) in MSEV proteins reveals their possible location (i.e.,  
325 potentially secreted into the venom gland lumen or associated with MSEV membrane).  
326 We therefore searched the 246 *Lh* 14 proteins for SP and TM domains. Of the 246  
327 proteins, 55 (22.35%) have a predicted SP domain, 37 (15.04%) have a predicted TM  
328 domain, while 6 (2.44%) have both a predicted SP and TM domain.

329         After annotation, we found that a majority (183/246 or 74.39%) of the 246  
330 proteins can be classified as core eukaryotic cell biology proteins (Class 1); 13/246  
331 (5.28%) proteins as virulence- and immunity-related based on associated GO terms  
332 (Ashburner et al. 2000; Carbon et al. 2019) (Class 2); and 50/246 or 20.32 % as novel  
333 sequences without high confidence annotation (Class 3) (Table S1). A  
334 presence/absence analysis of these 246 proteins in published transcriptomes (Goecks  
335 et al. 2013) of *Lb* or a more distantly related wasp, *G. hookeri* (for thresholds see  
336 Methods) revealed the following: only 43/246 (17.47%) *Lh* MSEV proteins are expected



337 to be found in *Lb* and/or *G. hookeri* (Table S1). Of these, 33/43 (76.74%) proteins were  
338 in Class 1 but only 7/43 (16.27%) and 3/43 (6.98%) were in Class 2 and Class 3  
339 categories, respectively. These results support the idea that, multiple but different,  
340 infection strategies and/or host evasion strategies might exist among different wasps  
341 infecting the same hosts.

342 While most of the Class 1 proteins were annotated as ribosomal or  
343 mitochondrial-related, a few were described as integral membrane proteins, vesicle  
344 trafficking protein SEC22b (E-value: 6.22E-145), and the ion channels sideroflexin 1  
345 and 2 (E-value: 0). We also identified an apolipoprotein (E-value: 1.02E-7) (Table S1).  
346 The presence of these membrane-associated proteins reinforces the vesicular nature of  
347 MSEVs.

348 Examples of Class 2 proteins include the neural/ectodermal development factor  
349 IMP-L2 (E-value: 5.29x10<sup>-50</sup>) and a protein involved in pain reception, CG9231 (E-value:  
350 4.39x10<sup>-15</sup>). A viral-like Diel protein (E-value: 1.77x10<sup>-7</sup>), viral Enhancin (E-value:  
351 6.02x10<sup>-5</sup>), I(2)37Cc (E-value: 3.39x10<sup>-165</sup>), odorant binding protein 56d-like (E-value:  
352 5.64x10<sup>-50</sup>), major royal jelly protein (E-value: 8.59x10<sup>-135</sup>), and two venom acid  
353 phosphatases Acph-1 (E-value: 4.12x10<sup>-5</sup>) were also found in the Class 2 category;  
354 their cDNA sequences were published previously (Heavner et al. 2013) (Table S1). It is  
355 possible that these MSEV proteins modulate the hosts' immune responses and/or  
356 influence host development to facilitate successful parasitism.

357 Within Class 3, 45 proteins (90%) lacked BLASTp and InterProScan results.  
358 However, Conserved Domain Database (CDD) (Marchler-Bauer et al. 2017) searches  
359 returned 9 hits identifying potentially functional domains (Table 1). This included (a) a

360 CD99L2 like antigen (%ID: 24%, E-value:  $3 \times 10^{-3}$ ), (b) a DEAD-like helicases superfamily  
361 member (%ID: 22%, E-value  $2 \times 10^{-4}$ ) and (c) a herpes outer envelope glycoprotein 350  
362 (gp350), (%ID: 28%, E-value:  $4 \times 10^{-3}$ ) (Table 1).

363 A BLASTp DELTA-BLAST of the potential gp350 domain against the nr database  
364 specifying “Vira” (taxid: 10239) under organism resulted in Crimean-Congo hemorrhagic  
365 fever orthonairovirus envelope glycoprotein (%ID: 30%, E-value:  $2.5 \times 10^{-1}$ ),  
366 Lymphocryptovirus Macaca gp350 (%ID: 29%, E-value:  $7.2 \times 10^{-1}$ ), and Gallid  
367 Alphaherpesvirus 1 envelope glycoprotein J (%ID: 26%, E-value: 1.2) as top hits.

368 BLASTp DELTA-BLAST of the potential gp350 domain against the nr database for  
369 Hymenoptera yielded an uncharacterized protein as the best hit (%ID: 24%, E-value:  
370  $8 \times 10^{-6}$ ) in the ant *Vollenhovia emeryi*. This ant protein is predicted to contain calcium-  
371 binding EGF domains. The second hit in this search is from *N. vitripennis* for a predicted  
372 mucin-3A like glycoprotein (Gendler and Spicer 1995) (%ID: 24%, E-value:  $2 \times 10^{-4}$ ).

373 Interestingly, transcripts related to the potential *Lh* gp350-like protein are not found in  
374 the *Lb* or *G. hookeri* transcriptomes (Table S1) (Goecks et al. 2013). Presence of this  
375 gp350-like protein in *Lh* MSEVs, but its absence in *Lb* MSEVs, suggests that it  
376 somehow contributes to differences in *Lh/Lb* host-parasite interactions and is therefore  
377 worthy of future studies. Complement receptor type 2 (CR2) in human B lymphocytes  
378 interacts with gp350 during Epstein-Barr infection (Young et al. 2008) and finding a  
379 verified homologue of CR2 in *Drosophila* hosts would be interesting in future research.

380 Because more than 200 proteins have now been added to the previously  
381 described MSEV proteome (Heavner et al. 2017), we re-evaluated our previous  
382 enrichment analysis. In an ortholog-based comparison of the superset to human

383 extracellular vesicle (EV) proteomes in Vesiclepedia (the most current and robust  
384 source of EV data (Kalra et al. 2012)), we found that the largest proportion of superset  
385 proteins (49%) are proteins specifically associated with exosomes (Fig. 1B). In human  
386 and mouse EV proteomes, mitochondrial and ribosomal proteins are enriched (Kalra et  
387 al. 2012). Accordingly, protein components of mitochondria (e.g., respiratory chain) and  
388 ribosomes (e.g., large and small subunit proteins) are found to be highly enriched in the  
389 *Lh* MSEV superset. However, we found that the significance of the enrichment was  
390 higher between the superset and exosomal proteins than mitochondrial or ribosomal  
391 proteins (Fig. 1C). These results demonstrate the similarities in the protein profiles of  
392 MSEVs and EVs.

393

394 **Do *Lh* MSEVs contain homologs of PDV or viral proteins?** Even though figitid  
395 *Leptopilina* wasps are distantly related to PDV-containing Ichneumonid and Braconid  
396 wasps (Misof et al. 2014; Strand and Burke 2015), an association of PDV-like viruses in  
397 figitid wasps cannot be discounted because of shared evolutionary history. Recent  
398 publications have identified capsid-less VLPs in Ichneumonidae wasps (Volkoff et al.  
399 2010; Pichon et al. 2015; Burke et al. 2018) and it is possible that *Lh* MSEVs have a  
400 similar viral origin. We therefore analyzed the *Lh* MSEV proteome superset against the  
401 GenBank PDV database, and then against its entire Viridae database.

402 To identify false positives, MSEV proteins with positive PDV hits (E-values were  
403 less than  $1.0 \times 10^{-3}$ , %ID was 20% or greater, and query coverage was 30% or higher)  
404 were also searched against the unrestricted nr database to compare relatedness. If an  
405 MSEV protein is similar to a viral or virus-related PDV protein, we expected that, in the

406 unrestricted nr database search, the MSEV query sequence would align again with the  
407 same viral subject sequences, but with a lower E-value (Table S2).

408 For PDV searches (Taxid: 10482 and Taxid: 40273), four proteins returned hits  
409 with E-values better than  $1.0 \times 10^{-20}$  and query coverage greater than 30%. Three of  
410 these hits are conserved proteins (cytochrome P450 and histone 4) while the fourth  
411 result identified an uncharacterized *Cotesia congregata* bracovirus (CcBV) protein (%ID:  
412 31.08%, E-value:  $1.38 \times 10^{-17}$ , query coverage: 77%) (Table S2). The unbiased BLASTp  
413 search against the entire nr database however had better results against eukaryotic  
414 proteins (E values: 0 to  $2.0 \times 10^{-7}$  and %ID from 100 to 56.25) (Table S2). In fact, the  
415 query that yielded the CcBV protein was better matched to a eukaryotic ribonuclease  
416 (%ID: 26.06%, E-value  $1.14 \times 10^{-16}$ , query coverage: 84%) (Table S2). These results  
417 suggest that MSEV sequence similarities with PDV proteins may not be significant.

418 We also searched the *Lh* MSEV superset for presence of *L. bouleardi* Filamentous  
419 Virus (LbFV) homologs (LbFV is a behavior manipulating virus of *Lb* (Varaldi et al. 2006;  
420 Patot et al. 2012)). Of LbFV's 108 genes, 13 are present in genomes of *Lb*, *Lh* and  
421 related species, and the 13 transcripts are expressed in the *Lb* venom gland (Di  
422 Giovanni et al. 2019). Within our thresholds, we obtained only three (of 13) sequences  
423 with similarity to LbFV ORFs. However, these three *Lh* MSEV proteins, with hits for  
424 LbFV sequences obtained better scoring hits in the unrestricted nr database, suggesting  
425 that the *Lh* MSEV proteins are not highly related to the LbFV proteins (Table S2).

426 When comparing MSEV proteins to the entirety of Viridae, a total of 35 MSEV  
427 proteins had hits with %IDs ranging from 30% to 71% and E-values ranging from  
428  $1.0 \times 10^{-22}$  to  $1.0 \times 10^{-178}$  (Table S2). However, a BLASTp search against the entire nr

429 database found that proteins with results for viral hits had better scores when searched  
430 against the entire nr database, indicating that while viral hits are possible, they are not  
431 the best match (Table S2). This result, in addition to the fact that 372 other MSEV  
432 proteins (including the Dieldel and Enhancin (Table S1)) did not return viral hits, would  
433 indicate that a majority of MSEV proteins are not closely related to viral proteins.

434

435 **The whole-body transcriptome contains expressed MSEV transcripts:** We  
436 performed a mixed-gender whole-body transcriptome sequencing and *de novo*  
437 assembly of *Lh* transcripts. This assembly generated 104,066 transcripts. This dataset  
438 is more than three times larger than the published data derived from abdomens of  
439 female wasps that has 31,400 transcripts (Goecks et al. 2013). A BLAST analysis of the  
440 female abdominal transcripts against the male/female whole-body transcripts showed  
441 that a majority (21,493/31,400, 68.44%) were present in the latter data set.

442 We searched the whole-body transcriptome for MSEV protein coding sequences  
443 using tBLASTn. Of 407 MSEV superset proteins, we identified transcripts for 371  
444 (91.15%) proteins. Despite the ~9% discrepancy (likely due to differences in expression  
445 levels due to different experimental conditions), these results largely verify the transcript  
446 data from (Goecks et al. 2013) that we have based our proteomic analyses on. Of the  
447 371 MSEV transcripts identified, 233 (62.8%) encode Class 1, 44 (11.86%) encode  
448 Class 2, and 94/371 (25.34%) encode Class 3 proteins.

449

450 **Assembly of the *Lh* genome:** We separately sequenced *Lh* 14 male and female  
451 genomic DNA and assembled the paired-end reads *de novo*, using ABySS (Jackman et

452 al. 2017). These assemblies have a modest scaffold N50 of 4,800 with more than  
453 100,000 scaffolds and an average coverage of 87% (Table 2). Assembly with MaSurCa  
454 (Zimin et al. 2013) provided similar results (data not shown), indicating that our  
455 assembly quality is limited likely due to factors such as large genome size and repetitive  
456 sequence regions (Dominguez Del Angel et al. 2018). Although the N50 values and  
457 large number of scaffolds indicate that the genome is not highly assembled, we found at  
458 least 80% of BUSCOs shared in the Insecta set in both assemblies (Table 2, BUSCOs  
459 in Table S3).

460         While still fragmented, a *de novo* joint assembly of male and female sequences  
461 using Platanus-allee and AGOUTI improved assembly and scaffolding statistics (N50:  
462 11,906, average coverage 91.1%). The number of found BUSCOs in the joint assembly  
463 rose to 90% (Table 2, BUSCOs in Table S3).

464         Analysis of K-mer multiplicity versus GC content in the genome sequencing  
465 reads using the K-mer analysis tool, KAT (Mapleson et al. 2016) showed three possible  
466 clusters, although they are difficult to distinguish (Fig. 2A). Cluster 1 has high multiplicity  
467 (450-650), Cluster 2 has lower multiplicity and a wide range of GC content, and Cluster  
468 3 has the lowest multiplicity and the highest GC content. Cluster 3 overlaps with Cluster  
469 2 making them hard to fully separate. BLAST searches of a random sample  
470 (1,672/4,482) from Cluster 1 contigs hit insect homologs 73% (1,220/1,672) of the time,  
471 *Acetobacter* homologs 13% (216/1,672) of the time, and then a variety of mostly  
472 Eukaryotic hits. Cluster 2 represents a majority of the wasp genome (>94%), and blast  
473 hits of a random sample (316/68,173) of its contigs almost exclusively had homologs in  
474 Hymenoptera (311/316; 98%) and mostly in *L. heterotoma* (227/316; 71%). Cluster 3 is

475 the smallest of the three and contigs from Cluster 3 had homologs exclusively in  
476 *Acetobacter* (110/110). There was no evidence for contamination from a viral source or  
477 discrete MSEV-specific set of nucleic acid sequence.

478 Furthermore, K-mer multiplicity versus GC content for the joint assembly of the  
479 *Lh 14* genome (Fig. 2B) showed a very similar heat map to that using the published  
480 assembly of *L. clavipes* (Fig. 2C; Bioproject: PRJNA84205 (Kraaijeveld et al. 2016)).  
481 The two genomes have a highly similar 27-mer/GC profiles that do not differ statistically  
482 (multivariate Cramér test statistic = 114,119,  $P = 0.73$ , number bootstrap-replicates =  
483 1000). The *L. heterotoma* assembly has 27-mers with approximately twice the  
484 multiplicity of those found in *L. clavipes*, which may represent increased repeat content  
485 in *L. heterotoma* and is supported by an assembly size of over 200 Mb larger than the *L.*  
486 *clavipes* genome (463 Mb vs 255 Mb) (Kraaijeveld et al. 2016).

487

488 **MSEV genes are encoded in the wasp genome:** Using our annotation pipeline,  
489 28,481 predicted genes were annotated. Within the annotated genes, we found 8 genes  
490 for the body color *yellow*, 3 *major royal jelly protein (mrjp)* genes, 25 odorant  
491 receptor/odorant binding protein coding genes, and 94 gene predictions for *cytochrome*  
492 *P450*. Some of these nuclear genes are not only involved in development and cellular  
493 processes, but are also included in the MSEV proteome (Table S1 and (Heavner et al.  
494 2017)). A search of gene predictions for MSEV proteins via tBLASTn identified 325 of  
495 407 (79.85%) MSEV sequences (Table 3). Of these, 153/407 (37.6%) had a percent  
496 identity of 95% or greater. Presence or removal of scaffolds with bacterial DNA

497 sequences from either the separate male/female or the joint assembly did not affect this  
498 number, supporting the nuclear location of a majority of the MSEV genes.

499 As gene prediction software can potentially miss genes (Wang et al. 2004), we  
500 searched the genomic scaffolds directly for MSEV-coding sequence regions using  
501 known protein sequences as queries via tBLASTn before and after removal of bacterial  
502 sequences. In both cases, 375/407 (92.13%) MSEV sequences were at least 70%  
503 complete as determined by query coverage (Table 3). Of these, 191/407 (46.9%) had a  
504 percent identity of 95% or greater.

505 The scaffolds containing MSEV genes (Fig. S1A) were also compared to a  
506 random subset of scaffolds without MSEV genes (Fig S1B) for their 27-mer/GC profiles.  
507 These appeared to not differ statistically (multivariate Cramér test statistic = 3755,  $P =$   
508 0.80, number bootstrap-replicates = 1000), indicating that the MSEV genes lie on  
509 scaffolds that resemble the rest of the genome.

510

511 **Characterization of select MSEV genes:** We spot checked small portions of the  
512 genome for gene structure predictions of MSEV virulence protein genes *SmGTPse01*  
513 (Class 2) and *p40* (Class 3). For this, we sequenced PCR products of gDNA  
514 corresponding to these genes.

515 The MSEV SmGTPase01 has prokaryotic-like GTPase domains and its gene is  
516 expected to lack introns (Heavner et al. 2017). The predicted *SmGTPase01* CDS spans  
517 936 bp, which contains the functional GTPase domain (Heavner 2018). Scaffolds from  
518 male and female genomes confirmed the absence of coding region introns (data not  
519 shown). We hypothesized that primers in 5' and 3' untranslated regions (UTRs) should



520 amplify the exact fragment from cDNA/gDNA as template based on manual  
521 characterization of the *SmGTPase01* locus (Heavner 2018) (Fig. 3A). This prediction  
522 was borne out and we amplified an 873 bp fragment only from female cDNA and from  
523 both male and female gDNA (Fig. 3C). The sequenced PCR products were identical to  
524 corresponding sequence within the assembly and the published transcript sequence  
525 from Goecks et al. (data not shown).

526 The *p40* gene encodes a protein that is structurally similar to T3SS bacterial  
527 needle tip proteins IpaD/SipD from *Shigella* and *Salmonella*. However, *p40*'s genomic  
528 sequence is expected to have introns (Heavner et al. 2017). The full *p40* gene was  
529 computationally assembled and predicted within both male and female genomes.  
530 Primers designed for *p40*'s 5' and 3' UTRs (Goecks et al. 2013; Ramroop 2016) (Fig.  
531 3B), allowed amplification of *p40*'s 939 bp cDNA only in preparations from female wasp  
532 extracts, but gDNA bands at 1,630 bp were detected from reactions when either male or  
533 female genome was used as template, indicating the presence of introns (Fig. 3D).  
534 Sequencing the cloned cDNA product from females confirmed the published cDNA  
535 sequence (Heavner et al., 2017). We also cloned and sequenced the gDNA products  
536 from male and female wasps and found the sequences to be identical (data not shown).

537 Unlike the well-characterized *Drosophila* hosts, the biology and molecular-  
538 genetics of their parasitic wasps have remained relatively obscure with only recent  
539 characterizations of *Leptopilina* and *Ganaspis* spp. (Melk and Govind 1999; Colinet et  
540 al. 2013; Goecks et al. 2013; Heavner et al. 2013; Mortimer et al. 2013; Heavner et al.  
541 2017; Di Giovanni et al. 2019; Khan et al. 2018). Our proteomic, transcriptomic and  
542 genomic results here expand the available information on *L. heterotoma*. Bioinformatics

543 analysis of the additional MSEV proteins does not alter the initial interpretation of the  
544 original 161 proteomic data. Genomic sequencing and analysis of scaffolds reveals that  
545 more than 92% of the MSEV genes reside on the wasp genome. We did not find  
546 evidence for MSEV gene association with endosymbiont or commensal bacterial DNA.  
547 We suspect that the remaining ~8% are also nuclear genes and this association will be  
548 confirmed in higher quality assemblies. Altogether, these results strongly suggest that,  
549 like other subcellular structures, MSEVs are encoded in the wasp nuclear genome.

550         The cellular nature of *Lh* vesicles is likely to be shared by closely related *Lv* and  
551 *Lb* wasps. Our previous work has shown that the overall morphologies of *Lh* and *Lv*  
552 MSEVs are similar (Morales et al. 2005; Chiu et al. 2006). However, this is not the case  
553 for *Lb* MSEVs; different *Lb* strains have varying MSEV morphologies (Dupas et al.  
554 1996; Gueguen et al. 2011; Wan et al. 2019). Interpretation of their identity also varies.  
555 For example, Di Giovanni et al. (Di Giovanni et al. 2019) contend that MSEVs/VLPs are  
556 derived from a virus ancestral to the LbFV. Our analysis of the expanded proteomic  
557 superset does not lend strong support to this line of thinking.

558         We did not find convincing evidence of PDV or other viral structural proteins in  
559 the *Lh* MSEV proteome. However, we cannot discount that MSEVs have a viral origin  
560 as our analysis is limited by the fragmentation of the genome. It is also possible that a  
561 virus related to MSEVs may not have been identified to date. Mechanistically,  
562 eukaryotic viruses and vesicles share cellular pathways involving the endomembrane  
563 systems of their cells of origin or their target cells (Nolte-'t Hoen et al. 2016), leading to  
564 overlap in protein functionality, but not necessarily origin. Thus, at least some of the  
565 Class 1 proteins in the MSEV proteome may be central to MSEV biogenesis in the wasp

566 or for their interactions with the host hemocytes' endomembrane machinery despite  
567 potentially being related to viruses. It is noteworthy that energy metabolism genes  
568 appear to be involved in rapid speciation and adaptation to new environments (Gershoni  
569 et al. 2009; Lane 2009), raising the possibility that MSEV mitochondrial proteins might  
570 contribute to this process. How *Lh* MSEVs are functionally similar to other insect or  
571 mammalian EVs remains to be explored experimentally. Functional characterization of  
572 predicted infection and immunity Class 2 proteins should explain the immune-  
573 suppressive strategies of these wasps. RNA interference, infection assays, and other  
574 experimental strategies should make this line of inquiry feasible.

575         Functional assignments are difficult for the unannotated Class 3 proteins. These  
576 are likely to be quite interesting, due to their different expression profiles in *Lh* versus *Lb*  
577 and *G. hookeri* species. This difference in expression may stem either from *cis* changes  
578 in their regulatory sequences, or from absence of these genes in the *Lb* or *G. hookeri*  
579 genomes. Recent comparative genomics analysis has shown that over 40% of venom  
580 genes in the closely-related species *N. vitripennis* and *N. giraulti* have diverged  
581 significantly and up to 25% of venom genes are specific to a species (Martinson et al.  
582 2017). A proteomic analysis of the venom genes of *Leptopilina spp.* and a molecular  
583 understanding of their expression will provide insights into how key activities within  
584 MSEVs evolved to parasitize the range of fruit fly hosts.

585         A key question regarding *Lh* virulence proteins critical to wasp success is  
586 whether their genes reside in a discrete region of the genome like a "virulence island"  
587 found in some microbial genomes (Dobrindt et al. 2004; Gal-Mor and Finlay 2006), or  
588 whether some genes are dispersed within the genome, while others occur in one or

589 more clusters as in wasps with PDVs (Volkoff et al. 2010; Pichon et al. 2015). More  
590 complete assemblies, scaffolded to the level of chromosomes, will describe the  
591 genome-wide distribution of these genes in *Lh* and related wasps. Key MSEV genes  
592 could serve as genetic markers in future studies. Comparative genomics will uncover  
593 additional gene family members of MSEV proteins in other *Leptopilina* wasps and  
594 enable the development of new functional genomics tools such as CRISPR-disrupted  
595 mutant alleles made in *N. vitripennis* (Werren et al. 2009; Siebert et al. 2015; Li et al.  
596 2017b; Li et al. 2017a). These approaches will open new avenues for understanding the  
597 biology of this host-parasite model.

598

599

## **ACKNOWLEDGEMENTS**

600 We thank Drs. W. Qiu and S. Singh, and J. Chou for discussions and critical comments,  
601 and our reviewers for insightful feedback. We are grateful to A. Corvelo at the New York  
602 Genome Center for help with genome assemblies. Bioinformatics work was conducted  
603 in-house and with the BIOMIX Shared Computing Cluster at Delaware Biotechnology  
604 Institute, University of Delaware. This work was supported by grants from NASA  
605 (NNX15AB42G), NSF (IOS-1121817), NIH (1F31GM111052-01A1, 5G12MD007603-30,  
606 and GM103446).

## LITERATURE CITED

- 607  
608  
609 Almagro Armenteros, J.J., K.D. Tsirigos, C.K. Sonderby, T.N. Petersen, O. Winther *et*  
610 *al.*, 2019 SignalP 5.0 improves signal peptide predictions using deep neural  
611 networks. *Nat Biotechnol* 37 (4):420-423.
- 612 Altschul, S.F., W. Gish, W. Miller, E.W. Myers, and D.J. Lipman, 1990 Basic local  
613 alignment search tool. *J Mol Biol* 215 (3):403-410.
- 614 Ashburner, M., C.A. Ball, J.A. Blake, D. Botstein, H. Butler *et al.*, 2000 Gene ontology:  
615 tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet* 25  
616 (1):25-29.
- 617 Baringhaus, L., and C. Franz, 2004 On a new multivariate two-sample test. 88 (1):190-  
618 206.
- 619 Buchfink, B., C. Xie, and D.H. Huson, 2015 Fast and sensitive protein alignment using  
620 DIAMOND. *Nature Methods* 12 (1):59-60.
- 621 Burke, G.R., T.J. Simmonds, B.J. Sharanowski, and S.M. Geib, 2018 Rapid Viral  
622 Symbiogenesis via Changes in Parasitoid Wasp Genome Architecture. *Mol Biol*  
623 *Evol* 35 (10):2463-2474.
- 624 Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos *et al.*, 2009 BLAST+:  
625 architecture and applications. *BMC Bioinformatics* 10:421.
- 626 Carbon, S., E. Douglass, N. Dunn, B. Good, N.L. Harris *et al.*, 2019 The Gene Ontology  
627 Resource: 20 years and still GOing strong. *Nucleic Acids Research* 47  
628 (D1):D330-D338.

629 Chiu, H., and S. Govind, 2002 Natural infection of *D. melanogaster* by virulent parasitic  
630 wasps induces apoptotic depletion of hematopoietic precursors. *Cell Death Differ*  
631 9 (12):1379-1381.

632 Chiu, H., J. Morales, and S. Govind, 2006 Identification and immuno-electron  
633 microscopy localization of p40, a protein component of immunosuppressive  
634 virus-like particles from *Leptopilina heterotoma*, a virulent parasitoid wasp of  
635 *Drosophila*. *J Gen Virol* 87 (Pt 2):461-470.

636 Cock, P.J.A., T. Antao, J.T. Chang, B.A. Chapman, C.J. Cox *et al.*, 2009 Biopython:  
637 freely available Python tools for computational molecular biology and  
638 bioinformatics. *Bioinformatics* 25 (11):1422-1423.

639 Colinet, D., E. Deleury, C. Anselme, D. Cazes, J. Poulain *et al.*, 2013 Extensive inter-  
640 and intraspecific venom variation in closely related parasites targeting the same  
641 host: the case of *Leptopilina* parasitoids of *Drosophila*. *Insect Biochem Mol Biol*  
642 43 (7):601-611.

643 Conesa, A., S. Götz, J.M. García-Gómez, J. Terol, M. Talón *et al.*, 2005 Blast2GO: a  
644 universal tool for annotation, visualization and analysis in functional genomics  
645 research. *Bioinformatics* 21 (18):3674-3676.

646 Consortium, U., 2015 UniProt: a hub for protein information. *Nucleic Acids Res* 43  
647 (Database issue):D204-212.

648 Crotti, E., A. Rizzi, B. Chouaia, I. Ricci, G. Favia *et al.*, 2010 Acetic Acid Bacteria, Newly  
649 Emerging Symbionts of Insects. *Applied and Environmental Microbiology* 76  
650 (21):6963-6970.

651 Di Giovanni, D., D. Lepetit, M. Boulesteix, Y. Couté, M. Ravallec *et al.*, 2019 A behavior-  
652 manipulating virus relative as a source of adaptive genes for parasitoid wasps.  
653 *bioRxiv:342758*.

654 Dobrindt, U., B. Hochhut, U. Hentschel, and J. Hacker, 2004 Genomic islands in  
655 pathogenic and environmental microorganisms. *Nat Rev Microbiol* 2 (5):414-424.

656 Dominguez Del Angel, V., E. Hjerde, L. Sterck, S. Capella-Gutierrez, C. Notredame *et*  
657 *al.*, 2018 Ten steps to get started in Genome Assembly and Annotation.  
658 *F1000Res* 7.

659 Drezen, J.M., M. Leobold, A. Bezier, E. Huguet, A.N. Volkoff *et al.*, 2017 Endogenous  
660 viruses of parasitic wasps: variations on a common theme. *Curr Opin Virol* 25:41-  
661 48.

662 Dupas, S., M. Brehelin, F. Frey, and Y. Carton, 1996 Immune suppressive virus-like  
663 particles in a *Drosophila* parasitoid: significance of their intraspecific  
664 morphological variations. *Parasitology* 113 ( Pt 3):207-212.

665 Engel, P., and N.A. Moran, 2013 The gut microbiota of insects - diversity in structure  
666 and function. *Fems Microbiology Reviews* 37 (5):699-735.

667 Franz, C., 2019 cramer: Multivariate Nonparametric Cramer-Test for the Two-Sample-  
668 Problem. R package version 0.9-3. <https://CRAN.R-project.org/package=cramer>.

669 Gal-Mor, O., and B.B. Finlay, 2006 Pathogenicity islands: a molecular toolbox for  
670 bacterial virulence. *Cell Microbiol* 8 (11):1707-1719.

671 Gauthier, J., J.M. Drezen, and E.A. Herniou, 2018 The recurrent domestication of  
672 viruses: major evolutionary transitions in parasitic wasps. *Parasitology* 145  
673 (6):713-723.



674 Gendler, S.J., and A.P. Spicer, 1995 Epithelial mucin genes. *Annu Rev Physiol* 57:607-  
675 634.

676 Gershoni, M., A.R. Templeton, and D. Mishmar, 2009 Mitochondrial bioenergetics as a  
677 major motive force of speciation. *Bioessays* 31 (6):642-650.

678 Goecks, J., N.T. Mortimer, J.A. Mobley, G.J. Bowersock, J. Taylor *et al.*, 2013  
679 Integrative approach reveals composition of endoparasitoid wasp venoms. *PLoS*  
680 *One* 8 (5):e64125.

681 Grabherr, M.G., B.J. Haas, M. Yassour, J.Z. Levin, D.A. Thompson *et al.*, 2011 Full-  
682 length transcriptome assembly from RNA-Seq data without a reference genome.  
683 *Nature Biotechnology* 29 (7):644-U130.

684 Gueguen, G., B. Onemola, and S. Govind, 2012 Association of a new Wolbachia strain  
685 with, and its effects on, *Leptopilina victoricae*, a virulent wasp parasitic to  
686 *Drosophila* spp. *Appl Environ Microbiol* 78 (16):5962-5966.

687 Gueguen, G., R. Rajwani, I. Paddibhatla, J. Morales, and S. Govind, 2011 VLPs of  
688 *Leptopilina boulardi* share biogenesis and overall stellate morphology with VLPs  
689 of the heterotoma clade. *Virus Res* 160 (1-2):159-165.

690 Götz, S., J.M. García-Gómez, J. Terol, T.D. Williams, S.H. Nagaraj *et al.*, 2008 High-  
691 throughput functional annotation and data mining with the Blast2GO suite.  
692 *Nucleic Acids Res* 36 (10):3420-3435.

693 Heavner, M.E., 2018 Evidence for Organelle-like Extracellular Vesicles From a Parasite  
694 of *Drosophila* and Their Function in Suppressing Host Immunity, pp. 236 in  
695 *Biochemistry*. City University of New York, Graduate Center.

696 Heavner, M.E., G. Gueguen, R. Rajwani, P.E. Pagan, C. Small *et al.*, 2013 Partial  
697 venom gland transcriptome of a *Drosophila* parasitoid wasp, *Leptopilina*  
698 *heterotoma*, reveals novel and shared bioactive profiles with stinging  
699 Hymenoptera. *Gene* 526 (2):195-204.

700 Heavner, M.E., A.D. Hudgins, R. Rajwani, J. Morales, and S. Govind, 2014 Harnessing  
701 the natural *Drosophila*-parasitoid model for integrating insect immunity with  
702 functional venomics. *Curr Opin Insect Sci* 6:61-67.

703 Heavner, M.E., J. Ramroop, G. Gueguen, G. Ramrattan, G. Dolios *et al.*, 2017 Novel  
704 Organelles with Elements of Bacterial and Eukaryotic Secretion Systems  
705 Weaponize Parasites of *Drosophila*. *Curr Biol* 27 (18):2869-2877.e2866.

706 Ihaka, R., and R. Gentleman, 1996 R: A Language for Data Analysis and Graphics.  
707 *Journal of Computational and Graphical Statistics* 5 (3):299-314.

708 Jackman, S.D., B.P. Vandervalk, H. Mohamadi, J. Chu, S. Yeo *et al.*, 2017 ABySS 2.0:  
709 resource-efficient assembly of large genomes using a Bloom filter. *Genome Res*  
710 27 (5):768-777.

711 Johnson, M., I. Zaretskaya, Y. Raytselis, Y. Merezuk, S. McGinnis *et al.*, 2008 NCBI  
712 BLAST: a better web interface. *Nucleic Acids Res* 36 (Web Server issue):W5-9.

713 Jones, P., D. Binns, H.Y. Chang, M. Fraser, W. Li *et al.*, 2014 InterProScan 5: genome-  
714 scale protein function classification. *Bioinformatics* 30 (9):1236-1240.

715 Kajitani, R., D. Yoshimura, M. Okuno, Y. Minakuchi, H. Kagoshima *et al.*, 2019  
716 *Platanus*-allee is a de novo haplotype assembler enabling a comprehensive  
717 access to divergent heterozygous regions. *Nature Communications* 10.

718 Kall, L., A. Krogh, and E.L. Sonnhammer, 2004 A combined transmembrane topology  
719 and signal peptide prediction method. *J Mol Biol* 338 (5):1027-1036.

720 Kall, L., A. Krogh, and E.L.L. Sonnhammer, 2007 Advantages of combined  
721 transmembrane topology and signal peptide prediction - the Phobius web server.  
722 *Nucleic Acids Research* 35:W429-W432.

723 Kalra, H., R.J. Simpson, H. Ji, E. Aikawa, P. Altevogt *et al.*, 2012 Vesiclepedia: a  
724 compendium for extracellular vesicles with continuous community annotation.  
725 *PLoS Biol* 10 (12):e1001450.

726 Kanehisa, M., Y. Sato, and K. Morishima, 2016 BlastKOALA and GhostKOALA: KEGG  
727 Tools for Functional Characterization of Genome and Metagenome Sequences. *J*  
728 *Mol Biol* 428 (4):726-731.

729 Keebaugh, E.S., and T.A. Schlenke, 2014 Insights from natural host-parasite  
730 interactions: The Drosophila model. *Developmental and Comparative*  
731 *Immunology* 42 (1):111-123.

732 Keller, O., M. Kollmar, M. Stanke, and S. Waack, 2011 A novel hybrid gene prediction  
733 method employing protein multiple sequence alignments. *Bioinformatics* 27  
734 (6):757-763.

735 Khan, S., D.T. Sowpati, and R.K. Mishra, 2018 Long-read genome sequence and  
736 assembly of *Leptopilina boulardi*: a specialist *Drosophila* parasitoid. *bioRxiv*.

737 Kim, D., B. Landmead, and S.L. Salzberg, 2015 HISAT: a fast spliced aligner with low  
738 memory requirements. *Nature Methods* 12 (4):357-U121.

739 Kraaijeveld, K., S.Y. Anvar, J. Frank, A. Schmitz, J. Bast *et al.*, 2016 Decay of Sexual  
740 Trait Genes in an Asexual Parasitoid Wasp. *Genome Biology and Evolution* 8  
741 (12):3685-3695.

742 Krogh, A., B. Larsson, G. von Heijne, and E.L. Sonnhammer, 2001 Predicting  
743 transmembrane protein topology with a hidden Markov model: application to  
744 complete genomes. *J Mol Biol* 305 (3):567-580.

745 Labrosse, C., K. Staslak, J. Lesobre, A. Grangeia, E. Huguet *et al.*, 2005 A RhoGAP  
746 protein as a main immune suppressive factor in the *Leptopilina boulardi*  
747 (*Hymenoptera*, *Figitidae*) - *Drosophila melanogaster* interaction. *Insect*  
748 *Biochemistry and Molecular Biology* 35 (2):93-103.

749 Lane, N., 2009 On the origin of bar codes. *Nature* 462 (7271):272-274.

750 Li, H., and R. Durbin, 2009 Fast and accurate short read alignment with Burrows-  
751 Wheeler transform. *Bioinformatics* 25 (14):1754-1760.

752 Li, M., L.Y.C. Au, D. Douglah, A. Chong, B.J. White *et al.*, 2017a Generation of heritable  
753 germline mutations in the jewel wasp *Nasonia vitripennis* using CRISPR/Cas9.  
754 *Sci Rep* 7 (1):901.

755 Li, M., M. Bui, and O.S. Akbari, 2017b Embryo Microinjection and Transplantation  
756 Technique for *Nasonia vitripennis* Genome Manipulation. *J Vis Exp* (130).

757 Li, W., A. Cowley, M. Uludag, T. Gur, H. McWilliam *et al.*, 2015 The EMBL-EBI  
758 bioinformatics web and programmatic tools framework. *Nucleic Acids Res* 43  
759 (W1):W580-584.

760 Mapleson, D., G. Garcia Accinelli, G. Kettleborough, J. Wright, and B.J. Clavijo, 2016  
761 KAT: a K-mer analysis toolkit to quality control NGS datasets and genome  
762 assemblies. *Bioinformatics* 33 (4):574-576.

763 Marchler-Bauer, A., Y. Bo, L. Han, J. He, C.J. Lanczycki *et al.*, 2017 CDD/SPARCLE:  
764 functional classification of proteins via subfamily domain architectures. *Nucleic  
765 Acids Res* 45 (D1):D200-D203.

766 Martinson, E.O., Mrinalini, Y.D. Kelkar, C.H. Chang, and J.H. Werren, 2017 The  
767 Evolution of Venom by Co-option of Single-Copy Genes. *Curr Biol* 27 (13):2007-  
768 2013.e2008.

769 Melk, J.P., and S. Govind, 1999 Developmental analysis of *Ganaspis xanthopoda*, a  
770 larval parasitoid of *Drosophila melanogaster*. *J Exp Biol* 202 (Pt 14):1885-1896.

771 Mikheenko, A., G. Valin, A. Prjibelski, V. Saveliev, and A. Gurevich, 2016 Icarus:  
772 visualizer for de novo assembly evaluation. *Bioinformatics* 32 (21):3321-3323.

773 Misof, B., S. Liu, K. Meusemann, R.S. Peters, A. Donath *et al.*, 2014 Phylogenomics  
774 resolves the timing and pattern of insect evolution. *Science* 346 (6210):763-767.

775 Morales, J., H. Chiu, T. Oo, R. Plaza, S. Hoskins *et al.*, 2005 Biogenesis, structure, and  
776 immune-suppressive effects of virus-like particles of a *Drosophila* parasitoid,  
777 *Leptopilina victoriae*. *J Insect Physiol* 51 (2):181-195.

778 Mortimer, N.T., J. Goecks, B.Z. Kacsoh, J.A. Mobley, G.J. Bowersock *et al.*, 2013  
779 Parasitoid wasp venom SERCA regulates *Drosophila* calcium levels and inhibits  
780 cellular immunity. *Proc Natl Acad Sci U S A* 110 (23):9427-9432.

781 Narendran, T.C., 2001 Parasitic Hymenoptera and Biological Control, pp. 1-12 in  
782 *Biocontrol Potential and its Exploitation in Sustainable Agriculture: Volume 2:*

783 *Insect Pests*, edited by R.K. Upadhyay, K.G. Mukerji and B.P. Chamola. Springer  
784 US, Boston, MA.

785 Nielsen, H., 2017 Predicting Secretory Proteins with SignalP. *Methods Mol Biol*  
786 1611:59-73.

787 Nolte-'t Hoen, E., T. Cremer, R.C. Gallo, and L.B. Margolis, 2016 Extracellular vesicles  
788 and viruses: Are they close relatives? *Proceedings of the National Academy of*  
789 *Sciences* 113 (33):9155-9161.

790 Pannebakker, B.A., L.P. Pijnacker, B.J. Zwaan, and L.W. Beukeboom, 2004 Cytology of  
791 Wolbachia-induced parthenogenesis in *Leptopilina clavipes* (Hymenoptera:  
792 Figitidae). *Genome* 47 (2):299-303.

793 Pathan, M., P. Fonseka, S.V. Chitti, T. Kang, R. Sanwani *et al.*, 2019 Vesiclepedia  
794 2019: a compendium of RNA, proteins, lipids and metabolites in extracellular  
795 vesicles. *Nucleic Acids Res* 47 (D1):D516-D519.

796 Pathan, M., S. Keerthikumar, C.S. Ang, L. Gangoda, C.Y. Quek *et al.*, 2015 FunRich:  
797 An open access standalone functional enrichment and interaction network  
798 analysis tool. *Proteomics* 15 (15):2597-2601.

799 Pathan, M., S. Keerthikumar, D. Chisanga, R. Alessandro, C.S. Ang *et al.*, 2017 A novel  
800 community driven software for functional enrichment analysis of extracellular  
801 vesicles data. *J Extracell Vesicles* 6 (1):1321455.

802 Patot, S., R. Allemand, F. Fleury, and J. Varaldi, 2012 An inherited virus influences the  
803 coexistence of parasitoid species through behaviour manipulation. *Ecol Lett* 15  
804 (6):603-610.

805 Pichon, A., A. Bézier, S. Urbach, J.M. Aury, V. Jouan *et al.*, 2015 Recurrent DNA virus  
806 domestication leading to different parasite virulence strategies. *Sci Adv* 1  
807 (10):e1501150.

808 Pritchard, L., J.A. White, P.R.J. Birch, and I.K. Toth, 2006 GenomeDiagram: a python  
809 package for the visualization of large-scale genomic data. *Bioinformatics* 22  
810 (5):616-617.

811 Ramroop, J., 2016 Mechanisms of Immune Activation and Suppression by Parasitic  
812 Wasps of *Drosophila* in *Biology*. Graduate Center, City University of New York.

813 Rizki, T.M., R.M. Rizki, and Y. Carton, 1990 *Leptopilina heterotoma* and *L. boulardi*:  
814 strategies to avoid cellular defense responses of *Drosophila melanogaster*. *Exp*  
815 *Parasitol* 70 (4):466-475.

816 Rodriguez, J.J., J.L. Fernandez-Triana, M.A. Smith, D.H. Janzen, W. Hallwachs *et al.*,  
817 2013 Extrapolations from field studies and known faunas converge on  
818 dramatically increased estimates of global microgastrine parasitoid wasp species  
819 richness (Hymenoptera: Braconidae). *Insect Conservation and Diversity* 6  
820 (4):530-536.

821 Schlenke, T.A., J. Morales, S. Govind, and A.G. Clark, 2007 Contrasting infection  
822 strategies in generalist and specialist wasp parasitoids of *Drosophila*  
823 *melanogaster*. *PLoS Pathog* 3 (10):1486-1501.

824 Siebert, A.L., D. Wheeler, and J.H. Werren, 2015 A new approach for investigating  
825 venom function applied to venom calreticulin in a parasitoid wasp. *Toxicon*.

826 Simão, F.A., R.M. Waterhouse, P. Ioannidis, E.V. Kriventseva, and E.M. Zdobnov, 2015  
827 BUSCO: assessing genome assembly and annotation completeness with single-  
828 copy orthologs. *Bioinformatics* 31 (19):3210-3212.

829 Small, C., I. Paddibhatla, R. Rajwani, and S. Govind, 2012 An introduction to parasitic  
830 wasps of *Drosophila* and the antiparasite immune response. *J Vis Exp*  
831 (63):e3347.

832 Smith-Unna, R., C. Boursnell, R. Patro, J.M. Hibberd, and S. Kelly, 2016 TransRate:  
833 reference-free quality assessment of de novo transcriptome assemblies.  
834 *Genome Research* 26 (8):1134-1144.

835 Sonnhammer, E.L., G. von Heijne, and A. Krogh, 1998 A hidden Markov model for  
836 predicting transmembrane helices in protein sequences. *Proc Int Conf Intell Syst*  
837 *Mol Biol* 6:175-182.

838 Stanke, M., and B. Morgenstern, 2005 AUGUSTUS: a web server for gene prediction in  
839 eukaryotes that allows user-defined constraints. *Nucleic Acids Res* 33 (Web  
840 Server issue):W465-467.

841 Stanke, M., R. Steinkamp, S. Waack, and B. Morgenstern, 2004 AUGUSTUS: a web  
842 server for gene finding in eukaryotes. *Nucleic Acids Res* 32 (Web Server  
843 issue):W309-312.

844 Strand, M.R., and G.R. Burke, 2015 Polydnviruses: From discovery to current insights.  
845 *Virology* 479-480:393-402.

846 Trapnell, C., A. Roberts, L. Goff, G. Pertea, D. Kim *et al.*, 2012 Differential gene and  
847 transcript expression analysis of RNA-seq experiments with TopHat and  
848 Cufflinks. *Nat Protoc* 7 (3):562-578.



849 Varaldi, J., S. Petit, M. Boulétreau, and F. Fleury, 2006 The virus infecting the parasitoid  
850 *Leptopilina boulardi* exerts a specific action on superparasitism behaviour.  
851 *Parasitology* 132 (Pt 6):747-756.

852 Volkoff, A.N., V. Jouan, S. Urbach, S. Samain, M. Bergoin *et al.*, 2010 Analysis of virion  
853 structural components reveals vestiges of the ancestral ichnovirus genome.  
854 *PLoS Pathog* 6 (5):e1000923.

855 Wan, B., E. Goguet, M. Ravallec, O. Pierre, S. Lemauf *et al.*, 2019 Venom Atypical  
856 Extracellular Vesicles as Interspecies Vehicles of Virulence Factors Involved in  
857 Host Specificity: The Case of a *Drosophila* Parasitoid Wasp. *Frontiers in*  
858 *Immunology* 10 (1688).

859 Wang, Z., Y. Chen, and Y. Li, 2004 A brief review of computational gene prediction  
860 methods. *Genomics Proteomics Bioinformatics* 2 (4):216-221.

861 Waterhouse, R.M., M. Seppey, F.A. Simão, M. Manni, P. Ioannidis *et al.*, 2017 BUSCO  
862 applications from quality assessments to gene prediction and phylogenomics.  
863 *Mol Biol Evol.*

864 Werren, J.H., L. Baldo, and M.E. Clark, 2008 *Wolbachia*: master manipulators of  
865 invertebrate biology. *Nat Rev Microbiol* 6 (10):741-751.

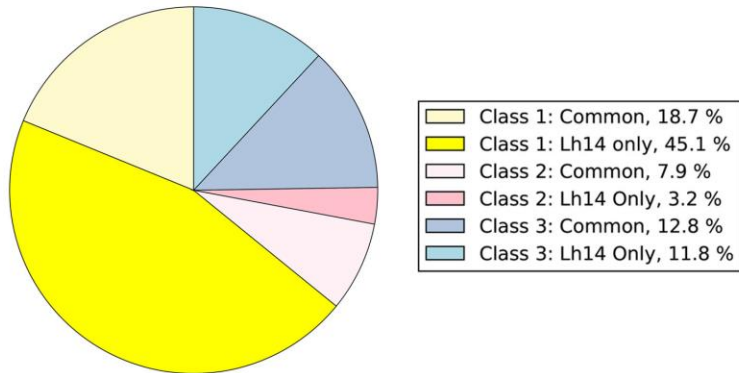
866 Werren, J.H., D.W. Loehlin, and J.D. Giebel, 2009 Larval RNAi in *Nasonia* (parasitoid  
867 wasp). *Cold Spring Harb Protoc* 2009 (10):pdb.prot5311.

868 Young, K.A., A.P. Herbert, P.N. Barlow, V.M. Holers, and J.P. Hannan, 2008 Molecular  
869 basis of the interaction between complement receptor type 2 (CR2/CD21) and  
870 Epstein-Barr virus glycoprotein gp350. *J Virol* 82 (22):11217-11227.

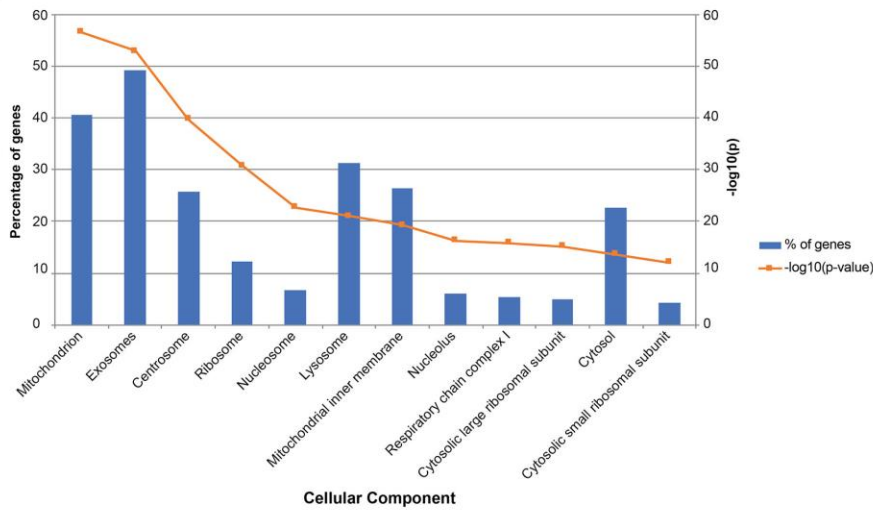
871 Zhang, S.V., L.T. Zhuo, and M.W. Hahn, 2016 AGOUTI: improving genome assembly  
872 and annotation using transcriptome data. *Gigascience* 5.  
873 Zimin, A.V., G. Marcais, D. Puiu, M. Roberts, S.L. Salzberg *et al.*, 2013 The MaSuRCA  
874 genome assembler. *Bioinformatics* 29 (21):2669-2677.  
875  
876

**FIGURES**

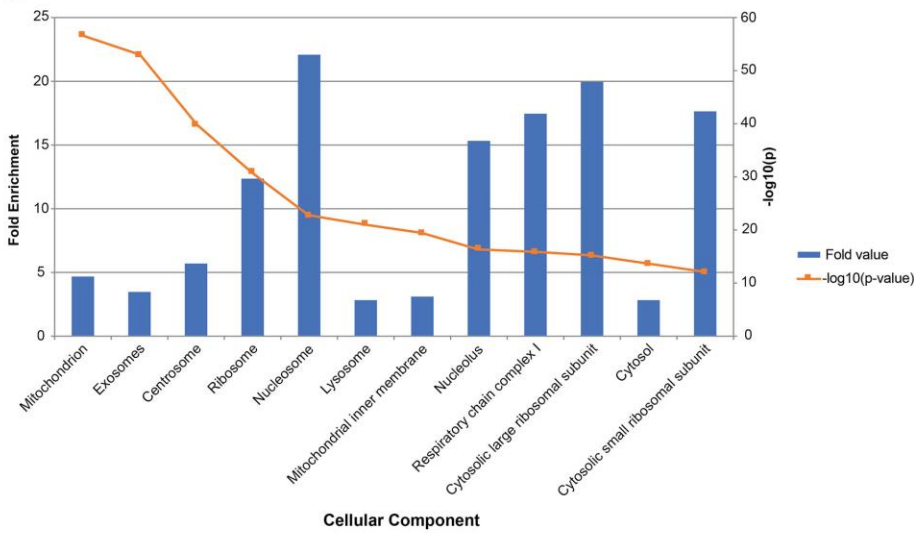
A.



B.

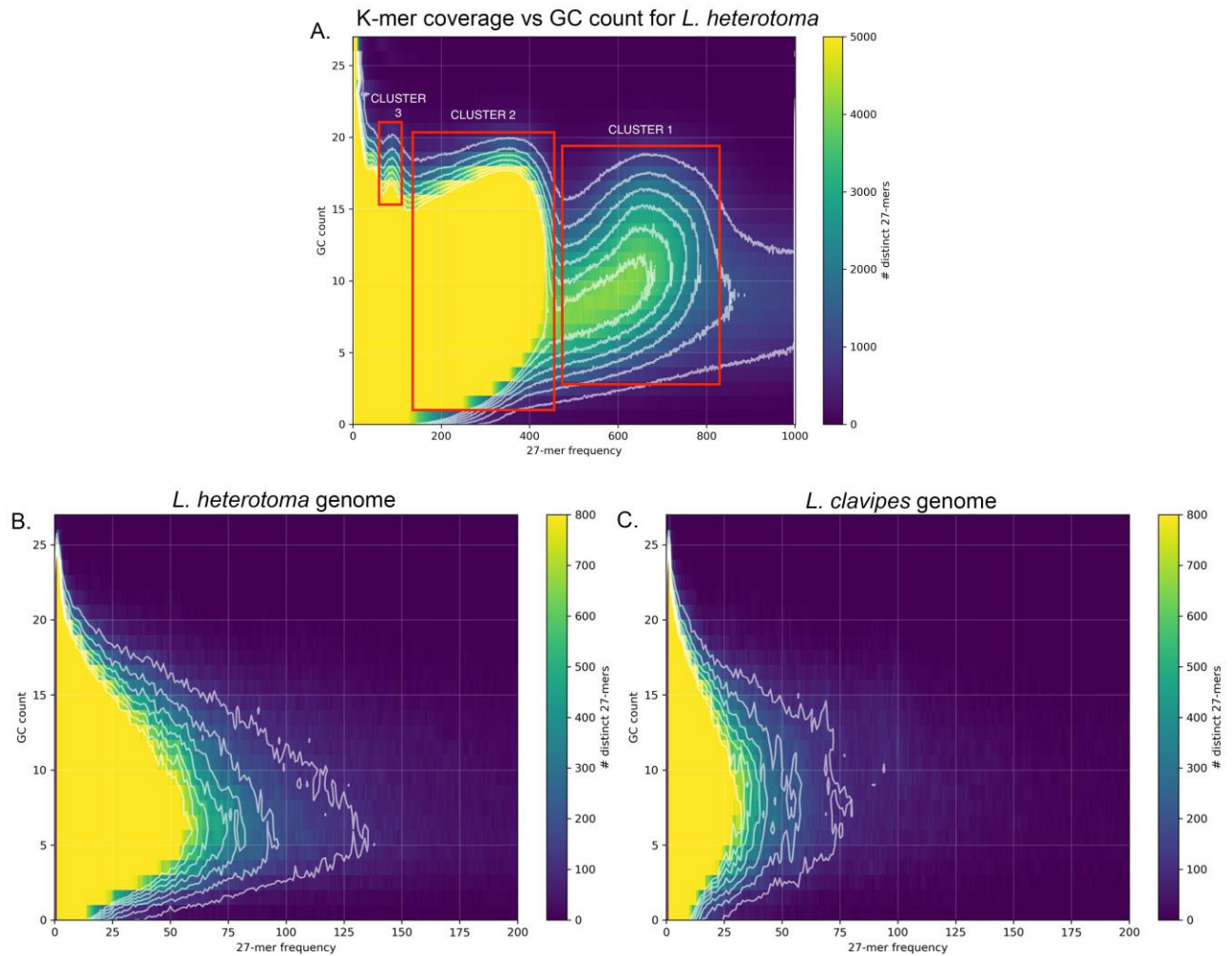


C.



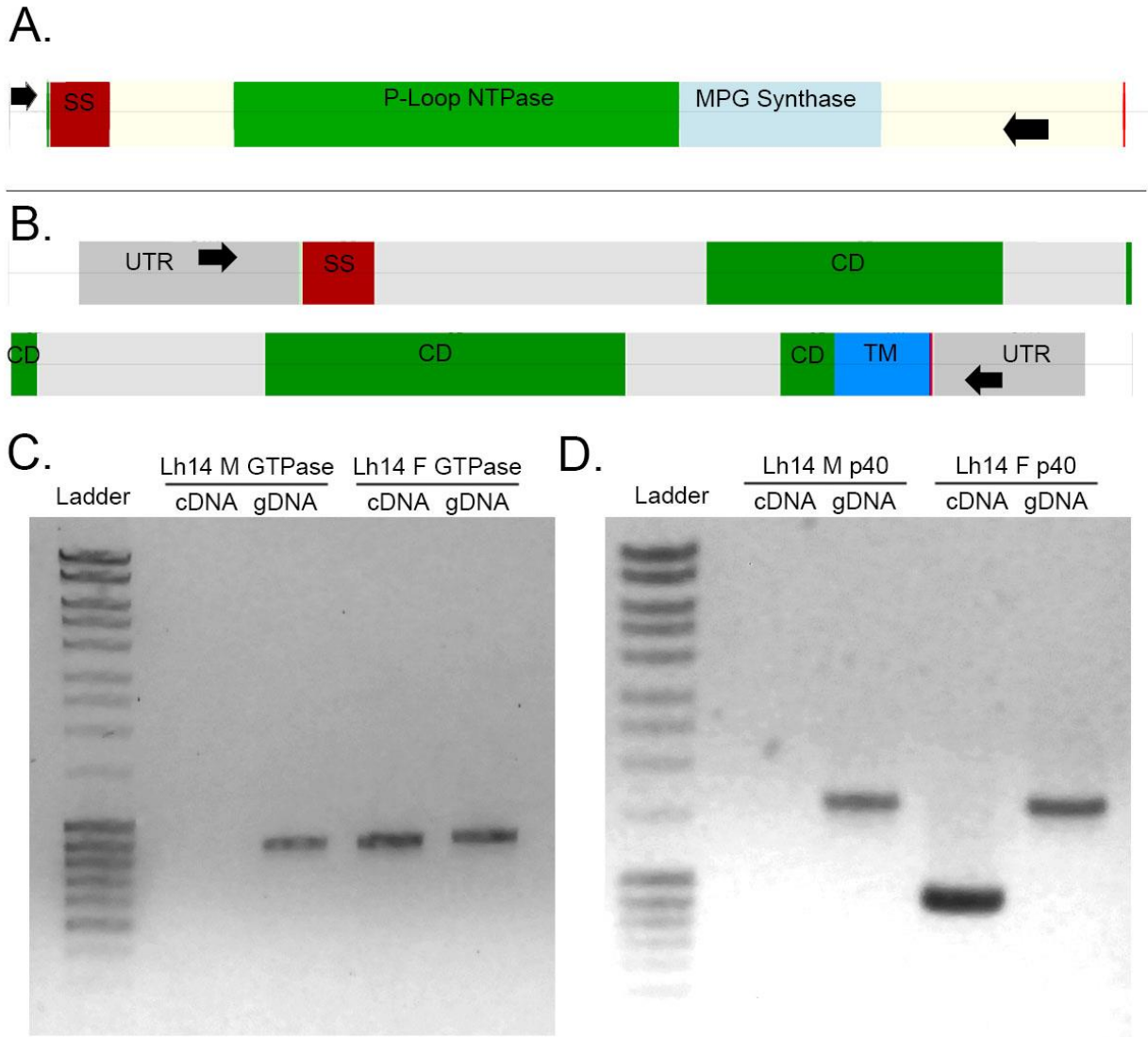
879 Figure 1. The superset of MSEV proteins: (A) *Lh* 14 MSEV proteins were annotated  
880 using BLAST2GO prior to class sorting via annotation and GO Terms. Wedges denoted  
881 as “Common,” were previously published in (Heavner et al. 2017) and represent  
882 proteins found in both *Lh* 14 and *Lh* NY MSEV proteomes. New proteins analyzed in  
883 this work are in wedges labeled “*Lh* 14 Only.” A majority of proteins belong to Class 1.  
884 Table S1 lists 246 proteins added to the superset *Lh* 14 proteome. (B) and (C)  
885 Enrichment analysis of MSEV superset shows high association with exosomes and  
886 mitochondria compared to other cellular organelles according to Vesiclepedia.  $-\log_{10}$   
887 (p-value) trend shown in orange for both graphs. The p-values were calculated with the  
888 Bonferroni method. (B) Percentage of MSEV genes associated with specific cellular  
889 compartments found in Vesiclepedia, relative to all MSEV genes. Of the superset  
890 proteins, 41 and 49% are associated with mitochondria and exosomes, respectively ( $p =$   
891  $3 \times 10^{-57}$ ;  $1 \times 10^{-53}$ ). (C) Fold-enrichment of the MSEV dataset in specific cellular  
892 compartments. Although many protein classes are present in the proteome, exosomal  
893 and mitochondrial proteins show more significant enrichments.

894



895  
 896 Figure 2. Analysis of K-mer coverage versus GC count. (A) Analysis of genomic reads.  
 897 27-mers generated from the cleaned Illumina reads used to assemble the *L. heterotoma*  
 898 genome binned by their GC count vs multiplicity (total counts among the reads). Bins  
 899 are colored by the number of distinct K-mers. Different clusters are identified as shown  
 900 and described in the text. (B and C): A map of 27-mer multiplicity versus GC content of  
 901 the joint assembly of the *Lh 14* genome (B) to a map from the published *L. clavipes*  
 902 genome (Bioproject: PRJNA84205) (C).

903



904

905 Figure 3. Predicted gene structures verified by PCR amplification experiments (A, B).  
 906 Diagrams showing primer locations and predicted gene structures of *SmGTPase01* (A)  
 907 and *p40* (B). Black arrows indicate primer locations, light gray indicates introns, UTR  
 908 regions are dark gray and labeled, exons encoding potential protein domains are  
 909 labeled as shown. Cream colored regions in panel A do not have a specified domain.  
 910 Diagrams were drawn using GenomeDiagram as part of the Biopython (v. 1.6) package  
 911 (Pritchard et al. 2006; Cock et al. 2009). Each row in the panels A and B diagrams

912 corresponds to approximately 1,000 bp. For primer sequences, see methods. (C and D)  
913 Ladder is Thermo Fisher MassRuler ladder. (C) PCR products for *SmGTPase01* from  
914 male or female cDNA and gDNA. All products are 873 bp long. Male cDNA PCR was  
915 negative. (D) PCR products for *p40* from male or female cDNA and gDNA. The  
916 expected band for *p40* cDNA is 939 bp and for gDNA is 1,630 bp. Male cDNA PCR was  
917 negative. Sequence analysis of PCR amplification products confirmed gene prediction  
918 results.

919

920

**TABLES**

921

<b>MSEV SUPERSET UNKNOWNNS CDD-SEARCH RESULTS</b>						
Query (in-house ID)	PSSM-ID	From	To	E-Value	Accession	Short name
GAJC01013214.1_14	<b>331760</b>	<b>25</b>	<b>98</b>	<b>0.000176</b>	cl26939	DEXDc superfamily
GAJC01012558.1_12	<b>330317</b>	<b>39</b>	<b>205</b>	<b>0.003987</b>	cl25496	Herpes_BLLF1 superfamily/gp350
GAJC01011863.1_13	<b>311912</b>	<b>86</b>	<b>187</b>	<b>0.003653</b>	cl07006	RNA_poll_A34 superfamily
GAJC01011463.1_48	<b>315064</b>	<b>234</b>	<b>335</b>	<b>0.002964</b>	cl13702	CD99L2 superfamily
GAJC01010930.1_16	<b>328726</b>	<b>32</b>	<b>61</b>	<b>0.001252</b>	cl21457	ICL_KPHMT superfamily
GAJC01010353.1_14	<b>331876</b>	<b>31</b>	<b>121</b>	<b>0.001483</b>	cl27055	MutS_III superfamily
GAJC01009713.1_25	<b>311628</b>	<b>138</b>	<b>225</b>	<b>0.000133</b>	cl06688	TSGP1 superfamily
GAJC01009493.1_4	<b>328724</b>	<b>79</b>	<b>96</b>	<b>0.001983</b>	cl21455	P-loop_NTPase superfamily
GAJC01002124.1_43	<b>330572</b>	<b>4</b>	<b>269</b>	<b>0.0073146</b>	cl25751	DUF4045 superfamily

922

923 Table 1. CDD-search results of MSEV “un-annotated” proteins in the super-set. MSEV  
924 ORFs that completed the BLAST2GO pipeline and did not return any results were run  
925 through the NCBI CDD-Search Version 3.16 (Accessed: Aug. 2018). Of 45 queries, only  
926 9 returned hits with threshold set to  $1 \times 10^{-2}$ . The ninth result came from a search with E-  
927 value threshold set to 1. Results listed are all unique, high scoring hits for each ORF  
928 that returned hits from the search.

929



930

<b>ASSEMBLY STATISTICS</b>				
		Male	Female	Joint
Assembly	N50 (bp)	<b>4,779</b>	<b>4,843</b>	<b>11906</b>
	No. scaffolds	<b>147,558</b>	<b>147,549</b>	<b>83,487</b>
	Largest scaffold (bp)	<b>306,667</b>	<b>176,371</b>	<b>375,275</b>
	Total length (bp)	<b>474,383,205</b>	<b>472,302,230</b>	<b>462,564,754</b>
	GC%	<b>27.54</b>	<b>27.28</b>	<b>27.84</b>
	Coverage (%)	<b>87.7</b>	<b>86.8</b>	<b>91.1</b>
BUSCOs (Insecta)	Complete	<b>80.20%</b>	<b>81.30%</b>	<b>90.90%</b>
	Single	<b>69.40%</b>	<b>71.80%</b>	<b>89.20%</b>
	Duplicated	<b>10.80%</b>	<b>9.50%</b>	<b>1.70%</b>
	Fragmented	<b>15.90%</b>	<b>15.10%</b>	<b>6.50%</b>
	Missing	<b>3.90%</b>	<b>3.60%</b>	<b>2.60%</b>
	n	<b>1658</b>	<b>1658</b>	<b>1658</b>

931

932 Table 2. Assembly statistics: Statistics of male, female, and combined (male plus  
933 female) *Lh* genomes as assessed by QUASTv4.0 and BUSCOv9.0. Percent coverage  
934 was found by mapping sequencing reads back to assembly using HISAT2. The  
935 identified BUSCOs can be found in Table S1. The QUAST program was run with  
936 parameters set for eukaryotic genomes and scaffolds. The BUSCO program was run  
937 with species set to 'Nasonia.' Contigs smaller than 500 bp were excluded.

938

939

<b>MSEV GENES FOUND IN GENOME ANALYSIS</b>				
	MSEV BLASTn scaffold results		AUGUSTUS prediction results	
	Found	Percentage	Found	Percentage
Female	<b>278</b>	<b>68.3</b>	<b>169</b>	<b>41.52</b>
Male	<b>275</b>	<b>67.58</b>	<b>166</b>	<b>40.78</b>
Shared in M+F	<b>265</b>		<b>159</b>	
Joint Assembly	<b>375</b>	<b>92.13%</b>	<b>325</b>	<b>79.85%</b>

940

941 Table 3. MSEV genes found in scaffolds and predictions: Gene predictions from  
942 genome assembly scaffolds and AUGUSTUS gene predictions were searched for  
943 MSEV genes using tBLASTn. Results better than %ID >70%, E-value <  $1 \times 10^{-50}$ , and  
944 query coverage > 70% were retained.

