

# Lawrence Berkeley National Laboratory

## LBL Publications

### Title

Molecular evolutionary trends and feeding ecology diversification in the Hemiptera, anchored by the milkweed bug genome.

### Permalink

<https://escholarship.org/uc/item/3dc8n4n3>

### Journal

Genome biology, 20(1)

### ISSN

1474-7596

### Authors

Panfilio, Kristen A  
Vargas Jentzsch, Iris M  
Benoit, Joshua B  
[et al.](#)

### Publication Date

2019-04-01

### DOI

10.1186/s13059-019-1660-0

Peer reviewed

RESEARCH

Open Access



# Molecular evolutionary trends and feeding ecology diversification in the Hemiptera, anchored by the milkweed bug genome

Kristen A. Panfilio<sup>1,2\*</sup> , Iris M. Vargas Jentzsch<sup>1</sup> , Joshua B. Benoit<sup>3</sup> , Deniz Ereyilmaz<sup>4,43</sup> , Yuichiro Suzuki<sup>5</sup> , Stefano Colella<sup>6,7</sup> , Hugh M. Robertson<sup>8</sup>, Monica F. Poelchau<sup>9</sup> , Robert M. Waterhouse<sup>10,11</sup> , Panagiotis Ioannidis<sup>10</sup> , Matthew T. Weirauch<sup>12</sup> , Daniel S. T. Hughes<sup>13</sup>, Shwetha C. Murali<sup>13,14,15</sup>, John H. Werren<sup>16</sup> , Chris G. C. Jacobs<sup>17,18</sup> , Elizabeth J. Duncan<sup>19,20</sup> , David Armisen<sup>21</sup>, Barbara M. I. Vreede<sup>22</sup> , Patrice Baa-Puyoulet<sup>6</sup>, Chloé S. Berger<sup>21</sup>, Chun-che Chang<sup>23,45</sup> , Hsu Chao<sup>13</sup>, Mei-Ju M. Chen<sup>9</sup> , Yen-Ta Chen<sup>1</sup> , Christopher P. Childers<sup>9</sup> , Ariel D. Chipman<sup>22</sup> , Andrew G. Cridge<sup>19</sup> , Antonin J. J. Crumière<sup>21</sup> , Peter K. Dearden<sup>19</sup> , Elise M. Didion<sup>3</sup> , Huyen Dinh<sup>13</sup> , Harsha Vardhan Doddapaneni<sup>13</sup> , Amanda Dolan<sup>16,24</sup>, Shannon Dugan<sup>13</sup> , Cassandra G. Extavour<sup>25,26</sup> , Gérard Febvay<sup>6</sup> , Markus Friedrich<sup>27</sup> , Neta Ginzburg<sup>22</sup>, Yi Han<sup>13</sup> , Peter Heger<sup>28</sup>, Christopher J. Holmes<sup>3</sup> , Thorsten Horn<sup>1</sup> , Yi-min Hsiao<sup>23,45</sup> , Emily C. Jennings<sup>3</sup> , J. Spencer Johnston<sup>29</sup> , Tamsin E. Jones<sup>25</sup> , Jeffery W. Jones<sup>27</sup>, Abderrahman Khila<sup>21</sup> , Stefan Koelzer<sup>1</sup>, Viera Kovacova<sup>30</sup> , Megan Leask<sup>19</sup>, Sandra L. Lee<sup>13</sup>, Chien-Yueh Lee<sup>9</sup> , Mackenzie R. Lovegrove<sup>19</sup>, Hsiao-ling Lu<sup>23,45</sup> , Yong Lu<sup>31</sup>, Patricia J. Moore<sup>32</sup> , Monica C. Munoz-Torres<sup>33</sup> , Donna M. Muzny<sup>13</sup> , Subba R. Palli<sup>34</sup> , Nicolas Parisot<sup>6</sup> , Leslie Pick<sup>31</sup> , Megan L. Porter<sup>35</sup> , Jiaxin Qu<sup>13</sup> , Peter N. Refki<sup>21,36</sup>, Rose Richter<sup>16,37</sup>, Rolando Rivera-Pomar<sup>38</sup> , Andrew J. Rosendale<sup>3</sup> , Siegfried Roth<sup>1</sup> , Lena Sachs<sup>1</sup>, M. Emília Santos<sup>21</sup>, Jan Seibert<sup>1</sup>, Essia Sghaier<sup>21</sup>, Jayendra N. Shukla<sup>34,39</sup> , Richard J. Stancliffe<sup>40,44</sup> , Olivia Tidswell<sup>19,41</sup>, Lucila Traverso<sup>42</sup>, Maurijn van der Zee<sup>17</sup> , Séverine Viala<sup>21</sup>, Kim C. Worley<sup>13</sup> , Evgeny M. Zdobnov<sup>10</sup>, Richard A. Gibbs<sup>13</sup> and Stephen Richards<sup>13</sup> 

## Abstract

**Background:** The Hemiptera (aphids, cicadas, and true bugs) are a key insect order, with high diversity for feeding ecology and excellent experimental tractability for molecular genetics. Building upon recent sequencing of hemipteran pests such as phloem-feeding aphids and blood-feeding bed bugs, we present the genome sequence and comparative analyses centered on the milkweed bug *Oncopeltus fasciatus*, a seed feeder of the family Lygaeidae.

(Continued on next page)

\* Correspondence: [kristen.panfilio@alum.swarthmore.edu](mailto:kristen.panfilio@alum.swarthmore.edu)

<sup>1</sup>Institute for Zoology: Developmental Biology, University of Cologne, Zùlpicher Str. 47b, 50674 Cologne, Germany

<sup>2</sup>School of Life Sciences, University of Warwick, Gibbet Hill Campus, Coventry CV4 7AL, UK

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



(Continued from previous page)

**Results:** The 926-Mb *Oncopeltus* genome is well represented by the current assembly and official gene set. We use our genomic and RNA-seq data not only to characterize the protein-coding gene repertoire and perform isoform-specific RNAi, but also to elucidate patterns of molecular evolution and physiology. We find ongoing, lineage-specific expansion and diversification of repressive C2H2 zinc finger proteins. The discovery of intron gain and turnover specific to the Hemiptera also prompted the evaluation of lineage and genome size as predictors of gene structure evolution. Furthermore, we identify enzymatic gains and losses that correlate with feeding biology, particularly for reductions associated with derived, fluid nutrition feeding.

**Conclusions:** With the milkweed bug, we now have a critical mass of sequenced species for a hemimetabolous insect order and close outgroup to the Holometabola, substantially improving the diversity of insect genomics. We thereby define commonalities among the Hemiptera and delve into how hemipteran genomes reflect distinct feeding ecologies. Given *Oncopeltus*'s strength as an experimental model, these new sequence resources bolster the foundation for molecular research and highlight technical considerations for the analysis of medium-sized invertebrate genomes.

**Keywords:** Phytophagy, Transcription factors, Gene structure, Lateral gene transfer, RNAi, Gene family evolution, Evolution of development

## Background

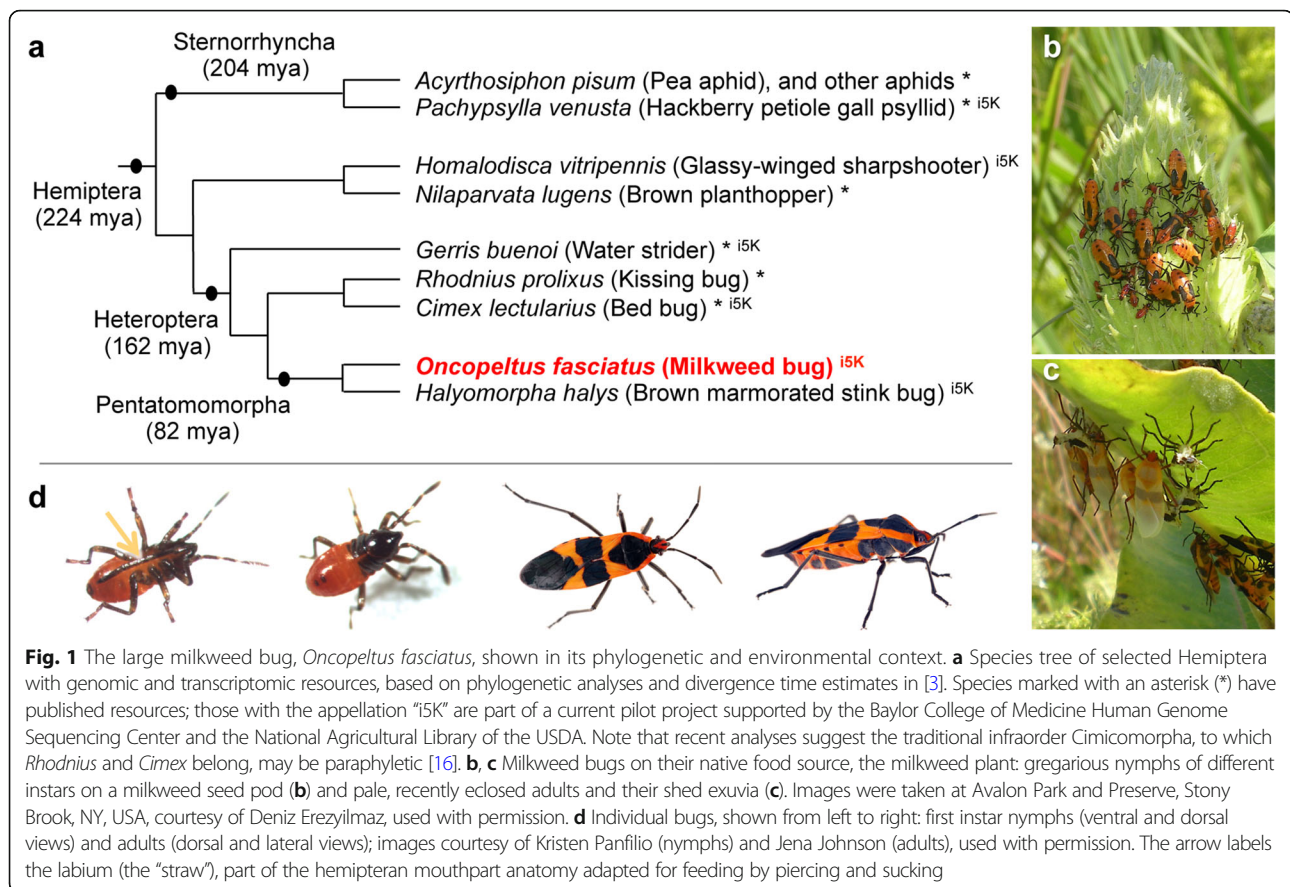
The number of animals with sequenced genomes continues to increase dramatically, and there are now over 100 insect species with assembled and annotated genomes [1]. However, the majority belong to the Holometabola (e.g., flies, beetles, wasps, butterflies), the group characterized by a biphasic life history with distinct larval and adult phases separated by dramatic metamorphosis during a pupal stage. The Holometabola represent only a fraction of the full morphological and ecological diversity across the Insecta: over half of all orders are hemimetabolous. Imbalance in genomic resources limits the exploration of this diversity, including the environmental and developmental requirements of a hemimetabolous lifestyle with a progression of flightless nymphal (juvenile) instars. Addressing this paucity, we report comparative analyses based on genome sequencing of the large milkweed bug, *Oncopeltus fasciatus*, as a hemimetabolous representative of the larger diversity of insects.

*Oncopeltus* is a member of the Hemiptera, the most species-rich hemimetabolous order. Together with the Thysanoptera and, traditionally, the Psocodea, the Hemiptera form the hemipteroid assemblage (or Acercaria), a close outgroup to the Holometabola [2, 3]. All Hemiptera share the same piercing and sucking mouthpart anatomy [4], yet they have diversified to exploit food sources ranging from seeds and plant tissues (phytophagy) to phloem sap (mucivory) and vertebrate blood (hematophagy). For this reason, many hemipterans are agricultural pests or human disease vectors, and genome sequencing efforts to date have focused on these species (Fig. 1, [5]), including phloem-feeding aphids [6–8], psyllids [9], and planthoppers [10], and the hematophagous kissing bug, *Rhodnius prolixus* [11], a vector of Chagas disease, and bed bug, *Cimex lectularius* [12, 13]. Building on transcriptomic data, genome projects are also in

progress for other pest species within the same infra-order as *Oncopeltus*, such as the stink bug *Halyomorpha halys* [14, 15].

The milkweed bug has feeding ecology traits that are both conservative and complementary to those of previously sequenced hemipterans. Its phytophagy is ancestral for the large infraorder Pentatomomorpha and representative of most extant Hemiptera [16]. Moreover, as a seed feeder, *Oncopeltus* has not undergone the marked lifestyle changes associated with fluid feeding (mucivory or hematophagy), including dependence on endosymbiotic bacteria to provide nutrients lacking in the diet. Gene loss in the pea aphid, *Acyrtosiphon pisum*, makes it reliant on the obligate endosymbiont *Buchnera aphidicola* for synthesis of essential amino acids [6, 17]. Although hematophagy arose independently in *Rhodnius* and *Cimex* [16], their respective endosymbionts, *Rhodococcus rhodnii* and *Wolbachia*, must provide vitamins lacking in a blood diet [18]. In contrast, the seed-feeding subfamily Lygaeinae, including *Oncopeltus*, is notable for the absence of prominent endosymbiotic anatomy: these bugs lack both the midgut crypts that typically house bacteria and the bacteriomes and endosymbiotic balls seen in other Lygaeidae [19].

As the native food source of *Oncopeltus* is the toxic milkweed plant, its own feeding biology has a number of interesting implications regarding detoxification and sequestration of cardenolide compounds. A prominent consequence of this diet is the bright red-orange aposomatic (warning) coloration seen in *Oncopeltus* embryos, nymphs, and adults [20, 21]. Thus, diet, metabolism, and body pigmentation are functionally linked biological features for which one may expect changes in gene repertoires to reflect the diversity within an order, and the Hemiptera provide an excellent opportunity to explore this.



Furthermore, *Oncopeltus* has been an established laboratory model organism for over 60 years, with a rich experimental tradition in a wide range of studies from physiology and development to evolutionary ecology [21–23]. It is among the few experimentally tractable hemimetabolous insect species, and it is amenable to a range of molecular techniques (e.g., [24–26]). In fact, it was one of the first insect species to be functionally investigated by RNA interference (RNAi, [27]). RNAi in *Oncopeltus* is highly effective across different life history stages, which has led to a resurgence of experimental work over the past 15 years, with a particular focus on the evolution of developmentally important regulatory genes (reviewed in [23]).

Here, we focus on these two themes—feeding biology diversity within the Hemiptera and *Oncopeltus* as a research model for macroevolutionary genetics. Key insights derive from a combination of global comparative genomics and detailed computational analyses that are supported by extensive manual curation, empirical data for gene expression, sequence validation, and new isoform-specific RNAi. We thereby identify genes with potentially restricted life history expression in *Oncopeltus* and that are unique to the Hemiptera, clarify evolutionary patterns of zinc finger protein family expansion,

categorize predictors of insect gene structure, and identify lateral gene transfer and amino acid metabolism features that correlate with feeding biology.

## Results and discussion

### The genome and its assembly

*Oncopeltus fasciatus* has a diploid chromosome number ( $2n$ ) of 16, comprised of seven autosomal pairs and two sex chromosomes with the XX/XY sex determination system [28, 29]. To analyze this genetic resource, we sequenced and assembled the genome using next-generation sequencing approaches (Table 1, see also the “Methods” section and Additional file 1: Supplemental Notes Sections 1–4). We measure the genome size to be 923 Mb in females and 928 Mb in males based on flow cytometry data (Additional file 1: Supplemental Note 2.1.a). The assembly thus contains 84% of the expected sequence, which is comparable to other recent, medium-sized insect genomes [12, 30]. However, our analyses of the  $k$ -mer frequency distribution in raw sequencing reads yielded ambiguous estimates of genome size and heterozygosity rate, which is suggestive of high heterozygosity and repetitive content ([31], Additional file 1: Supplemental Note 2.1.b). In further analyses, we indeed obtained high estimates of repetitive content, although heterozygosity does not unduly influence gene

**Table 1** *Oncopeltus fasciatus* genome metrics

Feature	Value	
2n chromosomes	16	
Genome size	926 Mb (mean between males and females)	
Assembly size	1099 Mb (contigs only, 774 Mb)	
Coverage	106.9× raw coverage, 83.7% of reads in final assembly	
Contig N50	4047 bp	
Scaffold N50	340.0 kb	
# scaffolds	17,222	
GC content	genome, 32.7%; protein-coding sequence (OGS v1.2), 42%	
OGS v1.1 (curated fraction)	19,690 models <sup>1</sup> (1426 models, 7.2%)	19,465 genes (1201 genes, 6.2%)
OGS v1.2 (curated fraction)	19,809 models <sup>1</sup> (1697 models, 8.7%)	19,616 genes (1518 genes, 7.7%)

<sup>1</sup>Individual genes may be represented by multiple models in cases of curated alternative isoforms or if exons of the gene are split across scaffolds

prediction (see below, based on protein orthology assessments). These computationally challenging features may be increasingly relevant as comparative genomics extends to insect species with larger genomes (> 1 Gb)—a common feature among hemimetabolous insects [5, 32].

As template DNA was prepared from dissected adults from which the gut material was removed, the resulting assembly is essentially free of contamination. Only five small scaffolds had high bacterial homology, each to a different, partial bacterial genome (Additional file 1: Supplemental Note 2.2).

#### The official gene set and conserved gene linkage

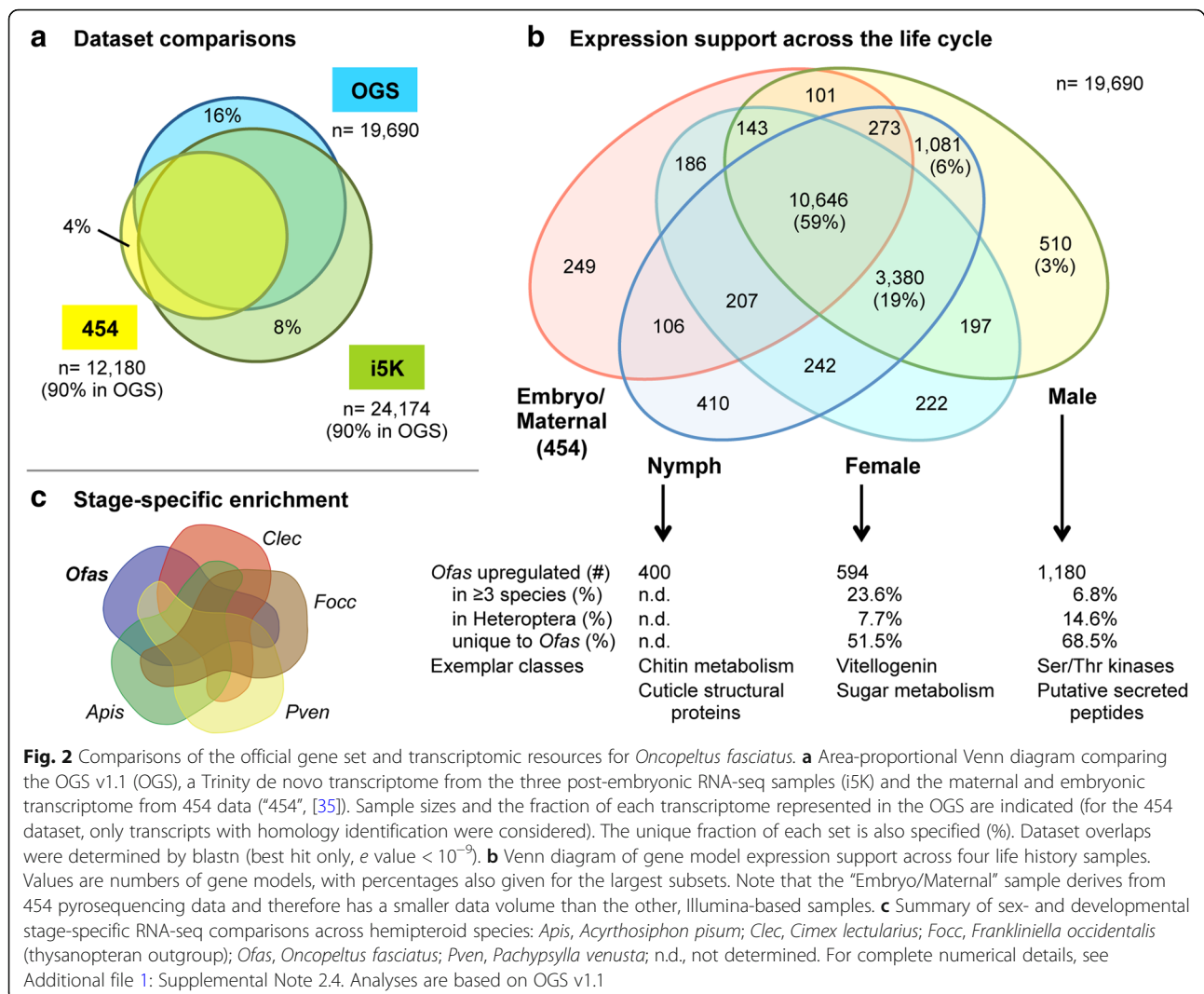
The official gene set (OGS) was generated by automatic annotation followed by manual curation in a large-scale effort by the research community (Additional file 1: Supplemental Notes Sections 3–4). Curation revised automatic models, added alternative isoforms and de novo models, and documented multiple models for genes split across scaffolds. We found that automatic predictions were rather conservative for hemipteran gene structure (see below). Thus, manual curation often extended gene loci as exons were added, including merging discrete automatic models (Additional file 1: Supplemental Note 4, and Table S4.4). The OGS v1.1 was generated for global analyses to characterize the gene repertoire. The latest version, OGS v1.2, primarily adds chemoreceptor genes of the ionotropic and odorant receptor classes and genes encoding metabolic enzymes. Altogether, the research community curated 1697 gene models (8.7% of OGS v1.2), including 316 de novo models (Additional file 2: Table S4.1, Additional file 1: Supplemental Notes Section 5). The majority of curated models are for genes encoding cuticular proteins (11%), chemoreceptors (19%), and developmental regulators such as transcription factors and signaling pathway components (40%, including the BMP/TGF- $\beta$ , Toll/NF- $\kappa$ B, Notch, Hedgehog, Torso RTK, and Wnt pathways).

In addition to assessing gene model quality, manual curation of genes whose orthologs are expected to occur in syntenic clusters also validates assembly scaffolding. Complete loci could be found for single orthologs of all Hox cluster genes, where *Hox3/zen* and *Hox4/Dfd* are linked in the current assembly and have  $\geq 99.9\%$  nucleotide identity with experimentally validated sequences ([27, 33, 34], Additional file 1: Supplemental Note 5.1.b). Conserved linkage was also confirmed for the homeobox genes of the Iroquois complex, the Wnt ligands *wingless* and *wnt10*, and two linked pairs from the Runt transcription factor complex (Additional file 1: Supplemental Notes 5.1.a, 5.1.c, 5.1.i, 5.1.j). Further evidence for correct scaffold assembly comes from the curation of large, multi-exonic loci. For example, the cell polarity and cytoskeletal regulator encoded by the conserved *furry* gene includes 47 exons spanning a 437-kb locus, which were all correctly assembled on a single scaffold.

#### Gene expression profiles across the milkweed bug life cycle

To augment published transcriptomic resources [35, 36], we sequenced three different post-embryonic samples (“i5K” dataset, see the “Methods” section). We then compared the OGS to the resulting de novo transcriptome and to a previously published embryonic and maternal (ovary) transcriptome (“454” pyrosequencing dataset, [35]). Our OGS is quite comprehensive, containing 90% of transcripts from each transcriptomic dataset (Fig. 2a). The OGS also contains an additional 3146 models (16% of OGS) not represented in either transcriptome, including 163 de novo models encoding chemoreceptors. Such genes are known for lineage-specific expansions and highly tissue- and stage-specific expression ([37, 38], and see below), and our OGS captures these genes with rare transcripts.

The OGS also incorporates many partial and unidentified 454 transcripts, nearly trebling the transcripts with



an assigned gene model or homology compared to the original study (from 9 to 26%, by blastn,  $e < 10^{-9}$ ). This included 10,130 transcripts that primarily mapped to UTRs and previously lacked recognizable coding sequence, such as for the *Oncopeltus brinker* ortholog, a BMP pathway component ([39], Additional file 1: Supplemental Note 5.1.f), and the enzyme-encoding genes *CTP synthase* and *roquin*. At the same time, the transcriptomes provided expression support for the identification of multiple isoforms in the OGS, such as for the germline determinant *nanos* [35]. More generally, most OGS gene models have expression support (91% of 19,690), with 74% expressed broadly in at least 3 of 4 samples (Fig. 2b). The inclusion of a fifth dataset from a published adult library [36] provided only a 1% gain in expression support, indicating that with the current study the expression data volume for *Oncopeltus* is quite complete.

RNA-seq studies were further conducted to establish male-, female-, and nymph-specific gene sets (Fig. 2b, c, Additional file 1: Supplemental Note 2.4), from which we also infer that the published adult dataset of unspecified sex is probably male. Moreover, most genes with stage-restricted or stage-enriched expression are in our male sample (Fig. 2b, c). For example, gustatory receptor (GR) genes show noticeable restriction to the adult male and published adult (probable male) samples ( $n = 169$  GRs: 40% no expression, 27% only expressed in these two samples), with half of these expressed in both biological replicates (52%). Interestingly, the nymphal sample is enriched for genes encoding structural cuticular proteins (94%, which is  $> 56\%$  more than any other sample). This likely reflects the ongoing molting cycles, with their cyclical upregulation of chitin metabolism and cuticular gene synthesis [40], that are experienced by the different instars and molt cycle stages of individuals pooled in this sample. Lastly, gene sets with sex-specific

enrichment across several hemipteroid species substantiate known aspects of male and female reproduction (Fig. 2c: serine-threonine kinases [41] or vitellogenin and other factors associated with oocyte generation, respectively). Some of these enriched genes have unknown functions and could comprise additional, novel factors associated with reproduction in *Oncopeltus*.

### Protein orthology and hemipteran copy number comparisons

To further assay protein-coding gene content, we compared *Oncopeltus* with other arthropods. A phylogeny based on strictly conserved single-copy orthologs correctly reconstructs the hemipteran and holometabolan clades' topologies (Fig. 3a, compare with Fig. 1a), although larger-scale insect relationships remain challenging [3].

We then expanded our appraisal to the Benchmarking Universal Single-Copy Orthologs dataset of 1658 Insecta genes (BUSCO v3, [42]). Virtually all BUSCO genes are present in the *Oncopeltus* OGS (98.9%, Fig. 3b, Additional file 1: Supplemental Note 6.1). Although some genes are fragmented, the assembly has a high level of BUSCO completeness (94.6%), independent of the annotation prediction limitations that missed some exons from current gene models. Furthermore, BUSCO assessments can elucidate potential consequences of high heterozygosity, which could result in the erroneous inclusion of multiple alleles for a single gene. In fact, the fraction of duplicated BUSCO genes in *Oncopeltus* (1.4%) is low, compared to both the well-assembled bed bug genome (2.2%, [12]) and the pea aphid (4.8%), which is known to have lineage-specific duplications [6, 43]. Thus, by these quality metrics, the *Oncopeltus* OGS and assembly are comparable to those of fellow hemipterans, strongly supporting the use of these resources in further comparisons.

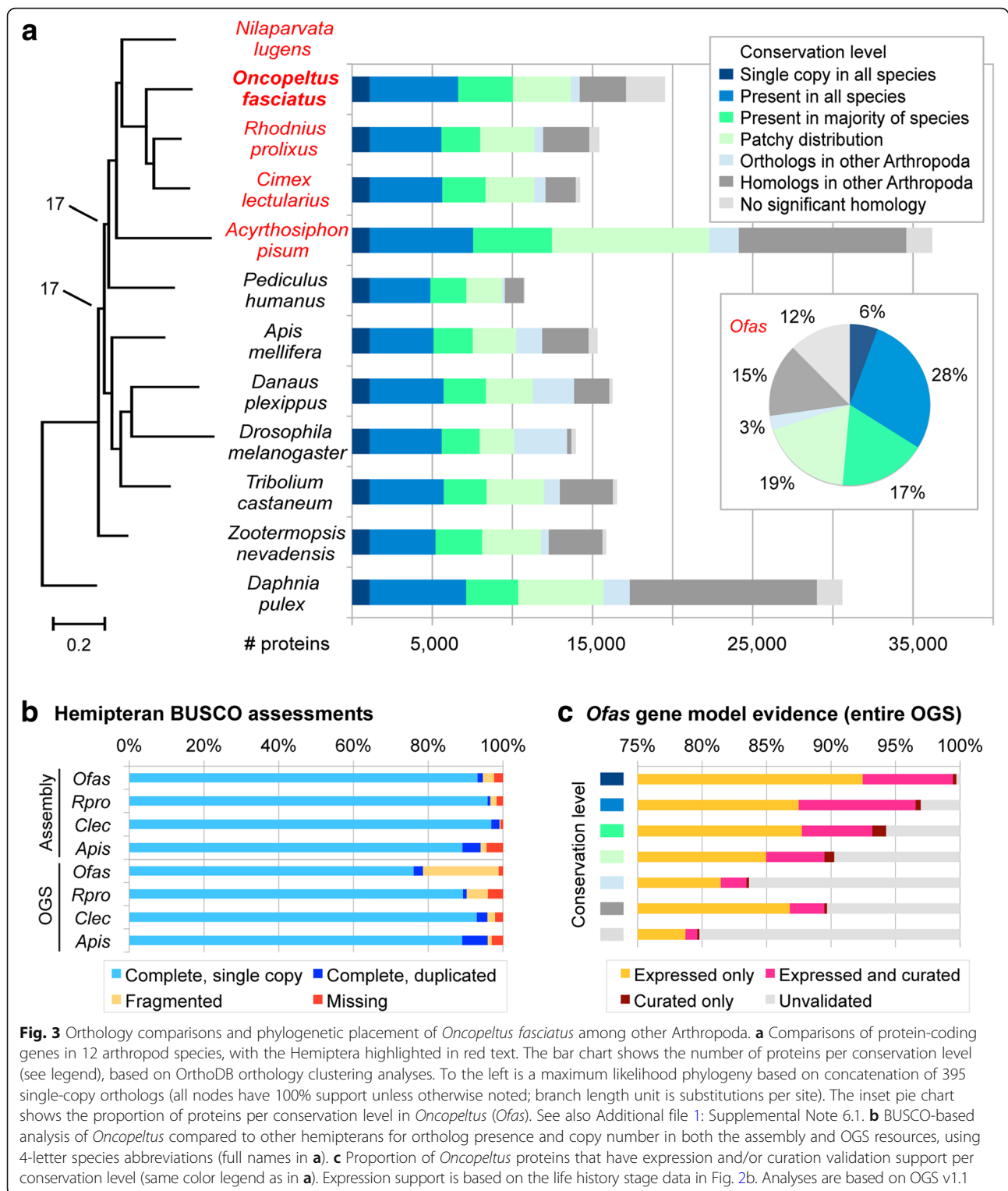
We next categorized all proteins by conservation in global, clustering-based orthology analyses (OrthoDB, [1, 44]). As in most species, half of *Oncopeltus* proteins are highly conserved (Fig. 3a). Moreover, 98% of all *Oncopeltus* protein-coding genes have homology, expression, and/or curation support (Fig. 3c). Proteins without homology include species-specific chemoreceptors and antimicrobial peptides (Additional file 1: Supplemental Note 5.1.h), as well as potentially novel or partial models. Overall, we estimate that the *Oncopeltus* protein repertoire is comparable to that of other insects in size and conservation. For the Hemiptera, *Oncopeltus* also has fewer missing orthology groups than either the kissing bug or pea aphid (Additional file 1: Table S6.1). Indeed, the pea aphid is a notable outlier, with its long branch in the phylogeny and for its large protein-coding gene content with low conservation (Fig. 3a). As more hemipteran genomes are sequenced,

other species now offer less derived alternatives for phylogenomic comparisons.

Compared to the pea aphid [43], *Oncopeltus* is more conservative in presence and copy number for several signaling pathway components. In contrast to gene absences in the pea aphid, *Oncopeltus* retains orthologs of the EGF pathway component *sprouty*, the BMP receptor *wishful thinking*, and the hormone nuclear receptor *Hr96* (Additional file 1: Supplemental Note 5.1.e). Also, whereas multiple copies were reported for the pea aphid, we find a single *Oncopeltus* ortholog for the BMP pathway components *decapentaplegic* and *Medea* and the Wnt pathway intracellular regulator encoded by *shaggy/GSK-3*, albeit with five potential isoforms of the latter (Additional file 1: Supplemental Notes 5.1.f, 5.1.j). Duplications of miRNA and piRNA gene silencing factors likewise seem to be restricted to the pea aphid, even compared to other aphid species ([45], Additional file 1: Supplemental Note 5.4.a). However, our survey of *Oncopeltus* and other hemimetabolous species reveals evidence for frequent, independent duplications of the Wnt pathway component *armadillo/β-catenin* ([46], Additional file 1: Supplemental Note 5.1.j). Curiously, *Oncopeltus* appears to encode fewer histone loci than any other arthropod genome and yet exhibits a similar, but possibly independent, pattern of duplications of histone acetyltransferases to those previously identified in *Cimex* and the pea aphid (Additional file 1: Supplemental Note 5.4.c).

On the other hand, we documented several notable *Oncopeltus*-specific duplications. For the BMP transducer *Mad*, we find evidence for three paralogs in *Oncopeltus*, where two occur in tandem and may reflect a particularly recent duplication (Additional file 1: Supplemental Note 5.1.f). Similarly, a tandem duplication of *wnt8* appears to be unique to *Oncopeltus* (Additional file 1: Supplemental Note 5.1.j). More striking is the identification of six potential paralogs of *cactus*, a member of the Toll/NF-κB signaling pathway for innate immunity, whereas the bed bug and kissing bug each retain only a single copy ([47], Additional file 1: Supplemental Note 5.1.g).

Lastly, we explored hemipteran-specific orthology groups against a backdrop of 107 other insect species [1]. What makes a bug a bug in terms of protein-coding genes? Several orthogroups contain potentially novel genes that show no homology outside the Hemiptera and await direct experimental analysis, for which the Hemiptera are particularly amenable (e.g., [5, 48–51]). Secondly, there are hemipteran-specific orthogroups of proteins with recognized functional domains and homologs in other insects, but where evolutionary divergence has led to lineage-specific subfamilies. One example is a heteropteran-specific cytochrome P450 (CYP) enzyme (EOG090W0V4B), which in *Oncopeltus* is expressed in all life history stages (Fig. 2b). The expansion of CYP protein families is associated with potential insecticide



resistance, as specific P450s can confer resistance to specific chemicals (e.g., [52, 53]; Additional file 1: Supplemental Notes 5.3.b, 5.3.c). Hence, the identification of lineage-specific CYP enzymes can suggest potential targets for integrated pest management approaches.

### Transcription factor repertoires and homeobox gene evolution

Having explored the global protein repertoire, we next focused specifically on transcription factors (TFs), which comprise a major class of proteins that has been



extensively studied in *Oncopeltus*. This is a class of key regulators of development whose functions can diverge substantially during evolution and for which RNAi-based experimental investigations have been particularly fruitful in the milkweed bug (e.g., [27, 33, 54–56], Additional file 1: Supplemental Notes 5.1.a–e).

To systematically evaluate the *Oncopeltus* TF repertoire, we used a pipeline to scan all predicted proteins and assign them to TF families, including orthology assignments where DNA binding motifs could be predicted (see the “Methods” section, [57]). We identified 762 putative TFs in *Oncopeltus*, which is similar to other insects for total TF count and for the size of each TF family (Fig. 4a: note that the heatmap also reflects the large, duplicated repertoire in the pea aphid, see also Additional file 2: Tables S6.3–S6.5).

We were able to infer DNA binding motifs for 25% ( $n = 189$ ) of *Oncopeltus* TFs, mostly based on data from *Drosophila melanogaster* (121 TFs) but also from distantly related taxa such as mammals (56 TFs). Such high conservation is further reflected in explicit orthology assignments for most proteins within several large TF families, including the homeodomain (53 of 85, 62%), basic helix-loop-helix (bHLH, 35 of 45, 78%), and fork-head box (16 of 17, 94%) families. In contrast, most C2H2 zinc finger proteins lack orthology assignment (only 22 of 360, 6%). Across species, the homeodomain and C2H2 zinc finger proteins are the two largest TF superfamilies (Fig. 4a). Given their very different rates of orthology assignment, we probed further into their pipeline predictions and the patterns of evolutionary diversification.

The number of homeodomain proteins identified by the pipeline displays a narrow normal distribution across species (Fig. 4b, mean  $\pm$  standard deviation  $97 \pm 9$ ), consistent with a highly conserved, slowly evolving protein family. Supporting this, many *Oncopeltus* homeodomain proteins that were manually curated also received a clear orthology assignment (Fig. 4c: pink), with only 4 exceptions (Fig. 4c: yellow). Only 1 case suggests a limitation of a pipeline that is not specifically tuned to hemipteran proteins (Goosecoid). Manual curation of partial or split models identified 11 further genes encoding homeodomains, bringing the actual tally in *Oncopeltus* to 96. Overall, we find the TF pipeline results to be a robust and reasonably comprehensive representation of these gene classes in *Oncopeltus*.

These analyses also uncovered a correction to the published *Oncopeltus* literature for the developmental patterning proteins encoded by the paralogs *engrailed* and *invected*. These genes arose from an ancient tandem duplication prior to the hexapod radiation. Their tail-to-tail orientation enables ongoing gene conversion [58], making orthology discrimination particularly

challenging. For *Oncopeltus*, we find that the genes also occur in a tail-to-tail orientation and that *invected* retains a diagnostic alternative exon [58]. These new data reveal that the purported *Oncopeltus engrailed* ortholog in previous developmental studies (e.g., [54, 59–62]) is in fact *invected* (Additional file 1: Supplemental Note 5.1.a).

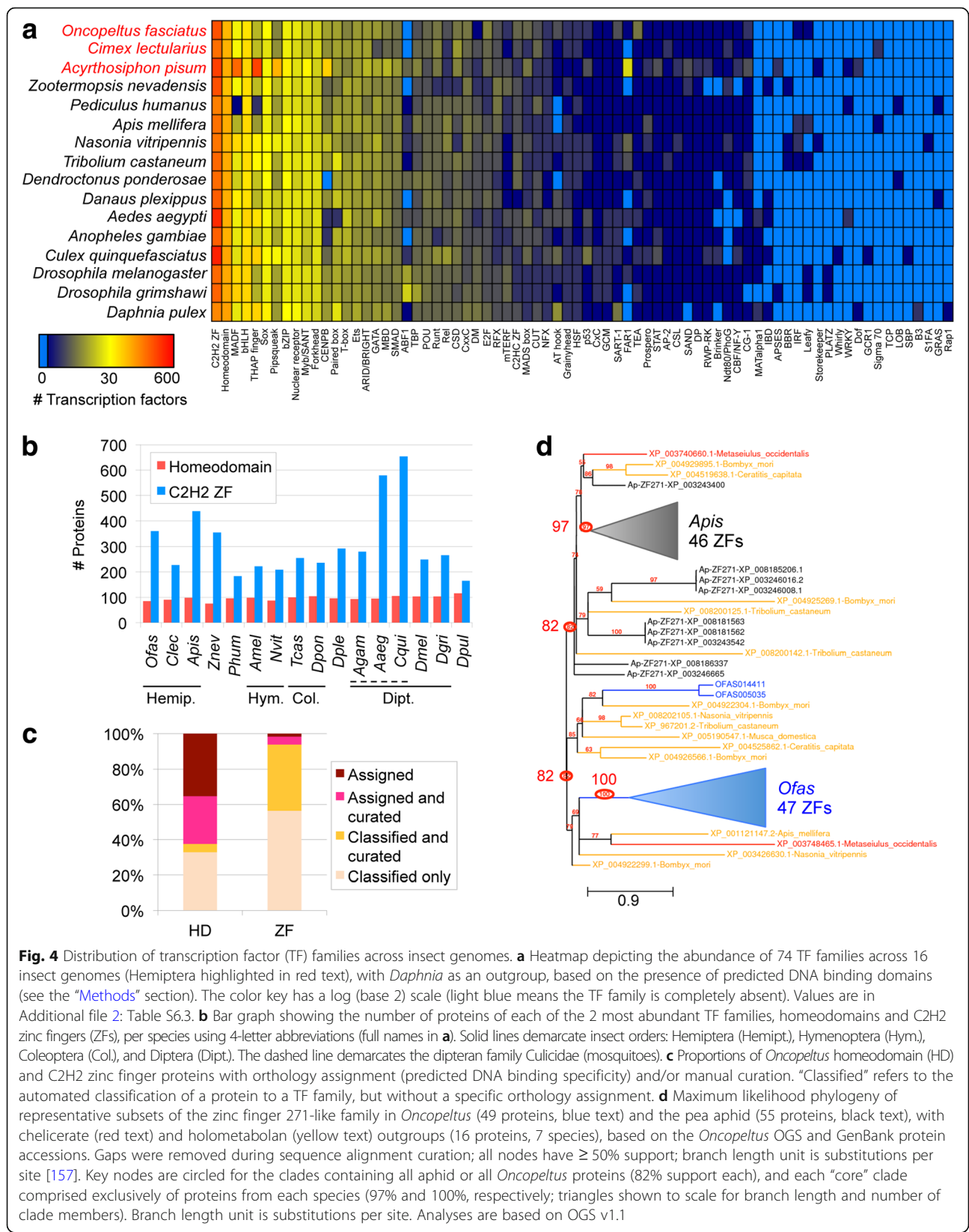
### Independent expansions of C2H2 zinc fingers within the Hemiptera

Unlike homeodomain proteins, C2H2 zinc finger (C2H2-ZF) repertoires are prominent for their large family size and variability throughout the animal kingdom [63], and this is further supported by our current analysis in insects. With  $> 350$  C2H2-ZFs, *Oncopeltus*, the pea aphid, termite, and some mosquito species have  $1.5\times$  more members than the insect median (Fig. 4b). This is nearly half of all *Oncopeltus* TFs. While the expansion in mosquitoes could have a single origin within the Culicinae, the distribution in the Hemiptera, where *Cimex* has only 227 C2H2-ZFs, suggests that independent expansions occurred in *Oncopeltus* and the pea aphid. Prior to the sequencing of other hemipteran genomes, the pea aphid’s large C2H2-ZF repertoire was attributed to the expansion of a novel subfamily, APEZ, also referred to as zinc finger 271-like [43].

In fact, manual curation in *Oncopeltus* confirms the presence of a subfamily with similar characteristics to APEZ (Fig. 4c: yellow fraction). In *Oncopeltus*, we find  $> 115$  proteins of the ZF271 class that are characterized by numerous tandem repeats of the C2H2-ZF domain and its penta-peptide linker, with 3–45 repeats per protein.

Intriguingly, we find evidence for ongoing evolutionary diversification of this subfamily. A number of *Oncopeltus* ZF271-like genes occur in tandem clusters of 4–8 genes—suggesting recent duplication events. Yet, clustered genes differ in gene structure (number and size of exons), and we identified a number of probable ZF271-like pseudogenes whose open reading frames have become disrupted—consistent with high turnover. *Oncopeltus* ZF271-like proteins also differ in the sequence and length of the zinc finger domains among themselves and compared to aphid proteins (WebLogo analysis, [64]), similar to zinc finger array shuffling seen in humans [65]. Furthermore, whole-protein phylogenetic analysis supports independent, rapid expansions in the pea aphid and *Oncopeltus* (Fig. 4d).

Clustered zinc finger gene expansion has long been recognized in mammals, with evidence for strong positive selection to increase both the number and diversity of zinc finger domains per protein as well as the total number of proteins [66]. This was initially found to reflect an arms race dynamic of co-evolution between selfish transposable elements and the C2H2-ZF proteins



**Fig. 4** Distribution of transcription factor (TF) families across insect genomes. **a** Heatmap depicting the abundance of 74 TF families across 16 insect genomes (Hemiptera highlighted in red text), with *Daphnia* as an outgroup, based on the presence of predicted DNA binding domains (see the “Methods” section). The color key has a log (base 2) scale (light blue means the TF family is completely absent). Values are in Additional file 2: Table S6.3. **b** Bar graph showing the number of proteins of each of the 2 most abundant TF families, homeodomains and C2H2 zinc fingers (ZFs), per species using 4-letter abbreviations (full names in **a**). Solid lines demarcate insect orders: Hemiptera (Hemipt), Hymenoptera (Hym), Coleoptera (Col), and Diptera (Dipt). The dashed line demarcates the dipteran family Culicidae (mosquitoes). **c** Proportions of *Oncopeltus* homeodomain (HD) and C2H2 zinc finger proteins with orthology assignment (predicted DNA binding specificity) and/or manual curation. “Classified” refers to the automated classification of a protein to a TF family, but without a specific orthology assignment. **d** Maximum likelihood phylogeny of representative subsets of the zinc finger 271-like family in *Oncopeltus* (49 proteins, blue text) and the pea aphid (55 proteins, black text), with chelicerate (red text) and holometabolan (yellow text) outgroups (16 proteins, 7 species), based on the *Oncopeltus* OGS and GenBank protein accessions. Gaps were removed during sequence alignment curation; all nodes have  $\geq 50\%$  support; branch length unit is substitutions per site [157]. Key nodes are circled for the clades containing all aphid or all *Oncopeltus* proteins (82% support each), and each “core” clade comprised exclusively of proteins from each species (97% and 100%, respectively; triangles shown to scale for branch length and number of clade members). Branch length unit is substitutions per site. Analyses are based on OGS v1.1

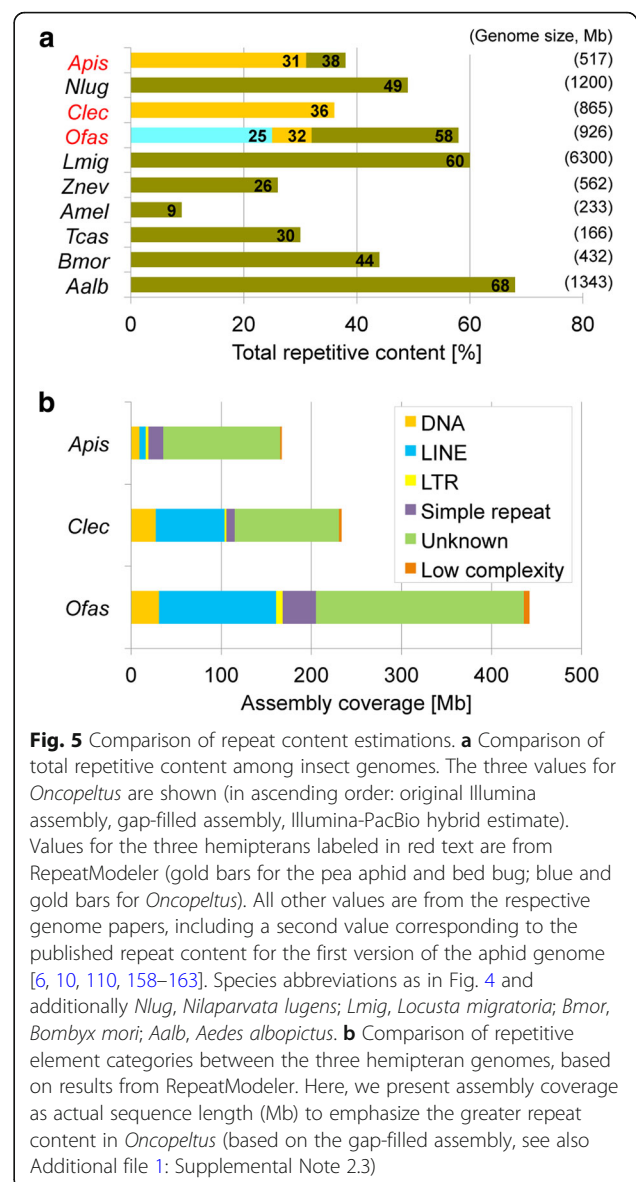
that would repress them [67]. In vertebrates, these C2H2-ZF proteins bind to the promoters of transposable elements via their zinc finger arrays and use their Krüppel-associated box (KRAB) domain to bind the chromatin-remodeling co-repressor KAP-1, which in turn recruits methyltransferases and deacetylases that silence the targeted promoter [68].

Insects do not have a direct ortholog of vertebrate KAP-1 (Additional file 1: Supplemental Note 5.4.d), and neither the aphid nor *Oncopeltus* ZF271-like subfamilies possess a KRAB domain or any other domain besides the zinc finger arrays. However, close molecular outgroups to this ZF271-like subfamily include the developmental repressor Krüppel [69] and the insulator protein CTCF [70] (data not shown). Like these outgroups, the *Oncopeltus* ZF271-like genes are strongly expressed: 98% have expression support, with 86% expressed in at least three different life history stages (Fig. 2b). Thus, the insect ZF271-like proteins may also play prominent roles in repressive DNA binding. Indeed, we find evidence for a functional methylation system in *Oncopeltus* (Additional file 1: Supplemental Note 5.4.c), like the pea aphid, which would provide a means of gene silencing by chromatin remodeling, albeit via mediators other than KAP-1.

However, an arms race model need not be the selective pressure that favors insect ZF271-like family expansions. Recent analyses in vertebrates identified sophisticated, additional regulatory potential by C2H2-ZF proteins, building upon original transposable element binding for new, lineage-specific and even positive gene regulation roles [65, 71, 72]. Moreover, although *Cimex* has half as many long terminal repeat (LTR) repetitive elements as *Oncopeltus* and the pea aphid, overall, we do not find a correlation between relative or absolute repetitive content and ZF271-like family expansion within the Hemiptera (see the next section).

### Proportional repeat content across hemipterans

With the aim of reducing assembly fragmentation and to obtain a better picture of repeat content, we performed low-coverage, long-read PacBio sequencing in *Oncopeltus* (Additional file 1: Supplemental Note 2.3). Using PacBio reads in a gap-filling assay on the Illumina assembly raised the total detected repetitive content from 25 to 32%, while repeat estimations based on simultaneous assessment of Illumina and PacBio reads nearly doubled this value to 58%. As expected, the capacity to identify repeats is strongly dependent on the assembly quality and sequencing technology, with the *Oncopeltus* repetitive content underrepresented in the current (Illumina-only) assembly. Furthermore, as increasing genome size compounds the challenge of assembling repeats, the repeat content of the current assembly is lower than in species with smaller genome sizes (Fig. 5a,



with the sole exception of the honey bee), and we therefore used our gap-filled dataset as a more accurate basis for further comparisons.

To support direct comparisons among hemipterans, we also performed our RepeatModeler analysis on the bed bug and pea aphid assemblies. Repeats comprised 36% and 31% of the respective assemblies, similar to the gap-filled value of 32% in *Oncopeltus*. Nevertheless, given the smaller sizes of these species' assemblies—651 Mb in the bed bug and 542 Mb in the pea aphid—the absolute repeat content is much higher in *Oncopeltus* (Fig. 5b). Excluding unknown repeats, the most abundant transposable elements in *Oncopeltus* are LINE retrotransposons, covering 10% of the assembly (Additional file 2: Table S2.5). This is also the case in the bed bug (12%), while in the pea aphid DNA transposons with terminal inverted repeats

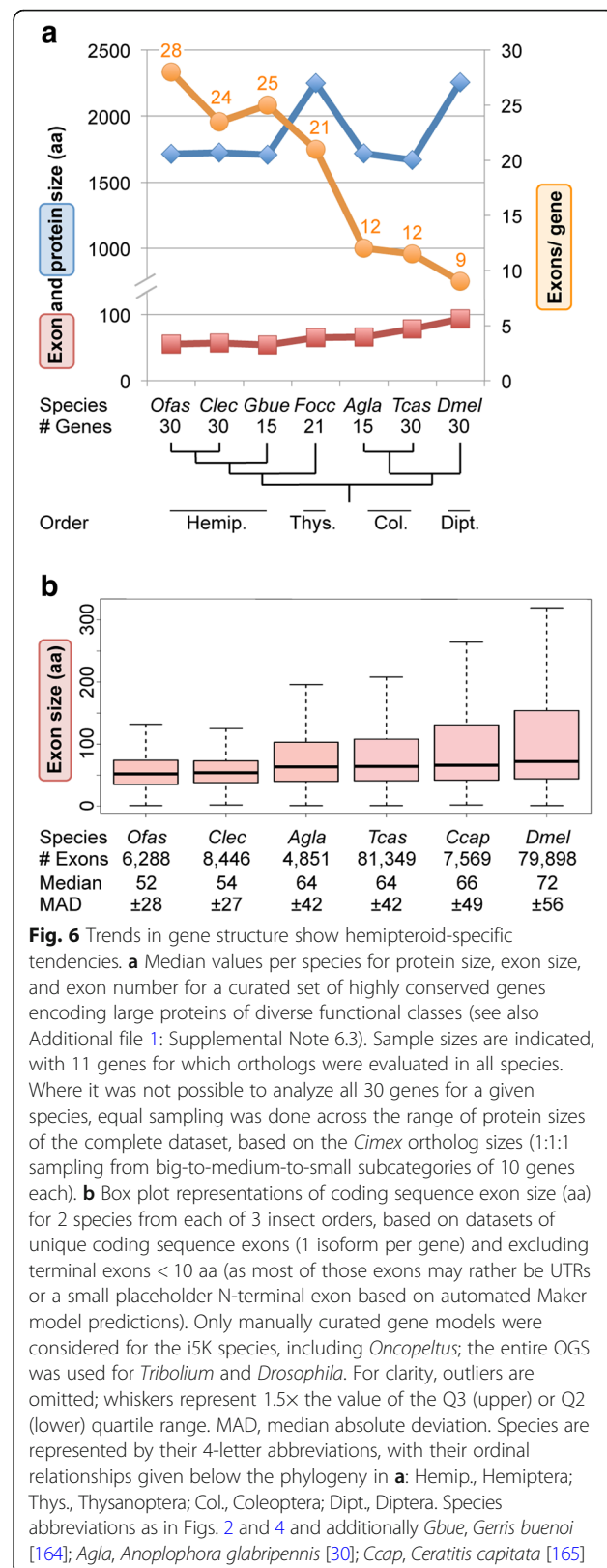
(TIRs) are the most abundant (2% of the assembly identified here and 4% reported from manual curation in the pea aphid genome paper, [6]). Across species, the remaining repeat categories appear to grow proportionally with assembly size, except for simple repeats, which were the category with the largest relative increase in size after gap filling in *Oncopeltus* (Additional file 1: Supplemental Note 2.3). However, given the mix of data types (Illumina [12] and Sanger [6]), these patterns should be treated as hypotheses for future testing.

### Lineage- and genome size-related trends in insect gene structure

Both our manual curation work and BUSCO analyses highlighted the fact that *Oncopeltus* genes are often comprised of many, small exons. We thus undertook a comparative analysis to determine whether this is a general feature to be considered for structural annotation of hemipteran genomes. We find that both lineage and genome size can serve as predictors of gene structure.

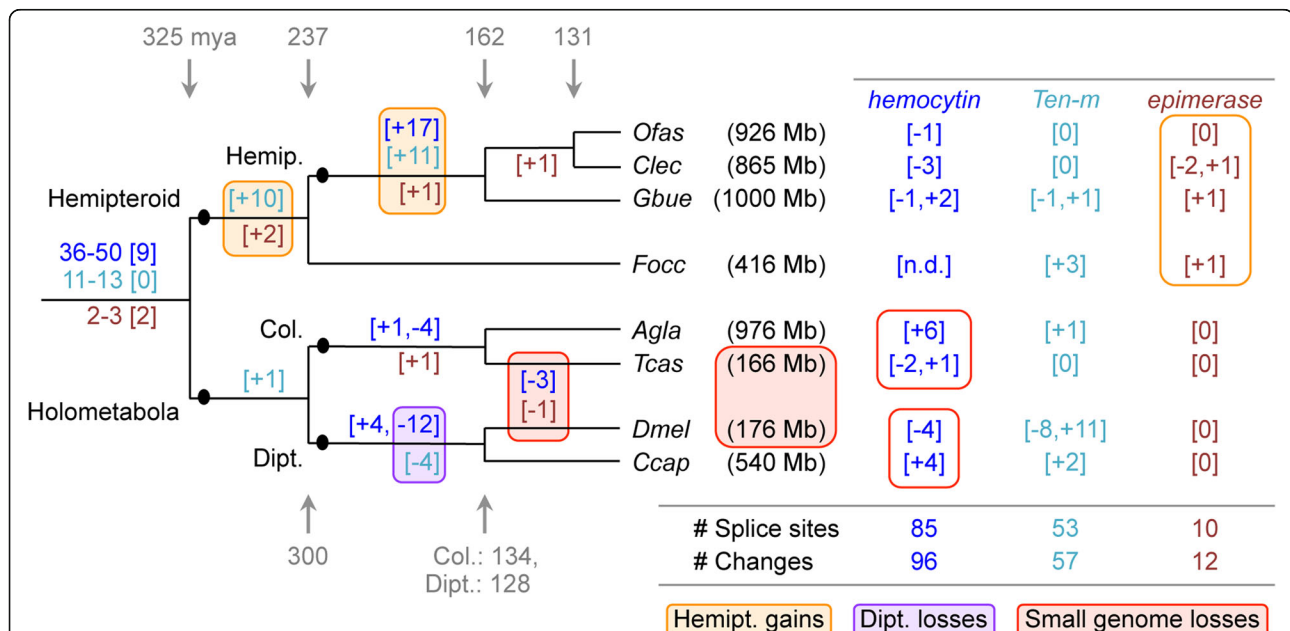
Firstly, we created a high-quality dataset of 30 functionally diverse, large genes whose manual curation could reasonably ensure complete gene models across 7 species from 4 insect orders (Fig. 6a, Additional file 1: Supplemental Note 6.3). Most species encode the same total number of amino acids for these conserved proteins, with the thrips *Frankliniella occidentalis* and the fruit fly being notable exceptions with larger proteins (Fig. 6a: blue plot line). However, the means of encoding this information differs between lineages, with hemipteroid orthologs comprised of twice as many exons as their holometabolous counterparts (Fig. 6a: orange plot line). Thus, there is an inverse correlation between exon number and exon size (Fig. 6a: orange vs. red plot lines). This analysis corroborates and extends previous probabilistic estimates of intron density, where the pea aphid as a sole hemipteran representative had the highest intron density of 10 insect species [73].

To test these trends, we next expanded our analysis to all manually curated exons in two species from each of three orders (Hemiptera, Coleoptera, Diptera). Here, we expect that curated exon sizes are accurate, without the need to assume that the entire gene models are complete. This large dataset corroborates our original findings, with bugs having small exons while both the median and Q3 quartile reflect larger exon sizes in beetles and flies (Fig. 6b). Notably, the median and median absolute deviation are highly similar between species pairs within the Hemiptera and Coleoptera. Meanwhile, the exon metrics within the Diptera support large protein sizes as a drosophilid-specific, rather than dipteran-wide, feature.



Does the high exon count in the Hemiptera reflect an ancient, conserved increase at the base of this lineage or ongoing remodeling of gene structure with high turnover? To assess the exact nature of evolutionary changes, we annotated intron positions within multiple sequence alignments of selected proteins and plotted gains and losses onto the phylogeny, providing a total sample of 165 evolutionary changes at 148 discrete splice sites (Fig. 7, see also Additional file 1: Supplemental Note 6.3 for gene selection and method). These data reveal several major correlates with intron gain or loss. The bases of both the hemipteroid and hemipteran radiations show the largest gains, while most losses occur in the dipteran lineage (Fig. 7: orange and purple shading, respectively). Furthermore, we find progressive gains across hemipteroid nodes, and it is only in this lineage that we additionally find species-specific splice changes for the highly conserved *epimerase* gene (Fig. 7: orange outline). Thus, we find evidence for both ancient intron gain and ongoing gene structure remodeling in this lineage.

Surprisingly, both *hemocytin* and *epimerase*—our exemplar genes with many (up to 74) and few exons (3–8 per species), respectively—show independent losses of the same splice sites in *Drosophila* and *Tribolium*. One feature these species share is a genome size 2.4–6.0× smaller than in the other species examined here (Fig. 7: red shading). Pairwise comparisons within orders also support this trend, as the beetle and fly species with larger genomes exhibit species-specific gains compared to intron loss in their sister taxa (Fig. 7: red outlines). Thus, genome size seems to positively correlate with intron number. However, lineage is a stronger predictor of gene structure: the coleopteran and dipteran species pairs have highly similar exon size metrics despite differences in genome size (Fig. 6b). A global computational analysis over longer evolutionary distances also supports a link between genome size and intron number in arthropods, although chelicerates and insects may experience different rates of evolutionary change in these features [74]. It will be interesting to see if the correlation with genome size is borne out in other invertebrate taxa.



**Fig. 7** Splice site evolution correlates with both lineage and genome size. Splice site changes are shown for *hemocytin* (blue text), *Tenascin major* (*Ten-m*, turquoise text), and *UDP-galactose 4'-epimerase* (brown text), mapped onto a species tree of eight insects. Patterns of splice site evolution were inferred based on the most parsimonious changes that could generate the given pattern within a protein sequence alignment of all orthologs (see also Additional file 1: Supplemental Note 6.3 for methodology and data sources). If inferred gains or losses were equally parsimonious, we remained agnostic and present a range for the ancestral number of splice sites present at the base of the tree, where the bracketed number indicates how many ancestral positions are still retained in all species. Along each lineage, subsequent changes are indicated in brackets, with the sign indicating gains (+) or losses (-). Values shown to the right are species-specific changes. The values shown between the *D. melanogaster* and *T. castaneum* lineages denote changes that have occurred independently in both species. Colored boxes highlight the largest sources of change, as indicated in the legend. Species are represented by their four-letter abbreviations (as in Fig. 6), and estimated genome sizes are indicated parenthetically (measured size [12, 30, 162, 165, 166]; draft assembly size: GenBank Genome IDs 14741 and 17730). Divergence times are shown in gray and given in millions of years [3]. Abbreviations as in Figs. 4 and 6, and also: Hemipt., hemipteroid assemblage (including *F. occidentalis*); n.d., no data

The selective pressures and mechanisms of intron gain in the Hemiptera will be a challenge to uncover. While median exon size (Fig. 6b) could reflect species-specific nucleosome sizes [75, 76], this does not explain why only the Hemiptera seldom exceed this (Fig. 6b: Q3 quartile). Given the gaps in draft genome assemblies, we remain cautious about interpreting (large) intron lengths but note that many hemipteran introns are too small to have harbored a functional transposase gene (e.g., median intron size of 429 bp,  $n = 69$  introns in *hemocytin* in *Cimex*). Such small introns could be consistent with the proliferation of non-autonomous short interspersed nuclear elements (SINEs). However, characterization of such highly divergent non-coding elements would require curated SINE libraries for insects, comparable to those generated for vertebrates and plants [75, 76]. Meanwhile, it appears that hemipteran open reading frames  $\geq 160$  bp are generally prevented by numerous in-frame stop codons just after the donor splice site. Most stop codons are encoded by the triplet TAA in both *Oncopeltus* and *Cimex* (data not shown), although these species' genomes are not particularly AT rich (Table 1).

Even if introns are small, having gene loci comprised of numerous introns and exons adds to the cost of gene expression in terms of both transcription duration and mRNA processing. One could argue that a gene like *hemocytin*, which encodes a clotting agent, would require rapid expression in the case of wounding—a common occurrence in adult *Cimex* females due to the traumatic insemination method of reproduction [12]. Thus, as our molecular understanding of comparative insect and particularly hemipteran biology deepens, we will need to increasingly consider how life history traits are manifest in genomic signatures at the structural level (e.g., Figs. 5, 6, and 7), as well as in terms of protein repertoires (Figs. 3 and 4).

#### Expansion after a novel lateral gene transfer event in phytophagous bugs

In addition to the need for cuticle repair, traumatic insemination may be responsible for the numerous lateral gene transfer (LGT) events predicted in the bed bug [12]. In contrast, the same pipeline analyses [77] followed by manual curation predicted very few LGTs in *Oncopeltus*, which lacks this unusual mating behavior. Here, we have identified 11 strong LGT candidates, and we confirmed the incorporation of bacterial DNA into the milkweed bug genome for all 5 candidates chosen for empirical testing (Additional file 2: Table S2.4). Curiously, we find several LGTs potentially involved in bacterial or plant cell wall metabolism that were acquired from different bacterial sources at different times during hemipteran lineage evolution, including 2 distinct LGTs

that are unique to *Oncopeltus* and implicated in the synthesis of peptidoglycan, a bacterial cell wall constituent (Additional file 1: Supplemental Note 2.2).

Conversely, two other validated LGT candidates are implicated in cell wall degradation. We find two strongly expressed, paralogous copies in *Oncopeltus* of a probable bacterial-origin gene encoding an endo-1,4-beta-mannosidase enzyme (MAN4, EC 3.2.1.78). Inspection of genome assemblies and protein accessions reveals that this LGT event occurred after the infraorder Pentatomomorpha, including the stink bug *Halyomorpha halys*, diverged from other hemipterans, including the bed bug (Fig. 8a). Independent duplications then led to the two copies in *Oncopeltus* and an astonishing nine tandem copies in *Halyomorpha* (Fig. 8b, Additional file 1: Figure S2.6). Since the original LGT event, the *mannosidase* genes have gained introns that are unique to each species and to subsets of paralogs (Fig. 8c). Thus, the “domestication” [78] of *mannosidase* homologs as multi-exonic genes further illustrates the hemipteran penchant for intron introduction and maintenance of small exons. The retention and subsequent expansion of these genes imply their positive selection, consistent with the phytophagous diet of these species. It is tempting to speculate that copy number proliferation in the stink bug correlates with the breadth of its diet, as this agricultural pest feeds on a number of different tissues in a range of host plants [79].

#### Cuticle development, structure, and warning pigmentation

The distinctive cuticle of *Oncopeltus* is produced through the combined action of genes that encode structural and pigmentation proteins, and the gene products that regulate their secretion at each life stage. Furthermore, the milkweed bug has been a powerful model for endocrine studies of hemimetabolous molting and metamorphosis since the 1960s [22, 80–83]. Therefore, we next focused on the presence and function of genes involved in these processes.

Molting is triggered by the release of ecdysteroids, steroid hormones that are synthesized from cholesterol by cytochrome P450 enzymes of the Halloween family [84], and we were able to identify these in the *Oncopeltus* genome (Additional file 1: Supplemental Notes 5.2.b, 5.3.b). From the ecdysone response cascade defined in *Drosophila* [85], we identified *Oncopeltus* orthologs of both early- and late-acting factors, including ecdysteroid hormones and their receptors. It will be interesting to see if the same regulatory relationships are conserved in the context of hemimetabolous molting in *Oncopeltus*. For example, *E75A* is required for reactivation of ecdysteroid production during the molt cycle in *Drosophila* larvae [86] and likely operates similarly in *Oncopeltus*, since *Of-E75A* RNAi



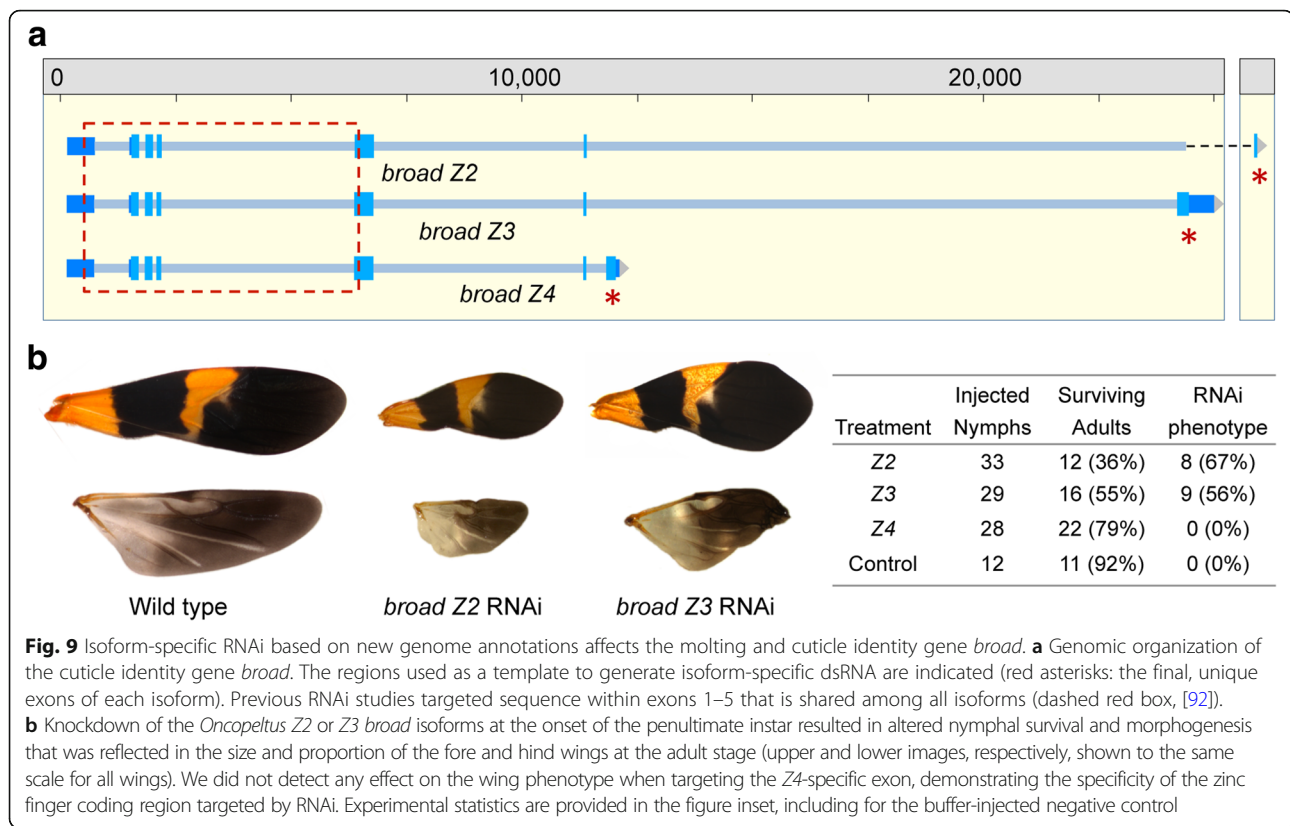


Table S5.12), we found a slight expansion in the *Oncopeltus* CPF family (Additional file 1: Figure S5.14). For cuticle production, similar to the bed bug and the Asian longhorned beetle [12, 30], we identified a single *chitin synthase* gene with conserved alternative splice isoforms, which suggests that *chitin synthase 2* is a duplication specific to only certain beetle and fly lineages within the Holometabola [93].

A major characteristic of the milkweed bug is the distinctive red-orange and black aposematic (warning) coloration within the cuticle and epidermis that deters predators (e.g., Figs. 1 and 9, [20, 21]). For black coloration, the melanin synthesis pathway known from holometabolous insects (e.g., [94, 95]) is conserved at the sequence (Additional file 1: Figure S5.15) and functional [96, 97] level in *Oncopeltus*, supporting conservation in hemimetabolous lineages as well. In contrast, production of the primary warning coloration, pteridine red erythropterin [98], has not been as extensively studied and remains an open avenue for hemimetabolous research. Pterin pigments are synthesized from GTP through a series of enzymatic reactions [99]. Thus far in *Oncopeltus*, we could identify orthologs of *punch*, which encodes a GTP cyclohydrolase [100], and *sepia*, which is required for the synthesis of the red eye pigment drospterin [101]. The bright red color of *Oncopeltus* eggs may in part reflect chemical protection transmitted parentally [102]. Thus, further identification of pigmentation genes will provide fitness indicators for

maternal contributions to developmental success under natural conditions (i.e., the presence of egg predators).

### Chemoreception and metabolism in relation to feeding biology

Aposematic pigmentation advertises the fact that toxins in the milkweed diet are incorporated into the insects themselves, a metabolic feat that was independently acquired in *Oncopeltus* and the monarch butterfly (*Danaus plexippus*), which shares this food source and body coloration [36, 103]. Moreover, given the fundamental differences between phytophagous, mucivorous, and hematophagous diets, we investigated to what extent differences in feeding ecology across hemipterans are represented in their chemoreceptor and metabolic enzyme repertoires.

Insects must smell and taste their environment to locate and identify food, mates, oviposition sites, and other essential cues. Perception of the enormous diversity of environmental chemicals is primarily mediated by the odorant (OR), gustatory (GR), and ionotropic (IR) families of chemoreceptors, which each encode tens to hundreds of distinct proteins [104–107]. Chemoreceptor family size appears to correlate with feeding ecology. *Oncopeltus* retains a moderate complement of each family, while species with derived fluid nutrition diets (sap or blood) have relatively depauperate OR and GR families (Table 2, Additional file 1: Supplemental Note



**Table 2** Numbers of chemoreceptor genes/proteins per family in selected insect species. In some cases, the number of proteins is higher than the number of genes due to an unusual form of alternative splicing, which is particularly notable for the *Oncopeltus* GRs. Data are shown for four Hemiptera as well as *Drosophila melanogaster*, the body louse *Pediculus humanus*, and the termite *Zootermopsis nevadensis* [11, 12, 104, 108–110, 168]

Species	Odorant	Gustatory	Ionotropic
<i>Oncopeltus fasciatus</i> <sup>1</sup>	120/121	115/169	37/37
<i>Cimex lectularius</i> <sup>1,2</sup>	48/49	24/36	30/30
<i>Rhodnius prolixus</i> <sup>1,2</sup>	116/116	28/30	33/33
<i>Acyrtosiphon pisum</i> <sup>3</sup>	79/79	77/77	19/19
<i>Pediculus humanus</i> <sup>2</sup>	12/13	6/8	14/14
<i>Zootermopsis nevadensis</i>	70/70	87/90	150/150
<i>Drosophila melanogaster</i>	60/62	60/68	65/65

<sup>1</sup>Hemiptera: Heteroptera

<sup>2</sup>Independent acquisitions of hematophagy [16]

<sup>3</sup>Hemiptera, phloem feeding

5.3.f, Additional file 3). In detail, a few conserved orthologs such as the OrCo protein and a fructose receptor are found across species, but other subfamilies are lineage specific. *Oncopeltus* and *Acyrtosiphon* retain a set of sugar receptors that was lost independently in the blood-feeding bugs (*Rhodnius* [11], *Cimex* [12]) and body louse (*Pediculus* [108]). Conversely, *Oncopeltus* and *Cimex* retain a set of candidate carbon dioxide receptors, a gene lineage lost from *Rhodnius*, *Acyrtosiphon*, and *Pediculus* [11, 12, 109], but which is similar to a GR subfamily expansion in the more distantly related hemimetabolous termite (Isoptera [110]). Comparable numbers of IRs occur across the Heteroptera. In addition to a conserved set of orthologs primarily involved in sensing temperature and certain acids and amines, *Oncopeltus* has a minor expansion of IRs distantly related to those involved in taste in *Drosophila*. The major expansions in each insect lineage are the

candidate “bitter” GRs ([111], Additional file 1: Supplemental Note 5.3.f and Figure S5.19). In summary, *Oncopeltus* exhibits moderate expansion of specific subfamilies likely to be involved in host plant recognition, consistent with it being a preferentially specialist feeder with a potentially patchy food source [21, 112].

As host plant recognition is only the first step, we further explored whether novel features of the *Oncopeltus* gene set are directly associated with its diet. We therefore used the CycADS annotation pipeline [113] to reconstruct the *Oncopeltus* metabolic network. The resulting BioCyc metabolism database for *Oncopeltus* (OncaCyc) was then compared with those for 26 other insect species ([114], <http://arthropodacyc.cycadsys.org/>), including 3 other hemipterans: the pea aphid, the green peach aphid, and the kissing bug (Tables 3 and 4). For a global metabolism analysis, we detected the presence of 1085 Enzyme Commission (EC) annotated reactions with at least 1 protein in the *Oncopeltus* genome (Additional file 1: Supplemental Note 6.4, Additional file 2: Table S6.10). Among these, 10 enzyme classes (represented by 17 genes) are unique and 17 are missing when compared to the other insects (Table 4, Additional file 2: Table S6.11).

We then specifically compared amino acid metabolism in the four hemipterans representing the three different diets. Eight enzymes are present only in *Oncopeltus* (Table 4), including the arginase that degrades arginine (Arg) into urea and ornithine, a precursor of proline (Pro). Given this difference, we extended our analysis to assess species’ repertoires for the entire urea cycle (Fig. 10a, Additional file 2: Table S6.13). *Oncopeltus* and six other species can degrade Arg but cannot synthesize it (Fig. 10b). Only the other three hemipterans can neither synthesize nor degrade Arg via this cycle (Fig. 10c), while most species have an almost complete cycle (Fig. 10d). This suggests that the ability to synthesize Arg was lost at the base of the Hemiptera, with subsequent, independent loss of Arg degradation capacity in the aphid and *Rhodnius* lineages. Retention of Arg

**Table 3** Hemipteran ArthropodaCyc database summaries. Overview statistics for the newly created database for *Oncopeltus fasciatus* (*Ofas*) in comparison with public databases for *Rhodnius prolixus* (*Rpro*), *Acyrtosiphon pisum* (*Apis*), and *Myzus persicae* (*Mper*) available from [114]. Based on OGS v1.1

Species ID	<i>Ofas</i>	<i>Rpro</i>	<i>Apis</i>	<i>Mper</i>	<i>Mper</i>
Gene set ID	OGS v1.1	RproC1.1 (Built on RproC1 assembly)	OGS v2.1b (Built on Acyr_2.0 assembly)	Clone G006 v1.0	Clone O v1.0
CycADS database ID	OncaCyc	RhoprCyc	AcypiCyc v2.1b	Myzpe_G006 Cyc	Myzpe_O Cyc
Total mRNA <sup>1</sup>	19,673	15,437	36,195	24,814	24,770
Pathways	294	312	307	319	306
Enzymatic reactions	2192	2366	2339	2384	2354
Polypeptides	19,820	15,471	36,228	24,849	24,805
Enzymes	3050	2660	5087	4646	4453
Compounds	1506	1665	1637	1603	1655

<sup>1</sup>In the BioCyc databases, all splice variants are counted in the summary tables for genes

**Table 4** Hemipteran ArthropodaCyc annotations of metabolic genes. Taxonomic abbreviations are as in Table 3

	<i>Ofas</i>	<i>Rpro</i>	<i>Apis</i>	<i>Mper</i>
Global metabolism				
EC <sup>1</sup> present in the genome	1085	1241	1288	1222
EC unique to this genome <sup>2</sup>	10	13	23	5
EC missing only in this genome <sup>2</sup>	17 <sup>4</sup>	8	2	6
Amino acid metabolism (KEGG)				
EC present in the genome	169	188	195	185
EC unique to this genome <sup>2</sup>	2	1	6	1
EC missing only in this genome <sup>2</sup>	5	2	0	2
EC unique to this genome <sup>3</sup>	8	10	12	8
EC missing only in this genome <sup>3</sup>	14	5	0	2

<sup>1</sup>“EC” refers to the number of proteins, as represented by their unique numerical designations within the Enzyme Commission (EC) classification system for enzymes and their catalytic reactions

<sup>2</sup>In comparison with all other insects from ArthropodaCyc

<sup>3</sup>in comparison among the four hemipterans

<sup>4</sup>Includes three EC categories added in OGS v1.2 (see also Additional file 2: Table S6.11)

degradation in *Oncopeltus* might be linked to the milkweed seed food source, as most seeds are very rich in Arg [115], and Arg is indeed among the metabolites detected in *Oncopeltus* [116]. However, the monarch butterfly is one of only a handful of species that retains the complete Arg pathway (Fig. 10: blue text). Despite a shared food source, these species may therefore differ in their overall Arg requirements or—in light of a possible group benefit of *Oncopeltus* aggregation during feeding ([21]; e.g., Fig. 1b)—in their efficiency of Arg uptake.

Other enzymes are also present only in the milkweed bug compared to the other examined hemipterans (Additional file 2: Table S6.12). Like other insects [114], *Oncopeltus* retains the ability to degrade tyrosine (Tyr). This pathway was uniquely lost in the aphids, where this essential amino acid is jointly synthesized—and consumed—by the aphid host and its endosymbiotic bacteria [6, 7, 17, 117]. Conversely, a gain specific to the milkweed bug lineage was the duplication of the Na<sup>+</sup>/K<sup>+</sup> ATPase alpha subunits whose amino acid substitutions confer resistance to milkweed cardenolides [36, 118]. In the *Oncopeltus* genome, we find support for the recent nature of these duplications: the genes encoding subunits ATPα1B and ATPα1C occur as a tandem duplication, notably on a scaffold that also harbors one of the clustered ZF271-like gene expansions (see above).

## Conclusions

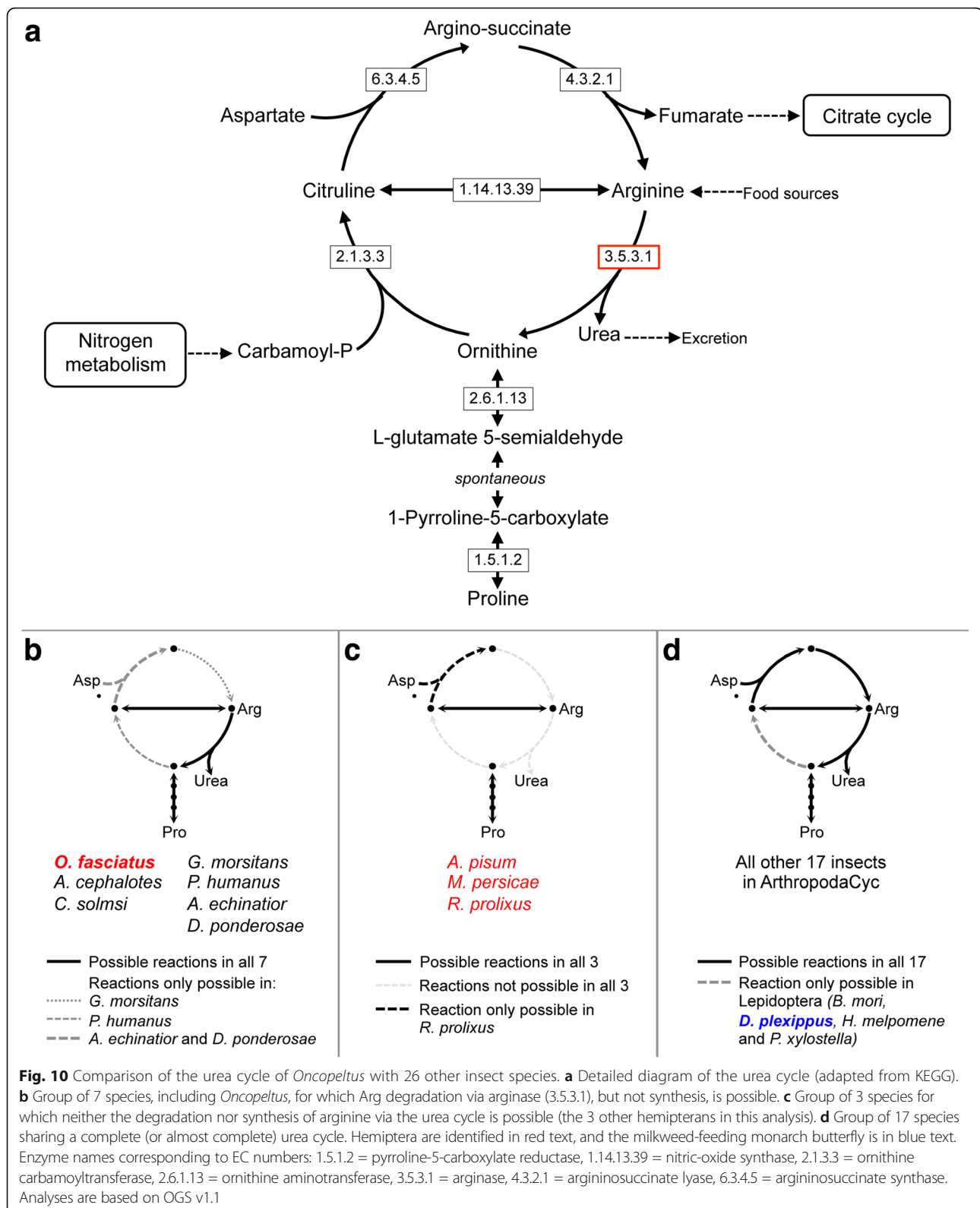
The integrated genomic and transcriptomic resources presented here for the milkweed bug *Oncopeltus fasciatus* (Figs. 2 and 5) underpin a number of insights into evolutionary and developmental genomics. Our

macroevolutionary comparisons across insect orders, now extended to the hemimetabolous Hemiptera, reveal unexpected patterns of molecular evolution. We also show how hemipteran feeding ecology and suites of related biological characters are reflected in the genome.

The gene structure trends we identified, with lineage predominating over genome size as a predictor and with many intron gains in the hemipteroid lineage (Figs. 6 and 7), offer initial parameters and hypotheses for the Hemiptera, Coleoptera, and Diptera. Such ordinal-level parameters can be evaluated against new species’ data and also inform customized pipelines for automated gene model predictions. At the same time, it will be interesting to explore the ramifications of hemipteroid intron gains, as there are few documented lineages with episodic intron gain [76]. For example, possessing more, small exons may provide greater scope to generate protein modularity via isoforms based on alternative exon usage [119]. Furthermore, with the larger genome sizes and lower gene densities of hemipteroids compared to the well-studied Hymenoptera, it remains open whether hemipteroid gene and intron size may also correlate with recombination rates [120].

Our analyses also highlight new directions for future experimental research, building on *Oncopeltus*’s long-standing history as a laboratory model and its active research community in the modern molecular genetics era (e.g., Fig. 9, [25–27]). Functional testing will clarify the roles of genes we have identified as unique to the Hemiptera, including those implicated in chemical protection, bacterial and plant cell wall metabolism, or encoding wholly novel proteins (Figs. 3 and 8, see also Additional file 1: Supplemental Note 2.2). Meanwhile, the prominent and species-specific expansions specifically of ZF271-like zinc fingers (Fig. 4), combined with the absence of the co-repressor KAP-1 in insects, argues for investigation into alternative interaction partners, which could clarify the nature of these zinc fingers’ regulatory roles and their binding targets.

One key output of this study is the generation of a metabolism database for *Oncopeltus*, contributing to the ArthropodaCyc collection (Table 3). In addition to comparisons with other species (Fig. 10), this database can also serve as a future reference for studies that use *Oncopeltus* as an ecotoxicology model (e.g., [121]). While we have primarily focused on feeding ecology in terms of broad comparisons between phytophagy and fluid feeding, *Oncopeltus* is also poised to support future work on nuances among phytophagous species. Despite its milkweed diet in the wild, the lab strain of *Oncopeltus* has long been adapted to feed on sunflower seeds, demonstrating a latent capacity for more generalist phytophagy [112]. This potential may also be reflected in a larger gustatory receptor repertoire than would be expected for



an obligate specialist feeder (Table 2). Thus, *Oncopeltus* can serve as a reference species for promiscuously phytophagous pests such as the stink bug. Finally, we have identified a number of key genes implicated in life history trade-offs in *Oncopeltus*, for traits such as cardenolide tolerance, pigmentation, and plasticity in reproduction under environmental variation. The genome data thus represent an important tool to elucidate the proximate mechanisms of fundamental aspects of life history evolution in both the laboratory and nature.

## Methods

(More information is available in Additional file 1: Supplemental Notes.)

### Milkweed bug strain, rearing, and DNA/RNA extraction

The milkweed bug *Oncopeltus fasciatus* (Dallas), Carolina Biological Supply strain (Burlington, NC, USA), was maintained in a laboratory colony under standard husbandry conditions (sunflower seed and water diet, 25 °C, 12:12 light-dark photoperiod). Voucher specimens for an adult female (record # ZFMK-TIS-26324) and adult male (record # ZFMK-TIS-26325) have been preserved in ethanol and deposited in the Biobank of the Centre for Molecular Biodiversity Research, Zoological Research Museum Alexander Koenig, Bonn, Germany (<https://www.zfmk.de/en/biobank>).

Genomic DNA was isolated from individual, dissected adults using the Blood & Cell Culture DNA Midi Kit (G/100) (Qiagen Inc., Valencia, CA, USA). Total RNA was isolated from individual, dissected adults and from pooled, mixed-instar nymphs with TRIzol Reagent (Invitrogen/Thermo Fisher Scientific, Waltham, MA, USA). Dissection improved the accessibility of muscle tissue by disrupting the exoskeleton, and the gut material was removed.

### Genome size calculations (flow cytometry, *k*-mer estimation)

Genome size estimations were obtained by flow cytometry with Hare and Johnston's protocol [122]. Four to five females and males each from the Carolina Biological Supply lab strain and a wild strain (collected from Athens, GA, USA; GPS coordinates: 33° 56' 52.8216" N, 83° 22' 38.3484" W) were measured (see also Additional file 1: Supplemental Note 2.1.a). At the bioinformatic level, we attempted to estimate the genome size by *k*-mer spectrum distribution analysis for a range of *k* = 15 to 34 counted with Jellyfish 2.1.4 [123] and bbmap [124], graphing these counts against the frequency of occurrence of *k*-mers (depth) and calculating genome size based on the coverage at the peak of the distribution (Additional file 1: Supplemental Note 2.1.b).

### Genome sequencing, assembly, annotation, and official gene set overview

Library preparation, sequencing, assembly, and automatic gene annotation were conducted at the Baylor College of Medicine Human Genome Sequencing Center (as in [12, 30]). About 1.1 billion 100-bp paired-end reads generated on an Illumina HiSeq2000s machine were assembled using ALLPATHS-LG [125], from two paired-end (PE) and two mate pair (MP) libraries specifically designed for this algorithm (Additional file 1: Supplemental Note 1). Three libraries were sequenced from an individual adult male (180- and 500-bp PE, 3-kb MP), with the fourth from an individual adult female (8–10-kb MP). The final assembly, "Ofas\_1.0" (see metrics in Table 1), has been deposited in GenBank (assembly accession GCA\_000696205.1).

Automated annotation of protein-coding genes was performed using a Maker 2.0 annotation pipeline [126] tuned specifically for arthropods (Additional file 1: Supplemental Note 3). These gene predictions were used as the starting point for manual curation via the Apollo v.1.0.4 web browser interface [127], and automatic and manual curations were compiled to generate the OGS (see also Additional file 1: Supplemental Note 4). The current version of the gene set, OGS v1.2, is deposited at GenBank as an "annotation-only" update to the Whole Genome Shotgun project (accession JHQO00000000). Here, we describe version JHQO02000000. The annotations can be downloaded from NCBI's ftp site, [ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/696/205/GCA\\_000696205.1\\_Ofas\\_1.0/](ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/696/205/GCA_000696205.1_Ofas_1.0/). The annotations are also available through the i5K Workspace@NAL [128], [https://i5k.nsl.usda.gov/data/Arthropoda/oncfas-\(Oncopeltus\\_fasciatus\)/Ofas\\_1.0/2.Official%20or%20Primary%20Gene%20Set/GCA\\_000696205.1\\_Ofas\\_1.0/](https://i5k.nsl.usda.gov/data/Arthropoda/oncfas-(Oncopeltus_fasciatus)/Ofas_1.0/2.Official%20or%20Primary%20Gene%20Set/GCA_000696205.1_Ofas_1.0/).

Databases of the genome assembly (definitive Illumina-only: Table 1, Additional file 1: Supplemental Note 1.3; provisional hybrid Illumina-PacBio: see below, Additional file 1: Supplemental Note 2.3), Maker automatic gene predictions (Additional file 1: Supplemental Note 3), and OGS v1.1 (Table 1, Additional file 1: Supplemental Note 4) are available through the i5K Workspace@NAL, and the Ag Data Commons data access system of the US Department of Agriculture's (USDA) National Agricultural Library as individual citable databases [129–132].

### Repeat content analysis

Repetitive regions were identified in the *Oncopeltus* genome assembly with RepeatModeler Open-1.0.8 [133] based on a species-specific repeat library generated de novo with RECON [134], RepeatScout [135], and Tandem Repeats Finder [136]. Then, RepeatMasker Open-4.0 [137] was used to mask the repeat sequences based on the

RepeatModeler library. Given the fragmented nature of the assembly, we attempted to fill and close the assembly gaps by sequencing additional material, generating long reads with single molecule real-time sequencing on a Pac-Bio RS II machine (estimated coverage of 8x). Gap filling on the Illumina assembly scaffolds was performed with PBJelly version 13.10.22, and the resulting assembly [132] was used for repeat content estimation and comparison with *Cimex lectularius* and *Acyrtosiphon pisum* (see also Additional file 1: Supplemental Note 2.3).

### Transcriptome resources

Total RNA from three distinct life history samples (pooled, mixed-instar nymphs; an adult male; an adult female) was also sequenced on an Illumina HiSeq2000s machine, producing a total of 72 million 100-bp paired-end reads (Additional file 1: Supplemental Note 1.3, Additional file 2: Table S1.1; GenBank Bioproject: PRJNA275739). These expression data were used to support the generation of the OGS at different stages of the project: as input for the evidence-guided automated annotation with Maker 2.0 (Additional file 1: Supplemental Note 3), as expression evidence tracks in the Apollo browser to support the community curation of the OGS, and, once assembled into a de novo transcriptome, as a point of comparison for quality control of the OGS.

The raw RNA-seq reads were pre-processed by filtering out low-quality bases (phred score < 30) and Truseq adapters with Trimmomatic-0.30. Further filtering removed ribosomal and mitochondrial RNA sequences with Bowtie 2 [138], based on a custom library built with all hemipteran ribosomal and mitochondrial RNA accessions from NCBI as of 7th February 2014 (6069 accessions). The pooled, filtered reads were mapped to the genome assembly with Tophat2-PE on CyVerse [139]. A second set of RNA-seq reads from an earlier study (“published adult” dataset, [36]) was also filtered and mapped in the same fashion, and both datasets were loaded into the *Oncopeltus* Apollo instance as evidence tracks (under the track names “pooled RNA-seq - cleaned reads” and “RNA-seq raw PE reads Andolfatto et al”, respectively).

Additionally, a de novo transcriptome was generated from our filtered RNA-seq reads (pooled from all three samples prepared in this study) using Trinity [140] and TransDecoder [141] with default parameters. This transcriptome is referred to as “i5K,” to distinguish it from a previously published maternal and early embryonic transcriptome for *Oncopeltus* (referred to as “454”, [35]). Both the i5K and 454 transcriptomes were mapped to the genome assembly with GMAP v. 2014-05-15 on CyVerse. These datasets were also loaded into the Apollo browser as evidence tracks to assist in manual curation.

### Life history stage-specific and sex-specific expression analyses in hemipteroids

Transcript expression of the OGS v1.1 genes was estimated by running RSEM2 [142] on the filtered RNA-seq datasets for the three i5K postembryonic stages against the OGS v1.1 cDNA dataset. Transcript expression was then based on the transcripts per million (TPM) value. The TPM values were processed by adding a value of 1 (to avoid zeros) and then performing a log<sub>2</sub> transformation. The number of expressed genes per RNA-seq library was compared for TPM cutoffs of > 1, > 0.5, and > 0.25. A > 0.25 cutoff was chosen, which reduced the number of expressed genes by 6.6% compared to a preliminary analysis based on a simple cutoff of  $\geq 10$  mapped reads per transcript, while the other TPM cutoffs were deemed too restrictive (reducing the expressed gene set by > 10%). This analysis was also applied to the “published adult” dataset [36]. To include the embryonic stages in the comparison, transcripts from the 454 transcriptome were used as blastn queries against the OGS v1.1 cDNA dataset (cutoff *e* value < 10<sup>-5</sup>). The results from all datasets were converted to a binary format to generate Venn diagrams (Fig. 2b).

Statistically significant sex-specific and developmental stage-specific gene enrichment was determined from RNA-seq datasets according to published methods [143, 144], with modifications. Data from *Oncopeltus* (see the previous methods section, Bioproject: PRJNA275739) were compared between stages and pairwise with the hemipterans *Cimex lectularius*, PRJNA275741; *Acyrtosiphon pisum*, PRJNA209321; and *Pachypsylla venusta*, PRJNA275248; as well as with the hemipteroid *Frankliniella occidentalis* (Thysanoptera), PRJNA203209 (see also Fig. 2c, Additional file 1: Supplemental Note 2.4).

### Protein orthology assessments via OrthoDB and BUSCO analyses

These analyses follow previously described approaches and with the current database and pipeline versions [1, 42, 44, 145], see Additional file 1: Supplemental Note 6.1 for further details.

### Global transcription factor identification

Likely transcription factors (TFs) were identified by scanning the amino acid sequences of predicted protein-coding genes for putative DNA binding domains (DBDs), and when possible, the DNA binding specificity of each TF was predicted using established procedures [57]. Briefly, all protein sequences were scanned for putative DBDs using the 81 Pfam [146] models listed in Weirauch and Hughes [147] and the HMMER tool [148], with the recommended detection thresholds of Per-sequence Eval < 0.01 and Per-domain conditional Eval < 0.01. Each

protein was classified into a family based on its DBDs and their order in the protein sequence (e.g., bZIPx1, AP2x2, Homeodomain+Pou). The resulting DBD amino acid sequences were then aligned within each family using Clustal Omega [149], with default settings. For protein pairs with multiple DBDs, each DBD was aligned separately. From these alignments, the sequence identity was calculated for all DBD sequence pairs (i.e., the percent of amino acid residues that are identical across all positions in the alignment). Using previously established sequence identity thresholds for each family [57], the predicted DNA binding specificities were mapped by simple transfer. For example, the DBD of OFAS001246-RA is 98% identical to the *Drosophila melanogaster* Bric a Brac 1 (Bab1) protein. Since the DNA binding specificity of Bab1 has already been experimentally determined, and the cutoff for the Pipsqueak family TFs is 85%, we can infer that OFAS001246-RA will have the same binding specificity as *Drosophila* Bab1.

### RNA interference

Double-stranded RNA (dsRNA) was designed to target the final, unique exon of the *broad* isoforms *Z2*, *Z3*, and *Z4*. A portion of the coding sequence for the zinc finger region from these exons (179 bp, 206 bp, and 216 bp, respectively) was cloned into a plasmid vector and used as a template for in vitro RNA synthesis, using the gene-specific primer pairs: *Of-Z2\_fwd*: 5'-ATGTGGCAGACAAGCATGCT-3', *Of-Z2\_rev*: 5'-CTAAAATTTGACATCAGTAGGC-3'; *Of-Z3\_fwd*: 5'-CCTTCTCCTGTTACTACTCAC-3', *Of-Z3\_rev*: 5'-TTATATGGGCGGCTGTCCAA-3'; and *Of-Z4\_fwd*: 5'-AACACTGACCTTGGTTACACA-3', *Of-Z4\_rev*: 5'-TAGGTGGAGGATTGCTAAAATT-3'. Two separate transcription reactions (one for each strand) were performed using the Ambion MEGAscript kit (Ambion, Austin, TX, USA). The reactions were purified by phenol/chloroform extraction followed by precipitation as described in the MEGAscript protocol. The separate strands were re-annealed in a thermocycler as described previously [27]. Nymphs were injected with a Hamilton syringe fitted with a 32-gauge needle as described [54]. The concentration of *Of-Z2*, *Of-Z3*, and *Of-Z4* dsRNA was 740 ng/μl, 1400 ng/μl, and 1200 ng/μl, respectively. All nymphs were injected within 8 h of the molt to the fourth (penultimate juvenile) instar ( $n \geq 12$  per treatment, see Fig. 9). Fore- and hindwings were then dissected from adults and photographed at the same scale as wings from wild type, uninjected controls.

### CycADS annotation and OncfaCyc database generation

We used the Cyc Annotation Database System (CycADS [113]), an automated annotation management system, to integrate protein annotations from different sources into a Cyc metabolic network reconstruction that was integrated into the ArthropodaCyc database.

Using our CycADS pipeline, *Oncopeltus fasciatus* proteins from the official gene set OGS v1.1 were annotated using different methods—including KAAS [150], PRIAM [151], Blast2GO [152, 153], and InterProScan with several approaches [154]—to obtain EC and GO numbers. All annotation information data were collected in the CycADS SQL database and automatically extracted to generate appropriate input files to build or update BioCyc databases [155] using the Pathway Tools software [156]. The OncfaCyc database, representing the metabolic protein-coding genes of *Oncopeltus*, was thus generated and is now included in the ArthropodaCyc database, a collection of arthropod metabolic network databases ([114], <http://arthropodacyc.cycadsys.org/>).

### Additional files

**Additional file 1:** Supplementary notes, figures, and small tables. (PDF 6142 kb)

**Additional file 2:** Large supporting tables. (XLSX 2222 kb)

**Additional file 3:** Chemoreceptor sequences in FASTA format. (FASTA 128 kb)

### Acknowledgements

We thank Dorith Rotenberg (Kansas State University, currently North Carolina State University, USA) and Michael Sparks (Agricultural Research Service, United States Department of Agriculture, USA) for generously making available the unpublished genome assemblies of the fellow hemipteroid i5K species *Frankliniella occidentalis* and *Halyomorpha halys*, respectively, for use in specific analyses presented here. Similarly, we thank Hans Kelstrup and Lynn Riddiford (Janelia Farm Research Campus, HHMI, USA) for sharing unpublished data on *Of-E75A* RNAi. We thank George Coupland (Max Planck Institute for Plant Breeding Research, Cologne, Germany) as well as Lisa Czaja, Kurt Steuber, and Bruno Huettel (Max Planck Genome Centre Cologne, Germany) for conducting the PacBio sequencing and providing support with data handling. We also thank Oliver Niehuis (Albert Ludwig University, Freiburg, Germany) and Alexander Klassmann (University of Cologne, Germany) for discussions on *k*-mer and gene structure analyses, respectively, Sarah Kingan (University of Rochester, USA) for assistance with LGT phylogenies, as well as Jeanne Wilbrandt (Zoologisches Forschungsmuseum Alexander Koenig, Bonn, Germany) for comments on the manuscript.

### Funding

Funding for genome sequencing, assembly and automated annotation was provided by the National Institutes of Health (NIH) grant U54 HG003273 (NHGRI) to RAG. The i5K pilot project (<https://www.hgsc.bcm.edu/arthropods>) assisted in sequencing of the *Oncopeltus fasciatus* genome. We also acknowledge funding for the project from German Research Foundation (DFG) grants PA 2044/1-1 and SFB 680 project A12 to KAP. Support for specific analyses was provided by the Swiss National Science Foundation with grant 31003A\_143936 to EMZ and PP00P3\_170664 to RMW; the European Research Council grant ERC-CoG #616346 to AK; DFG grant SFB 680 project A1 to SiR; the National Science Foundation with grant US NSF DEB1257053 to JHW; NIH grant R01GM113230 (NIGMS) to LP; and by NIH grants 5R01GM080203 (NIGMS) and 5R01HG004483 (NHGRI) and by the Director, Office of Science, Office of Basic Energy Sciences, U.S. Department of Energy, Contract No. DE-AC02-05CH11231 to MCMT.

### Availability of data and materials

All sequence data are publically available at the NCBI, bioproject number PRJNA229125 and in the USDA Ag Data Commons data access system [133]. In addition, assembled scaffolds, gene models, and a browser are available at the National Agricultural Library ([130–133], <https://i5k.nal.usda.gov/>)

*Oncopeltus fasciatus*). The OncoCyc metabolism database is available within the ArthropodaCyc collection (<http://arthropodacyc.cycadsys.org/>).

#### Authors' contributions

KAP and StR conceived the project. KAP managed and coordinated the project. KAP and SK provided the specimens for sequencing and performed the DNA and RNA extractions. StR, SD, SLL, HC, HVD, HD, YH, JQ, SCM, DSTH, KCW, DMM, and RAG constructed the libraries and performed the sequencing. StR, SCM, and DSTH performed the genome assembly and automated gene prediction. IMVJ, JSJ, and PJM analyzed the genome size. IMVJ, VK, PH, and KAP contributed to the repetitive content analyses. AD, RR, JHW, KAP, and SK performed the bacterial scaffold detection and LGT analyses. MCMT developed the Apollo software. KAP, IMVJ, MCMT, CPC, C-YL, and MFP implemented the Apollo-based manual curation. KAP, IMVJ, JBB, DE, YS, HMR, DA, CGCJ, BMIV, EJD, CSB, C-CC, Y-TC, ADC, AGC, AJJC, PKD, EMD, CGE, MF, NG, TH, Y-MH, ECJ, TEJ, JWJ, AK, ML, MRL, H-LL, YL, SRP, LP, MLP, PNR, RR-P, SiR, LS, MES, JS, ES, JNS, OT, LT, MVDZ, SV, and AJR participated in the manual curation and contributed to the supplemental notes. IMVJ, KAP, DSTH, M-JMC, CPC, C-YL, and MFP performed the curation quality control and generated the OGS. IMVJ, KAP, CJH, and JBB generated the de novo transcriptomes and performed the life history stage expression analyses. RMW, PI, KAP, and EMZ performed the orthology and phylogenomic analyses. MTW, KAP, IMVJ, PH, and BMIV performed the transcription factor analyses. EJD conducted the analyses of DNA methylation. KAP, PH, and RJS contributed to the comparative analyses of gene structure. DE conducted the RNAi experiments. SC, PB-P, GF, and NP generated and performed the comparative analyses on the OncoCyc database. KAP, IMVJ, JBB, DE, YS, SC, HMR, and MTW wrote the manuscript. KAP, IMVJ, JBB, DE, YS, SC, HMR, MFP, RMW, PI, MTW, StR, PJM, and AK edited the manuscript. IMVJ and KAP organized the supplementary materials. All authors approved the final manuscript.

#### Ethics approval and consent to participate

Not applicable.

#### Competing interests

The authors declare that they have no competing interests.

#### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

#### Author details

<sup>1</sup>Institute for Zoology: Developmental Biology, University of Cologne, Zùlpicher Str. 47b, 50674 Cologne, Germany. <sup>2</sup>School of Life Sciences, University of Warwick, Gibbet Hill Campus, Coventry CV4 7AL, UK. <sup>3</sup>Department of Biological Sciences, University of Cincinnati, Cincinnati, OH 45221, USA. <sup>4</sup>Department of Biochemistry and Cell Biology and Center for Developmental Genetics, Stony Brook University, Stony Brook, NY 11794, USA. <sup>5</sup>Department of Biological Sciences, Wellesley College, 106 Central St., Wellesley, MA 02481, USA. <sup>6</sup>Univ Lyon, INSA-Lyon, INRA, BF21, UMR0203, F-69621 Villeurbanne, France. <sup>7</sup>Present address: LSTM, Laboratoire des Symbioses Tropicales et Méditerranéennes, INRA, IRD, CIRAD, SupAgro, University of Montpellier, Montpellier, France. <sup>8</sup>Department of Entomology, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA. <sup>9</sup>National Agricultural Library, Beltsville, MD 20705, USA. <sup>10</sup>Department of Genetic Medicine and Development and Swiss Institute of Bioinformatics, University of Geneva, 1211 Geneva, Switzerland. <sup>11</sup>Present address: Department of Ecology and Evolution, University of Lausanne, 1015 Lausanne, Switzerland. <sup>12</sup>Center for Autoimmune Genomics and Etiology, Division of Biomedical Informatics, and Division of Developmental Biology, Cincinnati Children's Hospital Medical Center, Department of Pediatrics, College of Medicine, University of Cincinnati, Cincinnati, OH 45229, USA. <sup>13</sup>Human Genome Sequencing Center, Department of Human and Molecular Genetics, Baylor College of Medicine, One Baylor Plaza, Houston, TX 77030, USA. <sup>14</sup>Present address: Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA 98195, USA. <sup>15</sup>Present address: Howard Hughes Medical Institute, University of Washington, Seattle, WA 98195, USA. <sup>16</sup>Department of Biology, University of Rochester, Rochester, NY 14627, USA. <sup>17</sup>Institute of Biology, Leiden University, Sylviusweg 72, 2333 BE Leiden, Netherlands. <sup>18</sup>Max Planck Institute for Chemical Ecology, Hans-Knöll Strasse 8, 07745 Jena, Germany. <sup>19</sup>Department of Biochemistry and Genomics

Aotearoa, University of Otago, Dunedin 9054, New Zealand. <sup>20</sup>School of Biology, Faculty of Biological Sciences, University of Leeds, Leeds LS2 9JT, UK. <sup>21</sup>Institut de Génomique Fonctionnelle de Lyon, Université de Lyon, Université Claude Bernard Lyon 1, CNRS UMR 5242, École Normale Supérieure de Lyon, 46 Allée d'Italie, 69364 Lyon, France. <sup>22</sup>Department of Ecology, Evolution and Behavior, The Alexander Silberman Institute of Life Sciences, The Hebrew University of Jerusalem, Edmond J. Safra Campus, Givat Ram, 91904 Jerusalem, Israel. <sup>23</sup>Department of Entomology/Institute of Biotechnology, College of Bioresources and Agriculture, National Taiwan University, Taipei, Taiwan. <sup>24</sup>Present address: School of Life Sciences, Rochester Institute of Technology, Rochester, NY 14623, USA. <sup>25</sup>Department of Organismic and Evolutionary Biology, Harvard University, 26 Oxford Street, Cambridge, MA 02138, USA. <sup>26</sup>Department of Molecular and Cellular Biology, Harvard University, 26 Oxford Street, Cambridge, MA 02138, USA. <sup>27</sup>Department of Biological Sciences, Wayne State University, Detroit, MI 48202, USA. <sup>28</sup>Institute for Genetics, University of Cologne, Zùlpicher Straße 47a, 50674 Cologne, Germany. <sup>29</sup>Department of Entomology, Texas A&M University, College Station, TX 77843, USA. <sup>30</sup>CECAD, University of Cologne, Cologne, Germany. <sup>31</sup>Department of Entomology and Program in Molecular & Cell Biology, University of Maryland, College Park, MD 20742, USA. <sup>32</sup>Department of Entomology, University of Georgia, 120 Cedar St., Athens, GA 30602, USA. <sup>33</sup>Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA. <sup>34</sup>Department of Entomology, College of Agriculture, Food and Environment, University of Kentucky, Lexington, KY 40546, USA. <sup>35</sup>Department of Biology, University of Hawai'i at Mānoa, Honolulu, HI 96822, USA. <sup>36</sup>Present address: Department of Evolutionary Genetics, Max-Planck-Institut für Evolutionsbiologie, August-Thienemann-Straße 2, 24306 Plön, Germany. <sup>37</sup>Present address: Earthworks Institute, 185 Caroline Street, Rochester, NY 14620, USA. <sup>38</sup>Centro de Bioinvestigaciones, Universidad Nacional del Noroeste de Buenos Aires, Pergamino, Argentina. <sup>39</sup>Present address: Department of Biotechnology, Central University of Rajasthan (CURAJ), NH-8, Bandarsindri, Ajmer 305801, India. <sup>40</sup>Argelander-Institut für Astronomie, Universität Bonn, Auf dem Hügel 71, 53121 Bonn, Germany. <sup>41</sup>Present address: Department of Zoology, University of Cambridge, Cambridge CB2 3DT, UK. <sup>42</sup>Centro Regional de Estudios Genómicos, Facultad de Ciencias Exactas, Universidad Nacional de La Plata, La Plata, Argentina. <sup>43</sup>Present address: Department of Physiology, Anatomy and Genetics and Centre for Neural Circuits and Behaviour, University of Oxford, Oxford OX1 3SR, UK. <sup>44</sup>Present address: E. A. Milne Centre for Astrophysics, Department of Physics and Mathematics, University of Hull, Hull HU6 7RX, UK. <sup>45</sup>Research Center for Developmental Biology and Regenerative Medicine, National Taiwan University, Taipei, Taiwan.

Received: 2 November 2018 Accepted: 21 February 2019

Published online: 02 April 2019

#### References

- Zdobnov EM, Tegenfeldt F, Kuznetsov D, Waterhouse RM, Simão FA, Ioannidis P, Seppey M, Loetscher A, Kriventseva EV. OrthoDB v9.1: cataloging evolutionary and functional annotations for animal, fungal, plant, archaeal, bacterial and viral orthologs. *Nucleic Acids Res.* 2017;45:D744–9.
- Huang DY, Bechly G, Nel P, Engel MS, Prokop J, Azar D, Cai CY, van de Kamp T, Staniczek AH, Garrouste R, et al. New fossil insect order Permopsocida elucidates major radiation and evolution of suction feeding in hemimetabolous insects (Hexapoda: Acercaria). *Sci Rep.* 2016;6:23004.
- Misof B, Liu S, Meusemann K, Peters RS, Donath A, Mayer C, Frandsen PB, Ware J, Flouri T, Beutel RG, et al. Phylogenomics resolves the timing and pattern of insect evolution. *Science.* 2014;346:763–7.
- Grimaldi D, Engel MS. *Evolution of the insects.* Cambridge: Cambridge University Press; 2005.
- Panfilio KA, Angelini DR. By land, air, and sea: hemipteran diversity through the genomic lens. *Curr Opin Insect Sci.* 2018;25:106–15.
- The International Aphid Genomics Consortium. Genome sequence of the pea aphid *Acyrtosiphon pisum*. *PLoS Biol.* 2010;8:e1000313.
- Mathers TC, Chen Y, Kaithakottil G, Legeai F, Mugford ST, Baa-Puyoulet P, Bretaudeau A, Clavijo B, Colella S, Collin O, et al. Rapid transcriptional plasticity of duplicated gene clusters enables a clonally reproducing aphid to colonise diverse plant species. *Genome Biol.* 2017;18:27.
- Wenger JA, Cassone BJ, Legeai F, Johnston JS, Bansal R, Yates AD, Coates BS, Pavinato VA, Michel A. Whole genome sequence of the soybean aphid, *Aphis glycines*. *Insect Biochem Mol Biol.* 2017; in press.

9. Sloan DB, Nakabachi A, Richards S, Qu J, Murali SC, Gibbs RA, Moran NA. Parallel histories of horizontal gene transfer facilitated extreme reduction of endosymbiont genomes in sap-feeding insects. *Mol Biol Evol.* 2014;31:857–71.
10. Xue J, Zhou X, Zhang C-X, Yu L-L, Fan H-W, Wang Z, Xu H-J, Xi Y, Zhu Z-R, Zhou W-W, et al. Genomes of the rice pest brown planthopper and its endosymbionts reveal complex complementary contributions for host adaptation. *Genome Biol.* 2014;15:521.
11. Mesquita RD, Vionette-Amaral RJ, Lowenberger C, Rivera-Pomar R, Monteiro FA, Minx P, Spieth J, Carvalho AB, Panzera F, Lawson D, et al. Genome of *Rhodnius prolixus*, an insect vector of Chagas disease, reveals unique adaptations to hematophagy and parasite infection. *Proc Natl Acad Sci U S A.* 2015;112:14936–41.
12. Benoit JB, Adelman ZN, Reinhardt K, Dolan A, Poelchau M, Jennings EC, Szuter EM, Hagan RW, Gujar H, Shukla JN, et al. Unique features of a global human ectoparasite identified through sequencing of the bed bug genome. *Nat Commun.* 2016;7:10165.
13. Rosenfeld JA, Reeves D, Brugler MR, Narechania A, Simon S, Durrett R, Foox J, Shianna K, Schatz MC, Gandara J, et al. Genome assembly and geospatial phylogenomics of the bed bug *Cimex lectularius*. *Nat Commun.* 2016;7:10164.
14. Sparks ME, Shelby KS, Kuhar D, Gundersen-Rindal DE. Transcriptome of the invasive brown marmorated stink bug, *Halyomorpha halys* (Stal) (Heteroptera: Pentatomidae). *PLoS One.* 2014;9:e111646.
15. Ioannidis P, Lu Y, Kumar N, Creasy T, Daugherty S, Chibucos MC, Orvis J, Shetty A, Ott S, Flowers M, et al. Rapid transcriptome sequencing of an invasive pest, the brown marmorated stink bug, *Halyomorpha halys*. *BMC Genomics.* 2014;15:738.
16. Li H, Leavengood JM Jr, Chapman EG, Burkhardt D, Song F, Jiang P, Liu J, Zhou X, Cai W. Mitochondrial phylogenomics of Hemiptera reveals adaptive innovations driving the diversification of true bugs. *Proc Biol Sci.* 2017;284: 20171223.
17. Wilson ACC, Ashton PD, Charles H, Colella S, Febvay G, Jander G, Kushlan PF, Macdonald SJ, Schwartz JF, Thomas GH, Douglas AE. Genomic insight into the amino acid relations of the pea aphid, *Acyrtosiphon pisum*, with its symbiotic bacterium *Buchnera aphidicola*. *Insect Mol Biol.* 2010;19(Suppl 2):249–58.
18. Eichler S, Schaub GA. Development of symbionts in triatomine bugs and the effects of infections with trypanosomatids. *Exp Parasitol.* 2002;100:17–27.
19. Matsuura Y, Kikuchi Y, Hosokawa T, Koga R, Meng X-Y, Kamagata Y, Nikoh N, Fukatsu T. Evolution of symbiotic organs and endosymbionts in lygaeid stinkbugs. *The ISME Journal.* 2012;6:397–409.
20. Berenbaum MR, Miliczky E. Mantids and milkweed bugs - efficacy of aposomatic coloration against invertebrate predators. *Am Midl Nat.* 1984;111:64–8.
21. Burdfield-Steel ER, Shuker DM. The evolutionary ecology of the Lygaeidae. *Ecol Evol.* 2014;4:2278–301.
22. Lawrence PA. Mitosis and the cell cycle in the metamorphic moult of the milkweed bug *Oncopeltus fasciatus*; a radioautographic study. *J Cell Sci.* 1968;3:391–404.
23. Chipman AD. *Oncopeltus fasciatus* as an evo-devo research organism. *Genesis.* 2017;55:e23020.
24. Panfilio KA. Late extraembryonic development and its *zen*-RNAi-induced failure in the milkweed bug *Oncopeltus fasciatus*. *Dev Biol.* 2009;333: 297–311.
25. Panfilio KA, Roth S. Epithelial reorganization events during late extraembryonic development in a hemimetabolous insect. *Dev Biol.* 2010; 340:100–15.
26. Sharma AI, Yanes KO, Jin L, Garvey SL, Taha SM, Suzuki Y. The phenotypic plasticity of developmental modules. *Evodevo.* 2016;7:15.
27. Hughes CL, Kaufman TC. RNAi analysis of *Deformed*, *proboscipedia* and *Sex combs reduced* in the milkweed bug *Oncopeltus fasciatus*: novel roles for Hox genes in the hemipteran head. *Development.* 2000;127:3683–94.
28. Wolfe SL, John B. The organization and ultrastructure of male meiotic chromosomes in *Oncopeltus fasciatus*. *Chromosoma.* 1965;17:85–103.
29. Messthaler H, Traut W. Phases of sex chromosome inactivation in *Oncopeltus fasciatus* and *Pyrrhocoris apterus* (Insecta, Heteroptera). *Caryologia.* 1975;28:501–10.
30. McKenna DD, Scully ED, Pauchet Y, Hoover K, Kirsch R, Geib SM, Mitchell RF, Waterhouse RM, Ahn SJ, Arslan D, et al. Genome of the Asian longhorned beetle (*Anoplophora glabripennis*), a globally significant invasive species, reveals key functional and evolutionary innovations at the beetle-plant interface. *Genome Biol.* 2016;17:227.
31. Simpson JT. Exploring genome characteristics and sequence quality without a reference. *Bioinformatics.* 2014;30:1228–35.
32. Hanrahan SJ, Johnston JS. New genome size estimates of 134 species of arthropods. *Chromosom Res.* 2011;19:809–23.
33. Panfilio KA, Liu PZ, Akam M, Kaufman TC. *Oncopeltus fasciatus zen* is essential for serosal tissue function in katatrepsis. *Dev Biol.* 2006;292:226–43.
34. Tian X, Xie Q, Li M, Gao C, Cui Y, Xi L, Bu W. Phylogeny of pentatomomorph bugs (Hemiptera-Heteroptera:Pentatomomorpha) based on six Hox gene fragments. *Zootaxa.* 2011;2888:57–68.
35. Ewen-Campen B, Shaner N, Panfilio KA, Suzuki Y, Roth S, Extavour CG. The maternal and early embryonic transcriptome of the milkweed bug *Oncopeltus fasciatus*. *BMC Genomics.* 2011;12:61.
36. Zhen Y, Aardema ML, Medina EM, Schumer M, Andolfatto P. Parallel molecular evolution in an herbivore community. *Science.* 2012;337:1634–7.
37. Robertson HM. The insect chemoreceptor superfamily in *Drosophila pseudoobscura*: molecular evolution of ecologically-relevant genes over 25 million years. *J Insect Sci.* 2009;9:18.
38. Robertson HM. Taste: independent origins of chemoreception coding systems? *Curr Biol.* 2001;11:R560–2.
39. Jazwinska A, Rushlow C, Roth S. The role of *brinker* in mediating the graded response to Dpp in early *Drosophila* embryos. *Development.* 1999;126:3323–34.
40. Togawa T, Dunn WA, Emmons AC, Nagao J, Willis JH. Developmental expression patterns of cuticular protein genes with the R&R Consensus from *Anopheles gambiae*. *Insect Biochem Mol Biol.* 2008;38:508–19.
41. Karr TL. Fruit flies and the sperm proteome. *Hum Mol Genet* 2007, 16 Spec No. 2:R124–R133.
42. Waterhouse RM, Seppey M, Simao FA, Manni M, Ioannidis P, Klioutchnikov G, Kriventseva EV, Zdobnov EM. BUSCO applications from quality assessments to gene prediction and phylogenomics. *Mol Biol Evol.* 2017;35: 543–8.
43. Shigenobu S, Bickel RD, Brisson JA, Butts T, Chang CC, Christiaens O, Davis GK, Duncan EJ, Ferrier DE, Iga M, et al. Comprehensive survey of developmental genes in the pea aphid, *Acyrtosiphon pisum*: frequent lineage-specific duplications and losses of developmental genes. *Insect Mol Biol.* 2010;19(Suppl 2):47–62.
44. Kriventseva EV, Tegenfeldt F, Petty TJ, Waterhouse RM, Simão FA, Pozdnyakov IA, Ioannidis P, Zdobnov EM. OrthoDB v8: update of the hierarchical catalog of orthologs and the underlying free software. *Nucl Acids Res.* 2015;43:D250–6.
45. Bansal R, Michel AP. Core RNAi machinery and *Sid1*, a component for systemic RNAi, in the hemipteran insect, *Aphis glycines*. *Int J Mol Sci.* 2013; 14:3786–801.
46. Bao R, Fischer T, Bolognesi R, Brown SJ, Friedrich M. Parallel duplication and partial subfunctionalization of beta-catenin/armadillo during insect evolution. *Mol Biol Evol.* 2012;29:647–62.
47. Sachs L, Chen YT, Drechsler A, Lynch JA, Panfilio KA, Lassig M, Berg J, Roth S. Dynamic BMP signaling polarized by Toll patterns the dorsoventral axis in a hemimetabolous insect. *eLife.* 2015;4:e05502.
48. Armisen D, Refki PN, Crumiere AJ, Viala S, Toubiana W, Khila A. Predator strike shapes antipredator phenotype through new genetic interactions in water striders. *Nat Commun.* 2015;6:8153.
49. Konopova B, Smykal V, Jindra M. Common and distinct roles of juvenile hormone signaling genes in metamorphosis of holometabolous and hemimetabolous insects. *PLoS One.* 2011;6:e28728.
50. Vellichirammal NN, Gupta P, Hall TA, Brisson JA. Ecdysone signaling underlies the pea aphid transgenerational wing polyphenism. *Proc Natl Acad Sci U S A.* 2017;114:1419–23.
51. Wulff JP, Sierra I, Sterkel M, Holtorf M, Van Wielendaele P, Francini F, Broeck JV, Ons S. Orcokinin neuropeptides regulate ecdysis in the hemimetabolous insect *Rhodnius prolixus*. *Insect Biochem Mol Biol.* 2017;81:91–102.
52. Chiu TL, Wen Z, Rupasinghe SG, Schuler MA. Comparative molecular modeling of *Anopheles gambiae* CYP6Z1, a mosquito P450 capable of metabolizing DDT. *Proc Natl Acad Sci U S A.* 2008;105:8855–60.
53. Gong Y, Li T, Feng Y, Liu N. The function of two P450s, CYP9M10 and CYP6AA7, in the permethrin resistance of *Culex quinquefasciatus*. *Sci Rep.* 2017;7:587.
54. Liu PZ, Kaufman TC. *hunchback* is required for suppression of abdominal identity, and for proper germband growth and segmentation in the intermediate germband insect *Oncopeltus fasciatus*. *Development.* 2004;131:1515–27.
55. Schaeper ND, Pechmann M, Damen WGM, Prpic N-M, Wimmer EA. Evolutionary plasticity of *collier* function in head development of diverse arthropods. *Dev Biol.* 2010;344:363–76.



56. Aspiras AC, Smith FW, Angelini DR. Sex-specific gene interactions in the patterning of insect genitalia. *Dev Biol.* 2011;360:369–80.
57. Weirauch MT, Yang A, Albu M, Cote AG, Montenegro-Montero A, Drewe P, Najafabadi HS, Lambert SA, Mann I, Cook K, et al. Determination and inference of eukaryotic transcription factor sequence specificity. *Cell.* 2014;158:1431–43.
58. Peel AD, Telford MJ, Akam M. The evolution of hexapod engrailed-family genes: evidence for conservation and concerted evolution. *Proc Biol Sci.* 2006;273:1733–42.
59. Ben-David J, Chipman AD. Mutual regulatory interactions of the trunk gap genes during blastoderm patterning in the hemipteran *Oncopeltus fasciatus*. *Dev Biol.* 2010;346:140–9.
60. Erezylmaz DF, Kelstrup HC, Riddiford LM. The nuclear receptor E75A has a novel pair-rule-like function in patterning the milkweed bug, *Oncopeltus fasciatus*. *Dev Biol.* 2009;334:300–10.
61. Liu PZ, Kaufman TC. *even-skipped* is not a pair-rule gene but has segmental and gap-like functions in *Oncopeltus fasciatus*, an intermediate germband insect. *Development.* 2005;132:2081–92.
62. Weisbrod A, Cohen M, Chipman AD. Evolution of the insect terminal patterning system—insights from the milkweed bug, *Oncopeltus fasciatus*. *Dev Biol.* 2013;380:125–31.
63. Albertin CB, Simakov O, Mitros T, Wang ZY, Pungor JR, Edsinger-Gonzales E, Brenner S, Ragsdale CW, Rokhsar DS. The octopus genome and the evolution of cephalopod neural and morphological novelties. *Nature.* 2015;524:220–4.
64. Crooks GE, Hon G, Chandonia J-M, Brenner SE. WebLogo: a sequence logo generator. *Genome Res.* 2004;14:1188–90.
65. Najafabadi HS, Mnaimneh S, Schmitges FW, Garton M, Lam KN, Yang A, Albu M, Weirauch MT, Radovani E, Kim PM, et al. C2H2 zinc finger proteins greatly expand the human regulatory lexicon. *Nat Biotechnol.* 2015;33:555–62.
66. Emerson RO, Thomas JH. Adaptive evolution in zinc finger transcription factors. *PLoS Genet.* 2009;5:e1000325.
67. Thomas JH, Schneider S. Coevolution of retroelements and tandem zinc finger genes. *Genome Res.* 2011;21:1800–12.
68. Garcia-Perez JL, Widmann TJ, Adams IR. The impact of transposable elements on mammalian development. *Development.* 2016;143:4101–14.
69. Liu PZ, Kaufman TC. *Krüppel* is a gap gene in the intermediate insect *Oncopeltus fasciatus* and is required for development of both blastoderm and germband-derived segments. *Development.* 2004;131:4567–79.
70. Heger P, Marin B, Bartkuhn M, Schierenberg E, Wiehe T. The chromatin insulator CTCF and the emergence of metazoan diversity. *Proc Natl Acad Sci U S A.* 2012;109:17507–12.
71. Liu H, Chang L-H, Sun Y, Lu X, Stubbs L. Deep vertebrate roots for mammalian zinc finger transcription factor subfamilies. *Genome Biol Evol.* 2014;6:510–25.
72. Imbeault M, Hellebood P-Y, Trono D. KRAB zinc-finger proteins contribute to the evolution of gene regulatory networks. *Nature.* 2017;543:550–4.
73. Csuros M, Rogozin IB, Koonin EV. A detailed history of intron-rich eukaryotic ancestors inferred from a global survey of 100 complete genomes. *PLoS Comput Biol.* 2011;7:e1002150.
74. Hoy MA, Waterhouse RM, Wu K, Estep AS, Ioannidis P, Palmer WJ, Pomerantz AF, Simao FA, Thomas J, Jiggins FM, et al. Genome sequencing of the phytoseiid predatory mite *Metaseiulus occidentalis* reveals completely atomized Hox genes and superdynamic intron evolution. *Genome Biol Evol.* 2016;8:1762–75.
75. Seibt KM, Wenke T, Muders K, Truberg B, Schmidt T. Short interspersed nuclear elements (SINEs) are abundant in Solanaceae and have a family-specific impact on gene structure and genome organization. *Plant J.* 2016;86:268–85.
76. Huff JT, Zilberman D, Roy SW. Mechanism for DNA transposons to generate introns on genomic scales. *Nature.* 2016;538:533–6.
77. Wheeler D, Redding AJ, Werren JH. Characterization of an ancient lepidopteran lateral gene transfer. *PLoS One.* 2012;8:e59262.
78. Da Lage JL, Binder M, Hua-Van A, Janecsek S, Casane D. Gene make-up: rapid and massive intron gains after horizontal transfer of a bacterial alpha-amylase gene to Basidiomycetes. *BMC Evol Biol.* 2013;13:40.
79. Lee DH, Short BD, Joseph SV, Bergh JC, Leskey TC. Review of the biology, ecology, and management of *Halymorphia halsys* (Hemiptera: Pentatomidae) in China, Japan, and the Republic of Korea. *Environ Entomol.* 2013;42:627–41.
80. Lawrence PA. Cellular differentiation and pattern formation during metamorphosis of the milkweed bug *Oncopeltus*. *Dev Biol.* 1969;19:12–40.
81. Riddiford LM. Prevention of metamorphosis by exposure of insect eggs to juvenile hormone analogs. *Science.* 1970;167:287.
82. Willis JH, Lawrence PA. Deferred action of juvenile hormone. *Nature.* 1970;225:81–3.
83. Masner P, Bowers WS, Kalin M, Muhle T. Effect of precocene II on the endocrine regulation of development and reproduction in the bug, *Oncopeltus fasciatus*. *Gen Comp Endocrinol.* 1979;37:156–66.
84. Rewitz K, O'Connor M, Gilbert L. Molecular evolution of the insect Halloween family of cytochrome P450s: phylogeny, gene organization and functional conservation. *Insect Biochem Mol Biol.* 2007;37:741–53.
85. Huet F, Ruiz C, Richards G. Sequential gene activation by ecdysone in *Drosophila melanogaster*: the hierarchical equivalence of early and early late genes. *Development.* 1995;121:1195–204.
86. Bialecki M, Shilton A, Fichtenberg C, Segraves WA, Thummel CS. Loss of the ecdysteroid-inducible E75A orphan nuclear receptor uncouples molting from metamorphosis in *Drosophila*. *Dev Cell.* 2002;3:209–20.
87. Charles JP, Iwema T, Epa VC, Takaki K, Rynes J, Jindra M. Ligand-binding properties of a juvenile hormone receptor, methoprene-tolerant. *Proc Natl Acad Sci U S A.* 2011;108:21128–33.
88. Minakuchi C, Zhou X, Riddiford L. Kruppel homolog 1 (Kr-h1) mediates juvenile hormone action during metamorphosis of *Drosophila melanogaster*. *Mech Dev.* 2008;125:91–105.
89. Minakuchi C, Namiki T, Shinoda T. Kruppel homolog 1, an early juvenile hormone-response gene downstream of methoprene-tolerant, mediates its anti-metamorphic action in the red flour beetle *Tribolium castaneum*. *Dev Biol.* 2009;352:341–50.
90. DiBello PR, Withers DA, Bayer CA, Fristrom JW, Guild GM. The *Drosophila Broad-Complex* encodes a family of related proteins containing zinc fingers. *Genetics.* 1991;129:385–97.
91. Karim F, Guild G, Thummel C. The *Drosophila Broad-Complex* plays a key role in controlling ecdysone-regulated gene expression at the onset of metamorphosis. *Development.* 1993;118:977–88.
92. Erezylmaz DF, Riddiford LM, Truman JW. The pupal specifier broad directs progressive morphogenesis in a direct-developing insect. *Proc Natl Acad Sci U S A.* 2006;103:6925–30.
93. Arakane Y, Hogenkamp DG, Zhu YC, Kramer KJ, Specht CA, Beeman RW, Kanost MR, Muthukrishnan S. Characterization of two chitin synthase genes of the red flour beetle, *Tribolium castaneum*, and alternate exon usage in one of the genes during development. *Insect Biochem Mol Biol.* 2004;34:291–304.
94. True JR. Insect melanism: the molecules matter. *Trends Ecol Evol.* 2003;18:640–7.
95. Zhan SA, Guo QH, Li MH, Li MW, Li JY, Miao XX, Huang YP. Disruption of an N-acetyltransferase gene in the silkworm reveals a novel role in pigmentation. *Development.* 2010;137:4083–90.
96. Liu J, Lemonds TR, Popadic A. The genetic control of aposomatic black pigmentation in hemimetabolous insects: insights from *Oncopeltus fasciatus*. *Evol Dev.* 2014;16:270–7.
97. Liu J, Lemonds TR, Marden JH, Popadic A. A pathway analysis of melanin patterning in a hemimetabolous insect. *Genetics.* 2016;203:403–13.
98. Lawrence PA. Some new mutants of large milkweed bug *Oncopeltus fasciatus* Dall. *Genet Res.* 1970;15:347–50.
99. Morgan ED. *Biosynthesis in insects: advanced edition.* London: Royal Society of Chemistry; 2010.
100. McLean JR, Krishnakumar S, O'Donnell JM. Multiple mRNAs from the *Punch* locus of *Drosophila melanogaster* encode isoforms of GTP cyclohydrolase I with distinct N-terminal domains. *J Biol Chem.* 1993;268:27191–7.
101. Wiederrecht GJ, Paton DR, Brown GM. Enzymatic conversion of Dihydroneopterin triphosphate to the pyrimidodiazepine intermediate involved in the biosynthesis of the Drosoterins in *Drosophila melanogaster*. *J Biol Chem.* 1984;259:2195–200.
102. Newcombe D, Blount JD, Mitchell C, Moore AJ. Chemical egg defence in the large milkweed bug, *Oncopeltus fasciatus*, derives from maternal but not paternal diet. *Entomologia Experimentalis et Applicata.* 2013;149:197–205.
103. Zhan S, Merlin C, Boore JL, Reppert SM. The monarch butterfly genome yields insights into long-distance migration. *Cell.* 2011;147:1171–85.
104. Robertson HM, Warr CG, Carlson JR. Molecular evolution of the insect chemoreceptor gene superfamily in *Drosophila melanogaster*. *Proc Natl Acad Sci U S A.* 2003;100(Suppl 2):14537–42.
105. Joseph RM, Carlson JR. *Drosophila* chemoreceptors: a molecular interface between the chemical world and the brain. *Trends Genet.* 2015;31:683–95.
106. Benton R. Multigene family evolution: perspectives from insect chemoreceptors. *Trends Ecol Evol.* 2015;30:590–600.
107. Rytz R, Crosset V, Benton R. Ionotropic receptors (IRs): chemosensory ionotropic glutamate receptors in *Drosophila* and beyond. *Insect Biochem Mol Biol.* 2013;43:888–97.

108. Kirkness EF, Haas BJ, Sun W, Braig HR, Perotti MA, Clark JM, Lee SH, Robertson HM, Kennedy RC, Elhaik E, et al. Genome sequences of the human body louse and its primary endosymbiont provide insights into the permanent parasitic lifestyle. *Proc Natl Acad Sci U S A*. 2010;107:12168–73.
109. Smadja C, Shi P, Butlin RK, Robertson HM. Large gene family expansions and adaptive evolution for odorant and gustatory receptors in the pea aphid, *Acyrtosiphon pisum*. *Mol Biol Evol*. 2009;26:2073–86.
110. Terrapon N, Li C, Robertson HM, Ji L, Meng X, Booth W, Chen Z, Childers CP, Glastad KM, Gokhale K, et al. Molecular traces of alternative social organization in a termite genome. *Nat Commun*. 2014;5:3636.
111. Xu W, Papanicolaou A, Zhang HJ, Anderson A. Expansion of a bitter taste receptor family in a polyphagous insect herbivore. *Sci Rep*. 2016;6:23666.
112. Feir D. *Oncopeltus fasciatus*: a research animal. *Annu Rev Entomol*. 1974;19:81–96.
113. Vellozo AF, Véron AS, Baa-Puyoulet P, Huerta-Cepas J, Cottret L, Febvay G, Calevro F, Rahbe Y, Douglas AE, Gabaldón T, et al. CycADS: an annotation database system to ease the development and update of BioCyc databases. *Database*. 2011;2011:bar008.
114. Baa-Puyoulet P, Parisot N, Febvay G, Huerta-Cepas J, Vellozo AF, Gabaldón T, Calevro F, Charles H, Colella S. ArthropodaCyc: a CycADS powered collection of BioCyc databases to analyse and compare metabolism of arthropods. *Database* (Oxford). 2016;2016:baw081.
115. Hojilla-Evangelista MP, Evangelista RL. Characterization of milkweed (*Asclepias* spp.) seed proteins. *Ind Crops Prod*. 2009;29: 275–80.
116. Dean CAE, Teets NM, Košťál V, Šimek P, Denlinger DL. Enhanced stress responses and metabolic adjustments linked to diapause and onset of migration in the large milkweed bug *Oncopeltus fasciatus*. *Physiol Entomol*. 2016;41:152–61.
117. Rabatel A, Febvay G, Gaget K, Duport G, Baa-Puyoulet P, Sapountzis P, Bendridi N, Rey M, Rahbé Y, Charles H, et al. Tyrosine pathway regulation is host-mediated in the pea aphid symbiosis during late embryonic and early larval development. *BMC Genomics*. 2013;14:235.
118. Dobler S, Petschenka G, Wagschal V, Flacht L. Convergent adaptive evolution – how insects master the challenge of cardiac glycoside-containing host plants. *Entomologia Experimentalis et Applicata*. 2015;157:30–9.
119. Grau-Bové X, Ruiz-Trillo I, Irimia M. Origin of exon skipping-rich transcriptomes in animals driven by evolution of gene architecture. *Genome Biol*. 2018;19:135.
120. Niehuis O, Gibson JD, Rosenberg MS, Pannebakker BA, Koevoets T, Judson AK, Desjardins CA, Kennedy K, Duggan D, Beukeboom LW, et al. Recombination and its impact on the genome of the haplodiploid parasitoid wasp *Nasonia*. *PLoS One*. 2010;5:e8597.
121. Ferrero A, Torreblanca A, Garcera MD. Assessment of the effects of orally administered ferrous sulfate on *Oncopeltus fasciatus* (Heteroptera: Lygaeidae). *Environ Sci Pollut Res Int*. 2017;24:8551–61.
122. Hare EE, Johnston JS. Genome size determination using flow cytometry of propidium iodide-stained nuclei. *Methods Mol Biol*. 2011;772:3–12.
123. Marçais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics*. 2011;27:764–70.
124. Bushnell B. BBMap short read aligner; 2016.
125. Gnerre S, Maccallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, Sharpe T, Hall G, Shea TP, Sykes S, et al. High-quality draft assemblies of mammalian genomes from massively parallel sequence data. *Proc Natl Acad Sci U S A*. 2011;108:1513–8.
126. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database management tool for second-generation genome projects. *BMC Bioinformatics*. 2011;12:491.
127. Lee E, Helt G, Reese J, Munoz-Torres M, Childers C, Buels R, Stein L, Holmes I, Elsik C, Lewis S. Web Apollo: a web-based genomic annotation editing platform. *Genome Biol*. 2013;14:R93.
128. Poelchau M, Childers C, Moore G, Tsavatapalli V, Evans J, Lee CY, Lin H, Lin JW, Hackett K. The i5k Workspace@NAL—enabling genomic data access, visualization and curation of arthropod genomes. *Nucleic Acids Res*. 2015;43:D714–9.
129. Murali SC, The i5k genome assembly team (29 additional authors), Han Y, Richards S, Worley K, Muzny D, Gibbs R, Koelzer S, Panfilio KA: *Oncopeltus fasciatus* genome assembly 1.0. *Ag Data Commons (Database)* 2015:<https://doi.org/10.15482/USDA.ADC/1173238>.
130. Hughes DST, Koelzer S, Panfilio KA, Richards S: *Oncopeltus fasciatus* genome annotations v0.5.3. *Ag Data Commons (Database)* 2015:<https://doi.org/10.15482/USDA.ADC/1173237>.
131. Vargas Jentzsch IM, Hughes DST, Poelchau M, Robertson HM, Benoit JB, Rosendale AJ, Armisén D, Duncan EJ, Vreede BMI, Jacobs CGC, et al: *Oncopeltus fasciatus* Official Gene Set v1.1. *Ag Data Commons (Database)* 2015:<https://doi.org/10.15482/USDA.ADC/1173142>.
132. Vargas Jentzsch IM, Kovacova V, Stueber K, Koelzer S, Panfilio KA: *Oncopeltus fasciatus* hybrid genome assembly 1.0. *Ag Data Commons (Database)* 2019:<https://doi.org/10.15482/USDA.ADC/1503405>.
133. RepeatModeler Open-1.0.8 [<http://www.repeatmasker.org>]. Accessed 5 June 2015.
134. Bao Z, Eddy SR. Automated de novo identification of repeat sequence families in sequenced genomes. *Genome Res*. 2002;12:1269–76.
135. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large genomes. *Bioinformatics*. 2005;21(Suppl 1):i351–8.
136. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res*. 1999;27:573–80.
137. RepeatMasker Open-4.0. [<http://www.repeatmasker.org>]. Accessed 5 June 2015.
138. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*. 2012;9:357–9.
139. Goff S, Vaughn M, McKay S, Lyons E, Stapleton A, Gessler D, Matasci N, Wang L, Hanlon M, Lenards A, et al. The iPlant Collaborative: cyberinfrastructure for plant biology. *Front Plant Sci*. 2011;2:34.
140. Grabherr MG, Haas BJ, Levin M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q, et al. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotechnol*. 2011;29:644–52.
141. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D, Li B, Lieber M, et al. De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nat Protoc*. 2013;8:1494–512.
142. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics*. 2011;12:323.
143. Schoville SD, Chen YH, Andersson MN, Benoit JB, Bhandari A, Bowsher JH, Brevik K, Cappelle K, Chen MM, Childers AK, et al. A model species for agricultural pest genomics: the genome of the Colorado potato beetle, *Leptinotarsa decemlineata* (Coleoptera: Chrysomelidae). *Sci Rep*. 2018;8:1931.
144. Scolari F, Benoit JB, Michalkova V, Aksoy E, Takac P, Abd-Alla AM, Malacrida AR, Aksoy S, Attardo GM. The spermatophore in *Glossina morsitans morsitans*: insights into male contributions to reproduction. *Sci Rep*. 2016;6:20334.
145. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*. 2015;31:3210–2.
146. Finn RD, Mistry J, Tate J, Coggill P, Heger A, Pollington JE, Gavin OL, Gunasekaran P, Ceric G, Forslund K, et al. The Pfam protein families database. *Nucleic Acids Res*. 2010;38:D211–22.
147. Weirauch MT, Hughes TR. A catalogue of eukaryotic transcription factor types, their evolutionary origin, and species distribution. *Subcell Biochem*. 2011;52:25–73.
148. Eddy SR. A new generation of homology search tools based on probabilistic inference. *Genome Inform*. 2009;23:205–11.
149. Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H, Remmert M, Soding J, et al. Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Mol Syst Biol*. 2011;7:539.
150. Moriya Y, Itoh M, Okuda S, Yoshizawa AC, Kanehisa M. KAA5: an automatic genome annotation and pathway reconstruction server. *Nucleic Acids Res*. 2007;35:W182–5.
151. Claudel-Renard C, Chevalet C, Faraut T, Kahn D. Enzyme-specific profiles for genome annotation: PRIAM. *Nucleic Acids Res*. 2003;31:6633–9.
152. Conesa A, Götz S. Blast2GO: a comprehensive suite for functional analysis in plant genomics. *Int J Plant Genomics*. 2008;2008:619832.
153. Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, Robles M. Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics (Oxford, England)*. 2005;21:3674–6.
154. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka G, et al. InterProScan 5: genome-scale protein function classification. *Bioinformatics*. 2014;30:1236–40.
155. Karp PD, Ouzounis CA, Moore-Kochlacs C, Goldovsky L, Kaipa P, Ahrén D, Tsoka S, Darzentas N, Kunin V, López-Bigas N. Expansion of the BioCyc collection of pathway/genome databases to 160 genomes. *Nucleic Acids Res*. 2005;33:6083–9.
156. Karp PD, Paley SM, Krummenacker M, Latendresse M, Dale JM, Lee TJ, Kaipa P, Gilham F, Spaulding A, Popescu L, et al. Pathway Tools version 13.0: integrated software for pathway/genome informatics and systems biology. *Brief Bioinform*. 2010;11:40–79.
157. Dereeper A, Guignon V, Blanc G, Audic S, Buffet S, Chevenet F, Dufayard JF, Guindon S, Lefort V, Lescot M, et al. Phylogeny.fr: robust phylogenetic analysis for the non-specialist. *Nucleic Acids Res*. 2008;36:W465–9.

158. Wang X, Fang X, Yang P, Jiang X, Jiang F, Zhao D, Li B, Cui F, Wei J, Ma C, et al. The locust genome provides insight into swarm formation and long-distance flight. *Nat Commun*. 2014;5:2957.
159. The International Silkworm Genome Consortium. The genome of a lepidopteran model insect, the silkworm *Bombyx mori*. *Insect Biochem Mol Biol*. 2008;38:1036–45.
160. Elsik CG, Worley KC, Bennett AK, Beye M, Camara F, Childers CP, de Graaf DC, Debyser G, Deng J, Devreese B, et al. Finding the missing honey bee genes: lessons learned from a genome upgrade. *BMC Genomics*. 2014;15:86.
161. Honeybee Genome Sequencing Consortium. Insights into social insects from the genome of the honeybee *Apis mellifera*. *Nature*. 2006;443:931–49.
162. Richards S, Gibbs RA, Weinstock GM, Brown SJ, Denell R, Beeman RW, Gibbs R, Bucher G, Friedrich M, Gimmelikhuijzen CJ, et al. The genome of the model beetle and pest *Tribolium castaneum*. *Nature*. 2008;452:949–55.
163. Chen XG, Jiang X, Gu J, Xu M, Wu Y, Deng Y, Zhang C, Bonizzoni M, Dermauw W, Vontas J, et al. Genome sequence of the Asian tiger mosquito, *Aedes albopictus*, reveals insights into its biology, genetics, and evolution. *Proc Natl Acad Sci U S A*. 2015;112:E5907–15.
164. Armisen D, Rajakumar R, Friedrich M, Benoit JB, Robertson HM, Panfilio KA, Ahn S-J, Poelchau MF, Chao H, Dinh H, et al. The genome of the water strider *Gerris buenoi* reveals expansions of gene repertoires associated with adaptations to life on the water. *BMC Genomics* in press:acceptance e-mail 12 Oct. 2018.
165. Papanicolaou A, Schetelig MF, Arensburger P, Atkinson PW, Benoit JB, Bourtzis K, Castañera P, Cavanaugh JP, Chao H, Childers C, et al. The whole genome sequence of the Mediterranean fruit fly, *Ceratitidis capitata* (Wiedemann), reveals insights into the biology and adaptive evolution of a highly invasive pest species. *Genome Biol*. 2016;17:192.
166. Ellis LL, Huang W, Quinn AM, Ahuja A, Alfrejd B, Gomez FE, Hjelmén CE, Moore KL, Mackay TF, Johnston JS, Tarone AM. Intrapopulation genome size variation in *D. melanogaster* reflects life history variation and plasticity. *PLoS Genet*. 2014;10:e1004522.
167. Kumar S, Stecher G, Tamura K. MEGA7: molecular evolutionary genetics analysis version 7.0 for bigger datasets. *Mol Biol Evol*. 2016;33:1870–4.
168. Croset V, Rytz R, Cummins SF, Budd A, Brawand D, Kaessmann H, Gibson TJ, Benton R. Ancient protostome origin of chemosensory ionotropic glutamate receptors and the evolution of insect taste and olfaction. *PLoS Genet*. 2010;6:e1001064.

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

