

## Biocuration as an undergraduate training experience: Improving the annotation of the insect vector of Citrus greening disease

Surya Saha<sup>α,β,1,\*</sup>, Prashant S Hosmani<sup>α,β,1,\*</sup>, Krystal Villalobos-Ayala<sup>α,β,2</sup>, Sherry Miller<sup>α,β,3</sup>, Teresa Shippy<sup>α,β,3</sup>, Andrew Rosendale<sup>α,β,4</sup>, Chris Cordola<sup>β,2</sup>, Tracey Bell<sup>β,2</sup>, Hannah Mann<sup>β,2</sup>, Gabe DeAvila<sup>β,2</sup>, Daniel DeAvila<sup>β,2</sup>, Zachary Moore<sup>β,4</sup>, Kyle Buller<sup>β,4</sup>, Kathryn Ciolkevich<sup>β,4</sup>, Samantha Nandyal<sup>β,4</sup>, Robert Mahoney<sup>β,4</sup>, Joshua Van Voorhis<sup>β,4</sup>, Megan Dunlevy<sup>β,4</sup>, David Farrow<sup>β,4</sup>, David Hunter<sup>β,3</sup>, Taylar Morgan<sup>β,3</sup>, Kayla Shore<sup>β,3</sup>, Victoria Guzman<sup>β,3</sup>, Allison Izsak<sup>β,5</sup>, Danielle E Dixon<sup>β,1,6</sup>, Liliana Cano<sup>β,7</sup>, Andrew Cridge<sup>β,8</sup>, Shannon Johnson<sup>β,9</sup>, Brandi L Cantarel<sup>β,10</sup>, Stephen Richardson<sup>β,11,12</sup>, Adam English<sup>β,11,12</sup>, Nan Leng<sup>β,13</sup>, Xiaolong Cao<sup>γ,14</sup>, Haobo Jiang<sup>γ,15</sup>, Chris Childers<sup>α,16</sup>, Mei-Ju Chen<sup>α,17</sup>, Mirella Flores<sup>α,γ,1</sup>, Wayne Hunter<sup>β,γ,18</sup>, Michelle Cilia<sup>γ,19,20</sup>, Lukas A Mueller<sup>γ,1,21</sup>, Monica Munoz-Torres<sup>α,β,22</sup>, David Nelson<sup>β,23</sup>, Monica F. Poelchau<sup>α,16</sup>, Josh Benoit<sup>β,γ,4</sup>, Helen Wiersma-Koch<sup>α,β,2</sup>, Tom D'Elia<sup>β,γ,2</sup>, Susan J Brown<sup>β,γ,3</sup>

Correspondence should be addressed to Surya Saha (ss2489@cornell.edu).

\* Surya Saha and Prashant S. Hosmani contributed equally to this work.

Roles: <sup>α</sup> Key personnel; <sup>β</sup> Annotator; <sup>γ</sup> Genome Sequencing; <sup>γ</sup> Transcriptome; <sup>γ</sup> Principal Investigator

Affiliations: <sup>1</sup> Boyce Thompson Institute, Ithaca, NY 14853, USA; <sup>2</sup> Indian River State College, Fort Pierce, FL 34981, USA; <sup>3</sup> Kansas State University, Manhattan, KS 66506, USA; <sup>4</sup> University of Cincinnati, Cincinnati, OH 45220, USA; <sup>5</sup> Cornell University, Ithaca, NY 14853, USA; <sup>6</sup> University of Puget Sound, Tacoma, WA 98416, USA; <sup>7</sup> University of Florida/ IFAS Indian River Research and Education Center, Plant Pathology, Ft. Pierce, FL 34945, USA; <sup>8</sup> University of Otago, North Dunedin, Dunedin 9016, New Zealand; <sup>9</sup> Los Alamos National Laboratory, NM 87544, USA; <sup>10</sup> UT Southwestern Medical Center, Bioinformatics Core Facility, Department of Bioinformatics, Dallas, TX 75390, USA; <sup>11</sup> Baylor College of Medicine, i5K Arthropod Genomics, Houston, TX 77030, USA; <sup>12</sup> Baylor College of Medicine, Human Genome Sequencing Center, Houston, TX 77030, USA; <sup>13</sup> Illumina Inc., San Diego, CA 92122, USA; <sup>14</sup> Oklahoma State University, Department of Biochemistry and Molecular Biology, Stillwater, OK 74074, USA; <sup>15</sup> Oklahoma State University, Department of Entomology and Plant Pathology, Stillwater, OK 74074, USA; <sup>16</sup> USDA Agricultural Research Service, National Agricultural Library, Beltsville, MD 20705, USA; <sup>17</sup> National Taiwan University, Graduate Institute of Biomedical Electronics and Bioinformatics, Taipei 10617, TW; <sup>18</sup> USDA ARS, U. S. Horticultural Research Laboratory, Ft. Pierce, FL 34945, USA; <sup>19</sup> USDA ARS, Emerging Pests and Pathogens Research Unit, Ithaca, NY 14853, USA; <sup>20</sup> Cornell University, Plant Pathology and Plant-Microbe Biology Section, School of Integrative Plant Science, Ithaca, NY 14853, USA; <sup>21</sup> Cornell University, Plant Breeding and Genetics Section, School of Integrative Plant Science, Ithaca, NY 14853, USA; <sup>22</sup> Lawrence Berkeley National Laboratory, Environmental Genomics and Systems Biology, Berkeley, CA 94720, USA; <sup>23</sup> The University of Tennessee Health Science Center, Department of Microbiology, Immunology and Biochemistry, Memphis, TN 38163, USA;

## ABSTRACT

The Asian citrus psyllid (*Diaphorina citri* Kuwayama) is the insect vector of the bacterium *Candidatus Liberibacter asiaticus* (CLas), the pathogen associated with citrus Huanglongbing (HLB, citrus greening). HLB threatens citrus production worldwide. Suppression or reduction of the insect vector using chemical insecticides has been the primary method to inhibit the spread of citrus greening disease. Accurate structural and functional annotation of the Asian citrus psyllid genome, as well as a clear understanding of the interactions between the insect and CLas, are required for development of new molecular-based HLB control methods. A draft assembly of the *D. citri* genome has been generated and annotated with automated pipelines. However, knowledge transfer from well-curated reference genomes such as that of *Drosophila melanogaster* to newly sequenced ones is challenging due to the complexity and diversity of insect genomes. To identify and improve gene models as potential targets for pest control, we manually curated several gene families with a focus on genes that have key functional roles in *D. citri* biology and CLas interactions. This community effort produced 530 manually curated gene models across developmental, physiological, RNAi regulatory, and immunity-related pathways. As previously shown in the pea aphid, RNAi machinery genes putatively involved in the microRNA pathway have been specifically duplicated. A comprehensive transcriptome enabled us to identify a number of gene families that are either missing or misassembled in the draft genome. In order to develop biocuration as a training experience, we included undergraduate and graduate students from multiple institutions, as well as experienced annotators from the insect genomics research community. The resulting gene set (OGS v1.0) combines both automatically predicted and manually curated gene models. All data are available on <https://citrusgreening.org/>.

## KEYWORDS

Asian citrus psyllid, biocuration, i5k workspace at NAL, training, annotation, insect immunity

## ACKNOWLEDGEMENTS

Kascha Bohnenblust from Kansas State University assisted with the organization of meetings for the community curation effort. The following personnel at Los Alamos National Laboratory, Los Alamos, NM contributed to the Illumina and Pacbio sequencing: Kimberly McMurphy, Cheryl D Gleasner, Krista Reitenga, Shunsheng (Cliff) Han and Goutam Gupta. Christian Haudenschild from Illumina, Inc. supported the genome assembly. The following personnel at USDA-ARS, U.S. Horticultural Research Laboratory, Fort Pierce, FL contributed to the sample preparation for genome sequencing: Maria T. Gonzalez, Belkis Diego and Kathy Moulton.

## AUTHOR CONTRIBUTIONS

SS, PSH, KVA, SM, TS, AR, CC, TB, HM, GD, DD, ZM, KB, KC, SN, RM, JVV, MD, DF, DH, TM, KS, VG, AI, DED, LC, AC, CC, MC, MF, WH, MMT, DN, MFP, JB, HKW, TD and SJB contributed to the community curation. XC, HJ and MF generated the MCOT v1.0 transcriptome. WH, SJ, BLC, SR, AE and NL were involved in the genome sequencing. WH, MC, LAM, JB, TD and SJB were the Principal Investigators involved in this work. SS, PSH, SM, TS, MF, WH, MC, MMT, DN, MFP, JB, TD and SJB wrote the final manuscript.

## INTRODUCTION

The Asian citrus psyllid (ACP), *Diaphorina citri* Kuwayama (Hemiptera:Liviidae), is a phloem-feeding insect native to Southeastern and Southwestern Asia with a host range limited to plants in the citrus genus and related Rutaceae spp. (1). Accidental anthropogenic introductions of psyllid-infested citrus combined with the ability of psyllids to disperse rapidly has allowed *D. citri* to extend its distribution to most of southern and eastern Asia, the Arabian Peninsula, the Caribbean, and South, Central and North America(1–6). For years, ACP has been classified as a global pest that is capable of devastating citrus crops through transmission of the bacterial agent, *Candidatus Liberibacter asiaticus*, CLas, which is associated with Huanglongbing (HLB) or citrus greening disease. The psyllid alone has little economic importance and causes only minor plant damage while feeding (7,8).

HLB is the most destructive and economically important disease of citrus, with practically all commercial citrus species and cultivars susceptible to CLas infection(9). Infected trees yield premature, bitter and misshapen fruit that is unmarketable. In addition, tree death follows 5-10 years after initial infection (2,9,10). Furthermore, HLB drastically suppresses economic progress in southern and eastern Asia by impeding viable commercial citrus agriculture within those regions (11). Florida is one of the top citrus-producing regions in the world and the largest in the United States, with nearly double the output of California, the second largest citrus-producing state (12). HLB puts the 9 billion dollar Florida citrus industry, with an annual net value of 1.5 billion dollars, at tremendous risk [USDA 2009]. In 2008, the HLB infection rate within central Florida was low (1.4% to 3.6%), but reaching 100% in the southern and eastern portions of the state (13,14). In 2005, when HLB was first detected in Florida, 9.3 million tons of oranges were harvested, but production has declined steadily to 5.3 million tons in 2016 as ACP and HLB have spread (12).

Primary management strategies focus on disrupting the HLB transmission pathway by suppressing psyllid populations and impeding interactions between CLas and psyllids. These strategies currently rely on extensive chemical application, which has broad environmental impact and high costs, and are ultimately unsustainable. To develop molecular methods that exploit current gene-targeting technologies, detailed genetic and genomic knowledge, including a high quality official gene set (OGS), is required (15,16). Early efforts focused on *D. citri* transcript expression(3,17–19), analysis of the full transcriptome (20,21) and, more recently, analysis of the *D. citri* proteome (15). Arp et al. (22) performed a BLAST-based inventory of NCBI-predicted immune genes in *D. citri* (v100, see Methods). In contrast, we have conducted

broad structural and functional annotation with the aid of a comprehensive transcriptome and created an official gene set for *D. citri* with a focus upon completing the repertoire of immune genes.

Manual curation improves the quality of gene annotation and establishing a ‘version controlled’ official gene set (OGS) provides a set of high quality, well-documented genes for the entire research community. Although ACP is a significant agricultural pest, it is not a model organism and the size of the research community does not warrant “museum” or “jamboree” annotation strategies (23). To maximize the number of genes annotated in a relatively short time, we augmented the “cottage industry” strategy (24) by training undergraduates to perform basic annotation tasks. The dispersed ACP annotation community agreed on a set of standard operating procedures and defined a set of primary gene targets. Starting with automated gene predictions, we used several additional types of evidence including RNAseq and proteomics data including comparisons to other insects to generate the first official gene set (OGS v1.0). Using an independent transcriptome enabled us to identify a number of gene families that are either missing or misassembled in the draft genome.

Our manual curation efforts focused on genes of potential use in vector control including immunity-related genes and pathways, RNAi machinery genes, multiple clans of cytochrome P450 genes and other genes relevant to insect development and physiology. We speculate that targeted analysis of these genes in *D. citri* will provide the foundation for a better understanding of the interactions between psyllid host and CLas pathogen, and will open the possibilities for research that can eventually find solutions to manage the dispersion of this very destructive pest and HLB.

## RESULTS AND DISCUSSION

### NCBI-Diaci 1.1 draft genome assembly

The Diaci1.1 draft genome assembly was generated using Illumina paired-end and mate-pair data with low coverage Pacbio for scaffolding and uploaded to NCBI (PRJNA251515) after filtering out bacterial contamination. This genome contains 161,988 scaffolds with an N50 of 109.8kb. Given the high degree of fragmentation, we performed a Benchmarking sets of Universal Single-Copy Orthologs (BUSCO) version 1 (25) analysis using a set of conserved single-copy markers, which showed that a significant number of these genes were missing (26%) or fragmented (22%) as described in the curation section below.

### MCOT transcriptome

To generate a more comprehensive set of gene models, we used the MCOT pipeline (26) to produce a *D. citri* MCOT transcriptome assembly that combines transcripts from genome-based Maker (27,28) and Cufflinks (29) pipelines with RNAseq-based (*de novo*) Trinity (30) and Oases (31) assemblers to select the best models. *D. citri* MCOT v1.0 contains 30,562 CDS, transcripts

and proteins, some of which are based solely on transcript evidence from RNAseq ([ftp://ftp.citrusgreening.org/genomes/Diaphorina\\_citri/genome/diaci1.1/](ftp://ftp.citrusgreening.org/genomes/Diaphorina_citri/genome/diaci1.1/)).

The completeness of *D. citri* MCOT v1.0 was assessed based on comparison to the BUSCO v2 beta (25). BUSCO v2 is based on a set of 1,066 single-copy orthologs from 133 arthropods species ORTHODB v9 (32). *D. citri* MCOT v1.0 contains 1039/1066 (97.5%) complete BUSCO orthologs, 808 (75.8%) of which are single copy and 231 (21.7%) of which appear to be duplicated. Four additional BUSCO orthologs (0.4%) are fragmented, and only 23 (2.2%) were not present. BUSCO analysis was also performed on the other resources used for annotation including three stage-specific *de novo* assembled transcriptomes from egg, nymph and adult tissue (20), the NCBI-Diaci1.1 genome assembly itself, the Maker v1.1 and NCBI v100 predicted gene models, and the *D. citri* MCOT v1.0 gene models that could be mapped to the genome assembly (Figure 1 and Supplementary Table 1). The Maker v1.1 annotation set contains 18,205 protein-coding gene models. NCBI *D. citri* Annotation Release 100 (NCBI v100, [https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/Diaphorina\\_citri/100/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Diaphorina_citri/100/)) contains a total of 19,311 protein-coding gene models along with annotations of non-coding RNAs (776) and pseudogenes (207). As shown in Figure 1, none of these data sets proved to be as complete as *D. citri* MCOT v1.0, which contained a higher percentage of complete BUSCO orthologs and fewer fragmented or missing orthologs.

The functional annotation was generated using InterproScan5, BLAST and AHRD (33), which assigned function descriptions to 23,098 genes (75.6%), GO annotations to 15,314 genes (50%) and Pfam domains to 18,170 genes (59%). MCOT is available at <ftp://ftp.citrusgreening.org/annotation/MCOT/>.

## Manual curation workflow

Manual curation of the *D. citri* (NCBI-Diaci1.1) draft genome assembly was undertaken to improve the quality of automated annotations (Supplementary Table 2) produced by the Maker (Maker v1.1) and NCBI pipelines (NCBI v100). This community-based manual annotation was focused on immunity-related genes as targets for ACP control.

The Apollo Genome Annotation Editor hosted at i5k Workspace@NAL (<https://i5k.nal.usda.gov/>) was implemented for community-based curation of gene models (34). Multiple evidence tracks (Supplementary Table 3) were added to Apollo to assist in manual curation. Standard operating procedures were outlined at the beginning and refined based on feedback from annotators and availability of evidence resources (see Methods). A typical workflow for manual gene curation involved selection of orthologs from related species, search of the ACP genome and assessment, followed by correction, of the gene models based on evidence tracks to generate the final model. Any exceptions to this general workflow are specified in the respective gene reports (Supplementary Notes 1-39). We used ImmunoDB (35) as a primary source of orthologs for curation of immune genes (Supplementary Notes 1-31). However, we also report other gene families of functional and evolutionary importance (Supplementary Notes 32-39).



including aquaporins, cuticle proteins and secretory proteins. Following correction, final gene models were verified using reciprocal BLAST analysis. With this community-based curation effort, we annotated a total of 530 genes, the majority of which include genes predicted to function in immunity, development and physiology.

### **Annotation Edit Distance (AED) as a measure of quality for different annotations**

We employed the Annotation Edit Distance (AED) (28,36) metric to evaluate the quality of the different annotation pipelines based on evidence from expression data. The AED had a two-fold application here, selection of the best predicted gene set and quantification of improvement after manual curation. To generate AED scores, we compared the annotations to known insect proteins (NCBI taxonomy:"Hexapoda [6960]") and to a comprehensive genome-guided transcriptome assembled from the latest data (See Methods and Supplementary Table 4) and the NCBI-Diaci1.1 draft assembly using the Maker pipeline (27). Most of the RNAseq data used to create this transcriptome was not used in the prediction of genes by other pipelines (NCBI v100, Maker v1.1 and MCOT v1.0) since it had not yet been produced and therefore provides independent validation. RNAseq data used for building the transcriptome included adult, nymph, egg, CLas exposed and healthy (adult and nymph tissue) as well as gut-specific expression data (see Supplementary Table 4). We calculated AED scores for NCBI v100, Maker v1.1 and mapped MCOT v1.0 annotation sets.

The plot for the AED cumulative fraction of transcripts (Figure 2) shows significant expression evidence support for NCBI v100 genes compared to the Maker v1.1 and mapped MCOT v1.0 annotation sets. Therefore the NCBI v100 annotations were selected to create the official gene set v1 (OGS v1.0). The AED plot also shows that there are many gene models in all the annotation sets that do not have any expression evidence support (AED = 1.0). This could be due either to lack of RNAseq data for some of the loci or incorrect annotations. The predicted annotations may have been affected by misassemblies in the NCBI-Diaci1.1 draft genome.

To quantify improvements in the gene structure by manual curation, we calculated AED scores for only the curated genes (530 genes). Cumulative fractions of transcripts of AED for curated genes show improvements indicating that the intron-exon structures in the curated genes have been corrected by manual curation (Figure 2).

### **Training and curation strategy**

#### **Curation Strategy: Harnessing the Crowd**

A single individual or computer system is currently unable to fully curate a genome with precise biological fidelity. A growing number of genome sequencing projects are the combined efforts of global consortia, (for example parasitoid wasps (37), centipede (38) and bed bug (39)). This indicates that a centralized model of genome annotation, the design used in earlier sequencing projects for the fruit fly (40) and the human genome (41), is giving way to a global and

collaborative communities to generate high quality genome annotation. Beyond the problem of scale, curators require insights from others with expertise in specific gene families, which makes the process of curation inherently collaborative. Mobilizing groups of researchers to focus on these specific and manageable areas is more likely to distill the most pertinent and valuable knowledge from genome analysis.

We formed a team of 36 curators by enlisting collaborators primarily distributed across seven academic institutions (Indian River State College, Cornell University / Boyce Thompson Institute, Kansas State University, University of Cincinnati, University of Florida, University of Otago and University of Tennessee Health Center), including undergraduate and graduate students, postdocs, staff researchers, and faculty. Training was provided at in-person workshop sessions and online tutorials. Training sessions were structured to provide a review of principles associated with identification of specific genes by comparison to orthologs, gene structural aspects, how to search the genome for targets of interest, and using the combination of gene predictions and additional evidence tracks (Supplementary Table 3) to correct gene models with the Apollo platform. Lastly, secondary analyses, such as BLAST comparison to other insects and phylogenetic analyses, were used to examine specific genes or gene sets.

After initial in-person training, we continued to closely follow the team's progress via bi-weekly video conferences, collaborative documents on Google Drive and an online project management website (basecamp.com). The online forum facilitated discussions outside of the video conferences and allowed annotators to interact at their convenience. We established standard operating procedures for minimum evidence required for annotating a gene model, gene naming conventions and quality control checks (See Methods and Supplementary Table 5). This was valuable in ensuring that annotators followed consistent guidelines throughout annotation. Curation of a gene family by an annotator was followed by a presentation to summarize their results during the video conference and peer-review.

The undergraduate annotators at each site were mentored by a local faculty who coordinated separate in-person meetings. A few experienced annotators from the i5k community also participated in this curation effort. The video conferences facilitated a healthy discussion about curation methods among different groups of annotators. The combination of video conferences and online forum provided curators with the tools required to efficiently share data and information as they worked in teams across institutional boundaries. Moreover, expert community curators volunteered their time to multiple projects and this model allowed them to contribute according to their availability.

Manual curation was incorporated into a bioinformatics class in 2016 by one of the authors (Benoit) at the University of Cincinnati. Specifically, one week in the class was utilized for genome annotation through the *D. citri* Apollo instance for twenty-two students. This focused on how RNAseq datasets contribute to gene prediction and why these models need to be corrected manually. Each student was responsible for correcting two gene models. As part of this course, students were given an assessment test before and after the class. Three questions (Supplementary Table 6) of this test focused on gene prediction and genome annotation, which

were the major educational focus of the genome annotation week. The average score on these questions before the class was only ~38%, which improved to ~90% at the completion of this course.

Two of the participating institutions (Benoit, University of Cincinnati; D'Elia, Indian River State College) integrated the *D. citri* genome annotation into senior capstone courses. Students that participated in capstone projects at UC covered eight topic areas; aquaporins, acidic amino acid transporters, glycosphingolipid metabolism, glycolysis, histone binding, vitamin metabolism, vitamin transport, and Hox genes. With the addition of the capstone course, a total of 28 UC students participated in the manual curation process. Six students participated in capstone projects at IRSC during which they annotated 15 gene families. In total, 25 students directly participated in the manual curation process, contributing to 39 gene reports. This strategy reinforces the use of undergraduate students in gene annotation, as students show an increase in learning and produce scientific reports which contribute to peer-reviewed publications.

### **Immune pathway in *D. citri***

Identification of the pathogen-induced immune components in ACP is critical for understanding and influencing the interaction between *D. citri* and CLAs. The repertoire of immune genes is known to be a very diverse functional group and includes proteins that recognize infectious agents and initiate a signal, members of signal transduction pathways that relay the message to the nucleus, and the genes that are transcribed in response to infection. Below, we have briefly summarized our findings from manual structural and functional curation of immunity-related genes. Additional details are provided in gene reports for specific gene families and in the supplementary notes 1- 31.

### **Pathogen recognition molecules**

Recognition of infectious agents, the first step in immune defense, relies on a variety of cellular receptors including C-type lectins, Galectins, fibrinogen related proteins (FREPs), Peptidoglycan recognition proteins (PGRP), Beta-1, 3-Glucan recognition proteins ( $\beta$ GRP) and thioester-containing proteins (TEP). These proteins recognize pathogen associated molecular patterns (PAMPs) which are associated with microbial cells (42). Most of these recognition molecules are widely conserved in insects, although the copy number often varies.

#### **C-type Lectins**

Ten C-type lectins (CTLs) carbohydrate binding receptors (43,44), were identified in *D. citri* (Supplementary Note 1). These include three oxidized low density lipoprotein receptor genes and genes encoding C-type Lectin 3, C-type Lectin 5, C-type Lectin 8, E selectin, Perlucin, Agglutinin subunit alpha and an selectin-like osteoblast derived protein. The number of CTLs in *D. citri* is comparable to the number found in bed bugs (11), pea aphids (6) and honeybees (10).

#### **Galactoside-Binding Lectins**



A total of three Galactoside-Binding Lectins (galectins) and one partial galectin were identified and manually curated within the ACP genome (Supplementary Note 2). Galectins bind  $\beta$ -galactose with their structurally similar carbohydrate-recognition domains (45), which can function alone or in clusters creating a  $\beta$ -sandwich structure without  $\text{Ca}^{2+}$  binding sites (46–48). Although we found more galectins in ACP than has been previously reported in the pea aphid (2 genes, (49)) and bed bug (1 gene, (39)), we did not observe a substantial lineage specific expansion as seen in dipterans.

#### Fibrinogen related proteins

Similar to other hemipterans (pea aphid and bed bug), few fibrinogen related proteins (FREPs) have been identified in ACP. Three complete FREPs were manually annotated (Scabrous, Angiopoietin and Tenascin) although partial un-annotatable FREP gene models were also detected. Like several of the other recognition molecule classes, the FREPs appear to have expanded in mosquitoes (50)(Supplementary Note 3). The suggestion that this expansion is related to blood feeding is consistent with the apparent absence in ACP of ficolin, tachylectins and aslectin, which are likely involved in detecting blood-borne parasites (50).

#### PGRP and $\beta$ GRP

We identified one PGRP gene in *D. citri*. Insects have two classes of PGRPs: large (L) and small (S) (51). PGRP-L proteins recognize Gram-negative bacteria and activate the Imd pathway. PGRP-S genes interact with  $\beta$ GRPs such as GGBP to recognize components of Gram-positive bacteria and then activate the Toll pathway. Based on sequence similarity to other insect proteins, the *D. citri* PGRP protein seems to belong to the S class. We did not find any GGBP genes in *D. citri* (Supplementary Note 4). This is somewhat surprising, since GGBPs have been found in several hemipterans including pea aphids (49), bed bugs (39) and brown planthoppers (52).

#### Thioester containing proteins

Only two Thioester containing proteins (TEP) were identified in ACP (Supplementary Note 5), which is comparable to the number found in *Acyrthosiphon pisum* and *Nasonia vitripennis*. TEPs are members of an ancient protein family that includes vertebrate C complement and alpha-2-macroglobulin proteins (53). Insect TEPs seem to play a similar role to their vertebrate homologs, binding to invaders such as parasites or microbes, marking them for degradation, and they are upregulated by the JAK/STAT pathway during innate immune response (54).

### Signaling cascades associated with pathogenesis

Once a potential infection has been detected, a cellular response is initiated by signaling cascade. Typically, Gram-positive bacteria and fungi cause activation of the Toll pathway, while the Imd pathway responds to Gram-negative bacteria (55,56). The JAK/STAT pathway plays a role in several immune functions, including antiviral defense (57). These signaling pathways are highly conserved, so the discovery that pea aphids (49), are missing many components of the Imd pathway was quite surprising. More recently, missing Imd pathway genes have been

reported in several other hemipterans(58–61), leading to speculation that association with Gram-negative endosymbionts may have favored the loss of these genes (61,62).

### Toll Pathway

We identified four Toll receptors in *D. citri* (Supplementary Note 6). Comparison of the Toll receptors found in various insects suggests that there were six ancestral Toll receptors: *Toll-1*, *Toll-6*, *Toll-2/7*, *Toll-8*, *Toll-9* and *Toll-10* (63,64). Phylogenetic analysis indicates that the *D. citri* genes are orthologs of *Toll-1*, *Toll-6*, *Toll-7* and *Toll-8*, but orthologs of *Toll-9* and *Toll-10* were not found. Pea aphids and bed bugs have Toll receptors from every class but *Toll-9* (39,49). We found orthologs of five of the six Spätzle (Spz) ligand classes, including Spz1, Spz3, Spz4, Spz5 and Spz6 (Supplementary Note 7). The lack of Spz2 is not surprising since it has only been reported in Diptera and Hymenoptera. The downstream Toll pathway components are represented by single genes in most insects (65). Consistent with this, we identified single copy orthologs of *tube*, *pelle*, *MyD88*, TRAF6, *cactus* and *dorsal* (Supplementary Notes 8-13). Taken together, our findings suggest that the Toll pathway is largely conserved in *D. citri*, as it is in other insects.

### Imd Pathway

As has been observed for several other hemipterans (49,58–61), many components of the Imd pathway appear to be missing in *D. citri* (Supplementary Note 14). We were unable to identify orthologs of Dredd, FADD, Imd, IKKG, and Relish. We did, however, find orthologs of pathway components IKKB, TAK1 and TAB, as well as FAF1/Caspar, a negative regulator of the pathway. Several Gram negative bacteria have been identified as *D. citri* symbionts, including *Wolbachia*, *Candidatus Carsonella*, *Candidatus Profftella armaturae*, and an as yet unidentified enteric bacteria closely related to *Klebsiella variicola* and *Salmonella enterica* (66–70). Given this information, it is tempting to speculate that the loss of many Imd pathway genes may, in fact, be associated with the ability of *D. citri* to acquire and harbor these Gram negative symbionts and might also be important for its ability to act as a carrier of CLAs.

### JAK/STAT Pathway

The Janus kinase/signal transducer of activators of transcription (JAK/STAT) pathway is a signaling pathway that provides direct communication between the membrane and nucleus (71). We identified genes encoding the major components of the JAK/STAT pathway, namely the orthologs of *domeless*, *hopscotch* and *marelle*/STAT92E (Supplementary Note 15). The JAK/STAT pathway is involved in many developmental processes, in addition to its role in immunity, and has been found in all sequenced insects to date, including other hemipterans.

## Response to pathogens and pathogen-associated stress

In response to infection, insect cells employ microbicidal compounds such as antimicrobial peptides (AMP), lysozymes and reactive oxygen species (ROS) to destroy invading cells and also activate tissue repair, wound healing and haematopoiesis processes. We searched the ACP genome for antimicrobial compounds (AMPs and lysozymes), the melanization-inducing

Clip-domain serine proteases (CLIP), the protective superoxide dismutases (SOD), and autophagy-related genes.

#### Antimicrobial peptides

Although more than 250 antimicrobial peptides (AMPs) have been identified in insects, we searched the ACP genome for ten classes of known AMPs without success. The AMPs investigated included attacin, cecropin, defensin, dipteracin, drosocin, drosomycin, gambicin, holotricin, metchnikowin and thaumatin. Although defensins are one of the most widely conserved, ancient groups of AMPs (72), the absence of defensin in ACP is not unprecedented as its absence has also been reported the hemipteran *A. pisum* (49). While the pea aphid is lacking defensin (as well as most other previously identified insect AMPs) it does contain six thaumatin (antifungal) homologs. Despite its presence in the closely related pea aphid, we were unable to identify thaumatin in the ACP genome. It must be noted that absence of previously identified AMPs does not necessarily suggest absence of all AMPs. AMPs are an extremely large, diverse group of molecules often defined by structure and function rather than conserved motifs, making identification through comparative sequence analysis difficult.

#### Lysozymes

Five genes encoding lysozymes were found in the *D. citri* genome (Supplementary Note 16). Lysozymes hydrolyze bacterial peptidoglycan, disrupting cell walls and causing cell lysis. Many insects produce lysozymes, particularly c-type lysozymes, and secrete them into the hemolymph following bacterial infection. C-type lysozymes that commonly defend against Gram-positive bacteria have been reported in many different insect orders including Diptera, Hemiptera, and Lepidoptera (73). Although c-type lysozymes were not found in the initial search of the ACP genome, two c-type lysozyme transcripts were found in *D. citri* MCOT v1.0 and were subsequently used to identify these genes in the ACP genome. Additionally, three i-type lysozymes were identified in the ACP assembly.

#### Superoxide dismutases

Insect hemocytes can produce a burst of reactive oxygen species (ROS) to kill pathogens (74). Since ROS are also damaging to host cells, superoxide dismutases are necessary to detoxify ROS. We found a total of four superoxide dismutase (SOD) genes in *D. citri* (Supplementary Note 17). Similar to other insects (75–77), *D. citri* contains both CuZn and Mn SODs. One of the *D. citri* genes is an Mn SOD and the other three are CuZn SODs.

#### CLIP

Eleven Clip-domain serine proteases (CLIP), from four distinct evolutionary clades (CLIPA, CLIPB, CLIPC and CLIPD), were manually annotated in the *D. citri* genome and corresponding models identified in MCOT v1.0 (Supplementary Note 18). These clades are present as multigene families in insect genomes and function in the hemolymph in innate immune responses (78). In *Drosophila*, CLIPs are involved in melanization and the activation of the Toll pathway (79).

#### Autophagy

Using the *D. citri* genome and the MCOT gene set, we identified *D. citri* orthologs of autophagy-related genes known in *Drosophila* (Supplementary Note 19). Autophagy is the regulated breakdown of unnecessary or dysfunctional components of the cell. This process is highly conserved among all animals and is critical to the regulation of cell degradation and recycling of cellular components. The main pathway is macroautophagy, where specific cytoplasmic components are isolated from the remaining cell in a double-membraned vesicle called the autophagosome (80–82). We identified 17 out of 20 autophagy-related genes (Supplementary Note 19). There is only a single autophagy-related 8 gene in ACP gene sets compared to two for the *Drosophila* gene set, but this is common for non-dipteran insects (81). Thus, as expected, psyllids have the required repertoire of autophagy-related genes to undergo macroautophagy.

In summary, there is a reduction in the number of immune recognition, signaling and response genes in *D. citri* compared to insects from diptera. The reduction in the immunity genes is also observed in other hemipteran insect genomes such as *A. pisum* (49), *P. humanus* (58), *Bactericera cockerelli* (59), *R. prolixus* (60) and *B. tabaci* (61). The reduction of immunity-related genes in these insects has been attributed to association with their endosymbionts, which sometimes complement the immunity of the insects (83,84). In addition, this reduction in immune genes may be associated with insects that feed on nutritionally poor and relatively sterile food sources, such as blood and fluid from the xylem/phloem (58,83). However, blood-feeding mosquitoes actually show an increase in immune genes and this expansion has been attributed to the likelihood of encountering pathogens in their food source. Arp et al. (22) pointed out additional inconsistencies with the diet hypothesis, including the presence of a full immune system in an insect species that develops in a sterile environment.

### RNA interference pathway in *D. citri*

The RNA interference (RNAi) pathway is a highly conserved, complex method of endogenous gene regulation and viral control mediated through short interfering RNAs (siRNAs), microRNAs (miRNAs), and piwiRNAs (piRNAs). While all of these small RNA molecules function to modulate or silence gene expression, the method of gene silencing and the biogenesis differs (85). In *Drosophila*, it appears that genes in the RNAi machinery have subfunctionalized to have roles in specific small RNA silencing pathways (86–92). While the RNAi machinery genes have been shown to be conserved across major taxa, functional studies in insects have been limited to a handful of diptera. Investigating the complement of RNAi genes in *D. citri* may provide insight into the role that RNAi has on the immune response of phloem-feeding insects and could aid in better use of RNAi as a tool for pest management (93–95).

#### Core machinery

Class II (Drosha type) and class III (Dicer type) RNase III enzymes play an essential role in the biogenesis of small RNA molecules with Drosha and Dicer1 functioning to produce miRNAs and Dicer2 functioning to produce siRNAs (96–98). Our analysis of the *D. citri* genome revealed four possible loci with identity to insect Dicer proteins (Supplementary Note 20). However, further

analysis of the MCOT transcriptome suggests that *D. citri* contains only one gene orthologous to *Dicer1* (MCOT05108.0.CO) and one gene orthologous to *Dicer2* (MCOT13562.0.CO). The remaining two loci are likely the result of genome fragmentation and misassemblies. BLAST analysis of Droscha also identified multiple loci with homology to other insect Droscha proteins. MCOT transcriptome analysis indicates at least one or possibly two *droscha* homologs are present in *D. citri* (Supplementary Note 21).

dsRNA binding proteins act in concert with RNase III-type enzymes to bind and process precursor dsRNA molecules into small effector molecules (88,92,99,100). In some cases, these dsRNA binding proteins also function to load small RNA molecules into the RISC (101–103). In *Drosophila*, Pasha partners with Droscha (100), Loquacious (Loqs) partners with Dicer1 (92,99), and R2D2 partners with Dicer2 (88). In the *D. citri* genome, we identified two *pasha* homologs (Supplementary Note 22) and two *loqs* homologs but were initially unable to identify a true *r2d2* ortholog (Supplementary Note 23) in the genome. The apparent absence of *r2d2* was consistent with previous reports (22,104). However, a search of the MCOT (MCOT18647.0.CO) transcriptome identified a gene with similarity to R2D2 orthologs from bed bug, *Tribolium castaneum* and mosquitoes. While *r2d2* is likely to be present in the *D. citri* genome, it is not annotatable given the limitations of the current assembly. Alternatively, if *r2d2* is missing from *D. citri*, it is possible that one of the *Loqs* proteins identified functions in the RNAi pathway (104), as *Loqs* has been shown to associate with Dicer2 in both *Drosophila* and *Aedes aegypti* (105,106).

Argonaute (AGO) proteins present small RNA guide molecules to their complementary targets through silencing complexes and provide the ‘Slicer’ catalytic activity that is required for mRNA cleavage in some RNA silencing pathways (107–110). In *Drosophila* AGO1 is involved in the miRNA pathway, AGO2 is involved in silencing by siRNAs (86,91) and AGO3, PIWI and Aubergine (Aub) function in the piRNA pathway (87,89,111). In the *D. citri* genome, we have identified 4 AGO genes, AGO1, AGO2, AGO3 and one gene corresponding to the PIW/Aub class of proteins (Supplementary Note 24).

#### Auxiliary (RISC and other) factors

A subset of other genes known to be involved in the function or regulation of the RNA-induced silencing complex (RISC) in *Drosophila* and other organisms have been identified and annotated in the *D. citri* genome. The genes identified include two *Tudor Staphylococcal Nucleases* (TSN, Supplementary Note 25), one *vasa-intronic gene* (*vig-1*, Supplementary Note 25), one *armitage* (*Arm*) gene (Supplementary Note 27), and one *Fragile X Mental Retardation 1* (*FMR1*, Supplementary Note 28) gene. Additionally, several more genes known to be involved in the biogenesis or function of small RNA molecules have been identified. These include two *spindle-E* genes (Supplementary Note 29), one *Rm62* gene (Supplementary Note 30) and one *Ran* gene (Supplementary Note 31).

In summary, the *D. citri* genome has a full complement of RNAi machinery genes. Duplications are more frequent in genes that have previously been associated with the miRNA pathway (*droscha*, *pasha* and *loqs*) as opposed to the RNAi or piRNA pathways. This is an interesting



finding as the same result was found upon analysis of the pea aphid genome (112) but was not seen in the whitefly genome (61).

## Building the foundation for P450/Halloween genes targeting to Reduce Insect Pests

Cytochrome P450s (CYPs) in eukaryotes are heme containing membrane bound enzymes that activate molecular oxygen via a mechanism involving a thiolate ligand to the heme iron. Usually this requires an electron donor protein, the NADPH cytochrome P450 reductase, in the ER or ferredoxin and ferredoxin reductase in the mitochondria (113,114). Insects have four deep branching clades on phylogenetic trees and this represents some losses during evolution as up to 11 clades are found in other animals. These are termed CYP2, CYP3, CYP4 and mitochondrial clans in P450 nomenclature (115). Most species have a tendency to expand P450s in one or more clans via tandem duplications. One interpretation of these P450 “blooms” is diversification to handle many related compounds from the environment that may be toxic or potential carbon sources.

*D. citri* in its current assembly has 60 P450 genes that are identified and named as distinct P450s. There are also numerous fragments named as partials. *Diaphorina* has a P450 bloom in the clusters CYP3172, CYP3174, CYP3175, CYP3176, CYP3178 in the CYP4 clan. There is another in the CYP3167 family in the CYP2 clan and a third that includes CYP6KA, CYP6KC and CYP6KD in the CYP3 clan. *D. citri* has 3 CYP4G genes.

CYP2 and mito clans have many 1:1 orthologs but these are rare in the CYP3 and CYP4 clan (Figure 3). One exception in the CYP3 clan is CYP3087A1 and two neighbors CYP6DB1 and CYP6KB1 as they may be orthologs and probably should be in the same family. *R. prolixus* and *A. pisum* CYP3 clan genes have undergone gene blooms that were not found in *D. citri*. This may be interpreted as the common ancestor having few CYP3 clan P450s. The cluster at arc A (Figure 3) on the tree consisting of CYP6KB1, CYP6DB1 and CYP3087A1 may be evidence that these three are orthologs and should be in the same family. The large aphid clade of CYP6CY at arc C has no members from *Rhodnius* or *Diaphorina*, so it seems to be aphid specific. At arc B (Figure 3) the CYP395 family has four subfamilies C, D, E, F. There are many CYP395 genes in other hemiptera species, including *C. lectularius* (bedbug CYP395A,B), *Apolygus lucorum* (Hemiptera, a Mirid bug, CYP395G, H, J, K, L, M)), and *Cyrtorhinus lividipennis* (Hemiptera, green mirid bug, CYP395H, J, N). The fact they have been placed in different subfamilies suggests they are diverging from their common ancestor. The CYP3084 family with subfamilies A, B, C, D is only found in *Rhodnius* so far (116). The families CYP3088, CYP3089, CYP3090 and CYP3091 are also Hemiptera specific with some members in the same species noted earlier. The number of P450s varies with arthropod species from a low of 25 in the mite *Aculops lyoperscii* and 36 in *Pediculus humanus* (body louse) to over 200 in *Ixodes scapularis* (black-legged tick) and up to 158 in some mosquitos (117).

## CONCLUSION

We report the first draft assembly for the *D. citri* genome and the corresponding official gene set (OGS v1.0) which includes 530 manually curated genes and about 20,000 genes predicted by the NCBI Eukaryotic Genome Annotation Pipeline (NCBI v100). The community curation effort involved undergraduate students at multiple locations who were trained, individually or in a class setting, in gene curation as a part of this initiative. These students were supported by contributions from expert annotators in the insect genomics community. The major advantage of having both expert curators and undergraduate students work together in the annotation project was the training, exchange of ideas and community building. We also present standard operating procedures that can be used to guide and coordinate annotation by large virtual teams. This community annotation process will be continued to improve structural and functional characterization of the ACP genome as new versions are generated. The MCOT transcriptome reported in this paper and also available at [citrusgreening.org](http://citrusgreening.org) offers a genome-independent and comprehensive representation of the gene repertoire of *D. citri* that was used to improve the genome annotation. The MCOT transcriptome allowed us to identify lineage specific-gene models in *D. citri* and curate them. This gene set will support efforts in other hemipteran species that transmit bacterial pathogens.

In summary, we curated and described genes related to immunity and the RNAi pathway in addition to the cytochrome P450 genes. We report blooms in P450 genes in the CYP4, CYP2 and CYP3 clans which may be an evolutionary response to environmental stresses. Other important gene families that were curated as a part of the official gene set include aquaporins, cathepsins, cuticle and secretory proteins. We found the number of immunity-related genes to be reduced, even after direct targeting for improvement, in the *D. citri* genome similar to pea aphid and whitefly, which may reflect the association with microbial symbionts that have coevolved in both insects and the consumption of relatively sterile plant derived fluids. The genomic resources from this project will provide critical information underlying ACP biology that can be used to improve control of this pest.

## MATERIAL AND METHODS

### DNA extraction and library preparation

High-molecular weight DNA was extracted using the BioRad AquaPure Genomic DNA isolation kit from fresh intact *D. citri* collected from a citrus grove in Ft. Pierce, FL and reared at the USDA, ARS, U.S. Horticultural Research Laboratory, Ft. Pierce, FL. To generate PacBio libraries, DNA was sheared using a Covaris g-Tube and SMRT-bell library was prepared using the 10Kb protocol (PacBio DNA template prep kit 2.0; 3-10Kb), cat #001-540-835.

### Genome sequencing and assembly

Samples were prepared for sequencing using the TruSeq DNA library preparation kits for paired-end as well as long-insert mate-pair libraries. All were sequenced on the Illumina HiSeq2000 using 100bp or longer reads. Seven libraries were sequenced, with inserts ranging from “short” (ca. 275bp) to 10Kb. These are available in NCBI SRA and included 99.7 million

paired-end reads (NCBI SRA:SRX057205), 35.1 million 2kb mate-pair reads (NCBI SRA: SRX057204), 30 million 5kb mate-pair reads (NCBI SRA: SRX058250) and 30 million 10kb mate-pair reads (NCBI SRA: SRX216330). A second round of DNA sequencing was done with PacBio at 12X coverage (NCBI SRA: SRX218985) for scaffolding the Diaci1.0 Illumina assembly to create the Diaci1.1 version of the *D. citri* genome. Thirty-nine SMRTcells of the library were sequenced, all with 2x45 minute movies. A total of 2,750,690 post-filter reads were generated, with an average of 70,530 reads per SMRTcell. The post-filter mean read length was 2,504 bp with an error rate of 15%.

Velvet (118) was used with kmer 59 for generating the original assembly. PacBio long reads were mapped to the draft assembly using blasr (119) with the following parameters: -minMatch 8 -minPctIdentity 70 -bestn 5 -nCandidates 30 -maxScore -500 -nproc 8 -noSplitSubreads. These alignments were parsed using PBJelly (120) with default parameters.

### Maker and NCBI annotation

The maker control files are included in supplementary data. BLAST 2.2.27+ was used with augustus version 2.5.5 (121) and exonerate version 2.2.0 (122). Only contigs longer than 10kb were selected for annotation with psyllid specific transcriptome sequences from Reese et al. (20). Details of ACP genome annotation by the NCBI Eukaryotic Genome Annotation pipeline are available at [https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/Diaphorina\\_citri/100/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Diaphorina_citri/100/).

### Annotation edit distance

We mapped transcripts from MCOT v1.0 to the genome using GMAP (123). Out of 30,562 transcripts, 19,744 were mapped with at least 90% query coverage and 90% identity.

A genome guided transcriptome was generated to validate all annotation sets using all available RNAseq data (Supplementary Table 4) and insect proteins (NCBI taxonomy:"Hexapoda [6960]") from Swiss-Prot were used as sources of evidence. 622 million paired-end reads were mapped to NCBI-Diaci1.1 assembly using hisat2 (124) with a mapping rate of 81.78%. Mapped files were sorted with samtools rocksort. After sorting, a genome-guided transcriptome assembly was performed using StringTie (125) and the resulting transcriptome contained 210,890 transcripts (N50: 1,691bp).

Annotation edit distance was calculated for all gene models from NCBI v100, mapped MCOT, Maker v1.1 and curated gene sets. The Maker genome annotation (v2.31.8) pipeline was used for calculating AED (28,36).

### MCOT transcriptome assembly

The MCOT v1.0 set was generated with the MCOT pipeline (26). Maker v1.1 gene models and RNAseq data from adult, nymph and egg tissue (20) were used to generate a genome-based transcriptome assembly using Cufflinks (29). *Denovo* transcriptome assemblies of the adult, nymph and egg RNAseq data were performed with Trinity (30) and Oases (31). These are available at <ftp://ftp.citrusgreening.org/annotation/MCOT/>. Transcripts from Maker (27,28),

Cufflinks (29), Trinity (30) and Oases (31) were translated to proteins with Transdecoder version 2.0.1(126) and unique proteins were kept.

Reads were assembled with Trinity in two runs, one used reads as single end reads, and the other used them as paired end reads. Reads were trimmed based on fastq quality score (the --trimmomatic option was enabled and run under the default setting of Trinity). The transcripts of both runs were combined to make the final Trinity assembly.

Velvet-Oases (31) assemblies were performed for trimmed reads (trimmed in Trinity run, the read quality control step) from the egg, nymph and adult separately, with kmer length of 23, 25, 27 and 29 as single end reads, and kmer length 25 as paired end reads. The outputs from kmer 15 and kmer 27 were combined using the Oases merge function (--long and -min\_trans\_lgth 200) to generate the final assembly.

Reads from the egg, nymph and adult were first aligned to the genome with Tophat (127) with the insert length parameter based on each library (53, 24 and 90). The parameter --read-realign-edit-dist were set to 0 and -r 90 to ensure better alignment results. Gene models were generated by Cufflinks (128) with default settings, with the -frag-bias-correct and -multi-read-correct function (-b, -u) enabled to give the most accurate gene models.

Furthermore, protein sequences from each program (Maker, Oases, Trinity, Cufflinks) were compared with BLASTP, with a special scoring matrix (matching score of non-identical amino acids setting to -100 of the BLOSUM62 matrix), and compared with proteins from other arthropod species by normal BLASTP alignment. The best protein models from each source were selected to make the final MCOT v1.0 protein set, and the corresponding transcript set. MCOT v1.0 set has 30,562 genes and is available at <ftp://ftp.citrusgreening.org/annotation/MCOT/>.

MCOT v1.0 was analyzed using Interproscan 5 (129) and InterPro database, which contains predictive information about protein function with the options -goterms to get GO terms, -iprlookup to switch on look-up of corresponding InterPro annotation and -pa option to switch on lookup of corresponding Pathway annotation. We also performed BLAST analysis on the MCOT set using the NCBI nr, Swiss-Prot and the TrEMBL protein databases. These BLAST and InterproScan results were used as input for AHRD (Automated assignment of Human Readable Descriptions) (33), to assign descriptions, Pfam domains and Gene Ontology terms. This functional annotation was performed using a filter of bit score more than 50 and e-value less than e-10.

#### Community curation for structural and functional annotation of genes

Gene lists and orthologs were made available via the project management website (<https://basecamp.com/2923904/projects/9184795>) so that the annotators could volunteer to analyze genes of interest. We used ImmunoDB (35) as a primary source of immune gene orthologs, which provided expert-curated immunity genes for *Aedes aegypti*, *Anopheles gambiae*, *Drosophila melanogaster* and *Culex quinquefasciatus*. Other closely related

organisms used as sources of annotated orthologs include bed bug (*Cimex lectularius*) (39), pea aphid (*Acyrtosiphon pisum*) (62) and milkweed bug (130). All communication was done via emails and the project management website, which was also used to store presentations, gene reports, meeting minutes, working documents and data files.

Candidate gene models were identified on the *D. citri* genome by using orthologous proteins as query in Apollo blat and i5k BLAST. The NCBI conserved domains database (131) was used to identify the conserved domains in the orthologs and candidate genes. Multiple sequence alignments were generated using MUSCLE (132), tcofee (133) and clustal (134) to compare the ACP gene model to the query gene set. The final model was refined in Apollo using homology, RNAseq and proteomics evidence tracks. MEGA7 (135) was used to construct phylogenetic trees. Please see individual gene reports in supplementary notes 1-39 for detailed methods. When available, published literature was used to putatively assign molecular functions, participation in biological processes, and cellular localization for an annotated gene, associating term identifiers from the Gene Ontology Database (136) and PubMed identifiers from NCBI. Curated genes were assigned names and descriptions based on the function and domain structures available in published literature or NCBI.

Another strategy for identification of candidate genes involved preprocessing the query before searching the ACP genome on Apollo. A set of representative genes was BLASTN searched against a database of all contigs in the ACP genome. These contigs were BLASTX searched against all the database of named insect genes. Analysis of the match helped to identify missing exons and extend partial exons to achieve the best possible manually curated version of the ACP gene. Once gene models were reconstructed from the genomic DNA, they were located on the *D. citri* genome using the i5k BLAST server. The models were then manually edited based on evidence tracks to produce the final gene model.

We performed multiple cycles of internal review of curated gene models to identify errors and suggest improvements to annotators. Annotations were ranked a scale of A-D (Supplementary Table 5) during review depending on completeness and support from evidence tracks. Standard gene naming conventions were defined and agreed upon to ensure consistency. The curated gene models were exported from Apollo and validated using the i5k quality control pipeline ([https://github.com/NAL-i5K/I5KNAL\\_OGS/wiki/QC-phase](https://github.com/NAL-i5K/I5KNAL_OGS/wiki/QC-phase)) which identifies intra-model, inter-model and single-feature errors.

## REFERENCES

1. Halbert, S. E. and Núñez, C. A. (2004) *Florida Entomol.*, **87**, 401–402, DISTRIBUTION OF THE ASIAN CITRUS PSYLLID, DIAPHORINA CITRI KUWAYAMA (RHYNCHOTA: PSYLLIDAE) IN THE CARIBBEAN BASIN.
2. Halbert, S. E. and Manjunath, K. L. (2004) *Florida Entomol.*, **87**, 330–353, ASIAN CITRUS PSYLLIDS (STERNORRHYNCHA: PSYLLIDAE) AND GREENING DISEASE OF CITRUS: A LITERATURE REVIEW AND ASSESSMENT OF RISK IN FLORIDA.
3. Boykin, L. M., De Barro, P., Hall, D. G., et al. (2012) *Bull. Entomol. Res.*, **102**, 573–582,



- Overview of worldwide diversity of *Diaphorina citri* Kuwayama mitochondrial cytochrome oxidase 1 haplotypes: two Old World lineages and a New World invasion.
4. French, J. V., Kahlke, C. J. and Da Graça, J. V. (2001) *Subtrop. Plant Sci.*, **53**, 14–15, First record of the Asian citrus psylla, *Diaphorina citri* Kuwayama (Homoptera: Psyllidae) in Texas.
5. Pluke, R. W. H., Qureshi, J. A. and Stansly, P. A. (2008) *Florida Entomol.*, **91**, 36–42, CITRUS FLUSHING PATTERNS, DIAPHORINA CITRI (HEMIPTERA: PSYLLIDAE) POPULATIONS AND PARASITISM BY TAMARIXIA RADIATA (HYMENOPTERA: EULOPHIDAE) IN PUERTO RICO.
6. Tsai, J. H. and Liu, Y. H. (2000) *J. Econ. Entomol.*, **93**, Biology of *Diaphorina citri* (Homoptera: Psyllidae) on Four Host Plants.
7. Teixeira, D. do C., Saillard, C., Eveillard, S., et al. (2005) *Int. J. Syst. Evol. Microbiol.*, **55**, 1857–62, “Candidatus *Liberibacter americanus*”, associated with citrus huanglongbing (greening disease) in São Paulo State, Brazil.
8. Capoor, S. P., Rao, D. G., Viswanath, S. M., et al. (1967) *Indian J. Agric. Sci.*, **37**, 572–575, *Diaphorina citri* Kuway., a vector of the greening disease of citrus in India.
9. Bové, J. M. (2006) *J. Plant Pathol.*, **88**, 7–37, INVITED REVIEW HUANGLONGBING : A DESTRUCTIVE , NEWLY-EMERGING , CENTURY-OLD DISEASE OF CITRUS 1.
10. Manjunath, K. L., Halbert, S. E., Ramadugu, C., et al. (2008) *Phytopathology*, **98**, 387–396, Detection of ‘Candidatus *Liberibacter asiaticus*’ in *Diaphorina citri* and its importance in the management of citrus huanglongbing in Florida.
11. Leong, S. C. T., Abang, F., Beattie, A., et al. (2012) *Sci. World J.*, **2012**, Impacts of horticultural mineral oils and two insecticide practices on population fluctuation of *Diaphorina citri* and spread of huanglongbing in a citrus orchard in Sarawak.
12. Honig, L. October Crop Production Executive Summary  
[https://www.nass.usda.gov/Newsroom/Executive\\_Briefings/2016/10\\_12\\_2016.pdf](https://www.nass.usda.gov/Newsroom/Executive_Briefings/2016/10_12_2016.pdf).
13. Tiwari, S., Lewis-Rosenblum, H., Pelz-Stelinski, K., et al. (2010) *J. Econ. Entomol.*, **103**, Incidence of *Candidatus Liberibacter asiaticus* Infection in Abandoned Citrus Occurring in Proximity to Commercially Managed Groves.
14. Tabachnick, W. J. (2015) *J. Econ. Entomol.*, **108**, 839–848, *Diaphorina citri* (Hemiptera: Liviidae) vector competence for the citrus greening pathogen *Candidatus Liberibacter asiaticus*.
15. Ramsey, J. S., Johnson, R. S., Hoki, J. S., et al. (2015) *PLoS One*, **10**, e0140826, Metabolic Interplay between the Asian Citrus Psyllid and Its Proffella Symbiont: An Achilles’ Heel of the Citrus Greening Insect Vector.
16. de Andrade, E. C. and Hunter, W. B. In *RNA Interference*; InTech, 2016.
17. Marutani-Hert, M., Hunter, W. B. and Hall, D. G. (2010) *Florida Entomol.*, **93**, 519–525, Gene response to stress in the Asian citrus psyllid (Hemiptera: Psyllidae).
18. Hunter, W. B., Dowd, S. E., Katsar, C. S., et al. (2009) *Gene*, 18–29, Psyllid Biology : Expressed Genes in Adult Asian Citrus Psyllids , *Diaphorina citri* Kuwayama.
19. Hunter, W. B., Hail, D., Tipping, C., et al. In *Symposium Proceedings*; 2010.
20. Reese, J., Christenson, M. K., Leng, N., et al. (2014) *J Genomics*, **2**, 54–58, Characterization of the Asian Citrus Psyllid Transcriptome.
21. Fisher, T., Vyas, M., He, R., et al. (2014) *Pathogens*, **3**, 875–907, Comparison of Potato and Asian Citrus Psyllid Adult and Nymph Transcriptomes Identified Vector Transcripts with Potential Involvement in Circulative, Propagative *Liberibacter* Transmission.
22. Arp, A. P., Pelz-Stelinski, K. and Hunter, W. (2016) *Front. Physiol.*, **7**, 570, Annotation of the Asian citrus psyllid genome reveals a reduced innate immune system.
23. Stein, L. (2001) *Nat. Rev. Genet.*, **2**, 493–503, Genome annotation: from sequence to biology.
24. Elsik, C. G., Worley, K. C., Zhang, L., et al. (2006) *Genome Res.*, **16**, 1329–33,

- Community annotation: procedures, protocols, and supporting tools.
25. Simao, F. A., Waterhouse, R. M., Ioannidis, P., et al. (2015) *Bioinformatics*, btv351-, BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs.
26. Cao, X. and Jiang, H. (2015) *Insect Biochem. Mol. Biol.*, Integrated modeling of protein-coding genes in the *Manduca sexta* genome using RNA-Seq data from the biochemical model insect.
27. Cantarel, B. L., Korf, I., Robb, S. M. C., et al. (2008) *Genome Res.*, **18**, 188–96, MAKER: an easy-to-use annotation pipeline designed for emerging model organism genomes.
28. Yandell, M. and Ence, D. (2012) *Nat. Rev. Genet.*, **13**, 329–342, A beginner's guide to eukaryotic genome annotation.
29. Trapnell, C., Williams, B. A., Pertea, G., et al. (2010) *Nat. Biotechnol.*, **28**, 511–5, Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation.
30. Grabherr, M. G., Haas, B. J., Yassour, M., et al. (2011) *Nat. Biotechnol.*, **29**, Full-length transcriptome assembly from RNA-Seq data without a reference genome.
31. Schulz, M. H., Zerbino, D. R., Vingron, M., et al. (2012) *Bioinformatics*, **28**, 1086–1092, Oases : robust de novo RNA-seq assembly across the dynamic range of expression levels.
32. Waterhouse, R. M., Zdobnov, E. M., Tegenfeldt, F., et al. (2011) *Nucleic Acids Res.*, **39**, D283–8, OrthoDB: the hierarchical catalog of eukaryotic orthologs in 2011.
33. Schoof, H. In *Plant and Animal Genome XXIV Conference*.
34. Poelchau, M., Childers, C., Moore, G., et al. (2014) *Nucleic Acids Res.*, gku983-, The i5k Workspace@NAL--enabling genomic data access, visualization and curation of arthropod genomes.
35. Waterhouse, R. M., Kriventseva, E. V, Meister, S., et al. (2007) *Science*, **316**, 1738–43, Evolutionary dynamics of immune-related genes and pathways in disease-vector mosquitoes.
36. Eilbeck, K., Moore, B., Holt, C., et al. *Quantitative measures for the management and comparison of annotated genomes.*; 2009; Vol. 10.
37. Werren, J. H., Richards, S., Desjardins, C. A., et al. (2010) *Science*, **327**, 343–8, Functional and evolutionary insights from the genomes of three parasitoid *Nasonia* species.
38. Chipman, A. D., Ferrier, D. E. K., Brena, C., et al. (2014) *PLoS Biol.*, **12**, e1002005, The First Myriapod Genome Sequence Reveals Conservative Arthropod Gene Content and Genome Organisation in the Centipede *Strigamia maritima*.
39. Benoit, J. B., Adelman, Z. N., Reinhardt, K., et al. (2016) *Nat. Commun.*, **7**, 10165, Unique features of a global human ectoparasite identified through sequencing of the bed bug genome.
40. Adams, M. D., Celniker, S. E., Holt, R. A., et al. (2000) *Science (80-. )*, **287**, 2185–2195, The genome sequence of *Drosophila melanogaster*.
41. Venter, J. C., Adams, M. D., Myers, E. W., et al. (2001) *Science (80-. )*, **291**, 1304–1351, The Sequence of the Human Genome.
42. Christophides, G. K., Vlachou, D. and Kafatos, F. C. (2004) *Immunol Rev*, **198**, 127–148, Comparative and functional genomics of the innate immune system in the malaria vector *Anopheles gambiae*.
43. Dodd, R. B. and Drickamer, K. (2001) *Glycobiology*, **11**, 71R--79R, Lectin-like proteins in model organisms: implications for evolution of carbohydrate-binding activity.
44. Cambi, A. and Figdor, C. G. (2003) *Curr. Opin. Cell Biol.*, **15**, 539–546, Dual function of C-type lectin-like receptors in the immune system.
45. Cummings, R. D. and Liu, F. T. *Galectins Essentials of Glycobiology*. 2. Vol. Chapter 33

## 2009.

46. Leffler, H., Carlsson, S., Hedlund, M., et al. (2002) *Glycoconj. J.*, **19**, 433–440, Introduction to galectins.
47. Mitchell, D. A., Fadden, A. J. and Drickamer, K. (2001) *J. Biol. Chem.*, **276**, 28939–28945, A novel mechanism of carbohydrate recognition by the C-type lectins DC-SIGN and DC-SIGNR Subunit organization and binding to multivalent ligands.
48. Wang, L., Wang, L., Yang, J., et al. (2012) *Dev. Comp. Immunol.*, **36**, 591–601, A multi-CRD C-type lectin with broad recognition spectrum and cellular adhesion from *Argopecten irradians*.
49. Gerardo, N. M., Altincicek, B., Anselme, C., et al. (2010) *Genome Biol.*, **11**, R21, Immunity and other defenses in pea aphids, *Acyrtosiphon pisum*.
50. Wang, X., Zhao, Q. and Christensen, B. M. (2005) *BMC Genomics*, **6**, 1, Identification and characterization of the fibrinogen-like domain of fibrinogen-related proteins in the mosquito, *Anopheles gambiae*, and the fruitfly, *Drosophila melanogaster*, genomes.
51. Dziarski, R. and Gupta, D. (2006) *Genome Biol.*, **7**, 1, The peptidoglycan recognition proteins (PGRPs).
52. Bao, Y.-Y., Qu, L.-Y., Zhao, D., et al. (2013) *BMC Genomics*, **14**, 1, The genome-and transcriptome-wide analysis of innate immunity in the brown planthopper, *Nilaparvata lugens*.
53. Blandin, S. and Levashina, E. A. (2004) *Mol. Immunol.*, **40**, 903–908, Thioester-containing proteins and insect immunity.
54. Agaisse, H. and Perrimon, N. (2004) *Immunol. Rev.*, **198**, 72–82, The roles of JAK/STAT signaling in *Drosophila* immune responses.
55. Lindsay, S. A. and Wasserman, S. A. (2014) *Dev. Comp. Immunol.*, **42**, 16–24, Conventional and non-conventional *Drosophila* Toll signaling.
56. Myllymäki, H., Valanne, S. and Rämet, M. (2014) *J. Immunol.*, **192**, 3455–3462, The *Drosophila* imd signaling pathway.
57. Myllymäki, H. and Rämet, M. (2014) *Scand. J. Immunol.*, **79**, 377–385, JAK/STAT pathway in *Drosophila* immunity.
58. Kim, J. H., Min, J. S., Kang, J. S., et al. (2011) *Insect Biochem. Mol. Biol.*, **41**, 332–339, Comparison of the humoral and cellular immune responses between body and head lice following bacterial challenge.
59. Nachappa, P., Levy, J. and Tamborindéguy, C. (2012) *Mol. Genet. genomics*, **287**, 803–817, Transcriptome analyses of *Bactericera cockerelli* adults in response to *Candidatus Liberibacter solanacearum* infection.
60. Mesquita, R. D., Vionette-Amaral, R. J., Lowenberger, C., et al. (2015) *Proc. Natl. Acad. Sci.*, **112**, 14936–14941, Genome of *Rhodnius prolixus*, an insect vector of Chagas disease, reveals unique adaptations to hematophagy and parasite infection.
61. Chen, W., Hasegawa, D. K., Kaur, N., et al. (2016) *BMC Biol.*, **14**, 110, The draft genome of whitefly *Bemisia tabaci* MEAM1, a global crop pest, provides novel insights into virus transmission, host adaptation, and insecticide resistance.
62. Consortium, T. I. A. G. (2010) *PLoS Biol.*, **8**, e1000313, Genome sequence of the pea aphid *Acyrtosiphon pisum*.
63. Evans, J. D., Aronstein, K., Chen, Y. P., et al. (2006) *Insect Mol. Biol.*, **15**, 645–656, Immune pathways and defence mechanisms in honey bees *Apis mellifera*.
64. Benton, M. A., Pechmann, M., Frey, N., et al. (2016) *Curr. Biol.*, Toll Genes Have an Ancestral Role in Axis Elongation.
65. Viljakainen, L. (2015) *Brief. Funct. Genomics*, elv002-, Evolutionary genetics of insect innate immunity.
66. Nakabachi, A., Yamashita, A., Toh, H., et al. (2006) *Science (80-. )*, **314**, 267, The 160-kilobase genome of the bacterial endosymbiont *Carsonella*.

67. Hilgenboecker, K., Hammerstein, P., Schlattmann, P., et al. (2008) *FEMS Microbiol Lett*, **281**, 215–220, How many species are infected with Wolbachia?--A statistical analysis of current data.
68. Saha, S., Hunter, W. B., Reese, J., et al. (2012) *PLoS One*, **7**, e50067, Survey of Endosymbionts in the *Diaphorina citri* Metagenome and Assembly of a Wolbachia wDi Draft Genome.
69. Sloan, D. B. and Moran, N. A. (2012) *Mol Biol Evol*, **29**, 3781–3792, Genome reduction and co-evolution between the primary and secondary bacterial symbionts of psyllids.
70. Nakabachi, A., Ueoka, R., Oshima, K., et al. *Defensive Bacteriome Symbiont with a Drastically Reduced Genome*; 2013; Vol. 23.
71. O'Shea, J. J., Schwartz, D. M., Villarino, A. V., et al. (2015) *Annu. Rev. Med.*, **66**, 311–328, The JAK-STAT pathway: impact on human disease and therapeutic intervention\*.
72. Zhang, L. and Gallo, R. L. (2016) *Curr. Biol.*, **26**, R14--R19, Antimicrobial peptides.
73. Callewaert, L. and Michiels, C. W. (2010) *J. Biosci.*, **35**, 127–160, Lysozymes in the animal kingdom.
74. Lavine, M. D. and Strand, M. R. (2002) *Insect Biochem. Mol. Biol.*, **32**, 1295–1309, Insect hemocytes and their role in immunity.
75. Bordo, D., Djinoovic, K. and Bolognesi, M. (1994) *J. Mol. Biol.*, **238**, 366–386, Conserved patterns in the Cu, Zn superoxide dismutase family.
76. Parker, J. D., Parker, K. M., Sohal, B. H., et al. (2004) *Proc. Natl. Acad. Sci. U. S. A.*, **101**, 3486–3489, Decreased expression of Cu--Zn superoxide dismutase 1 in ants with extreme lifespan.
77. Colinet, D., Cazes, D., Belghazi, M., et al. (2011) *J. Biol. Chem.*, **286**, 40110–40121, Extracellular superoxide dismutase in insects characterization, function, and interspecific variation in parasitoid wasp venom.
78. Kanost, M. R. and Jiang, H. (2015) *Curr. Opin. insect Sci.*, **11**, 47–55, Clip-domain serine proteases as immune factors in insect hemolymph.
79. Veillard, F., Troxler, L. and Reichhart, J.-M. (2016) *Biochimie*, **122**, 255–269, *Drosophila melanogaster* clip-domain serine proteases: Structure, function and regulation.
80. Chang, Y.-Y. and Neufeld, T. P. (2010) *FEBS Lett.*, **584**, 1342–1349, Autophagy takes flight in *Drosophila*.
81. Malagoli, D., Abdalla, F. C., Cao, Y., et al. (2010) *Autophagy*, **6**, 575–588, Autophagy and its physiological relevance in arthropods: current knowledge and perspectives.
82. Zirin, J. and Perrimon, N. In *Seminars in immunopathology*; 2010; Vol. 32, pp. 363–372.
83. Altincicek, B., Gross, J. and Vilcinskis, A. (2008) *Insect Mol. Biol.*, **17**, 711–716, Wounding-mediated gene expression and accelerated viviparous reproduction of the pea aphid *Acyrtosiphon pisum*.
84. Ratzka, C., Gross, R. and Feldhaar, H. (2012) *Insects*, **3**, 553–572, Endosymbiont tolerance and control within insect hosts.
85. Ghildiyal, M. and Zamore, P. D. (2009) *Nat Rev Genet*, **10**, 94–108, Small silencing RNAs: an expanding universe.
86. Hammond, S. M., Bernstein, E., Beach, D., et al. (2000) *Nature*, **404**, 293–296, An RNA-directed nuclease mediates post-transcriptional gene silencing in *Drosophila* cells.
87. Pal-Bhadra, M., Bhadra, U. and Birchler, J. A. (2002) *Mol Cell Biol*, **9**, 315–327, RNAi related mechanisms affect both transcriptional and posttranscriptional transgene silencing in *Drosophila*.
88. Liu, Q., Rand, T. A., Kalidas, S., et al. (2003) *Science (80-. )*, **301**, 1921–1925, R2D2, a bridge between the initiation and effector steps of the *Drosophila* RNAi pathway.
89. Aravin, A. A., Klenov, M. S., Vagin, V. V., et al. (2004) *Mol Cell Biol*, **24**, 6742–6750, Dissection of a natural RNA silencing process in the *Drosophila melanogaster* germ line.
90. Lee, Y. S., Nakahara, K., Pham, J. W., et al. (2004) *Cell*, **117**, 69–81, Distinct roles for



- Drosophila Dicer-1 and Dicer-2 in the siRNA/miRNA silencing pathways.
91. Okamura, K., Ishizuka, A., Siomi, H., et al. (2004) *Genes Dev*, **18**, 1655–1666, Distinct roles for Argonaute proteins in small RNA-directed RNA cleavage pathways.
92. Saito, K., Ishizuka, A., Siomi, H., et al. (2005) *PLoS Biol*, **3**, e235, Processing of pre-microRNAs by the Dicer-1-Loquacious complex in *Drosophila* cells.
93. Scott, J. G., Michel, K., Bartholomay, L., et al. (2013) *J. Insect Physiol.*, Towards the elements of successful insect RNAi.
94. Christiaens, O. and Smagghe, G. (2014) *Curr. Opin. Insect Sci.*, **6**, 15–21, The challenge of RNAi-mediated control of hemipterans.
95. Kola, V. S. R., Renuka, P., Madhav, M. S., et al. (2015) *Front. Physiol.*, **6**, 119, Key enzymes and proteins of crop insects as candidate for RNAi based gene silencing.
96. Bernstein, E., Caudy, A. A., Hammond, S. M., et al. (2001) *Nature*, **409**, 363–366, Role for a bidentate ribonuclease in the initiation step of RNA interference.
97. Knight, S. W. and Bass, B. L. (2001) *Science (80-. )*, **293**, 2269–2271, A role for the RNase III enzyme DCR-1 in RNA interference and germ line development in *Caenorhabditis elegans*.
98. Lee, Y., Ahn, C., Han, J., et al. (2003) *Nature*, **425**, 415–419, The nuclear RNase III Drosha initiates microRNA processing.
99. Leuschner, P. J., Obernosterer, G. and Martinez, J. (2005) *Curr Biol*, **15**, R603-5, MicroRNAs: Loquacious speaks out.
100. Yeom, K. H., Lee, Y., Han, J., et al. (2006) *Nucleic Acids Res*, **34**, 4622–4629, Characterization of DGCR8/Pasha, the essential cofactor for Drosha in primary miRNA processing.
101. Liu, X., Jiang, F., Kalidas, S., et al. (2006) *RNA*, **12**, 1514–1520, Dicer-2 and R2D2 coordinately bind siRNA to promote assembly of the siRISC complexes.
102. Liu, X., Park, J. K., Jiang, F., et al. (2007) *RNA*, **13**, 2324–2329, Dicer-1, but not Loquacious, is critical for assembly of miRNA-induced silencing complexes.
103. Okamura, K., Robine, N., Liu, Y., et al. (2011) *Mol Cell Biol*, **31**, 884–896, R2D2 organizes small regulatory RNA pathways in *Drosophila*.
104. Taning, C. N. T., Andrade, E. C., Hunter, W. B., et al. (2016) *Sci. Rep.*, **6**, 38082, Asian Citrus Psyllid RNAi Pathway – RNAi evidence.
105. Czech, B., Malone, C. D., Zhou, R., et al. (2008) *Nature*, **453**, 798–802, An endogenous small interfering RNA pathway in *Drosophila*.
106. Haac, M. E., Anderson, M. A., Eggleston, H., et al. (2015) *Nucleic Acids Res*, **43**, 3688–3700, The hub protein loquacious connects the microRNA and short interfering RNA pathways in mosquitoes.
107. Hammond, S. M., Boettcher, S., Caudy, A. A., et al. (2001) *Science (80-. )*, **293**, 1146–1150, Argonaute2, a link between genetic and biochemical analyses of RNAi.
108. Liu, J., Carmell, M. A., Rivas, F. V., et al. (2004) *Science (80-. )*, **305**, 1437–1441, Argonaute2 is the catalytic engine of mammalian RNAi.
109. Meister, G. and Tuschl, T. (2004) *Nature*, **431**, 343–349, Mechanisms of gene silencing by double-stranded RNA.
110. Verdel, A., Jia, S., Gerber, S., et al. (2004) *Science (80-. )*, **303**, 672–676, RNAi-mediated targeting of heterochromatin by the RITS complex.
111. Kennerdell, J. R., Yamaguchi, S. and Carthew, R. W. (2002) *Genes Dev*, **16**, 1884–1889, RNAi is activated during *Drosophila* oocyte maturation in a manner dependent on aubergine and spindle-E.
112. Jaubert-Possamai, S., Rispe, C., Tanguy, S., et al. (2010) *Mol Biol Evol*, **27**, 979–987, Expansion of the miRNA pathway in the hemipteran insect *Acyrtosiphon pisum*.
113. Zhang, Y., Wang, Y., Wang, L., et al. (2016) *Pestic. Biochem. Physiol.*, **127**, 21–27, Knockdown of NADPH-cytochrome P450 reductase results in reduced resistance to



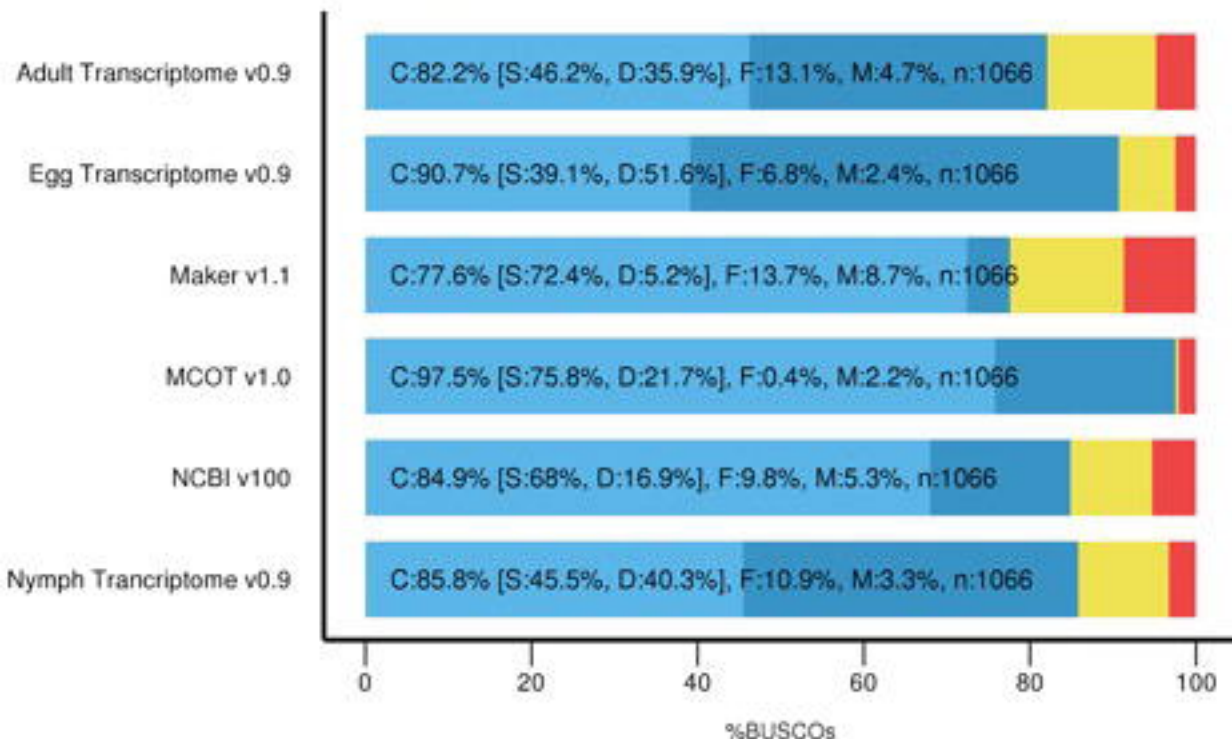
- buprofezin in the small brown planthopper, *Laodelphax striatellus* (fall{é}n).
114. McLean, K. J., Luciakova, D., Belcher, J., et al. In *Monooxygenase, Peroxidase and Peroxygenase Properties and Mechanisms of Cytochrome P450*; Springer, 2015; pp. 299–317.
115. Feyereisen, R. (2006) *Biochem. Soc. Trans.*, **34**, 1252–1255, Evolution of insect P450.
116. Schama, R., Pedrini, N., Juárez, M. P., et al. (2016) *Insect Biochem. Mol. Biol.*, **69**, 91–104, *Rhodnius prolixus* supergene families of enzymes potentially associated with insecticide resistance.
117. Richards, S., Gibbs, R. A., Weinstock, G. M., et al. (2008) *Nature*, **452**, 949–55, The genome of the model beetle and pest *Tribolium castaneum*.
118. Zerbino, D. R. and Birney, E. (2008) *Genome Res.*, **18**, 821–9, Velvet: algorithms for de novo short read assembly using de Bruijn graphs.
119. Chaisson, M. J. and Tesler, G. (2012) *BMC Bioinformatics*, **13**, 238, Mapping single molecule sequencing reads using Basic Local Alignment with Successive Refinement (BLASR): Theory and Application.
120. English, A. C., Richards, S., Han, Y., et al. (2012) *PLoS One*, **7**, e47768, Mind the Gap: Upgrading Genomes with Pacific Biosciences RS Long-Read Sequencing Technology.
121. Stanke, M., Steinkamp, R., Waack, S., et al. (2004) *Nucleic Acids Res.*, **32**, W309–W312, AUGUSTUS: a web server for gene finding in eukaryotes.
122. Slater, G. S. C. and Birney, E. (2005) *BMC Bioinformatics*, **6**, 31, Automated generation of heuristics for biological sequence comparison.
123. Wu, T. D. and Watanabe, C. K. (2005) *Bioinformatics*, **21**, 1859–1875, GMAP: a genomic mapping and alignment program for mRNA and EST sequences.
124. Kim, D., Langmead, B. and Salzberg, S. L. (2015) *Nat. Methods*, **advance on**, HISAT: a fast spliced aligner with low memory requirements.
125. Perte, M., Perte, G. M., Antonescu, C. M., et al. (2015) *Nat. Biotechnol.*, **advance on**, StringTie enables improved reconstruction of a transcriptome from RNA-seq reads.
126. Haas, B. J., Papanicolaou, A., Yassour, M., et al. (2013) *Nat. Protoc.*, **8**, 1494–1512, De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis.
127. Kim, D., Perte, G., Trapnell, C., et al. (2013) *Genome Biol.*, **14**, R36, TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions.
128. Trapnell, C., Roberts, A., Goff, L., et al. (2012) *Nat. Protoc.*, **7**, 562–578, Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks.
129. Jones, P., Binns, D., Chang, H.-Y., et al. (2014) *Bioinformatics*, **30**, 1236–1240, InterProScan 5: genome-scale protein function classification.
130. Vargas Jentzsch, I. M., Hughes, D. S. T. and Poelchau, M. F. T. The *O. fasciatus* curation community, Richards S, Panfilio KA. 2015. *Oncopeltus fasciatus* official gene set v1. 1.
131. Marchler-Bauer, A., Derbyshire, M. K., Gonzales, N. R., et al. (2015) *Nucleic Acids Res.*, **43**, D222–6, CDD: NCBI's conserved domain database.
132. Edgar, R. C. (2004) *Nucleic Acids Res.*, **32**, 1792–1797, MUSCLE: multiple sequence alignment with high accuracy and high throughput.
133. Notredame, C., Higgins, D. G. and Heringa, J. (2000) *J. Mol. Biol.*, **302**, 205–217, T-Coffee: A novel method for fast and accurate multiple sequence alignment.
134. Sievers, F., Wilm, A., Dineen, D., et al. (2011) *Mol. Syst. Biol.*, **7**, 539, Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega.
135. Kumar, S., Stecher, G. and Tamura, K. (2016) *Mol. Biol. Evol.*, msw054, MEGA7: Molecular Evolutionary Genetics Analysis version 7.0 for bigger datasets.
136. Consortium, G. O. and others (2004) *Nucleic Acids Res.*, **32**, D258–D261, The Gene Ontology (GO) database and informatics resource.







# BUSCO Assessment Results



AED cumulative plots

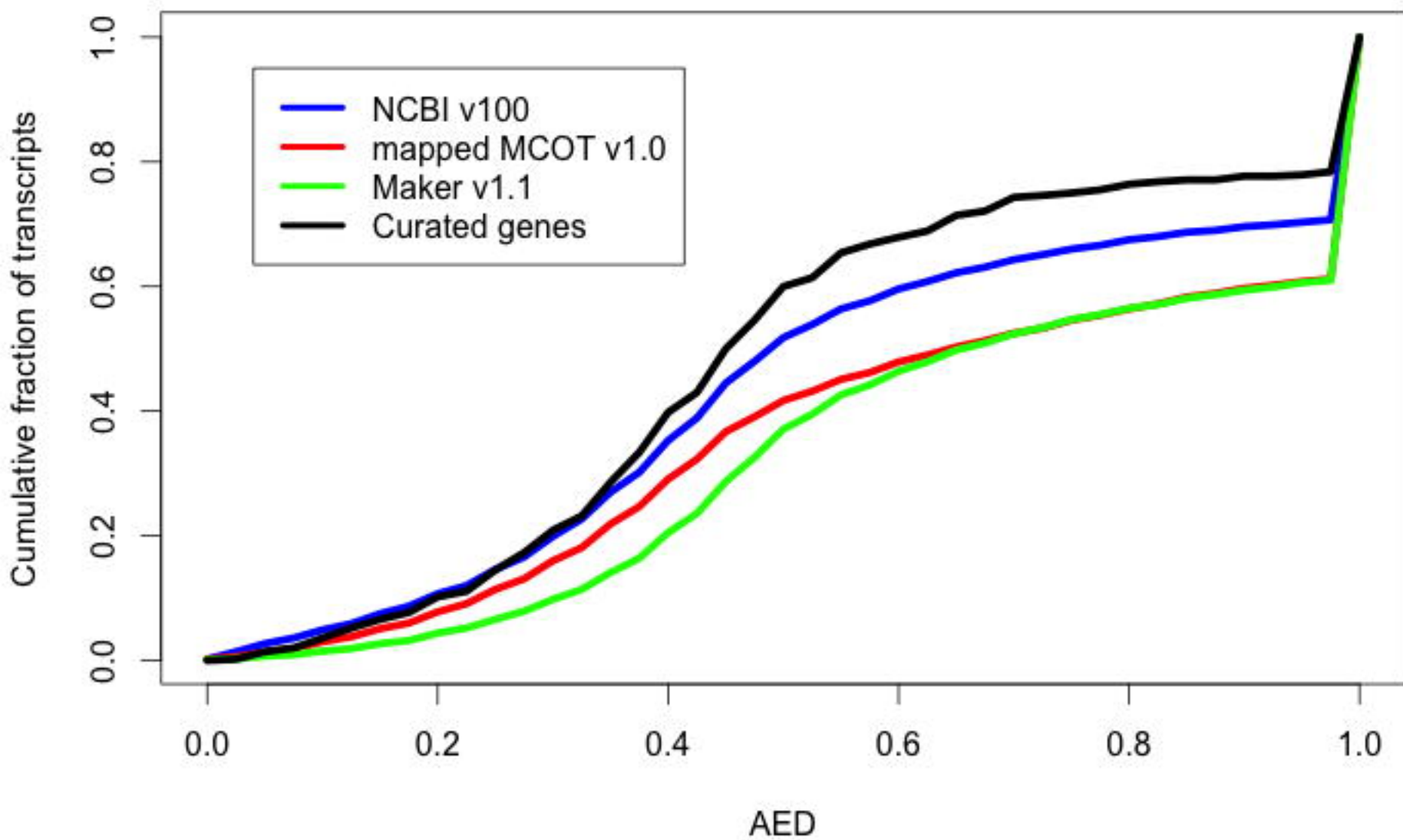


Table 1: Immune gene pathways and gene counts in ACP, Pea and Aphid Whitefly. ACP genes identified only in MCOT v1.0 are in ().

	Pathway/Genes	ACP	Pea aphid	Whitefly
<i>Pathogen recognition molecules</i>				
	CTLs: C-Type Lectins	10	5	5
	GALEs: Galactoside-Binding Lectins	4	1	4
	FREPs: Fibrinogen-Related Proteins	3	1	5
	PGRPs: Peptidoglycan Recognition Proteins	1	0	1
	BGBPs: 1,3-beta-D Glucan Binding Proteins	0	2	5
	TEPs: Thio-Ester Containing Proteins	2	1	1
<i>Signaling cascades associated with pathogenesis</i>				
Toll pathway and receptors	Toll receptors	5	7	5
	Spaetzle	5	6	9
	Tube	1	1	1
	Pelle	1	1	1
	MyD88	1	1	1
	CACT	1	1	0
	TRAF6	1	2	1
IMD pathway members	CASPAR	3	1	1
	FADD	0	0	0
	IKKB/ird5	1	0	0
	IMD	0	0	0
	TAK1	1	1	1
	TAB	1	1	1
JAK/STAT pathway members	DOME	1	3	1
	HOP	1	1	1
	STAT	1	3	1
<i>Response to pathogens and pathogen-associated stress</i>				
	AMPs: Anti-Microbial Peptides	0	7	4



	LYSs: Lysozymes	5	3	5
	SODs: Superoxide Dismutases	4	4	5
	CLIPs: CLIP-Domain Serine Proteases	14	3	6
	Autophagy	15 (2)	8	16
	PPOs: Prophenoloxidases	2	2	4
	IAPs: Inhibitors of Apoptosis	4	7	4

Table 2: Homolog number of core machinery and auxiliary RNAi components in insects. Proteins highlighted in orange are RNase Type III enzymes. Proteins highlighted in blue are dsRNA binding proteins. Proteins highlighted in green are AGO proteins. Proteins highlighted in yellow have been implicated in RISC or in small RNA biogenesis. \*indicates homolog was determined by BLAST and reciprocal BLAST analysis using NCBI's non redundant databases. Homolog number with an asterisk were determined by publications or reported in ImmunoDB. ? Indicates homolog is likely present but not annotated in the current genome assembly. Note: *Drosophila melanogaster* has six other RNA helicase genes related to Rm62. Some of the mosquito Rm62 proteins could be orthologous to these related RNA helicases.

	Dcr1	Dcr2	Drosha	Loqs	R2D2	Pasha	AGO1	AGO2	AGO3	PIWI/Aub	Armi	TSN	VIG-1	Spn-E
<i>D. melanogaster</i>	1	1	1	1	1	1	1	1	1	2	1	1	1	1
<i>A. gambiae</i>	1	1	1	1	1	1	1	1	1	2	1	1	1	1
<i>A. aegypti</i>	1	1	1	1	1	1	2	1	1	7	2	1	1	1
<i>C. quinquefasciatus</i>	1	1	1	1	1	1	1	2	1	6	2	1	0*	1
<i>T. castaneum</i>	1	1	1	1	2	1	1	2	1	1	3*	1*	2*	1*
<i>C. lectularius</i>	1	1	1	1	1	1	2	2	1	4	2*	1*	1*	1*
<i>D. citri</i>	1	1	1 or 2	2	1 <sup>ψ</sup>	2	1	1	1	1	1	2	1	2

d in orange are  
 en are AGO family  
 omolog number  
 mber with no  
 was unable to be  
 es with homology

Rm62	Ran	FMR1
1	1	1
6	1*	0
9	1*	1
10	1	1
2*	1*	2*
0*	1*	1*
1	1	1

Species	Number of P450s
<i>Apolygus lucorum</i> (mirid bug)	46
<i>Acyrtosiphon pisum</i> (pea aphid)	58
<i>Cyrtorhinus lividipennis</i> (green miridbug)	59
<i>Cimex lectularius</i> (bedbug)	60
<b><i>Diaphorina citri</i> (Asian citrus psyllid)</b>	60
<i>Laodelphax striatellus</i> (small brown planthopper)	63
<i>Nilaparvata lugens</i> (brown planthopper)	66
<i>Rhodnius prolixus</i> (kissing bug)	87
<i>Bemisia tabaci</i> (whitefly)	128
<i>Homalodisca vitripennis</i> (glassy winged sharpshooter)	142