# Genomic Analyses of New Genes and Their Phenotypic Effects Reveal Rapid Evolution of Essential Functions in Drosophila Development — Source link ⧉

Shengqian Xia, Nicholas W. VanKuren, C. T. Chen, Li Zhang ...+10 more authors

**Institutions:** University of Chicago, Chinese Academy of Sciences, Ohio State University, University of São Paulo ...+1 more institutions

Related papers:

- Origins, evolution, and phenotypic impact of new genes

- Newly evolved genes: moving from comparative genomics to functional studies in model systems. How important is genetic novelty for species adaptation and diversification?

- Novel genes exhibit distinct patterns of function acquisition and network integration

- New genes as drivers of phenotypic evolution.

- New genes drive the evolution of gene interaction networks in the human and mouse genomes

Share this paper: 📘 🐦 in ✉

View more about this paper here: https://typeset.io/papers/genomic-analyses-of-new-genes-and-their-phenotypic-effects-2fpfwmc8xu

1

2

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

6

7

Shengqian Xia[*1], Nicholas W. VanKuren[1*], Chunyan Chen[2,3*], Li Zhang[1], Clause Kemkemer[1], Yi Shao[2,3], Hangxing Jia[2,3], UnJin Lee[1,4], Alexander S. Advani[4], Andrea Gschwend[5], Maria Vibranovski[6], Sidi Chen[7], Yong E. Zhang[2,3,8] and Manyuan Long[1,4]

12

13

14

1. Department of Ecology and Evolution, The University of Chicago, Chicago, USA.
2. State Key Laboratory of Integrated Management of Pest Insects and Rodents & Key Laboratory of Zoological Systematics and Evolution, Institute of Zoology, Chinese Academy of Sciences, Beijing 100101, China.
3. University of Chinese Academy of Sciences, Beijing 100049, China.
4. Committee on Genetics, Genomics and Systems Biology, The University of Chicago, Chicago, USA.
5. Department of Horticulture & Crop Science, The Ohio State University, Columbus, Ohio, USA
6. Department of Genetics and Evolutionary Biology, University of São Paulo, Sao Paulo, Brazil.
7. Department of Genetics, Yale School of Medicine, West Haven, USA.
8. Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming 650223, China.

30

31

32

33

34

35

* Co-first authors

Corresponding author: mlong@uchicago.edu; zhangyong@ioz.ac.cn

38

39

40

41

42 **ABSTRACT**

43 It is a conventionally held dogma that the genetic basis underlying development is

44 conserved in a long evolutionary time scale. Ample experiments based on mutational,

45 biochemical, functional, and complementary knockdown/knockout approaches have

46 revealed the unexpectedly important role of recently evolved new genes in the

47 development of *Drosophila*. The recent progress in the analyses of gene effects and

48 improvements in the computational identification of new genes, which has led to large

49 sample sizes of new genes, open the door to investigate the evolution of gene

50 essentiality with a phylogenetically high resolution. These advancements also raised

51 interesting issues related to phenotypic effect analyses of genes, particularly of those

52 that recently originated. Here we reported our analyses of these issues, including the

53 dating of gene ages, the interpretation of RNAi data that may confuse false

54 positive/false negative rates, and the potential confounding impact of compensation

55 and developmental effects that were not considered during previous CRISPR

56 knockout experiments. We further analyzed new data from knockdowns of 702 new

57 genes (~66% of total 1,070 *Drosophila melanogaster* new genes), revealing a

58 similarly high proportion of essential genes from recent evolution, compared to those

59 found in distant ancestors of *D. melanogaster*. Knockout of a few young genes

60 detected analogous essentiality. Furthermore, our experimentally determined

61 distribution and comparison of knockdown efficiency in different RNAi libraries

62 provided valuable data for general functional analyses of genes. Taken together, these

63 data, along with an improved understanding of the phenotypic effect analyses of new

64 genes, provide further evidence to the conclusion that new genes in *Drosophila*

65 quickly evolved essential functions in viability during development.

66

67

68

69 [*Keywords*: new gene; *Drosophila*; RNAi; the off-target effect; false positive; false

70 negative; CRISPR; essentiality; the compensation effect]

71

72

73

74

75

## INTRODUCTION

77

78   The question of how often new genes evolve essential functions is a critical problem

79   in understanding the genetic basis of development and general phenotypic evolution.

80   New genes in evolution have widely attracted discussion (Long and Langley. 1993;

81   Long et al. 2003; Chen et al. 2013; Carvunis et al. 2012; Ding et al. 2013; McLysaght

82   and Hurst. 2015), supported by increasing studies with fulsome evidence in various

83   organisms (e.g. Ruiz-Orera et al. 2018; Xie et al. 2019; Vakirlis et al. 2020; Witt et al.

84   2019; Jiang and Assis. 2017; Rogers et al. 2014; Schroeder et al. 2020). The detected

85   large number of new genes with unexpected rate of new gene evolution (e.g. Zhang et

86   al. 2019; Shao et al. 2019; Zhang et al. 2010a) and the revealed important functions of

87   new genes (Kasinathan et al. 2020; Lee et al. 2019; Long et al. 2013) challenged a

88   widely held dogma that the genetic basis in control of development is conserved in a

89   long time scale of evolution (Ashburner et al. 1999; Gould. 2002; Carroll. 2005;

90   Krebs et al. 2013). Our previous work used the RNAi knockdown in a smaller sample

91   showing that new genes may quickly become essential in *Drosophila* and that

92   potential for a gene to develop an essential function is independent of its age (Chen et

93   al. 2010). This work suggests a tremendous and quickly evolving genetic diversity,

94   which had not been previously anticipated.  Since then, genomes of better quality

95   from more species have allowed for more reliable new gene annotation (Shao et al.

96   2019).  In addition, technical progress in the detection of gene effects has increased

97   with better equipped knockdown libraries and direct CRISPR knockout methods.

98   Related scientific discoveries and technical development in knockdown and knockout

99   techniques -- e.g., Green et al (2014) and Kondo et al (2017), respectively -- can be

100  considered when investigating the evolution of gene essentiality.

101

102  We will present in this report our recent experiments and computational analyses,

103  examining a few important issues raised in recent years (e.g. by Kondo et al. (2017)

104  and Green et al. (2014)) that we find to be generally relevant for the detection of the

105  phenotypic effects of genes, particularly of those that recently originated. Our

106  investigations include the following: 1) the estimation of new gene ages; 2) an

107  evaluation of the knockdown efficiency distribution in RNAi experiments; 3) an

108  understanding of the differences between different RNAi libraries in phenotyping

3

109    large samples of new genes for viability effects; 4) an interpretation of knockout data

110    regarding the compensation effect. Our analyses, with additional evidence published

111    recently by our group and others, provide ample and strong evidence to further

112    support a notion suggested by the fitness effect analysis of new genes in *Drosophila*:

113    new genes have quickly evolved essential functions in viability during development.

114

115    **RESULTS AND DISCUSSION**

116

117    **Identification of *Drosophila*-specific genes by the age dating.**

118

119    Two new gene datasets are available for *D. melanogaster*, which include the dataset

120    of Kondo *et al* (2017; the K-dataset, the underlying pipeline as the K-pipeline) and the

121    dataset we recently reported (Zhang et al, 2010b; Shao et al, 2019 called as the

122    GenTree Fly dataset, the G-dataset). In order to determine which dataset is more

123    accurate and thus could be used in the downstream analyses, we estimated their

124    qualities by performing systematic comparison. Kondo *et al.* (2017) identified 1,182

125    new genes that postdated the split of *D. melanogaster* and *D. pseudoobscura*

126    (Branches 3~6, ~40 Mya(million years ago); Fig. 1A and 1B). They inferred the ages

127    of these genes by incorporating the UCSC DNA-level synteny information, homology

128    information based on comparison of annotated proteins and RNA expression

129    profiling. By contrast, we identified 654 new genes in this same evolutionary period

130    (Fig. 1A and 1B) using the same UCSC synteny information (Rhead et al. 2010) and

131    our maximum parsimony-based pipelines (Zhang et al. 2010b).

132

133    We investigated why the K-dataset was almost twice the G-dataset over the same

134    evolutionary period. The K-dataset contained 471 new gene candidates in the G-

135    dataset (471/654 = 72%) (Fig. 1B, red), among which 313 are of the exact same ages

136    while 158 show either younger (123 genes) or older (35 genes) (Supplemental Table

137    S1). For the remaining 183 genes absent in the K-dataset, we found 101 as authentic

138    new genes (Fig. 1B, deep blue; Supplemental Table S2) after extensive validation by

139    manually checking UCSC synteny information and four additional resources

140    including the FlyBase ortholog annotation, Ensembl Metazoa homolog annotation,

141    protein prediction in outgroup species and published literatures (see also Materials

4

142 and Methods). This result indicates a high false negative rate in the K-dataset. By

143 contrast, only 19 genes are old genes, which represent false positives in the G-dataset

144 (Fig. 1B, light purple). Then, for 45 genes (Fig. 1B, sky blue), they are located in

145 clusters of tandemly amplified genes or transposon-rich regions, where synteny is

146 often ambiguous or difficult to build especially when outgroup species genomes are

147 poorly assembled. Their ages are difficult to infer. Analogously, the final remaining

148 18 candidates (Fig. 1B, green) are dubious, either with inconsistent topology between

149 UCSC synteny information and Ensembl tree, marginal protein similarity between

150 species or gene structure model changes.

151

152 We next examined the 711 (1182-471=711) new gene candidates unique to the K-

153 dataset by manually examining phylogenetic distribution and their syntenic

154 relationship with genes in various species. We could only confirm 49 authentic new

155 genes, which represent false negatives of the G-dataset. By contrast, 318 out of 711

156 genes were incorrectly dated as new genes due to four problematic practices (Fig. 1B;

157 Supplemental Table S2 lists the 318 false positives): 1) neglecting 275 that have

158 orthologs in outgroup species; 2) taking 32 noncoding or pseudogene models as

159 protein coding genes; 3) treating 6 redundant entries of same genes as different genes;

160 and 4) misdating 5 polycistronic coding genes reported by the literatures. In addition,

161 242 (Fig. 1B, sky blue) genes are located in repetitive regions. To be conservative, the

162 G-dataset excluded these genes from dating. Then, the remaining 102 candidates (Fig.

163 1B, green) are dubious.

164

165 We conclude, based on above exhaustive manual evaluation, that the G-dataset is of

166 much higher quality compared to the K-dataset: 1) the false negative rate and false

167 positive rate of the G-dataset is estimated as 7.5% (49) and 2.9% (19), respectively; 2)

168 The both parameters of the K-dataset are higher, 8.5% (101) and 26.9% (318),

169 respectively; 3) the G-dataset only contains 9.6% (63) low-quality candidates (genes

170 in repetitive or dubious categories), while the K-dataset consists of 29.1% (344) such

171 candidates. Overall, 56% (662/1,182) new gene candidates in the K-dataset are either

172 false positive or dubious.

173

174 **Measuring reproducibility and efficiency of knockdown.**

175

5

176  We investigated the consistency of RNAi experiments with the same lines and the
177  same drivers in different laboratories, conditions, and years. Zeng et al. (2015)
178  screened 16,562 transgenic RNAi lines using an Act5C-Gal4 driver to detect the
179  lethality of 12,705 protein-coding genes (~90% of all annotated coding genes) in their
180  study of intestinal stem cell development and maintenance. Their dataset included
181  RNAi lines targeting the same 103 genes that were measured for lethality by Chen et
182  al. (2010). Chen et al. (2010) and Zeng et al. (2015) obtained the same phenotypes for
183  88 (85.4%) genes, including 30 (29.1%) of the lethal phenotype and 58 (56.3%) of
184  non-lethal phenotype (Fig. 2A, Supplemental Table S3). These data suggest that
185  despite differences in independent observers, lab environments, and years to conduct
186  experiments, the vast majority of RNAi knockdown experiments are reproducible for
187  phenotyping lethality and non-lethality.

188

189  We also tested consistency between RNAi lines with different RNAi drivers (called
190  new drivers) or same drivers in different genome positions. Specifically, the datasets
191  of Chen et al (2010) and Zeng et al (2015) shared 86 new genes in knockdown
192  experiments, mostly (81.4%, 70) with different RNAi drivers and fewer (18.6%, 16)
193  same drivers in different genome positions (Supplemental Table S4).  This dataset
194  showed that: 7 genes were consistently lethal; 42 genes were consistently non-lethal;
195  and 37 genes have different phenotypes (Fig. 2B). Thus, the two groups with different
196  drivers or same drivers with different positions show that more genes (57.0%, 49)
197  have the same phenotypes.

198

199  We considered an additional factor in RNAi knockdown, sensitivity, in the two
200  widely used RNAi libraries: the Vienna Drosophila Resource Center's (VDRC's) GD
201  and KK libraries (Dietzl et al., 2007).  The GD libraries were constructed using P-
202  elements to randomly insert hairpin RNAs (average 321bp) into the genome targeting
203  individual genes, while the KK library inserted constructs carrying hairpin RNAs
204  (average 357bp) into a specific landing site by ΦC31-mediated homologous
205  recombination. All KK lines carry an insertion at 30B3, but a proportion (23-25%)
206  also carry an insertion at 40D3 (*tio* locus) that results in pupal lethality when using
207  constitutive drivers like Act5C-GAL4 (Green et al. 2014; Vissers et al. 2016). Unless
208  specified, no lines discussed below contain 40D3 insertions.

209

6

210   Given the intrinsic different designs of GD and KK libraries, we hypothesized that

211   they have different false negative or false positive rates, which cause the

212   inconsistency shown in Fig. 2B. Only GD lines were examined previously, and they

213   have a high false negative rate (39.9%) but low false positive rate (<2%) (Dietzl et al.

214   2007). The high false negative rate is likely caused by insufficient target gene

215   knockdown, while false positives may be due to off target effects (Dietzl et al. 2007).

216   We thus tested the knockdown efficiency of 75 KK lines targeting randomly selected

217   75 young genes (Supplemental Table S5, Fig. 3A). We found that the knockdown

218   efficiency of KK lines is generally lower than the efficiency of 64 GD lines as

219   previously reported (Dietzl et al. 2007). Specifically, using the same driver (Act5C),

220   we found that in general, GD lines have significantly higher knockdown efficiency

221   than KK lines, as shown by the knockdown expression as the percentage of the

222   control expression (Fig. 3A). That is, the KK lines have an average knockdown

223   efficiency as 48.6% of control expression while the GD lines show an average

224   efficiency as 38.1% (Fig. 3B and 3C, $t$-test $P$ = 0.031). Notably, the expression

225   reduction to 50~60% level of the wide-type level was observed to have no significant

226   fitness loss due to widespread haplosufficiency (Huang et al. 2010; VanKuren and

227   Long, 2018). Detecting any fitness effect may be expected when the expression drops

228   to a lower level, for example, 20~30% or lower of the control expression.  In this

229   range of knockdown efficiency, we observed that only 29% of KK lines but 41% of

230   GD lines reduced target expression levels to ≤20% of control levels; 37% of KK lines

231   but 53% of GD lines were seen to reduce target expression levels to ≤ 30% of control

232   levels (Fig. 3A). Thus, it is expected that GD lines have a significantly higher power

233   in detecting lethal phenotypes as shown in the next section.

234

235   To estimate false positive rate of KK lines, we constructed 10 randomly chosen new

236   KK lines targeting one member of a young duplicate gene pair, in addition to one KK

237   lines and 3 TRiP lines (Transgenic RNAi Project, BDSC, Materials and Methods).

238   The rationale is that for each gene of interest its paralog is the most likely off target.

239   The same rationale was also followed by (Dietzl et al. 2007) when false positive rates

240   of GD lines were estimated. We measured the knockdown efficiency and estimated

241   off-target effects using these 14 lines with qPCR experiments in adult whole bodies

242   (Fig. 4).  We found that two lines likely produce off-target effects (*NV-CG31958*,

7

243    *34008* (the TRiP line)), for both of which the expression of paralog is down-regulated

244    to similar or even lower level compared to the corresponding gene of interest. Twelve

245    other lines have significantly higher target effects than off-target effects, among

246    which 10 genes reduced activity to 20-80% expression level of the control (7 genes

247    reduced activity to 20-40%) and only two genes (*CG32164*, *CG7046*) reach≤20% of

248    control levels. Thus, if we take 20% as the cutoff of efficient knockdown, only

249    *CG31958* could be counted as the false positive, and *CG32164* and *CG7046* be

250    counted as the true positives. Collectively speaking, the off-target effects are rare

251    while insufficient knockdowns are pervasive.

252

253    These experiments detected a variation of knockdown efficiency among different

254    drivers where newer KK lines have lower efficiency and thus higher false negatives

255    compared to older GD lines. Therefore, these observations offer an alternative

256    interpretation of the incongruence than the false-positive-only rationale of Kondo et al

257    (2017): when new RNAi drivers were added to the analysis, insufficient knockdown

258    was also introduced with a high probability. This would create incongruence between

259    old and new drivers if the old and new triggers have significantly different sensitivity.

260

261    **Phenotyping essentiality of new genes in RNAi libraries.**

262

263    We first investigated differences in measured lethality between the KK and GD

264    libraries used in Chen et al (2010). To control for the confounding effect of *tio*

265    insertion in the KK lines, we genotyped these lines using PCR-amplification and

266    found that out of 153 KK lines we collected, 47 (30.7%) had two landing sites and 6

267    (3.9%) had only 40D3 landing site (the confounding site) (Green et al. 2014), which

268    all showed lethal phenotypes  (Supplemental Table S6). Using the recombination

269    approach (Green et al. 2014), we recovered 41 of the 47 lines into the lines that have

270    only the 30B3 site.  The RNAi knockdown of 140 KK lines carrying insertions only at

271    30B3 identified 12 genes (8.6%) with lethal phenotypes (Supplemental Table S6).

272    Meanwhile, 12 genes in 59 GD lines (20.3%) were detected to have lethal knockdown

273    effects (Chen et al. 2010), significantly higher than the KK libraries ($P = 0.0112$,

274    Fisher's Exact Test). As aforementioned, this difference is likely due to higher false

275    negative rate of KK lines (Fig. 3).

8

276

277    By using the essentiality data of 10,652 old genes provided by VDRC

278    (https://stockcenter.vdrc.at/control/library_rnai) that were in branch 0 (Shao et al.

279    2019), we characterized the statistical distribution of essential old genes (Fig. 5). We

280    independently sampled 1000 times, with each randomly sampling 150 old genes and

281    calculating the proportion of essential ones. We found that in the GD library, the

282    probability to obtain a proportion of essential new genes equal or lower than 20.3% is

283    0.780. Meanwhile, in the KK library, the probability to observe a proportion of

284    essential new genes equal or lower than 8.6% is 0.867. These analyses of GD and KK

285    libraries reveal similarly that the proportions of new and old genes with lethal

286    phenotypes are not statistically different.

287

288    Further analysis of gene essentiality data in a recent version of VDRC libraries

289    (retrieved online in April 2019) detected with increased resolution the proportions of

290    essential genes in six detectable ancestral stages of *D. melanogaster*. We reported the

291    analysis of the GD library, which has a significantly higher knockdown efficiency

292    than the KK library. In total, 11,354 genes (72% of 15,682 genes in the species,

293    Ensembl 73) have been phenotyped for their lethality or nonlethality, including 702

294    *Drosophila* genus specific genes (66% of 1,070 detected *Drosophila*-specific genes)

295    (Long et al. 2013; Shao et al. 2019) and 10,652 genes that predated the *Drosophila*

296    divergence 40 Mya.

297

298    We parsimoniously mapped the 702 *Drosophila*-specific genes on the six ancestral

299    stages by examining their species distribution in the *Drosophila* phylogeny (Shao et

300    al. 2019) (Fig. 6A). Of the 702 genes, 19.7% (138) are directly observed to be

301    essential, similar to the proportion of essential old genes, 18.9% ($P = 0.6212$, Fisher's

302    exact test). We considered a low knockdown efficiency as shown by the 47% of GD

303    lines whose knockdowns are expressed at the level of 30% or higher of the control

304    (Fig. 3A), suggesting that 47% of RNA lines are invalid for the testing and should be

305    subtracted from the total tested lines.

306

307    Thus, the actual proportion of essential genes can be estimated by correcting for the

308    bias of false positives (Fp) and false negatives (Fn) by following formula:

309

9

310          Corrected proportion of essential genes $= [E - (T \cdot F_p)] / [T - (T \cdot F_n)]$

311

312   Where E and T are observed number of essential genes and total number of genes

313 examined, respectively. $F_p$ was measured as 1.6% (Dietzl et al., 2007) while $F_n$ as

314 47% as estimated above or 39.9% as measured previously (Dietzl et al., 2007). Thus,

315 the estimated proportion of essential genes after correcting false positives and false

316 negatives can be as high as 36.5% for the estimated false negative rate of 47% in this

317 study. The corrected proportion can be also as high as 32.2% given the previously

318 measured false negative rate of 39.9%. Furthermore, all six stages show a stable

319 proportion of essential genes; none of the proportions is statistically different from the

320 proportion of old genes (Fig. 6A). Meanwhile, lethal rates of new genes which belong

321 to three origin mechanism categories (DNA-based duplication, RNA-based

322 duplication and orphan genes, Shao et al., 2019) also show no significant difference

323 (Fig. 6B). Interestingly, 21.7% of orphan genes, some of which might be *de novo*

324 genes (Long et al., 2013), are essential. These data add new insight into the evolution

325 of essentiality in all ancestral stages: soon after genes originated and fixed in *D.*

326 *melanogaster*, a stable proportion of new genes is essential throughout entire

327 evolutionary process from ancient ancestors to the speciation of *D. melanogaster*.

328

329 These data of knockdown experiments on a large number of new genes further

330 supported what we proposed before: *Drosophila* new genes rapidly evolve essential

331 functions within the divergence of *Drosophila* genus; knockdown of these genes leads

332 to death of flies.

333

334 **Analyses of mutants identified young essential genes**

335

336 Kondo et al (2017) recommended and used CRISPR/Cas9-mediated mutagenesis to

337 create small frameshift indel mutations in targeted genes. This method has two

338 potential issues. First, it is now well documented that vertebrate cells such as

339 mammalian cells or zebrafish cells recognize such aberrant mRNAs and compensate

340 for their loss by increasing expression of genes with high sequence similarity, such as

341 paralogs in zebrafish, worm and other organisms (Rossi et al. 2015; El-Brolosy and

342 Stainier 2017; El-Brolosy et al. 2019); Ma et al, 2019; Serobyan et al, 2020). This has

343 the effect of producing false negatives especially for recent duplicates. We confirmed

10

344    that a similar compensation effect exists in *Drosophila*. Specifically, when we

345    induced a one-nucleotide deletion using CRISPR/Cas9 into the ORF region of *vismay*

346    (*vis*), a *D. melanogaster*-specific gene duplicated from a parental gene, *achintya*

347    (*achi*), 0.8 Mya, with a nucleotide similarity of 92% between the two copies. We

348    found that *achi* in the *vis* mutant was significantly upregulated whereas a randomly

349    selected unrelated gene *CG12608* and the distantly related gene *hth* (nucleotide

350    similarity of 45%) to *vis*, did not show such an effect (Fig. 7). Second, CRISPR/Cas9-

351    mediated mutagenesis cannot detect the effects of maternal and paternal effect genes,

352    which can be common in *Drosophila* (Perrimon et al. 1989; Raices et al. 2019) and

353    can be detected by RNAi knockdown. Therefore, the two approaches of knockdown

354    and knockout/mutagenesis are complementary to each other given their technical

355    characteristics.

356

357    Actually, in-depth analyses of several cases already provided further evidence

358    supporting essentiality of new genes in development. First, Ross et al (2013) reported

359    a stepwise neofunctionalization evolution in which a centromere-targeting gene in

360    *Drosophila*, *Umbrea*, was generated less than 15 Mya. Both RNAi knockdown,

361    rescue experiments and P-element mediated gene knockout revealed that *Umbrea*

362    evolved a species-specific essentiality to target centromere in chromosome

363    segregation (Chen et al. 2010; Ross et al. 2013). Second, Lee et al (2019) recently

364    detected stage-specific (embryos/larvae/pupa) lethality associated with RNAi

365    knockdown and CRISPR knockout in *Cocoon,* a gene emerged 4 Mya in the common

366    ancestor of the clade of *D. melanogaster-simulans*. These data show that *Cocoon* is

367    essential for the survival at multiple developmental stages, including the critical

368    embryonic stage. Third, P-element insertion/excision experiments show the

369    essentiality of *K81* as a paternal element in early development. This gene only exists

370    in the *Drosophila melanogaster*-subgroup species that diverged 6 Mya (Loppin et al.

371    2005). Fourth, Zeus, a gene that duplicated from the highly conserved transcription

372    factor *CAF40* 4 Mya in the common ancestor of *D. melanogaster* and *D. simulans*

373    rapidly evolved new essential functions in male reproductive functions, as detected in

374    the null mutants and knockdown (Chen et al. 2012; Ventura,  2019).  Fifth, A pair of

375    extremely young duplicates, *Apollo* (*Apl*) and *Artemis* (*Arts*), was found to have been

376    fixed 200,000 years ago in *D. melanogaster* populations (VanKuren and Long, 2018).

11

377    CRISPR-created gene deletions of these genes showed that both evolved distinct

378    essential functions in gametogenesis and *Apl* critical function in development. Sixth,

379    in a comprehensive functional and evolutionary analysis of the ZAD-ZNF gene

380    family in *Drosophila* (Kasinathan et al, 2020), 86 paralogous copies were identified

381    with phenotypic effects detected by knockdown and knockout in *D. melanogaster*. It

382    was found that the proportion (17/58 = 29.3%) of lethal copies in old duplicates (>40

383    Mya) and the proportion (11/28 = 39.3%) of lethal copies in *Drosophila*-specific

384    duplicates (<40Mya) are statistically similar. Further functional analyses of one of the

385    non-essential young copies (*CG17802, Nicknack*) reported by Kondo et al (2017)

386    clearly unveiled an essential function for larval development. These pieces of

387    evidence strongly support the notion that new genes can quickly evolve essential

388    functions in development.

389

390    **Concluding Remarks**

391

392    We appreciate the extensive experiment, computation and data-compilation by Kondo

393    et al (2017) and their interests in the evolution of gene essentiality. However, we

394    found that the K-pipeline and K-dataset were associated with a high false positive

395    rate. Moreover, their interpretation of RNAi data is problematic due to confusing the

396    false negative and false positive, while they applied an incorrect CRISPR mutagenesis

397    that neglected compensation and parental effects. The data we created in this study,

398    while revealing their errors and technical insufficiencies, increase understanding of

399    technical subtleties for analyzing effects of young duplicate genes and other genes.

400    More data and additional analyses of related scientific issues for the testing of fitness

401    and functional effects of new genes from the two complementary approaches, RNAi

402    knockdown and CRISPR knockout, provided a strong support for the concept: the

403    new genes rapidly evolved essential functions in development in *Drosophila*. This

404    challenges a conventional belief in the antiquity of important gene functions in

405    general (Jacob, 1977; Mayr, 1982; Ashburner et al. 1999; Krebs et al. 2013) and in

406    development process in specific (Gould. 2002; Carroll. 2005).

407

408

409    **MATERIALS AND METHODS**

410    ***Comparison of the K-dataset and the G-dataset for Drosophila new genes.***

411    *Overall comparison scheme*

412    To our knowledge, there is no published genome-wide evaluation of gene ages in
413    *Drosophila*. Specifically, around a decade ago, we took advantage of the syntenic
414    genomic (DNA) alignment generated by the UCSC group and performed the genome-
415    wide age dating for the first time in *Drosophila* (Zhang et al. 2010b). At that moment,
416    we compared our data to previous studies based on limited number of cases and
417    discovered the general reproducibility across studies. The genome-wide dataset by
418    Kondo *et al.* enabled a systematic large-scale comparison. We began with repeating
419    our pipeline on FlyBase annotation v6.02 and identified 654 new genes originated on
420    the branch toward *D. melanogaster* after the species split of *D. melanogaster* and *D.*
421    *pseudoobscura* (Fig. 1). Concurrently, based on similar FlyBase release v6.13, Kondo
422    *et al.* performed dating according to the same syntenic alignment of UCSC, which is
423    further complemented by v6.02 protein-level BLAST search and filter with testis-
424    specific expression (Kondo et al. 2017). With additional Dollo-parsimonious searches,
425    they identified 1,182 new genes originated in the same period including
426    *melanogaster*-group to *melanogaster*-only, which correspond to age group 3 to 6 in
427    our analysis, respectively (Supplemental Fig. S1). Since the FlyBase annotation
428    version is similar (v6.02 vs. v6.13), only 12 entries out of the G-dataset and K-dataset
429    are not comparable due to "Gene model change" (Supplemental Fig. S1). They
430    represent either expired models or new models in v6.13. Except them, all other genes
431    can be compared across two datasets.

432

433    We found that 471 out of the 654 new gene candidates in the G-dataset are covered in
434    the K-database by comparing the Ensembl IDs of these databases.  Moreover, 313
435    (66%) genes show the exact same ages. Since manual curation needs extensive
436    efforts,  we did not examine why the remaining 158 genes show minor age difference.
437    Instead, we subsequently only focused on those genes which show conflicting dating
438    results, *i.e.*, included or excluded in new gene dataset across two studies. As a result,
439    we classified the conflicting cases into six major categories, which can be further
440    divided into around 20 more specific sub-categories (Supplemental Fig. S1). We
441    documented how we performed classification as below.

442

443     *Four independent information sources facilitate evaluations of two age datasets*

444

445     The challenges in the dating of gene age largely lie in the ambiguity of calling

446     orthologs across outgroup species (Liebeskind et al. 2016). We found that the conflict

447     of age dating was often due to the difference of DNA-level synteny and protein-level

448     homology search. Specifically, for a gene of interest, *A*, the UCSC best-to-best

449     synteny information shows that its ortholog is present in one outgroup species, *B*.

450     However, the protein-level information may reveal an absence. The opposite scenario

451     can occur too. In these conflicts, we turned to independent resources including

452     FlyBase ortholog annotation, the homolog annotation and gene family tree provided

453     by Ensembl Metazoa (St Pierre et al. 2014; Kersey et al. 2015), protein prediction in

454     outgroup species based on gene models of *D. melanogaster* and literatures.

455     Specifically, FlyBase provided AAA (Assembly/Alignment/Annotation) syntenic

456     ortholog annotation. If species *B* encodes a FlyBase annotated ortholog, gene *A* likely

457     predated the species-split of *D. melanogaster* and *B*. Similarly, Ensembl Metazoa

458     provided one-to-one best-to-best ortholog annotation. We used it like FlyBase.

459     Finally, for some cases where synteny predicted orthologous regions of *B* do not

460     harbor an annotated gene, we conceptually translated this region with the protein of

461     *D. melanogaster* as the template. BLAST (Tblastn) was used here. We have two

462     reasons to perform additional annotation: 1) recently evolved genes are often poorly

463     annotated; 2) annotation quality of outgroup species is presumably worse compared to

464     *D. melanogaster* and we need to correct this bias. For particularly interesting cases

465     (*e.g.* polycistronic coding genes), we searched literatures describing their evolutionary

466     history.

467

468     *Conflicting cases could be classified into six categories*

469

470     We implemented a series of customized rules to call the presence of ortholog of gene

471     *A* in species *B*. The first set of rules are used to call presence of ortholog based on

472     gene prediction. For a synteny-predicted candidate orthologous region in *B*, we ran

473     Tblastn to predict whether this region encodes an orthologous protein of *A*. If Tblastn

474     could align the protein of *D. melanogaster* beyond the following thresholds (identity

475     cutoff > 70% & coverage cutoff > 30%, identity cutoff > 30% & coverage cutoff >

476     70%, identity cutoff > 50% & coverage cutoff > 50%), we believed that the ortholog

14

477   is present. If the alignment meets with the threshold (identity cutoff < 30% &

478   coverage cutoff < 30%), the ortholog of *A* is absent in *B*. For all other cases, we called

479   them as "boundary cases" if there is also no ortholog annotated by FlyBase and

480   Ensembl. A total of 80 candidate new genes fall in this category including 65 cases in

481   the list of Kondo *et al.* and 15 cases in our list (Supplemental Fig. S1).

482

483   Secondly, for 275 out of 318 new genes dated by Kondo *et al* but not by us

484   (Supplemental Fig. S1), we identified orthologs as supported by at least two

485   independent sources (FlyBase, Ensembl homolog, Ensembl gene family tree, and/or

486   prediction). For example, in case of *FBgn0027589*, the ortholog is present across all

487   12 *Drosophila* species, which is supported by both Ensembl and FlyBase. The

488   remaining 43 new genes are misidentified due to other types of problems (*e.g.*

489   annotation problem due to polycistronic structure such as *tal-1A/tal-2A/tal-3A*). All

490   these 318 genes are marked as "Dating problem in Kondo *et al*" (Supplemental Fig.

491   S1). In the opposite scenario, 101 new genes are only identified by us, which could be

492   divided into four cases: 1) for 50 genes called as old genes in the K-dataset, their new

493   gene calling were also supported by lack of one-to-one orthologs annotated by

494   Ensembl or FlyBase; 2) for 5 genes called as old genes in the K-dataset, they are

495   subject to complex evolutionary trajectories (pseudogenization of parental copies),

496   such as *FBgn0032740* and *Cyp6t1*; 3) for 10 cases excluded by the K-dataset, we

497   examined phylogenetic trees provided by Ensembl and confirmed that new gene are

498   derived as suggested by longer branch length; 4) for 36 genes excluded in the K-

499   dataset, FlyBase and Ensembl do not annotate orthologous genes in the outgroup

500   species for most genes, which are also consistent with the lack of Tblastn hits.

501

502   Compared to the K-dataset, the false positives and false negatives are much fewer in

503   the G-dataset (Fig. 1, Supplemental Fig. S1). In the G-dataset, we misidentified only

504   19 new genes with 14 cases caused by double or triple losses in the outgroups (*e.g.*

505   *CG2291*). In our pipeline, we first searched against closely-related species and then

506   went for remotely-related species. If at least two independent losses are needed to

507   explain the phylogenetic distribution of orthologs, we will assign a young age to gene

508   *A* by following the maximum parsimony. The underlying assumption is: 1) the

509   possibility of double or triple losses should be low; 2) the genomic alignment between

510   *D. melanogaster* and remotely-related species is less reliable compared to that

15

511    between *D. melanogaster* and closely-related species. Consistently, we only identified

512    14 cases with support by at least one additional source (FlyBase, Ensembl). The

513    remaining 5 cases are caused by lack of sensitivity of UCSC genome alignment. Both

514    FlyBase and Ensembl annotate orthologs in the outgroup species, but synteny does

515    not cover the corresponding regions. In opposite, genome alignment built spurious

516    alignment in remotely-related species for 49 new genes identified by Kondo *et al*.

517    However, our Tblastn search could not identify a protein at all. Thus, we referred

518    them as false negatives. All these 68 genes are put into the third category entitled with

519    "Dating problem in this work" (Supplemental Fig. S1).

520

521    A fourth category ("Not applicable or difficult for dating") consists of genes which

522    are most resistant for dating due to their sequence features undesirable for new gene

523    identification. 242 specific new genes claimed by Kondo *et al.* belongs to this

524    category (Supplemental Fig. S1). For 178 out of 242 genes, we found that synteny is

525    in conflict for outgroup species sharing the same phylogenetic relationship relative to

526    *D. melanogaster* (e.g. *D. simulans* and *D. sechellia*). These genes are generally

527    located in repetitive regions (*e.g.* tandem amplification). It is thus likely that the

528    orthologous regions may not be equally well assembled or be subject to species-

529    specific gene conversion across these outgroup species. We thus excluded these

530    genes. For example, we masked 19 out of 242 genes including 12 *Ste* genes, 5 Y-

531    linked genes and 2 genes encoded by contigs but not anchored to five major

532    chromosome arms. *Ste* is the X-linked tandem gene families each with redundant

533    copies (Supplemental Fig. S2A). In the UCSC Net track, the most assembles can only

534    reach level 2 of one-way syntenic mapping, rather than the adequate reciprocal

535    syntenic mapping as level 1. The closely related species (*e.g. D. simulans*) also

536    encodes multiple copies, but the corresponding region is not fully assembled and

537    filled with lots of gaps and many of these copies even cannot be assigned to

538    chromosomes (Supplemental Fig. S2B). The size contrast of assemblies between *D.*

539    *sechellia* and *D. melanogaster* suggests that the region in *D. sechellia* might not be

540    properly assembled due to its repetitive structure (Supplemental Fig. S2A and S2C).

541

542    In order to date each member correctly, high-quality outgroup genome must be

543    available first. As for 5 Y-linked genes, Koerich *et al.* (2008) assigned all of them to

544    be old genes (Koerich et al. 2008). The remaining 45 genes consist of three subtypes:

16

545    1) 21 fast-evolving small proteins (<100 amino acids); 2) 20 polycistronic genes

546    without previous literature support; 3) 4 tandem duplicates. Different from the

547    aforementioned 178 genes, the syntenic information is consistent across outgroup

548    species suggesting that they are old genes. The small proteins or polycistronic genes

549    are poorly annotated across outgroups. For 4 duplicates, Ensembl phylogenetic tree

550    could not provide diagnostic information to infer the duplication order. Thus, we

551    believed that all these three subtypes are difficult to date as of now. In our private new

552    gene list, 45 candidates also belong to gene families. Similar to the above 4 tandem

553    duplicates, the synteny is consistent across outgroup species showing that these 45

554    genes are derived copies in their respective families. However, Ensembl phylogenetic

555    tree could not provide additional support. So, we also put these genes into the same

556    fourth category.

557

558    The fifth category ("Different ortholog definition") only consists of 28 private new

559    genes identified by Kondo *et al* (Supplemental Fig. S1). For 24 out of 28 cases, the

560    UCSC syntenic chain only covers a small portion (<30%) of coding regions or mainly

561    corresponds to untranslated regions (UTRs) in outgroup species. Since the dating of

562    Kondo *et al* is protein-centric, they called these genes as new genes. By contrast, our

563    dating pipeline works on DNA-level and identified these genes as old genes. The

564    reason why we took the age of most conserved exons to represent the age of whole

565    genes is that these exons usually represent most important functional regions.

566    Moreover, by performing dating on DNA-level, our dating does not depend on

567    annotation quality of outgroup species. Actually, for all 24 cases, whether the coding

568    region is accurately annotated is unknown due to the lack of protein evidence. For the

569    remaining 4 cases, they all represent translocated genes. In our terminology, we only

570    referred derived duplicate or orphan genes as new genes. By contrast, Kondo *et al.*

571    interpreted incorrectly translocated genes as new genes.

572

573    The sixth and final category refers to the aforementioned 12 entries not comparable

574    due to "Gene model change" (Supplemental Fig. S1).

575

576    In the main text, we merged "Gene model change", "boundary cases" and "different

577    ortholog definition" as one dubious (in green color of Fig. 1B) category to simplify.

578

579    *D. melanogaster-specific gene identification*

580

581    Candidate new genes were initially collected from previous studies (Zhou et al. 2008;

582    Zhang et al. 2010b; Chen et al. 2012). We removed from this list of 233 candidates: 1)

583    any genes whose *D. melanogaster*-specific release 6.05 (http://flybase.org) annotation

584    status is 'withdrawn', 2) genes not located on the major chromosome arms 2L, 2R ,

585    3L, 3R, or X, and 3) members of large tandem arrays, including the *Sperm dynein*

586    *intermediate chain* (Nurminsky et al. 1998; Yeh et al. 2012), *Stellate* (*Ste*), and X:

587    19,900,000-19,960,000 arrays that are *D. melanogaster*-specific but are impossible to

588    specifically study. We checked syntenic whole-genome alignments of the remaining

589    84 genes manually using our multi-species alignments at the UCSC Genome Browser.

590    To be conservative, we required that all outgroups including the *D. simulans*,

591    *D.sechellia*, *D. yakuba*, and *D. erecta* genome assemblies contained no assembly

592    gaps, transposable elements, or repeats corresponding to the flanking regions of the

593    putative *D. melanogaster*-specific gene. *D. melanogaster*-specific gene origination

594    mechanisms and parental genes were taken from the original studies and confirmed

595    using BLAT and BLASTp. If a gene had multiple significant (e < $10^{-10}$) full-length

596    BLASTp hits *in D. melanogaster*, the hit that was most similar to the *D.*

597    *melanogaster*-specific gene was assumed to be the parent. We used available *D.*

598    *simulans* and *D. yakuba* next-generation sequencing reads to test the presence of

599    putative *D. melanogaster*-specific tandem duplications in these two species (Green et

600    al. 2014; Rogers et al. 2014). We found no breakpoint spanning read pairs supporting

601    *D. melanogaster* tandem duplications in any of 20 *D. simulans* or 20 *D. yakuba*

602    genomes. Thus, these tandem duplications are specifically found in *D. melanogaster*

603    and are not simply missing from the *D. yakuba* and *D. simulans* reference genome

604    assemblies. We checked if any of the duplications in our final set are segregating

605    rather than being fixed within *D. melanogaster* by analyzing 17 whole genome re-

606    sequencing data from the DPGP2 core Rwanda (RG) genomes (Ni et al. 2008). We

607    required tandem duplications to have at least one read uniquely mapped to each of the

608    three unique breakpoints in order to be called as 'present' in a particular line. Ten

609    genes are not found in any of 17 additional *D. melanogaster* genomes we analyzed,

610    suggesting that they are found specifically in the reference stock. Finally, we curated

611    10 *D. melanogaster*-specific genes. This dataset is actually a subset of G_K new gene

612    data list.

613

*RNAi strain construction*

615

Since species-specific new genes are under-represented in public RNAi lines, we generated new RNAi lines following Dietzl et al. (2007). Briefly speaking, we designed RNAi reagents using the E-RNAi server (http://www.dkfz.de/signaling/e-rnai3/) and kept constructs with all possible 19-mers uniquely matching the intended target gene and excluded designs with >1 CAN repeat (simple tandem repeats of the trinucleotide with N indicates any base) (Ma et al. 2006). Constructs were cloned into pKC26 following the Vienna *Drosophila* Resource Center's (VDRC's) KK library strategy (http://stockcenter.vdrc.at, last accessed 2 February 2016). We introgressed the X chromosome from Bloomington *Drosophila* Stock Center line 34772, which expresses ΦC31 integrase in ovary under control of the *nanos* promoter, into the VDRC 60100 strain. Strain 60100 carries attP sites at 2L:22,019,296 and 2L:9,437,482 (Green et al. 2014). We ensured that our RNAi constructs were inserted only at the 2L:9,437,482 site using PCR following Green et al. (2014). RNAi constructs were injected into the 60100-ΦC31 at 250 ng/μL. Surviving adult flies were crossed to sna$^{Sco}$/CyO balancer flies (BDSC 9325) and individual insertion strains were isolated by backcrossing.

632

*RNAi screen*

634

We knocked down target gene expression using driver lines constitutively and ubiquitously expressing GAL4 under the control of either the *Actin5C* or *αTubulin84B* promoter. We replaced driver line's balancer chromosomes with GFP-marked chromosomes to track non-RNAi progeny. Control crosses used flies from the background strains 60100-ΦC31, 25709, or 25710 crossed to driver strains. Five males and five virgin driver females were used in each cross. Crosses were grown at 25°C, 40% - 60% humidity, and a 12h:12h light:dark cycle. F1 progeny were counted at day 19 after crossing, after all pupae had emerged. We screened F1 RNAi flies for visible morphological defects in 1) wings: vein patterning and numbers, wing periphery; 2) notum: general bristle organization and number, structure and smoothness; 3) legs: number of segments. We monitored survival of RNAi F1s by counting GFP and non-GFP L1, L3 larvae and pupae. We tested RNAi F1 sterility by

647 crossing individual RNAi F1 flies to 60100-ΦC31 and monitoring vials for L1

648 production. Ten replicates for each sex for each line were performed.

649

650 *RNAi knockdown specificity and sensitivity*

651

652 We sought to address two known problems of RNAi technology using RT-qPCR.

653 First, since off-target effects are often discussed in RNAi experiments(Dietzl et al.

654 2007) we need to test whether target gene expression are specifically knocked down,

655 although our constructs are computationally predicted to be specific. Second, since

656 the RNAi knockdown is often incomplete (Dietzl et al. 2007), we need to estimate

657 how many genes are adequately knockdown in expression. We targeted a random

658 dataset of 14 *D. melanogaster*-specific genes. We collected qPCR primers from

659 FlyPrimerBank (Hu et al. 2013). For those genes not found in FlyPrimerBank, we

660 took Primer-BLAST to design primers by specifically targeting a ~100 bp region of

661 the gene (Supplemental Table S7). We confirmed primer specificity with PCR and

662 Sanger sequencing.

663

664 We randomly selected 75 KK RNAi lines (no *tio* site insertion) to analyze their knock

665 down efficiency. We cross these 75 KK RNAi lines with same driver which was used

666 in Dietzl et al., 2007 for GD RNAi line knock down efficiency test. We extracted

667 RNA from sets of 8 adult males (2~4 day old) in triplicate from each RNAi cross

668 using TRIzol (Catalog# 15596-026, Invitrogen, USA), treated ~2 μg RNA with

669 RNase-free DNase I (Catalog# M0303S NEW ENGLAND Biolabs, USA), then used

670 1 μL treated RNA in cDNA synthesis with SuperScript III Reverse Transcriptase

671 (Invitrogen, USA) using oligo(dT)$_{20}$ primers. cDNA was diluted 1:40 in water before

672 using 1 μL as template in 10 μL qPCRs with iTaq$^{TM}$ Universal SYBR Green

673 Supermix (Catalog# 1725121, Bio-Rad, USA) and 400 nM each primer. Reactions

674 were run on a Bio-Rad C1000 Touch thermal cycler with CFX96 detection system

675 (BioRad, CA). Cycling conditions were 95°C for 30 sec, then 45 cycles of 95°C for 5

676 sec, 60°C for 30 sec, and 72°C for 15 sec. We normalized gene expression levels

677 using the $\Delta\Delta C_T$ method and *RpL32* as the control (Livak and Schmittgen 2001; Dietzl

678 et al. 2007). We tested the specificity and efficiency (90%< qPCR Efficiency<110%)

679 of qPCR primers using an 8-log$_2$ dilution series for each primer pair (VanKuren and

680 Long 2018).

20

681

*Testing Compensation Effects of New Gene Duplicates*

683

684 We generated the frameshift mutation line of *vis* using the CRISPR protocol

685 previously developed (VanKuren and Long 2018) but with one single sgRNA for one

686 gene as Kondo et al (2017) did. The sgRNA-*vis* primer below was synthesized (the

687 underlined sequence):

688 5'-GAAATTAATACGACTCACTATAGG<u>ATGTACGGCAGAACATAA</u>GTTTAA

689 GAGCTATGCTGGAA-3';

690 We used the following sequence-specific qRT-PCR primers to test the compensatory

691 expression of *achi*, the duplicate of *vis*. Two control genes including *CG12608* and

692 *hth* were examined too. Since *vis*'s expression is largely testis-specific, we extracted

693 RNAs from testis of mated 4-day males and used qRT-PCR with 3 replicates to assess

694 the expression, as developed previously (VanKuren and Long 2018).

695 248bp:

696 Achi-RT1F: 5'-AAAGTGACAGGTTTCTCTGTTTG-3';

697 Achi-RT1R: 5'-CTGATCCTCCTCCACGATGAC-3'.

698 237bp:

699 CG12608-RT1F: 5'-CATAGTGGGCACCTACGAG-3';

700 CG12608-RT1R: 5'-TGCGAGAGTATGATCTGCGAC-3'.

701 92bp:

702 hth-RT1F:5'-CCTAGTCATGTATCGCCGGTC-3';

703 hth-RT1R:5'-AGCGGATGTTCATAAATCGCA-3'.

704 Internal control:

705 113bp:

706 RpL32-RT1F: 5'-AGCATACAGGCCCAAGATCG-3';

707 RpL32-RT1R: 5'-TGTTGTCGATACCCTTGGGC-3'.

708

## ACKNOWLEDGMENTS

714 results for KK lines of the correct insertion site. Y.E. Zhang was supported by the
715 National Key R&D Program of China (2019YFA0802600, 2018YFC1406902) and
716 the National Natural Science Foundation of China (91731302, 31771410, 31970565).
717 M. Long was supported by NSF1026200 and NIH R01GM116113.

718
719
720
721



722
723

724 **Figure 1.** Summary of new gene candidates in the K-dataset and G-dataset. A.
725  phylogenetic distribution of gene origination identified by the K-pipeline and the
726  G-pipeline as shown in the two datasets. B. Evaluation of the two datasets based
727  on individual gene analyses. The two datasets share 471 candidates (red). The G-
728  dataset consist of 101 authentic candidates (deep blue) undetected in the K-
729  dataset, 19 false positives (light purple), 18 dubious cases (green) and 45 cases
730  not applicable for dating (sky blue). By contrast, the K-dataset includes 49 bona
731  fide new gene candidate, 318 false positives, 102 dubious cases and 242 difficult
732  cases. Note, the K-dataset mentions 1,182 genes in the main text, however its
733  associated supplemental table includes 1,176 genes with 6 genes listed more than
734  once.

735
736

22

737

738

**Figure 2.** The reproducibility analysis of RNAi experiments by comparing two groups of independent experiments by Chen et al (2010) and Zeng et al (2015). A. Phenotypes of same 103 RNAi lines analyzed by Chen et al (2010) and Zeng et al (2015) using same lines; B. Phenotypes of 86 same new genes knocked down by two different drivers or the same drivers with different insertion sites. The old drivers detected 29 genes as lethal while 57 non-lethal; the new drivers detected 20 genes as lethal while 66 non-lethal.

746

747

748

749

750

751

752

753

754

23

**Figure 3.** Knockdown efficiency in the KK and GD libraries revealed GD lines have significantly higher knockdown efficiency than the KK lines. A. The knockdown efficiency of the 75 KK lines was measured, compared to the expression of the wild-type control and the standard deviation is calculated from the measurement of three repeats; *P* refers to proportion of genes with the expression lower than a certain threshold while the values of KK lines are generated in this work and that of GD lines are extracted from Dietzl et al. (2007). B. The distributions of knockdown efficiency of KK and GD lines. C. The Q-Q Plots between KK and GD lines.

24

**Figure 4.** Experimental comparison of the efficiency and off-target effects explain the conservative nature of RNAi knockdown experiments and limited off-targets propensity. For each young duplicate gene pair specific for *D. melanogaster* and *melanogaster* species complex, we examined their expression intensity relative to the wide type control in whole body flies with qPCR. The standard deviation is calculated based on three replicates.
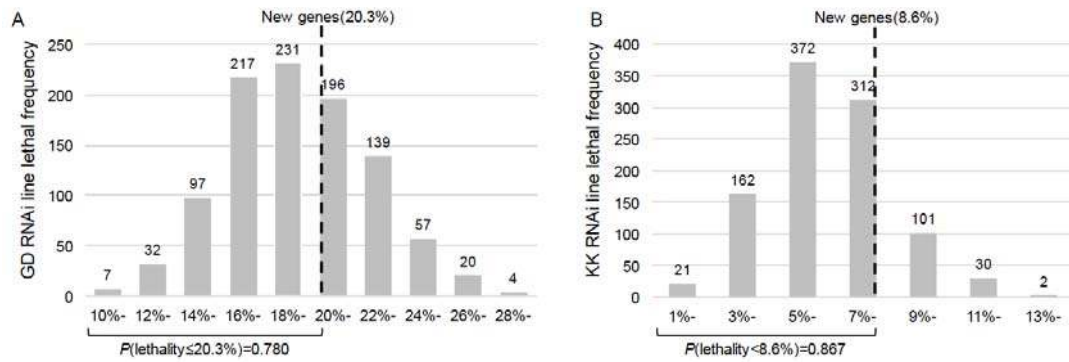
**Figure 5.** Comparison of proportions of lethality between new genes and old genes in GD lines (**A**) and KK lines (**B**) suggests that in both GD and KK lines, new genes have an equally high probability to be lethal as old genes. Since old genes are much more abundant than new genes, we generated 1000 random sample of old genes with the same number of new genes and then plotted the distribution of proportion of essential genes as histograms.
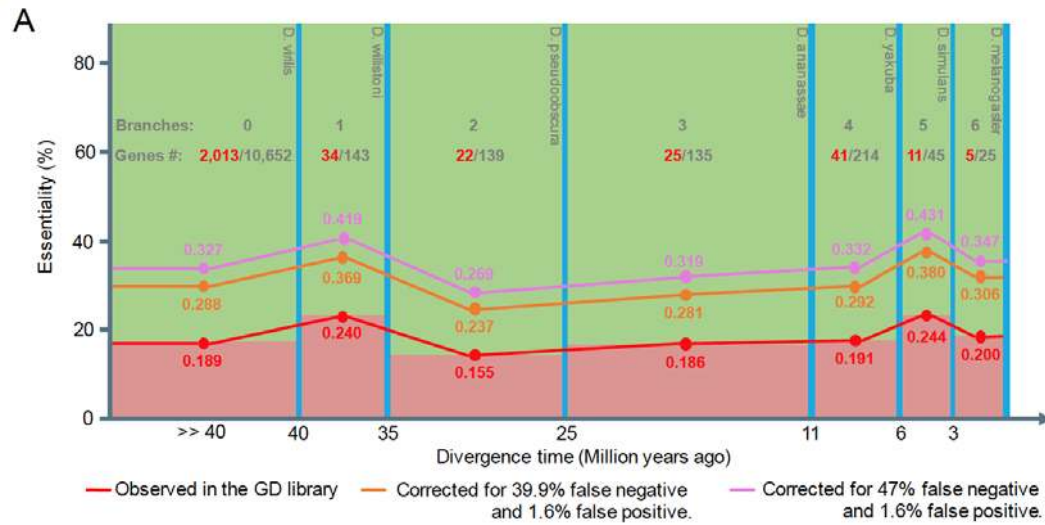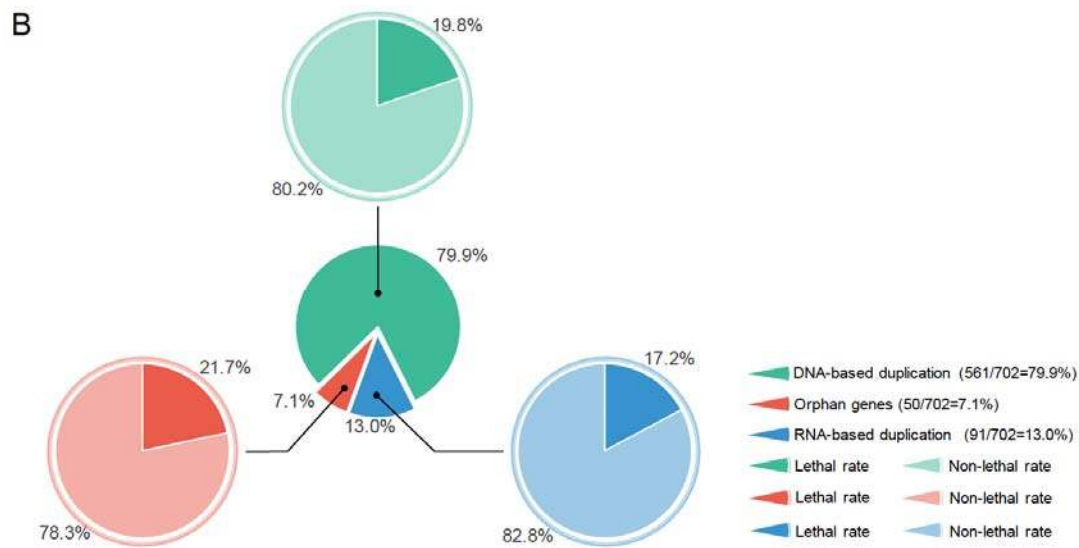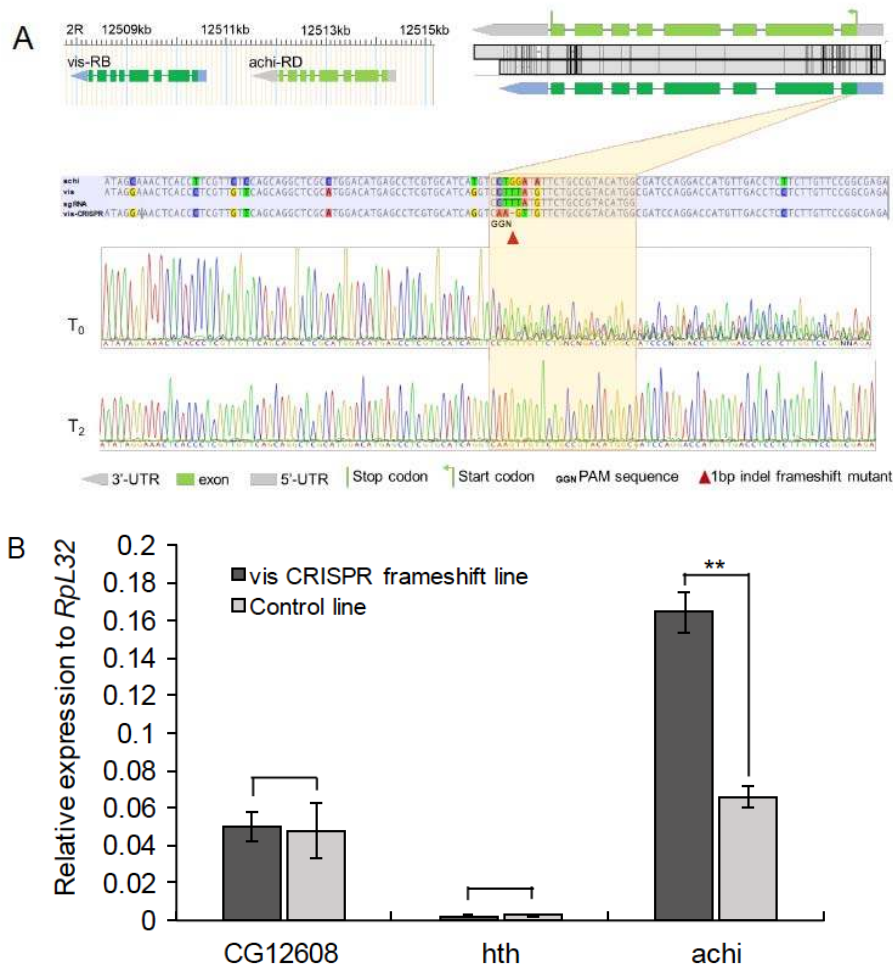
812



813

**Figure 6.** Lethality proportion of 702 *Drosophila*-specific genes. A. Lethality proportion of 702 *Drosophila*-specific genes in 6 ancestral stages of extant *D. melanogaster*, compared to the lethality proportion of 10,652 genes older than 40 Mya. No stages show an essentiality proportion significantly different from that of old genes (0.189). B. Lethality proportion of 702 *Drosophila*-specific genes based on three origin mechanism catalogs. No catalog shows a lethality proportion significantly different from that of old genes (0.189).

821

**Figure 7.** CRISPR/Cas9 frameshift mutant could induce compensatory effect in *Drosophila*. A. Design of CRISPR/Cas9 mutant. We targeted a randomly chosen young gene, *vis*, which emerged via duplication of *achi* in the common ancestor of *melanogaster* species complex. The genomic arrangement of two genes are shown in the upper left panel with the boxes referring to exons and connecting lines as introns. The pair shares a high sequence identity (0.92) in their 9 exons, which is schematically shown in the upper right panel. The middle panel shows the diverged site between *vis* and *achi*, which was chosen to design a short guide RNA (sgRNA) specifically targeting *vis*. The mutation (CTTTA→AAGT) was marked with a red triangle. The raw sanger sequencing data for the initial generation (T0) and the second generation of offspring (T2) was shown. B. The compensation effect of *achi*. In the frameshift mutant of *vis*, *achi*'s expression is

28

836    significantly increased (*P*=0.0003). By contrast, the unrelated *CG12608* and the

837    remotely related *hth* did not show any significant upregulation. *RpL32* was used

838    as a control as in (VanKuren and Long 2018).

839

840    **Supplementary Figures**

841

842    **Figure S1.** Age dating between this work and Kondo *et al.* This figure, following Fig.

843    1 in the main text, adds specific information on how we classified genes into six

844    major categories or dozens of subcategories. For more details, please refer to

845    Materials and Methods.

846

847    **Figure S2.** A representative difficult-to-date locus in the *K*-dataset. A. The syntenic

848    view of *Ste* locus between *D. melanogaster* and *D. simulans* shows fragmented

849    continuity. Due to its multiplicative nature, *Ste* locus is difficult to assemble. In

850    the UCSC Net track, the most assembles can only reach level 2 of one-way

851    syntenic mapping, rather than a better reciprocal syntenic mapping as level 1. B.

852    Some orthologous region in *D. simulans* (lifted from *D. melanogaster*) is not

853    anchored to the chromosome (X) and they are arbitrarily assembled as chrU. C.

854    In *D. sechellia*, two scaffolds are assembled with the major scaffold super_20

855    spanning 200 kb, in contrast to the assembly of 15 kb for the orthologous region

856    of *D. melanogaster*.

857

858    **Supplementary Tables**

859

860    **Table S1**. The list of genes with the exact ages across the G-dataset and the K-dataset,

861    genes with slightly younger ages in the K-dataset and genes with slightly older

862    ages in the K-dataset, respectively.

863

864    **Table S2**. The list of false negatives and false positives in the K-dataset. Since the

865    *Pan-Drosophilid* age group in the K-dataset corresponds to the age group 0, 1 or

866    2 in the G-dataset (Fig. 1A), we simply replaced the *Pan-Drosophilid* age group

867    as 0/1/2 in the table if applicable.

868

869    **Table S3**. 103 knockdown experiments repeated by two independent works (Chen et

870     al. 2010; Zeng et al. 2015). Note, Chen et al (2010) works classified phenotypes

871     as lethal, semi-lethal and viable. Since there are only few genes deemed as semi-

872     lethal, we merged them into lethal gene groups to simplify.

873

874     **Table S4**. For 86 new genes with different RNAi drivers, the consistency between

875     different drivers in Chen et al (2010) and Zeng et al (2015) is listed.

876

877     **Table S5**. The knockdown efficiency data of KK library and GD library.

878

879     **Table S6**. The genotyping results of 153 KK lines, the corrected lines by

880     recombination and knockdown results.

881

882     **Table S7**.  Primers for 75 KK lines knockdown efficiency tests.

883

884     **REFERENCES**

885     1. Ashburner M, Misra S, Roote J, Lewis SE, Blazej R, Davis T, Doyle C, Galle
886        R, George R, Harris N et al. 1999. An exploration of the sequence of a 2.9-Mb
887        region of the genome of Drosophila melanogaster: the Adh region. *Genetics*
888        **153**: 179-219.
889     2. Carroll SB. 2005. *Endless Forms Most Beautiful*: The New Science of Evo
890        Devo. W. W. Norton & Company
891     3. Carvunis AR, Rolland T, Wapinski I, Calderwood MA, Yildirim MA, Simonis
892        N, Charloteaux B, Hidalgo CA, Barbette J, Santhanam B et al. 2012. Proto-
893        genes and de novo gene birth. *Nature* **487**: 370-374.
894     4. Chen SD, Krinsky BH, Long MY. 2013. New genes as drivers of phenotypic
895        evolution. *Nat Rev Genet* **14**: 645-660.
896     5. Chen SD, Spletter M, Ni XC, White KP, Luo L, Long M. 2012. Frequent
897        recent origination of brain genes shaped the evolution of foraging behavior in
898        Drosophila. *Cell Rep* **1**: 118-132.
899     6. Chen SD, Zhang YE, Long MY. 2010. New genes in Drosophila quickly
900        become essential. *Science* **330**: 1682-1685.
901     7. Dietzl G, Chen D, Schnorrer F, Su KC, Barinova Y, Fellner M, Gasser B,
902        Kinsey K, Oppel S, Scheiblauer S et al. 2007. A genome-wide transgenic
903        RNAi library for conditional gene inactivation in Drosophila. *Nature* **448**:
904        151-156.
905     8. Ding Y, Zhou Q, Wang W. 2013. Origins of new genes and evolution of their
906        novel functions. *Annu Rev Ecol Evol Syst* **43**: 345-363.
907     9. El-Brolosy MA, Kontarakis Z, Rossi A, Kuenne C, Gunther S, Fukuda N,
908        Kikhi K, Boezio GLM, Takacs CM, Lai SL et al. 2019. Genetic compensation
909        triggered by mutant mRNA degradation. *Nature* **568**: 193-197.
910     10. El-Brolosy MA, Stainier DYR. 2017. Genetic compensation: A phenomenon
911        in search of mechanisms. *Plos Genet* **13**: e1006780.

11. Gould SJ. 2002. The structure of evolutionary theory. Kelknap Press of Harvard University Press. Cambridge, Massashusetts and London, England.

12. Green EW, Fedele G, Giorgini F, Kyriacou CP. 2014. A Drosophila RNAi collection is subject to dominant phenotypic effects. *Nat Methods* **11**: 222.

13. Hu Y, Sopko R, Foos M, Kelley C, Flockhart I, Ammeux N, Wang X, Perkins L, Perrimon N, Mohr SE. 2013. FlyPrimerBank: an online database for Drosophila melanogaster gene expression analysis and knockdown evaluation of RNAi reagents. *G3* **3**: 1607-1616.

14. Huang N, Lee I, Marcotte EM, Hurles ME. 2010. Characterising and predicting haploinsufficiency in the human genome. *Plos Genet* **6**: e1001154.

15. Jacob F. 1977. Evolution and tinkering. *Science* **196**: 1161-1166.

16. Jiang X, Assis R. 2017. Natural selection drives rapid functional evolution of young Drosophila duplicate genes. *Mol Biol Evol* **34**:3089-3098.

17. Kasinathan B, Colmenares SU, McConnell H, Young JM, Karpen GH, Malik HS, 2020. Innovation of heterochromatin functions drives rapid evolution of essential ZAD-ZNF genes in *Drosophila* . bioRxiv doi.org/10.1101/2020.07.08.192740.

18. Kersey PJ, Allen JE, Armean I, Boddu S, Bolt BJ, Carvalho-Silva D, Christensen M, Davis P, Falin LJ, Grabmueller C. 2015. Ensembl Genomes 2016: more genomes, more complexity. *Nucleic Acids Res* **44**: D574-D580.

19. Koerich LB, Wang XY, Clark AG, Carvalho AB. 2008. Low conservation of gene content in the Drosophila Y chromosome. *Nature* **456**: 949-951.

20. Kondo S, Vedanayagam J, Mohammed J, Eizadshenass S, Kan LJ, Pang N, Aradhya R, Siepel A, Steinhauer J, Lai EC. 2017. New genes often acquire male-specific functions but rarely become essential in Drosophila. *Genes Dev* **31**: 1841-1846.

21. Krebs JE, Gildstein ES and Kilpatrick ST. 2013. Lewin's essential genes. Jones & Bartlett Publishers.

22. Lee YCG, Ventura IM, Rice GR, Chen D-Y, Colmenares SU, and Long M. 2019. Rapid Evolution of Gained Essential Developmental Functions of a Young Gene via Interactions with Other Essential Genes. *Mol Biol Evol* **36**: 2212–2226.

23. Liebeskind BJ, McWhite CD, Marcotte EM. 2016. Towards Consensus Gene Ages. *Genome Biol Evol* **8**: 1812-1823.

24. Livak KJ, Schmittgen TD. 2001. Analysis of relative gene expression data using real-time quantitative PCR and the 2(T)(-Delta Delta C) method. *Methods* **25**: 402-408.

25. Long MY, Langley CH. 1993. Natural selection and the origin of jingwei, a chimeric processed functional gene in Drosophila. *Science* **260**: 91-95.

26. Long, MY, Betrán E, Thornton K, and Wang W. 2003. The origin of new genes: glimpses from the young and old. *Nature Reviews Genetics* **4**: 865-875.

27. Long MY, VanKuren NW, Chen SD, Vibranovski MD. 2013. New gene evolution: Little did we know. *Annu Rev Genet* **47**: 307-333.

28. Loppin B, Lepetit D, Dorus S, Couble P, Karr TL. 2005. Origin and neofunctionalization of a Drosophila paternal effect gene essential for zygote viability. *Current Biology* **15**: 87-93.

29. Ma Y, Creanga A, Lum L, Beachy PA. 2006. Prevalence of off-target effects in Drosophila RNA interference screens. *Nature* **443**: 359-363.

30. Ma Z, Zhu P, Shi H, Guo L, Zhang Q, Chen Y, Chen S, Zhang Z, Peng J, Chen J. 2019. PTC-bearing mRNA elicits a genetic compensation response via

31

962          Upf3a and COMPASS components. *Nature* **568**:259–263.

963 31. Mayr EJ. 1982. The Growth of Biological Thought - Diversity, Evolution, and
964          Inheritance. *New York Rev Books* **29**: 41-42.

965 32. Ni JQ, Markstein M, Binari R, Pfeiffer B, Liu LP, Villalta C, Booker M,
966          Perkins L, Perrimon N. 2008. Vector and parameters for targeted transgenic
967          RNA interference in Drosophila melanogaster. *Nat Methods* **5**: 49-51.

968 33. Nurminsky DI, Nurminskaya MV, De Aguiar D, Hartl DL. 1998. Selective
969          sweep of a newly evolved sperm-specific gene in Drosophila. *Nature* **396**:
970          572-575.

971 34. Perrimon N, Engstrom L, Mahowald AP. 1989. Zygotic Lethals with Specific
972          Maternal Effect Phenotypes in Drosophila-Melanogaster .1. Loci on the X-
973          Chromosome. *Genetics* **121**: 333-352.

974 35. Raices JB, Otto PA, Vibranovski MD. 2019. Haploid selection drives new
975          gene male germline expression. *Genome Res* **29**: 1115-1122.

976 36. Rhead B, Karolchik D, Kuhn RM, Hinrichs AS, Zweig AS, Fujita PA,
977          Diekhans M, Smith KE, Rosenbloom KR, Raney BJ et al. 2010. The UCSC
978          Genome Browser database: update 2010. *Nucleic Acids Res* **38**: D613-D619.

979 37. Rogers RL, Cridland JM, Shao L, Hu TT, Andolfatto P, Thornton KR. 2014.
980          Landscape of standing variation for tandem duplications in Drosophila yakuba
981          and Drosophila simulans. *Mol Biol Evol* **31**: 1750-1766.

982 38. Ross BD, Rosin L, Thomae AW, Hiatt MA, Vermaak D, de la Cruz AFA,
983          Imhof A, Mellone BG, Malik HS. 2013. Stepwise evolution of essential
984          centromere function in a Drosophila neogene. *Science* **340**: 1211-1214.

985 39. Rossi A, Kontarakis Z, Gerri C, Nolte H, Holper S, Kruger M, Stainier DYR.
986          2015. Genetic compensation induced by deleterious mutations but not gene
987          knockdowns. *Nature* **524**: 230-235.

988 40. Ruiz-Orera J, Verdaguer-Grau P, Villanueva-Cañas JL, Messeguer X, Mar
989          Albà M. 2018. Translation of neutrally evolving peptides provides a basis for
990          de novo gene evolution. *Nature Ecol Evol* **2**, 890–896.

991 41. Schroeder, CM, Tomlin, SA, Valenzuela, JR and Malik, HS, 2020. A rapidly
992          evolving actin mediates fertility and developmental tradeoffs in Drosophila.
993          bioRxiv.

994 42. Serobyan V, Kontarakis Z, El-Brolosy MA, Welker JM, Tolstenkov O,
995          Saadeldein AM, Retzer N, Gottschalk A, Wehman AM, Stainier DYR, 2020.
996          Transcriptional adaptation in Caenorhabditis elegans *eLife* **9**: e50014.

997 43. Shao Y, Chen C, Shen H, He BZ, Yu D, Jiang S, Zhao S, Gao Z, Zhu Z, Chen
998          X et al. 2019. GenTree, an integrated resource for analyzing the evolution and
999          function of primate-specific coding genes. *Genome Res* **29**: 682-696.

1000 44. St Pierre SE, Ponting L, Stefancsik R, McQuilton P, FlyBase C. 2014.
1001          FlyBase 102--advanced approaches to interrogating FlyBase. *Nucleic Acids*
1002          *Res* **42**: D780-D788.

1003 45. Vakirlis N, Acar O, Hsu B, Castilho Coelho N, Van Oss SB, Wacholder A,
1004          Medetgul-Ernar K, Bowman RW, 2nd, Hines CP, Iannotta J et al. 2020. De
1005          novo emergence of adaptive membrane proteins from thymine-rich genomic
1006          sequences. *Nat Commun* **11**: 781.

1007 46. VanKuren NW, Long MY. 2018. Gene duplicates resolving sexual conflict
1008          rapidly evolved essential gametogenesis functions. *Nat Ecol Evol* **31**: 705-712.

1009 47. Ventura IM. 2019. Functional Evolution of Young Retrogenes with
1010          Regulatory Roles in Drosophila. The University of Chicago Ph.D. dissertation
1011          10.6082/uchicago.1799.

48. Vissers JHA, Manning SA, Kulkarni A, Harvey KF. 2016. A Drosophila RNAi library modulates Hippo pathway-dependent tissue growth. *Nat Commun* **7**: 10368.

49. Witt E, Benjamin S,Svetec N, Zhao L. 2019. Testis single-cell RNA-seq reveals the dynamics of de novo gene transcription and germline mutational bias in Drosophila. *eLife* **8**: e47138.

50. Xie C, Bekpen C, Kunzel S, Keshavarz M, Krebs-Wheaton R, Skrabar N, Ullrich KK, Tautz D. 2019. A de novo evolved gene in the house mouse regulates female pregnancy cycles. *elife* **8**.

51. Yeh SD, Do T, Chan C, Cordova A, Carranza F, Yamamoto EA, Abbassi M, Gandasetiawan KA, Librado P, Damia E. 2012. Functional evidence that a recently evolved Drosophila sperm-specific gene boosts sperm competition. *Proc Natl Acad Sci* **109**: 2043-2048.

52. Zeng XK, Han LL, Singh SR, Liu HH, Neumuller RA, Yan D, Hu YH, Liu Y, Liu W, Lin XH et al. 2015. Genome-wide RNAi screen identifies networks involved in intestinal stem cell regulation in Drosophila. *Cell Rep* **10**: 1226-1238.

53. Zhang L, Ren Y, Yang T, Li G, Chen J, Gschwend AR, Yu Y, Hou G, Zi J, Zhou R et al. 2019. Rapid evolution of protein diversity by de novo origination in Oryza. *Nat Ecol Evol* **3**: 679-690.

54. Zhang, YE, Vibranovski, MD, Landback P, Marais GA and Long MY. 2010a. Chromosomal redistribution of male-biased genes in mammalian evolution with two bursts of gene gain on the X chromosome. *PLoS Biol* **8**: e1000494.

55. Zhang YE, Vibranovski MD, Krinsky BH, Long M. 2010b. Age-dependent chromosomal distribution of male-biased genes in Drosophila. *Genome Res* **20**: 1526-1533.

56. Zhou Q, Zhang GJ, Zhang Y, Xu SY, Zhao RP, Zhan ZB, Li X, Ding Y, Yang S, Wang W. 2008. On the origin of new genes in Drosophila. *Genome Res* **18**: 1446-1455.