# Research gaps and new insights in the intriguing evolution of Drosophila seminal proteins — Source link ⧉

Juan Hurtado, Almeida F, Belliard S, Santiago Revale ...+1 more authors

Institutions: University of Buenos Aires, International Trademark Association, Wellcome Trust Centre for Human Genetics

Related papers:

- New genes as drivers of phenotypic evolution.

- Genomic Analyses of New Genes and Their Phenotypic Effects Reveal Rapid Evolution of Essential Functions in Drosophila Development

- Origins, evolution, and phenotypic impact of new genes

- Comparative expression profiling reveals widespread coordinated evolution of gene expression across eukaryotes

- Proteomics in evolutionary ecology.

1  **Research gaps and new insights in the intriguing evolution of *Drosophila* seminal**
2  **proteins**

3

4  Hurtado J[*,1,2], Almeida FC[1,2], Belliard SA[3], Revale S[4] and Hasson E[1,2]

5

6  [1] Departamento de Ecología, Genética y Evolución, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos
7  Aires (UBA), CABA, Argentina.

8  [2] Instituto de Ecología, Genética y Evolución de Buenos Aires, Consejo Nacional de Investigaciones Científicas y
9  Técnicas (CONICET), CABA, Argentina.

10  [3] Laboratorio de Insectos de Importancia Agronómica, IGEAF (INTA), GV-IABIMO (CONICET), Hurlingham, Buenos
11  Aires, Argentina.

12  [4] Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, OX3 7BN, UK.

13  [*] Author for Correspondence: Juan Hurtado, jhurtado@ege.fcen.uba.ar

14

15  **Abstract**

16  While the striking effects that seminal fluid proteins (SFPs) exert on females are fairly conserved
17  among Diptera, they exhibit remarkable evolutionary lability. Consequently, most SFPs lack
18  detectable homologs among the repertoire of SFPs of phylogenetically distant species. How such
19  a rapidly changing proteome "manages" to conserve functions across taxa is a fascinating
20  question. However, this and other pivotal aspects of SFPs' evolution remain elusive because
21  discoveries on these proteins have been mainly restricted to the model *D. melanogaster*. Here,
22  we provide an overview of the current knowledge on the inter-specific divergence of *Drosophila*
23  SFPs and compile the increasing amount of relevant genomic information from multiple species.
24  Capitalizing the accumulated knowledge in *D. melanogaster*, we present novel sets of high-
25  confidence SFP candidates and transcription factors presumptively involved in regulating the
26  expression of SFPs. We also address open questions by performing comparative genomic
27  analyses that failed to support the existence of conserved SFPs shared by most dipterans and
28  indicated that gene co-option is the most frequent mechanism accounting for the origin of
29  *Drosophila* SFP-coding genes. We hope our update establishes a starting point to integrate, as
30  more species are assayed for SFPs, further data and thus, to widen the understanding of the
31  intricate evolution of these proteins.

32

36

**Introduction**

During mating, spermatozoa expelled from the testes travel throughout the ejaculatory duct into the female reproductive tract accompanied by a rich repertoire of proteins and peptides known as Seminal Fluid Proteins (SFPs) (reviewed in, e.g., Avila et al. 2011; Avila et al. 2016; Chapman 2008; Ramm 2020). These proteins, likely adapted to sperm competition and fertilization, have been highly studied in *Drosophila melanogaster* (e.g., Civetta & Ranz 2019; Hopkins, Sepil, Bonham et al. 2019; Hopkins, Sepil, Thézénas et al. 2019; Misra & Wolfner 2020; Ravi Ram et al. 2005; Ravi Ram & Ramesh 2003; Ravi Ram & Wolfner 2007; Wigby et al. 2020; Wolfner 2007). Once inside the female, some of these proteins will remain bound to spermatozoa, contributing to sperm functions, and some may even interact with the already stored sperm from previous mates (e.g., Avila et al. 2011; Holman 2009; Misra & Wolfner 2020; Ravi Ram & Wolfner 2007; Singh et al. 2018; Wolfner 2007). Many others instead will interact intimately with female biomolecules in the reproductive tract and other organs, and are capable of changing drastically her physiology and behavior (e.g., Avila et al. 2011; Avila et al. 2016; Avila & Wolfner 2017; Lung & Wolfner 1999; Ravi Ram et al. 2005; Ravi Ram & Wolfner 2007).

In *Drosophila*, decrease of female receptivity to mating, increase of egg production, and conformational modification of the female reproductive organs stand out among the profound changes that SFPs trigger in the female (reviewed in Avila et al. 2016). Given the conflicts of interest between males and females (and between competing males), some of the SFPs effects, while beneficial to the last-mating male, can be detrimental to the female (Chapman et al. 1995; Lung et al. 2002; Mueller et al. 2007; Wigby & Chapman 2005). Thus, rapid antagonistic coevolution is expected between some SFPs and female-derived proteins that interact with them (e.g., Sirot et al. 2014). Nevertheless, other SFPs work synergistically with female biomolecules to facilitate fertilization or progeny production for the mutual benefit of males and females (Avila et al. 2016; Wolfner 2009). Therefore, they are expected to diverge more slowly. In fact, sequence comparisons between closely related *Drosophila* species revealed that some SFPs have evolved extremely fast by positive selection while others are conserved by purifying selection (e.g., Almeida & Desalle 2008; Haerty et al. 2007; Turner & Hoekstra 2008; Wong et al. 2012).

65　The biochemical classes into which SFPs typically fall (e.g., proteases, protease inhibitors, lectins,
66　lipases, and cysteine-rich secretory proteins) seem quite conserved among Diptera, even among
67　animals from different classes (reviewed in, e.g., Avila et al. 2016; Wigby et al. 2020). This
68　suggests that the functional spectrum of SFPs is adaptively restricted at the molecular level.
69　Nonetheless, a striking pattern for the vast majority of SFPs is the lack of detectable homologs
70　among SFPs of phylogenetically distant species (Ahmed-Braimah et al. 2017; Almeida & Desalle
71　2009; Davies & Chapman 2006; Haerty et al. 2007; Mueller et al. 2005). Even though the rapid
72　divergence of some of these proteins may hinder homology detection, the main reason behind
73　this pattern seems to be the rapid turnover (gain and loss) of genes encoding SFPs (seminal
74　genes) (Sirot 2019; Sirot et al. 2014). It remains unknown, however, whether a core of a particular
75　SFPs, playing essential reproductive roles, has been conserved over long evolutionary periods.
76　Neither do we know how new seminal genes arise so frequently, or to what extent regulatory
77　elements of seminal genes are conserved across species.

78　Addressing these broad evolutionary questions requires performing multi-species comparative
79　analyses which, in turn, requires extensive omic information on the seminal proteome of several
80　related species. While most of the achieved findings on SFPs have been restricted to *D.*
81　*melanogaster*, in recent years, the seminal proteome has been characterized in many other
82　species, including Drosophilids. This brings up an opportunity to use the *Drosophila* model to
83　address open questions on SFPs evolution and capitalize the accumulated knowledge in *D.*
84　*melanogaster*.

85　Here, to elucidate some answers, we review the current knowledge on the evolution of
86　*Drosophila* SFPs, compiled genomic data from multiple species, and performed molecular
87　evolutionary analyses using bioinformatic tools. We structured the text into sections, each of
88　which tackles a specific topic by presenting knowledge gaps, new insights, and future
89　perspectives.

90

91　**Identification**

92　In *D. melanogaster*, as in many other dipteran species, the main secretory tissues of the male
93　reproductive system are the accessory glands, a pair of merocrine glands attached to the anterior
94　region of the ejaculatory duct (Avila et al. 2016; Chen 1984; Gillott 1996). While mutant males
95　without accessory glands cannot elicit the normal postmating responses in their female mates
96　(Kalb et al. 1993), it has long been known that ACcessory glands Proteins (ACPs) alone are
97　sufficient for triggering these responses in virgin females (reviewed in Ravi Ram & Wolfner 2007).
98　In fact, the first studies on male reproductive proteins aimed to identify SFPs focusing on the
99　male accessory glands.

100    The very first SFP to be identified was 'Sex Peptide' (SP, also known as Acp70A). It was purified
101    from an HPLC fraction of accessory gland extracts that proved, after being injected into virgin
102    females, to reproduce the well-known postmating responses (Chen et al. 1988). The authors also
103    showed that SP gene is transcribed specifically in the male accessory glands. Afterwards, diverse
104    methods such as Expressed Sequence Tags screening, RT-PCR, subtracting hybridization, and
105    cDNA microarray hybridization allowed the identification of many other genes specifically
106    expressed in the male accessory glands (reviewed in Chapman & Davies 2004). Among those
107    genes, the ones encoding proteins or peptides with a predicted signal peptide—that permits
108    canonical merocrine secretion—were considered as candidate seminal genes (Ravi Ram &
109    Wolfner 2007; Swanson et al. 2001). By 2005, using this double criterion, accessory gland-specific
110    expression and capacity to encode secretory proteins, it was possible to identify ~90 putative
111    seminal genes. Five additional seminal genes—or presumptive seminal genes—were found in
112    other organs of the male reproductive tract: the testes, the ejaculatory duct, and the ejaculatory
113    bulb (Cavener & MacIntyre 1983; Dyanov & Dzitoeva 1995; Kopantseva et al. 1990; Ludwig et al.
114    1991; Lung & Wolfner 2001; Richmond et al. 1980; Saudan et al. 2002; Sheehan et al. 1979).
115    Seven additional candidate genes were identified by mass spectrometry of tryptic peptides from
116    accessory glands secretions (Walker et al. 2006).

117    Until 2008, only 22 of the predicted seminal genes were confirmed, mainly by means of
118    immunological techniques, to be transferred to females during mating (e.g., Aigaki et al. 1991;
119    Bertram et al. 1996; Cho et al. 1999; Coleman et al. 1995; Kopantseva et al. 1990; Lung & Wolfner
120    1999, 2001; Meikle et al. 1990; Ravi Ram et al. 2005; Wong et al. 2008). In 2008, Findlay et al.
121    conducted a proteomic screen that largely extended the list of proven SFPs. The authors used
122    isotopic labeling of the female to distinguish, among proteins isolated from the reproductive tract
123    of newly mated females, between female proteins and proteins transferred from unlabeled
124    males. In this way, they confirmed 75 of the previously predicted SFPs and revealed 63 novel
125    ones. More recently, Sepil et al. (2019) applied quantitative proteomics to identify proteins that
126    after mating become significantly less abundant in male reproductive tissues but more abundant
127    in the female reproductive tract, as expected precisely for SFPs. They also cross-referenced their
128    quantification results with transcriptomic and sequence databases to obtain a list of high-
129    confidence candidate SFPs meeting stringent multiple criteria. Some of these candidates were
130    already known as predicted or confirmed SFPs, while nine were novel discoveries (Sepil et al.
131    2019). While we were concluding this report, Wigby et al. (2020) combined data from these and
132    other proteomic studies to provide a list of 292 *D. melanogaster* SFPs. However, the conditions
133    they evaluated may have been too lax; according to modENCODE [implemented in FlyBase
134    r2020_03 (Graveley et al. 2010; Thurmond et al. 2019)] and FlyAtlas2 (Leader et al. 2018), some
135    of the genes they proposed as novel candidates are not expressed in the male reproductive

136    tissues but in the female (e.g., *FBgn0262536*, *FBgn0262484*, and *FBgn0261989*), and thus, it is
137    not clear that all these genes encode SFPs.

138    According to our bibliographic search, the current number of confirmed—or high-confidence
139    candidate—non-sperm SFPs in *D. melanogaster* [hereafter Known Seminal Proteins (KSPs)] is 173
140    (see source studies in supplementary table S1). Our list includes 1) genes encoding proteins
141    previously confirmed to be transferred by males to females during mating, 2) genes meeting the
142    stringent multiple criteria adopted by Sepil et al. (2019), or 3) those genes more expressed in
143    male reproductive tissues than in any other tissue (according to modENCODE and FlyAtlas2) also
144    encoding secretory proteins found in the mating plug [according to Avila et al. (2015) and Wigby
145    et al. (2020)]. Nonetheless, due to current methodological limitations, some other SFPs probably
146    remain to be discovered. Given the leading role of accessory glands as suppliers of SFPs through
147    merocrine secretion, genes that 1) are strongly expressed in the accessory glands and 2) encode
148    secretory proteins can be considered seminal genes. Based on this expression/secretion (double)
149    criterion, a suitable way of finding new candidate seminal genes may be to search in
150    transcriptomic databases for genes expressed in the male accessory glands and to assess which
151    of those genes encode secretory proteins using *in silico* prediction approaches.

152    Before searching for new candidate seminal genes, we explored to what extent *D. melanogaster*
153    Known Seminal Genes (KSGs) meet the expression/secretion criterion by evaluating two
154    conditions. First, we used the RNA-seq databases modENCODE (implemented in FlyBase
155    r2020_03) and FlyAtlas2 to check which seminal genes are strongly expressed in the accessory
156    glands. Second, we used SignalP-5.0—a deep neural network-based tool that identifies signal
157    peptides and their cleavage sites (Almagro Armenteros et al. 2019; Nielsen et al. 1997)—to
158    evaluate which SFPs have signal peptide required for secretion. Among the 173 KSGs, 159 (93.0%)
159    showed relatively high expression in the accessory glands [> 25 Reads/Fragment Per Kilobase of
160    transcript per Million mapped reads (R/FPKM), which is within the 60-70th percentile] according
161    to one or both databases; 156 (90.2%) encoded a protein with a predicted signal peptide; 151
162    (87.2%) meet both conditions (supplementary table S1), and; 165 (95.3%) meet at least one of
163    them. Most of the few genes not meeting any of these conditions are expressed specifically in
164    the testes. These numbers not only confirm that the vast majority of SFPs are expressed in the
165    accessory glands but also show that their secretion is mainly merocrine (but see Corrigan et al.
166    2014; Leiblich et al. 2012).

167    However, the two conditions we evaluated in the KSPs may be too lax for finding new candidate
168    genes. For instance, accessory glands expression level could be inflated in modENCODE or
169    FlyAtlas2, or SignalP could wrongly predict the presence of a signal peptide. Moreover, a signal
170    peptide would only guarantee translocation into the endoplasmic reticulum followed by signal
171    sequence cleavage. Thus, even if a gene truly meets both conditions, the protein may be retained,

172    for instance, in the endoplasmic reticulum or the Golgi apparatus of accessory glands cells. For
173    these reasons, we decided to evaluate *D. melanogaster* genes for a more restrictive set of six
174    conditions that also relies on the expression/secretion criterion:

175    1) At least 'Very High' expression in the accessory glands (> 100 RPKM, which is within the ~90th
176    percentile) according to modENCODE.

177    2) At least 'Moderately High' expression (> 25 RPKM) and expression enrichment in the accessory
178    glands (relative to other adult tissues) according to modENCODE.

179    3) At least 'Very High' expression in the accessory glands (> 100 FPKM, which is within the ~90th
180    percentile) according to FlyAtlas2.

181    4) At least 'Moderately High' expression (> 25 FPKM) and expression enrichment in the accessory
182    glands (relative to whole adult male flies) according to FlyAtlas2.

183    5) Ability to encode a protein with a signal peptide according to SignalP.

184    6) Ability to encode a secretory protein according to DeepLoc, a prediction algorithm that uses
185    deep neural networks to predict protein localization relying on sequence information (Almagro
186    Armenteros et al. 2017). Unlike SignalP, this software differentiates between 10 subcellular
187    localizations and distinguishes proteins of the extracellular space from proteins of the secretory
188    pathway that are retained in the cell.

189    Genes fulfilling conditions 1 (or 2) and 3 (or 4) are highly (or differentially) expressed in the
190    accessory glands according to different databases, while genes fulfilling conditions 5 and 6 are
191    predicted to encode secretory proteins by different software programs. Therefore, we
192    recognized 219 *D. melanogaster* genes that met conditions 1 (or 2), 3 (or 4), 5, and 6 as seminal
193    gene candidates. These 219 genes included 122 KSGs, 43 previously predicted but unconfirmed
194    seminal genes, and 54 newly identified candidates (fig. 1, supplementary table S1). From the 97
195    candidates that are not among the KSGs, 46 (22 previously predicted seminal genes and 24 novel
196    discoveries) met all six conditions and were dubbed Unconfirmed High Confident Candidates
197    (UHCCs) (fig. 1, table 1).

198    As previously noticed, *D. melanogaster* seminal genes share other quite singular features: a
199    significantly biased location on autosomes, particularly on the second chromosome (Findlay et
200    al. 2008; Ravi Ram & Wolfner 2007), and, on average, high *Ka/Ks* ratios (Ahmed-Braimah et al.
201    2017; Almeida & Desalle 2008; Haerty et al. 2007; Holloway & Begun 2004). The UHCCs resemble
202    KSGs regarding chromosomal location (fig. 2) and *Ka/Ks* ratio (fig. 3). In addition, using the
203    functional annotation tool DAVID (Huang et al. 2009), we performed gene-enrichment analyses

204 for molecular function of both UHCCs and KSGs. These analyses also revealed similarities
205 between these groups of genes: eight out of the nine (89%) Gene Ontology (GO) terms annotated
206 to UHCCs are among the terms annotated to KSGs, and the two most represented GO terms in
207 the UHCCs are among the over-represented terms in the KSGs (table 2). Thus, we will henceforth
208 refer to the 173 KSGs and the 46 UHCCs together (a total of 219 genes) as an updated list of *D.*
209 *melanogaster* seminal genes.

210 Aside from *D. melanogaster*, the only *Drosophila* species in which seminal genes were extensively
211 identified are *D. mojavensis* (Almeida & Desalle 2009; Kelleher et al. 2009; Wagstaff & Begun
212 2005), *D. pseudoobscura* (Karr et al. 2019), *D. simulans* (Begun & Lindfors 2005; Findlay et al.
213 2008; Swanson et al. 2001), *D. virilis* (Ahmed-Braimah et al. 2017), and *D. yakuba* (Begun et al.
214 2006; Findlay et al. 2008). Some (or a few) putative seminal genes were also identified in *D.*
215 *biarmipes* (Imamura et al. 1998), *D. erecta* (Begun et al. 2006), *D. funebris* (Baumann et al. 1975;
216 Schmidt et al. 1989), *D. mayaguana* (Almeida & Desalle 2009), and *D. suzukii* (Ohashi et al. 1991;
217 Schmidt et al. 1993). Given the good recall of the stringent criteria we used here to identify
218 candidates, we think that other *Drosophila* species could be assayed for seminal genes using
219 similar criteria. Thus, further research on transcriptomic data generated from accessory glands
220 would provide enough starting information to identify at low cost seminal genes in many species.

221 However, identifying SFPs in multiple species is only part of the equation. The evolution of the
222 seminal proteome may also diverge through changes in the expression level of seminal genes.
223 Begun and Lindfors (2005) found that transcript abundance of the seminal gene *Acp24A4*
224 (*FBgn0051779*) differs drastically between *D. melanogaster* and its sibling *D. simulans*. Findlay et
225 al. (2009) reported differences between *D. melanogaster*, *D. simulans*, and *D. yakuba* in the
226 expression level and sex-specificity of several seminal genes. Similarly, Ahmed-Braimah et al.
227 (2017) uncovered large differences in seminal transcripts abundance between members of the
228 *virilis* subgroup. Although these studies documented divergence between closely related species
229 for seminal genes at the regulatory level, neither the cis nor the trans regulatory elements have
230 been studied in depth.

231 Transcription is a key control point of gene expression, thus the evolution of transcription factors
232 (TFs) that are expressed in the male accessory glands may explain much of the changes in
233 expression of seminal genes across species. However, most of the accessory glands TFs have yet
234 to be identified. To our knowledge, the only known accessory glands' TFs are the hox gene *Abd-*
235 *B* (*FBgn0000015*), the homeodomain transcription repressor *dve* (*FBgn0020307*), and the paired-
236 rule gene *prd* (*FBgn0003145*), which are required for the normal development of accessory
237 glands and the production of functional ACPs (Gligorov et al. 2013; Minami et al. 2012; Xue &
238 Noll 2002). Nevertheless, these genes encode pleiotropic master regulators involved in the
239 morphogenesis of several organs and may be subjected to strong evolutionary constraints.

240    Therefore, future research focused on the identification of accessory glands TFs will advance our
241    understanding of how seminal genes' expression has evolved.

242    It can be argued that TFs implicated in the regulation of seminal genes' expression (seminal TFs)
243    correlate with seminal genes in transcript abundance. Ayroles et al. (2011) found 224 *D.*
244    *melanogaster* genes that, besides being expressed in male reproductive tissues, showed
245    correlated expression patterns to at least seven KSGs. Therefore, we updated this list to the
246    current release (FlyBase r2020_03) and searched it for accessory glands TFs using an online
247    prediction tool implemented in AnimalTFDB3.0, a comprehensive database of animal TFs (Hu et
248    al. 2019). This first search led to the identification of eight putative seminal TFs, including the
249    known *prd* and genes with unknown function (e.g., *FBgn0034870*, *FBgn0030933*, and
250    *FBgn0028480*). We confirmed that all these candidates are distinctly expressed in the male
251    accessory glands according to both modENCODE (implemented in FlyBase r2020_03) and
252    FlyAtlas2.

253    Expression pattern does not necessarily correlate between seminal genes and seminal TFs. Thus,
254    we made a second search of TFs in a more extensive list of genes including all those whose
255    expression is enriched in the male accessory glands according to modENCODE (no less than
256    'Moderately High' in accessory glands but no more than 'Moderate' in any non-reproductive adult
257    tissues) and FlyAtlas2 (accessory glands enrichment higher than 1). This second search retrieved
258    most of the genes found in the first search plus six new candidates that have not been implicated
259    in reproduction (table 3).

260    Next, we explored whether the candidate TFs we identified in *D. melanogaster* are also expressed
261    in the male accessory glands of *D. virilis*, where accessory glands-biased transcripts were recently
262    identified by RNA-seq (Ahmed-Braimah et al. 2017). Seven of the 14 *D. melanogaster* candidates
263    showed clear homology to *D. virilis* genes with accessory glands-biased transcripts that were also
264    predicted to encode TFs (table 3). This contrasts with the low proportion (16.9%) of *D.*
265    *melanogaster* seminal genes having homologs among *D. virilis* seminal genes. In addition, *Ka/Ks*
266    ratios estimated for the candidate seminal TFs (0.10 on average, range: 0.03–0.28,) were lower
267    than those estimated for seminal genes (0.27 on average, range: 0.02–1.51) (fig. 3). These results
268    suggest that the high turnover rate and the rapidly adaptive evolution of SFPs do not have a
269    strong correlate in the evolution of seminal TFs.

270    The evolution of seminal genes' regulatory networks may follow the evolution of cis elements
271    rather than that of TFs. However, enhancers, insulators, and promoters that are active in the
272    male accessory glands have not been thoroughly investigated. Thus, the study of seminal TFs and
273    their binding sites is an important area for future research.

274   Besides TFs and their binding sites, post-transcriptional factors such as microRNAs (miRNAs) are
275   also involved in the regulation of seminal genes' expression. Recently, Mohorianu et al. (2018)
276   made an important contribution to the understanding of seminal regulatory networks by
277   assessing the role of miRNAs in the modulation of ejaculate composition. The authors found
278   evidence for the presence of several regulatory miRNAs that bind to a given sequence of the 3'
279   untranslated region (UTR) of seminal transcripts, likely repressing translation. Each miRNA
280   targets a specific group of seminal genes that share the corresponding 3' UTR target site, which
281   provides males with a mechanism to adjust ejaculate composition (Mohorianu et al. 2018). These
282   findings indicate that seminal genes UTRs and accessory glands miRNAs may have been involved
283   in the evolution of the seminal proteome.

284   Beyond the regulatory elements identified in *D. melanogaster*, causes underlying the divergence
285   of seminal genes at the regulatory level remain mostly unknown. Certainly, comparative
286   genomics will help to address this problem, however, we first need to identify the involved
287   elements in other species. Therefore, future research studying accessory glands transcriptome in
288   different *Drosophila* species will likely benefit this unexplored field.

289

290   **Turnover Rate**

291   One of the most striking characteristics of SFPs is their rich diversity, which seems to be causally
292   related, at least in part, to sexual conflict (Chapman 2008, 2018). In theory, postmating sexual
293   selection can escalate the evolutionary tension between the fitness interests of males and
294   females because male adaptations to sperm competition can be harmful to females (Chapman
295   et al. 1995; Lung et al. 2002; Mueller et al. 2007). Selection will then favor both female traits that
296   counteract detrimental male adaptations and male traits that respond to female resistance,
297   potentially leading to coevolutionary arms races between male persistence and female
298   resistance (Arnqvist 2004; Chapman et al. 2003). SFPs, by affecting female physiology and
299   behavior, clearly influence fertilization success and sperm competitiveness. Therefore, sexual
300   antagonistic coevolution between SFPs and the female counterparts likely accounts for the rapid
301   divergence of seminal proteomes (Sirot et al. 2014).

302   As sperm competition and sexual conflict can lead to rapid adaptive divergence of orthologous
303   SFPs, they may also promote divergence of the seminal protein repertoire through the gain of
304   novel seminal genes as well as through seminal gene loss. On one hand, females will not be
305   adapted to resist the action of novel SFPs. On the other hand, the expression of ancient SFPs—
306   whose action has been neutralized by females' counter-adaptations—will not be sustained by
307   natural selection. According to this hypothesis, turnover of seminal genes would be adaptive for
308   males because it would provide males with resources to "stay ahead" of female resistance

309   (Chapman 2018; Sirot et al. 2014). Evidence supporting sexual conflict as a driver of seminal
310   protein evolution abounds and comes from diverse sources (reviewed in Chapman 2018; Hollis
311   et al. 2019; Sirot et al. 2014).

312   High turnover rate of seminal gene sets was first noted by Wagstaff and Begun (2005). Assaying
313   the just released *D. pseudoobscura* genome for orthologs of *D. melanogaster* ACP-coding genes,
314   the authors noticed an unexpectedly high proportion of absences, suggesting that an important
315   number of seminal genes are lineage-specific. Later that year, Begun and Lindfors explored the
316   presence/absence patterns of three *D. simulans* ACP-coding genes across closely related species
317   of the *melanogaster* subgroup, to which *D. simulans* belongs. They found that two of these genes
318   (*Acp23D4* and *Acp54A1*) were absent in at least one species but had one to three copies in the
319   rest. Mueller et al. (2005), by performing comparative sequence analysis on 52 ACP-coding genes
320   of the *melanogaster* subgroup, found that 22 of them were not conserved in *D. pseudoobscura*.
321   Overall, these studies introduced the idea that the fraction of the genome encoding SFPs is, by
322   means of gene gain and loss, unusually dynamic.

323   With the release of the genomes of 12 *Drosophila* species (Drosophila 12 Genomes Consortium
324   2007), several comparative studies confirmed this pattern (e.g., Ahmed-Braimah et al. 2017;
325   Findlay et al. 2008, 2009; Haerty et al. 2007; Zhang et al. 2007). However, since too few dipteran
326   species were assayed for extensive identification of seminal genes, a comprehensive analysis to
327   trace the origin and loss of seminal genes in a phylogenetic context is lacking. Currently, we do
328   not know, for instance, to what extent orthologs of *D. melanogaster* seminal genes also encode
329   SFPs in other species of the genus. We do not know either how long ago these genes have
330   encoded SFPs in the *D. melanogaster* lineage. Identifying seminal genes/proteins in other
331   *Drosophila* species would allow to not only survey the evolutionary history of SFPs but also study
332   how new SFPs arise and how regulatory elements of seminal genes diverge between species. So
333   far, these questions have been barely explored.

334   Another question that arises is whether a core of SFPs playing essential reproductive roles has
335   been conserved throughout evolution. In such a case, these "essential SFPs", critical for
336   reproduction, should be present in a vast number of taxa. They could be searched by recognizing
337   the SFPs shared not only by closely related species but also by several phylogenetically distant
338   taxa; those shared only by closely related species would include both essential and non-essential
339   ones.

340   Intending to survey this hypothesis in Diptera, here we compiled a list of SFPs of the
341   *melanogaster* subgroup (those identified in *D. melanogaster*, *D. simulans* and/or *D. yakuba*) and
342   search it for homologs among SFPs identified in other dipteran taxa with well-known seminal
343   genes/proteins [see methodological procedures in Methods (Orthology of SFPs among Diptera)].

344   Taking into account that identification studies are hardly exhaustive, we only considered the
345   three outgroup taxa for which SFPs or seminal genes were identified in no less than two species
346   by independent extensive searches. These taxa were the *virilis-repleta* radiation of the *Drosophila*
347   subgenus (that split from *D. melanogaster* ~35 mya), tephritid fruit flies (that split from *D.*
348   *melanogaster* ~120 mya), and mosquitoes (that split from *D. melanogaster* ~250 mya). We
349   clustered all annotated proteins of 19 *Drosophila* species, including the *melanogaster* subgroup
350   and the *virilis-repleta* radiation, in 23782 groups of orthologs (orthogroups), 196 of which have
351   at least one seminal gene of the *melanogaster* subgroup. Among these 196 orthogroups 41
352   contain seminal genes of the *virilis-repleta* radiation, 11 have at least one homolog of tephritid
353   seminal genes, and 25 have at least one homolog of mosquitoes' seminal genes (fig. 4). Caution
354   should be taken when comparing these numbers because they relied on different homology
355   criteria, some applied by different previous studies (supplementary table S2). However,
356   considering that 11,298 out of the 13,969 (81%) protein-coding genes have certainly clear
357   homologs in mosquitoes (blastp bit score > 50), the number of SFPs shared by the four evaluated
358   taxa seemed to be remarkably low: only two orthogroups had seminal genes of the four taxa.
359   One of these orthogroups contains only one *D. melanogaster* seminal gene (*FBgn0034753*),
360   which encodes a peptidyl-prolyl cis-trans isomerase. The other contains five *D. melanogaster*
361   seminal paralogs that encode protease inhibitors with Kazal domains and belong to a tandem
362   gene cluster located in the left arm of the second chromosome. Within this gene family, we found
363   *FBgn0266364*, which was identified as a novel candidate in the present report, and *FBgn0051704*,
364   which is reported in FlyBase r2020_03 as ortholog of *SPINK2*, a human gene implicated in male
365   infertility.

366   Although the number of taxa included in our analysis is low, the results indicate that most SFPs
367   in Diptera are lineage-specific, which strongly suggests that most SFPs have a short evolutionary
368   life (or diverges rapidly beyond detectable homology) and that not many—if any—have been
369   critical for reproduction throughout Diptera evolution. Still, even if the seminal protein
370   repertoires of the taxa we analyzed were fairly complete, our results would be far from being
371   conclusive because homology detection across dipteran families can be inefficient for rapidly
372   evolving seminal genes. In this sense, it would be more feasible to search for "essential SFPs"
373   within specific groups of the *Drosophila* genus. However, the repertoire of SFPs is currently
374   known for too few species. Thus, the search for "essential SFPs" within *Drosophila* must await
375   more studies assaying SFPs in a wider spectrum of species.

376   Despite those observations and claims, gene birth and death rates were never estimated for SFPs.
377   To obtain these estimates, we pruned the 196 orthogroups containing *D. melanogaster* seminal
378   genes, leaving only the nine species for which genomic annotations were updated at least once
379   [see Methods (Gene Birth and Death Rates)]. Then, duplications, losses, and orthogroup gains
380   were identified in the gene trees of each orthogroup (fig. 5) and each event rate was estimated

381    from the obtained figures. Taking into account divergence dates reported in Obbard et al. (2012),
382    the estimated duplication rate was 0.0097 duplications per gene per million years (/gene/my)
383    and the loss rate was 0.0122 losses/gene/my (0.0133 duplications/gene/my and 0.0212
384    losses/gene/my considering only the species of the *melanogaster* group). The species with the
385    greatest gene loss rate was *D. sechellia* (49 losses), which could be an artifact of genome
386    sequencing, assembly, and annotation. However, the number of protein-coding genes annotated
387    for this species is the highest in the *melanogaster* group and a similar pattern of high gene loss
388    was previously observed for olfactory genes in this species (Almeida et al. 2014; McBride 2007).
389    The authors associated this with *D. sechellia* specialization and endemism, which could also have
390    implications for the mating system and reproductive proteins. Regarding orthogroup gains in the
391    *D. melanogaster* lineage, the estimated rate was 0.0047 gains/gene/my and the total number of
392    identified events was 87. The acquisitions were inferred in the ancestors of the *Sophophora*
393    subgenus (25), the *melanogaster* group (22), the *melanogaster* subgroup (35), and the
394    *melanogaster* complex (*D. melanogaster*, *D. simulans*, and *D. sechellia*) (5). Interestingly, the
395    latter figure accounts for more than half the number of putative *de novo* genes identified by Zhou
396    et al. (2008) in the *melanogaster* complex.

397    Using 12 *Drosophila* genomes, Hahn et al. (2007) estimated a total event (gene duplications +
398    losses) rate of 0.0013 events/gene/my based on Tamura et al. (2004) divergence dates. Using the
399    same dates, we estimated for the *D. melanogaster* seminal genes an event rate of 0.0096
400    events/gene/my (0.0111 events/gene/my considering only the species of the *melanogaster*
401    group). This suggests that seminal genes' families, though they may not contain seminal genes of
402    non-*D. melanogaster* species, are approximately seven times more dynamic than the average
403    gene family in *Drosophila*.

404

**Mechanisms of Origin**

406    The high turnover rate in seminal genes/proteins repertoires implies a high proportion of novel
407    seminal genes/proteins restricted to young lineages or unique species. This facilitates studying
408    the evolution of novel genes in a common cellular background (i.e., accessory glands) in groups
409    of closely related species, where the molecular routes of gene origin are more likely traceable.
410    Thus, seminal genes provide an excellent opportunity to investigate how novel proteins and
411    biological functions emerge. Four mechanisms have been reported or proposed so far as
412    responsible for the origin of seminal genes in *Drosophila*: duplication of seminal genes,
413    duplication of non-seminal genes, gene co-option into the  male reproductive tract, and *de novo*
414    evolution (reviewed in Sirot 2019).

415   The first mechanism proposed was duplication of preexisting seminal genes (e.g., Almeida &
416   Desalle 2009; Findlay et al. 2008; Holloway & Begun 2004; Mueller et al. 2005; Wagstaff & Begun
417   2005). When a seminal gene is entirely duplicated so that both copies, the new and the old,
418   encode the same SFP, ensuing mutations may lead to subfunctionalization or
419   neofunctionalization, giving rise to novel SFPs with similar amino acid sequences. Most of the
420   seminal genes encoding these proteins are located in clusters of nearby genes on the second
421   chromosome (fig. 2), showing that tandem duplication followed by mutation has played an
422   important role in the divergence of the seminal proteome. For instance, *FBgn0043825*,
423   *FBgn0051872*, and *FBgn0265264* are three paralogs located in tandem on the left arm of the
424   second chromosome, which encode SFPs with triglyceride lipase activity (Mueller et al. 2005).

425   Duplication of genes that are not expressed in the male reproductive system and do not encode
426   SFPs may also be a source of novel seminal genes (Sirot 2019); if a duplicate ends up placed under
427   the control of regulatory elements driving its expression in the accessory glands, it may become
428   a new seminal gene. Genes encoding proteins that already have secretion signals are likely
429   sources for this mechanism. An example of this is the origin of the seminal gene *FBgn0052833*,
430   which resulted from a duplication-mediated co-option of a female-expressed gene whose original
431   copy encodes a secretory protein of the sperm storage organs (Sirot et al. 2014). Another
432   example comes from odorant binding proteins (OBPs), a highly dynamic family of olfactory genes
433   that are usually expressed in the antennae. Four OBP genes, however, have been co-opted into
434   the accessory glands exclusively in the lineage leading to the *melanogaster* group (Almeida et al.
435   2014). Interestingly, the rates of protein evolution of these genes were the highest among OBPs.

436   Although duplication may facilitate sequence or expression evolution because of initial
437   redundancy (one copy can change, while the other maintains the original function), some
438   *Drosophila* seminal genes seem to have arisen via gene co-option in the absence of a previous
439   gene duplication event (Findlay et al. 2008). *FBgn0262571*, a *D. melanogaster* seminal gene
440   exclusively expressed in the male accessory glands, belongs to a single-copy gene family (Sepil et
441   al. 2019). Its orthologs, despite encoding proteins with secretion signal, are not within the
442   repertoire of seminal genes in either *D. mojavensis*, *D. pseudoobscura*, or *D. virilis* (the only non-
443   *melanogaster* group species of the genus in which seminal genes were extensively identified).
444   Therefore, despite not being duplicated, this gene was potentially co-opted into the accessory
445   glands in the *D. melanogaster* lineage, during the evolution of the *melanogaster* group.

446   Some other seminal genes may have emerged *de novo* from ancestrally noncoding DNA (Begun
447   et al. 2006; Findlay et al. 2008; Haerty et al. 2007). While sperm competition and sexual conflict
448   may steadily select for innovation in the male ejaculate, "fitness valleys" limit the paths available
449   for the evolution of preexisting proteins (Camps et al. 2007). In this sense, young *de novo* seminal
450   genes may be less constrained and may have more opportunities to fill the emerging functional

451 niches. Curiously, the first evidence consistent with *de novo* gene birth comes from studies aimed
452 to identify genes specifically expressed in *Drosophila* male accessory glands (Begun et al. 2006)
453 or testes (Begun et al. 2007; Zhao et al. 2014). Given the high proportion of insect seminal genes
454 without identified orthologs, *de novo* gene birth is believed to account for the origin of many
455 seminal genes (reviewed in Sirot 2019). So far, however, no *Drosophila* seminal genes have yet
456 been identified as *de novo* genes with high confidence, possibly because distinguishing *de novo*
457 birth from horizontal transfer or rapid protein divergence (which is common among seminal
458 proteins) is challenging (Zile et al. 2020).

459 Despite particular cases, a broad-scale analysis to determine the relative contribution of the
460 alternative mechanisms of origin has yet to be completed. In an attempt to discern which of the
461 mentioned mechanisms were responsible for the origin of young *D. melanogaster* SFPs [those
462 that have arisen during the evolution of the *melanogaster* species group, i.e., less than ~25
463 million years ago (mya)], we identified gene families that included *melanogaster* group's seminal
464 genes. Given that homology detection power banishes with divergence, evaluating alternative
465 mechanisms of origin for older genes would be much more uncertain. Gene families were
466 obtained by clustering the proteins of reference proteomes of 19 *Drosophila* species [see
467 Methods (Seminal Gene Families)]. This analysis revealed that our set of 219 *D. melanogaster*
468 SFPs belong to 168 gene families. To determine which seminal genes have likely emerged after
469 the origin of the *melanogaster* group (which were dubbed young seminal genes), and to infer the
470 most likely mechanism of origin, we manually inspected the gene family tree of all these 168
471 gene families. Specifically, we explored the presence/absence of orthologs and paralogs, and
472 whether they had been classified as SFPs. We then applied the parsimony principle to determine,
473 according to the observed pattern, which mechanism was most likely responsible for the origin
474 of each young *D. melanogaster* SFP (fig. 6 illustrates our criteria). See Methods (Seminal Gene
475 Families) for a more detailed description of the applied criteria. In cases where *n* mechanisms
476 were equally likely, we assigned "*1/n* genes" to each mechanism.

477 In this way, we estimated that 76 *D. melanogaster* seminal genes existed as seminal genes
478 (before the split from the lineage leading to *D. pseudoobscura* (~25 mya). For 13 seminal genes,
479 we could not determine whether the origin was before or after that split since they exhibited
480 uncertain homology to sequences of outgroup or distant species. Among the remaining 130 *D.*
481 *melanogaster* seminal genes (i.e., the tentatively young ones), we classified ~27 (20.6%) as
482 duplicates of preexisting seminal genes, ~7 (5.3%) as co-opted duplicates (duplicates of non-
483 seminal genes), ~47 (36.5%) as co-opted without duplication, and ~49 (37.6%) as putative
484 orphans (fig. 7).

485 These results may give the impression that *de novo* emergence was responsible for the origin of
486 many *D. melanogaster* seminal genes. However, our approach did not contemplate all possible

487    mechanisms of gene origin and may have confounded some. For instance, a non-orphan seminal
488    gene showing fast evolution may have diverged beyond detectable homology and be construed
489    as an orphan gene. Some of the proteomes we used may be incomplete due to potentially
490    defective genomic annotations, which may also have led to the overestimation of taxonomically
491    restricted genes. In consequence, the actual number of orphans among seminal genes of the
492    *melanogaster* group is surely lower than the one we estimated. In fact, we could not ensure *de*
493    *novo* status for any of the identified putative orphans [see applied criteria in Methods (*De Novo*
494    Status Validation)]. Briefly, after examining several *Drosophila* annotated genomes, we failed to
495    find taxonomically restricted seminal genes with syntenic homologous, reliably noncoding
496    sequences in any outgroup species. This means that these gene families, which were initially
497    identified as taxonomically restricted to the *melanogaster* group, may be classified as originating
498    through rapid evolution (among other mechanisms) rather than *de novo* emergence. Therefore,
499    the relative contribution of *de novo* emergence to the origin of *Drosophila* seminal genes may be
500    more limited than previously thought. Gene co-option, on the other hand, appears to be the most
501    frequent mechanism of origin.

502    To uncover the possible ancestral expression pattern of those few seminal genes that, according
503    to our analysis, appear to have arisen via duplication-mediated co-option, we checked the
504    expression pattern of the respective non-seminal paralogs. According to modENCODE
505    (implemented in FlyBase r2020_03), these paralogs are expressed in the larval salivary gland, the
506    adult female spermatheca, the pupal fat body, or the adult digestive system. Whether these
507    tissues represent common sources for co-option into the seminal fluid will require further cross-
508    species exploration of co-opted seminal genes (for examples in other insects see Martinson et al.
509    2017; Meslin et al. 2015).

510    Alternative mechanisms of seminal genes' origin—such as exon/domain shuffling, gene
511    fission/fusion, horizontal gene transfer, and reading-frame shift—should be explored in the
512    future. Also, further identification of SFPs in more *Drosophila* species will allow for more accurate
513    discrimination between alternative mechanisms, for dating gene origin more precisely, and for
514    exploring gene origin in other groups.

515

**Conclusions**

517    Here, we provided an overview of the inter-specific divergence of *Drosophila* SFPs summarizing
518    the current state of knowledge and emphasizing the intriguing aspects that are less understood.
519    We focused on the conservation of SFPs across the order Diptera and the mechanisms of origin
520    of *Drosophila* seminal genes. We not only inspected some of the main contributions to these

521 topics but also compiled genomic information from multiple species and performed molecular
522 evolutionary analyses to address some broad questions that remain open.

523 Using reviewed criteria, we presented a novel set of high-confidence seminal protein candidates
524 for *D. melanogaster* and generated a database of *Drosophila* SFPs. We also provided, for the first
525 time, a list of accessory glands (putative or confirmed) TFs presumptively controlling the
526 expression of SFPs.

527 Two interesting patterns derive from our comparative genomic analyses. First, given the low
528 number of common SFPs found among the three inspected dipteran families, the hypothesis that
529 there is a core of indispensable, "essential SFPs" conserved across Diptera seems unlikely.
530 Second, gene co-option appears to be the most frequent mechanism accounting for the origin of
531 *Drosophila* seminal genes. As *de novo* evolution could not be ensured for any seminal gene, our
532 analysis failed to support the hypothesis that *de novo* emergence is a frequent mechanism of
533 origin for seminal genes.

534 Despite the insights we have gained, it is evident that characterizing the seminal proteome in
535 more species, especially in those outside the *melanogaster* group, is imperative to fill important
536 knowledge gaps. While proteomics on isotopic labeled flies and quantitative proteomics have
537 proven to be useful to carry out this task, our searches suggest that RNA-seq on accessory glands,
538 which is less challenging and cheaper, would provide valuable starting information.

539

540 **Methods**

541 Orthology of SFPs among Diptera

542 Supplementary table S2 summarizes the sources of the list of SFPs for each considered taxa (lists
543 are available upon request). To identify the orthologs of the SFPs identified in the *melanogaster*
544 group (ingroup), we employed the following strategy. First, we gathered the proteomes of 19
545 *Drosophila* species (see below) and used Orthofinder, a platform for comparative genomics
546 (Emms & Kelly 2015, 2019), to cluster the proteins in groups of orthologs (orthogroups). Then,
547 we searched for the orthogroups that had any SFP of the *melanogaster* subgroup [i.e., the 219 of
548 *D. melanogaster* or those of *D. simulans* and/or *D. yakuba* identified by Findlay et al. (2008)]. The
549 input protein sequences were obtained from reference proteomes available in FlyBase, NCBI, or
550 specific genome projects' sites. The *Drosophila* species of the *melanogaster* group included in the
551 analysis were *D. ananassae* [dana_r1.06 (FlyBase r2020_03)], *D. biarmipes* [Dbia_2.0 (Richards et
552 al. unpublished, NCBI)], *D. bipectinata* [Dbip_2.0 (Richards et al. unpublished, NCBI)], *D. elegans*
553 [Dele_2.0 (Richards et al. unpublished, NCBI)], *D. erecta* [dere_r1.05 (FlyBase r2020_03)], *D.*

554  *eugracilis* [Deug_2.0 (Richards et al. unpublished, NCBI)], *D. ficusphila* [Dfic_2.0 (Richards et al.
555  unpublished, NCBI)], *D. kikkawai* [Dkik_2.0 (Richards et al. unpublished, NCBI)], *D. mauritiana*
556  [dmauMS17_r1.0 (Nolte et al. 2013)], *D. melanogaster* [dmel_r6.34 (FlyBase r2020_03)], *D.*
557  *rhopaloa* [Drho_2.0 (Richards et al. unpublished, NCBI)], *D. sechellia* [dsec_r1.3 (FlyBase
558  r2020_03)], *D. simulans* [dsim_r2.02 (FlyBase r2020_03)], *D. suzukii* (Joanna C. Chiu 2020,
559  personal communication), *D. takahashii* [Dtak_2.0 (Richards et al. unpublished, NCBI)], and *D.*
560  *yakuba* [dyak_r1.05 (FlyBase 2017_03) re-annotated by Yang et al. (2018)]. Species belonging to
561  other species groups (outgroups) were *D. mojavensis* [dmoj_r1.04 (FlyBase r2017_03) re-
562  annotated by Yang et al. (2018)], *D. pseudoobscura* [UCI_Dpse_MV25 (Liao et al. unpublished,
563  NCBI)], and *D. virilis* [dvir_r1.06 (FlyBase 2017_03) re-annotated by Yang et al. (2018)]. These
564  three species were chosen because they were the only ones outside the *melanogaster* group in
565  which seminal genes were extensively studied. As Yang et al. (2018) did not annotate CDSs, we
566  predicted for *D. mojavensis*, *D. virilis*, and *D. yakuba* one protein per gene with RefProt pipeline
567  (Revale & Hurtado, available upon request), which is based on TransDecoder (Haas et al. 2013),
568  Blast (Altschul et al. 1990), HMMER (hmmer.org), and several inhouse R scripts (R-project.org).
569  In our experience, Orthofinder has limited recall when clustering sequences of very distantly
570  related species. Therefore, to recognize orthogroups with SFPs of species outside *Drosophila*
571  (*Aedes aegypti*, *Aedes albopictus*, *Anopheles gambiae*, *Bactrocera dorsalis*, and *Ceratitis capitata*)
572  we relied on previous orthology assignments based on Blast (supplementary table S2). We
573  considered a SFP to be shared between *melanogaster* subgroup and any given outgroup if the
574  protein was clustered together with an outgroup SFP in the same orthogroup.

575

576  Molecular Evolutionary Analyses

577  Estimates of the ratio between the rate of non-synonymous substitution (*Ka*) and the rate of
578  synonymous substitutions (*Ks*) can be used as a proxy to investigate the evolutionary forces that
579  shape the evolution of proteins. Close to zero ratios are associated with purifying selection,
580  whereas ratios close or higher than one mean that the gene evolves under neutrality or that
581  some codons are positively selected. We employed PAML-4.8 (Yang 2007) to obtain $\omega$, a
582  likelihood-based estimator of *Ka*/*Ks*, for each orthogroup.

583

584  Gene Birth and Death Rates

585  We pruned the 196 orthogroups containing *D. melanogaster* SFP-coding genes (see above) to
586  include only those species with updated genome annotations, leaving in this way the orthologs
587  of *D. melanogaster*, *D. simulans*, *D. sechellia*, *D. erecta*, *D. yakuba*, *D. ananassae*, *D.*

588 *pseudoobscura*, *D. mojavensis*, and *D. virilis*. Then we employed the program Notung-2.9.1.5
589 (Chen et al. 2000; Darby et al. 2017) to identify gene duplications, losses, and *de novo* gains in
590 each orthogroup by comparing gene trees with the species tree. To be conservative and avoid
591 overestimation, we edited the Notung results to remove duplications and losses when there was
592 an even number of genes per species. With the total number of each of these events for each
593 branch of the *Drosophila* phylogeny, we estimated per gene rates by dividing the number of
594 events by the number of genes in the ancestral branches. These events were summed across all
595 branches and the sum was divided by the total phylogeny time to obtain the rates using the
596 formulas described in Vieira et al. (2007). A gene gain was identified for each orthogroup
597 exclusive of a monophyletic clade.

598

599 Seminal Gene Families

600 Since Orthofinder inference relies on reciprocal best alignment hits, some paralogous sequences
601 ended up grouped in separate orthogroups. Thus, with the aim of identifying paralogous
602 orthogroups, we compared *D. melanogaster* sequences clustered in different orthogroups using
603 Blastp. We then merged orthogroups with aligned sequences into more inclusive gene families.
604 Since we used a conservative bit score cutoff of 80 for filtering hits, the number of recognized
605 gene families probably represent an upper bound of the actual number. Our objective was to
606 determine the origin of *D. melanogaster* seminal genes that had emerged during the evolution
607 of the *melanogaster* group (i.e., after the split from the lineage leading to *D. pseudoobscura*), so
608 we considered the species belonging to other groups as outgroups. We then used the gene trees
609 generated by Orthofinder to investigate the origins of the *melanogaster* group SFPs. Within each
610 orthogroup, the last common ancestor gene between an outgroup seminal gene and a *D.*
611 *melanogaster* seminal gene was considered as a seminal gene. Similarly, the last common
612 ancestor gene at the root of any orthogroup containing homologs to seminal genes of tephritids
613 or mosquitoes was also considered as a seminal gene. With these considerations, we inferred the
614 most likely mechanism of origin of each *D. melanogaster* seminal gene by manually inspecting
615 the respective gene family tree. Specifically, we explored the presence/absence of orthologs and
616 paralogs among species of the *melanogaster* group and outgroups applying the parsimony
617 principle over gene gain/loss events (fig. 6). In this way, we first distinguished between "ancient"
618 (those that had emerged before the split from the lineage leading to *D. pseudoobscura*, ~25 mya)
619 and tentatively young (those lacking homologs among outgroup seminal genes, that have likely
620 emerged after the split from the lineage leading to *D. pseudoobscura*) *D. melanogaster* seminal
621 genes. Then, we classified tentatively young seminal genes into the following four categories:
622 duplicated, co-opted after being duplicated, co-opted without being duplicated, and orphan. In
623 those cases where *n* mechanisms were equally likely, we assigned "*1/n* genes" to each

624    mechanism. Some *D. melanogaster* proteins may have evolved very rapidly, hindering homology

625    detection. Thus, in the case of SFPs classified as orphan with our approach, we evaluated distant

626    homology by comparing *D. melanogaster* SFPs to non-redundant proteins sequences from NCBI

627    databases using Blastp (blast.ncbi.nlm.nih.gov). In this case, we admitted hits (bit score > 39)

628    against sequences of any Diptera: those with any bit score higher than 50 were considered to

629    reflect homology while those with bit scores between 39 and 50 were considered uncertain. Also,

630    for each apparent orphan seminal gene, we checked manually the absence of syntenic open

631    reading frames encoding similar proteins (Blastp: bit score > 39 or positives > 60%) in the *D.*

632    *pseudobscura* genome by using the Ensembl Metazoa genome browser (Howe et al. 2019).

633

634    *De Novo* Status Validation

635    To validate the *de novo* status of the putative orphans, we used the conservative criteria applied

636    by Zile et al. (2020). Briefly, as *de novo* genes should have syntenic, homologous noncoding

637    sequences in closely related outgroup species, we inspected each orphan candidate for syntenic,

638    homologous noncoding sequences in well-annotated genomes of outgroup species. Particularly,

639    we examined the latest public assemblies for *D. anananassae* [DanaRS2.1 (Zhang et

640    al.unpublished, NCBI)], *D. elegans* [Dele_2.0 (Richards et al. unpublished, NCBI)], *D. erecta*

641    [DereRS2 (Zhang et al.unpublished, NCBI)], *D. pseudoobscura* [UCI_Dpse_MV25 (Liao et al.

642    unpublished, NCBI)], *D. simulans* [Prin_Dsim_3.0 (Pinharanda et al. unpublished, NCBI)], *D.*

643    *suzukii* [LBDM_Dsuz_2.1.pri (Paris et al. unpublished, NCBI)], and *D. yakuba*

644    [Prin_Dyak_Tai18E2_2.0 (Reilly et al. unpublished, NCBI)]. For instance, for a gene family

645    restricted to the *melanogaster* complex (*D. melanogaster*, *D. sechellia* and *D. simulans*), any

646    species outside this complex (i.e., *D. ananassae*, *D. elegans*, *D. erecta*, *D. pseudoobscura*, *D.*

647    *suzukii* and *D. yakuba*) was considered an outgroup. Thus, for each gene family having orphan

648    candidates, Blastn searches were applied to search the syntenic genomic regions of the outgroup

649    genomes for homologous sequences (bit score > 39 or identities > 60%). The found homologous

650    syntenic sequences showing evidence of being transcribed (i.e., evidence from RNA-Seq

651    alignment data) were searched—employing Blastp searches—for the absence of homologous

652    open reading frames (bit score < 39 and positives < 60%).

653

657  *suzukii* 2.0 genome. This work was supported by the Agencia Nacional de Promoción Científica y

658  Técnica through grants awarded to JH and EH.

659

660  **Authors' Contributions**

661  JH conceived and designed the study, compiled and analyzed the data, and took the lead in

662  writing the manuscript. FCA was involved in planning the work and analysis design; she also

663  estimated rates of molecular evolution and gene gain/loss. SAB performed functional

664  annotations and designed the figures. SR helped integrate genomic information and predict

665  protein sequences. EH was involved in planning the work and supervised the project. All authors

666  discussed the results and contributed to the final manuscript.

667

668  **Supplementary Material**

669  Table S1. List of *D. melanogaster* seminal genes. As KSGs we included genes encoding proteins

670  previously confirmed to be transferred by males into females during mating, those meeting

671  stringent multiple criteria that indicate so according to Sepil et al. (2019), or those expressed in

672  male reproductive tissues more than in any other tissue (according to modENCODE and FlyAtlas2)

673  also encoding secretable proteins found in the mating plug [according to Avila et al. (2015) and

674  Wigby et al. (2020)]. As candidates, we included our novel candidates as well as previously

675  predicted seminal genes. We excluded genes expressed specifically in the testes (according to

676  FlyAtlas2) that encode sperm proteins (Wigby et al. 2020), those candidates proposed only by

677  Wigby et al. (2020) that show low expression in male reproductive tissues and higher expression

678  in other male and female tissues (according to modENCODE and FlyAtlas2), and those proposed

679  only by Ayroles et al. (2011) that do not encode secretable proteins (signalP). The evaluated

680  conditions for the expression/secretion criterion and sources that previously identified the gene

681  as seminal are shown for each gene (see supplementary references).

682  Table S2. SFPs of the *melanogaster* subgroup, the *virilis-repleta* radiation, tephritids, and

683  mosquitoes. Sources and methods used to compile the list are summarized for each considered

684  species (see supplementary references).

685

686  **Data Availability**

687 Despite no new data were generated in support of this research, the compiled information and
688 data underlying our analyses are available in the article, in its online supplementary material,
689 and/or at the open-access databases duly mentioned in the text.

690

691 **References**

692 Ahmed-Braimah YH, Unckless RL, Clark AG. 2017. Evolutionary dynamics of male reproductive
693 genes in the *Drosophila virilis* subgroup. G3 Genes, Genomes, Genet. 7:3145–3155. doi:
694 10.1534/g3.117.1136.

695 Aigaki T, Fleischmann I, Chen PS, Kubli E. 1991. Ectopic expression of sex peptide alters
696 reproductive behavior of female *D. melanogaster*. Neuron. 7:557–563. doi: 10.1016/0896-
697 6273(91)90368-A.

698 Almagro Armenteros JJ et al. 2019. SignalP 5.0 improves signal peptide predictions using deep
699 neural networks. Nat. Biotechnol. 37:420–423. doi: 10.1038/s41587-019-0036-z.

700 Almagro Armenteros JJ, Sønderby CK, Sønderby SK, Nielsen H, Winther O. 2017. DeepLoc:
701 prediction of protein subcellular localization using deep learning. Bioinformatics. 33:3387–3395.
702 doi: 10.1093/bioinformatics/btx431.

703 Almeida FC, Desalle R. 2008. Evidence of adaptive evolution of accessory gland proteins in closely
704 related species of the *Drosophila repleta* group. Mol. Biol. Evol. 25:2043–2053. doi:
705 10.1093/molbev/msn155.

706 Almeida FC, Desalle R. 2009. Orthology, function and evolution of accessory gland proteins in the
707 *Drosophila repleta* group. Genetics. 181:235–245. doi: 10.1534/genetics.108.096263.

708 Almeida FC, Sánchez-Gracia A, Campos JL, Rozas J. 2014. Family size evolution in *Drosophila*
709 chemosensory gene families: A comparative analysis with a critical appraisal of methods.
710 Genome Biol. Evol. 6:1669–1682. doi: 10.1093/gbe/evu130.

711 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. J.
712 Mol. Biol. 215:403–410. doi: 10.1016/S0022-2836(05)80360-2.

713 Arnqvist G. 2004. Sexual conflict and sexual selection: lost in the chase. Evolution. 58:1383–1393.
714 doi: 10.1111/j.0014-3820.2004.tb01716.x.

715 Avila FW et al. 2015. Retention of ejaculate by *Drosophila melanogaster* females requires the
716 male-derived mating plug protein PEBme. Genetics. 200:1171–1179. doi:

717    10.1534/genetics.115.176669.

718    Avila FW et al. 2016. Nature and functions of glands and ducts in the *Drosophila* reproductive
719    tract. In: Extracellular Composite Matrices in Arthropods. Cohen, E & Moussian, B, editors.
720    Switzerland: Springer International Publishing pp. 411–444. doi: 10.1007/978-3-319-40740-1.

721    Avila FW, Sirot LK, LaFlamme BA, Rubinstein CD, Wolfner MF. 2011. Insect seminal fluid proteins:
722    identification and function. Annu. Rev. Entomol. 56:21–40. doi: 10.1146/annurev-ento-120709-
723    144823.

724    Avila FW, Wolfner MF. 2017. Cleavage of the *Drosophila* seminal protein Acp36DE in mated
725    females enhances its sperm storage activity. J Insect Physiol. 101:66–72.
726    doi:10.1016/j.jinsphys.2017.06.015.

727    Ayroles JF, Laflamme BA, Stone EA, Wolfner MF, Mackay TFC. 2011. Functional genome
728    annotation of *Drosophila* seminal fluid proteins using transcriptional genetic networks. Genet.
729    Res. 93:387–395. doi: 10.1017/S0016672311000346.

730    Baumann H, Wilson KJ, Chen PS, Humbel RE. 1975. The amino-acid sequence of a peptide (PS-1)
731    from *Drosophila funebris*: a paragonial peptide from males which reduces the receptivity of the
732    female. Eur. J. Biochem. 52:521–529. doi: 10.1111/j.1432-1033.1975.tb04023.x.

733    Begun DJ, Lindfors HA. 2005. Rapid evolution of genomic Acp complement in the *melanogaster*
734    subgroup of *Drosophila*. Mol. Biol. Evol. 22:2010–2021. doi: 10.1093/molbev/msi201.

735    Begun DJ, Lindfors HA, Kern AD, Jones CD. 2007. Evidence for de novo evolution of testis-
736    expressed genes in the *Drosophila yakuba*/*Drosophila erecta* clade. Genetics. 176:1131–1137.
737    doi: 10.1534/genetics.106.069245.

738    Begun DJ, Lindfors HA, Thompson ME, Holloway AK. 2006. Recently evolved genes identified from
739    *Drosophila yakuba* and *D. erecta* accessory gland expressed sequence tags. Genetics. 172:1675–
740    1681. doi: 10.1534/genetics.105.050336.

741    Bertram MJ, Neubaum DM, Wolfner MF. 1996. Localization of the *Drosophila* male accessory
742    gland protein Acp36DE in the mated female suggests a role in sperm storage. Insect Biochem.
743    Mol. Biol. 26:971–980. doi: 10.1016/s0965-1748(96)00064-1.

744    Camps M, Herman A, Loh E, Loeb LA. 2007. Genetic constraints on protein evolution. Crit. Rev.
745    Biochem. Mol. Biol. 42:313–326. doi: 10.1080/10409230701597642.

746    Cavener DR, MacIntyre RJ. 1983. Biphasic expression and function of glucose dehydrogenase in

747    *Drosophila melanogaster*. Proc. Natl. Acad. Sci. U. S. A. 80:6286–6288. doi:
748    10.1073/pnas.80.20.6286.

749    Chapman T. 2018. Sexual conflict: Mechanisms and emerging themes in resistance biology. Am.
750    Nat. 192:217–229. doi: 10.1086/698169.

751    Chapman T et al. 2003. The sex peptide of *Drosophila melanogaster*: Female post-mating
752    responses analyzed by using RNA interference. Proc. Natl. Acad. Sci. U. S. A. 100:9923–9928. doi:
753    10.1073/pnas.1631635100.

754    Chapman T. 2008. The soup in my fly: Evolution, form and function of seminal fluid proteins. PLoS
755    Biol. 6:1379–1382. doi: 10.1371/journal.pbio.0060179.

756    Chapman T, Davies SJ. 2004. Functions and analysis of the seminal fluid proteins of male
757    *Drosophila melanogaster* fruit flies. Peptides. 25:1477–1490. doi:
758    10.1016/j.peptides.2003.10.023.

759    Chapman T, Liddle LF, Kalb JM, Wolfner MF, Partridge L. 1995. Cost of mating in *Drosophila*
760    *melanogaster* females is mediated by male accessory gland products. Nature. 373:241–244. doi:
761    10.1038/373241a0.

762    Chen K, Durand D, Farach-Colton M. 2000. NOTUNG: a program for dating gene duplications and
763    optimizing gene family trees. J. Comput. Biol. 7:429–447. doi: 10.1089/106652700750050871.

764    Chen PS et al. 1988. A male accessory gland peptide that regulates reproductive behavior of
765    female *D. melanogaster*. Cell. 54:291–298. doi: 10.1016/0092-8674(88)90192-4.

766    Chen S. 1984. Biochemistry of Insect Male Accessory Glands. Annu. Rev. Entomol. 29:233–255.

767    Cho KS et al. 1999. A 45-kDa cAMP-dependent phosphoprotein which is related to the product of
768    Mst57Dc in Drosophila melanogaster. Insect Biochem. Mol. Biol. 29:701–710. doi:
769    10.1016/S0965-1748(99)00046-6.

770    Civetta A, Ranz JM. 2019. Genetic factors influencing sperm competition. Front. Genet. 10:820.
771    doi: 10.3389/fgene.2019.00820.

772    Coleman S, Drähn B, Petersen G, Stolorov J, Kraus K. 1995. A *Drosophila* male accessory gland
773    protein that is a member of the serpin superfamily of proteinase inhibitors is transferred to
774    females during mating. Insect Biochem. Mol. Biol. 25:203–207. doi: 10.1016/0965-
775    1748(94)00055-M.

776    Corrigan L et al. 2014. BMP-regulated exosomes from *Drosophila* male reproductive glands

777      reprogram female behavior. J. Cell Biol. 206:671–688. doi: 10.1083/jcb.201401072.

778      Darby CA, Stolzer M, Ropp PJ, Barker D, Durand D. 2017. Xenolog classification. Bioinformatics.
779      33:640–649. doi: 10.1093/bioinformatics/btw686.

780      Davies SJ, Chapman T. 2006. Identification of genes expressed in the accessory glands of male
781      Mediterranean Fruit Flies (*Ceratitis capitata*). Insect Biochem. Mol. Biol. 36:846–856. doi:
782      10.1016/j.ibmb.2006.08.009.

783      Drosophila 12 Genomes Consortium. 2007. Evolution of genes and genomes on the *Drosophila*
784      phylogeny. Nature. 450:203–218. doi: 10.1038/nature06341.

785      Dyanov HM, Dzitoeva SG. 1995. Method for attachment of microscopic preparations on glass for
786      in situ hybridization, PRINS and in situ PCR studies. Biotechniques. 18:822–826.

787      Emms DM, Kelly S. 2019. OrthoFinder: Phylogenetic orthology inference for comparative
788      genomics. Genome Biol. 20:1–14. doi: 10.1186/s13059-019-1832-y.

789      Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole genome comparisons
790      dramatically improves orthogroup inference accuracy. Genome Biol. 16:1–14. doi:
791      10.1186/s13059-015-0721-2.

792      Findlay GD, MacCoss MJ, Swanson WJ. 2009. Proteomic discovery of previously unannotated,
793      rapidly evolving seminal fluid genes in *Drosophila*. Genome Res. 19:886–895. doi:
794      10.1101/gr.089391.108.

795      Findlay GD, Yi X, Maccoss MJ, Swanson WJ. 2008. Proteomics reveals novel *Drosophila* seminal
796      fluid proteins transferred at mating. PLoS Biol. 6:e178. doi: 10.1371/journal.pbio.0060178.

797      Gillott C. 1996. Male insect accessory glands: Functions and control of secretory activity.
798      Invertebr. Reprod. Dev. 30:199–205. doi: 10.1080/07924259.1996.9672546.

799      Gligorov D, Sitnik JL, Maeda RK, Wolfner MF, Karch F. 2013. A novel function for the hox gene
800      Abd-b in the male accessory gland regulates the long-term female post-mating response in
801      *Drosophila*. PLoS Genet. 9(3):e1003395. doi: 10.1371/journal.pgen.1003395.

802      Graveley BR et al. 2010. The *D. melanogaster* transcriptome: modENCODE RNA-Seq data.
803      http://www.modencode.org/celniker/.

804      Haas BJ et al. 2013. De novo transcript sequence reconstruction from RNA-seq using the Trinity
805      platform for reference generation and analysis. Nat. Protoc. 8:1494–1512. doi:
806      10.1038/nprot.2013.084.

807    Haerty W et al. 2007. Evolution in the fast lane: Rapidly evolving sex-related genes in *Drosophila*.
808    Genetics. 177:1321–1335. doi: 10.1534/genetics.107.078865.

809    Hahn MW, Han MV, Han S-G. 2007. Gene family evolution across 12 *Drosophila* genomes. PLoS
810    Genet. 3:e197. doi: 10.1371/journal.pgen.0030197.

811    Hollis B et al. 2019. Sexual conflict drives male manipulation of female postmating responses in
812    *Drosophila melanogaster*. Proc. Natl. Acad. Sci. U. S. A. 116:8437–8444. doi:
813    10.1073/pnas.1821386116.

814    Holloway AK, Begun DJ. 2004. Molecular evolution and population genetics of duplicated
815    accessory gland protein genes in *Drosophila*. Mol. Biol. Evol. 21:1625–1628. doi:
816    10.1093/molbev/msh195.

817    Holman L. 2009. *Drosophila melanogaster* seminal fluid can protect the sperm of other males.
818    Funct. Ecol. 23:180–186. doi: 10.1111/j.1365-2435.2008.01509.x.

819    Hopkins BR, Sepil I, Bonham S et al. 2019. BMP signaling inhibition in *Drosophila* secondary cells
820    remodels the seminal proteome and self and rival ejaculate functions. Proc. Natl. Acad. Sci. U. S.
821    A. 116(49):24719–24728. doi: 10.1073/pnas.1914491116.

822    Hopkins BR, Sepil I, Thézénas ML et al. 2019. Divergent allocation of sperm and the seminal
823    proteome along a competition gradient in *Drosophila melanogaster*. Proc. Natl. Acad. Sci. U. S. A.
824    116(36):17925-17933. doi: 10.1073/pnas.1906149116.

825    Howe KL et al. 2019. Ensembl Genomes 2020—enabling non-vertebrate genomic research.
826    Nucleic Acids Res. 48:D689–D695. doi: 10.1093/nar/gkz890.

827    Hu H et al. 2019. AnimalTFDB 3.0: a comprehensive resource for annotation and prediction of
828    animal transcription factors. Nucleic Acids Res. 47:D33–D38. doi: 10.1093/nar/gky822.

829    Huang DW, Sherman BT, Lempicki RA. 2009. Systematic and integrative analysis of large gene lists
830    using DAVID bioinformatics resources. Nat. Protoc. 4:44–57. doi: 10.1038/nprot.2008.211.

831    Imamura M, Haino-Fukushima K, Aigaki T, Fuyama Y. 1998. Ovulation stimulating substances in
832    *Drosophila biarmipes* males: Their origin, genetic variation in the response of females, and
833    molecular characterization. Insect Biochem. Mol. Biol. 28:365–372. doi: 10.1016/S0965-
834    1748(98)00004-6.

835    Kalb JM, Dibenedetto AJ, Wolfnert MF. 1993. Probing the function of *Drosophila melanogaster*
836    accessory glands by directed cell ablation. Proc. Natl. Acad. Sci. U. S. A. 90:8093–8097.

837    Karr TL, Southern H, Rosenow MA, Gossmann TI, Snook RR. 2019. The old and the new: Discovery
838    proteomics identifies putative novel seminal fluid proteins in *Drosophila*. Mol. Cell. Proteomics.
839    18:S23–S33. doi: 10.1074/mcp.RA118.001098.

840    Kelleher ES, Watts TD, LaFlamme BA, Haynes PA, Markow TA. 2009. Proteomic analysis of
841    *Drosophila mojavensis* male accessory glands suggests novel classes of seminal fluid proteins.
842    Insect Biochem. Mol. Biol. 39:366–371. doi: 10.1016/j.ibmb.2009.03.003.

843    Kopantseva MR et al. 1990. The proteins of the ejaculatory bulbs in different species of
844    *Drosophila*. Zh. Obshch. Biol. 51:125–140.

845    Leader DP, Krause SA, Pandit A, Davies SA, Dow JAT. 2018. FlyAtlas 2: a new version of the
846    *Drosophila melanogaster* expression atlas with RNA-Seq, miRNA-Seq and sex-specific data.
847    Nucleic Acids Res. 46:D809–D815. doi: 10.1093/nar/gkx976.

848    Lefevre G. 1976. A photographic representation and interpretation of the polytene chromosomes
849    of *Drosophila melanogaster* salivary glands. In: The Genetics and Biology of *Drosophila*.
850    Ashburner, M & Novitski, E, editors. Vol. Ia London: Academic Press pp. 31–66.

851    Leiblich A et al. 2012. Bone morphogenetic protein- and mating-dependent secretory cell growth
852    and migration in the *Drosophila* accessory gland. Proc. Natl. Acad. Sci. U. S. A. 109:19292–19297.
853    doi: 10.1073/pnas.1214517109.

854    Ludwig MZ et al. 1991. Genetic control and expression of the major ejaculatory bulb protein PEB-
855    me in *Drosophila melanogaster*. Biochem. Genet. 29:215–240.

856    Lung O et al. 2002. The *Drosophila melanogaster* seminal fluid protein Acp62F is a protease
857    inhibitor that is toxic upon ectopic expression. Genetics. 160:211–224.

858    Lung O, Wolfner MF. 1999. *Drosophila* seminal fluid proteins enter the circulatory system of the
859    mated female fly by crossing the posterior vaginal wall. Insect Biochem. Mol. Biol. 29:1043–1052.
860    doi: 10.1016/S0965-1748(99)00078-8.

861    Lung O, Wolfner MF. 2001. Identification and characterization of the major *Drosophila*
862    *melanogaster* mating plug protein. Insect Biochem. Mol. Biol. 31:543–551. doi: 10.1016/S0965-
863    1748(00)00154-5.

864    Martinson EO, Mrinalini,  Kelkar YD, Chang C, Werren JH. 2017. The evolution of venom by co-
865    option of single copy genes. Curr. Biol. 27(13):2007–2013.e8. doi:10.1016/j.cub.2017.05.032.

866    McBride CS. 2007. Rapid evolution of smell and taste receptor genes during host specialization

867    in *Drosophila sechellia*. Proc. Natl. Acad. Sci. U. S. A. 104:4996–5001. doi:
868    10.1073/pnas.0608424104.

869    Meikle DB, Sheehan KB, Phillis DM, Richmond RC. 1990. Localization and longevity of seminal-
870    fluid esterase 6 in mated female *Drosophila melanogaster*. J. Insect Physiol. 36:93–101. doi:
871    10.1016/0022-1910(90)90179-J.

872    Meslin C et al. 2015. Digestive organ in the female reproductive tract borrows genes from
873    multiple organ systems to adopt critical functions. Mol. Biol. Evol. 32(6):1567–1580. doi:
874    10.1093/molbev/msv048.

875    Minami R et al. 2012. The homeodomain protein defective proventriculus is essential for male
876    accessory gland development to enhance fecundity in *Drosophila*. PLoS One. 7:e32302. doi:
877    10.1371/journal.pone.0032302.

878    Misra S, Wolfner MF. 2020. *Drosophila seminal* sex peptide associates with rival as well as own
879    sperm, providing SP function in polyandrous females. Elife. 9:e58322. doi: 10.7554/eLife.58322.

880    Mohorianu I, Fowler EK, Dalmay T, Chapman T. 2018. Control of seminal fluid protein expression
881    via regulatory hubs in *Drosophila melanogaster*. Proc. R. Soc. B Biol. Sci. 285:20181681. doi:
882    10.1098/rspb.2018.1681.

883    Mueller JL et al. 2005. Cross-species comparison of *Drosophila* male accessory gland protein
884    genes. Genetics. 171:131–143. doi: 10.1534/genetics.105.043844.

885    Mueller JL, Page JL, Wolfner MF. 2007. An ectopic expression screen reveals the protective and
886    toxic effects of *Drosophila* seminal fluid proteins. Genetics. 175:777–783. doi:
887    10.1534/genetics.106.065318.

888    Nielsen H, Engelbrecht J, Brunak S, von Heijne G. 1997. Identification of prokaryotic and
889    eukaryotic signal peptides and prediction of their cleavage sites. Protein Eng. Des. Sel. 10:1–6.
890    doi: 10.1093/protein/10.1.1.

891    Nolte V, Pandey RV, Kofler R, Schlötterer C. 2013. Genome-wide patterns of natural variation
892    reveal strong selective sweeps and ongoing genomic conflict in *Drosophila mauritiana*. Genome
893    Res. 23:99–110. doi: 10.1101/gr.139873.112.

894    Obbard DJ et al. 2012. Estimating divergence dates and substitution rates in the *Drosophila*
895    phylogeny. Mol. Biol. Evol. 29:3459–3473. doi: 10.1093/molbev/mss150.

896    Ohashi YY, Haino-Fukushima K, Fuyama Y. 1991. Purification and characterization of an ovulation

897    stimulating substance from the male accessory glands of *Drosophila suzukii*. Insect Biochem.
898    21:413–419. doi: 10.1016/0020-1790(91)90008-3.

899    Ramm SA. 2020. Seminal fluid and accessory male investment in sperm competition. Phil. Trans.
900    R. Soc. B. 375: 20200068. doi: 10.1098/rstb.2020.0068.

901    Ravi Ram K, Ji S, Wolfner MF. 2005. Fates and targets of male accessory gland proteins in mated
902    female *Drosophila melanogaster*. Insect Biochem. Mol. Biol. 35:1059–1071. doi:
903    10.1016/j.ibmb.2005.05.001.

904    Ravi Ram K, Ramesh SR. 2003. Male accessory gland proteins in *Drosophila*: A multifacet field.
905    Indian J. Exp. Biol. 41:1372–1383.

906    Ravi Ram K, Wolfner MF. 2007. Seminal influences: *Drosophila* Acps and the molecular interplay
907    between males and females during reproduction. Integr. Comp. Biol. 47:427–445. doi:
908    10.1093/icb/icm046.

909    Richmond RC, Gilbert DG, Sheehan KB, Gromko MH, Butterworth FM. 1980. Esterase 6 and
910    reproduction in *Drosophila melanogaster*. Science. 207:1483–1485. doi:
911    10.1126/science.6767273.

912    Saudan P et al. 2002. Ductus ejaculatorius peptide 99B (DUP99B), a novel *Drosophila*
913    *melanogaster* sex-peptide pheromone. Eur. J. Biochem. 269:989–997. doi: 10.1046/j.0014-
914    2956.2001.02733.x.

915    Schmidt T et al. 1993. *Drosophila suzukii* contains a peptide homologous to the *Drosophila*
916    *melanogaster* sex-peptide and functional in both species. Insect Biochem. Mol. Biol. 23:571–579.
917    doi: 10.1016/0965-1748(93)90030-V.

918    Schmidt T, Stumm-Zollinger E, Chen PS, Böhlen P, Stone SR. 1989. A male accessory gland peptide
919    with protease inhibitory activity in *Drosophila funebris*. J. Biol. Chem. 264:9745–9749.

920    Sepil I et al. 2019. Quantitative proteomics identification of seminal fluid proteins in male
921    *Drosophila melanogaster*. Mol. Cell. Proteomics. 18:S46–S58. doi: 10.1074/mcp.RA118.000831.

922    Sheehan K, Richmond RC, Cochrane BJ. 1979. Studies of esterase 6 in *Drosophila melanogaster*.
923    III. The developmental pattern and tissue distribution. Insect Biochem. 9:443–450. doi:
924    10.1016/0020-1790(79)90062-3.

925    Singh A et al. 2018. Long-term interaction between *Drosophila* sperm and sex peptide is mediated
926    by other seminal proteins that bind only transiently to sperm. Insect Biochem. Mol. Biol. 102:43–

927     51. doi: 10.1016/j.ibmb.2018.09.004.

928     Sirot LK. 2019. On the evolutionary origins of insect seminal fluid proteins. Gen. Comp.
929     Endocrinol. 278:104–111. doi: 10.1016/j.ygcen.2019.01.011.

930     Sirot LK, Wong A, Chapman T, Wolfner MF. 2014. Sexual conflict and seminal fluid proteins: A
931     dynamic landscape of sexual interactions. Cold Spring Harb. Perspect. Biol. 7:a017533. doi:
932     10.1101/cshperspect.a017533.

933     Swanson WJ, Clark AG, Waldrip-Dail HM, Wolfner MF, Aquadro CF. 2001. Evolutionary EST
934     analysis identifies rapidly evolving male reproductive proteins in *Drosophila*. Proc. Natl. Acad. Sci.
935     U. S. A. 98:7375–7379. doi: 10.1073/pnas.131568198.

936     Tamura K, Subramanian S, Kumar S. 2004. Temporal patterns of fruit fly (*Drosophila*) evolution
937     revealed by mutation clocks. Mol. Biol. Evol. 21:36–44. doi: 10.1093/molbev/msg236.

938     Thurmond J et al. 2019. FlyBase 2.0: the next generation. Nucleic Acids Res. 47:D759–D765. doi:
939     10.1093/nar/gky1003.

940     Turner LM, Hoekstra HE. 2008. Causes and consequences of the evolution of reproductive
941     proteins. Int. J. Dev. Biol. 52:769–780. doi: 10.1387/ijdb.082577lt.

942     Vieira FG, Sánchez-Gracia A, Rozas J. 2007. Comparative genomic analysis of the odorant-binding
943     protein family in 12 Drosophila genomes: purifying selection and birth-and-death evolution.
944     Genome Biol. 8:R235. doi: 10.1186/gb-2007-8-11-r235.

945     Wagstaff BJ, Begun DJ. 2005. Comparative genomics of accessory gland protein genes in
946     *Drosophila melanogaster* and *D. pseudoobscura*. Mol. Biol. Evol. 22:818–832. doi:
947     10.1093/molbev/msi067.

948     Walker MJ et al. 2006. Proteomic identification of *Drosophila melanogaster* male accessory gland
949     proteins, including a pro-cathepsin and a soluble γ-glutamyl transpeptidase. Proteome Sci. 4:1–
950     10. doi: 10.1186/1477-5956-4-9.

951     Wigby S et al. 2020. The *Drosophila* seminal proteome and its role in postcopulatory sexual
952     selection. Trans R Soc Lond B Biol Sci. 375(1813):20200072. doi: 10.1098/rstb.2020.0072.

953     Wigby S, Chapman T. 2005. Sex peptide causes mating costs in female *Drosophila melanogaster*.
954     Curr. Biol. 15(4):316–21. doi: 10.1016/j.cub.2005.01.051.

955     Wolfner MF. 2009. Battle and ballet: Molecular interactions between the sexes in *Drosophila*. J.
956     Hered. 100:399–410. doi: 10.1093/jhered/esp013.

957  Wolfner MF. 2007. 'S.P.E.R.M.' (seminal proteins (are) essential reproductive modulators): the
958  view from Drosophila. Soc. Reprod. Fertil. Suppl. 65:183–199.

959  Wong A et al. 2008. A role for Acp29AB, a predicted seminal fluid lectin, in female sperm storage
960  in *Drosophila melanogaster*. Genetics. 180:921–931. doi: 10.1534/genetics.108.092106.

961  Wong A, Turchin MC, Wolfner MF, Aquadro CF. 2012. Temporally variable selection on
962  proteolysis-related reproductive tract proteins in Drosophila. Mol. Biol. Evol. 29:229–238. doi:
963  10.1093/molbev/msr197.

964  Xue L, Noll M. 2002. Dual role of the Pax gene paired in accessory gland development of
965  *Drosophila*. Development. 129:339–346.

966  Yang H et al. 2018. Re-annotation of eight *Drosophila* genomes. Life Sci. Alliance. 1:1–14. doi:
967  10.26508/lsa.201800156.

968  Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. Mol. Biol. Evol. 24:1586–
969  1591. doi: 10.1093/molbev/msm088.

970  Zhang Y, Sturgill D, Parisi M, Kumar S, Oliver B. 2007. Constraint and turnover in sex-biased gene
971  expression in the genus *Drosophila*. Nature. 450:233–237. doi: 10.1038/nature06323.

972  Zhao L, Saelao P, Jones CD, Begun DJ. 2014. Origin and spread of de novo genes in *Drosophila*
973  *melanogaster* populations. Science. 343(6172):769-72. doi: 10.1126/science.1248286.

974  Zhou Q et al. 2008. On the origin of new genes in *Drosophila*. Genome Res. 18:1446–1455. doi:
975  10.1101/gr.076588.108.

976  Zile K, Dessimoz C, Wurm Y, Masel J. 2020. Only a single taxonomically restricted gene family in
977  the *Drosophila melanogaster* subgroup can be identified with high confidence. Genome Biol.
978  Evol. 12(8):1355–1366. doi:10.1093/gbe/evaa127.

979

980  **Figure and Table Legends**

981  Fig. 1. Venn diagram representing the overlap between the candidate seminal genes we
982  identified (Candidates) and other sets of putative or confirmed *D. melanogaster* seminal genes.
983  Candidates are those genes we identified (1) to be highly (or differentially) expressed in the
984  accessory glands according to two transcriptomic databases and also (2) to encode secretory
985  proteins with two software programs. Known Seminal Genes (KSGs) are those encoding proteins
986  previously confirmed to be transferred by males into females during mating or those meeting

987    stringent multiple criteria that indicate so. Unconfirmed High Confident Candidates (UHCCs) are
988    those Candidates, not included among KSGs, that are both highly and differentially expressed in
989    the accessory glands according to the two consulted transcriptomic databases. Predicted but
990    unconfirmed seminal genes are previously predicted seminal genes not included among KSGs.

991    Fig. 2. Chromosomal location of *D. melanogaster* seminal genes. Drawings of polytene
992    chromosomes were modified from Lefevre's photographic maps (Lefevre 1976) and gene
993    locations were obtained from FlyBase.

994    Fig. 3. Mean *Ka/Ks* ($\omega$) across the *melanogaster* group for Known Seminal Genes (KSGs),
995    Unconfirmed High Confident Candidates (UHCCs), and candidate transcription factors driving the
996    expression of seminal genes in the accessory glands (TFs). TFs searches are described in the
997    Identification section and estimation procedures in Methods (Molecular Evolutionary Analyses).
998    The horizontal discontinuous line represents the mean value for all protein-coding genes
999    [according to Haerty et al. (2007)]. Different letters above boxes indicate differences between
1000   groups and * indicates differences between the group and the mean value (GLM followed by
1001   Tukey comparisons; $p < 0.05$).

1002   Fig. 4. Seminal genes shared between the *melanogaster* subgroup and other Diptera. Numbers
1003   refer to the 196 *Drosophila* orthogroups (generated with Orthofinder) having at least one seminal
1004   gene of the *melanogaster* subgroup. Orthogroups having seminal genes of various taxa are
1005   represented by overlapped areas.

1006   Fig. 5. Duplication (blue), loss (magenta), and *de novo* emergence (black) events among
1007   orthogroups containing *D. melanogaster* seminal genes. The numbers of events are shown per
1008   branch. Since orthogroups without *D. melanogaster* SFPs were not considered, *de novo* gains for
1009   branches outside the *D. melanogaster* lineage, which are zero, are not shown. Divergence times
1010   were obtained from Obbard et al. (2012).

1011   Fig. 6. Expected gene family topology for each considered mechanism of seminal gene origin.
1012   Ingroup genes represent *melanogaster* genes, while outgroup genes represent genes of any non-
1013   *melanogaster* group species for which seminal genes are known. Magenta branches correspond
1014   to seminal genes, while black branches correspond to non-seminal genes. Grey discontinuous
1015   branches stand for the absent of homologs.

1016   Fig. 7. Most likely mechanisms of origin of *D. melanogaster* seminal genes. Mechanisms were
1017   proposed according to our analysis of seminal gene families only for tentatively young seminal
1018   genes, i.e., those that have likely emerged after the split from the lineage leading to *D.*
1019   *pseudobscura*. Uncertain genes represent those we could not determine whether they are young
1020   or ancient.

1021    Table 1. List of Unconfirmed High Confident Candidates (UHCCs). Name, chromosomal location,
1022    and molecular function (taken from FlyBase r2020_03) are shown for each gene.
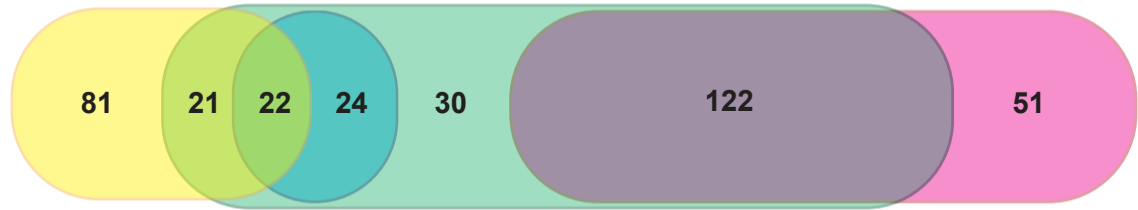
1023    Table 2. Molecular function annotation of Known Seminal Genes (KSGs) and Unconfirmed High
1024    Confident Candidates (UHCCs). For each group, count (and percentage) and false discovery rate
1025    (FDR) are shown for each GO term found with DAVID with more than one gene.

1026    Table 3. *D. melanogaster* seminal transcription factors candidates. Aligment e-value and the
1027    assigned DNA-binding domain family are shown for each candidate found with AnimalTFDB3.0.
1028    The first search was performed on genes whose expression strongly correlates to KSGs expression
1029    according to Ayroles et al. (2011). The second search was performed on genes whose expression
1030    is enriched in the male accessory glands according to modENCODE and FlyAtlas2 *D. virilis* search,
1031    which was performed using Blastp (alignment bit score > 80), shows the presence/absence of
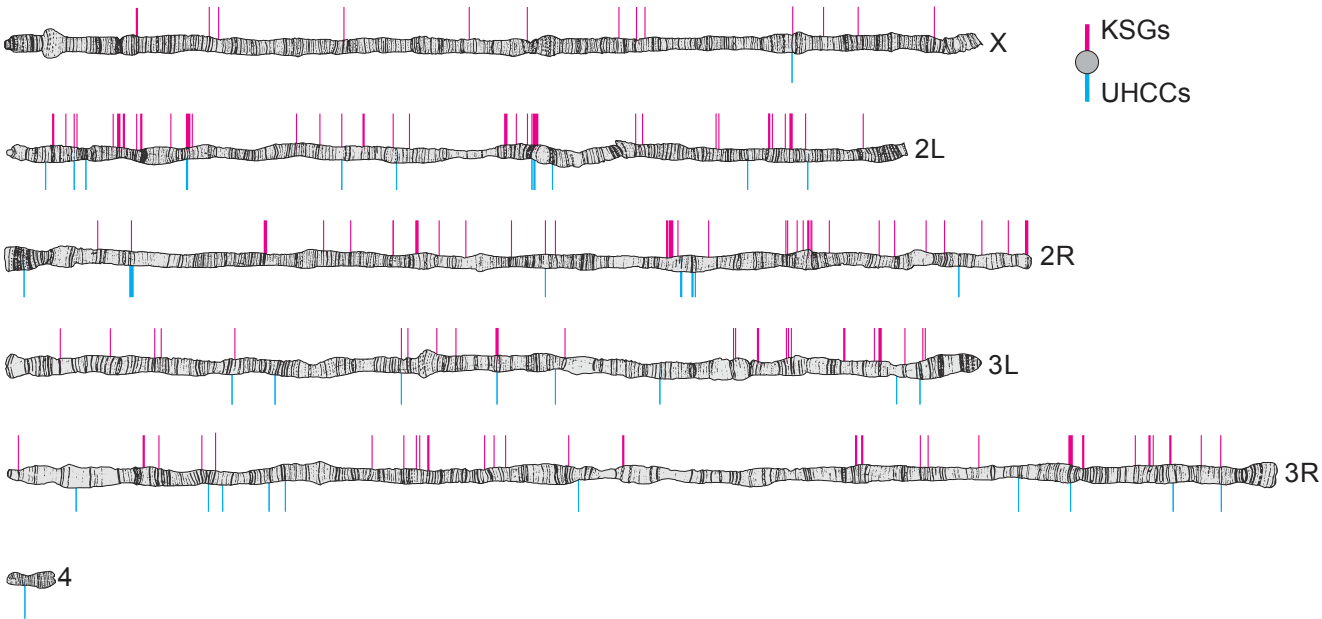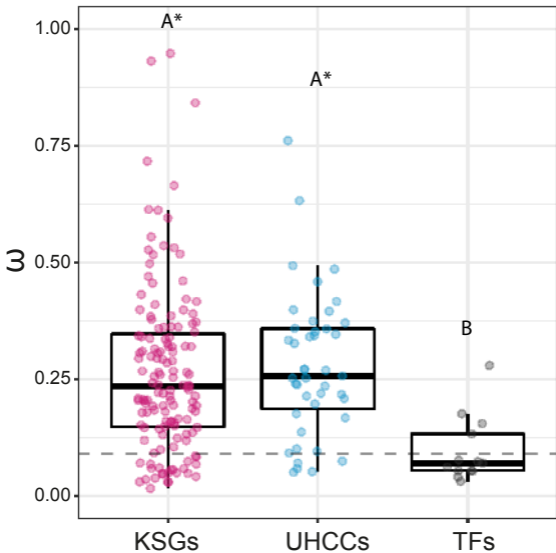1032    homologs among the *D. virilis* putative seminal TFs.
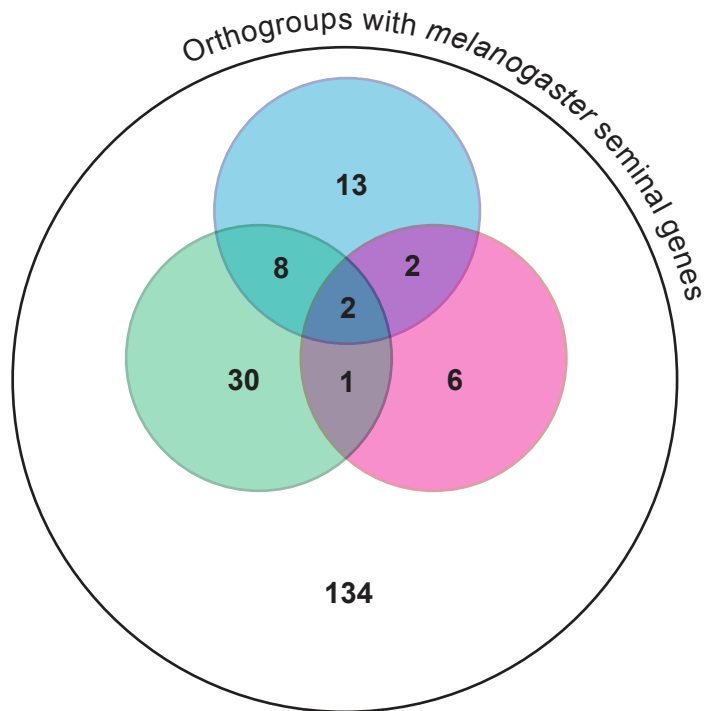
1033

Candidates (219)  UHCCs (46)
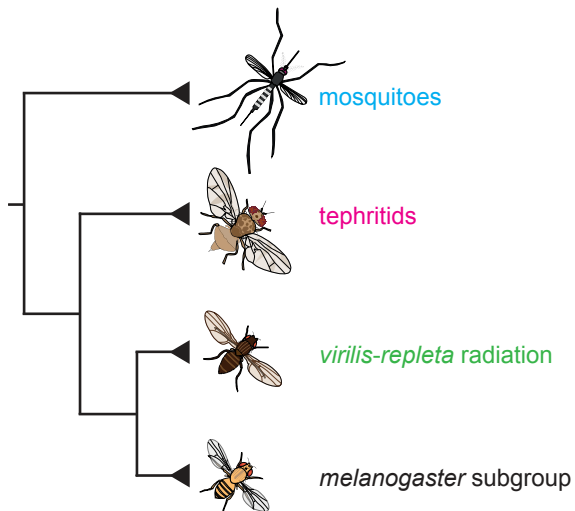
| 81 | 21 | 22 | 24 | 30 | 122 | 51 |

Predicted but unconfirmed (124)  KSGs (173)

X

2L

2R

3L

3R

4

KSGs

UHCCs

Orthogroups with *melanogaster* seminal genes

mosquitoes

tephritids

*virilis-repleta* radiation
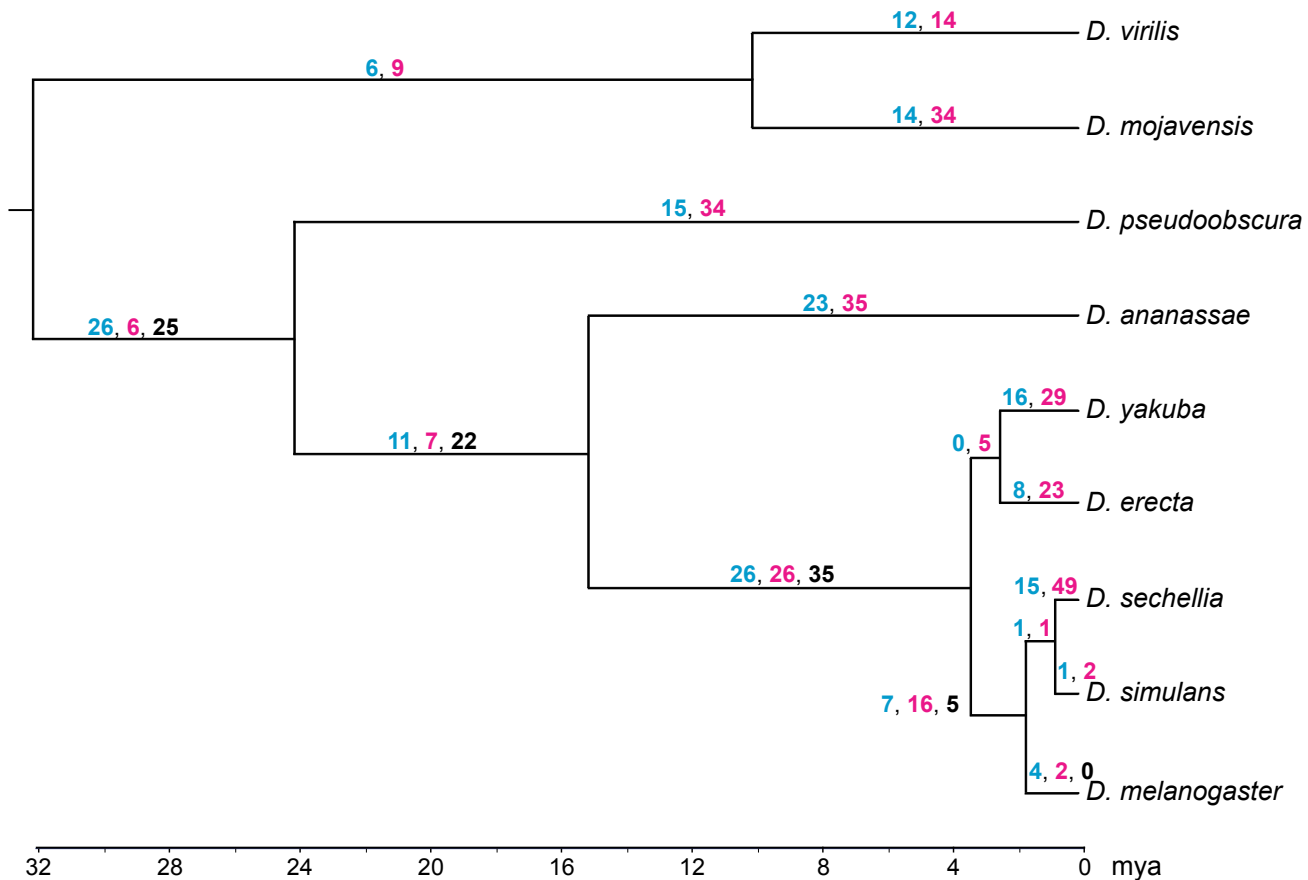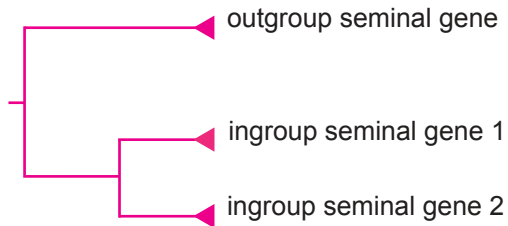
*melanogaster* subgroup

13

8

2

2

30

1

6

134

also containing (homologs to) seminal genes of:

mosquitoes    *virilis-repleta* radiation    tephritids

**Duplication**

- outgroup seminal gene
- ingroup seminal gene 1
- ingroup seminal gene 2

**Duplication-mediated co-option**

- outgroup non-seminal gene
- ingroup non-seminal gene
- ingroup seminal gene

**Co-option without duplication**

- outgroup non-seminal gene
- ingroup seminal gene

*De novo* emergence

- ingroup seminal gene

| Name or symbol | FlyBase ID | Novel candidate | Chromosomal location | Molecular function (GO) |
|---|---|---|---|---|
| Manf | FBgn0027095 | Yes | 3R | unknown |
| CG4271 | FBgn0031409 | Yes | 2L | serine-type endopeptidase/hydrolase activity |
| atilla | FBgn0032422 | Yes | 2L | unknown |
| CG17549 | FBgn0032774 | No | 2L | unknown |
| CG9336 | FBgn0032897 | Yes | 2L | unknown |
| CG11112 | FBgn0033164 | No | 2R | unknown |
| CG11113 | FBgn0033165 | No | 2R | unknown |
| Gbp1 | FBgn0034199 | Yes | 2R | cytokine activity |
| CG13557 | FBgn0034867 | Yes | 2R | unknown |
| CG12310 | FBgn0036467 | Yes | 3L | unknown |
| CG11977 | FBgn0037650 | No | 3R | unknown |
| CG8420 | FBgn0037664 | No | 3R | unknown |
| SPH202 | FBgn0039599 | No | 3R | serine-type endopeptidase activity |
| Lectin-21Ca | FBgn0040107 | No | 2L | carbohydrate binding |
| BG642312 | FBgn0047334 | No | 3L | unknown |
| CG31997 | FBgn0051997 | Yes | 4 | unknown |
| CG32382 | FBgn0052382 | No | 3L | serine-type endopeptidase/hydrolase activity |
| CG33290 | FBgn0053290 | No | 3L | unknown |
| Acp54A1 | FBgn0083936 | No | 2R | unknown |
| CG34299 | FBgn0085328 | Yes | 3R | unknown |
| CG34103 | FBgn0250831 | No | 3R | unknown |
| CG15394 | FBgn0250835 | No | 2L | unknown |
| CG42471 | FBgn0259961 | No | 2L | unknown |
| CG42481 | FBgn0259971 | Yes | 3L | unknown |
| CG42521 | FBgn0260396 | Yes | 3L | unknown |
| CG12163 | FBgn0260462 | Yes | 3R | cysteine-type peptidase/hydrolase activity |
| CG42852 | FBgn0262099 | Yes | 3L | unknown |
| CG43057 | FBgn0262359 | No | 2L | unknown |
| CG43061 | FBgn0262363 | No | 3R | unknown |
| CG43101 | FBgn0262547 | No | 2R | unknown |
| CG43123 | FBgn0262583 | No | 2R | unknown |
| CG43185 | FBgn0262814 | Yes | 2L | unknown |
| CG43254 | FBgn0262899 | Yes | 3R | unknown |
| CG43267 | FBgn0262948 | Yes | 2R | unknown |
| CG43350 | FBgn0263082 | Yes | 2L | serine-type endopeptidase inhibitor activity |
| CG43392 | FBgn0263249 | Yes | 3L | unknown |
| CG43679 | FBgn0263762 | Yes | 3L | unknown |
| CG43788 | FBgn0264329 | Yes | 2R | unknown |
| CG43789 | FBgn0264330 | Yes | 2R | unknown |
| CG44102 | FBgn0264911 | Yes | 2R | unknown |
| CG13639 | FBgn0265266 | No | 3R | unknown |
| CG18258 | FBgn0265267 | No | X | carboxylic ester hydrolase activity |
| CG44388 | FBgn0265538 | Yes | 2R | unknown |

| | | | | |
|---|---|---|---|---|
| CG44574 | FBgn0265785 | No | 2L | unknown |
| CG45011 | FBgn0266363 | No | 2L | unknown |
| CG45012 | FBgn0266364 | Yes | 2L | unknown |

| GO term | | KSGs | | UHCCs | |
|---|---|---|---|---|---|
| | | Count | FDR | Count | FDR |
| serine-type endopeptidase inhibitor activity | | 18 (10.4%) | 1.76E-16 | 0 | – |
| hormone activity | | 6 (3.5%) | 3.09E-04 | 0 | – |
| galactose binding | | 5 (2.9%) | 3.09E-04 | 0 | – |
| lipase activity | | 5 (2.9%) | 0.00414 | 0 | – |
| serine-type endopeptidase activity | | 11 (6.4%) | 0.00804 | 3 (6.5%) | 0.07748 |
| odorant binding | | 7 (4.0%) | 0.01045 | 0 | – |
| flavin-linked sulfhydryl oxidase activity | | 3 (1.7%) | 0.01045 | 0 | – |
| peptidase inhibitor activity | | 3 (1.7%) | 0.01362 | 0 | – |
| carbohydrate binding | | 6 (3.5%) | 0.01389 | 0 | – |
| hydrolase activity | acting on ester bonds | 4 (2.3%) | 0.02990 | 3 (6.5%) | 0.07739 |
| | carboxyesterase activity | 4 (2.3%) | 0.19229 | | |
| protein disulfide isomerase activity | | 3 (1.7%) | 0.09327 | 0 | – |
| thiol oxidase activity | | 2 (1.2%) | 0.18233 | 0 | – |
| unannotated | | 77 (44.5%) | – | 37 (80.4%) | – |

| Name or symbol | FlyBase ID | TF family | e-value | First search | Second search | *D. virilis* search |
|---|---|---|---|---|---|---|
| retn | FBgn0004795 | ARID | 3.10E-22 | Yes | No | No |
| CG7556 | FBgn0030990 | MYB | 5.00E-16 | Yes | Yes | No |
| prd | FBgn0003145 | PAX | 1.10E-71 | Yes | Yes | Yes |
| toe | FBgn0036285 | PAX | 1.00E-33 | Yes | Yes | Yes |
| CG13559 | FBgn0034870 | zf-LITAF-like | 5.30E-17 | Yes | Yes | No |
| CG6470 | FBgn0030933 | zf-C2H2 | 0.00020 | Yes | Yes | No |
| CG17841 | FBgn0028480 | TRAM_LAG1_CLN8 | 2.60E-63 | Yes | No | No |
| Myc | FBgn0262656 | bHLH | 5.90E-11 | Yes | No | Yes |
| CrebA | FBgn0004396 | TF_bZIP | 3.20E-15 | No | Yes | Yes |
| stc | FBgn0001978 | zf-NF-X1 | 1.10E-10 | No | Yes | Yes |
| CG3065 | FBgn0034946 | zf-H2C2_2 | 4.60E-22 | No | Yes | Yes |
| Bap111 | FBgn0030093 | HMG | 1.30E-16 | No | Yes | No |
| pzg | FBgn0259785 | zf-C2H2 | 7.50E-09 | No | Yes | Yes |
| CG11414 | FBgn0035024 | zf-C2H2 | 7.00E-05 | No | Yes | No |