

 Open access • Posted Content • DOI:10.1101/2021.06.10.447990

A putative de novo evolved gene required for spermatid chromatin condensation in *Drosophila melanogaster* — [Source link](#)

Emily L Rivard, Emily L Rivard, Andrew G. Ludwig, Prajal H. Patel ...+11 more authors

Institutions: Harvard University, College of the Holy Cross, University of Münster, Ohio State University ...+1 more institutions

Published on: 11 Jun 2021 - bioRxiv (Cold Spring Harbor Laboratory)

Topics: Chromatin, Comparative genomics, Spermatid, Gene and Drosophila melanogaster

Related papers:

- [A putative de novo evolved gene required for spermatid chromatin condensation in *Drosophila melanogaster*.](#)
- [Chromosomal rearrangements as a source of new gene formation in *Drosophila yakuba*](#)
- [A Continuum of Evolving De Novo Genes Drives Protein-Coding Novelty in *Drosophila*](#)
- [On the origin of new genes in *Drosophila*](#)
- [Heterochromatic genes in *Drosophila*: a comparative analysis of two genes.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/a-putative-de-novo-evolved-gene-required-for-spermatid-5fn8lyq7k0>

A putative *de novo* evolved gene required for spermatid chromatin condensation in *Drosophila melanogaster*

Emily L. Rivard^{1,4,*}, Andrew G. Ludwig^{1,*}, Prajal H. Patel^{1,*}, Anna Grandchamp², Sarah E. Arnold^{1,5}, Alina Berger², Emilie M. Scott¹, Brendan J. Kelly^{1,6}, Grace C. Mascha¹, Erich Bornberg-Bauer^{2,3}, Geoffrey D. Findlay^{1,**}

1. Department of Biology, College of the Holy Cross, Worcester, MA, USA

2. Institute for Evolution and Biodiversity, University of Münster, Münster, Germany

3. Department of Protein Evolution, Max Planck Institute for Developmental Biology, Tübingen, Germany

4. Current address: Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA, USA

5. Current address: Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA, USA

6. Current address: Department of Entomology, The Ohio State University, Columbus, OH, USA

*equal contribution

**Correspondence to: gfindlay@holycross.edu

25 **Abstract**

26

27 Comparative genomics has enabled the identification of genes that potentially evolved *de novo*
28 from non-coding sequences. Many such genes are expressed in male reproductive tissues, but
29 their functions remain poorly understood. To address this, we conducted a functional genetic
30 screen of over 40 putative *de novo* genes with testis-enriched expression in *Drosophila*
31 *melanogaster* and identified one gene, *atlas*, required for male fertility. Detailed genetic and
32 cytological analyses show that *atlas* is required for proper chromatin condensation during the
33 final stages of spermatogenesis. Atlas protein is expressed in spermatid nuclei and facilitates
34 the transition from histone- to protamine-based chromatin packaging. Complementary
35 evolutionary analyses revealed the complex evolutionary history of *atlas*. The protein-coding
36 portion of the gene likely arose at the base of the *Drosophila* genus on the X chromosome but
37 was unlikely to be essential, as it was then lost in several independent lineages. Within the last
38 ~15 million years, however, the gene moved to an autosome, where it fused with a conserved
39 non-coding RNA and evolved a non-redundant role in male fertility. Altogether, this study
40 provides insight into the integration of novel genes into biological processes, the links between
41 genomic innovation and functional evolution, and the genetic control of a fundamental
42 developmental process, gametogenesis.

43

44 Introduction

45

46 The evolution of new genes is integral to the extensive genotypic and phenotypic
47 diversity observed across species. The best-characterized mechanism of novel gene
48 emergence is gene duplication (Zhang 2003; Lipinski et al. 2011); however, rapid expansion in
49 high-quality genomic resources has provided mounting evidence of lineage-specific sequences
50 and the existence of alternative modes of new gene origination. One such mechanism is *de*
51 *novo* evolution, the birth of new genes from previously non-genic or intronic regions, which is
52 now a widely acknowledged source of protein-coding and RNA genes (McLysaght and Hurst
53 2016; Bornberg-Bauer and Schmitz 2017; Van Oss and Carvunis 2019). Although *de novo*
54 origination was once considered an unlikely event, catalogs of *de novo* genes have now been
55 published for an expansive range of species (Zhao et al. 2014; Guerzoni and McLysaght 2016;
56 Li et al. 2016; Ruiz-Orera et al. 2016; Lu et al. 2017; Zhang et al. 2019; Chamakura et al. 2020;
57 Puntambekar et al. 2020). Multiple models explain how protein-coding *de novo* genes may
58 acquire both an open reading frame (ORF) and regulatory sequences permitting transcription
59 (Carvunis et al. 2012; Schlötterer 2015; Wilson et al. 2017; Schmitz et al. 2018). Interrogation
60 of the biochemical and biophysical properties of the proteins encoded by *de novo* genes has
61 offered initial insight into the mechanisms of emergence and functional potential of these genes
62 (Schmitz et al. 2018; Vakirlis et al. 2018; Heames et al. 2020; Lange et al. 2021).

63 The capacity of protein-coding *de novo* genes to evolve important functions is a topic of
64 interest from evolutionary, physiological and molecular perspectives (Keeling et al. 2019). In the
65 last couple of decades, the products of *de novo* genes have been shown to play diverse roles in
66 a variety of organisms. For example, *de novo* genes function in fundamental molecular
67 processes in yeast, such as *BSC4*, a gene implicated in DNA repair, and *MDF1*, which mediates
68 crosstalk between reproduction and growth (Cai et al. 2008; Li et al. 2014). *De novo* genes also
69 evolve roles in organismal responses to disease and changing environmental factors. A
70 putatively *de novo* evolved gene in rice regulates the plant's pathogen resistance response to
71 strains causing bacterial blight (Xiao et al. 2009). Antifreeze glycoprotein genes, essential for
72 survival in frigid ocean temperatures, evolved *de novo* in the ancestor of Arctic codfishes to
73 coincide with cooling oceans in the Northern Hemisphere (Baalsrud et al. 2018; Zhuang et al.
74 2019). *De novo* genes are additionally implicated in the development and physiology of
75 mammals. In house mice, a *de novo* evolved gene expressed in the oviduct functions in female
76 fertility by regulating pregnancy cycles (Xie et al. 2019). A *de novo* gene found in humans and
77 chimpanzees regulates the oncogenesis and growth of neuroblastoma, revealing the relevance
78 of novel genes to human disease (Suenaga et al. 2014). These studies have started to
79 demonstrate the significance of *de novo* genes, thereby challenging previous assumptions that
80 only ancient, highly conserved genes can be essential.

81 Across multicellular animals, male reproductive tissues serve as hubs for new gene
82 emergence via numerous mechanisms, including *de novo* evolution (Marques et al. 2005;
83 Levine et al. 2006; Begun et al. 2007; Kaessmann 2010; Baker et al. 2012; Cui et al. 2015;
84 Ruiz-Orera et al. 2016). Proposed causes of this "out of the testis" phenomenon include the
85 high level of promiscuous transcription in testis cells (Soumillon et al. 2013; Necsulea and
86 Kaessmann 2014), the relative simplicity of promoter regions driving expression in the testis
87 (Sorourian et al. 2014), and preferential retention of novel genes with male-biased expression
88 (Palmieri et al. 2014). Sexual selection also drives rapid evolution of reproductive proteins
89 (Wilburn and Swanson 2016) and could drive the emergence of new genes as a mechanism of
90 improving male reproductive ability (Levine et al. 2006). The testis-biased expression of novel
91 genes, combined with growing evidence for new genes acting across a variety of tissue
92 contexts, suggests that many novel genes may function in male reproduction. For example, a
93 pair of young duplicate genes in *Drosophila*, *apollo* and *artemis*, are essential for male and
94 female fertility, respectively (VanKuren and Long 2018). Continued efforts to identify and

95 characterize testis-expressed novel genes are consequently critical for understanding the
96 genetic basis of male reproductive phenotypes.

97 *Drosophila* serves as an ideal system for interrogating the prevalence, sequence
98 attributes, expression patterns, and functions of testis-expressed *de novo* genes. The
99 availability of well-annotated genomes for numerous *Drosophila* species, the tractability of flies
100 to molecular genetics techniques, and our thorough understanding of *Drosophila* reproductive
101 processes facilitate comprehensive analyses of novel fly reproductive proteins (Demarco et al.
102 2014; Hales et al. 2015). As observed in other biological systems, *Drosophila de novo* genes
103 retained by selection demonstrate enriched expression in the testis (Levine et al. 2006; Begun
104 et al. 2007; Zhao et al. 2014; Heames et al. 2020). The expression patterns of emerging *de*
105 *novo* genes in the *Drosophila* testis were recently analyzed at single cell resolution (Witt et al.
106 2019), thereby providing insight into the dynamics of novel gene expression throughout
107 spermatogenesis. In addition to bioinformatic screens that have started to identify *de novo*
108 genes and large-scale expression analyses of testis-expressed genes, RNAi (Reinhardt et al.
109 2013) and CRISPR/Cas9-based (Kondo et al. 2017) functional screens have identified putative,
110 testis-expressed *de novo* genes required for fertility. However, a need remains for in-depth
111 experimental and evolutionary characterization of the genes identified in such screens. Detailed
112 examination of the function of *de novo* proteins will enable us to understand how these proteins
113 might integrate into existing gene networks and become essential.

114 We previously conducted a pilot functional screen of *de novo* genes with testis-enriched
115 expression in *D. melanogaster* and identified two novel genes, *goddard* and *saturn*, that are
116 required for full fertility (Gubala et al. 2017). *Goddard* knockdown males failed to produce any
117 sperm. *Saturn* knockdown males produced fewer sperm, which were inefficient at migrating to
118 female sperm storage organs. Subsequent characterization of *Goddard* using null deletion
119 alleles and a biochemically tagged rescue construct showed that the protein localizes to
120 elongating sperm axonemes and that, in its absence, individualization complexes associate less
121 efficiently with spermatid nuclei and do not successfully progress along sperm tails (Lange et al.
122 2021). These data suggested that putative *de novo* genes can evolve essential roles in a
123 rapidly evolving reproductive process, spermatogenesis.

124 Here, we expand this functional screen by evaluating whether any of 42 putative *de novo*
125 genes that show testis-enriched expression in *D. melanogaster* are required for male fertility.
126 We identified one gene, which we named *atlas*, whose knockdown or knockout results in nearly
127 complete male sterility. We show that *atlas* encodes a transition protein that facilitates
128 spermatid chromatin condensation. The *atlas* gene in *D. melanogaster* arose when a likely *de*
129 *novo* evolved protein-coding sequence moved off of the X chromosome and was inserted
130 upstream of a well-conserved non-coding RNA. While the *atlas* protein-coding sequence has
131 undergone multiple, independent gene loss events since its apparent origin at the base of the
132 *Drosophila* genus, the gene has evolved a critical function in *D. melanogaster*. These results
133 underscore the importance of detailed functional and evolutionary characterization in
134 understanding the origins of new protein-coding genes and the selective forces that affect their
135 subsequent evolution.

136

137 **Results**

138

139 **An RNAi screen identifies a putative *de novo* gene essential for *Drosophila* male fertility**

140

141 A previous pilot screen of 11 putative *de novo* evolved, testis-expressed genes identified

142 two genes that are critical for male fertility in *Drosophila melanogaster* (Gubala et al. 2017).

143 This result, and other recent work (e.g., Abrusán 2013; Zhang et al. 2015; VanKuren and Long

144 2018), suggested that lineage-specific, newly evolved genes can rapidly become important for

145 fertility, perhaps by gaining interactions with existing protein networks. To determine more

146 comprehensively the frequency with which potential *de novo* evolved genes become essential

147 for fertility, we identified *de novo* or putative *de novo* evolved genes with testis-biased

148 expression. A previous computational analysis identified genes that are detectable only within

149 the *Drosophila* genus, lack identifiable protein domains, and show no homology to other known

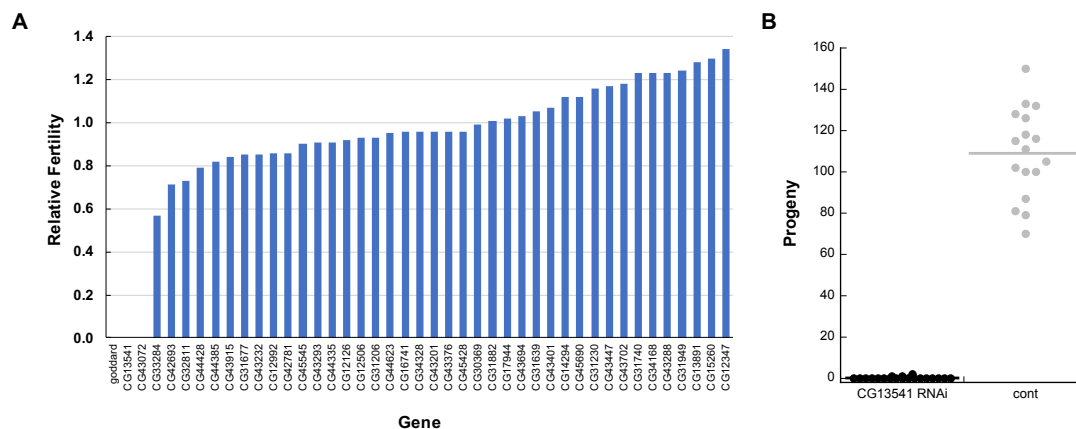
150 proteins through BLASTP and TBLASTN searches (Heames et al. 2020). We filtered these

151 genes to identify those expressed exclusively or predominantly in the testis, a common site of

152 *de novo* gene expression in animal species (Begun et al. 2007; Wu et al. 2011; Palmieri et al.

153 2014; Xie et al. 2019). This resulted in a set of 96 target genes.

154



155

156 **Figure 1. An RNAi screen of putative *de novo* genes identifies CG13541 as a major contributor to**

157 ***Drosophila melanogaster* male fertility.** A) All RNAi lines that showed at least partial knockdown of the

158 target gene were screened in group fertility assays (see Materials and Methods). Relative fertility was

159 calculated by dividing the average number of progeny produced per female mated to knockdown males

160 by the average number of progeny produced per female mated to control males in a contemporaneous

161 experiment. Relative fertility measurements lack error bars because each gene was tested in only 1-2

162 replicates. Knockdown of *goddard* was used as a positive control. B) A single-mating, single-pair fertility

163 assay confirms the observed defect when males are knocked down for *CG13541*, as knockdown males

164 showed significantly reduced fertility (control fertility (mean \pm SEM): 109.0 ± 5.3 ; knockdown fertility: $0.2 \pm$

165 0.1 ; two-sample *t*-test assuming unequal variances, $p = 5.6 \times 10^{-13}$).

166

167 We used testis-specific RNA interference to screen these genes for roles in male fertility.

168 We obtained RNAi lines from the Vienna *Drosophila* Resource Center (VDRC) and the

169 Transgenic RNAi Project (TRiP) and constructed additional lines using the TRiP-style

170 pValium20 vector (Ni et al. 2011), which is optimized for male germline expression. We

171 obtained and tested an RNAi line for 57 genes and induced knockdown for each with the *Bam-*

172 *GAL4* driver, which is expressed in the male germline, and enhanced with a copy of *UAS-*

173 *Dicer2*. RT-PCR confirmed at least partial knockdown in lines representing 42 genes (see

174 example in Fig. S1). We then screened knockdown males for fertility by allowing groups of 7

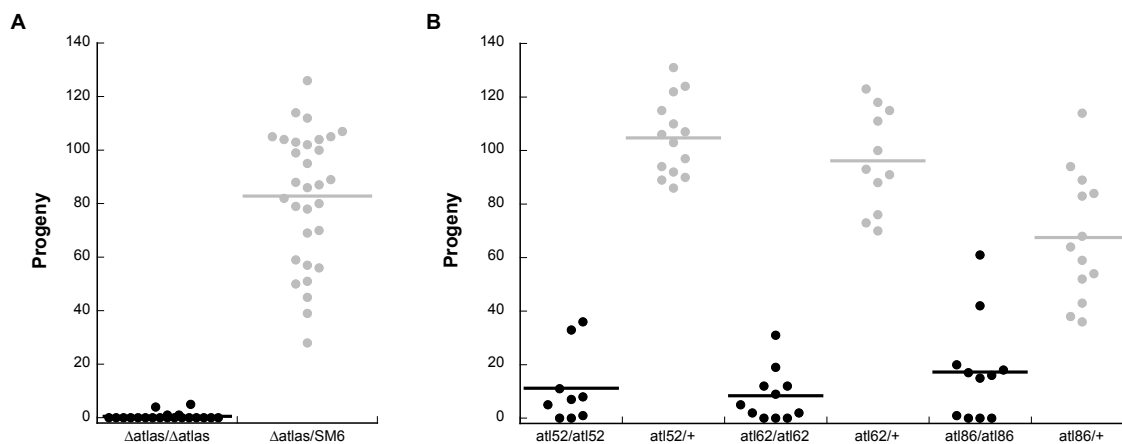
175 knockdown males to mate with 5 wild-type females for 2 days. Progeny counts were

176 standardized to the number of progeny produced by concurrently mated groups of 7 control

177 males and 5 wild-type females. The results are shown in Fig. 1A. This initial screen identified
178 *CG13541*, whose knockdown severely reduced male fertility. We confirmed the result for
179 *CG13541* by performing single-pair mating fertility assays (Fig. 1B). Knockdown of *CG43072*
180 and *CG33284* also showed consistent fertility defects, but these results did not replicate upon
181 testing with CRISPR-generated null alleles for each gene, so we do not consider them further
182 here. Consistent with our previous convention of naming testis-expressed genes after American
183 rocketry (Gubala et al. 2017), we will from here on refer to *CG13541* as *atlas*.

184 185 CRISPR-mediated gene mutation validates atlas RNAi results 186

187 We validated the observed fertility defect by using CRISPR/Cas9-based genome editing
188 to construct putative loss-of-function alleles for *atlas* (Fig. S2). The principal allele we used for
189 validation and further functional studies (described below) was a null allele that completely
190 deleted *atlas* from the genome. This allele was generated by targeting each end of the locus
191 with a gRNA. We made three additional frameshift alleles by inducing double-stranded breaks
192 at a gRNA target site just downstream of the *atlas* start codon, which induced non-homologous
193 end joining. Males homozygous for the *atlas* deletion allele have the same fertility defect as
194 knockdown males (Fig. 2A). Males homozygous for any of three frameshift alleles showed
195 significantly reduced, but non-zero, fertility (Fig. 2B). It is possible that residual *atlas* function
196 may be present in these animals, perhaps due to translation initiation at a downstream start
197 codon to create a shorter protein with partial function. Finally, we constructed a genomic rescue
198 construct carrying both the *atlas* coding region and its native regulatory sequences. *Atlas* null
199 males that carried a single copy of the rescue construct had fully restored fertility (Fig. S3).
200 Overall, these data demonstrate that *atlas* loss, and not an RNAi or CRISPR off-target, causes
201 nearly complete male sterility.
202



203
204
205 **Figure 2. CRISPR-generated deletion and frameshift alleles of *atlas* confirm the gene's**
206 **requirement for male fertility.** A) Single-pair fertility assay for males homozygous for the null ($\Delta atlas$)
207 allele or heterozygous controls ($\Delta atlas/SM6$). Flies homozygous for the deletion had significantly reduced
208 fertility (control fertility: 82.9 ± 4.5 ; null fertility: 0.3 ± 0.6 ; two-sample *t*-test assuming unequal variances, p
209 $= 5.4 \times 10^{-18}$). B) Single-pair fertility assays for males homozygous or heterozygous for three frameshift
210 alleles of *atlas* generated by imprecise non-homologous end joining at a CRISPR/Cas9 target site just
211 downstream of the start codon: *atlas*⁵² (control fertility: 104.7 ± 3.8 , mutant fertility: 11.2 ± 4.6 ; two-sample
212 *t*-test assuming unequal variances: $p = 8 \times 10^{-12}$), *atlas*⁶² (control fertility: 96.2 ± 5.7 ; mutant fertility: $8.4 \pm$
213 2.9 ; two-sample *t*-test assuming unequal variances: $p = 6.1 \times 10^{-10}$) and *atlas*⁸⁶ (control fertility: 67.5 ± 6.6 ;
214 mutant fertility: 17.3 ± 5.8 ; two-sample *t*-test assuming unequal variances: $p = 9.5 \times 10^{-6}$).
215

216 *Atlas is required for proper spermatid nuclear condensation*

217

218 We next examined how *atlas* loss of function impacted male fertility at the cellular level.

219 Dissection and phase-contrast imaging of *atlas* deletion null or knockdown male reproductive

220 tracts revealed that while the pre-meiotic and meiotic stages of spermatogenesis appeared

221 normal, sperm accumulated at the basal end of the testes, rather than in the seminal vesicles

222 (SVs), over the first week of adulthood (Fig. 3A and Fig. S4). To further characterize the fertility

223 defects in the absence of *atlas*, we examined the Mst35Bb-GFP (“protamine”-GFP) marker in

224 null or knockdown backgrounds (Manier et al. 2010). Mst35Bb encodes one of two protamine-

225 like proteins (highly similar paralogs of each other) that bind DNA in mature sperm. Its GFP

226 fusion construct thus allows visualization of nuclei in late stage spermiogenesis and mature

227 sperm. Consistent with the observed conglomeration of sperm tails at the basal testes, SVs

228 from either *atlas* null or knockdown males contain fewer mature sperm (Fig. 3B, S4C). The

229 nuclei of sperm from null males also appeared wider and less elongated than those of controls.

230 Together, these data suggest that *atlas* is required after meiosis, as developing spermatids take

231 on their final structures.

232 We next examined two post-meiotic processes: individualization of 64-cell spermatid

233 cysts into mature sperm, and spermatid nuclear condensation. Individualization initiates when

234 an actin-rich individualization complex (IC) associates with the bundle of spermatid nuclei. The

235 IC then proceeds down the sperm tails, expelling cytoplasmic waste and remodeling cell

236 membranes to form 64 individual sperm. We visualized this process in males 0-1 days old,

237 when spermatogenesis occurs at high levels, by staining whole mount testes for actin (Fig. 3C-

238 D). Although ICs could associate with nuclear bundles present at the basal end of the testes in

239 both control and *atlas* null males, we observed significantly fewer nuclear bundle-associated ICs

240 in nulls (Fig. 3C-D). While control testes typically had several ICs progressing down sperm tails,

241 we saw a significantly reduced proportion of progressed bundles in nulls (Fig. 3C-D). In some

242 null testes, we also observed individual investment cones dissociated away from progressing

243 ICs (Fig. 3C).

244 The ability of ICs to assemble at nuclear bundles and progress down sperm tails may be

245 reduced if nuclear condensation is aberrant (reviewed in Steinhauer 2015). During *Drosophila*

246 spermiogenesis, round spermatid nuclei undergo a series of stepwise, morphological changes

247 that are the product of two distinct, but related processes: changes in the chromatin packaging

248 of DNA, and changes in nuclear shape (Rathke et al. 2007; Fabian and Brill 2012; Rathke et al.

249 2014). The end result is thin, condensed nuclei. We quantified this process in testes dissected

250 from newly eclosed wild-type and *atlas* null males expressing Mst35Bb-GFP, which marks the

251 final stages of condensation. We shredded the post-meiotic region of the testes in the presence

252 of a fixative and counted the number of nuclear bundles that exhibited each of five stages of

253 condensation (Rathke et al. 2007): round nuclei, early canoe-stage (unmarked with Mst35Bb-

254 GFP), late canoe-stage (marked with Mst35Bb-GFP), elongated nuclei, and fully condensed

255 nuclei (Fig. 4). Condensation of the nuclear bundles in *atlas* null testes progressed at similar

256 rates to controls through the late canoe stage (Table S1). However, in *atlas* null males, all

257 nuclear bundles that progressed past the canoe stage (which included ~60% [range: 26-100%]

258 of all observed bundles) showed an aberrant “curled” phenotype (Fig. 4; Table S1). These data

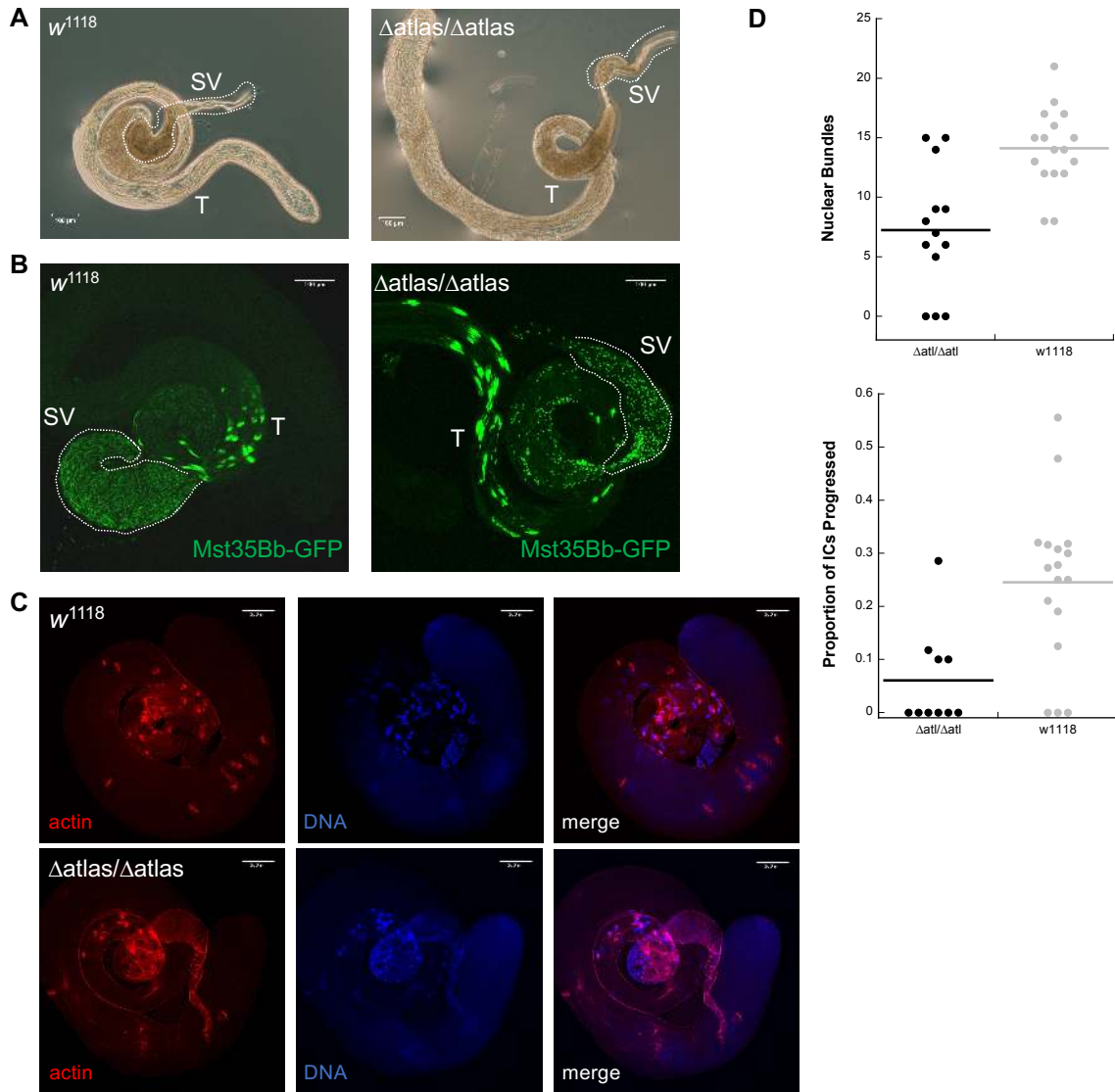
259 suggest that Atlas protein is required during the later stages of nuclear condensation and are

260 consistent with the idea that the loss of *atlas* affects nuclear shape in a way that reduces IC

261 assembly and sperm individualization (see Fig. 3C).

262

263



264

265

266

267

268

269

270

271

272

273

274

275

276

277

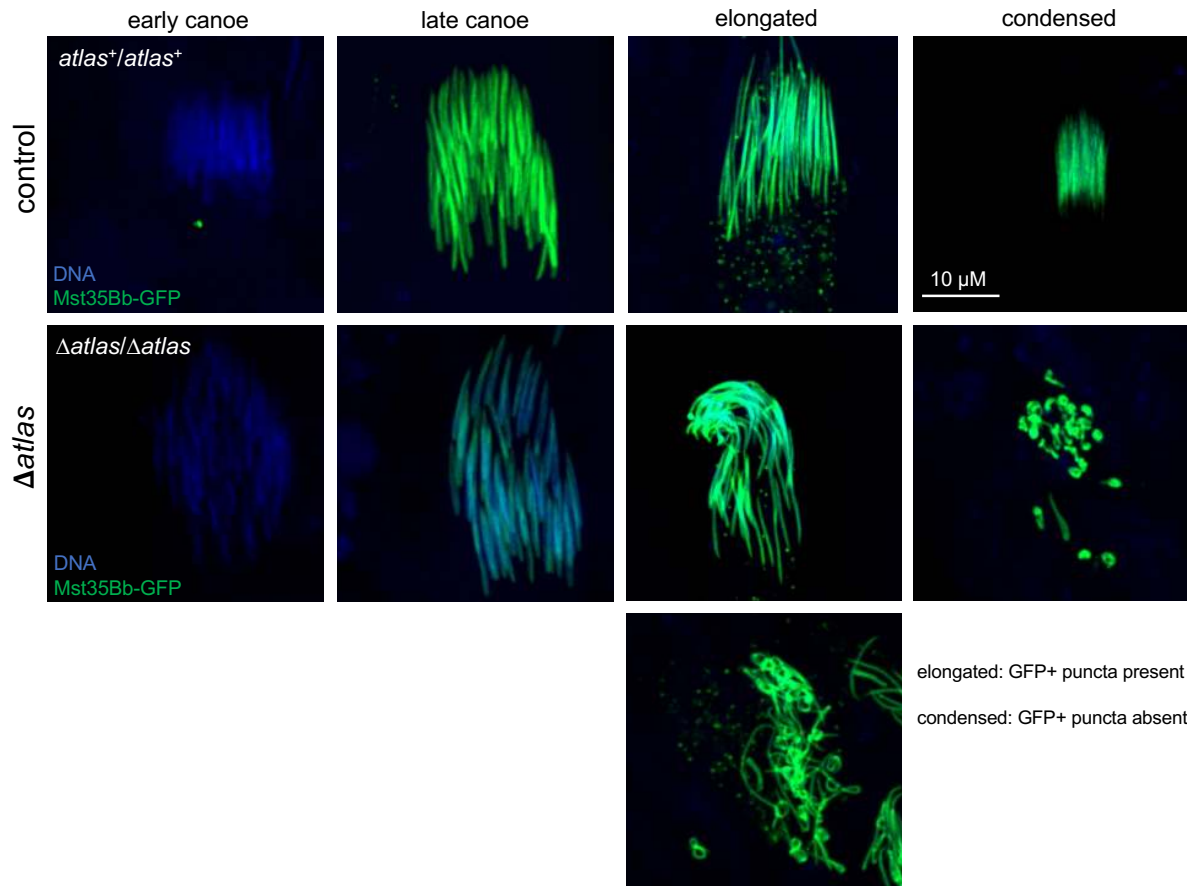
278

279

280

281

Figure 3. Cytological investigations of the *atlas* mutant fertility defect. A) Phase contrast microscopy of male reproductive tracts dissected from 7-day-old, unmated control (w^{1118}) or *atlas* null males. Control males show the expected accumulation of sperm in the seminal vesicle (SV), which appears here as a darker brown shading, while null male have an aberrant accumulation of sperm tails at the basal end of the testis (T). B) Visualization of Mst35Bb-GFP in 4-day old control and *atlas* null testes. While Mst35Bb is expressed in spermatid nuclei in the absence of *atlas*, the nuclei appear shorter and much less numerous in the outlined SV. C) Representative images from phalloidin staining of w^{1118} and *atlas* null testes used to assess the association of individualization complexes (ICs) with nuclear bundles and the progression of ICs down the length of sperm tails. D) At top, number of nuclear bundles with ICs associated in control and *atlas* null testes. Significantly more ICs were observed in control testes (control: $N = 17$, median = 14; mutant: $N = 13$, median = 7; Wilcoxon rank-sum test $W = 34$, $p = 0.0014$). At bottom, proportion of all observed ICs that were intact and that had progressed away from nuclear bundles. Three mutant testes with no observed ICs were excluded from the analysis. A significantly higher proportion of ICs progressed in control testes (control: $N = 17$, median = 0.27; mutant: $N = 10$, median = 0; Wilcoxon rank-sum test $W = 28$, $p = 0.0038$).



282
283
284
285
286
287
288
289

Figure 4. *Atlas* null males show aberrant nuclear shaping at and beyond the elongated stage of spermatid nuclear condensation. Early and late canoe stages were distinguished by the absence or presence of Mst35Bb-GFP, respectively. Elongated and condensed stages were distinguished by the presence or absence of GFP-positive puncta, respectively. As shown in Table S1, nuclear bundles from *atlas* null testes consistently took on a curved shape after the canoe stage, though the degree of curvature was variable, as exemplified above.

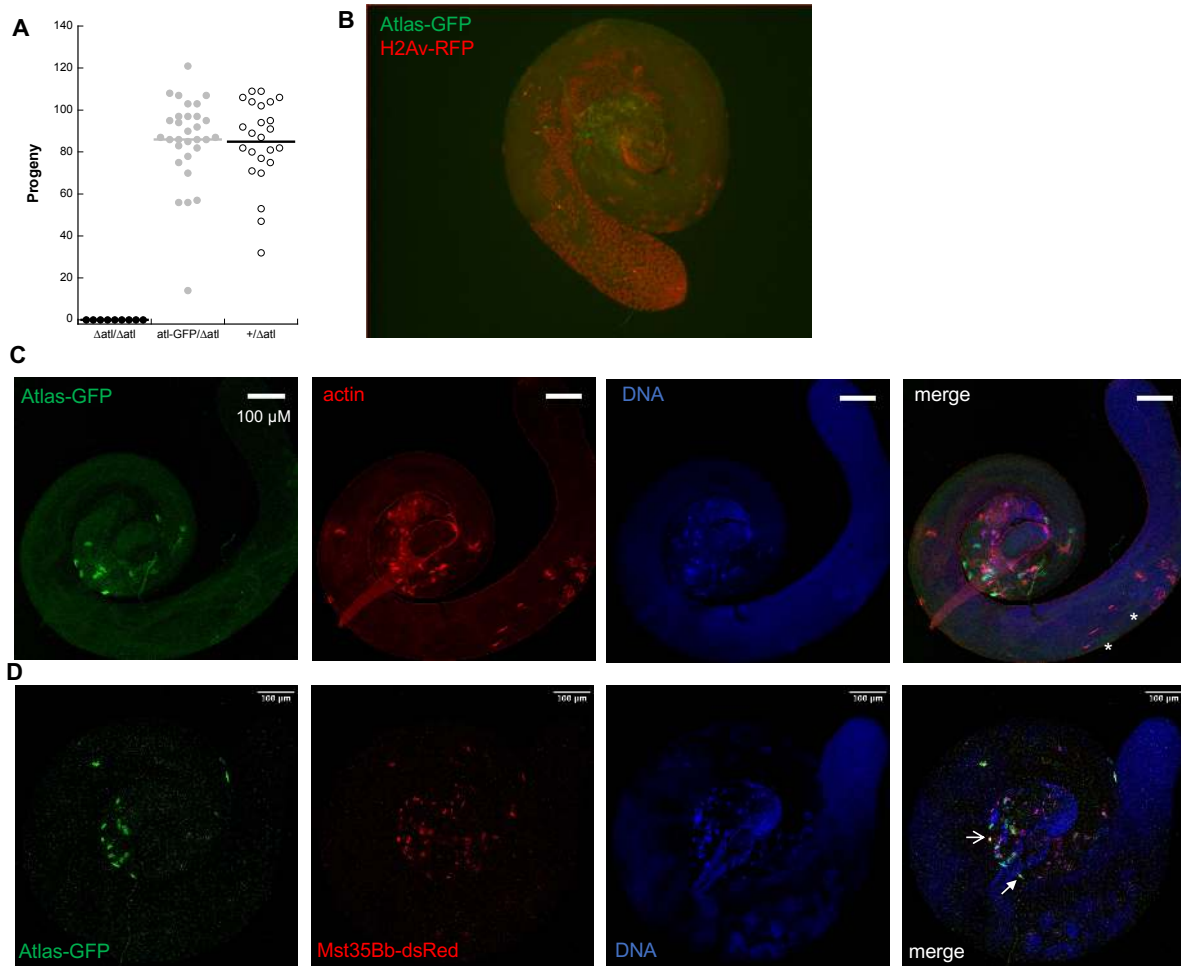
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305

That condensing spermatid nuclei are misshapen in the absence of *atlas* suggests the possibility that Atlas protein is critical for nuclear condensation. This idea is further supported by its predicted biochemical properties. Previously characterized spermatid chromatin binding proteins are small and highly basic (Jayaramaiah Raja and Renkawitz-Pohl 2005; Rathke et al. 2007; Kanippayoor et al. 2013), as the excess of positively charged amino acid side chains facilitates ionic interactions with negatively charged DNA. Many such proteins (i.e., Tpl94D, Mst35Ba, Mst35Bb, Prtl99C and Mst77F) also contain a conserved protein domain, the high-mobility-group box (HMG-box) domain (Dorus et al. 2008; Alvi et al. 2013; Rathke et al. 2014; Eren-Ghiani et al. 2015; Gärtner et al. 2015; Alvi et al. 2016; Gärtner et al. 2019), suggesting that this variety of chromatin binding proteins could have originated through gene duplication and divergence. Consistent with its putative *de novo* origin, Atlas lacks a detectable HMG-box domain. However, Atlas is otherwise similar to these other sperm chromatin binding proteins: the ~20 kDa protein has a highly basic predicted isoelectric point of 10.7, and its primary sequence contains the sequence KRDK, which matches the canonical consensus sequence for nuclear import, K(K/R)X(K/R) (Lange et al. 2007). To test the hypothesis that Atlas is nuclear localized, and could thus bind DNA, we generated an *atlas-GFP* transgene under UAS control

306 and expressed it ubiquitously using *tubulin*-GAL4 and in the early male germline using *Bam*-
307 GAL4. In both larval salivary glands and early male germline cells, atlas-GFP appeared to be
308 nuclear localized (Fig. S5).

309 While these results were consistent with Atlas protein localizing to the nucleus, they did
310 not allow us to visualize Atlas in the cells in which it is normally expressed. To do so, we used
311 CRISPR/Cas9-induced homology directed repair (<https://flycrispr.org/scarless-gene-editing/>)
312 (Bruckner et al. 2017; Bier et al. 2018; Hill et al. 2019) to create an *atlas*-GFP fusion at the
313 endogenous *atlas* locus (see Fig. S6 and Materials and Methods). We first confirmed the
314 functionality of the knock-in allele by showing that males with the *atlas* locus genotype *atlas*-
315 GFP/ Δ *atlas* had equivalent fertility to males of genotype *atlas*+/ Δ *atlas* (Fig. 5A). We then
316 visualized Atlas-GFP fusion protein in whole-mount testes in conjunction with phalloidin-stained
317 actin (Fig. 5C). Atlas-GFP was absent from seminal vesicles, consistent with its absence from
318 the proteome of mature *D. melanogaster* sperm (Dorus et al. 2006; Wasbrough et al. 2010).
319 Instead, Atlas-GFP colocalized with condensing nuclear bundles near the basal end of the
320 testes (Fig. 5C). Actin-based ICs were also observed in the basal testes, but generally did not
321 co-localize with Atlas-GFP, suggesting that Atlas-GFP is present in condensing nuclei before IC
322 association (Fig. 5C). This result, taken together with the aberrant nuclear condensation in the
323 absence of *atlas* (Fig. 4), is consistent with the idea that Atlas is a transition protein. Transition
324 proteins are chromatin components that act transiently during spermatid nuclear condensation.
325 A series of transition proteins first replace histones as the primary DNA binding proteins in the
326 nucleus and then give way to protamines, the proteins that package chromatin in mature sperm
327 (Rathke et al. 2007; Rathke et al. 2014; Gärtner et al. 2015).

328 To further elucidate the role of *atlas* in nuclear condensation, we next examined Atlas-
329 GFP localization in the presence of either an early spermatid nuclear marker, histone H2Av-
330 RFP (Schuh et al. 2007; Rathke et al. 2014), or Mst35Bb-dsRed (Manier et al. 2010), a marker
331 of nuclei from the late canoe stage through final condensation. Atlas-GFP showed no co-
332 localization with H2Av-RFP, suggesting that Atlas functions after histone removal (Fig. 5B). In
333 contrast, some GFP-positive bundles co-localized with Mst35Bb-dsRed, but others did not (Fig.
334 5D). These data suggest that Atlas may be one of the final transition proteins used in nuclear
335 condensation before the chromatin becomes fully condensed with protamines.
336



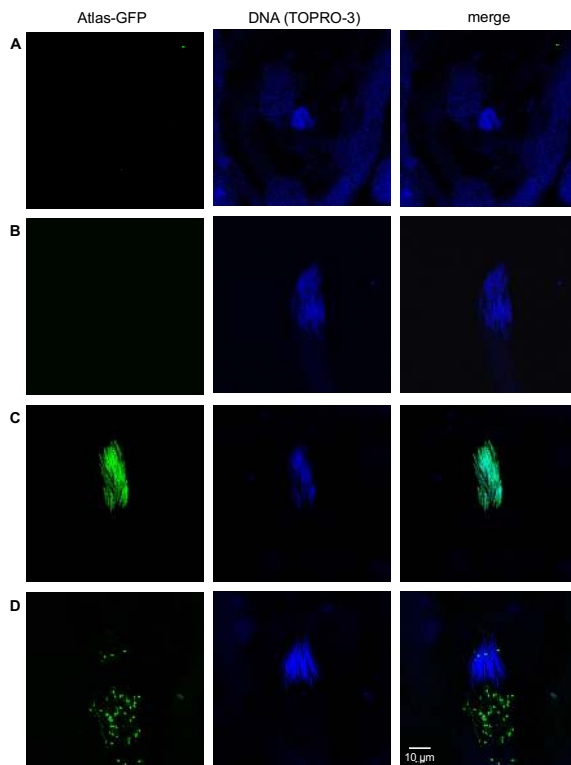
337
338

339 **Figure 5. An *atlas-GFP* allele generated at the endogenous *atlas* locus rescues the fertility defect**
340 **of *atlas* null flies and encodes a protein that localizes to condensing spermatid nuclei.** A) A single
341 copy of the *atlas-GFP* allele completely rescues the fertility defect caused by the $\Delta atlas$ allele and shows
342 equivalent fertility to males heterozygous for the wild-type *atlas* allele ($\Delta atlas/atlas-GFP$ fertility: $86.0 \pm$
343 3.2 ; $\Delta atlas/+$ fertility: 84.9 ± 4.1 ; two-sample t-test assuming unequal variances, $p = 0.85$). B) Atlas-GFP
344 does not co-localize with histone H2Av-RFP, a marker of the initial stages of spermatid nuclear
345 condensation. C) Visualization of Atlas-GFP in whole-mount testes from *atlas-GFP* homozygotes shows
346 that the fusion protein co-localizes with a subset of condensing spermatid nuclear bundles. While actin
347 associates with fully condensed nuclei at the basal testis, Atlas-GFP does not overlap and is also absent
348 from the seminal vesicle. Some Atlas-GFP is observed near progressing individualization complexes
349 toward the apical testis in the merged image (marked with asterisks; see also Fig. 6D). D) Atlas-GFP
350 partially colocalizes with Mst35Bb-dsRed, a marker of the final stage of nuclear condensation. Open
351 arrow: example of co-localization. Filled arrowhead: example of Atlas-GFP that does not co-localize with
352 Mst35Bb-dsRed. Collectively, these data suggest that *atlas* may serve as a transition protein involved in
353 the final stages of nuclear condensation.

354

355 To determine the stage(s) of nuclear condensation at which *atlas* functions, we analyzed
356 the shape of fixed nuclear bundles from shredded testes isolated from *atlas-GFP* males on the
357 day of eclosion. Based on the stage of the defect in *atlas* null males (Fig. 3-4) and the pattern
358 of Atlas-GFP-positive bundles in whole-mount testes (Fig. 5), we hypothesized that Atlas-GFP
359 would localize to the later stages of nuclear condensation. Consistent with this hypothesis, we

360 did not detect Atlas-GFP in round or early canoe stage bundles (Fig. 6A-B). Atlas-GFP co-
361 localized with DNA in late canoe stage bundles (Fig. 6C). Interestingly, when nuclei elongated
362 further, GFP was detected not in the nucleus, but as puncta basal to the nuclei (Fig. 6D; see
363 also Fig. 5C). Since Atlas-GFP is not observed in mature sperm in the SV (Fig. 5C), these data
364 suggest that Atlas may function as a transition protein that facilitates the condensation of
365 spermatid nuclei from histone-based DNA packaging to protamine-like-based DNA packaging
366 (Rathke et al. 2007) and is then removed from nuclei once protamines bind DNA. Indeed, the
367 appearance of Atlas in nuclei during the late canoe stage of condensation is similar to the
368 pattern observed for a previously characterized transition protein, Tpl94D (Rathke et al. 2007).
369 We hypothesize that the failure of *atlas* null sperm to form needle-like nuclei can be explained
370 by the absence of Atlas from the late canoe nucleus. It is also possible that the apparent
371 removal of Atlas-GFP from nuclei (Fig. 6D) represents a mechanism for removing transition
372 proteins from the nucleus after they exert their functions. We observed above that some
373 Mst35Bb-GFP also appears to be removed from the nucleus in puncta during the elongation
374 stage of nuclear condensation (see elongated stage of control nuclear bundles, Fig. 4), even
375 though other Mst35Bb-GFP molecules ultimately package DNA in mature, individualized sperm.
376 This could occur if Mst35Bb-GFP is present in excess of what is needed to package DNA.
377



378
379

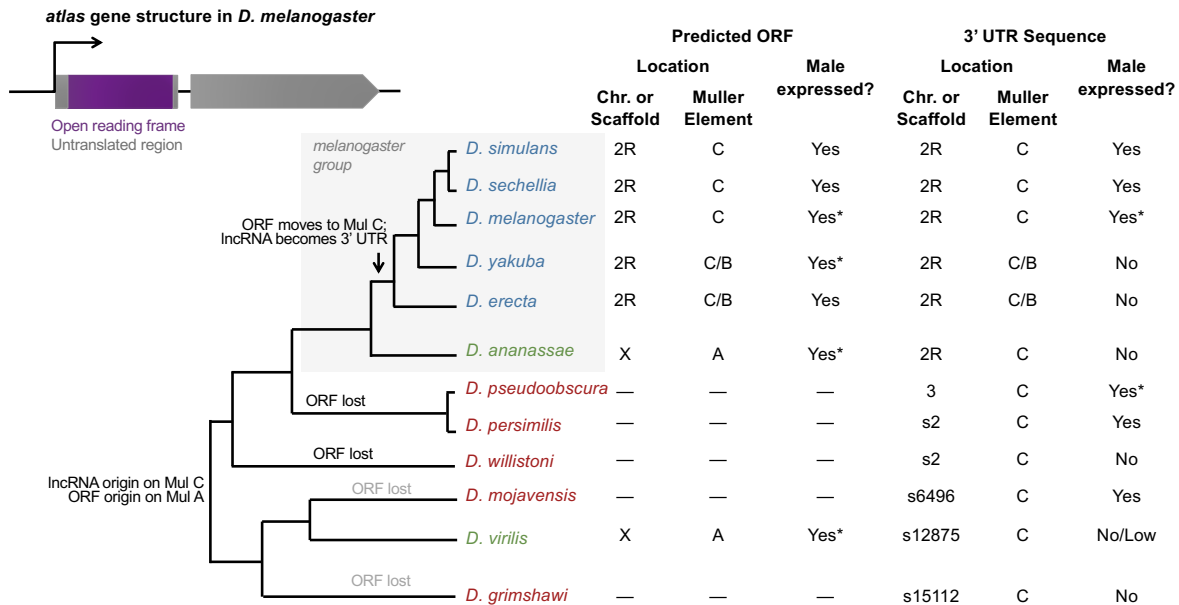
380 **Figure 6. Atlas-GFP is present in late canoe-stage spermatid nuclei and then appears to leave the**
381 **nucleus in puncta.** Staging of condensing spermatid nuclei fixed in paraformaldehyde from *atlas*-GFP
382 males stained with TO-PRO-3 DNA stain. Atlas-GFP is not detectable in (A) round stage or (B) early
383 canoe stage nuclei. Atlas-GFP is nuclear localized in the late canoe stage (C). When nuclei become fully
384 elongated (D), puncta of Atlas-GFP appear to be removed from the nucleus.

385
386

387 *Evolutionary origins of atlas*

388

389 To better understand the evolutionary origin of *atlas* and its evolution since emergence,
 390 we used a combination of BLAST- and synteny-based approaches to identify *atlas* orthologs
 391 throughout the genus (Gubala et al. 2017; Rele et al. 2020). One notable feature of this two-
 392 exon gene is that the protein-coding region (519 nucleotides) is contained entirely within the first
 393 exon (622 nt); the longer, second exon (910 nt) appears to be entirely non-coding (Fig. 7).
 394 Surprisingly, the second exon is more widely conserved. BLASTN detected significant matches
 395 to this region (range of hit length: 185-864 nt) in the same genomic location on Muller element C
 396 (see Schaeffer et al. 2008 for explanation of Muller elements), as assessed by synteny, in all
 397 *Drosophila* species examined, including distantly related species such as *D. virilis* and *D.*
 398 *grimshawi* (Fig. 7 and S7). The protein-coding first exon showed a more limited phylogenetic
 399 distribution. In most members of the *melanogaster* group of *Drosophila* (gray box in Fig. 7), this
 400 exon is found in a conserved position, adjacent to the non-coding region on the equivalent of *D.*
 401 *melanogaster* chromosome 2R (Table S2). In *D. ananassae*, however, the protein-coding
 402 region is found on the X chromosome (Muller element A). A putative ortholog for the protein-
 403 coding sequence is detectable by BLASTP in a partially syntenic region on the same Muller
 404 element in *D. virilis* (Table S2, Fig. S8). These data suggest that the *atlas* protein-coding
 405 sequence initially arose on Muller element A and then moved to Muller element C, giving rise to
 406 the gene structure observed in extant *D. melanogaster* and its sister species.
 407



408

409 **Figure 7. Molecular evolution and gene expression of *atlas* across the *Drosophila* genus.** The
 410 gene structure of *atlas* in *D. melanogaster* is shown at top left. The predicted protein-coding sequence is
 411 contained entirely within exon 1, while exon 2 encodes the presumed 3' UTR. The gene is located on
 412 chromosome 2R, equivalent to Muller element C. The phylogeny shows BLAST- and synteny-based
 413 detection of sequences orthologous to the protein-coding sequence and the 3' UTR sequence across
 414 *Drosophila* species. Sex-specific adult RNA-seq data were used to assess male expression across
 415 species, with RT-PCR verification performed in species marked with asterisks. RNA-seq data for the
 416 syntenic region of the 3' UTR in *D. virilis* were ambiguous; see Figs. S7 and S9.
 417

418

419 To confirm the lack of *atlas* protein-coding sequences identifiable by BLASTP or
 420 TBLASTN in most non-*melanogaster* group species, we identified the regions syntenic to those
 421 containing *atlas* in *D. ananassae* and *D. virilis* in 11 additional *Drosophila* species and used

421 more sensitive methods to search for potential orthologs (Rele et al. 2020). Specifically, we: a)
422 relaxed the BLAST cut-offs for detection, since default parameters can cause false-negative
423 results when searching for potential *de novo* genes in divergent species (Weisman et al. 2020);
424 b) used adult male RNA-seq data to detect transcribed areas within each syntenic region that
425 did not match annotated genes; and, c) predicted the isoelectric point of the potential proteins
426 encoded, under the hypothesis that Atlas orthologs would have conserved, DNA-binding
427 functions. The results are summarized in Fig. S8 and Table S2. These searches detected no
428 evidence for *atlas* orthologs in the following *Drosophila* species: *obscura*, *miranda*, *willistoni*,
429 *hydei*, *arizonae*, *mojavensis*, *navojoa*, and *grimshawi*. Some of these species had unannotated,
430 male-expressed transcripts in the regions syntenic to the Muller A location of *atlas* in *D.*
431 *ananassae* and *D. virilis*, but when each was compared with BLASTP to the *D. melanogaster*
432 proteome, all matched proteins other than Atlas, suggesting they may be lineage-specific
433 paralogs of other genes (Fig. S8). In sister species *D. pseudoobscura* and *D. persimilis*, we
434 detected a male-expressed transcript predicted to encode a protein with a pI > 10 in the region
435 syntenic to the location of *atlas* in *D. virilis*, but the predicted protein sequences showed no
436 significant BLASTP similarity to *atlas* orthologs (Table S2). While this predicted protein may
437 represent a divergent *atlas* ortholog, the abSENSE method predicts low probabilities of BLASTP
438 detection failure when searching for Atlas protein in these species (0.02 and 0.04, respectively),
439 so we favor the hypothesis of a lineage-specific, newly evolved gene in the region. Conversely,
440 in *D. busckii*, we detected, in the region syntenic to the *D. virilis atlas* locus, a male-expressed
441 transcript predicted to encode a protein with significant BLASTP identity to *D. melanogaster*
442 Atlas ($e = 4e-8$), but with a predicted pI of 5.1 and a ~50 percent shorter open-reading frame
443 (Table S2). The ortholog status of this predicted protein is also unclear, but because of its
444 dramatically altered size and pI, it is unlikely to have a functional role equivalent to that of *D.*
445 *melanogaster* Atlas.

446 To investigate whether the protein-coding region may have reproductive functions in
447 other species, we used sex-specific RNA-seq data from numerous *Drosophila* species curated
448 by the Genomics Education Partnership (Rele et al. 2020; thegep.org) and verified several of
449 these results by RT-PCR (Fig. 7, Fig. S9). In all species in which *atlas* was detected, the
450 protein-coding region was expressed specifically in males regardless of its genomic location
451 (Fig. 7). Interestingly, the non-coding region showed male-specific expression in species
452 lacking an unambiguous, orthologous coding region, such as *D. pseudoobscura* and *D.*
453 *mojavensis*. Conversely, while *D. yakuba* and *D. erecta* express the protein-coding region
454 robustly, we found no RNA-seq evidence to support expression of the non-coding second exon,
455 in spite of its sequence conservation (Fig. S7). Based on its high level of sequence
456 conservation, consistent genomic location and expression in a variety of species, it is possible
457 that what we now consider to be the 3' untranslated region of *atlas* from *D. melanogaster* was,
458 ancestrally, a non-coding RNA.

459 The FlyBase database reports two transcript isoforms of *atlas* in *D. melanogaster*: the
460 *atlas*-RA isoform is 986 nucleotides, while the *atlas*-RB isoform is 1528 nt. These isoforms
461 differ in how much of the second, non-coding exon is included in the transcript. We used RT-
462 PCR of whole male cDNA to assess the presence of these isoforms and their relative
463 abundances. Primers designed to amplify a region present in both isoforms produced products
464 that appeared more abundant than primers designed to amplify only the long isoform, even
465 though both primer pairs appeared to amplify genomic DNA with equal efficiency. Based on
466 RT-PCR band intensities and controlling for product size and genomic PCR band intensities, we
467 estimated that the short isoform is about 3-fold more abundant. This difference in abundance is
468 mirrored in available RNA-seq data, which show approximately 3- to 4-fold higher levels of
469 expression in the upstream part of exon 2 (Fig. S7), a pattern that also appears in *D. simulans*
470 and *D. sechellia*. Evaluating the potential significance of this finding awaits functional
471 characterization of the non-coding region.

472 As we have observed for other putative *de novo* genes with essential male reproductive
473 functions (Gubala et al. 2017), the pattern of *atlas* protein-coding sequence presence/absence
474 across the phylogeny is difficult to explain parsimoniously. If we assume that gene birth events
475 are less frequent than gene deaths, since the latter can occur through many possible mutational
476 events and can happen separately along multiple phylogenetic lineages, our data support the
477 hypothesis of a single origin of the protein-coding sequence at the base of the genus, followed
478 by independent loss events on the lineages leading to *D. grimshawi*, *D. mojavensis* and *D.*
479 *willistoni*, and potentially also *D. pseudoobscura/persimilis* (Table S2). We summarize these
480 findings for 12 representative species of *Drosophila* in Fig. 7. The general patterns of loss do
481 not change when all species of Table S2 are considered, though an additional loss in the
482 *melanogaster* group is likely due to the absence of a detectable ortholog in *D. kikkawai* and *D.*
483 *serrata*. As noted above, the pattern of gene loss can also appear due to orthology detection
484 failure (Weisman et al. 2020), for which we tried to account with our additional search methods
485 described above. We also note, however, that the probability of BLASTP-based ortholog
486 detection failure is relatively low for some *Drosophila* species that lack *atlas*, including *D.*
487 *pseudoobscura* (probability of non-detection due to divergence = 0.02), *D. persimilis* ($p = 0.04$)
488 and *D. willistoni* ($p = 0.06$). The probability is higher for other species, *D. mojavensis* ($p = 0.33$)
489 and *D. grimshawi* ($p = 0.66$), underscoring the importance of our additional search strategies.
490 Overall, our data support the hypothesis of multiple, independent loss events within *Drosophila*.

491 AbSENSE produces a 1.00 probability of BLASTP-based *atlas* ortholog detection failure
492 outside of *Drosophila*, reflecting the protein's short length and relatively rapid divergence (see
493 below). Indeed, the protein-coding and non-coding transcriptomes from each species showed
494 no matches to *Atlas* protein or cDNA sequences by BLAST. We thus used another synteny-
495 based approach, summarized in Fig. S10, to look for the protein-coding gene in other Dipterans
496 with well-resolved genomes: *Musca domestica*, *Glossina morsitans*, *Lucilia cuprina*, *Aedes*
497 *aegypti*, *Anopheles darlingi*, *Anopheles gambiae*, *Culex quinquefasciatus* and *Mayetiola*
498 *destructor*. In none of these species was a putative homolog found in any potential syntenic
499 region.

500 Recognizing the limitation of even this approach, we also used HMMER (Potter et al.
501 2018) to search iteratively either all genomes in ENSEMBL, or all metazoan genomes in
502 ENSEMBL, for annotated proteins with identity to *Atlas* from *D. melanogaster* or *D. virilis*.
503 These searches initially identified significant hits to the *Atlas* orthologs we identified above from
504 other *Drosophila* species. When these collections of orthologs were used as queries, no further
505 proteins outside of *Drosophila* were a significant match. As a control, we performed the same
506 search strategy with *D. melanogaster* Mst35Bb, a protein whose length, amino acid
507 composition, and function are similar to *Atlas*. These searches readily identified orthologs
508 throughout Diptera, consistent with predictions of its conservation from the OrthoDB database
509 (Zdobnov et al. 2017). Thus, we conclude that *atlas* is a putative *de novo* evolved gene that is
510 limited phylogenetically to the *Drosophila* genus.

511 Finally, we used standard tests of molecular evolution to examine the selective
512 pressures that have shaped *Atlas* protein within the *melanogaster* group. We aligned the *atlas*
513 protein-coding sequences from 12 species and used PAML to ask whether a model (M8_
514 allowing for positive selection, as well as neutral evolution and purifying selection, explained the
515 data better than models (M7 and M8a) that allowed only neutral evolution and purifying
516 selection (Yang 2007; McGeary and Findlay 2020). These data showed that while the *atlas*
517 protein-coding sequence's rate of evolution was accelerated relative to most *Drosophila* proteins
518 (whole-gene estimated d_N/d_S , $\omega = 0.41$ by PAML model M0), there was no significant evidence
519 for positive selection acting to recurrently diversify a subset of sites within the protein (Table 2).

520
521

522 **Table 2. PAML sites tests for positive selection acting on atlas in the *melanogaster***
523 **group.**

524

525 Model M0 (uniform ω across all sites): $\omega = 0.41$, $\ln L = -4032.24$, $np = 23$

526 Model M7 (10 site classes, each with $0 \leq \omega \leq 1$): $\ln L = -3961.28$, $np = 24$

527 Model M8 (10 site classes as in M7, plus one class with $\omega \geq 1$): $\ln L = -3960.41$, $np = 26$, ω for
528 extra class of sites = 1.39 (9.1% of sites)

529 Model M8a (10 site classes as in M7, plus one class with $\omega = 1$): $\ln L = -3961.02$, $np = 25$,
530 16.6% of sites in the $\omega = 1$ class

531

532 M7 vs. M8 likelihood ratio test: $\chi^2 = 1.74$, $df = 2$, $p = 0.42$

533 M8 vs. M8a likelihood ratio test: $\chi^2 = 1.22$, $df = 1$, $p = 0.27$

534

535

536

537 Discussion

538
539 Across taxa, many *de novo* evolved genes are expressed in the male reproductive
540 system (Levine et al. 2006; Begun et al. 2007; Cui et al. 2015; Ruiz-Orera et al. 2016).
541 Identifying those genes that have evolved essential roles in reproduction will provide insight into
542 how newly evolved genes integrate with existing cellular networks (Abrusán 2013) and how
543 evolutionary novelties permit adaptation in the face of sexual selection. Here, we screened 42
544 putatively *de novo* evolved genes for major effects on male *D. melanogaster* reproduction. Our
545 primary screen identified three genes whose knockdown caused an apparent reduction in male
546 fertility. However, subsequent CRISPR-mediated gene deletion revealed that only one of these
547 genes, *atlas*, was truly essential. This result underscores the importance of validating genes
548 identified in RNAi screens through traditional loss-of-function genetics and other approaches.

549 Using such genetic tools, we then showed that *atlas* loss of function reduces fertility by
550 affecting mature sperm production. During spermiogenesis, *atlas* mutants show aberrant
551 nuclear condensation and an inability to individualize spermatid bundles successfully. GFP-
552 tagged Atlas protein localizes to condensing spermatid nuclei in the basal testis and partially co-
553 localizes with Mst35Bb, a protamine around which DNA is wrapped in mature, individualized
554 sperm. Evolutionary analysis showed that the *atlas* protein-coding sequence likely arose at the
555 base of the *Drosophila* genus but was unlikely to have played an essential role immediately
556 upon birth, as the gene was subsequently lost along several independent lineages. Within the
557 *melanogaster* group of *Drosophila*, however, the gene moved from the X chromosome to an
558 autosome, where it formed a single transcriptional unit with a conserved, non-coding sequence.
559 Since this point, the gene has encoded a protein with a conserved length, isoelectric point and
560 male-specific expression pattern, suggesting potential functional conservation over the last ~15
561 million years.

562 Several lines of evidence suggest that Atlas is a transition protein that facilitates the
563 change from histone-based to protamine-based chromatin packaging in spermatid nuclei. Atlas
564 localizes throughout spermatid nuclei (Fig. 6) and has biochemical properties consistent with
565 direct DNA interaction. The protein appears specifically at the late canoe stage of nuclear
566 compaction (Fig. 5B-D). Its lack of overlap with testis-specific histones (Fig. 5C), partial overlap
567 with Mst35Bb (Fig. 5D), likely removal from needle-stage nuclei (Fig. 6) and absence from
568 mature sperm (Figs. 5-6) are all consistent with the expression profile of a transition protein.
569 Several other transition proteins have been characterized in *D. melanogaster*, including Tpl94D,
570 thmg-1, thmg-2, and Mst84B (Rathke et al. 2007; Gärtner et al. 2015; Gärtner et al. 2019).
571 Collectively, the transition proteins vary in the stage of nuclear condensation at which they first
572 appear and the range of nuclear shapes over which they are found (Hundertmark et al. 2018),
573 but otherwise match Atlas in their biochemical properties, transient roles, and localization
574 throughout the nucleus. Compared to these other transition proteins, Atlas is present over a
575 fairly narrow range of nuclear condensation stages and reaches its peak expression just prior to
576 the onset of individualization. *Atlas* is also the only transition protein gene characterized to date
577 whose removal disrupts fertility, as *Tpl94D*, *thmg-1*, *thmg-2* and *Mst84B* mutants are all fertile
578 (Rathke et al. 2007; Gärtner et al. 2015; Gärtner et al. 2019). This may reflect the relatively
579 later timing of Atlas's expression in spermatid nuclei, reduced functional redundancy between
580 DNA-binding proteins at the later stages of condensation, a potential interaction between Atlas
581 and an essential protamine-like protein, and/or a more stringent requirement for DNA binding at
582 these stages.

583 Transition proteins give way in spermatid nuclei to protamine-like proteins, which bind
584 DNA in mature sperm and persist through fertilization. In this way, protamine-like proteins
585 function analogously to vertebrate protamines, though they are believed to be evolutionarily
586 independent (Jayaramaiah Raja and Renkawitz-Pohl 2005; Doyen et al. 2015). In *D.*
587 *melanogaster*, protamine-like proteins include Mst35Ba, Mst35Bb, PrtI99C and Mst77F (Rathke

588 et al. 2010; Eren-Ghiani et al. 2015; Kimura and Loppin 2016). Interestingly, while all
589 characterized protamine-like proteins are present in mature sperm, only some are essential for
590 fertility. Knockouts of *Mst35Ba*, *Mst35Bb*, or both show occasional nuclear shaping defects, but
591 male fertility is normal (Rathke et al. 2010; Tirmarche et al. 2014). In contrast, mutants of
592 *Prtl99C* or *Mst77F* are sterile. *Prtl99C* and *Mst35Ba/b* bind condensed DNA independently of
593 each other and contribute additively to the shortening of needle-stage nuclei, but *Prtl99C*'s
594 effect is ~3x greater (Doyen et al. 2015; Eren-Ghiani et al. 2015). This difference is apparently
595 great enough to reduce fertility only in *Prtl99C* mutants. In contrast, *Mst77F* and *Mst35Ba/b*
596 show a genetic interaction, as *Mst35Ba/b* null flies become nearly sterile in an *Mst77F*
597 heterozygous background (Kimura and Loppin 2016). Furthermore, while *Mst35Bb*-GFP is
598 expressed in *Mst77F* nulls, these flies show deformed spermatid nuclei that do not reach a
599 recognizable needle-like stage. Because *atlas* nulls show considerable phenotypic similarity to
600 *Mst77F* nulls, but not *Prtl99C* nulls, we hypothesize that Atlas may act in a pathway with
601 *Mst77F*. Our observation of inefficient IC movement down sperm tails in *atlas* null testes is
602 reminiscent of a similar phenotype in *Mst77F* nulls (Kimura and Loppin 2016), providing further
603 evidence that these proteins may act in a common pathway. In both cases, ICs can form at
604 misshapen canoe-stage nuclei, but fail at a subsequent step. While the exact relationship
605 between nuclear abnormalities in the late canoe stage and individualization is not entirely
606 understood, it is possible that nuclear shape and the organization of nuclear bundles impact the
607 ability of IC association and IC progression, as is also observed in mutants of another gene,
608 *dPSMG1*, which controls nuclear shape (Gärtner et al. 2019).

609 Because *de novo* genes emerge from non-coding sequences, they typically encode
610 proteins that are short and lack complex structure (Schlötterer 2015; Van Oss and Carvunis
611 2019). Indeed, expanding the length of the protein-coding region and evolving higher-level
612 protein structures are hypothesized to be among the final stages of new gene evolution
613 (Bornberg-Bauer and Schmitz 2017). In light of these constraints, what kinds of cellular
614 functions might be available to newly evolved proteins? Vakirlis et al. (2020) overexpressed
615 emerging proto-genes in *S. cerevisiae* and found that those encoding proteins with predicted
616 transmembrane (TM) domains were more likely to be adaptive, as assessed by the effect of
617 proto-gene overexpression on growth rate. Such proteins may arise when thymine-rich
618 intergenic regions undergo mutations that allow protein-coding gene birth and expression, since
619 many codons with multiple U nucleotides encode amino acids commonly found in TM domains
620 (Vakirlis et al. 2020). Our imaging data, in addition to the prediction tools employed by Vakirlis
621 et al. (2020), suggest that Atlas does not contain a TM domain. However, just as the amino
622 acid compositional requirements of a TM domain are not overly complex, neither are those of
623 DNA binding proteins. In essence, these proteins must simply be small, have a high
624 concentration of positively charged residues, and contain a nuclear localization signal, which
625 itself requires a small patch of positively charged residues (Lange et al. 2007). Thus, DNA
626 binding proteins may be a relatively easy class of protein to evolve *de novo*.

627 While many putative *de novo* genes are expressed in the *D. melanogaster* testis (Fig. 1A
628 and Heames et al. 2020), *atlas* was the only verified hit from our screen that was essential for
629 male fertility. This result raises two related questions. First, why has selection maintained the
630 expression of the other potential *de novo* genes in our screen? In general, it is common for the
631 knockdown of protein-coding genes expressed in reproductive tissues in *D. melanogaster* to
632 result in no detectable fertility defects (Ravi Ram and Wolfner 2007; Schnakenberg et al. 2011;
633 Findlay et al. 2014; Gubala et al. 2017). One hypothesis to explain this pattern is that while the
634 loss of function of such genes may cause small reductions in fertility that would be subject to
635 strong negative selection in nature, the conditions used to assay such knockdown animals in the
636 primary screens are rarely tailored to detect differences of this magnitude. Another possibility,
637 not mutually exclusive, is that while the genes may be expendable in non-competitive, non-
638 exhaustive mating conditions, their absence may result in lower fitness in sperm competitive

639 environment, environments in which males mate several times in quick succession, or
640 environments in which sperm must persist in storage for longer intervals or during less optimal
641 conditions (Wong et al. 2008; Yeh et al. 2012; Civetta and Finn 2014).

642 A second question raised by our finding that *atlas* encodes an essential transition protein
643 is: how might *atlas* have evolved to become essential for fertility in *D. melanogaster*, particularly
644 when other transition proteins appear functionally redundant? Other proteins involved in
645 spermatid chromatin compaction show variable levels of conservation across *Drosophila*. For
646 example, the protamines around which DNA is wrapped in mature sperm (Jayaramaiah Raja
647 and Renkawitz-Pohl 2005) are found across all sequenced *Drosophila* species (Alvi et al. 2013),
648 and orthologs are also reported in FlyBase from other Dipteran and non-Dipteran insects (Larkin
649 et al. 2021). However, transition protein Tpl94D is reported to be restricted to species ranging
650 from *D. melanogaster* to *D. pseudoobscura* (Alvi et al. 2016), as are the related proteins tHMG1
651 and tHMG2 with high-mobility group domains (Gärtner et al. 2015; Larkin et al. 2021). Results
652 like these suggest that while some protamine-like proteins (i.e., Mst35Ba and Mst35Bb) have
653 consistently been among the final chromatin-packaging proteins, the specific proteins facilitating
654 the transition from histones to protamines have likely varied over evolutionary time. Against this
655 backdrop, and based on our analyses of the protein's presence/absence, biochemical
656 properties, and expression patterns in extant species, we hypothesize that while the Atlas
657 protein likely had some DNA-binding ability and male-specific expression upon its origin, it was
658 only one of several proteins involved in spermatid chromatin compaction. Since *atlas* was lost
659 independently in several lineages after its birth (Fig. 7), Atlas was likely non-essential at its
660 outset, but rather evolved an essential function within the *melanogaster* group of species. Such
661 evolution of essentiality could have occurred because of the loss of a protein with a
662 complementary function and/or changes in the process of spermatogenesis that thrust Atlas into
663 a functionally unique role. It is also worth noting that species that have evidently lost *atlas* might
664 have undergone other compensatory changes in their repertoires of spermatid DNA binding
665 proteins. For example, *D. willistoni* lacks *atlas* but appears to have several additional paralogs
666 of the protamines found only in duplicate in *D. melanogaster*, which could have evolved
667 transition-protein-like roles.

668 While our study cannot establish whether the movement of the *atlas* protein-coding
669 sequence off of the X chromosome onto an autosome early in the evolution of the *melanogaster*
670 group (Fig. 7) affected the gene's essentiality, such movement remains noteworthy. Prior work
671 has found a significant dearth of testis-expressed genes on the X chromosome in *Drosophila*
672 (Parisi et al. 2003; Parisi et al. 2004; Dorus et al. 2006; Vibranovski, Zhang, et al. 2009) and
673 other species (Emerson et al. 2004; Reinke et al. 2004). Furthermore, *Drosophila* exhibit
674 suppression of X-linked testis-expressed genes, and transfer of such genes from the X
675 chromosome to autosomal loci results in higher expression levels (Kemkemer et al. 2014;
676 Argyridou et al. 2017; Argyridou and Parsch 2018). One of several proposed mechanisms for
677 both the paucity of X-linked testis-expressed genes and the suppression of their expression is
678 meiotic sex chromosome inactivation (MSCI), in which the X chromosome becomes
679 transcriptionally silenced earlier than autosomes (Vibranovski, Lopes, et al. 2009; Zhang,
680 Vibranovski, Krinsky, et al. 2010; Zhang, Vibranovski, Landback, et al. 2010; Vibranovski et al.
681 2012; Gao et al. 2014; Mahadevaraju et al. 2021). Thus, genes that affect meiotic or post-
682 meiotic processes, as *atlas* does, could exert beneficial effects more strongly and/or for a longer
683 period of time if they become encoded autosomally. While the *atlas* protein-coding sequence
684 appears to show male-specific expression regardless of its chromosomal location, it is possible
685 that the movement of *atlas* to chromosome 2 allowed it to evolve a broader or different
686 expression pattern that expanded or modified its role in spermiogenesis. The complex
687 molecular bases of both X suppression and "escape" from the X chromosome in *Drosophila*
688 continue to be actively investigated and debated (Meiklejohn et al. 2011; Mikhaylova and
689 Nurminsky 2011; Vibranovski et al. 2012; Gao et al. 2014; Vibranovski 2014; Landeen et al.

690 2016; Mahadevaraju et al. 2021), but continued research in this area might inform further
691 interrogation of the forces driving *atlas* off of the X chromosome.

692 The movement of the *atlas* protein-coding sequence to chromosome 2 also created the
693 two-exon gene observed in *D. melanogaster*, in which the longer second exon appears to be
694 entirely non-coding. This second exon is highly conserved across the genus in both sequence
695 and genomic location, and it shows male-specific expression in several species that lack the
696 protein-coding sequence upstream (Fig. 7 and Fig. S7). These patterns of conservation
697 suggest that the second exon might originally have been a non-coding RNA, a class of molecule
698 whose importance in *Drosophila* male reproduction has recently become recognized (Wen et al.
699 2016; Bouska and Bai 2021). While these previous examples of functional ncRNAs in
700 spermatogenesis have generally acted in *trans* to regulate other genes or affect the functions of
701 other proteins, it is also possible that the long 3' UTR of *atlas* in *D. melanogaster* could affect
702 the translation of *atlas* transcripts. Many genes functioning in spermatid differentiation are
703 transcribed early in spermatogenesis but translationally repressed until later in spermiogenesis,
704 a phenomenon that relies on various forms of post-transcriptional regulation (White-Cooper
705 2010; Lim et al. 2012). Future studies of the *atlas* protein-coding sequence in the absence of its
706 3' UTR, the expression patterns of Atlas protein in species in which it is encoded from the X
707 chromosome, or the genetic ablation of the conserved region in species lacking the protein-
708 coding sequence will provide additional insights.

709 A final issue raised by our results is the exact timing and mechanism of origin for the
710 *atlas* protein-coding sequence. The bioinformatic screen (Heames et al. 2020) that identified
711 *atlas* and the other genes tested in Fig. 1 was designed to identify both “*de novo*” genes,
712 defined as protein-coding regions in *Drosophila* that had recognizable, but non-ORF-
713 maintaining, TBLASTN hits in outgroup species, and “putative *de novo*” genes, which had no
714 TBLASTN hits in outgroup species. (Importantly, the screen also eliminated any protein with an
715 identifiable protein domain, thus reducing the chances of identifying divergent members of gene
716 families.) The vast majority of the genes we tested with RNAi, including *atlas*, fell into the
717 putative *de novo* category. The bioinformatic screen’s criteria were reasonable for a high-
718 throughput analysis, but BLAST-based methods have known limitations for detecting
719 orthologous sequences in diverged species (Moyers and Zhang 2015, 2018; Weisman et al.
720 2020). The lack of identifiable *atlas* protein-coding genes in several *Drosophila* species (e.g., *D.*
721 *pseudoobscura* and *D. willistoni*) is unlikely to be due to BLAST homology detection failure, and
722 extensive synteny-based searches confirmed the gene’s absence (Fig. S9). BLAST and
723 synteny-based searches for orthologs in non-*Drosophila* species also did not detect an ortholog,
724 though BLAST searches are not predicted to have adequate sensitivity, for a protein of this size
725 and evolutionary rate, at this level of species divergence (Weisman et al. 2020). Hence, in
726 addition to using synteny to search for orthologs, we used HMMER, which employs hidden
727 Markov models and builds a sequence profile of the target protein using information from
728 multiple orthologs. Since HMMER also did not detect orthologs outside of *Drosophila*, we
729 hypothesize that *atlas* evolved *de novo* at the base of the genus. However, since we remain
730 unable to identify the non-protein-coding sequence from which *atlas* arose, we continue to refer
731 to *atlas* as a putative *de novo* gene (McLysaght and Hurst 2016).

732 Overall, we find that while many putative *de novo* evolved genes are expressed in the *D.*
733 *melanogaster* testes, few have major, non-redundant effects on fertility. However, several such
734 genes have acquired critical roles, acting at distinct stages of spermatogenesis and sperm
735 function. We showed previously that the putative *de novo* gene *saturn* is required for maximal
736 sperm production, as well as for the ability of transferred sperm to migrate successfully to sperm
737 storage organs in females (Gubala et al. 2017). Another putative *de novo* gene, *goddard*, is
738 required for sperm production and encodes a cytoplasmic protein that appears to localize to
739 elongating axonemes (Gubala et al. 2017; Lange et al. 2021). Loss of *goddard* impairs the
740 individualization of spermatid bundles (Lange et al. 2021), thus exerting an effect that appears

741 to be upstream of those observed for *saturn* and *atlas*. Here, we report another novel function
742 for a putative *de novo* gene: encoding an essential transition protein that is necessary for proper
743 nuclear condensation in spermiogenesis. Taken together, these results demonstrate that while
744 many *de novo* genes may play subtle roles or share functional redundancy with other genes, *de*
745 *novo* genes can also become essential players in complex cellular processes that mediate
746 successful reproduction.
747

748 **Materials and Methods**

749

750 *RNA Interference Screen*

751

752 *De novo* and putative *de novo* genes inferred to be no older than the *Drosophila* genus
753 were identified previously (Heames et al. 2020). We filtered these genes with publicly available
754 RNA-seq data (Brown et al. 2014) to identify those expressed predominantly in the testes
755 (>50% of RPKM sum deriving from the testes from ModENCODE data; Brown et al. 2014),
756 giving a total of 96 genes. To assess each of these candidates for effects on male fertility, we
757 induced knockdown in the male germline by crossing UAS-RNAi flies to Bam-GAL4, UAS-
758 Dicer2 flies (White-Cooper 2012; Gubala et al. 2017). Control flies were generated by crossing
759 the genetic background into which UAS-RNAi was inserted crossed to the same GAL4 line.
760 Flies carrying UAS-RNAi were of two types. Roughly half of the genes had publicly available
761 lines from the Vienna *Drosophila* Resource Center (Dietzl et al. 2007) or the Transgenic RNAi
762 Project (Ni et al. 2011). For the other genes, no publicly available RNAi stock was available, so
763 we constructed TRiP-style stocks in the pValium20 vector as previously described (Findlay et al.
764 2014). These constructs were integrated into an AttP site in stock BL 25709 ($y^1 v^1 P\{\text{nos-}$
765 $\text{phiC31}\int\text{nt.NLS}\}X; P\{\text{CaryP}\}\text{attP40}$) from the Bloomington *Drosophila* Stock Center (injections
766 by Genetivision; Houston, TX, USA) and crossed into a $y v$ background to screen for v^+ . We
767 attempted at least two rounds of transgenic production for each gene. In total, we were able to
768 obtain and test RNAi lines for 57 of the 96 identified genes.

769 We initially screened males knocked down for each candidate gene for major fertility
770 defects by crossing groups of 7 knockdown or control males to 5 virgin Canton S females,
771 letting the adults lay eggs for ~48 hours, and then discarding adults and quantifying the resulting
772 progeny by counting the pupal cases, as previously described (Gubala et al. 2017). To assess
773 the degree of knockdown achieved, 10 whole males of each line were homogenized in TRIzol
774 reagent (Life Technologies, Carlsbad, CA). RNA isolation, DNase treatment, cDNA synthesis
775 and semi-quantitative RT-PCR with gene-specific primers were performed as previously
776 described; amplification of *RpL32* was used as a positive control (Gubala et al. 2017). Any
777 gene that did not show at least partial knockdown was discarded from further analysis, leaving a
778 total of 42 genes successfully screened. Table S3 lists all lines used and results of tests for
779 effective target gene knockdown.

780

781 *CRISPR Genome Editing*

782

783 To validate RNAi results for *atlas*, *CG43072* and *CG33284*, we used CRISPR/Cas9
784 genome-editing to generate null alleles that could be used for further analysis, as described
785 previously (Lange et al. 2021). Briefly, our general strategy was to design gRNAs in the pU6.3
786 vector (*Drosophila* Genome Resource Center (DGRC) #1362) that targeted each end of a locus.
787 These plasmids, along with plasmids encoding gRNAs that targeted the w^+ locus, were co-
788 injected by Rainbow Transgenics (Camarillo, CA) into embryos laid by *vasa*-Cas9 females in a
789 w^+ background, Bloomington stock #51323 (Ge et al. 2016). G_0 animals were crossed to w^-
790 flies, and members of G_1 broods with a higher-than-expected fraction of w^- progeny were
791 individually crossed to an appropriate balancer line and then PCR-screened for the desired
792 deletion of the targeted locus.

793 We also constructed three frameshift, expected loss-of-function alleles for *atlas* by using
794 CRISPR to induce non-homologous end joining at a single PAM site just downstream of the
795 *atlas* start codon. *Vasa*-Cas9 embryos were co-injected and screened for w^- progeny as
796 described above. We then used squish preps to isolate DNA from G_1 flies and used a PCR-
797 RFLP assay to detect mutations. PCR products spanning the gRNA-targeted site were digested
798 with *Bfal* (New England Biolabs (NEB), Ipswich, MA); undigested products in which the

799 expected *Bfal* site was lost indicated a mutation, which was balanced and then confirmed by
800 PCR and sequencing of homozygous mutant lines.

801 We used scarless CRISPR editing and homology-directed repair (HDR) to insert the
802 GFP protein-coding sequence in-frame at the end of the *atlas* protein sequence (see Fig. S6;
803 (<https://flycrispr.org/scarless-gene-editing/>) (Bruckner et al. 2017; Bier et al. 2018; Hill et al.
804 2019). We first generated an *atlas*-GFP DNA construct by cloning the *atlas* protein-coding
805 sequence into pENTR and using LR Clonase II (Thermo Fisher Scientific, Waltham, MA) to
806 recombine the sequence with pTWG (DGRC #1076; T. Murphy), generating a C-terminally
807 tagged *atlas*:GFP construct. We amplified the *atlas* fragment from *vasa*-Cas9 strain #51323
808 genomic DNA. Once *atlas*-GFP was obtained in a plasmid, we amplified it with primers that
809 contained 5' *Esp3I* sites and overhangs designed for Golden Gate Assembly (GGA) and that, in
810 the case of the reverse primer, also added on 42 nucleotides downstream of the *atlas* stop
811 codon to reach a PAM site identified by FlyCRISPR TargetFinder (Gratz et al. 2014) as being
812 optimal for Cas9/gRNA recognition and cleavage. The primer also introduced a mutation in the
813 PAM site so that insertion of the designed piece of DNA into the genome *in vivo* would not be
814 subject to re-cutting. We also used the NEB Q5 Site-Directed Mutagenesis kit to introduce a
815 silent mutation into the *atlas* protein-coding sequence to eliminate an internal *Esp3I* site. The
816 resulting construct was used as the “left” homology arm for homology-directed repair (HDR)
817 editing. We constructed a “right” homology arm by using NEB Q5 PCR to amplify a 982-bp
818 fragment downstream of the PAM site, using primers modified to contain *Esp3I* sites and
819 overhangs compatible with GGA. We performed GGA by combining these left and right arms, a
820 plasmid containing a PiggyBac transposase-excisable 3xP3-dsRed flanked by *Esp3I* sites, and
821 backbone plasmid pXZ13, with *Esp3I* and T4 DNA ligase (NEB). A combination of colony PCR,
822 restriction digestion and sequencing identified properly assembled plasmids suitable for HDR.

823 *Vasa*-Cas9 embryos were co-injected with the assembled plasmid and a pU6.3 plasmid
824 encoding a gRNA targeting the region just downstream of the *atlas* stop codon. G0 flies were
825 crossed to *w*¹¹¹⁸ adults, and G1 flies were screened for red fluorescent eyes using the
826 NIGHTSEA system (NIGHTSEA LLC, Lexington, MA). Six balanced lines from two independent
827 G1 broods were established. To remove the dsRed from the *atlas* locus, we crossed these lines
828 to a PiggyBac transposase line (BDSC #8285) and then selected against pBac and dsRed in
829 the following generation. PCR and sequencing confirmed the expected “scarless” insertion of
830 GFP at the *atlas* locus.

831 832 *Atlas Rescue Line*

833
834 We constructed an HA-tagged *atlas* rescue line that contained the *atlas* gene flanked by
835 1345 bp of sequence upstream of the start codon (but excluding the coding sequence of
836 upstream gene CG3124) and 3000 bp of sequence downstream of the stop codon (including the
837 full 3' UTR) as follows. Genomic sequences were PCR amplified using Q5 High fidelity
838 Polymerase (NEB), purified Canton S genomic DNA (Gentra Puregene Tissue Kit, Qiagen,
839 Germantown, MD), and the *atlas* rescue F1/R1 and *atlas* rescue F3/R3 primer sets (see Table
840 S4). The 3x-HA tag was likewise amplified from pTWH (DGRC 1100; T. Murphy) using *atlas*
841 rescue F2/R2 primers. These DNA fragments were subsequently assembled into a XbaI/Ascl-
842 linearized w+attB plasmid (Addgene, Watertown, MA, plasmid 30326, deposited by J.
843 Sekelsky). The assembled construct was then phiC31 integrated into the PBac{y⁺-attP-
844 3B}VK00037 (Bloomington *Drosophila* Stock Center (BDSC) stock #24872) docking site
845 (Rainbow Transgenics) and crossed into the *atlas* null background to assess rescue.

846 847 *Fertility Assays and Sperm Visualization*

848

849 To validate the finding of reduced fertility for *atlas* knockdown males in the group fertility
850 assay described above, we performed single-pair fertility assays in which knockdown or mutant
851 males or their controls were mated individually to Canton S virgin females. Based on previous
852 experience analyzing genes that resulted in sterility or near-sterility (Gubala et al. 2017; Lange
853 et al. 2021), we designed assays with $N = 20$ -30 flies per male genotype. Matings were
854 observed, and males were discarded after copulation. Females were allowed to lay eggs into
855 the vials for 4 days and then discarded. Pupal cases were counted as a measure of fertility.
856 Crosses to generate and mating assays involving RNAi flies were maintained at 25° to optimize
857 knockdown. Before all assays, flies were reared to sexual maturity (3-7 days) in single-sex
858 groups on cornmeal-molasses food supplemented with dry yeast grains (Gubala et al. 2017).

859 To assess the ability of *atlas*-GFP to rescue the fertility defect of the $\Delta atlas$ allele, we
860 crossed *atlas*-GFP and *w1118* flies to $\Delta atlas/SM6$. Males with genotypes *atlas*-GFP/ $\Delta atlas$ and
861 $+\Delta atlas$ were compared using the single-pair fertility assay described above.

862 To observe the production of sperm in knockdown or mutant males, we introduced the
863 Mst35Bb-GFP marker into these males, which labels mature sperm and late-stage spermatid
864 nuclei with GFP (Manier et al. 2010). Samples were prepared, imaged and analyzed as
865 described previously (Gubala et al. 2017); because the large differences in testis shape and
866 sperm production observed in initial phase contrast imaging would be obvious to any
867 experimenter, we were not blind to male genotypes.

868 *Atlas-GFP Ectopic Expression*

870
871 We used the Gateway cloning system (Thermo) to construct an *atlas*-GFP transgene
872 expressed under UAS control (primers in Table S4). The *atlas* protein-coding sequence in
873 pENTR was recombined with pTWG (*Drosophila* Genomics Resource Center, T. Murphy) as
874 described above. The resulting plasmid was then inserted into *w*-flies using P-element-
875 mediated transposition (Rainbow Transgenics), *w*+ G1s were selected, and several
876 independent insertions were balanced. We crossed male UAS-*atlas*:GFP flies to females from
877 two different driver lines: *tubulin*-GAL4 (to drive ubiquitous expression) and *Bam*-GAL4 (to drive
878 expression in the early germline). We dissected larval salivary glands of the *tub>atlas*:GFP
879 males, since these cells are exceptionally large and ideal for visualizing subcellular localization.
880 We then dissected the testes of *Bam>atlas*:GFP males to evaluate whether the localization
881 pattern observed in the salivary gland was consistent in testis tissue, albeit not the same cells in
882 which endogenous *atlas* appears to be expressed. Protein localization was visualized by
883 fluorescence confocal microscopy on a Leica SP5 microscope (Leica Microsystems, Wetzlar,
884 Germany) and images were captured with LASAF as described previously (Lange et al. 2021).

885 *Imaging Spermatogenesis and Spermatid Nuclear Condensation*

887
888 We used phase-contrast microscopy to examine the stages of spermatogenesis in whole
889 mount testes (White-Cooper 2004). To assess the processes of nuclear condensation and
890 individualization of 64-cell cysts of spermatids in the post-meiotic stages of spermatogenesis,
891 we used fluorescence and confocal microscopy to visualize actin-based individualization
892 complexes and nuclei. Samples were processed, and actin and nuclear DNA were visualized
893 with TRITC-phalloidin (Molecular Probes, Eugene, OR) and TOPRO-3 iodine (Thermo),
894 respectively, as described previously (Lange et al. 2021). The final stages of nuclear
895 condensation were visualized with the Mst35Bb-GFP marker described above, as well as an
896 equivalent marker, Mst35Bb-dsRed (Manier et al. 2010). Earlier nuclear stages were visualized
897 with histone H2AvD-RFP (BDSC stock #23651), which is present in round spermatid nuclei and
898 the earliest stages of nuclear elongation (Clarkson and Saint 1999; Rathke et al. 2007). Images

899 with H2AvD-RFP were obtained with epifluorescence microscopy, since we lack an appropriate
900 confocal laser for RFP.

901 To examine spermatid nuclei at various stages of condensation, we visualized nuclear
902 bundles using TOPRO-3. Testes of newly eclosed (<1 day old) *atlas* null and control males
903 were dissected in PBS. Testes were then transferred to a droplet of 4% paraformaldehyde on
904 poly-L-lysine treated glass slides and were gently shredded in the post-meiotic region to release
905 sperm bundles. Testes were gently squashed beneath coverslips coated in Sigmacote (Sigma
906 Aldrich, St. Louis, MO). We then froze slides in liquid nitrogen for a few seconds and popped off
907 of the siliconized coverslip with a razor. Slides were incubated in Coplin jars filled with 95%
908 ethanol at -20°C for 30 minutes and then mounted in VECTASHIELD (Vector Laboratories,
909 Burlingame, CA). Nuclear staging was performed by examining the shape of the nuclei. Early
910 and late canoe stages of condensation were distinguished by the absence or presence of
911 Mst35Bb-GFP, respectively. Elongated and condensed stages were distinguished by the
912 presence or absence, respectively, of vesicles of GFP-tagged nuclear proteins (Atlas-GFP or
913 Mst35Bb-GFP) located basal to the nuclei. Examples of stages are given in Fig. 4. Confocal
914 stacks were taken on a Leica SP5 microscope, images were captured by LASAF, and ImageJ
915 was used to flatten stacks into a single, two-dimensional image. All intact nuclear bundles were
916 counted for each dissection.

917 For the experiments measuring nuclear condensation stage (Table S1), a sample size of
918 $N = 10$ for each genotype was selected based on the magnitude of the *atlas* null phenotype and
919 the consistent differences observed in previous dissections of these genotypes with Mst35Bb-
920 GFP. Likewise, for the IC-nuclear bundle association and IC progression analysis (Fig. 3C-D),
921 we selected sample sizes of $N = \sim 15$ per genotype based on pilot experiments showing that
922 aberrant actin phenotypes were highly consistent in null testes and previous experience with
923 such quantification (Lange et al. 2021).

924 *Evolutionary and Gene Expression Analysis of atlas*

925 We searched for orthologs of the *D. melanogaster* Atlas protein in the original
926 sequenced *Drosophila* species with BLASTP searches in FlyBase (Clark et al. 2007). We also
927 used TBLASTN searches to identify orthologs in species lacking complete protein annotations.
928 We identified syntenic regions for each species by looking for conserved neighboring genes,
929 such as *ord* and *CG3124*. In addition to analyzing the *atlas* coding region, we conducted
930 separate BLASTN searches for the sequence of the *D. melanogaster* 3'UTR across *Drosophila*
931 species since it has a different conservation pattern than the coding sequence.

932 To test for sex-specific expression biases for both the ORF and the 3' UTR sequences,
933 we used adult male- and female-specific RNA-seq data from numerous *Drosophila* species
934 accessed through the Genomics Education Partnership version of the UCSC Genome Browser
935 (<http://gander.wustl.edu/>) and initially collected by Brown et al. (2014) and Chen et al. (2014).
936 We also confirmed these findings experimentally in several species by performing RT-PCR on
937 cDNA isolated from whole males and whole females, as previously described (Gubala et al.
938 2017).

939 To search for *atlas* orthologs in non-*Drosophila* Dipterans, we obtained from ENSEMBL
940 Metazoa the genomes of *Musca domestica*, *Glossina morsitans*, *Lucilia cuprina*, *Aedes aegypti*,
941 *Anopheles darlingi*, *Anopheles gambiae*, *Culex quinquefasciatus* and *Mayetiola destructor*. We
942 performed a synteny search (summarized in Fig. S10) in each species by identifying the nearest
943 neighbors of *atlas* in the *D. ananassae* and *D. virilis* genomes that had an identifiable homolog
944 in each species. In all cases, the homologs of the nearest neighbors on each side of *atlas* were
945 found on different contigs, suggesting synteny breakdown. We obtained up to 1 Mb of
946 sequence on each side of each identified homolog and queried it with BLASTN, TBLASTN, and
947 Exonerate (Slater and Birney 2005) for regions with significantly similarity to any portion of the

950 Atlas protein or cDNA sequences. No significant hits, and no hits better than what could be
951 found in other parts of the genome, were found. Finally, we used HMMER to search for
952 orthologs in all annotated proteomes and all metazoan proteomes. We first queried the
953 database with Atlas from either *D. melanogaster* or *D. virilis* and accepted hits that fell below an
954 e-value cutoff of 0.01 and a minimum hit length of 3%. These hits were then included iteratively
955 in subsequent searches until no new significant hits were found.

956 We analyzed the molecular evolution of the *atlas* protein-coding sequence by obtaining
957 orthologous protein-coding sequences from *melanogaster* group species. (Analysis out of this
958 group was not performed due to high sequence divergence and poor alignment quality.) We
959 used BLASTP to identify these sequences from GenBank and then extracted the coding DNA
960 sequence for each. Sequences were aligned, checked for recombination, used to construct a
961 gene tree, and analyzed with the PAML sites test as described previously (McGeary and Findlay
962 2020), except that alignment positions that included gaps were masked from the PAML
963 analysis. We initially analyzed a set of 13 species (*melanogaster*, *simulans*, *sechellia*, *yakuba*,
964 *erecta*, *suzukii*, *takahashii*, *biarmipes*, *hopaloea*, *ficusphila*, *elegans*, *eugracilis* and *ananassae*);
965 we excluded an ortholog detected in *D. bipunctinata* due to poor alignment. This initial analysis
966 detected a class of sites with significant evidence of positive selection, but closer inspection of
967 the alignment revealed that the site with the strongest evidence of selection, corresponding to
968 *D. melanogaster* residue 31R, may have been driven by a questionable alignment due to an
969 insertion in that region that was unique to *D. takahashii*. The results reported derive from an
970 analysis with *D. takahashii* excluded, which produced a more reliable alignment and showed no
971 evidence for any sites under positive selection.
972

973 **Acknowledgements**

974

975 We thank Dr. Alexis Hill for assistance with developing the *atlas*-GFP transgene; Gynesis
976 Vance, Elvis Perez, Ishanpepe Jagusah, Emily Gualdino, and students in the Biology 261 Lab at
977 Holy Cross for assistance with RNAi fertility screens; Dr. Rob Bellin for assistance with confocal
978 microscopy; Dr. Justin McAlister for use of the NIGHTSEA system; the Bloomington Stock
979 Center and the Vienna *Drosophila* Resource Center for fly strains; and Dr. Mariana Wolfner and
980 members of the Findlay and Bornberg-Bauer labs for helpful discussions about the project and
981 the manuscript. This work was supported by NSF CAREER Award #1652013 (to GDF) and a
982 Humboldt Fellowship (to AG). GDF and EBB are also grateful to the Evolution Think Tank at
983 the University of Münster, which funded GDF as a visiting fellow for in-person collaboration.
984

References

- 985
986
987 Abrusán G. 2013. Integration of New Genes into Cellular Networks, and Their Structural
988 Maturation. *Genetics* 195:1407-1417.
- 989 Alvi ZA, Chu T-C, Schawaroch V, Klaus AV. 2016. Genomic and expression analysis of
990 transition proteins in *Drosophila*. *Spermatogenesis* 5:e1178518.
- 991 Alvi ZA, Chu T-C, Schawaroch V, Klaus AV. 2013. Protamine-like proteins in 12 sequenced
992 species of *Drosophila*. *Protein and peptide letters* 20:17-35.
- 993 Argyridou E, Huylmans AK, Königer A, Parsch J. 2017. X-linkage is not a general inhibitor of
994 tissue-specific gene expression in *Drosophila melanogaster*. *Heredity* 119:27-34.
- 995 Argyridou E, Parsch J. 2018. Regulation of the X Chromosome in the Germline and Soma of
996 *Drosophila melanogaster* Males. *Genes* 9:242.
- 997 Baalsrud HT, Tørresen OK, Solbakken MH, Salzburger W, Hanel R, Jakobsen KS, Jentoft S.
998 2018. *De novo* gene evolution of antifreeze glycoproteins in codfishes revealed by whole
999 genome sequence data. *Molecular Biology and Evolution* 35:593-606.
- 1000 Baker RH, Narechania A, Johns PM, Wilkinson GS. 2012. Gene duplication, tissue-specific
1001 gene expression and sexual conflict in stalk-eyed flies (Diopsidae). *Philosophical*
1002 *Transactions of the Royal Society of London. Series B, Biological Sciences* 367:2357-
1003 2375.
- 1004 Begun DJ, Lindfors HA, Kern AD, Jones CD. 2007. Evidence for *de Novo* Evolution of Testis-
1005 Expressed Genes in the *Drosophila yakuba/Drosophila erecta* Clade. *Genetics*
1006 176:1131-1137.
- 1007 Bier E, Harrison MM, O'Connor-Giles KM, Wildonger J. 2018. Advances in Engineering the Fly
1008 Genome with the CRISPR-Cas System. *Genetics* 208:1-18.
- 1009 Bornberg-Bauer E, Schmitz JF. 2017. Fact or fiction: Updates on how protein-coding genes
1010 might emerge *de novo* from previously non-coding DNA. *F1000Research* 6:57.
- 1011 Bouska MJ, Bai H. 2021. Long noncoding RNA regulation of spermatogenesis via the spectrin
1012 cytoskeleton in *Drosophila*. *G3 Genes|Genomes|Genetics*.
- 1013 Brown JB, Boley N, Eisman R, May GE, Stoiber MH, Duff MO, Booth BW, Wen J, Park S,
1014 Suzuki AM, et al. 2014. Diversity and dynamics of the *Drosophila* transcriptome. *Nature*
1015 512:393-399.
- 1016 Bruckner JJ, Zhan H, Gratz SJ, Rao M, Ukken F, Zilberg G, O'Connor-Giles KM. 2017. Fife
1017 organizes synaptic vesicles and calcium channels for high-probability neurotransmitter
1018 release. *Journal of Cell Biology* 216:231-246.
- 1019 Cai J, Zhao R, Jiang H, Wang W. 2008. *De novo* origination of a new protein-coding gene in
1020 *Saccharomyces cerevisiae*. *Genetics* 179:487-496.
- 1021 Carvunis A-R, Rolland T, Wapinski I, Calderwood MA, Yildirim MA, Simonis N, Charlotiaux B,
1022 Hidalgo C, Barbette J, Santhanam B, et al. 2012. Proto-genes and *de novo* gene birth.
1023 *Nature* 487:370-374.
- 1024 Chamakura KR, Tran JS, O'Leary C, Lisciandro HG, Antillon SF, Garza KD, Tran E, Min L,
1025 Young R. 2020. Rapid *de novo* evolution of lysis genes in single-stranded RNA phages.
1026 *Nature Communications* 11.
- 1027 Chen Z-X, Sturgill D, Qu J, Jiang H, Park S, Boley N, Suzuki AM, Fletcher AR, Plachetzki DC,
1028 FitzGerald PC, et al. 2014. Comparative validation of the *D. melanogaster* modENCODE
1029 transcriptome annotation. *Genome Research* 24:1209-1223.
- 1030 Civetta A, Finn S. 2014. Do candidate genes mediating conspecific sperm precedence affect
1031 sperm competitive ability within species? A test case in *Drosophila*. *G3 (Bethesda, Md.)*
1032 4:1701-1707.
- 1033 Clark AG, Eisen MB, Smith DR, Bergman CM, Oliver B, Markow TA, Kaufman TC, Kellis M,
1034 Gelbart W, Iyer VN, et al. 2007. Evolution of genes and genomes on the *Drosophila*
1035 phylogeny. *Nature* 450:203-218.

- 1036 Clarkson M, Saint R. 1999. A His2AvDGFP fusion gene complements a lethal His2AvD mutant
1037 allele and provides an in vivo marker for *Drosophila* chromosome behavior. *DNA and*
1038 *Cell Biology* 18:457-462.
- 1039 Cui X, Lv Y, Chen M, Nikoloski Z, Twell D, Zhang D. 2015. Young genes out of the male: an
1040 insight from evolutionary age analysis of the pollen transcriptome. *Molecular Plant*
1041 8:935-945.
- 1042 Demarco RS, Eikenes TH, Haglund K, Jones DL. 2014. Investigating spermatogenesis in
1043 *Drosophila melanogaster*. *Methods* 68:218-227.
- 1044 Dietzl G, Chen D, Schnorrer F, Su K-C, Barinova Y, Fellner M, Gasser B, Kinsey K, Oettel S,
1045 Scheiblauer S, et al. 2007. A genome-wide transgenic RNAi library for conditional gene
1046 inactivation in *Drosophila*. *Nature* 448:151-156.
- 1047 Dorus S, Busby SA, Gerike U, Shabanowitz J, Hunt DF, Karr TL. 2006. Genomic and functional
1048 evolution of the *Drosophila melanogaster* sperm proteome. *Nature Genetics* 38:1440-
1049 1445.
- 1050 Dorus S, Freeman ZN, Parker ER, Heath BD, Karr TL. 2008. Recent Origins of Sperm Genes in
1051 *Drosophila*. *Molecular Biology and Evolution* 25:2157-2166.
- 1052 Doyen CM, Chalkley GE, Voets O, Bezstarosti K, Demmers JA, Moshkin YM, Verrijzer CP.
1053 2015. A Testis-Specific Chaperone and the Chromatin Remodeler ISWI Mediate
1054 Repackaging of the Paternal Genome. *Cell Reports* 13:1310-1318.
- 1055 Emerson JJ, Kaessmann H, Betrán E, Long M. 2004. Extensive gene traffic on the mammalian
1056 X chromosome. *Science* 303:537-540.
- 1057 Eren-Ghiani Z, Rathke C, Theofel I, Renkawitz-Pohl R. 2015. Prtl99C Acts Together with
1058 Protamines and Safeguards Male Fertility in *Drosophila*. *Cell Reports* 13:2327-2335.
- 1059 Fabian L, Brill JA. 2012. *Drosophila* spermiogenesis. *Spermatogenesis* 2:197-212.
- 1060 Findlay GD, Sitnik JL, Wang W, Aquadro CF, Clark NL, Wolfner MF. 2014. Evolutionary rate
1061 covariation identifies new members of a protein network required for *Drosophila*
1062 *melanogaster* female post-mating responses. *PLoS Genetics* 10:e1004108.
- 1063 Gao G, Vibranovski MD, Zhang L, Li Z, Liu M, Zhang YE, Li X, Zhang W, Fan Q, Vankuren NW,
1064 et al. 2014. A long-term demasculinization of X-linked intergenic noncoding RNAs in
1065 *Drosophila melanogaster*. *Genome Research* 24:629-638.
- 1066 Gärtner SMK, Hundertmark T, Nolte H, Theofel I, Eren-Ghiani Z, Tetzner C, Duchow TB,
1067 Rathke C, Krüger M, Renkawitz-Pohl R. 2019. Stage-specific testes proteomics of
1068 *Drosophila melanogaster* identifies essential proteins for male fertility. *European Journal*
1069 *of Cell Biology* 98:103-115.
- 1070 Gärtner SMK, Rothenbusch S, Buxa MK, Theofel I, Renkawitz R, Rathke C, Renkawitz-Pohl R.
1071 2015. The HMG-box-containing proteins tHMG-1 and tHMG-2 interact during the
1072 histone-to-protamine transition in *Drosophila* spermatogenesis. *European Journal of Cell*
1073 *Biology* 94:46-59.
- 1074 Ge DT, Tipping C, Brodsky MH, Zamore PD. 2016. Rapid screening for CRISPR-directed
1075 editing of the *Drosophila* genome using *white* coconversion. *G3 (Bethesda, Md.)* 6:3197-
1076 3206.
- 1077 Gratz SJ, Ukken FP, Rubinstein CD, Thiede G, Donohue LK, Cummings AM, O'Connor-Giles
1078 KM. 2014. Highly specific and efficient CRISPR/Cas9-catalyzed homology-directed
1079 repair in *Drosophila*. *Genetics* 196:961-971.
- 1080 Gubala AM, Schmitz JF, Kearns MJ, Vinh TT, Bornberg-Bauer E, Wolfner MF, Findlay GD.
1081 2017. The *goddard* and *saturn* genes are essential for *Drosophila* male fertility and may
1082 have arisen *de novo*. *Molecular Biology and Evolution* 34:1066-1082.
- 1083 Guerzoni D, McLysaght A. 2016. *De novo* genes arise at a slow but steady rate along the
1084 primate lineage and have been subject to incomplete lineage sorting. *Genome biology*
1085 and evolution 8:1222-1232.

- 1086 Hales KG, Korey CA, Larracuente AM, Roberts DM. 2015. Genetics on the fly: a primer on the
1087 *Drosophila* model system. *Genetics* 201:815-842.
- 1088 Heames B, Schmitz J, Bornberg-Bauer E. 2020. A continuum of evolving de novo genes drives
1089 protein-coding novelty in *Drosophila*. *Journal of Molecular Evolution* 88:382-398.
- 1090 Hill A, Jain P, Ben-Shahar Y. 2019. The *Drosophila* ERG channel seizure plays a role in the
1091 neuronal homeostatic stress response. *PLoS Genetics* 15:e1008288.
- 1092 Hundertmark T, Gärtner SMK, Rathke C, Renkawitz-Pohl R. 2018. Nejire/dCBP-mediated
1093 histone H3 acetylation during spermatogenesis is essential for male fertility in *Drosophila*
1094 *melanogaster*. *PLOS ONE* 13:e0203622.
- 1095 Jayaramaiah Raja S, Renkawitz-Pohl R. 2005. Replacement by *Drosophila melanogaster*
1096 protamines and Mst77F of histones during chromatin condensation in late spermatids
1097 and role of Sesame in the removal of these proteins from the male pronucleus.
1098 *Molecular and Cellular Biology* 25:6165-6177.
- 1099 Kaessmann H. 2010. Origins, evolution, and phenotypic impact of new genes. *Genome*
1100 *Research* 20:1313-1326.
- 1101 Kanippayoor RL, Alpern JHM, Moehring AJ. 2013. Protamines and spermatogenesis in
1102 *Drosophila* and *Homo sapiens*. *Spermatogenesis* 3:e24376.
- 1103 Keeling DM, Garza P, Nartey CM, Carvunis A-R. 2019. The meanings of 'function' in biology
1104 and the problematic case of de novo gene emergence. *eLife* 8:e47014.
- 1105 Kemkemer C, Catalán A, Parsch J. 2014. 'Escaping' the X chromosome leads to increased
1106 gene expression in the male germline of *Drosophila melanogaster*. *Heredity* 112:149-
1107 155.
- 1108 Kimura S, Loppin B. 2016. The *Drosophila* chromosomal protein Mst77F is processed to
1109 generate an essential component of mature sperm chromatin. *Open Biology* 6:160207.
- 1110 Kondo S, Vedanayagam J, Mohammed J, Eizadshenass S, Kan L, Pang N, Aradhya R, Siepel
1111 A, Steinhauer J, Lai EC. 2017. New genes often acquire male-specific functions but
1112 rarely become essential in *Drosophila*. *Genes & development* 31:1841-1846.
- 1113 Landeen EL, Muirhead CA, Wright L, Meiklejohn CD, Presgraves DC. 2016. Sex Chromosome-
1114 wide Transcriptional Suppression and Compensatory *Cis*-Regulatory Evolution Mediate
1115 Gene Expression in the *Drosophila* Male Germline. *PLOS Biology* 14:e1002499.
- 1116 Lange A, Mills RE, Lange CJ, Stewart M, Devine SE, Corbett AH. 2007. Classical nuclear
1117 localization signals: definition, function, and interaction with Importin α . *Journal of*
1118 *Biological Chemistry* 282:5101-5105.
- 1119 Lange A, Patel PH, Heames B, Damry AM, Saenger T, Jackson CJ, Findlay GD, Bornberg-
1120 Bauer E. 2021. Structural and functional characterization of a putative de novo gene in
1121 *Drosophila*. *Nature Communications* 12:1667.
- 1122 Larkin A, Marygold SJ, Antonazzo G, Attrill H, Gilberto, Garapati PV, Joshua, Millburn G,
1123 Strelets VB, Tabone CJ, et al. 2021. FlyBase: updates to the *Drosophila melanogaster*
1124 knowledge base. *Nucleic acids research* 49:D899-D907.
- 1125 Levine MT, Jones CD, Kern AD, Lindfors HA, Begun DJ. 2006. Novel genes derived from
1126 noncoding DNA in *Drosophila melanogaster* are frequently X-linked and exhibit testis-
1127 biased expression. *Proceedings of the National Academy of Sciences* 103:9935-9939.
- 1128 Li D, Yan Z, Lu L, Jiang H, Wang W. 2014. Pleiotropy of the *de novo*-originated gene *MDF1*.
1129 *Scientific reports* 4:7280.
- 1130 Li Z-W, Chen X, Wu Q, Hagmann J, Han T-S, Zou Y-P, Ge S, Guo Y-L. 2016. On the origin of
1131 *de novo* genes in *Arabidopsis thaliana* populations. *Genome biology and evolution*
1132 8:2190-2202.
- 1133 Lim C, Tarayrah L, Chen X. 2012. Transcriptional regulation during *Drosophila*
1134 spermatogenesis. *Spermatogenesis* 2:158-166.

- 1135 Lipinski KJ, Farslow JC, Fitzpatrick KA, Lynch M, Katju V, Bergthorsson U. 2011. High
1136 spontaneous rate of gene duplication in *Caenorhabditis elegans*. *Current biology : CB*
1137 21:306-310.
- 1138 Lu T-C, Leu J-Y, Lin W-C. 2017. A comprehensive analysis of transcript-supported de novo
1139 genes in *Saccharomyces sensu stricto* yeasts. *Molecular Biology and Evolution* 34:2823-
1140 2838.
- 1141 Mahadevaraju S, Fear JM, Akeju M, Galletta BJ, Pinheiro MMLS, Avelino CC, Cabral-De-Mello
1142 DC, Conlon K, Dell'Orso S, Demere Z, et al. 2021. Dynamic sex chromosome
1143 expression in *Drosophila* male germ cells. *Nature Communications* 12.
- 1144 Manier MK, Belote JM, Berben KS, Novikov D, Stuart WT, Pitnick S. 2010. Resolving
1145 mechanisms of competitive fertilization success in *Drosophila melanogaster*. *Science*
1146 328:354-357.
- 1147 Marques AC, Dupanloup I, Vinckenbosch N, Reymond A, Kaessmann H. 2005. Emergence of
1148 young human genes after a burst of retroposition in primates. *PLOS Biology* 3:e357.
- 1149 McGeary MK, Findlay GD. 2020. Molecular evolution of the sex peptide network in *Drosophila*.
1150 *Journal of Evolutionary Biology* 33:629-641.
- 1151 McLysaght A, Hurst LD. 2016. Open questions in the study of *de novo* genes: what, how and
1152 why. *Nature Reviews Genetics* 17:567.
- 1153 Meiklejohn CD, Landeen EL, Cook JM, Kingan SB, Presgraves DC. 2011. Sex Chromosome-
1154 Specific Regulation in the *Drosophila* Male Germline But Little Evidence for
1155 Chromosomal Dosage Compensation or Meiotic Inactivation. *PLOS Biology* 9:e1001126.
- 1156 Mikhaylova LM, Nurminsky DI. 2011. Lack of global meiotic sex chromosome inactivation, and
1157 paucity of tissue-specific gene expression on the *Drosophila* X chromosome. *BMC*
1158 *Biology* 9:29.
- 1159 Moyers BA, Zhang J. 2015. Phylostratigraphic bias creates spurious patterns of genome
1160 evolution. *Molecular Biology and Evolution* 32:258-267.
- 1161 Moyers BA, Zhang J. 2018. Toward reducing phylostratigraphic errors and biases. *Genome*
1162 *biology and evolution* 10:2037-2048.
- 1163 Necsulea A, Kaessmann H. 2014. Evolutionary dynamics of coding and non-coding
1164 transcriptomes. *Nature Reviews Genetics* 15:734-748.
- 1165 Ni J-Q, Zhou R, Czech B, Liu L-P, Holderbaum L, Yang-Zhou D, Shim H-S, Tao R, Handler D,
1166 Karpowicz P, et al. 2011. A genome-scale shRNA resource for transgenic RNAi in
1167 *Drosophila*. *Nature Methods* 8:405-407.
- 1168 Palmieri N, Kosiol C, Schlötterer C. 2014. The life cycle of *Drosophila* orphan genes. *eLife*
1169 3:e01311.
- 1170 Parisi M, Nuttall R, Edwards P, Minor J, Naiman D, Lü J, Doctolero M, Vainer M, Chan C,
1171 Malley J, et al. 2004. A survey of ovary-, testis-, and soma-biased gene expression in
1172 *Drosophila melanogaster* adults. *Genome biology* 5:R40.
- 1173 Parisi M, Nuttall R, Naiman D, Bouffard G, Malley J, Andrews J, Eastman S, Oliver B. 2003.
1174 Paucity of genes on the *Drosophila* X chromosome showing male-biased expression.
1175 *Science* 299:697-700.
- 1176 Potter SC, Luciani A, Eddy SR, Park Y, Lopez R, Finn RD. 2018. HMMER web server: 2018
1177 update. *Nucleic acids research* 46:W200-W204.
- 1178 Puntambekar S, Newhouse R, Navas JSM, Chauhan R, Vernaz G, Willis T, Wayland MT,
1179 Umrانيا Y, Miska EA, Prabakaran S. 2020. Evolutionary divergence of novel open
1180 reading frames in cichlids speciation. *Scientific reports* 10.
- 1181 Rathke C, Baarends WM, Awe S, Renkawitz-Pohl R. 2014. Chromatin dynamics during
1182 spermiogenesis. In: *Chromatin and epigenetic regulation of animal development*. p. 155-
1183 168.
- 1184 Rathke C, Baarends WM, Jayaramaiah-Raja S, Bartkuhn M, Renkawitz R, Renkawitz-Pohl R.
1185 2007. Transition from a nucleosome-based to a protamine-based chromatin

- 1186 configuration during spermiogenesis in *Drosophila*. Journal of cell science 120:1689-
1187 1700.
- 1188 Rathke C, Barckmann B, Burkhard S, Jayaramaiah-Raja S, Roote J, Renkawitz-Pohl R. 2010.
1189 Distinct functions of Mst77F and protamines in nuclear shaping and chromatin
1190 condensation during *Drosophila* spermiogenesis. European Journal of Cell Biology
1191 89:326-338.
- 1192 Ravi Ram K, Wolfner MF. 2007. Sustained post-mating response in *Drosophila melanogaster*
1193 requires multiple seminal fluid proteins. PLoS Genetics 3:e238.
- 1194 Reinhardt JA, Wanjiru BM, Brant AT, Saelao P, Begun DJ, Jones CD. 2013. *De novo* ORFs in
1195 *Drosophila* are important to organismal fitness and evolved rapidly from previously non-
1196 coding sequences. PLoS Genetics 9:e1003860.
- 1197 Reinke V, San Gil I, Ward S, Kazmer K. 2004. Genome-wide germline-enriched and sex-biased
1198 expression profiles in *Caenorhabditis elegans*. Development 131:311-323.
- 1199 Rele CP, Sandlin KM, Leung W, Reed LK. 2020. Manual Annotation of Genes within *Drosophila*
1200 Species: the Genomics Education Partnership protocol. bioRxiv 12.10.420521.
- 1201 Ruiz-Orera J, Hernandez-Rodriguez J, Chiva C, Sabidó E, Kondova I, Bontrop R, Marqués-
1202 Bonet T, Albà MM. 2016. Origins of *de novo* genes in human and chimpanzee. PLoS
1203 Genetics 11:e1005721.
- 1204 Schaeffer SW, Bhutkar A, McAllister BF, Matsuda M, Matzkin LM, O'Grady PM, Rohde C,
1205 Valente VLS, Aguadé M, Anderson WW, et al. 2008. Polytene Chromosomal Maps of 11
1206 *Drosophila* Species: The Order of Genomic Scaffolds Inferred From Genetic and
1207 Physical Maps. Genetics 179:1601-1655.
- 1208 Schlötterer C. 2015. Genes from scratch – the evolutionary fate of *de novo* genes. Trends in
1209 Genetics 31:215-219.
- 1210 Schmitz JF, Ullrich KK, Bornberg-Bauer E. 2018. Incipient *de novo* genes can evolve from
1211 frozen accidents that escaped rapid transcript turnover. Nature Ecology & Evolution
1212 2:1626-1632.
- 1213 Schnakenberg SL, Matias WR, Siegal ML. 2011. Sperm-storage defects and live birth in
1214 *Drosophila* females lacking spermathecal secretory cells. PLOS Biology 9:e1001192.
- 1215 Schuh M, Lehner CF, Heidmann S. 2007. Incorporation of *Drosophila* CID/CENP-A and CENP-
1216 C into Centromeres during Early Embryonic Anaphase. Current Biology 17:237-243.
- 1217 Slater G, Birney E. 2005. Automated generation of heuristics for biological sequence
1218 comparison. BMC Bioinformatics 6:31.
- 1219 Sorourian M, Kunte MM, Domingues S, Gallach M, Özdil F, Río J, Betrán E. 2014. Relocation
1220 facilitates the acquisition of short *cis*-regulatory regions that drive the expression of
1221 retrogenes during spermatogenesis in *Drosophila*. Molecular Biology and Evolution
1222 31:2170-2180.
- 1223 Soumillon M, Necsulea A, Weier M, Brawand D, Zhang X, Gu H, Barthès P, Kokkinaki M, Nef S,
1224 Gnirke A, et al. 2013. Cellular source and mechanisms of high transcriptome complexity
1225 in the mammalian testis. Cell Reports 3:2179-2190.
- 1226 Steinhauer J. 2015. Separating from the pack: molecular mechanisms of *Drosophila* spermatid
1227 individualization. Spermatogenesis 5:e1041345.
- 1228 Suenaga Y, Islam SMR, Alagu J, Kaneko Y, Kato M, Tanaka Y, Kawana H, Hossain S,
1229 Matsumoto D, Yamamoto M, et al. 2014. NCYM, a *cis*-antisense gene of MYCN,
1230 encodes a *de novo* evolved protein that inhibits GSK3 β resulting in the stabilization of
1231 MYCN in human neuroblastomas. PLoS Genetics 10:e1003996.
- 1232 Tirmarche S, Kimura S, Sapey-Triomphe L, Sullivan W, Landmann F, Loppin B. 2014.
1233 *Drosophila* Protamine-Like Mst35Ba and Mst35Bb Are Required for Proper Sperm
1234 Nuclear Morphology but Are Dispensable for Male Fertility. G3
1235 Genes|Genomes|Genetics 4:2241-2245.

- 1236 Vakirlis N, Acar O, Hsu B, Castilho Coelho N, Van Oss SB, Wacholder A, Medetgul-Ernar K,
1237 Bowman RW, Hines CP, Iannotta J, et al. 2020. De novo emergence of adaptive
1238 membrane proteins from thymine-rich genomic sequences. *Nature Communications* 11.
1239 Vakirlis N, Hebert AS, Opulente DA, Achaz G, Hittinger CT, Fischer G, Coon JJ, Lafontaine I.
1240 2018. A molecular portrait of de novo genes in yeasts. *Molecular Biology and Evolution*
1241 35:631-645.
- 1242 Van Oss SB, Carvunis A-R. 2019. *De novo* gene birth. *PLoS Genetics* 15:e1008160.
- 1243 VanKuren NW, Long M. 2018. Gene duplicates resolving sexual conflict rapidly evolved
1244 essential gametogenesis functions. *Nature Ecology & Evolution* 2:705-712.
- 1245 Vibranovski MD. 2014. Meiotic Sex Chromosome Inactivation in *Drosophila*. *Journal of*
1246 *Genomics* 2:104-117.
- 1247 Vibranovski MD, Lopes HF, Karr TL, Long M. 2009. Stage-specific expression profiling of
1248 *Drosophila* spermatogenesis suggests that meiotic sex chromosome inactivation drives
1249 genomic relocation of testis-expressed genes. *PLoS Genetics* 5:e1000731.
- 1250 Vibranovski MD, Zhang Y, Long M. 2009. General gene movement off the X chromosome in the
1251 *Drosophila* genus. *Genome Research* 19:897-903.
- 1252 Vibranovski MD, Zhang YE, Kemkemer C, Lopes HF, Karr TL, Long M. 2012. Re-analysis of the
1253 larval testis data on meiotic sex chromosome inactivation revealed evidence for tissue-
1254 specific gene expression related to the *Drosophila* X chromosome. *BMC Biology* 10:49.
- 1255 Wasbrough ER, Dorus S, Hester S, Howard-Murkin J, Lilley K, Wilkin E, Polpitiya A, Petritis K,
1256 Karr TL. 2010. The *Drosophila melanogaster* sperm proteome-II (DmSP-II). *Journal of*
1257 *Proteomics* 73:2171-2185.
- 1258 Weisman CM, Murray AW, Eddy SR. 2020. Many, but not all, lineage-specific genes can be
1259 explained by homology detection failure. *PLOS Biology* 18:e3000862.
- 1260 Wen K, Yang L, Xiong T, Di C, Ma D, Wu M, Xue Z, Zhang X, Long L, Zhang W, et al. 2016.
1261 Critical roles of long noncoding RNAs in *Drosophila* Spermatogenesis. *Genome*
1262 *Research* 26:1233-1244.
- 1263 White-Cooper H. 2010. Molecular mechanisms of gene regulation during *Drosophila*
1264 spermatogenesis. *REPRODUCTION* 139:11-21.
- 1265 White-Cooper H. 2004. Spermatogenesis: analysis of meiosis and morphogenesis. In:
1266 Henderson D, editor. *Drosophila Cytogenetics Protocols*. Totowa, NJ: Humana Press. p.
1267 45-75.
- 1268 White-Cooper H. 2012. Tissue, cell type and stage-specific ectopic gene expression and RNAi
1269 induction in the *Drosophila* testis. *Spermatogenesis* 2:11-22.
- 1270 Wilburn DB, Swanson WJ. 2016. From molecules to mating: Rapid evolution and biochemical
1271 studies of reproductive proteins. *Journal of Proteomics* 135:12-25.
- 1272 Wilson BA, Foy SG, Neme R, Masel J. 2017. Young genes are highly disordered as predicted
1273 by the preadaptation hypothesis of *de novo* gene birth. *Nature Ecology & Evolution*
1274 1:0146.
- 1275 Witt E, Benjamin S, Svetec N, Zhao L. 2019. Testis single-cell RNA-seq reveals the dynamics of
1276 *de novo* gene transcription and germline mutational bias in *Drosophila*. *eLife* 8.
- 1277 Wong A, Albright SN, Giebel JD, Ram KR, Ji S, Fiumera AC, Wolfner MF. 2008. A role for
1278 Acp29AB, a predicted seminal fluid lectin, in female sperm storage in *Drosophila*
1279 *melanogaster*. *Genetics* 180:921-931.
- 1280 Wu DD, Irwin DM, Zhang YP. 2011. De novo origin of human protein-coding genes. *PLoS*
1281 *Genetics* 7.
- 1282 Xiao W, Liu H, Li Y, Li X, Xu C, Long M, Wang S. 2009. A rice gene of *de novo* origin negatively
1283 regulates pathogen-induced defense response. *PLOS ONE* 4:e4603.
- 1284 Xie C, Bekpen C, Künzel S, Keshavarz M, Krebs-Wheaton R, Skrabar N, Ullrich KK, Tautz D.
1285 2019. A *de novo* evolved gene in the house mouse regulates female pregnancy cycles.
1286 *eLife* 8.

- 1287 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular Biology and*
1288 *Evolution* 24:1586-1591.
- 1289 Yeh SD, Do T, Chan C, Cordova A, Carranza F, Yamamoto EA, Abbassi M, Gandasetiawan
1290 KA, Librado P, Damia E, et al. 2012. Functional evidence that a recently evolved
1291 *Drosophila* sperm-specific gene boosts sperm competition. *Proceedings of the National*
1292 *Academy of Sciences* 109:2043-2048.
- 1293 Zdobnov EM, Tegenfeldt F, Kuznetsov D, Waterhouse RM, Simão FA, Ioannidis P, Seppey M,
1294 Loetscher A, Kriventseva EV. 2017. OrthoDB v9.1: cataloging evolutionary and
1295 functional annotations for animal, fungal, plant, archaeal, bacterial and viral orthologs.
1296 *Nucleic acids research* 45:D744-D749.
- 1297 Zhang J. 2003. Evolution by gene duplication: an update. *Trends in Ecology & Evolution* 18:292-
1298 298.
- 1299 Zhang L, Ren Y, Yang T, Li G, Chen J, Gschwend AR, Yu Y, Hou G, Zi J, Zhou R, et al. 2019.
1300 Rapid evolution of protein diversity by de novo origination in *Oryza*. *Nature Ecology &*
1301 *Evolution* 3:679-690.
- 1302 Zhang W, Landback P, Gschwend AR, Shen B, Long M. 2015. New genes drive the evolution of
1303 gene interaction networks in the human and mouse genomes. *Genome biology* 16.
- 1304 Zhang YE, Vibranovski MD, Krinsky BH, Long M. 2010. Age-dependent chromosomal
1305 distribution of male-biased genes in *Drosophila*. *Genome Research* 20:1526-1533.
- 1306 Zhang YE, Vibranovski MD, Landback P, Marais GAB, Long M. 2010. Chromosomal
1307 Redistribution of Male-Biased Genes in Mammalian Evolution with Two Bursts of Gene
1308 Gain on the X Chromosome. *PLOS Biology* 8:e1000494.
- 1309 Zhao L, Saelao P, Jones CD, Begun DJ. 2014. Origin and spread of *de novo* genes in
1310 *Drosophila melanogaster* populations. *Science* 343:769-772.
- 1311 Zhuang X, Yang C, Murphy KR, Cheng CHC. 2019. Molecular mechanism and history of non-
1312 sense to sense evolution of antifreeze glycoprotein gene in northern gadids.
1313 *Proceedings of the National Academy of Sciences* 116:4400-4405.
- 1314