


RESEARCH

Open Access



Comparative genomic analysis of six *Glossina* genomes, vectors of African trypanosomes

Geoffrey M. Attardo^{22*}, Adly M. M. Abd-Alla¹³, Alvaro Acosta-Serrano¹⁶, James E. Allen⁶, Rosemary Bateta², Joshua B. Benoit²⁴, Kostas Bourtzis¹³, Jelle Caers¹⁵, Guy Caljon²¹, Mikkel B. Christensen⁶, David W. Farrow²⁴, Markus Friedrich³³, Aurélie Hua-Van⁵, Emily C. Jennings²⁴, Denis M. Larkin¹⁹, Daniel Lawson¹⁰, Michael J. Lehane¹⁶, Vasileios P. Lenis³⁰, Ernesto Lowy-Gallego⁶, Rosaline W. Macharia^{27,12}, Anna R. Malacrida²⁹, Heather G. Marco²³, Daniel Masiga¹², Gareth L. Maslen⁶, Irina Matetovici¹¹, Richard P. Meisel²⁵, Irene Meki¹³, Veronika Michalkova^{7,20}, Wolfgang J. Miller¹⁷, Patrick Minx³², Paul O. Mireji^{2,14}, Lino Ometto^{8,29}, Andrew G. Parker¹³, Rita Rio³⁴, Clair Rose¹⁶, Andrew J. Rosendale^{18,24}, Omar Rota-Stabelli⁸, Grazia Savini²⁹, Liliane Schoofs¹⁵, Francesca Scolari²⁹, Martin T. Swain¹, Peter Takáč³¹, Chad Tomlinson³², George Tsiamis²⁸, Jan Van Den Abbeele¹¹, Aurelien Vigneron³⁵, Jingwen Wang⁹, Wesley C. Warren^{32,36}, Robert M. Waterhouse²⁶, Matthew T. Weirauch⁴, Brian L. Weiss³⁵, Richard K. Wilson³², Xin Zhao³ and Serap Aksoy^{35*} 

Abstract

Background: Tsetse flies (*Glossina* sp.) are the vectors of human and animal trypanosomiasis throughout sub-Saharan Africa. Tsetse flies are distinguished from other Diptera by unique adaptations, including lactation and the birthing of live young (obligate viviparity), a vertebrate blood-specific diet by both sexes, and obligate bacterial symbiosis. This work describes the comparative analysis of six *Glossina* genomes representing three sub-genera: *Morsitans* (*G. morsitans morsitans*, *G. pallidipes*, *G. austeni*), *Palpalis* (*G. palpalis*, *G. fuscipes*), and *Fusca* (*G. brevipalpis*) which represent different habitats, host preferences, and vectorial capacity.

Results: Genomic analyses validate established evolutionary relationships and sub-genera. Syntenic analysis of *Glossina* relative to *Drosophila melanogaster* shows reduced structural conservation across the sex-linked X chromosome. Sex-linked scaffolds show increased rates of female-specific gene expression and lower evolutionary rates relative to autosome associated genes. Tsetse-specific genes are enriched in protease, odorant-binding, and helicase activities. Lactation-associated genes are conserved across all *Glossina* species while male seminal proteins are rapidly evolving. Olfactory and gustatory genes are reduced across the genus relative to other insects. Vision-associated Rhodopsin genes show conservation of motion detection/tracking functions and variance in the Rhodopsin detecting colors in the blue wavelength ranges.

Conclusions: Expanded genomic discoveries reveal the genetics underlying *Glossina* biology and provide a rich body of knowledge for basic science and disease control. They also provide insight into the evolutionary biology underlying novel adaptations and are relevant to applied aspects of vector control such as trap design and discovery of novel pest and disease control strategies.

Keywords: Tsetse, Trypanosomiasis, Hematophagy, Lactation, Disease, Neglected, Symbiosis

* Correspondence: gmattardo@ucdavis.edu; serap.aksoy@yale.edu

²²Department of Entomology and Nematology, University of California, Davis, Davis, CA, USA

³⁵Department of Epidemiology of Microbial Diseases, Yale School of Public Health, New Haven, CT, USA

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



Background

Flies in the genus *Glossina* (tsetse flies) are vectors of African trypanosomes, which are of great medical and economic importance in Africa. Sleeping sickness (human African trypanosomiasis or HAT) is caused by two distinct subspecies of the African trypanosomes transmitted by tsetse. In East and Southern Africa, *Trypanosoma brucei rhodesiense* causes the acute *Rhodesiense* form of the disease, while in Central and West Africa *T. b. gambiense* causes the chronic *Gambiense* form of the disease, which comprises about 95% of all reported HAT cases. Devastating epidemics in the twentieth century resulted in hundreds of thousands of deaths in sub-Saharan Africa [1], but more effective diagnostics now indicate that data concerning sleeping sickness deaths are subject to gross errors due to underreporting [2]. With hindsight, it is thus reasonable to infer that in reality, millions may have died from sleeping sickness since the implementation of trypanosomiasis surveillance and record-keeping by African colonial powers at the beginning of the twentieth century. Loss of interest and funding for control programs within the endemic countries resulted in a steep rise in incidence after the post-independence period of the 1960s. In an ambitious campaign to control the transmission of trypanosomiasis in Africa, multiple groups came together in a public/private partnership. These include the WHO, multiple non-governmental organizations, Sanofi Aventis, and Bayer. The public sector groups developed and implemented multi-country control strategies, and the companies donated the drugs required for the treatment of the disease. The campaign reduced the global incidence of *Gambiense* HAT to <3000 cases in 2015 [3]. Based on the success of the control campaign, there are now plans to eliminate *Gambiense* HAT as a public health problem by 2030 [4]. In contrast, control of *Rhodesiense* HAT has been more complex as disease transmission involves domestic animals, which serve as reservoirs for the parasite. Hence, the elimination of the *Rhodesiense* disease will require treatment or elimination of domestic reservoirs and/or reduction of tsetse vector populations. These strategies play a key part while medical interventions are used largely for humanitarian purposes. In addition to the public health impact of HAT, animal African trypanosomiasis (AAT or nagana) limits the availability of meat and milk products in large regions of Africa. It also excludes effective cattle rearing from ten million square kilometers of Africa [5] with wide implications for land use, i.e., constraints on mixed agriculture and lack of animal labor for plowing [6]. Economic losses in cattle production are estimated at 1–1.2 billion US dollars, and total agricultural losses caused by AAT are estimated at 4.75 billion US dollars per year [7, 8].

Achieving disease control in the mammalian host has been difficult given the lack of vaccines. This is due to the process of antigenic variation the parasite displays in its host. Hence, accurate diagnosis of the parasite and staging of the disease are important. This is of particular importance due to the high toxicity of current drugs available for the treatment of late-stage disease although the introduction of a simpler and shorter nifurtimox and eflornithine combination therapy (NECT) [9] and discovery of new oral drugs, such as fexinidazole [10] and acoziborole, are exciting developments. Although powerful molecular diagnostics have been developed in research settings, few have yet to reach the patients or national control programs [11]. Further complicating control efforts, trypanosomes are showing resistance to available drugs for treatment [12, 13]. While vector control is essential for zoonotic *Rhodesiense* HAT, it has not played a major role in *Gambiense* HAT as it was considered too expensive and difficult to deploy in the resource-poor settings of HAT foci. However, modeling, historical investigations, and practical interventions demonstrate the significant role that vector control can play in the control of *Gambiense* HAT [14–16], especially given the possibility of long-term carriage of trypanosomes in both human and animal reservoirs [17, 18]. The African Union has made removal of trypanosomiasis via tsetse fly control a key priority for the continent [19].

Within the *Glossinidae*, 33 extant taxa are described from 22 species in 4 subgenera. The first three subgenera *Austenina* Townsend, *Nemorhina* Robineau-Desvoidy, and *Glossina* Wiedemann correspond to the *Fusca*, *Palpalis*, and *Morsitans* species groups, respectively [20]. The fourth subgenus *Machadomia* was established in 1987 to incorporate *G. austeni*. The relationship of *G. austeni* Newstead with respect to the *Palpalis* and *Morsitans* complex flies remains controversial [21]. While molecular taxonomy shows that *Palpalis* and *Morsitans* species groups are monophyletic, the *Fusca* species group emerges as a sister group to all remaining *Glossinidae* [22]. *Morsitans* group taxa are adapted to drier habitats relative to the other two subgenera [23]. *Palpalis* group flies tend to occur in riverine and lacustrine habitats. *Fusca* group flies largely inhabit moist forests of West Africa. The host specificity of the different species groups vary, with the *Palpalis* group flies displaying strong anthropophily while the others are more zoophilic in preference. The principal vectors of HAT include *G. palpalis* s.l., *G. fuscipes*, and *G. m. morsitans* s.l. The riverine habitats of *Palpalis* group flies and their adaptability to peridomestic environments along with human blood meal preferences make them excellent vectors for HAT. Other species belonging to the *Morsitans* group (such as *G. pallidipes*) can also

transmit human disease, but principally play an important role in AAT transmission. In particular, *G. pallidipes* has a wide distribution and a devastating effect in East Africa. Also, of interest is *G. brevipalpis*, an ancestral tsetse species within the *Fusca* species complex. This species exhibits poor vectorial capacity with *T. brucei* relative to *G. m. morsitans* in laboratory infection experiments using colonized fly lines [24]. Comparison of the susceptibility of *G. brevipalpis* to *Trypanosoma congolense* (a species that acts as a major causative agent of AAT) also showed it has a much lower rate of infection relative to *Glossina austeni* [25].

To expand the genetic/genomic knowledge and develop new and/or improved vector control tools, a consortium in 2004, the International Glossina Genome Initiative (IGGI), was established to generate genetic and molecular resources for the tsetse research community [26]. The first tsetse fly genome from the *Glossina m. morsitans* species was published in 2014 [27]. However, questions regarding the genetics underlying tsetse species-specific traits, such as host preference and vector competence, required additional context. As such, we have assembled genomes from four species representing the three major *Glossina* sub-genera: *Morsitans* (*G. m. morsitans*, *G. pallidipes*), *Palpalis* (*G. palpalis*, *G. fuscipes*), and *Fusca* (*G. brevipalpis*) as well as one species with conflicted phylogenetic associations *Morsitans/Machadomia* (*G. austeni*). These species represent flies with differences in geographical localization, ecological preferences, host specificity, and vectorial capacity (Fig. 1). Here, we report on the evolution and genetics underlying this genus by comparison of their genomic architecture and predicted protein-coding sequences as well as highlighting some of the genetic differences that hold clues to the differing biology between these species.

Results and discussion

Genome assemblies and global features of note

The genomic sequences for the tsetse species described here originated from mother and daughter lines for each respective *Glossina* species (Additional file 1: Table S1). Sequencing and assembly of the resulting reads produced scaffolds of varied sizes, contiguity, and coverage (Table 1). The total assembled sequencing coverage varied between 45 and 58× for each species. The average assembled size was 359 Mb with the greatest contiguity measured for *G. pallidipes*, which comprised the fewest contigs ($n = 7275$) with an N50 contig length of 167 kb. On average, the new *Glossina* assemblies resulted in fewer contigs (17,604 vs 24,071) at a higher level of contiguity (72 vs 49 kb) than the original *G. morsitans* assembly. This is likely due to the advancements in the sequencing technologies and software utilized to sequence and assemble these genomes relative to the

original *G. morsitans* genome. The *G. morsitans* genome also has fewer predicted genes relative to the more recently produced genomes, suggesting that additional sequencing on this species would be informative.

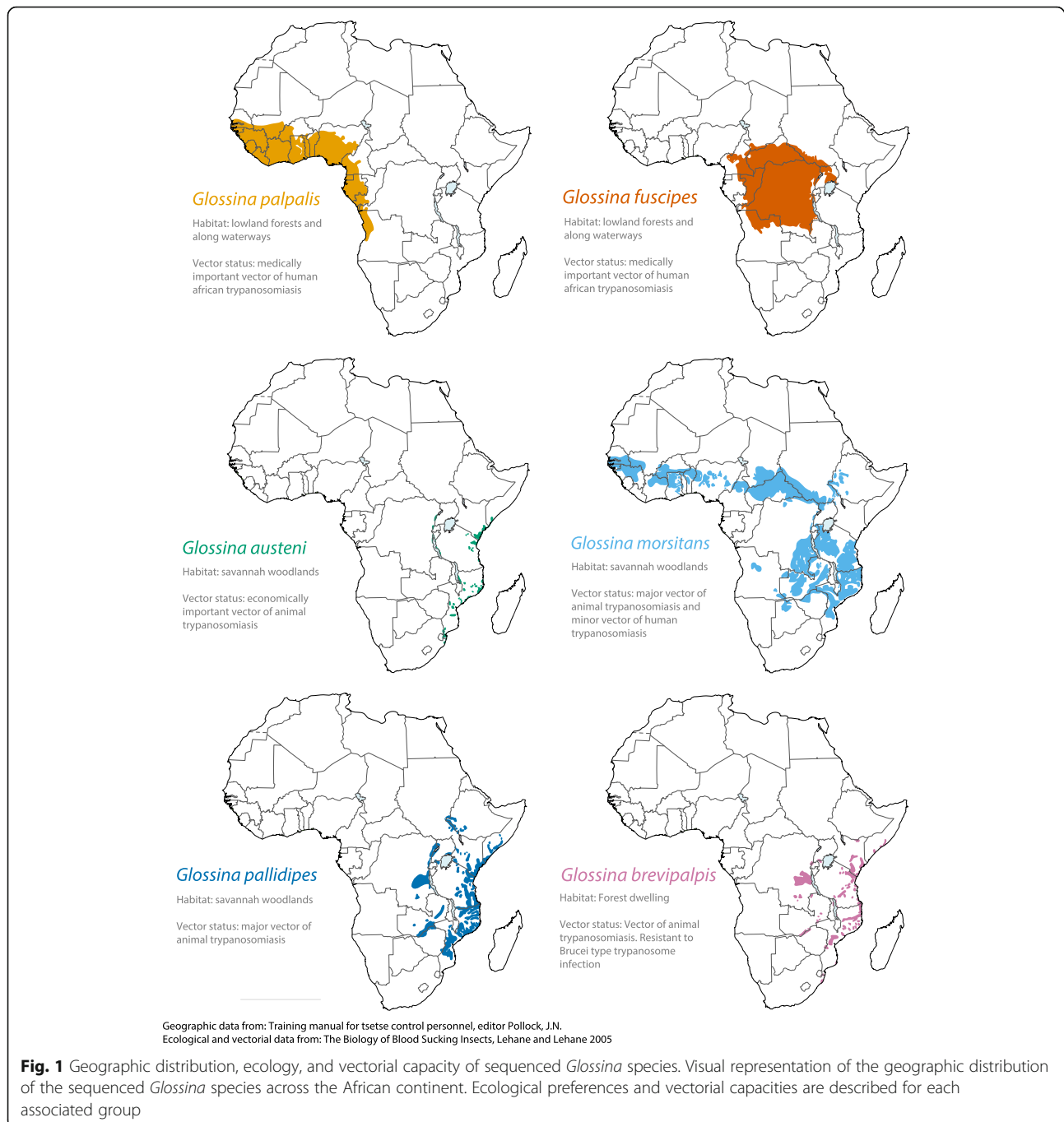
The GC content of these genomes ranges from 27% (*G. brevipalpis*) to 35% (*G. pallidipes*). Genomic regions with low GC content are associated with heterochromatic DNA which is often transcriptionally inactive [28]. The lower GC content in *G. brevipalpis* relative to other tsetse species could result in additional regions of lower transcriptional activity.

Completeness and accuracy of gene model predictions (Additional file 1: Table S2) within the genomes were determined by Benchmarking Universal Single-Copy Orthologs (BUSCO) analysis (Table 2). This analysis revealed high levels of representation of universal orthologs in all *Glossina* species. The scores for the genomes ranged from 92% representation (*G. morsitans*) to between 97 and 98% (the remaining five *Glossina* species). The lower level of representation within *G. morsitans* probably results from the fact that it was assembled from sequence data derived from multiple older technologies using the now unsupported Celera assembly software [29].

Repeat analysis and transposable element composition

A comparative analysis was performed on the quantities and types of repetitive elements contained within the six tsetse genomes (Fig. 2). The analysis reveals a similar content across the six genomes in terms of the number of consensus sequences and subclass diversity. The total percentage of masked repeats ranges from 34.95% (*G. brevipalpis*) to 39.99% (*G. pallidipes*) (Additional file 1: Table S3) consisting mainly of dispersed transposable elements (TEs) as well as simple repeats (tandem, satellite, and low-complexity sequences). *G. brevipalpis* contains the highest proportion of simple repeats and the lowest proportions and coverage of TEs. For all tsetse genomes, three subclasses of TEs predominate: DAN terminal inverted repeats (TIR) transposons (class II DNA), rolling circle Helitrons (Class II RC), and long interspersed nuclear elements (class I LINE). Other class I elements such as LTR retroelements (class I LTR) or small interspersed nuclear elements (class I SINE) are very scarce. In all genomes, a significant part of dispersed repetitive elements remains unknown and then unclassified (Fig. 2a).

After the clustering of the 7583 TE consensus into 2906 clusters, a distribution analysis reveals that most of the TE content is either shared between 5 species or is species-specific (Fig. 2b, Additional file 1: Table S4). The *G. brevipalpis* genome, containing the lowest overall repeat content, is substantially different from the other 5 species. For instance, whereas these genomes have very



similar proportions of LINE families, the most abundant one (LINE/CR1) is largely underrepresented in the *G. brevipalpis* genome (Fig. 2c). On the opposite, among the DNA subclass, the DNA/TcMar families (especially *mariner*) are very abundant in all genomes including *G. brevipalpis*, but this genome also contains a higher proportion of DNA/TcMar-Tc1 families (Fig. 2d). Based on these analyses, the DNA/*mariner* elements appear to have diversified and expanded prior to the split between *G. brevipalpis* and the rest of the *Glossina* species, but

also after the split (Additional file 2: Figure S1). This differs from the LINE/CR1 families which seem to have expanded and diversified mainly after the split, in the other 5 species, whereas the DNA/Tc1 family would have specifically expanded in *G. brevipalpis*.

The total assembled repeats did not correlate with assembly contiguity measures, meaning high repeat content did not equate to lower assembly contiguity. Nonetheless, given highly repeat-rich regions are largely inaccessible to short read length sequencing and

Table 1 *Glossina* species contig and scaffold assembly statistics

Scaffold length	<i>Glossina morsitans</i>	<i>Glossina pallidipes</i>	<i>Glossina austeni</i>	<i>Glossina fuscipes</i>	<i>Glossina palpalis</i>	<i>Glossina brevipalpis</i>
Total genomic coverage	100x	46x	50x	52x	58x	45x
Genome size (Mb)	366	357	370	374	380	315
> 1 Mb	13	102	78	70	63	81
250 kb–1 Mb	138	248	316	393	395	202
100–250 kb	605	184	248	330	326	136
10–100 kb	3663	290	379	496	709	257
5–10 kb	737	106	94	165	507	85
2–5 kb	1933	255	206	252	978	156
< 2 kb	6718	541	884	689	948	734
Total no. of contigs	24,071	7275	18,748	13,688	31,320	16,993
N50 contig length (kb)	49	167	46	64	24	62
Total no. of Scaffolds	13,807	1726	2205	2395	3926	1651
GC content (%)	33	35	34	34	34	27
N50 scaffold length (kb)	120	1038	812	561	575	1209
L50 (rank of N50 scaffold)	569	94	115	178	186	62
Repeat content (%)	34.95	35.49	38.64	37.09	35.49	37.67

N50 is defined as the minimum contig length needed to cover 50% of the genome. L50 is defined as the smallest number of contigs whose length sum makes up half of genome size

Table 2 Quantification of *Glossina* gene predictions and genomic completeness by Benchmarking Universal Single-Copy Orthologs (BUSCO) analysis

Species	Complete BUSCOs (%)	Complete and single-copy BUSCOs (%)	Complete and duplicated BUSCOs (%)	Fragmented BUSCOs (%)	Missing BUSCOs (%)	Total BUSCO groups searched (%)
BUSCO gene analysis results (percentage) (diptera_odb9 geneset)						
<i>G. morsitans</i>	93.53	88.00	5.54	3.22	3.25	100.00
<i>G. pallidipes</i>	95.53	90.78	4.75	2.72	1.75	100.00
<i>G. austeni</i>	97.11	93.00	4.11	2.18	0.71	100.00
<i>G. fuscipes</i>	96.50	91.14	5.36	2.32	1.18	100.00
<i>G. palpalis</i>	95.00	87.53	7.47	3.32	1.68	100.00
<i>G. brevipalpis</i>	95.14	89.03	6.11	2.97	1.89	100.00
BUSCO genomic analysis results (percentage) (diptera_odb9 geneset)						
<i>G. morsitans</i>	92.03	91.25	0.79	3.32	4.64	100.00
<i>G. pallidipes</i>	98.43	97.36	1.07	1.07	0.50	100.00
<i>G. austeni</i>	98.07	97.18	0.89	1.25	0.68	100.00
<i>G. fuscipes</i>	98.32	97.21	1.11	1.18	0.50	100.00
<i>G. palpalis</i>	97.07	92.85	4.22	1.86	1.07	100.00
<i>G. brevipalpis</i>	97.96	97.11	0.86	1.25	0.79	100.00

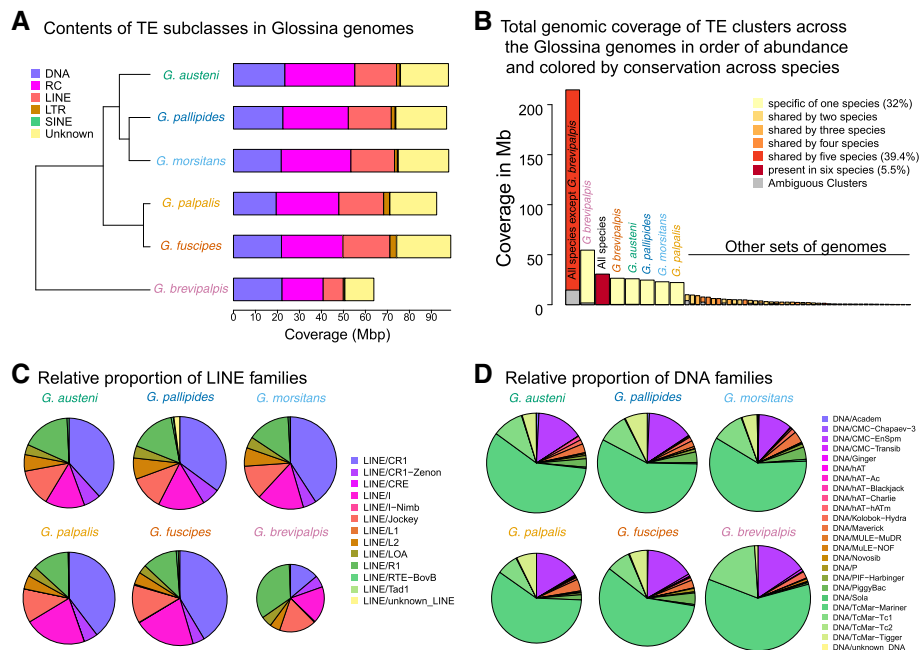


Fig. 2 Comparative analysis of repetitive elements within the *Glossina* genomes. **a** Graphical representation of the constitution and sequence coverage by the various classes of identified dispersed repetitive elements. **b** Coverage of TE families that are shared between species. More than 75% of the total coverage (eight first magnified bars) correspond to TE either specific to one species, shared by all species, or shared by the five closest. **c** Relative constitution of DNA terminal inverted repeat (TIR) families across the *Glossina* genomes. **d** Relative constitution of long interspersed nuclear elements (LINEs) across the *Glossina* genomes. For **c** and **d**, the size of the pie charts reflects the proportion of the subclass among the dispersed repetitive sequences

assembly methods, our approximation of repeat element content in *Glossina* is likely an underestimation and more detailed distribution and measures of transposable elements will require further experimentation.

Multiple genetic comparisons confirm *Glossina* phylogenetic relationships and the inclusion of *G. austeni* as a member of the *Morsitans* sub-genus

Sequence similarity between the genomes was analyzed using whole-genome nucleotide alignments of supercontigs and predicted coding sequences from the five new *Glossina* genomes as well as those from the *Musca domestica* genome using *G. m. morsitans* as a reference (Fig. 3a). The results indicate that *G. pallidipes* and *G. austeni* are most similar at the sequence level to *G. m. morsitans*. This is followed by the species in the *Palpalis* sub-genus (*G. fuscipes* and *G. palpalis*). The remaining species (*G. brevipalpis*) shows the least sequence conservation relative to *G. m. morsitans* followed by the out-group species *M. domestica*. The lower sequence similarity between *G. brevipalpis* and the other tsetse species reinforces its status as a sister group to the *Morsitans* and *Palpalis* sub-genera.

Alignment of the predicted coding sequences produced a similar result to that observed in the whole-genome alignment in terms of similarity to *G. m. morsitans*

(Fig. 2a). Of interest is that more than 25% of the *G. m. morsitans* exon sequences were not align-able with *G. brevipalpis*, indicating that they were either lost, have diverged beyond alignability, or were in an unsequenced region in *G. brevipalpis*. In addition, *G. brevipalpis* has on average ~ 5000 fewer predicted protein-coding genes than the other species. Given the low GC content of the *G. brevipalpis* sequenced genome, it is possible that some of the regions containing these sequences lie within heterochromatin. Difficulties associated with sequencing heterochromatic regions may have excluded these regions from our analysis; however, it also implies that if these protein-coding genes are indeed present, they are located in a region of the genome with low transcriptional activity.

We inferred the phylogeny and divergence times of *Glossina* using a concatenated alignment of 286 single-copy gene orthologs (478,000 nucleotide positions) universal to *Glossina* (Fig. 3b). The tree recovered from this analysis has support from both maximum likelihood and Bayesian analyses, using respectively homogeneous and heterogeneous models of replacement. A coalescent-aware analysis further returned full support, indicating a speciation process characterized by clear lineage sorting with no introgression between species (Additional file 2: Figure S2). These results suggest an allopatric speciation

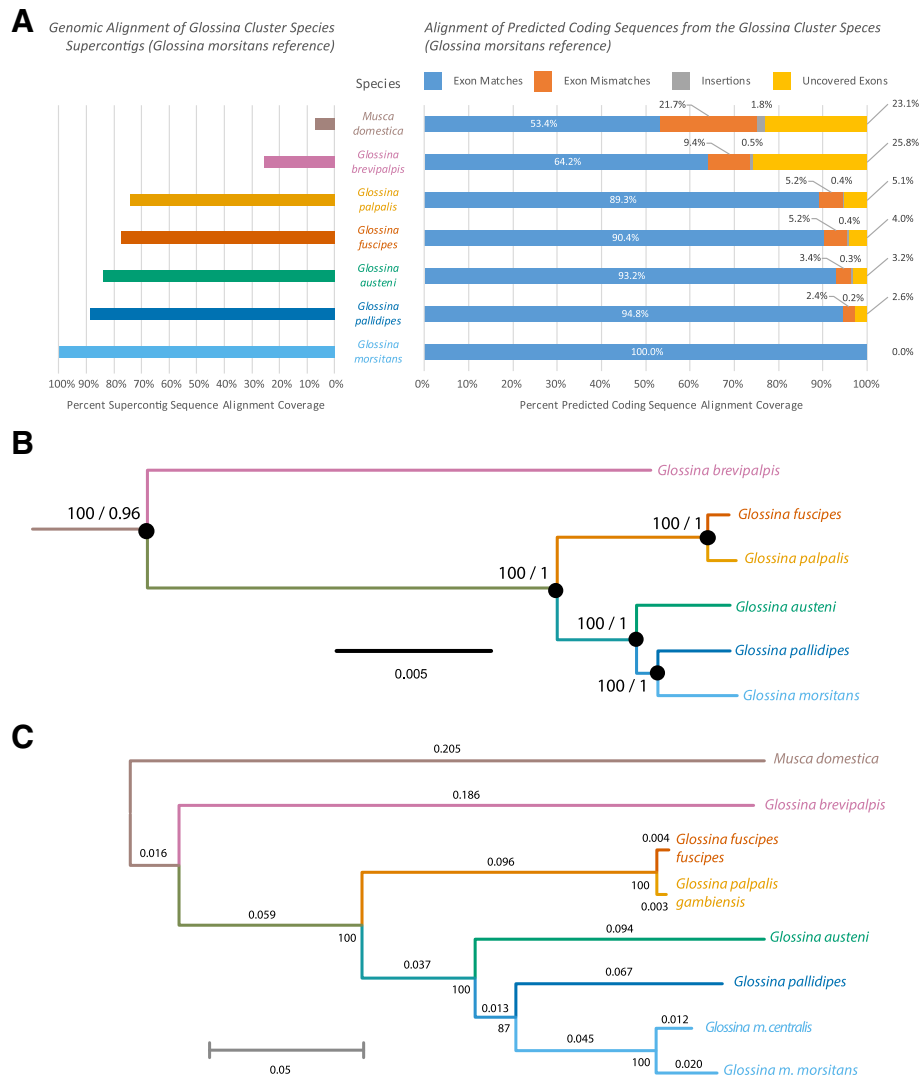


Fig. 3 *Glossina* whole-genome alignment, phylogenetic analysis of orthologous protein-coding nuclear genes, and phylogenetic analysis of mitochondrial sequences. **a** Analysis of whole-genome and protein-coding sequence alignment. The left graph reflects the percentage of total genomic sequence aligning to the *G. m. morsitans* reference. The right side of the graph represents the alignment of all predicted coding sequences from the genomes with coloration representing matches, mismatches, insertions, and uncovered exons. **b** Phylogenetic tree from conserved protein-coding sequences. Black dots at nodes indicate full support from maximum likelihood (Raxml), Bayesian (Phylobayes), and coalescent-aware (Astral) analyses. Raxml and Phylobayes analyses are based on an amino acid dataset of 117,782 positions from 286 genes from 12 species. The Astral analyses are based on a 1125-nucleotide dataset of 478,617 positions from the 6 *Glossina* (full trees are in Additional file 2: Figure S2A-C). The values at nodes represent the bootstrap supports and posterior probabilities from the maximum likelihood and Bayesian analyses, respectively (Bootstrap/posterior probability). **c** Molecular phylogeny derived from whole mitochondrial genome sequences. The analysis was performed using the maximum likelihood method with MEGA 6.0

process characterized by a small founder population size followed by little to no introgression among newly formed species.

Furthermore, we assembled complete mitochondrial (mtDNA) genome sequences for each species as well as *Glossina morsitans centralis* as references for use in distinguishing samples at the species, sub-species, or haplotype levels. All the mtDNA genomes encode large (16S rRNA) and small (12S rRNA) rRNAs, 22 tRNAs, and 13 protein-coding genes. Phylogenetic analysis of the

resulting sequences using the maximum likelihood method resulted in a tree with congruent topology to that produced by the analysis of the concatenated nuclear gene alignment (Fig. 3c). A comparative analysis of the mtDNA sequences identified variable marker regions with which to identify different tsetse species via traditional sequencing and/or high-resolution melt analysis (HRM) (Additional file 2: Figure S3). Analysis of the amplicons from this region using HRM facilitated the discrimination of these products based on their

composition, length, and GC content. The use of HRM on these variable regions successfully resolved the differences between test samples consisting of different tsetse species as well as individuals with different haplotypes or from different populations (Additional file 2: Figure S4). This method provides a rapid, cost-effective, and relatively low-tech way of identifying differences in the field-caught tsetse for the purposes of population genetics and measurement of population diversity.

The trees derived from the nuclear and mitochondrial phylogenetic analyses agree with previously published phylogenies for tsetse [22, 30, 31], and the species delineate into groups representing the defined *Fusca*, *Palpalis*, and *Morsitans* sub-genera.

A contentious issue within the taxonomy of *Glossina* is the placement of *G. austeni* within the *Machadomia* sub-genus. Comparative anatomical analysis of the male genitalia places *G. austeni* within the *morsitans* sub-genus. However, female *G. austeni* genitalia bear anatomical similarities to the members of the *Fusca* sub-genus. In addition, *G. austeni*'s habitat preferences and some external morphology resemble those of the *palpalis* sub-genus [30]. Recent molecular evidence suggests that *G. austeni* are closer to the *morsitans* sub-genus [22, 31]. The data generated via the three discrete analyses described above all support the hypothesis that *G. austeni* is a member of the *Morsitans* sub-genus rather than the *Palpalis* sub-genus and belongs as a member of the *Morsitans* group rather than its own discrete sub-genus.

Comparative analysis of *Glossina* with *Drosophila* reveals reduced synteny and female-specific gene expression on X-linked scaffolds

The scaffolds in each *Glossina* spp. genome assembly were assigned to chromosomal arms based on orthology and relative position to protein-coding sequences in the *D. melanogaster* genome (*Drosophila*) [32]. The *Glossina* and *Drosophila* genomes contain six chromosome arms (Muller elements A–F) [33–35]. We assigned between 31 and 52% of annotated genes in each species to a Muller element, which we used to assign >96% of scaffolds to Muller elements in each species (Additional file 1: Table S5). From these results, we inferred the relative size of each Muller element in each species by counting the number of annotated genes assigned to each element and calculating the cumulative length of all assembled scaffolds assigned to each element. Using either measure, we find that element E is the largest and element F is the shortest in all species, consistent with the observations in *Drosophila* [36] (Fig. 4).

Mapping of the *Glossina* scaffolds to the *Drosophila* Muller elements reveals differing levels of conservation of synteny (homologous genomic regions with

maintained orders and orientations) across these six species relative to *Drosophila*. In *G. m. morsitans*, the X chromosome is composed of Muller elements A, D, and F as opposed to the *Drosophila* X which only contains A and sometimes D [35], and all other *Glossina* species besides *G. brevipalpis* have the same karyotype [37]. We therefore assume that the same elements are X-linked in the other *Glossina* species (apart from *G. brevipalpis*). This analysis reveals that scaffolds mapping to *Drosophila* Muller element A show a reduced overall level of syntenic conservation relative to the other Muller elements while the scaffolds mapping to *Drosophila* Muller element D (part of the *Glossina* X chromosome, but not the *D. melanogaster* X) retain more regions of synteny conservation. We hypothesize that the lower syntenic conservation on element A reflects a higher rate of rearrangement because it has been X-linked for more time (both in the *Drosophila* and *Glossina* lineages) than element D (only in *Glossina*) and rearrangement rates are higher on the X chromosome (element A) in *Drosophila* [36].

To examine the relationship between gene expression and DNA sequence evolution, we compared the gene expression levels between the X chromosome and autosomes using sex-specific RNA-seq libraries derived from whole males, whole non-lactating females, and whole lactating females for all the *Glossina* species apart from *G. pallidipes*. Consistent with the previous results from *G. m. morsitans* [35], the ratio of female to male expression is greater on the X chromosome than on the autosomes across species (Additional file 2: Figure S5). In addition, there is a deficiency of genes with male-biased expression (upregulated in males relative to females) on the X-linked elements in all species (Additional file 2: Figure S6). Reduced levels of male-biased gene expression have also been observed in mosquitoes and is a conserved feature of the *Anopheles* genus [38]. The X chromosome is hemizygous in males, which exposes recessive mutations to natural selection and can accelerate the rate of adaptive substitutions and facilitate the purging of deleterious mutations on the X chromosome [39, 40]. Using dN/dS values for annotated genes, we fail to find any evidence for this faster-X effect across the entire phylogeny or along any individual lineages (Additional file 2: Figure S7). The faster-X effect is expected to be greatest for genes with male-biased expression because they are under selection in males [39], but we find no evidence for faster-X evolution of male-biased genes in any of the *Glossina* species. In contrast, there is some evidence for “slower-X” evolution among female-biased genes (Additional file 2: Figure S8), suggesting that purifying selection is more effective at purging deleterious mutations on the X chromosome [41]. Genes with female-biased expression tend to be broadly expressed

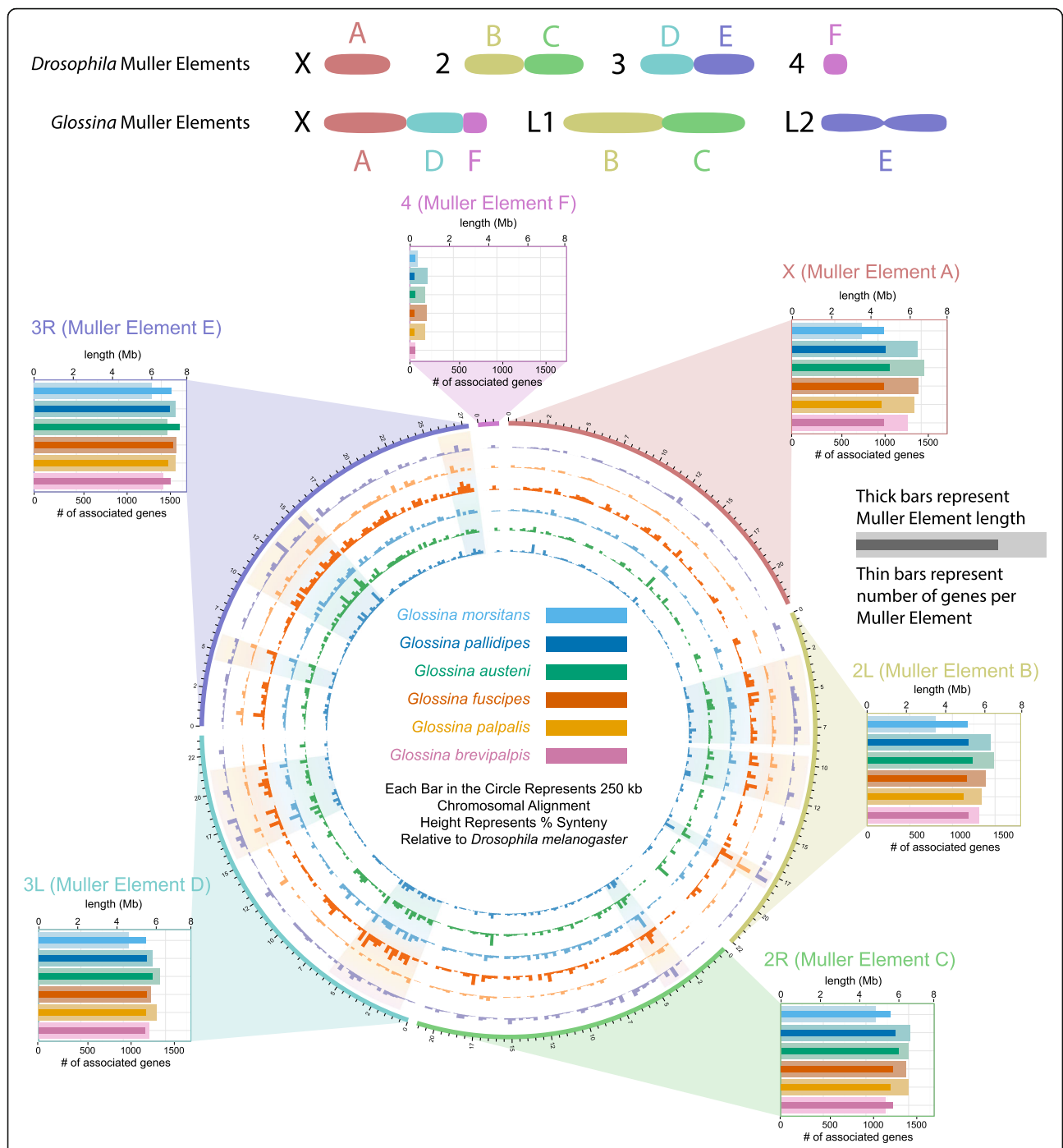


Fig. 4 Visualization of syntenic block analysis data and predicted Muller element sizes. Level of syntenic conservation between tsetse scaffolds and *Drosophila* chromosomal structures (Muller elements). The color-coded concentric circles consisting of bars represent the percent of syntenic conservation of orthologous protein-coding gene sequences between the *Glossina* genomic scaffolds and *Drosophila* Muller elements. Each bar represents 250 kb of aligned sequence, and bar heights represent the percent of syntenic conservation. The graphs on the periphery of the circle illustrate the combined predicted length and number of genes associated with the Muller elements for each tsetse species. The thin darkly colored bars represent the number of 1:1 orthologs between each *Glossina* species and *D. melanogaster*. The thicker lightly colored bands represent the predicted length of each Muller element for each species. This was calculated as the sum of the lengths of all scaffolds mapped to those Muller elements

[42], suggesting that pleiotropic constraints on female-biased genes increase the magnitude of purifying selection and produce the observed slower-X effect [43].

The exception to these observations is element F. Element F, the smallest X-linked element, has low female expression and an excess of genes with male-biased expression (Additional file 2: Figure S9). In contrast with the other X-linked Muller elements in *Glossina*, the dN/dS ratios of all Element F-associated genes (male-biased and unbiased) suggest that they are evolving faster than the rest of the genome across all tsetse lineages (Additional file 2: Figure S10). The F elements in *Drosophila* species, while not being X-linked, show similar properties in that they have lower levels of synteny, increased rates of inversion, and higher rates of protein-coding sequence evolution, suggesting that the F element is rapidly evolving in flies within Schizophora [44].

The *G. austeni* genome contains *Wolbachia*-derived chromosomal insertions

A notable feature of the *G. m. morsitans* genome was the integration of large segments of the *Wolbachia* symbiont genome via horizontal gene transfer (HGT). Characterization of the *G. m. morsitans* HGT events revealed that the chromosomal sequences with transferred material contained a high degree of nucleotide polymorphisms, coupled with insertions and deletions [45]. These observations were used in this analysis to distinguish cytoplasmic from chromosomal *Wolbachia* sequences during the in silico characterization of the tsetse genomes. Analysis of the six assemblies revealed that all contain sequences homologous to *Wolbachia*. However, in *G. pallidipes*, *G. fuscipes*, *G. palpalis*, and *G. brevipalpis*, the homologous sequences were limited to short fragments and likely represent artifacts. Additional screening of these lines by PCR with *Wolbachia*-specific primers yielded negative results, suggesting that this is the case. This is in agreement with negative PCR-based screening of *Wolbachia* infections in natural populations of these species indicating that these short segments could be artifacts or contaminants [46]. The exception to this is *G. austeni* which contains more extensive chromosomal integrations of *Wolbachia* DNA (Additional file 1: Table S6).

All *Wolbachia* sequences, chromosomal and cytoplasmic, identified in *G. austeni* were mapped against the reference genomes of *Wolbachia* strains *wMel*, *wGmm*, and the chromosomal insertions A and B in *G. m. morsitans* (Fig. 5). The *G. austeni* chromosomal insertions range in size from 500 to 95,673 bps with at least 812 DNA fragments identified in silico. Sequence homology between *wMel*, *wGmm*, and the chromosomal insertions A and B in *G. m. morsitans* varied between 98 and 63%, with the

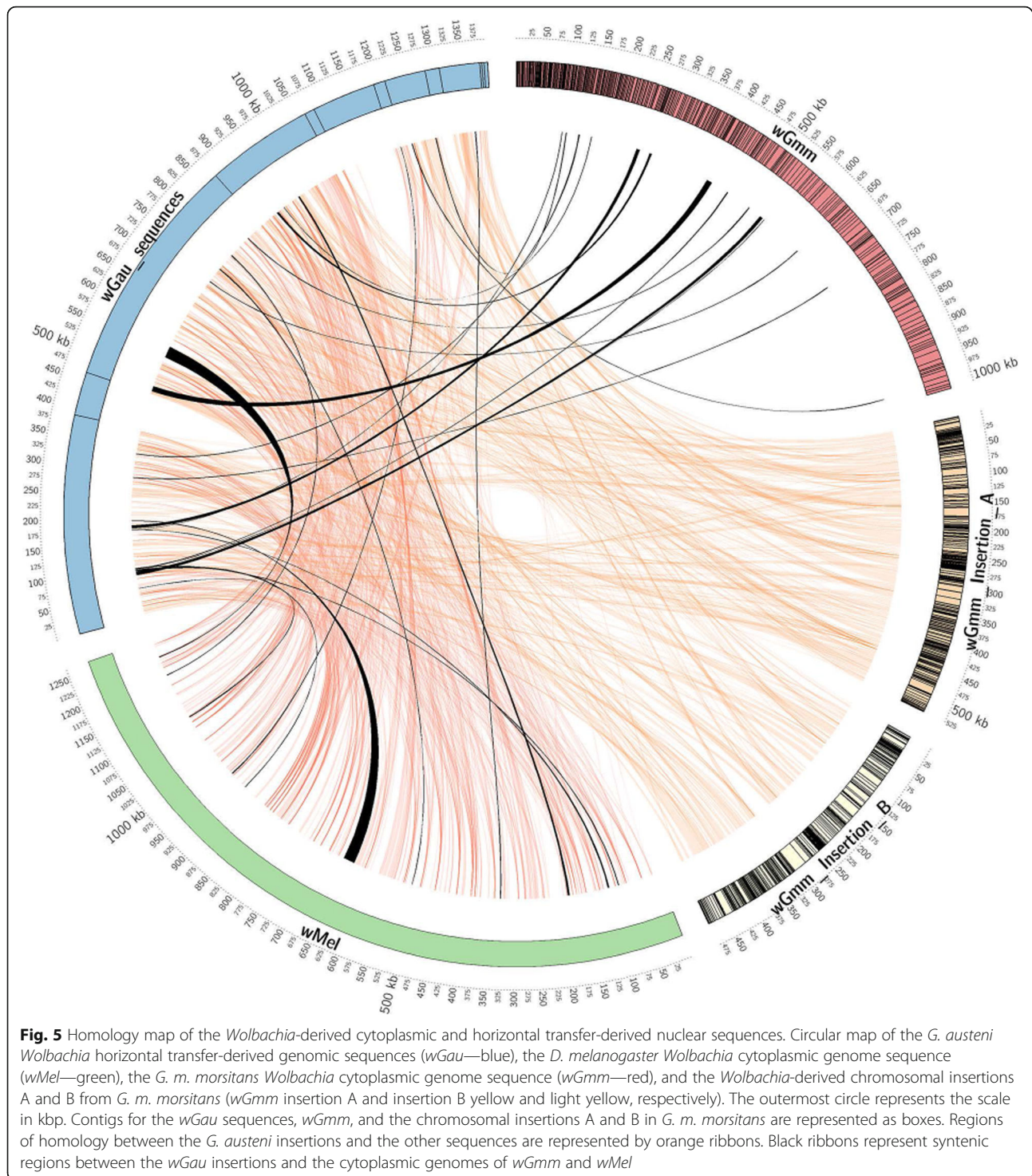
highest sequence homologies observed with chromosomal insertions A and B from *G. m. morsitans*. The similarity between the genomic insertions in *G. m. morsitans* and *G. austeni* relative to cytoplasmic *Wolbachia* sequences suggests they could be derived from an event in a common ancestor. However, the absence of comparable insertions in *G. pallidipes* (a closer relative to *G. m. morsitans*) indicate that either these insertions occurred independently or that the region containing the insertions was not assembled in *G. pallidipes*. Additional data from field-based *Glossina* species/sub-species is required to determine the true origin of these events.

The biological implications of the insertions in *G. morsitans* and *G. austeni* remain ambiguous. Prior gene expression analyses of *Wolbachia* insertions in *G. morsitans* using RNA-seq data found little to no evidence of gene expression from these insertions [45]. This suggests that these may be accidental transfer events associated with the long-term symbiosis between the species. Additional research is required to understand the origin, evolutionary history, and functionality of these HGT events.

Analysis of *Glossina* genus- and sub-genus-specific gene families reveals functional enrichments

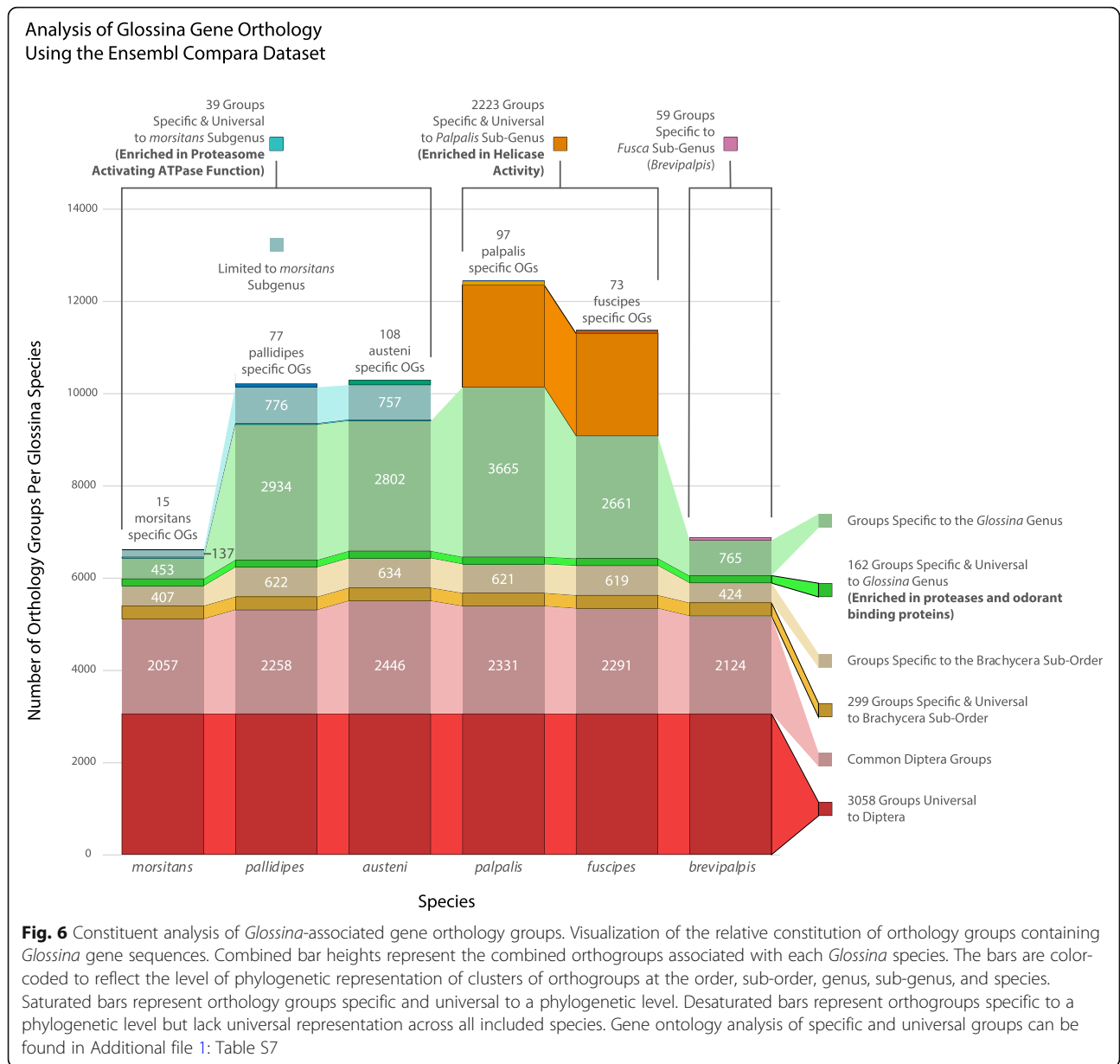
All annotated *Glossina* genes were assigned to groups (orthology groups (OGs)) containing predicted orthologs from other insect and arthropod species represented within VectorBase. A global analysis of all the groups containing *Glossina* genes was utilized to determine the gene composition of these flies relative to their Dipteran relatives and between the *Glossina* sub-genera. An array of 12 Diptera is represented within this analysis including *Anopheles gambiae* (Nematocera), *Aedes aegypti* (Nematocera), *Lutzomyia longipalpis* (Nematocera), *Drosophila melanogaster* (Brachycera), *Stomoxys calcitrans* (Brachycera), and *Musca domestica* (Brachycera).

The tsetse-associated OGs range from those containing sequences representing all the dipteran species included in the analysis to those with sequences specific to individual tsetse species. The composition of these OGs breaks down to a core of 3058 OGs with constituents universal to Diptera (93,430 genes), 299 OGs specific and universal to Brachyceran flies (4975 genes), and 162 OGs specific and universal to *Glossina* (1548 genes). A dramatic feature identified by this analysis is the presence of 2223 OGs specific and universal to the *Palpalis* sub-genus (*G. fuscipes* and *G. palpalis* 4948 genes). This contrasts with the members of the *Morsitans* sub-genus (*G. m. morsitans*, *G. pallidipes*, and *G. austeni*) in which there are 137 specific and universal OGs (153 genes) (Fig. 6, Additional file 3, Additional file 4).



To understand the functional significance of the *Glossina*-specific OGS, we performed an analysis of functional enrichment of Gene Ontology (GO) terms within these groups. Many of the *Glossina*-specific genes are not currently associated with GO annotations as they lack characterized homologs in other

species. As such, these sequences were not included in this analysis. However, ~60% of the genes within the combined *Glossina* gene repertoire are associated with GO annotations, which allowed for the analysis of a sizable proportion of the dataset (Additional file 1: Table S7, Additional file 5).



***Glossina* genus universal and specific genes are enriched in genes coding for proteases and odorant-binding proteins**

The orthology groups containing genes specific and universal to the *Glossina* genus are enriched in odorant-binding and serine-type endopeptidase activities. The universality of these genes within *Glossina* and their absence from the other surveyed Dipteran species suggest they are currently associated with tsetse-specific adaptations.

The ontology categories with the most significant *p* values across all six species represent proteolysis-associated genes (serine-type endopeptidase activity (GO:0004252) and proteolysis (GO:0006508)). This category

encompasses 92 *Glossina*-specific proteases with predicted serine-type endopeptidase activity. The abundance of this category may be an adaptation to the protein-rich blood-specific diet of both male and female flies. A similar expansion of serine proteases is associated with blood-feeding in mosquitoes, and the presence of an equivalent expansion in tsetse may represent an example of convergent evolution [47]. This class of peptidases is also associated with critical functions in immunity, development, and reproduction in Diptera [48–51]. Homology analysis of these proteases by BLAST against an insect-specific subset of the NCBI NR database reveals that most bear the closest homology to chymotrypsins and trypsin proteases in other Brachyceran

Diptera (Additional file 6). Many of these homologs remain undefined in terms of their function in other systems. Determination of the functions associated with these expansions will require further investigations into their expression patterns and analysis of their putative roles in digestion, development, reproduction, and immunity.

The other enriched GO term common to all *Glossina* is for genes encoding odorant-binding proteins (OBPs). Of the 370 OBPs annotated within *Glossina*, 55 lack orthologs in species outside of *Glossina*. The primary function of OBPs is to bind small hydrophobic molecules to assist in their mobilization in an aqueous environment. These proteins are often associated with olfaction functions as many are specifically expressed in chemosensory-associated tissues/organs where they bind small hydrophobic molecules and transport them to odorant receptors [52, 53]. However, functional analyses in *G. m. morsitans* have associated an OBP (OBP6) with developmental activation of hematopoiesis during larvigenesis in response to the mutualistic *Wigglesworthia* symbiont [54]. In addition, many of the OBPs identified in this analysis are characterized as *Glossina*-specific seminal proteins with male accessory gland-specific expression patterns. They are primary constituents of the spermatophore structure produced by the male tsetse during mating [55]. The genus-specific nature of these OBPs suggests that they are key components of reproductive adaptations of male tsetse.

The *Palpalis* sub-genus contains a large group of sub-genus-specific genes

A large group of genes specific and universal to the members of the *Palpalis* sub-genus (*G. palpalis* and *G. fuscipes*) was a defining feature of the orthology analysis. The expansion includes 2223 OGs and encompasses 4948 genes between *G. palpalis* and *G. fuscipes*. Homology-based analysis of these genes by comparison against the NCBI NR database revealed significant (e value $< 1 \times 10^{-10}$) results for 603 of the genes. Within this subset of genes, ~5% represent bacterial contamination from tsetse's obligate endosymbiont *Wigglesworthia*. Sequences homologous to another well-known bacterial symbiont *Spiroplasma* were found exclusively in *G. fuscipes*. This agrees with previous observations of *Spiroplasma* infection of colonized and field-collected *G. fuscipes* flies [56].

Four genes bear homology to viral sequences (GPP1051037/GFUI045295 and GPP1016422/GFUI028200). These sequences are homologous to genes from *Ichnoviruses*. These symbiotic viruses are transmitted by parasitic Ichneumonid wasps with their eggs to suppress the immune system of host insects [57]. These

genes may have originated from a horizontal transfer event during an attempted parasitization.

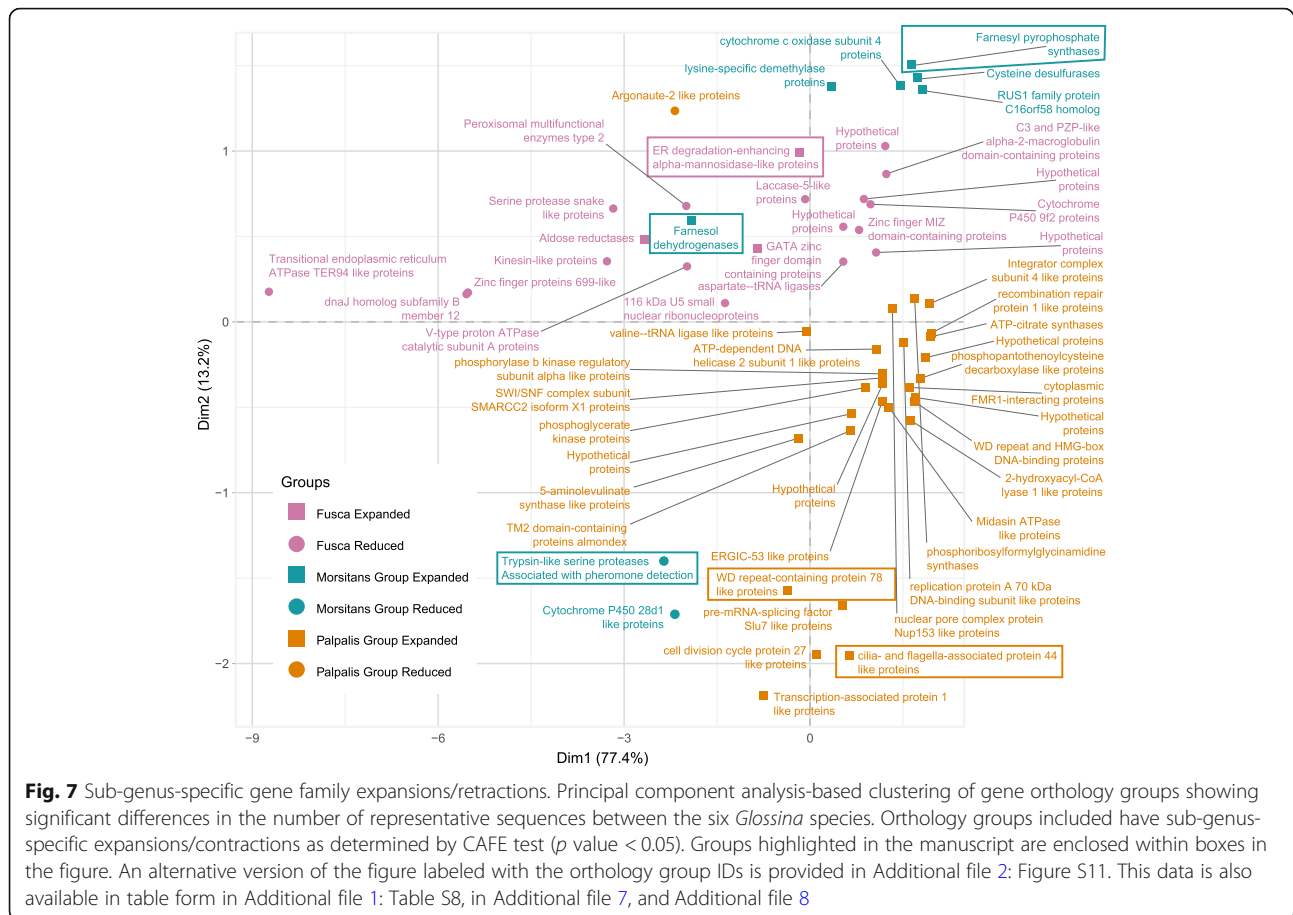
Another feature of note is the abundance of putative proteins with predicted helicase activity. Of the 603 genes with significant hits, 64 (10.5%) are homologous to characterized helicases. Functional enrichment analysis confirms the enrichment of helicase activity in this gene set. These proteins are associated with the production of small RNA's (miRNAs, siRNAs, and piRNAs) which mediate post-transcriptional gene expression and the defensive response against viruses and transposable elements. Of the 64 genes, 41 were homologous to the armitage (*armi*) helicase. Recent work in *Drosophila* shows that *armi* is a reproductive tissue-specific protein and is responsible for binding and targeting mRNAs for processing into piRNAs by the PIWI complex [58]. The reason for the accumulation of this class of genes within the *Palpalis* sub-genus is unknown. However, given the association of these proteins with small RNA production, they could be associated with a defensive response against viral challenges or overactive transposable elements. A similar phenomenon is seen in *Aedes aegypti* where components of the PIWI pathway have been amplified and function outside of the reproductive tissues to generate piRNAs against viral genes [59].

Analysis of gene family variations reveals sub-genus-specific expansions and contractions of genes involved in sperm production and chemosensation

In addition to unique gene families, we identified orthology groups showing significant variation in gene numbers between *Glossina* species. Of interest are groups showing significant sub-genus-specific expansions or contractions, which may represent lineage-specific adaptations. General trends that we observed in these groups show the largest number of gene family expansions within the *Palpalis* sub-genus and the largest number of gene family contractions within *G. brevipalpis* (a member of the *Fusca* sub-genus) (Fig. 7). A second version of the figure labeled with orthology group IDs is available in Additional file 2: Figure S11. The raw data from which this figure was derived can be found in Additional file 1: Table S8 and in Additional file 7 and Additional file 8 (CAFÉ and BLAST analyses data, respectively).

Palpalis sub-genus-specific expansion of sperm-associated genes

Members of the *Palpalis* sub-genus had a total of 29 gene family expansions and 1 contraction relative to the other 4 tsetse species. Of the three sub-genera, this represents the largest number of expansions and parallels with the large number of *Palpalis*-specific orthology groups.



Two gene families expanded within the *Palpalis* group (VBGT00770000031191 and VBGT00190000014373) encode WD repeat-containing proteins. The *Drosophila* orthologs contained within these families (*cg13930*, *dic61B*, *cg9313*, *cg34124*) are testis-specific and associated with cilia/flagellar biosynthesis and sperm production [60]. Alteration/diversification of sperm-associated proteins could explain the split of the *Palpalis* sub-genus from the other *Glossina* and the potential incipient speciation documented between *G. palpalis* and *G. fuscipes* [61].

The *Morsitans* sub-genera shows reductions in chemosensory protein genes

Within the *Morsitans* sub-genus, 6 gene families are expanded and 2 are contracted relative to the other tsetse species. Of interest, 1 of the contracted gene families encodes chemosensory proteins (VBGT00190000010664) orthologous to the CheB and CheA series of proteins in *D. melanogaster*. The genes encoding these proteins are expressed exclusively in the gustatory sensilla of the forelegs of male flies and are associated with the detection of low-volatility pheromones secreted by the female in higher flies

[62]. Of interest is that the number of genes in *G. palpalis* (14), *G. fuscipes* (15), and *G. brevipalpis* (14) are expanded within this family relative to *D. melanogaster* (12), *M. domestica* (10), and *S. calcitrans* (4). However, the *Morsitans* group flies *G. m. morsitans* (7), *G. pallidipes* (7), and *G. austeni* (5) all appear to have lost some members of this family. The functional significance of these changes is unknown. However, it could represent an optimization of the male chemosensory repertoire within the *Morsitans* sub-genus.

In terms of expanded gene families in *Morsitans*, we find two encoding enzymes associated with the terpenoid backbone biosynthesis pathway (VBGT00190000010926—farnesyl pyrophosphate synthase; VBGT00840000047886—farnesol dehydrogenase). This pathway is essential for the generation of precursors required for the synthesis of the insect hormone juvenile hormone (JH). In adult *G. m. morsitans*, JH levels play an important role in regulating nutrient balance before and during pregnancy. High JH titers activate lipid biosynthesis and accumulation in the fat body prior to lactation. During lactation, JH titers fall, resulting in the catabolism and mobilization of stored lipids for use in milk production [63].

Comparative analysis of the immune-associated genes in *Glossina* species reveals specific expansions, contractions, and losses relative to *Musca domestica* and *Drosophila melanogaster*

Tsetse flies are exposed to bacterial, viral, protozoan, and fungal microorganisms exhibiting a broad spectrum of beneficial, commensal, parasitic, and pathogenic phenotypes within their host. Yet, the diversity and intensity of the microbial challenge facing tsetse flies are limited relative to that of related Brachyceran flies such as *D. melanogaster* and *M. domestica* in terms of the level of exposure, microbial diversity, and host-microbe relationships. While tsetse larvae live in a protected environment (maternal uterus) feeding on maternally produced lactation secretions, larval *D. melanogaster* and *M. domestica* spend their entire immature development in rotting organic materials surrounded by and feeding on a diverse array of microbes. The adult stages also differ in that tsetse feed exclusively on blood which exposes

them to a distinct yet limited array of microbial fauna. The immune function and genetic complement of *D. melanogaster* are well characterized and provide the opportunity to compare the constitution of orthologous immune gene sequences between *M. domestica* and the *Glossina* species [64]. Orthology groups containing *Drosophila* genes associated with the “Immune System Process” GO tag (GO:0002376) were selected and analyzed to measure the presence/absence or variance in the number of orthologous sequences in *Glossina* (Fig. 8, Additional file 9, Additional file 10).

Several orthologs within this ontology group are highly conserved across all species and are confirmed participants with the fly’s antimicrobial immune response. These genes include the peptidoglycan recognition proteins (PGRPs) (with the exception of the PGRP SC1+2 genes) [65]; prophenoloxidase 1, 2, and 3 [54]; the reactive oxygen intermediates *dual oxidase* and *peroxiredoxin 5* [66, 67]; and antiviral (RNAi pathway associated) *dicer*

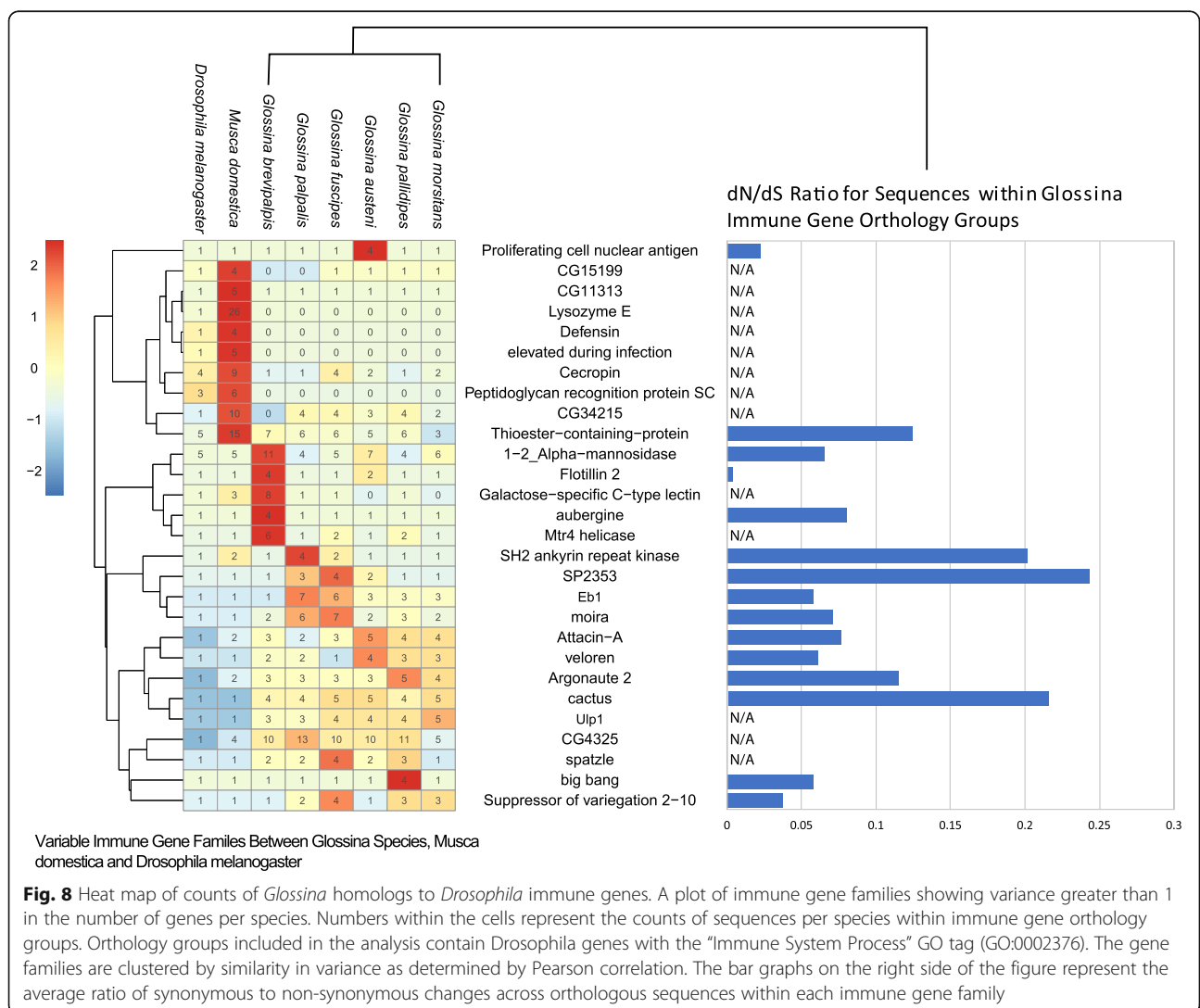


Fig. 8 Heat map of counts of *Glossina* homologs to *Drosophila* immune genes. A plot of immune gene families showing variance greater than 1 in the number of genes per species. Numbers within the cells represent the counts of sequences per species within immune gene orthology groups. Orthology groups included in the analysis contain *Drosophila* genes with the “Immune System Process” GO tag (GO:0002376). The gene families are clustered by similarity in variance as determined by Pearson correlation. The bar graphs on the right side of the figure represent the average ratio of synonymous to non-synonymous changes across orthologous sequences within each immune gene family

2 and *argonaute 2*. The antimicrobial peptide-encoding genes *attacin* (variants A and B) and *cecropin* (variants A1, A2, B, and C) are found within *Glossina* but have diverged significantly (the highest % identity based on blastx comparison = 84%) from closely related fly taxa [68–70]. Analysis of the variance in the numbers of immune gene ortholog/paralogs between the *Glossina* species relative to *M. domestica* and *D. melanogaster* revealed a number of interesting patterns (Fig. 7).

Glossina* species are missing immune gene families present in *D. melanogaster* and *M. domestica

Several gene families are missing within the *Glossina* species although expanded within *M. domestica* (Fig. 7). These include *lysozyme E*, *defensin*, *elevated during infection*, and the *PGRP-SC1+2* gene families. These may be adaptations to the microbe-rich diet and environment in which *M. domestica* larvae and adults exist. The expansion of immune gene families in *M. domestica* relative to *D. melanogaster* was previously documented in the publication of the *M. domestica* genome [71]. However, the added context of the *Glossina* immune gene complement highlights the significance of the expansion of these families relative to their loss in all *Glossina* species. The loss of these families may represent the reduced dietary and environmental exposure to microbial challenge associated with the dramatic differences in tsetse life history.

***Glossina* species show immune gene family expansions associated with the Toll and IMD pathways**

In contrast, we observed several *Glossina* immune-related gene families which are expanded relative to orthologous families in *Drosophila* and *M. domestica* (Fig. 7). Duplications of this nature often reflect evolutionarily important aspects of an organism's biology and, in the case of tsetse, may have resulted from the fly's unique association with parasitic African trypanosomes. Prominent among the expanded immune-related *Glossina* genes are those that encode *Attacin A* and *Attacin B*, which are IMD pathway-produced effector antimicrobial molecules, and *Cactus*, a negative regulator of the Toll signaling pathway. Analysis of the evolutionary rate of these gene families by dN/dS analysis reveals significant variability. We were not able to obtain dN/dS ratios for all families due to the large sequence differences in some family members making an accurate alignment difficult. Whether this is due to the rapid genetic changes or inaccuracies in the gene models remains to be determined and will require additional curation to establish.

However, families with high-quality alignments showed significant variability in their evolutionary rates. *Cactus* is expanded across all *Glossina* species and appears to be evolving rapidly relative to other immune

gene families. This could have significant implications on the regulation of the Toll pathway signaling in immunity and development. The SP2353 gene family is a Laminin G domain-containing protein associated with various binding functions and is associated with negative regulation of immune responses [72]. This gene is primarily expanded in the *Palpalis* sub-genus and is the most rapidly evolving gene family relative to the other representative families. As both of these gene families are associated with negative regulation of immune pathways, it is possible these could be associated with adaptations to obligate symbiosis.

The most highly expanded immune-related gene across *Glossina* species are the orthologs of *Drosophila* CG4325. RNAi-based studies in *Drosophila* indicate that CG4325 is a regulator of both the Toll and IMD signaling pathways [73]. Significant expansion of this gene family in *Glossina* substantiates data that demonstrated the functional importance of the Toll and IMD pathways in tsetse's response to trypanosome challenge [74, 75]. Finally, all six *Glossina* genomes encode multiple copies of *moira*. This gene, which is involved in cell proliferation processes [76], is differentially expressed upon trypanosome infection when compared to uninfected *G. m. morsitans* [77]. In an effort to eliminate parasite infections, tsetse flies produce reactive oxygen intermediates that cause collateral cytotoxic damage [66]. Additionally, trypanosome infection of tsetse's salivary glands induces the expression of fly genes that encode proteins associated with stress and cell division processes, further indicating that parasite infection results in extensive damage to host cells. Expansion of *moira* gene copy number in *Glossina*'s genome may reflect the fly's need to maintain epithelial homeostasis in the face of damage caused by trypanosome infections.

***G. brevipalpis* has a species-specific expansion of immune-associated proteins**

An interesting highlight from this analysis is the identification of a gene expansion associated with alpha-mannosidase activity (VBGT00190000009892). An orthologous *Drosophila* gene (*α-Man-Ia*) is an essential component in the encapsulation response by hemocytes to attack by parasitoid wasps. This enzyme modifies lamellocyte surface glycoproteins to facilitate the recognition and encapsulation of foreign bodies. As described in the *G. m. morsitans* genome paper and here, there is evidence of parasitization by parasitoid wasps in the genomes of these flies in the form of integrated gene sequences homologous to polydnvirus genes [27]. The expansion of these proteins could be an evolutionary response to pressure induced by parasitization although the current status of tsetse-associated parasitoids is unknown.

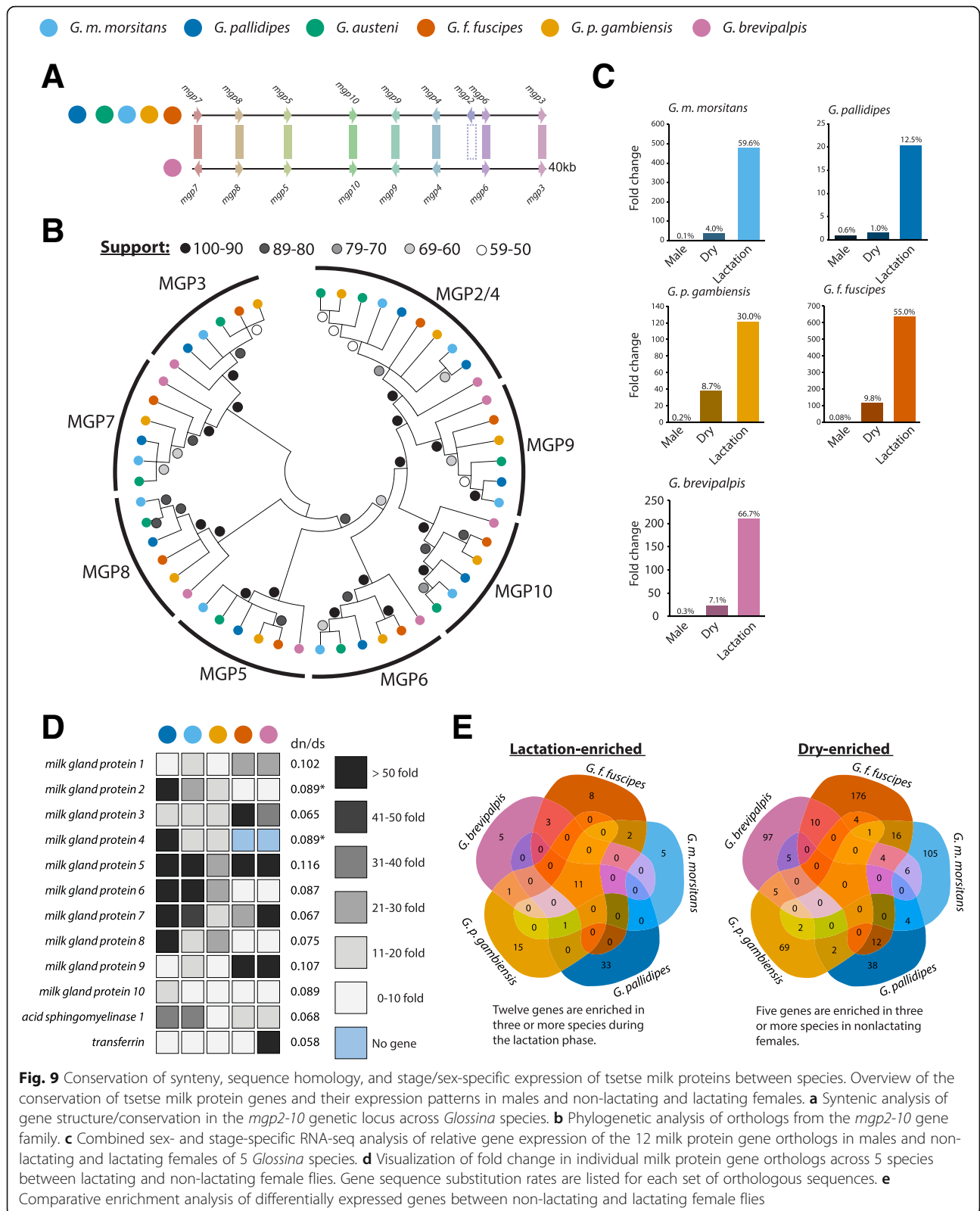


Fig. 9 Conservation of synteny, sequence homology, and stage/sex-specific expression of tsetse milk proteins between species. Overview of the conservation of tsetse milk protein genes and their expression patterns in males and non-lactating and lactating females. **a** Syntenic analysis of gene structure/conservation in the *mgp2-10* genetic locus across *Glossina* species. **b** Phylogenetic analysis of orthologs from the *mgp2-10* gene family. **c** Combined sex- and stage-specific RNA-seq analysis of relative gene expression of the 12 milk protein gene orthologs in males and non-lactating and lactating females of 5 *Glossina* species. **d** Visualization of fold change in individual milk protein gene orthologs across 5 species between lactating and non-lactating female flies. Gene sequence substitution rates are listed for each set of orthologous sequences. **e** Comparative enrichment analysis of differentially expressed genes between non-lactating and lactating female flies

Tsetse reproductive genetics

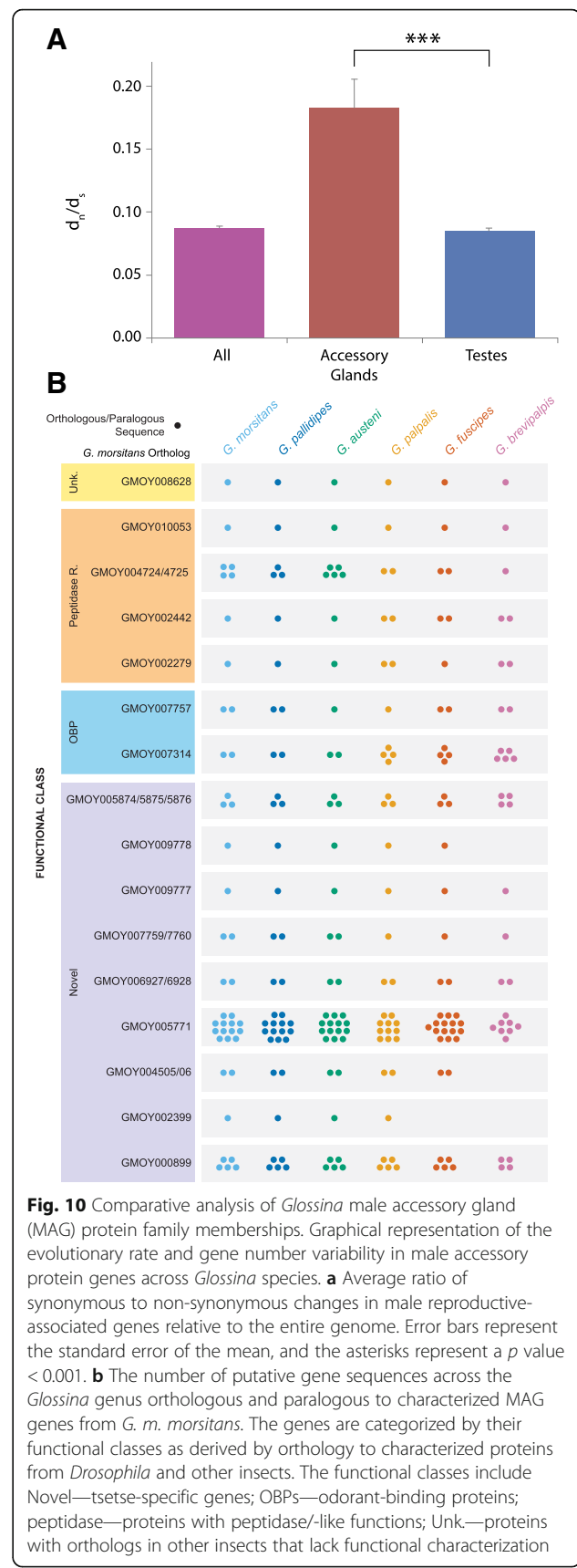
Milk protein genes are universal and tightly conserved in *Glossina* (Fig. 9, Additional file 1: Table S9)

The intrauterine development and nourishment of individual larval offspring are a defining characteristic of the *Hippoboscoidea* superfamily, which includes the *Glossinidae* (tsetse flies), *Hippoboscidae* (keds flies), *Nycteribiidae* (bat flies), and *Streblidae* (bat flies) families [78]. Nutrient provisioning is accomplished by the secretion of a milk-like substance from specialized glands into the uterus where the larval flies consume the milk. Dry weight of tsetse milk is roughly 50% protein and 50% lipids [79]. A compiled list of the milk protein orthologs from the six species of tsetse have been assembled (Additional file 1: Table S9).

Milk protein genes 2-10 (*mgp2-10*) in *G. m. morsitans* are the largest milk protein gene family. These genes are tsetse-specific, lack conserved functional protein domains, and their origin is currently unknown. However, experimental evidence suggests they act as lipid emulsification agents and possible phosphate carrier molecules in the milk [80]. Search for orthologous sequences to these genes revealed 1:1 orthologs to each of the nine genes in the five new *Glossina* species except for *G. brevipalpis* which lacks an orthologous sequence for the *mgp2* gene. These genes are conserved at the levels of both synteny and sequence (Fig. 9a, b). Comparative expression analysis of these genes (and the other characterized milk protein orthologs: *milk gland protein 1*, *acid sphingomyelinase*, and *transferrin* [81, 82]) in male and non-lactating and lactating females shows sex- and lactation-specific expression profiles across the five species for which sex-specific RNA-seq data was available (Fig. 9c, d). Comparison of sequence variation across species for these genes by dN/dS analysis indicates that they are under heavy negative selective pressure (Fig. 9d). Enrichment analysis based on comparison of lactation-based RNA-seq data confirms that these 12 orthologous sequences are enriched in lactating flies across all *Glossina* (Fig. 9e). The *mgp2-10* gene family is a unique and conserved adaptation that appears essential to the evolution of lactation in the *Glossina* genus. Determination of the origins of this protein family requires genomic analyses of other members of the *Hippoboscoidea* superfamily that exhibit viviparity along with other species closely related to this group.

Tsetse seminal protein genes are rapidly evolving and vary in number and sequence conservation between species (Fig. 10, Additional file 1: Table S10)

Recent proteomic analysis of male seminal proteins in *G. m. morsitans* revealed an array of proteins transferred from the male to the female as components of the spermatophore [55]. Cross-referencing of the proteomic



data with tissue-specific transcriptomic analyses of the testes and male accessory glands (MAGs) allowed us to identify the tissues from which these proteins are derived. Many of the MAG-associated proteins are *Glossina*-specific and are derived from gene families with multiple paralogs. These sequences were used to identify and annotate orthologous sequences in the other five *Glossina* species. In contrast to the milk proteins, sequence variance and differences in paralog numbers vary in male reproductive genes between the six *Glossina* species.

This is particularly evident in the genes with MAG-biased/specific expression. MAG-biased/specific genes are represented by 22 highly expressed gene families encoding characterized seminal fluid proteins (SFP). We investigated the evolutionary rate of reproductive genes over-expressed in the MAGs and testes, relative to a set of 5513 *G. m. morsitans* genes, orthologous between the six species (Fig. 10a). The average dN/dS ratio is higher in MAG-biased genes than in testis-biased genes or the entire *Glossina* ortholog gene set suggesting that the MAG genes are under relaxed selective constraints. In addition, we found high heterogeneity in the selective pressure across MAG genes. This is specifically evident in the tsetse-specific genes *GMOY002399*, *GMOY007759*, *GMOY004505*, and *GMOY005874* (a protein with OBP like conserved cysteine residues) as well as the OBP ortholog *GMOY007314*. All five genes encode seminal fluid proteins as confirmed by the proteomic analysis of the spermatophore [55].

In addition to sequence variability, the number of paralogs per species differs as well (Fig. 10b). This is similar to comparative analysis observations in *Anopheles* and *Drosophila* species [83, 84]. This variance is especially evident in *Glossina*-specific protein families (i.e., *GMOY002399*, *GMOY004505/4506*, *GMOY005771*). In particular, there are a large number of gene orthologs/paralogs to the *GMOY005771* gene across all *Glossina* species revealing a large family of MAG genes of unknown function. The number of orthologs/paralogs differs significantly between *Glossina* species. In addition, the two *G. m. morsitans* paralogs *GMOY004724* and *GMOY004725* (predicted peptidase regulators) appear to display a higher number of putative gene duplications in the *Morsitans* sub-genus relative to the *Palpalis* and *Fusca* sub-genera. Conservation appears instead to be more evident across testis genes that code for proteins associated with conserved structural and functional components of sperm. Overall, the comparison of the MAG-biased genes across *Glossina* reveals that this group shows substantial variability in terms of genomic composition and rate of evolution.

This is in agreement with other studies indicating that male accessory proteins evolve at high rates due to the intraspecific competition between males or sexually antagonistic coevolution between males and females [85].

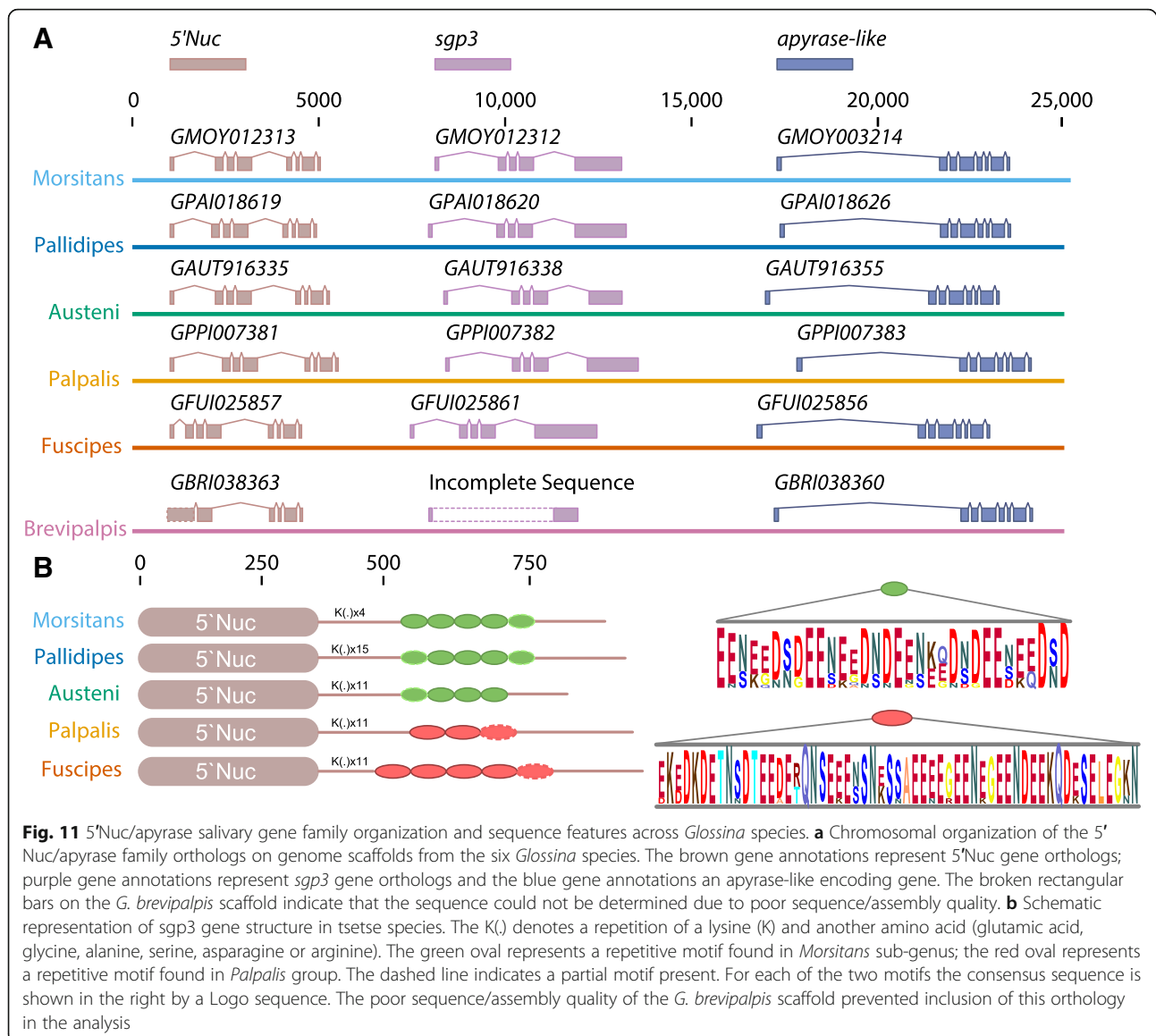
Olfactory-associated protein-coding genes are conserved and reduced in number relative to other Diptera

Comparative analyses of genes responsible for perireceptor olfaction activities revealed high conservation of the repertoire among the six species. The genes appear to scatter across their respective genomes with only a few duplicates occurring in clusters [86]. *Glossina* species expanded loci that include Gr21a (responsible for CO₂ detection) [87], Or67d (mediates *cis*-vaccenyl acetate reception), and Obp83a (thought to be olfactory specific) [88]. The expanded loci suggest the involvement of gene duplication and/or transposition in their emergence [86]. All six species lack sugar receptors likely as a result of tsetse's streamlined blood-feeding behavior. Although our analysis did not reveal major discrepancies among the species, *G. brevipalpis* has lost three key gustatory receptors (Gr58c, Gr66a, and Gr32a) compared to other species. In addition, *G. brevipalpis* showed higher structural gene rearrangements that could be attributed to its evolutionary distance relative to the other tsetse species [89].

A salivary protein gene shows sub-genus-specific repeat motifs (Fig. 11)

Efficient acquisition of a blood meal by tsetse relies on a broad repertoire of physiologically active saliva components inoculated at the bite site. These proteins modulate early host responses, which, in addition to facilitating blood-feeding can also influence the efficacy of parasite transmission [90, 91]. The differences in the competence of different tsetse fly species to develop mature *T. brucei* salivary gland infections may also be correlated with species-specific variations in saliva proteins. Tsetse saliva raises a species-specific IgG response in their mammalian hosts [92]. This response could potentially function as a biomarker to monitor the exposure of host populations to tsetse flies [93].

The *sgp3* gene [94] is characterized in all the tsetse species by two regions: a metallophosphoesterase/5' nucleotidase and a repetitive glutamate/aspartate/asparagine-rich region (Fig. 11a). The complete sequence for this gene from *G. brevipalpis* could not be obtained due to a gap in the sequence. The metallophosphoesterase/5' nucleotidase region is highly conserved between all tsetse species. However, the sequences contain sub-genus-specific (*Morsitans* and *Palpalis*) repeat motifs within the glutamate/aspartate/asparagine region. The motifs differ in size (32 amino acids in the *Morsitans*



group and 57 amino acids in the *Palpalis* group) and amino acid composition (Fig. 11b). Moreover, within each sub-genus, there are differences in the number of repetitive motifs. Within the *Morsitans* group, *G. m. morsitans* and *G. pallidipes* have 5 motifs while *G. austeni* has only 4. In the *Palpalis* group, *G. palpalis* has 3 repetitive motifs and *G. fuscipes* 5. Between the metallo-phosphoesterase/5' nucleotidase and the glutamate/aspartate/asparagine-rich regions, there are a series of amino acids doublets comprising a lysine at the first position followed on the second position by another amino acid (glutamic acid, glycine, alanine, serine, asparagine, or arginine). These differences may account for the differential immunogenic "sub-genus-specific" antibody response caused by Sgp3 in *Morsitans* and *Palpalis* group flies [92].

Comparison of vision-associated Rhodopsin genes reveals conservation of motion tracking receptors and variation in receptors sensitive to blue wavelengths (Fig. 12)

Vision plays an important role in host and mate-seeking by flies within the *Glossina* genus. This aspect of their biology is a critical factor in the optimization and development of trap/target technologies [95, 96]. Analysis of the light-sensitive Rhodopsin proteins across the *Glossina* species reveals orthologs to those described in the *G. m. morsitans* genome (Fig. 12a). The expanded analysis provided by these additional genomes corroborates observations made for the original *G. m. morsitans* genome, including the conservation of the blue-sensitive *Rh5* rhodopsin and the loss of one of the two dipteran UV-sensitive Rhodopsins, *Rh4* [27]. The availability of the new genomes provides complete sequences for an

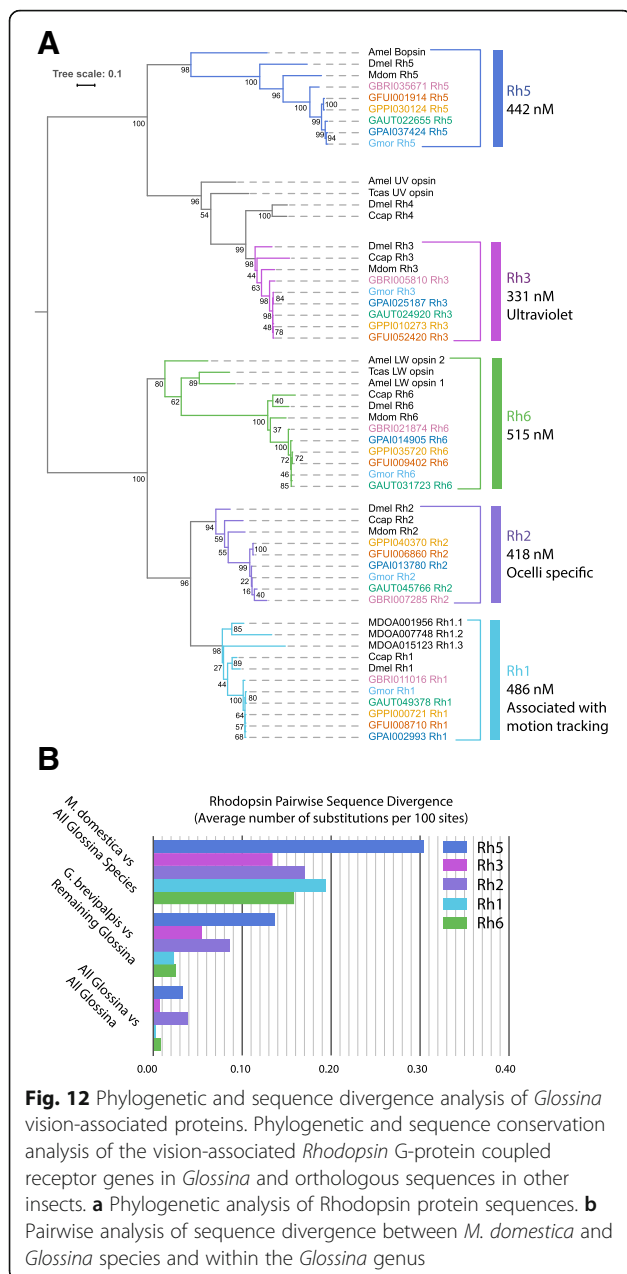


Fig. 12 Phylogenetic and sequence divergence analysis of *Glossina* vision-associated proteins. Phylogenetic and sequence conservation analysis of the vision-associated Rhodopsin G-protein coupled receptor genes in *Glossina* and orthologous sequences in other insects. **a** Phylogenetic analysis of Rhodopsin protein sequences. **b** Pairwise analysis of sequence divergence between *M. domestica* and *Glossina* species and within the *Glossina* genus

additional long wavelength-sensitive Rhodopsin gene, *Rh2*. Prior to this analysis, the recovery of a complete sequence from *G. m. morsitans* was not possible due to the poor sequence quality at its locus.

Rhodopsin protein sequence divergence among the 6 *Glossina* species and *M. domestica* (as an outgroup) was investigated by calculating pairwise sequence divergence. As expected, the average pairwise sequence divergence between *M. domestica* and any *Glossina* species is higher than the maximum sequence divergence among *Glossina* species for any of the 5 investigated Rhodopsin subfamilies, ranging between 0.13 and 0.3 substitutions per 100 sites. Average sequence divergence of *G. brevipalpis* to

other *Glossina* is consistently lower than *Musca* vs *Glossina* but also higher than the average pairwise distances between all other *Glossina*, suggesting the older evolutionary lineage of *G. brevipalpis* (Fig. 12b).

Three interesting aspects emerge in the comparison between subfamilies at the level of sequence divergence between *Glossina* species. The *Rh1* subfamily, which is deployed in motion vision, has the lowest average sequence divergence suggesting the strongest level of purifying selection. *Rh2*, which is expressed in the ocelli, and *Rh5*, which is expressed in color-discriminating inner photoreceptors, are characterized by conspicuously higher-than-average sequence divergence among *Glossina* species. This observation could account for the varying attractivity of trap and targets to different tsetse species.

Conclusions

The comparative genomic analysis of these six *Glossina* species highlights the important aspects of *Glossina* evolution and provides further insights into their unique biology. Additional documentation of other comparative analyses is included in Additional file 1. These include additional information on *Glossina*-specific gene enrichments/expansions/contractions, *Glossina* salivary protein genes, genes encoding neuropeptides and their receptors (Additional file 1: Table S11 and S12), cuticular protein genes (Additional file 1: Table S13, Additional file 11), *Glossina* transcription factor genes and their putative binding sites (Additional file 2: Figure S12, Additional file 12), and peritrophic matrix protein genes (Additional file 1: Table S14). The results derived from the analysis of these genomes are applicable to many aspects of tsetse biology including host seeking, digestion, immunity, metabolism, endocrine function, reproduction, and evolution. This expanded knowledge has important practical relevance. Indeed, tsetse control strategies utilize trapping as a key aspect of population management. These traps use both olfactory and visual stimuli to attract tsetse. The findings of a reduced contingent of olfactory-associated genes and the variability of color sensing Rhodopsin genes provide research avenues into improvements of trap efficacy. A deeper understanding of the important chemosensory and visual stimuli associated with the different species could facilitate the refinement of trap designs for specific species. The findings associated with *Glossina* digestive biology, including the enrichment of proteolysis-associated genes and identification of *Glossina*-specific expansions of immune-associated proteins, provide new insights and avenues of investigation into vector competence and vector/parasite relationships. Analysis of the female and male reproduction-associated genes reveals the differential evolutionary pressures on females and males. The

conservation of female milk proteins across species highlights the fact that this unique biology is optimized and under strong negative evolutionary pressure. In counterpoint, male accessory gland-derived seminal proteins appear to have evolved rapidly between *Glossina* species and with little conservation relative to other Diptera in gene orthology and functional conservation. Tsetse reproduction is slow due to their unique viviparous adaptations, making these adaptations a potential target for the development of new control measures. The knowledge derived from these comparisons provides context and new targets for functional analysis of the genetics and molecular biology of tsetse reproduction. In addition to the practical aspects of the knowledge derived from these analyses, they also provide a look at the genetics underlying the evolution of unique adaptive traits and the resources to develop a deeper understanding of these processes.

Materials and methods

Aim

The aim of these studies was to generate and mine the genomic sequences of six species of tsetse flies with different ecological niches, host preferences, and vectorial capacities. The goals of the analyses performed here are to identify the novel genetic features specific to tsetse flies and to characterize the differences between the *Glossina* species to correlate the genetic changes with phenotypic differences in these divergent species. This was accomplished by the analyses described below.

Glossina strains

All genomes were sequenced from DNA obtained from two to four lines of flies originating from individual pregnant females and their female offspring. Species collections were derived from laboratory strains with varied histories (Additional file 1: Table S1). The *G. pallidipes*, *G. palpalis*, and *G. fuscipes* flies were maintained in the laboratory at the Slovak Academy of Sciences in Bratislava, Slovakia. The *G. brevipalpis* strain was maintained in the Insect Pest Control Laboratory of the Joint FAO/IAEA Division of Nuclear Techniques in Food and Agriculture, Seibersdorf, Austria. Finally, *G. austeni* were obtained from the Tsetse Trypanosomiasis Research Institute in Tanga, Tanzania. Females were given two blood meals supplemented with 20 mg/ml tetracycline to cure them of symbionts to eliminate non-tsetse-derived DNA.

Genomic sequencing and assembly

Total genomic DNA was isolated from female pools for each species. High-quality/high-molecular weight DNA was isolated from individual flies using Genomic-tip purification columns (QIAGEN) and the associated

buffer kit. Samples were treated according to the protocol for tissue-based DNA extraction. The pooled individual DNA isolates were utilized for sequencing on Illumina HiSeq2000 instruments. The sequencing plan followed the recommendations provided in the ALL-PATHS-LG assembler [97]. Using this model, we targeted 45× sequence coverage each of fragments (overlapping paired reads ~180 bp length) and 3 kb paired-end (PE) sequences as well as 5× coverage of 8 kb PE sequences. The first draft assembly scaffold gaps of each species were closed where possible with the mapping of the same species assembly input sequences (overlapping paired reads ~180 bp length) and local gap assembly [98]. Contaminating sequences and contigs 200 bp or less were removed (Table 1). The genome assemblies for each of the six species are available at www.vectorbase.org [99–104].

Scaffold mapping to Muller elements and sex-specific Muller element expression biases

We mapped scaffolds in each *Glossina* spp. genome assembly to chromosomes using homology relationships with *D. melanogaster* (Additional file 1: Table S5). This method exploits the remarkable conservation of chromosome arm (Muller element) gene content across flies [35, 105, 106]. We used the 1:1 orthologs between each *Glossina* species and *D. melanogaster* from OrthoDB [107] to assign scaffolds from each species to Muller elements, applying an approach previously developed for house fly [32]. For each species, a gene was assigned to a Muller element if it was a 1:1 ortholog with a *D. melanogaster* gene. Then, each scaffold was assigned to a Muller element if the majority (>50%) of genes with 1:1 orthologs on that scaffold were assigned to a single Muller element.

We used the RNA-seq data (described below) to compare the gene expression in males and females. Expression comparisons were between male flies and either lactating (L) or non-lactating (NL) females.

Repeat feature annotation

Repeat libraries for each species were generated using RepeatModeler [108]. The resultant libraries were used to annotate the genome with RepeatMasker [109], alongside tandem and low-complexity repeats identified with TRF [110] and DUST [111]. The proportion of the genome covered by different repeat classes is shown in Table 1, Fig. 2, and Additional file 1: Tables S3 and S4. Comparative analysis of TE repeats between species was achieved by clustering the RepeatModeler sequences using Usearch5 [112] with an identity threshold of 80%.

Automated gene annotation

Gene annotation was performed with MAKER [113], using the first two rounds to iteratively improve the training of the ab initio gene predictions derived from the combined Benchmarking Universal Single-Copy Orthologs (BUSCO) [114] and Core Eukaryotic Genes Mapping Approach (CEGMA) [115] HMMs, which were aligned to the genome assemblies using GeneWise [116]. RNA-seq data for each species (described below) were used to build a reference-guided transcriptome assembly with Tophat [117] and Cufflinks [118]. The initial MAKER analysis produced unrealistically high numbers of gene models, so InterProScan [119] and OrthoMCL [120] were used to identify gene predictions which lacked strong evidence. Only the gene models that met one or more of the following criteria were retained: (a) an annotation edit distance < 1 [121], (b) at least one InterPro domain (other than simple coils or signal peptides), and (c) an ortholog in the *Glossina* species complex. This process resulted in a reduction of 12–25% in the number of gene models for each species (Additional file 1: Table S2). Genes from all six species were assigned to 15,038 orthology groups via the Ensembl Compara “GeneTrees” pipeline [122].

For all types of ncRNA except tRNA and rRNA genes, we predicted RNA gene models by aligning sequences from Rfam [123] against the genome using BLASTN [124]. The BLAST results were then used to seed Infernal [125] searches of the aligned regions with the corresponding Rfam covariance models. rRNA genes were predicted with RNAmmer [126] and tRNA genes with tRNAscan-SE [127].

Manual gene annotation

Glossina sequence data and annotation data were loaded into the Apollo [128] community annotation instances in VectorBase [129]. Manual annotations, primarily from a workshop held in Kenya in 2015, underwent both manual and automated quality control to remove incomplete and invalid modifications and then merged with the automated gene set. Gene set versions are maintained at www.vectorbase.org for each organism. All highlighted cells relate to the current gene set version indicated in the table. Statistics for older gene set versions are provided along with the relevant version number.

Genome completeness analysis (BUSCO and CEGMA analyses)

Quality of the genome assembly and training of the ab initio predictors used in the gene prediction pipeline was determined using the diptera_odb9 database which represents 25 Dipteran species and contains a total of 2799 BUSCO (Benchmarking Universal Single-Copy Orthologs)

genes derived from the OrthoDB v9 dataset [114] (Table 2).

Identification of horizontal gene transfer events

All genome sequence files for *G. pallidipes*, *G. palpalis*, *G. fuscipes*, *G. austeni*, and *G. brevipalpis* used for the whole-genome assembly were also introduced into a custom pipeline for the identification of putative horizontal gene transfer (HGT) events between *Wolbachia* and tsetse. *Wolbachia* sequences were filtered out from WGS reads using a combination of MIRA [130] and NextGenMap [131] mapping approaches. The reference sequences used were *wMel* (AE017196), *wRi* (CP001391), *wBm* (AE017321), *wGmm* (AWUH01000000), *wHa* (NC_021089), *wNo* (NC_021084), *wOo* (NC_018267), *wPip* (NC_010981), and the chromosomal insertions A and B in *G. morsitans morsitans*. All filtered putative *Wolbachia*-specific sequences were further examined using blast and custom-made databases.

To identify the chromosomal *Wolbachia* insertions, the following criteria were used: sequences that (relative to the reference genomes) (a) exhibit high homology to the insertion sequences A and B from *G. m. morsitans*, (b) exhibit a high degree of nucleotide polymorphisms (at least 10 polymorphisms/100 bp) with the reference genomes, and (c) contain a high degree of polymorphism coupled with insertions and/or deletions. *Wolbachia*-specific sequences for each *Glossina* species were assembled with MIRA using a de novo approach. For *G. pallidipes*, *G. palpalis*, *G. fuscipes*, and *G. brevipalpis* assembled sequences corresponding only to cytoplasmic *Wolbachia* were identified. Genomic insertions were only observed in assembled sequences from *G. austeni* (Additional file 1: Table S6). The statistics for the *G. austeni* assembled sequences are as follows: N50 4493, N90 1191, and mean contig length 2778 bps. During the process of identifying HGT events in *G. fuscipes*, we also recovered *Spiroplasma* sequences, but none of the recovered sequences was chromosomal.

Whole-genome pairwise alignment

We generated all possible pairwise alignments between the six *Glossina* species (including *G. m. morsitans*) and an outgroup, *M. domestica*, using the Ensembl Compara software pipeline [122]. LASTZ [132] was used to create pairwise alignments, which were then joined to create “nets” representing the best alignment with respect to a reference genome [133]. *G. m. morsitans* was always used as the reference for any alignment of which it was a member; otherwise, the reference genome was randomly assigned. Coverage statistics and configuration parameters for all alignments are available at https://www.vectorbase.org/compara_analyses.html.

Glossina phylogeny prediction

We identified orthologous genes across the six *Glossina* species and six outgroups (*M. domestica*, *D. melanogaster*, *D. ananassae*, *D. grimshawi*, *L. longipalpis*, and *A. gambiae*) by employing a reciprocal-best-hit (RBH) approach in which *G. m. morsitans* was used as a focal species. We identified 286 orthologs with a clear reciprocal relationship among the 12 species. All orthologs were aligned individually using MAFFT [134] and concatenated in a super-alignment of 478,617 nucleotide positions. The nucleotide alignment was translated in the corresponding amino acids and passed through Gblocks [135] (imposing “half allowed gap positions” and leaving the remaining parameters at default) to obtain a dataset of 117,783 amino acid positions. This dataset was used for a maximum likelihood analysis in RAxML [136] employing the LG+G+F model of replacement and for a Bayesian analysis using Phylobayes [137, 138] employing the heterogeneous CAT+G model of replacement. We further performed a coalescent-aware analysis using Astral [139] and the 286 single-gene trees obtained using Raxml [136] and analyzing the alignments at the nucleotide level with the GTR+G model of replacement.

Rate of molecular evolution and selective pressure

We used PAML 4.7 [140] to analyze the rate of molecular evolution and identify heterogeneity in the levels of selective pressure acting across the phylogenetic tree (*(G. morsitans, G. pallidipes), G. austeni*), (*G. fuscipes, G. palpalis*), *G. brevipalpis*). We aligned orthologous gene sequences with PRANK [141], without providing a guide tree, using the tool TranslatorX [142]. Subsequently, to minimize false signals of rapid evolution, we removed the problematic alignment regions using an approach similar to that proposed by [143] implemented in a custom perl script.

We estimated the rate of non-synonymous, dN, and synonymous, dS, substitution over all branches of the phylogenetic tree using the “free-ratio” model, which allows branch-specific levels of selective pressure (i.e., of $\omega = dN/dS$), an additional class of sites under positive selection (M8; model = 0 and NSsites = 8). In these cases, each comparison was tested using a χ^2 test with 2 degrees of freedom. To account for multiple testing, for each set of comparisons, we estimated the false discovery rate (FDR) using the qvalue [144] package implemented in R (R Development Core Team 2009).

Mitochondrial genome analysis and phylogeny

The mtDNA genomes of *G. m. centralis* and *G. brevipalpis* were sequenced using the Illumina HiSeq system, and about 15 kb of mitochondrial sequence of each species was obtained. These sequences were used to identify

the mtDNA sequences within the sequenced tsetse genomes (*G. pallidipes*, *G. m. morsitans*, *G. p. gambiense*, *G. f. fuscipes*, and *G. austeni*) from the available genomic data. Sanger sequencing confirmed the mtDNA genome sequence of each tsetse species. This involved PCR amplification of the whole mtDNA genome using 14 pairs of degenerate primers designed to cover the whole mitochondrial genomes of the sequence species (Additional file 1: Table S15). The PCR products were sent for Sanger sequencing. The sequences obtained by Sanger and Illumina sequencing for each species were assembled using the SegMan program from the laser-gene software package (DNASTar Inc., Madison, USA). The phylogenetic analysis based on these sequences was performed using the maximum likelihood method with the MEGA 6.0 [145].

Synteny analysis

The synteny analysis was derived from whole-genome alignments performed as follows using tools from the UCSC Genome Browser [146]. The LASTZ software package (version 1.02.00) generated the initial pairwise sequence alignments with the following parameters: $E = 30$, $H = 2000$, $K = 3000$, $L = 2200$, $O = 400$, and the default substitution matrix. From these alignments, Kent’s toolbox (version 349) [146] was used to generate chain and nets (higher-level abstractions of pairwise sequence alignments) with the following parameters: $-verbose = 0$, $-minScore = 3000$, and $-linearGap = \text{medium}$. The multiple alignment format (MAF) files were built with MULTIZ for TBA package (version 01.21.09) [147], using the chains and nets, along with the phylogenetic relationships and distances between species. Using the MAF files, pairwise homologous syntenic blocks (HSBs) were automatically defined using the SyntenyTracker software [148]. Briefly, the SyntenyTracker software defines an HSB as a set of two or more consecutive orthologous markers in homologous regions of the two genomes, such that no other defined HSB is within the region bordered by these markers. There are two exceptions to this rule: the first involves single orthologous markers not otherwise defined within HSBs, and the second involves two consecutive singleton markers separated by a distance less than the resolution threshold (10 kb for this analysis). As the 10-kb blocks were too small for visualization in Circos [149], they were aggregated into larger 250-kb histogram blocks, where each 250 kb Circos block shows the fraction of sequence identified as syntenic for a particular species when aligned to *D. melanogaster*. Synteny blocks are available for visualization from the Evolution Highway comparative chromosome browser: <http://eh-demo.ncsa.uiuc.edu/drosophila/>.

Orthology and paralogy inference and analysis

Phylogenetic trees were inferred with the Ensembl Compara “GeneTrees” pipeline [122] using all species from the VectorBase database of arthropod disease vectors [129]. The trees include 33 non-*Glossina* species, such as *D. melanogaster*, which act as outgroup comparators. All analyses are based on the VectorBase April 2016 version of the phylogenetic trees. Representative proteins from all genes were clustered and aligned, and trees showing orthologs and paralogs were inferred with respect to the NCBI taxonomy tree (<http://www.ncbi.nlm.nih.gov/taxonomy>).

The 15,038 predicted gene trees containing *Glossina* sequences were parsed to quantify the trees based on their constituent species. Raw tree files (Additional file 3) were parsed using a custom PERL script which is accessible via Github (https://github.com/attardog/Comp_Genomics_Scripts/releases/latest) to determine gene counts for representative Dipteran species for each gene tree [150]. Count data were imported into Excel and filtered using pivot tables to categorize orthology groups based upon species constitution (Additional file 4).

The orthology groups were broken into cohorts based on the phylogenetic composition of species within each group. The *Glossina* containing orthology groups were categorized as follows: common to Diptera (including the Nematocera sub-order), Brachycera sub-order-specific, *Glossina* genus-specific, *Glossina* sub-genus-specific (*Morsitans* and *Palpalis*), or *Glossina* species-specific. Each category is subdivided into two groups, universal groups that contain representative sequences from all species within the phylogenetic category or partial orthology groups containing sequences from some but not all members of the phylogenetic category. *Glossina* gene IDs and associated FASTA sequences associated with groups of interest were extracted using a custom Perl script for gene ontology analysis (https://github.com/attardog/Comp_Genomics_Scripts/releases/latest) [150].

Gene Ontology analysis

Gene-associated GO terms were obtained from the VectorBase annotation database via the BioMart interface. Genes from *Glossina* genus and sub-genus-specific orthology groups were isolated and tested for enrichment of GO terms. Analysis for GO terms for enrichment was performed with the R package “topGO.” The enriched genes were separated into species-specific lists compared against the entirety of predicted protein-coding genes from the respective species. Significance of enrichment was determined using Fisher’s exact test (Additional file 1: Table S7, and Additional file 5).

Identification and analysis of gene expansions/contractions

Gene trees containing orthologs/paralogs representing each of the six *Glossina* species were analyzed to identify sub-genus-associated gene expansions/contractions. Gene trees were considered for analysis if the variance in the number of orthologs/paralogs between the six species was greater than 2. Variable gene trees were tested for phylogenetic significance relative to the predicted *Glossina* phylogeny using the CAFE software package [151] to reject potentially inaccurate variance predictions due to erroneous gene annotations. Gene trees with a CAFE score of < 0.05 were considered significant (Additional file 1: Table S8, Additional file 7, Additional file 8).

Sequences from gene trees satisfying the variance and CAFE thresholds were extracted with a custom PERL script (https://github.com/attardog/Comp_Genomics_Scripts/releases/latest) and analyzed by BLASTP analysis [124] against an insect-specific subset of the NCBI NR database. Gene trees were annotated with the most common description associated with the top BLAST hits of its constituent sequences. Gene trees were subjected to PCA analysis in R using the FactoMineR and Factoextra packages using species-specific gene counts as input data. The results were plotted and annotated with their associated BLAST-derived descriptions.

Immune gene and dN/dS analysis

Orthology groups containing *Drosophila* genes associated with the GO term (GO:0002376) were queried from the dataset used in the orthology/paralogy analysis described above. Gene counts were visualized using the pHeatmap package in R (<https://cran.r-project.org/web/packages/pheatmap/index.html>). The associated dendrograms are generated based on the similarity of the gene counts per family as determined by Pearson correlation. The dN/dS values were derived from the molecular evolution and selective pressure analysis described above.

RNA-seq data

Total RNA was isolated for each of the six tsetse species from whole male and whole female (non-lactating and lactating) for RNA-seq library construction. Poly(A)+ RNA was isolated, then measured with an Agilent Bioanalyzer for quality. Samples were considered to be of high quality if they showed intact ribosomal RNA peaks and lacked signs of degradation. Samples passing quality control were used to generate non-normalized cDNA libraries with a modified version of the Nu-GEN Ovation® RNASeq System V2 (<https://www.nugen.com/products/ovation-rna-seq-system-v2>). We sequenced each cDNA library (0.125 lane) on an

Illumina HiSeq 2000 instrument (~ 36 Gb per lane) at 100 bp in length.

RNA-seq analyses were conducted based on methods described in Benoit et al. [80], Rosendale et al. [152], and Scolari et al. [55] with slight modifications. RNA-seq datasets were acquired from whole males, whole dry females, and whole lactating females. The SRA numbers for each of the libraries are listed in (Additional file 1: Table S16) [153–169].

RNA-seq datasets were quality controlled using the FastQC (Babraham Bioinformatics) software package. Each set was trimmed/cleaned with CLC Genomics (Qiagen), and quality was re-assessed with FastQC. Each dataset was mapped to the predicted genes from each *Glossina* genome with CLC Genomics. Each read required at least 95% similarity over 50% of length with three mismatches allowed. Transcripts per million (TPM) was used as a proxy for gene expression. Relative transcript abundance differences were determined as the TPM in one sample relative to the TPM of another dataset (e.g., male/lactating female). A proportion-based statistical analysis [170] followed by Bonferroni correction at 0.05 was used to identify genes with significant sex- and stage-specific transcript enrichment. This stringent statistical analysis was used as only one replicate was available for each treatment.

Enriched transcripts in lactating and dry transcriptomes from the species examined were compared to orthologous sequences in *G. m. morsitans* [27]. Overlap was determined by comparison of the enrichment status of orthologous sequences in the *Glossina* species tested. The results of this analysis are visualized in a Venn diagram (<http://bioinformatics.psb.ugent.be/webtools/Venn/>). Determination of dN/dS values and production of phylogenetic trees was conducted with the use of DataMonkey [171, 172] for dN/dS analyses and MEGA5 for alignment and tree construction [173].

Cuticular protein analysis

The predicted peptide sequences from each species were analyzed by BLASTp analysis [124] against characteristic sequence motifs derived from several families of cuticle proteins [174]. Predicted cuticle proteins were further analyzed with CutProtFam-Pred, a cuticle protein prediction tool described in Ioannidou et al. [175], to assign genes to specific families of cuticle proteins. To find the closest putative homolog to cuticle protein genes in *Glossina*, genes were searched (BLASTp) against Refseq protein database from the National Center for Biotechnology Information (NCBI). The protein sequences with the lowest *e* value were considered the closest putative homologs (Additional file 11).

Transcription factor identification and annotation

Likely transcription factors (TFs) were identified by scanning the amino acid sequences of predicted protein-coding genes for putative DNA binding domains (DBDs). When possible, we predicted the DNA-binding specificity of each TF using the procedures described in Weirauch et al. [176]. Briefly, we scanned all protein sequences for putative DBDs using the 81 Pfam [177] models listed in Weirauch and Hughes [178] and the HMMER tool [179], 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). We then aligned the resulting DBD sequences within each family using clustalOmega [180], with default settings. For protein pairs with multiple DBDs, each DBD was aligned separately. From these alignments, we calculated the sequence identity of all DBD sequence pairs (i.e., the percent of AA residues that are identical across all positions in the alignment). Using previously established sequence identity thresholds for each family [176], we mapped the predicted DNA binding specificities by simple transfer. For example, the DBD of the *G. austeni* GAUT024062-PA protein is identical to the DBD of *D. melanogaster* mirr (FBgn0014343). Since the DNA-binding specificity of mirr has already been experimentally determined, and the cutoff for Homeodomain family of TFs is 70%, we can infer that GAUT024062-PA will have the same binding specificity as mirr. All associated data can be found in Additional file 12.

Additional files

- Additional file 1:** Supplementary text and supplementary tables. (DOCX 186 kb)
- Additional file 2:** Supplemental Figures S1-S12 and associated captions. (PDF 2260 kb)
- Additional file 3:** Raw ensemble dipteran orthology data. (TXT 1842 kb)
- Additional file 4:** Orthology group species composition data. (XLSX 3807 kb)
- Additional file 5:** Top GO results of genus- and sub-genus-specific gene ID. (XLSX 21 kb)
- Additional file 6:** Results from the BLAST analysis of *Glossina*-specific serine proteases. (XLSX 22 kb)
- Additional file 7:** Results of a CAFÉ analysis of the variable orthology groups to identify their closest dipteran homologs. (XLSX 224 kb)
- Additional file 8:** Summary of the BLAST results of all members of the variable orthology groups to identify their closest dipteran homologs. (XLSX 679 kb)
- Additional file 9:** Counts of *Glossina* and *Musca* orthologs/paralogs of *Drosophila melanogaster* immunity-associated genes. (XLSX 37 kb)
- Additional file 10:** List of gene IDs and full names for *Drosophila* immunity-associated genes and the orthologous/paralogous genes identified in *Musca* and the *Glossina* species. (XLSX 122 kb)

Additional file 11: List of putative cuticle protein genes identified in all *Glossina* species. (XLSX 51 kb)

Additional file 12: Transcription Factor Data. (ZIP 883 kb)

Additional file 13: Review history. (DOXC 31 kb)

Abbreviations

AAT: Animal African trypanosomiasis; DBD: DNA-binding domain; GO: Gene Ontology; HAT: Human African trypanosomiasis; HGT: Horizontal gene transfer; MAG: Male accessory gland; MGP: Milk gland protein; miRNA: MicroRNA; mtDNA: Mitochondrial DNA; OBP: Odorant-binding protein; OG: Orthology group; piRNA: Piwi-interacting RNA; rRNA: Ribosomal RNA; SFP: Seminal fluid protein; siRNA: Small interfering RNA; tRNA: Transfer RNA

Acknowledgements

We thank the production sequencing group of McDonnell Genome Institute at Washington University for the library construction, sequencing, and data curation. Great thanks to the members of the Comparative Genomics workshop held at the Biotechnology Research Institute - Kenya Agricultural and Livestock Research Organization, Kikuyu, Kenya, including Muna Abry, Willis Adero, Erick Aroko, Joel Bargul, Tania Bishola, Lorna Jemosop Chebon, Appolinaire Djikeng, John Irungu, Evelyn Kamau, Christine Kamidi, Caleb Kibet, Esther Kimani, Kelvin Kimenyi, Mathuriin Koffi, Benard Kulohoma, Clarence Manger, Abraham Mayoke, David Mburu, Grace Murilla, Mary Murithi, Ramadhan Mwakubambanya, Sarah Mwangi, Nelly Ndungu, Joyce Njuguna, Benson Nyambega, Faith Obange, Samuel Ochieng, Edwin Ogola, Owallah (Martin) Ogwang, Sylvance Okoth, Luicer Olubayo, Irene Onyango, Fred Osowo, David Price, Martin Rono, Sharon Towett, Kelvin Wachiuri, Kevin Wamae, and Mark Wamalwa.

The workshop was sponsored by the D43 TW007391 award from the Fogarty International Center to SA and was facilitated by the Yale School of Public Health, Kenya Agricultural and Livestock Research Organization (KALRO), International Centre of Insect Physiology and Ecology (ICIPE), South African National Bioinformatics Institute (SANBI), International Livestock Research Institute (ILRI), and Biosciences Eastern and Central Africa.

Review history

The review history is available as Additional file 13

Authors' contributions

JBB, GMA, HGM, JC, AA-S, and RR are the annotation group leaders. GMA, WCW, SA, and MJL are the project leaders. DL, EL, GLM, and MB were responsible for the analysis of whole genomic sequences and database management. ECJ, JBB, and VM were responsible for the BUSCO and female reproductive gene analysis. DM, POM, and RWM were responsible for the chemosensory gene analysis. AJR and DWF were responsible for the cuticular protein gene analysis. GMA and JEA were responsible for the gene orthology and expansion analyses. WCW, CT, PM, and RKW were responsible for the genome sequencing, assembly, and analysis. GT and KB were responsible for the horizontal gene transfer analysis. AV, BLW, JW, and RB were responsible for the immune gene analysis. ARM, FS, and GS were responsible for the male reproductive gene analysis. AMMA, IM, and AGP were responsible for the mitochondrial DNA sequence analysis. LO and OR-S were responsible for the molecular evolution and phylogenetic analyses. HGM, JC, and LS were responsible for the neuropeptide and G protein-coupled receptor analysis. MF was responsible for the rhodopsin gene analysis. RMW was responsible for the orthology and comparative genomics advice and manuscript editing. AA-S and CR were responsible for the peritrophin gene analysis. SA and MJL were responsible for the project conception. SA was responsible for the project funding. SA and WCW were responsible for the project management. PT, SA, and AMMA were responsible for the provision of the experimental material. JVDA, IM, GC, and XZ were responsible for the salivary protein gene analysis. RR was responsible for the symbiont-associated gene analysis. MTS, DML, and VPPEL were responsible for the syntenic analysis of genomes. MTW was responsible for the transcription factor and DNA-binding motif prediction. AHV and WJM were responsible for the transposable element analysis. RPM was responsible for the X chromosome and sex-linked expression analysis. All authors read and approved the final manuscript.

Funding

This work was supported by NIH Grants D43 TW007391, U01AI115648, R01AI051584, R03TW008413, and R03TW009444 to SA; Grant R21AI109263 to GA and SA from NIH-NIAID; Grant U54HG003079 from NIH-NHGRI to RKW and SA, McDonnell Genome Institute at Washington University School of Medicine; partial funding from the National Research Foundation to HGM (Grant # 10924); and Swiss National Science Foundation grant PP00P3_170664 to RMW. This research was partially supported by the Slovak Research and Development Agency under contract no. APW-15-0604 entitled "Reduction of fecundity and trypanosomiasis control of tsetse flies by the application of sterile insect techniques and molecular methods."

Availability of data and materials

The genomes, transcriptomes, and predicted protein-coding sequences are available from VectorBase and are included within the references [100–105]. The raw RNA-seq datasets generated and/or analyzed during the current study are available from the NCBI SRA database repository at the following link <https://www.ncbi.nlm.nih.gov/sra/SRP158014> and are listed within the reference list [157–173]. All data generated during the analyses of these datasets are included in this published article and its supplementary information files.

Ethics approval and consent to participate

Not applicable

Consent for publication

Not applicable

Competing interests

All authors declare that they have no competing interests.

Author details

¹Institute of Biological, Environmental and Rural Sciences, Aberystwyth University, Aberystwyth, Ceredigion, UK. ²Department of Biochemistry, Biotechnology Research Institute - Kenya Agricultural and Livestock Research Organization, Kikuyu, Kenya. ³CAS Center for Influenza Research and Early-warning (CASCIRE), Chinese Academy of Sciences, Beijing, China. ⁴Center for Autoimmune Genomics and Etiology and Divisions of Biomedical Informatics and Developmental Biology, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA. ⁵Laboratoire Evolution, Genomes, Comportement, Ecologie, CNRS, IRD, Univ. Paris-Sud, Université Paris-Saclay, Gif-sur-Yvette, France. ⁶VectorBase, European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Cambridge, Cambridgeshire, UK. ⁷Department of Biological Sciences, Florida International University, Miami, Florida, USA. ⁸Department of Sustainable Ecosystems and Bioresources, Research and Innovation Centre, Fondazione Edmund Mach, San Michele all'Adige, TN, Italy. ⁹School of Life Sciences, Fudan University, Shanghai, China. ¹⁰Department of Life Sciences, Imperial College London, London, UK. ¹¹Biomedical Sciences, Institute of Tropical Medicine, Antwerp, Belgium. ¹²Molecular Biology and Bioinformatics Unit, International Center for Insect Physiology and Ecology, Nairobi, Kenya. ¹³Insect Pest Control Laboratory, Joint FAO/IAEA Division of Nuclear Techniques in Food & Agriculture, Vienna, Austria. ¹⁴Centre for Geographic Medicine Research Coast, Kenya Medical Research Institute, Kilifi, Kenya. ¹⁵Department of Biology - Functional Genomics and Proteomics Group, KU Leuven, Leuven, Belgium. ¹⁶Department of Vector Biology, Liverpool School of Tropical Medicine, Merseyside, Liverpool, UK. ¹⁷Department of Cell and Developmental Biology, Medical University of Vienna, Vienna, Austria. ¹⁸Department of Biology, Mount St. Joseph University, Cincinnati, OH, USA. ¹⁹Department of Comparative Biomedical Sciences, Royal Veterinary College, London, UK. ²⁰Institute of Zoology, Slovak Academy of Sciences, Bratislava, Slovakia. ²¹Laboratory of Microbiology, Parasitology and Hygiene, University of Antwerp, Antwerp, Belgium. ²²Department of Entomology and Nematology, University of California, Davis, Davis, CA, USA. ²³Department of Biological Sciences, University of Cape Town, Rondebosch, South Africa. ²⁴Department of Biological Sciences, University of Cincinnati, Cincinnati, OH, USA. ²⁵Department of Biology and Biochemistry, University of Houston, Houston, TX, USA. ²⁶Department of Ecology & Evolution, Swiss Institute of Bioinformatics, University of Lausanne, Lausanne, Switzerland. ²⁷Centre for Biotechnology and Bioinformatics, University of Nairobi, Nairobi, Kenya. ²⁸Department of Environmental and Natural Resources Management,

University of Patras, Agrinio, EtoIoakarnania, Greece. ²⁹Department of Biology and Biotechnology, University of Pavia, Pavia, Italy. ³⁰Schools of Medicine and Dentistry, University of Plymouth, Plymouth, UK. ³¹Department of Animal Systematics, Ústav zoologie SAV, Scientica, Ltd, Bratislava, Slovakia. ³²McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO, USA. ³³Department of Biological Sciences, Wayne State University, Detroit, MI, USA. ³⁴Department of Biology, West Virginia University, Morgantown, WV, USA. ³⁵Department of Epidemiology of Microbial Diseases, Yale School of Public Health, New Haven, CT, USA. ³⁶Bond Life Sciences Center, University of Missouri, Columbia, MO, USA.

Received: 31 January 2019 Accepted: 22 July 2019

Published online: 02 September 2019

References

- Lyons M. The colonial disease. A social history of sleeping sickness in northern Zaire, 1900–1940. Cambridge UK: Cambridge University Press; 1992.
- Odiit M, Coleman PG, Liu WC, McDermott JJ, Fevre EM, Welburn SC, Woolhouse ME. Quantifying the level of under-detection of *Trypanosoma brucei rhodesiense* sleeping sickness cases. *Tropical Med Int Health*. 2005;10:840–9.
- Franco JR, Cecchi G, Priotto G, Paone M, Diarra A, Grout L, Simarro PP, Zhao W, Argaw D. Monitoring the elimination of human African trypanosomiasis: update to 2016. *PLoS Negl Trop Dis*. 2018;12:e0006890.
- Franco JR, Simarro PP, Diarra A, Ruiz-Postigo JA, Jannin JG. The journey towards elimination of gambiense human African trypanosomiasis: not far, nor easy. *Parasitology*. 2014;141:748–60.
- Steelman CD. Effects of external and internal arthropod parasites on domestic livestock production. *Annu Rev Entomol*. 1976;21:155–78.
- Jordan A. Trypanosomiasis control and African rural development. London: Longman; 1986.
- Budd LT. Tsetse and Trypanosomiasis Research and Development since 1980: an economic analysis. DFID, Livestock Production Programme, Animal Health Programme/Natural Resources Systems Programme: UK; 1999.
- Alsam M. The effect of the Tsetse Fly on African development. *Am Econ Rev*. 2015;105:382–410.
- Opigo J, Woodrow C. NECT trial: more than a small victory over sleeping sickness. *Lancet*. 2009;374:7–9.
- Mesu V, Kalonji WM, Bardonneau C, Mordt OV, Blesson S, Simon F, Delhomme S, Bernhard S, Kuziena W, Lubaki JF, et al. Oral fexinidazole for late-stage African *Trypanosoma brucei* gambiense trypanosomiasis: a pivotal multicentre, randomised, non-inferiority trial. *Lancet*. 2018;391:144–54.
- Buscher P, Deborggraeve S. How can molecular diagnostics contribute to the elimination of human African trypanosomiasis? *Expert Rev Mol Diagn*. 2015;15:607–15.
- Anene BM, Onah DN, Nawa Y. Drug resistance in pathogenic African trypanosomes: what hopes for the future? *Vet Parasitol*. 2001;96:83–100.
- Geerts S, Holmes PH, Eisler MC, Diall O. African bovine trypanosomiasis: the problem of drug resistance. *Trends Parasitol*. 2001;17:25–8.
- Lehane M, Alfarouk I, Bucheton B, Camara M, Harris A, Kaba D, Lumbala C, Peka M, Rayaisse JB, Waiswa C, et al. Tsetse control and the elimination of Gambian sleeping sickness. *PLoS Negl Trop Dis*. 2016;10:e0004437.
- Solano P, Torr SJ, Lehane MJ. Is vector control needed to eliminate gambiense human African trypanosomiasis? *Front Cell Infect Microbiol*. 2013;3:33.
- Courtin F, Camara M, Rayaisse JB, Kagbadoune M, Dama E, Camara O, Traore IS, Rouamba J, Peyllhard M, Somda MB, et al. Reducing human-tsetse contact significantly enhances the efficacy of sleeping sickness active screening campaigns: a promising result in the context of elimination. *PLoS Negl Trop Dis*. 2015;9:e0003727.
- Ilboudo H, Jamonneau V, Camara M, Camara O, Dama E, Leno M, Ouendeno F, Courtin F, Sakande H, Sanon R, et al. Diversity of response to *Trypanosoma brucei* gambiense infections in the Forécariah mangrove focus (Guinea): perspectives for a better control of sleeping sickness. *Microbes Infect*. 2011;13:943–52.
- Molyneux DH. Animal reservoirs and Gambian trypanosomiasis. *Ann Soc Belg Med Trop*. 1973;53:605–18.
- Kabayo JP. Aiming to eliminate tsetse from Africa. *Trends Parasitol*. 2002;18:473–5.
- Krafsur ES. Tsetse flies: genetics, evolution, and role as vectors. *Infect Genet Evol*. 2009;9:124–41.
- Travassos Santos Dias J. Contribuição para o estudo da sistemática do género *Glossina* Wiedemann 1830 (Insecta, Brachycera, Cyclorhapha, Glossinidae) Proposta para a criação de um novo subgénero. *Garcia de Orta, Ser Zool, Lisboa*. 1987;14:67–78.
- Dyer NA, Lawton SP, Ravel S, Choi KS, Lehane MJ, Robinson AS, Okedi LM, Hall MJ, Solano P, Donnelly MJ. Molecular phylogenetics of tsetse flies (Diptera: Glossinidae) based on mitochondrial (COI, 16S, ND2) and nuclear ribosomal DNA sequences, with an emphasis on the palpalis group. *Mol Phylogenet Evol*. 2008;49:227–39.
- Rogers D, Robinson T. Tsetse distribution. In: Maudlin I, Holmes P, Miles M, editors. *The trypanosomiasis*. Oxford: CAB International; 2004. p. 139–79.
- Moloo SK, Kabata JM, Sabwa CL. A study on the maturation of procyclic *Trypanosoma brucei brucei* in *Glossina morsitans centralis* and *G. brevipalpis*. *Med Vet Entomol*. 1994;8:369–74.
- Motloang M, Masumu J, Mans B, Van den Bossche P, Latif A. Vector competence of *Glossina austeni* and *Glossina brevipalpis* for *Trypanosoma congolense* in KwaZulu-Natal, South Africa. *Onderstepoort J Vet Res*. 2012;79:E1–6.
- Aksoy S, Berriman M, Hall N, Hattori M, Hide W, Lehane MJ. A case for a *Glossina* genome project. *Trends Parasitol*. 2005;21:107–11.
- IGGI. Genome sequence of the tsetse fly (*Glossina morsitans*): vector of African trypanosomiasis. *Science*. 2014;344:380–6.
- Bernardi G. The neoselectionist theory of genome evolution. *Proc Natl Acad Sci U S A*. 2007;104:8385–90.
- Denisov G, Walenz B, Halpern AL, Miller J, Axelrod N, Levy S, Sutton G. Consensus generation and variant detection by Celera assembler. *Bioinformatics*. 2008;24:1035–40.
- Gooding RH, Krafsur ES. Tsetse genetics: contributions to biology, systematics, and control of tsetse flies. *Annu Rev Entomol*. 2005;50:101–23.
- Petersen FT, Meier R, Kutty SN, Wiegmann BM. The phylogeny and evolution of host choice in the Hippoboscoidea (Diptera) as reconstructed using four molecular markers. *Mol Phylogenet Evol*. 2007;45:111–22.
- Meisel RP, Scott JG, Clark AG. Transcriptome differences between alternative sex determining genotypes in the house fly, *Musca domestica*. *Genome Biol Evol*. 2015;7:2051–61.
- Brelsfoard C, Tsiamis G, Falchetto M, Gomulski L, Telleria E, Alam U, Ntountoumis E, Scolari F, Swain M, Takac P, et al. *Wolbachia* symbiont genome sequence and extensive chromosomal insertions present in the host *Glossina morsitans morsitans* genome. *PLoS Negl Trop Dis*. 2014;8:e2728.
- Muller HJ. Bearings of the 'Drosophila' work on systematics. In: Huxley J, editor. *The new systematics*. Oxford: Clarendon; 1940. p. 185–268.
- Vicoso B, Bachtrog D. Numerous transitions of sex chromosomes in Diptera. *PLoS Biol*. 2015;13:e1002078.
- Schaeffer SW, Bhutkar A, McAllister BF, Matsuda M, Matzkin LM, O'Grady PM, Rohde C, Valente VL, Aguade M, Anderson WW, et al. Polytenic chromosomal maps of 11 *Drosophila* species: the order of genomic scaffolds inferred from genetic and physical maps. *Genetics*. 2008;179:1601–55.
- Willhoelt U. Fluorescence in situ hybridization of ribosomal DNA to mitotic chromosomes of tsetse flies (Diptera: Glossinidae: Glossina). *Chromosom Res*. 1997;5:262–7.
- Papa F, Windbichler N, Waterhouse RM, Cagnetti A, D'Amato R, Persampieri T, Lawniczak MKN, Nolan T, Papathanos PA. Rapid evolution of female-biased genes among four species of *Anopheles* malaria mosquitoes. *Genome Res*. 2017;27:1536–48.
- Meisel RP, Connallon T. The faster-X effect: integrating theory and data. *Trends Genet*. 2013;29:537–44.
- Charlesworth B, Coyne JA, Barton NH. The relative rates of evolution of sex chromosomes and autosomes. *Am Nat*. 1987;130:113–46.
- Mank JE, Vicoso B, Berlin S, Charlesworth B. Effective population size and the faster-X effect: empirical results and their interpretation. *Evolution*. 2010;64:663–74.
- Meisel RP. Towards a more nuanced understanding of the relationship between sex-biased gene expression and rates of protein-coding sequence evolution. *Mol Biol Evol*. 2011;28:1893–900.
- Larracuent AM, Sackton TB, Greenberg AJ, Wong A, Singh ND, Sturgill D, Zhang Y, Oliver B, Clark AG. Evolution of protein-coding genes in *Drosophila*. *Trends Genet*. 2008;24:114–23.
- Leung W, Shaffer CD, Reed LK, Smith ST, Barshop W, Dirkes W, Dothager M, Lee P, Wong J, Xiong D, et al. *Drosophila* Muller F elements maintain a distinct set of genomic properties over 40 million years of evolution. *G3: Genes|Genomes|Genetics*. 2015;5:719.

45. Brelsfoard C, Tsiamis G, Falchetto M, Gomulski LM, Telleria E, Alam U, Doudoumis V, Scolari F, Benoit JB, Swain M, et al. Presence of extensive Wolbachia symbiont insertions discovered in the genome of its host *Glossina morsitans morsitans*. *PLoS Negl Trop Dis*. 2014;8:e2728.
46. Doudoumis V, Alam U, Aksoy E, Abd-Alla AM, Tsiamis G, Brelsfoard C, Aksoy S, Bourtzis K. Tsetse-Wolbachia symbiosis: comes of age and has great potential for pest and disease control. *J Invertebr Pathol*. 2013;112(Suppl):S94–103.
47. Wu DD, Wang GD, Irwin DM, Zhang YP. A profound role for the expansion of trypsin-like serine protease family in the evolution of hematophagy in mosquito. *Mol Biol Evol*. 2009;26:2333–41.
48. Gorman MJ, Paskewitz SM. Serine proteases as mediators of mosquito immune responses. *Insect Biochem Mol Biol*. 2001;31:257–62.
49. LaFlamme BA, Ram KR, Wolfner MF. The *Drosophila melanogaster* seminal fluid protease "seminase" regulates proteolytic and post-mating reproductive processes. *PLoS Genet*. 2012;8:e1002435.
50. Sirot LK, Findlay GD, Sitnik JL, Frasher D, Avila FW, Wolfner MF. Molecular characterization and evolution of a gene family encoding both female- and male-specific reproductive proteins in *Drosophila*. *Mol Biol Evol*. 2014;31:1554–67.
51. Hamilton JV, Munks RJ, Lehane SM, Lehane MJ. Association of midgut defensin with a novel serine protease in the blood-sucking fly *Stomoxys calcitrans*. *Insect Mol Biol*. 2002;11:197–205.
52. Larter NK, Sun JS, Carlson JR. Organization and function of *Drosophila* odorant binding proteins. *Elife*. 2016;5:1–22.
53. Leal WS. Odorant reception in insects: roles of receptors, binding proteins, and degrading enzymes. *Annu Rev Entomol*. 2013;58:373–91.
54. Benoit JB, Vigneron A, Broderick NA, Wu Y, Sun JS, Carlson JR, Aksoy S, Weiss BL. Symbiont-induced odorant binding proteins mediate insect host hematopoiesis. *Elife*. 2017;6:1–24.
55. Scolari F, Benoit JB, Michalkova V, Aksoy E, Takac P, Abd-Alla AM, Malacrida AR, Aksoy S, Attardo GM. The Spermatothore in *Glossina morsitans morsitans*: insights into male contributions to reproduction. *Sci Rep*. 2016;6:20334.
56. Doudoumis V, Blow F, Saridaki A, Augustinos A, Dyer NA, Goodhead I, Solano P, Rayaisse JB, Takac P, Mekonnen S, et al. Challenging the *Wigglesworthia*, *Sodalis*, *Wolbachia* symbiosis dogma in tsetse flies: *Spiroplasma* is present in both laboratory and natural populations. *Sci Rep*. 2017;7:4699.
57. Tschopp A, Riedel M, Kropf C, Nentwig W, Klopstein S. The evolution of host associations in the parasitic wasp genus *Ichneumon* (Hymenoptera: Ichneumonidae): convergent adaptations to host pupation sites. *BMC Evol Biol*. 2013;13:74.
58. Pandey RR, Homolka D, Chen KM, Sachidanandam R, Fauvarque MO, Pillai RS. Recruitment of Armitage and Yb to a transcript triggers its phased processing into primary piRNAs in *Drosophila* ovaries. *PLoS Genet*. 2017;13:e1006956.
59. Miesen P, Joosten J, van Rij RP. PIWIs go viral: arbovirus-derived piRNAs in vector mosquitoes. *PLoS Pathog*. 2016;12:e1006017.
60. Avidor-Reiss T, Maer AM, Koundakjian E, Polyanovsky A, Keil T, Subramaniam S, Zuker CS. Decoding cilia function: defining specialized genes required for compartmentalized cilia biogenesis. *Cell*. 2004;117:527–39 2018/02/04.
61. Ravel S, de Meeus T, Dujardin JP, Zeze DG, Gooding RH, Dsoufou I, Sane B, Cuny G, Solano P. The tsetse fly *Glossina palpalis palpalis* is composed of several genetically differentiated small populations in the sleeping sickness focus of Bonon, Cote d'Ivoire. *Infect Genet Evol*. 2007;7:116–25.
62. Starostina E, Xu A, Lin H, Pikieli CW. A *Drosophila* protein family implicated in pheromone perception is related to Tay-Sachs GM2-activator protein. *J Biol Chem*. 2009;284:585–94.
63. Baumann AA, Benoit JB, Michalkova V, Mireji PO, Attardo GM, Moulton JK, Wilson TG, Aksoy S. Juvenile hormone and insulin suppress lipolysis between periods of lactation during tsetse fly pregnancy. *Mol Cell Endocrinol*. 2013;372:30–41.
64. Buchon N, Silverman N, Cherry S. Immunity in *Drosophila melanogaster*—from microbial recognition to whole-organism physiology. *Nat Rev Immunol*. 2014;14:796–810.
65. Dziarski R, Gupta D. The peptidoglycan recognition proteins (PGRPs). *Genome Biol*. 2006;7:232.
66. Vigneron A, Aksoy E, Weiss BL, Bing X, Zhao X, Awuoch E, O'Neill MB, Wu Y, Attardo GM, Aksoy S. A fine-tuned vector-parasite dialogue in tsetse's cardia determines peritrophic matrix integrity and trypanosome transmission success. *PLoS Pathog*. 2018;14:e1006972.
67. MacLeod ET, Maudlin I, Darby AC, Welburn SC. Antioxidants promote establishment of trypanosome infections in tsetse. *Parasitology*. 2007;134:827–31.
68. Hao Z, Kasumba I, Lehane MJ, Gibson WC, Kwon J, Aksoy S. Tsetse immune responses and trypanosome transmission: implications for the development of tsetse-based strategies to reduce trypanosomiasis. *Proceedings of the National Academy of Sciences, USA*. 2001;98:12648–53.
69. Aksoy S, Weiss BL, Attardo GM. Trypanosome transmission dynamics in tsetse. *Curr Opin Insect Sci*. 2014;3:43–9.
70. Hu C, Aksoy S. Innate immune responses regulate trypanosome parasite infection of the tsetse fly *Glossina morsitans morsitans*. *Mol Microbiol*. 2006;60:1194–204.
71. Scott JG, Warren WC, Beukeboom LW, Bopp D, Clark AG, Giers SD, Hediger M, Jones AK, Kasai S, Leichter CA, et al. Genome of the house fly, *Musca domestica* L., a global vector of diseases with adaptations to a septic environment. *Genome Biol*. 2014;15:466.
72. Cronin SJ, Nehme NT, Limmer S, Liegeois S, Pospisilik JA, Schramek D, Leibbrandt A, Simoes Rde M, Gruber S, Puc U, et al. Genome-wide RNAi screen identifies genes involved in intestinal pathogenic bacterial infection. *Science*. 2009;325:340–3.
73. Valanne S, Myllymaki H, Kallio J, Schmid MR, Kleino A, Murumagi A, Airaksinen L, Kotipelto T, Kaustio M, Ulvila J, et al. Genome-wide RNA interference in *Drosophila* cells identifies G protein-coupled receptor kinase 2 as a conserved regulator of NF-kappaB signaling. *J Immunol*. 2010;184:6188–98.
74. Lehane MJ, Aksoy S, Gibson W, Kerhornou A, Berriman M, Hamilton J, Soares MB, Bonaldo MF, Lehane S, Hall N. Adult midgut expressed sequence tags from the tsetse fly *Glossina morsitans morsitans* and expression analysis of putative immune response genes. *Genome Biol*. 2003;4:R63.
75. Aksoy E, Vigneron A, Bing X, Zhao X, O'Neill M, Wu YN, Bangs JD, Weiss BL, Aksoy S. Mammalian African trypanosome VSG coat enhances tsetse's vector competence. *Proc Natl Acad Sci U S A*. 2016;113:6961–6.
76. Nakamura K, Ida H, Yamaguchi M. Transcriptional regulation of the *Drosophila moira* and *Osa* genes by the DREF pathway. *Nucleic Acids Res*. 2008;36:3905–15.
77. Gloria-Soria A, Dunn WA, Yu X, Vigneron A, Lee K-Y, Li M, Weiss BL, Zhao H, Aksoy S, Caccione A. Uncovering genomic regions associated with *Trypanosoma* infections in wild populations of the tsetse fly *Glossina fuscipes*. *G3: Genes|Genomes|Genetics*. 2018;8:887–97.
78. Meier R, Kotrba M, Ferrar P. Oviviviparity and viviparity in the Diptera. *Biol Rev Camb Philos Soc*. 1999;74:199–258.
79. Cmelik SHW, Bursell E, Slack E. Composition of gut contents of third-instar tsetse larvae (*Glossina morsitans* Westwood). *Comp Biochem Physiol*. 1969;29:447–53.
80. Benoit JB, Attardo GM, Michalkova V, Krause TB, Bohova J, Zhang Q, Baumann AA, Mireji PO, Takac P, Denlinger DL, et al. A novel highly divergent protein family identified from a viviparous insect by RNA-seq analysis: a potential target for tsetse fly-specific abortifacients. *PLoS Genet*. 2014;10:e1003874.
81. Benoit JB, Attardo GM, Michalkova V, Takac P, Bohova J, Aksoy S. Sphingomyelinase activity in mother's milk is essential for juvenile development: a case from lactating tsetse flies. *Biol Reprod*. 2012;87(17):1–10.
82. Guz N, Attardo GM, Wu Y, Aksoy S. Molecular aspects of transferrin expression in the tsetse fly (*Glossina morsitans morsitans*). *J Insect Physiol*. 2007;53:715–23.
83. Neafsey DE, Waterhouse RM, Abai MR, Aganezov SS, Alekseyev MA, Allen JE, Amon J, Arca B, Arensburg P, Artemov G, et al. Mosquito genomics. Highly evolvable malaria vectors: the genomes of 16 *Anopheles* mosquitoes. *Science*. 2015;347:1258522.
84. Wong A, Turchin MC, Wolfner MF, Aquadro CF. Evidence for positive selection on *Drosophila melanogaster* seminal fluid protease homologs. *Mol Biol Evol*. 2008;25:497–506.
85. Findlay GD, MacCoss MJ, Swanson WJ. Proteomic discovery of previously unannotated, rapidly evolving seminal fluid genes in *Drosophila*. *Genome Res*. 2009;19:886–96.
86. Macharia R, Mireji P, Murungi E, Murilla G, Christoffels A, Aksoy S, Masiga D. Genome-wide comparative analysis of chemosensory gene families in five tsetse fly species. *PLoS Negl Trop Dis*. 2016;10:e004421.
87. Obiero GFO, Mireji PO, Nyanjom SRG, Christoffels A, Robertson HM, Masiga DK. Odorant and gustatory receptors in the tsetse fly *Glossina morsitans morsitans*. *PLoS Negl Trop Dis*. 2014;8:e2663.
88. Liu R, Lehane S, He X, Lehane M, Hertz-Fowler C, Berriman M, Pickett JA, Field LM, Zhou JJ. Characterisations of odorant-binding proteins in the tsetse fly *Glossina morsitans morsitans*. *Cell Mol Life Sci*. 2010;67:919–29.

89. Rio RV, Symula RE, Wang J, Lohs C, Wu YN, Snyder AK, Bjornson RD, Oshima K, Biehl BS, Perna NT, et al. Insight into the transmission biology and species-specific functional capabilities of tsetse (Diptera: Glossinidae) obligate symbiont *Wigglesworthia*. *mBio*. 2012;3:1–13.
90. Caljon G, Van Reet N, De Trez C, Vermeersch M, Perez-Morga D, Van Den Abbeele J. The dermis as a delivery site of *Trypanosoma brucei* for tsetse flies. *PLoS Pathog*. 2016;12:e1005744.
91. Caljon G, Van Den Abbeele J, Stijlemans B, Coosemans M, De Baetselier P, Magez S. Tsetse fly saliva accelerates the onset of *Trypanosoma brucei* infection in a mouse model associated with a reduced host inflammatory response. *Infect Immun*. 2006;74:6324–30.
92. Zhao X, Silva TL, Cronin L, Savage AF, O'Neill M, Nerima B, Okedi LM, Aksoy S. Immunogenicity and serological cross-reactivity of saliva proteins among different tsetse species. *PLoS Negl Trop Dis*. 2015;9:e0004038.
93. Dama E, Cornelie S, Bienvenu Somda M, Camara M, Kambire R, Courtin F, Jamonneau V, Demetere E, Seveno M, Bengaly Z, et al. Identification of *Glossina palpalis gambiensis* specific salivary antigens: towards the development of a serologic biomarker of human exposure to tsetse flies in West Africa. *Microbes Infect*. 2013;15:416–27.
94. Van Den Abbeele J, Caljon G, Dierick JF, Moens L, De Ridder K, Coosemans M. The *Glossina morsitans* tsetse fly saliva: general characteristics and identification of novel salivary proteins. *Insect Biochem Mol Biol*. 2007;37:1075–85.
95. Lindh JM, Goswami P, Blackburn RS, Arnold SE, Vale GA, Lehane MJ, Torr SJ. Optimizing the colour and fabric of targets for the control of the tsetse fly *Glossina fuscipes fuscipes*. *PLoS Negl Trop Dis*. 2012;6:e1661.
96. Green CH, Cosens D. Spectral responses of the tsetse fly, *Glossina morsitans morsitans*. *J Insect Physiol*. 1983;29:795–800.
97. 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*. 2011; 108:1513–8.
98. Tsai IJ, Otto TD, Berriman M. Improving draft assemblies by iterative mapping and assembly of short reads to eliminate gaps. *Genome Biol*. 2010;11:R41.
99. Aksoy S, Warren WC, Lawson D, Attardo G: Genome assembly for *Glossina brevipalpis*. University W. Vectorbase; 2019. GCA_000671755.1, <https://www.vectorbase.org/organisms/glossina-brevipalpis>, Access Date - 22 Oct 2018.
100. Aksoy S, Warren WC, Lawson D, Attardo G: Genome assembly for *Glossina fuscipes*. University W. Vectorbase; 2019. GCA_000671735.1, <https://www.vectorbase.org/organisms/glossina-fuscipes>, Access Date - 22 Oct 2018.
101. Aksoy S, Warren WC, Lawson D, Attardo G: Genome assembly for *Glossina palpalis*. University W. Vectorbase; 2019. GCA_000818775.1, <https://www.vectorbase.org/organisms/glossina-palpalis>, Access Date - 22 Oct 2018.
102. Aksoy S, Warren WC, Lawson D, Attardo G: Genome assembly for *Glossina austeni*. University W. Vectorbase; 2019. GCA_000688735.1, <https://www.vectorbase.org/organisms/glossina-austeni>, Access Date - 22 Oct 2018.
103. Aksoy S, Warren WC, Lawson D, Attardo G: Genome assembly for *Glossina pallidipes*. University W. Vectorbase; 2019. GCA_000688715.1, <https://www.vectorbase.org/organisms/glossina-pallidipes>, Access Date - 22 Oct 2018.
104. Berriman M, Aksoy S, Lawson D: Genome assembly for *Glossina morsitans morsitans*. Institute WTS. Vectorbase; 2010. GCA_001077435.1, <https://www.vectorbase.org/organisms/glossina-morsitans>, 25 Jun 2018.
105. Weller GL, Foster GG. Genetic maps of the sheep blowfly *Lucilia cuprina*: linkage-group correlations with other dipteran genera. *Genome*. 1993;36:495–506.
106. Foster TJ, Davis MA, Roberts DE, Takeshita K, Kleckner N. Genetic organization of transposon Tn10. *Cell*. 1981;23:201–13.
107. Zdobnov EM, Tegenfeldt F, Kuznetsov D, Waterhouse RM, Simao FA, Ioannidis P, Seppely 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.
108. Smit A, Hubley R: RepeatModeler Open-1.0; 2008–2015 <http://www.repeatmasker.org>, Access date: 29 May 2014.
109. Smit A, Hubley R, Green P: RepeatMasker Open-3.0; 1996–2010. <http://www.repeatmasker.org>, Access date: 5 Feb 2014.
110. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res*. 1999;27:573–80.
111. Morgulis A, Gertz EM, Schaffer AA, Agarwala R. A fast and symmetric DUST implementation to mask low-complexity DNA sequences. *J Comput Biol*. 2006;13:1028–40.
112. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*. 2010;26:2460–1.
113. Cantarel BL, Korf I, Robb SM, Parra G, Ross E, Moore B, Holt C, Sanchez Alvarado A, Yandell M. MAKER: an easy-to-use annotation pipeline designed for emerging model organism genomes. *Genome Res*. 2008;18:188–96.
114. Waterhouse RM, Seppey M, Simao FA, Manni M, Ioannidis P, Kloutchnikov G, Kriventseva EV, Zdobnov EM. BUSCO applications from quality assessments to gene prediction and phylogenomics. *Mol Biol Evol*. 2018;35:543–8.
115. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in eukaryotic genomes. *Bioinformatics*. 2007;23:1061–7.
116. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. *Genome Res*. 2004;14:988–95.
117. Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics*. 2009;25:1105–11.
118. Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL, Rinn JL, Pachter L. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and cufflinks. *Nat Protoc*. 2012;7:562–78.
119. Jones P, Binns D, Chang HY, 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.
120. Li L, Stoeckert CJ Jr, Roos DS. OrthoMCL: identification of ortholog groups for eukaryotic genomes. *Genome Res*. 2003;13:2178–89.
121. Yandell M, Ence D. A beginner's guide to eukaryotic genome annotation. *Nat Rev Genet*. 2012;13:329–42.
122. Herrero J, Muffato M, Beal K, Fitzgerald S, Gordon L, Pignatelli M, Vilella AJ, Searle SM, Amode R, Brent S, et al. Ensembl comparative genomics resources. *Database (Oxford)*. 2016;2016:1–17.
123. Nawrocki EP, Burge SW, Bateman A, Daub J, Eberhardt RY, Eddy SR, Floden EW, Gardner PP, Jones TA, Tate J, Finn RD. Rfam 12.0: updates to the RNA families database. *Nucleic Acids Res*. 2015;43:D130–7.
124. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. BLAST+: architecture and applications. *BMC bioinformatics*. 2009;10:421.
125. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics*. 2013;29:2933–5.
126. Lagesen K, Hallin P, Rodland EA, Staerfeldt HH, Rognes T, Ussery DW. RNAmmer: consistent and rapid annotation of ribosomal RNA genes. *Nucleic Acids Res*. 2007;35:3100–8.
127. Lowe TM, Eddy SR. tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res*. 1997;25:955–64.
128. Lee E, Helt GA, Reese JT, Munoz-Torres MC, Childers CP, Buels RM, Stein L, Holmes IH, Elisk CG, Lewis SE. Web Apollo: a web-based genomic annotation editing platform. *Genome Biol*. 2013;14:R93.
129. Giraldo-Calderon GI, Emrich SJ, MacCallum RM, Maslen G, Dialynas E, Topalis P, Ho N, Gesing S, VectorBase C, Madey G, et al. VectorBase: an updated bioinformatics resource for invertebrate vectors and other organisms related with human diseases. *Nucleic Acids Res*. 2015;43:D707–13.
130. Chevreur B, Pfisterer T, Drescher B, Driesel AJ, Muller WE, Wetter T, Suhai S. Using the miraEST assembler for reliable and automated mRNA transcript assembly and SNP detection in sequenced ESTs. *Genome Res*. 2004;14: 1147–59.
131. Sedlaczek FJ, Rescheneder P, von Haeseler A. NextGenMap: fast and accurate read mapping in highly polymorphic genomes. *Bioinformatics*. 2013;29:2790–1.
132. Harris R. Improved pairwise alignment of genomic DNA. State College, PA: The Pennsylvania State University, College of Engineering; 2007.
133. Kent WJ, Baertsch R, Hinrichs A, Miller W, Haussler D. Evolution's cauldron: duplication, deletion, and rearrangement in the mouse and human genomes. *Proc Natl Acad Sci U S A*. 2003;100:11484–9.
134. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol*. 2013; 30:772–80.
135. Castresana J. Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Mol Biol Evol*. 2000;17:540–52.
136. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics*. 2014;30:1312–3.
137. Lartillot N, Philippe H. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Mol Biol Evol*. 2004; 21:1095–109.
138. Lartillot N, Brinkmann H, Philippe H. Suppression of long-branch attraction artefacts in the animal phylogeny using a site-heterogeneous model. *BMC Evol Biol*. 2007;7(Suppl 1):S4.

139. Chandonia JM, Hon G, Walker NS, Lo Conte L, Koehl P, Levitt M, Brenner SE. The ASTRAL compendium in 2004. *Nucleic Acids Res.* 2004;32:D189–92.
140. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol.* 2007;24:1586–91.
141. Loytynoja A, Goldman N. An algorithm for progressive multiple alignment of sequences with insertions. *Proc Natl Acad Sci U S A.* 2005;102:10557–62.
142. Abascal F, Zardoya R, Telford MJ. TranslatorX: multiple alignment of nucleotide sequences guided by amino acid translations. *Nucleic Acids Res.* 2010;38:W7–13.
143. Han XY, Sizer KC, Thompson EJ, Kabanja J, Li J, Hu P, Gomez-Valero L, Silva FJ. Comparative sequence analysis of *Mycobacterium leprae* and the new leprosy-causing *Mycobacterium lepromatosis*. *J Bacteriol.* 2009;191:6067–74.
144. Storey KB. Life in the slow lane: molecular mechanisms of estivation. *Comp Biochem Physiol A Mol Integr Physiol.* 2002;133:733–54.
145. Tamura K, Stecher G, Peterson D, Filipinski A, Kumar S. MEGA6: molecular evolutionary genetics analysis version 6.0. *Mol Biol Evol.* 2013;30:2725–9.
146. Kuhn RM, Haussler D, Kent WJ. The UCSC genome browser and associated tools. *Brief Bioinform.* 2013;14:144–61.
147. Blanchette M, Kent WJ, Riemer C, Elnitski L, Smit AF, Roskin KM, Baertsch R, Rosenbloom K, Clawson H, Green ED, et al. Aligning multiple genomic sequences with the threaded blockset aligner. *Genome Res.* 2004;14:708–15.
148. Donthu R, Lewin HA, Larkin DM. SyntenyTracker: a tool for defining homologous synteny blocks using radiation hybrid maps and whole-genome sequence. *BMC Res Notes.* 2009;2:148.
149. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA. Circos: an information aesthetic for comparative genomics. *Genome Res.* 2009;19:1639–45.
150. Attardo GM: Scripts for orthology group analysis and FASTA sequence extraction. 1.0 edition: GitHub; 2018. https://github.com/attardog/Comp_Genomics_Scripts/releases/latest, Access date: 15 May 2019.
151. De Bie T, Cristianini N, Demuth JP, Hahn MW. CAFE: a computational tool for the study of gene family evolution. *Bioinformatics.* 2006;22:1269–71.
152. Rosendale AJ, Romick-Rosendale LE, Watanabe M, Dunlevy ME, Benoit JB. Mechanistic underpinnings of dehydration stress in the American dog tick revealed through RNA-Seq and metabolomics. *J Exp Biol.* 2016;219:1808–19.
153. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Non-lactating whole Female *G. austeni* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRS686473, <https://www.ncbi.nlm.nih.gov/sra/SRX682983>, Access date: 24 Aug 2014.
154. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Whole Male *G. austeni* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRS686445, <https://www.ncbi.nlm.nih.gov/sra/SRX682955>, Access date: 25 Aug 2014.
155. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Lactating whole Female *G. morsitans* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRS430097, <https://www.ncbi.nlm.nih.gov/sra/SRS430097>, Access date: 22 Jul 2015.
156. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Non-lactating whole Female *G. morsitans* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRS430099, <https://www.ncbi.nlm.nih.gov/sra/SRS430099>, Access date: 22 Jul 2015.
157. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Male reproductive tract *G. morsitans* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRS2364381, <https://www.ncbi.nlm.nih.gov/sra/SRS2364381>, Access date: 18 Jul 2017.
158. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Lactating whole Female *G. palpalis gambiensi*s RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698159, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698159>, Access date: 1 Dec 2018.
159. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Non-lactating whole Female *G. palpalis gambiensi*s RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698158, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698158>, Access date: 1 Dec 2018.
160. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Whole Male *G. palpalis gambiensi*s RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698161, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698161>, Access date: 1 Dec 2018.
161. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Lactating whole Female *G. pallidipes* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698160, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698160>, Access date: 1 Dec 2018.
162. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Non-lactating whole Female *G. pallidipes* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698163, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698163>, Access date: 1 Dec 2018.
163. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Whole Male *G. pallidipes* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698162, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698162>, Access date: 1 Dec 2018.
164. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Lactating whole Female *G. fuscipes* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698165, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698165>, Access date: 1 Dec 2018.
165. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Non-lactating whole Female *G. fuscipes* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698167, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698167>, Access date: 1 Dec 2018.
166. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Whole Male *G. fuscipes* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698167, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698167>, Access date: 1 Dec 2018.
167. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Lactating whole Female *G. brevipalpis* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698166, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698166>, Access date: 1 Dec 2018.
168. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Non-lactating whole Female *G. brevipalpis* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698169, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698169>, Access date: 1 Dec 2018.
169. Attardo G, Benoit JB, Michalkova V, Takac P, Aksoy S: Whole Male *G. brevipalpis* RNA-seq. NCBI Sequence Read Archive Database (SRA); 2019. SRR7698168, <https://www.ncbi.nlm.nih.gov/Traces/sra/?run=SRR7698168>, Access date: 1 Dec 2018.
170. Baggerly KA, Deng L, Morris JS, Aldaz CM. Differential expression in SAGE: accounting for normal between-library variation. *Bioinformatics.* 2003;19:1477–83.
171. Pond SLK, Frost SDW. Datamonkey: rapid detection of selective pressure on individual sites of codon alignments. *Bioinformatics.* 2005;21:2531–3.
172. Delpont W, Poon AFY, Frost SDW, Pond SLK. Datamonkey 2010: a suite of phylogenetic analysis tools for evolutionary biology. *Bioinformatics.* 2010;26:2455–7.
173. Tamura K, Peterson D, Peterson N, Stecher G, Nei M, Kumar S. MEGA5: molecular evolutionary genetics analysis using maximum likelihood, evolutionary distance, and maximum parsimony methods. *Mol Biol Evol.* 2011;28:2731–9.
174. Willis JH. Structural cuticular proteins from arthropods: annotation, nomenclature, and sequence characteristics in the genomics era. *Insect Biochem Mol Biol.* 2010;40:189–204.
175. Ioannidou ZS, Theodoropoulou MC, Papandreou NC, Willis JH, Hamodrakas SJ. CutProtFam-Pred: detection and classification of putative structural cuticular proteins from sequence alone, based on profile hidden Markov models. *Insect Biochem Mol Biol.* 2014;52:51–9.
176. 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.
177. Finn RD, Mistry J, Tate J, Coghill 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.
178. Weirauch MT, Hughes TR. A catalogue of eukaryotic transcription factor types, their evolutionary origin, and species distribution. *Subcell Biochem.* 2011;52:25–73.
179. Eddy SR. A new generation of homology search tools based on probabilistic inference. *Genome Inform.* 2009;23:205–11.
180. 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.

Publisher's Note

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