

# 1 Intermolecular interactions drive protein adaptive and co-adaptive evolution 2 at both species and population levels

3 Junhui Peng, Li Zhao

4 Laboratory of Evolutionary Genetics and Genomics, The Rockefeller University, New York, NY  
5 10065, USA

6 \*Correspondence to: [lzhao@rockefeller.edu](mailto:lzhao@rockefeller.edu)

## 8 **Abstract**

9 Proteins are the building blocks for almost all the functions in cells. Understanding the molecular  
10 evolution of proteins and the forces that shape protein evolution is an essential step in  
11 understanding the basis of function and evolution. Previous studies have shown that adaptation  
12 occurs frequently at the protein surface, such as in genes involved in host-pathogen  
13 interactions. However, it remains unclear whether adaptive sites are distributed randomly or at  
14 regions that are associated with particular structural or functional characteristics across the  
15 genome, since many of the proteins lack structural or functional annotations. Here, we seek to  
16 tackle this question by combining large-scale bioinformatic prediction, structural analysis,  
17 phylogenetic inference, and population genomic analysis of *Drosophila* protein-coding genes.  
18 By estimating and comparing the rate of adaptive substitutions at protein and residue level, we  
19 showed that adaptation is more relevant to function-related rather than structure-related  
20 properties. Among the function-related properties, we found that molecular interactions in  
21 proteins contribute to adaptive evolution, and putative binding residues exhibit higher rates of  
22 adaptation. We observed that physical interactions might play a role in the co-adaptation of fast-  
23 adaptive proteins. We found that strongly differentiated amino acids in protein coding genes are  
24 mostly adaptive, which may contribute to the long-term adaptive evolution. Our results suggest  
25 important roles of intermolecular interactions and co-adaptation in the adaptive evolution of  
26 proteins both at the species and population levels.

27

## 28 **Introduction**

29 Natural selection plays an important role in molecular evolution of protein sequences. Recent  
30 advances in genome sequencing and reliable inference methods at both phylogenetic and  
31 population levels have enabled fast and robust estimation of evolutionary rates and adaptation  
32 driven by natural selection. In addition, the increased availabilities of structural and functional  
33 data of proteins have made it possible to study how structural and functional constraints affect

34 protein sequence evolution and adaptation. It is now well established that different proteins and  
35 different sites within a protein have varying rates of evolution and adaptation due to both  
36 structural and functional constraints (Echave et al., 2016; Kosiol et al., 2008; Lindblad-Toh et al.,  
37 2011; Zhang and Yang, 2015). For example, genes that are highly expressed or perform  
38 essential functions are under strong purifying selection and tend to evolve slowly (Drummond et  
39 al., 2005; Moutinho et al., 2019; Pál et al., 2001; Zhang and He, 2005; Zhang and Yang, 2015);  
40 genes involved in host-pathogen interactions, e.g., immune responses and antiviral responses,  
41 show exceptionally high rates of adaptive changes (Enard et al., 2016; Nielsen et al., 2005;  
42 Obbard et al., 2009; Palmer et al., 2018; Sackton et al., 2007; Sironi et al., 2015; Uricchio et al.,  
43 2019); and residues that are intrinsically disordered or at the protein surface are fast evolving  
44 and has been proved to be hotspots of adaptive evolution (Afanasyeva et al., 2018; Goldman et  
45 al., 1998; Lin et al., 2007; Moutinho et al., 2019; Ramsey et al., 2011). More recently,  
46 Slodkowitz & Goldman (Slodkowitz and Goldman, 2020) employed genomic-scale integrated  
47 structural and phylogenetic evolutionary analysis in mammals and showed that positively  
48 selected residues are clustered near ligand binding sites, especially in proteins that are  
49 associated with immune responses and xenobiotic metabolism.

50         Although evidence have shown that adaptation is more likely to occur at intrinsically  
51 disordered regions and clustered at the surface of proteins, the functional properties of  
52 adaptation in the genomic scale remains unclear. Moreover, due to lack of structural and  
53 functional information of many proteins in the genome, the underlying mechanism derived from  
54 current studies might be incomplete. Here, we systematically investigated the evolution and  
55 adaptation of protein-coding genes in *Drosophila melanogaster* by comparing it to its closely  
56 related species, in order to distinguish the main factors that impact the evolution and adaptation at  
57 the protein-coding level. We applied large-scale bioinformatic and structural analysis to obtain  
58 structural and functional properties of proteins. We then classified residues into different  
59 structural and functional sites. By comparing rates of sequence evolution and adaptation  
60 between different proteins and different sites, we were able to locate hotspots of adaptation at  
61 genome scale. We showed that, for *D. melanogaster* proteins, adaptation is more sensitive to  
62 functional properties rather than structural ones. Interestingly, we found that putative binding  
63 regions including allosteric sites at protein surface show higher rates of adaptation than other  
64 sites. For proteins that are under fast-adaptive evolution, we showed that they tend to interact  
65 with each other more frequently than random expectations and are often associated with  
66 reproduction, immunity, and environmental information processing in *D. melanogaster*. In  
67 addition, we showed that interacting proteins in *D. melanogaster* might undergo co-adaptive

68 evolution. Furthermore, we hypothesize that molecular interactions or physical interactions  
69 might be an important mechanism that contribute to the adaptive and co-adaptive evolution in *D.*  
70 *melanogaster* genome. At last, we showed that the accumulation of short-term adaptation to  
71 local environments could be a possible genetic mechanism that contribute to long-term adaptive  
72 evolution.

73

## 74 **Results**

### 75 **Impact of gene properties on evolution of protein-coding genes in *D. melanogaster***

76 To uncover the main factors that impact the evolutionary rates of genes, we analyzed 13528  
77 protein-coding genes in *D. melanogaster* using genome data from *melanogaster* subgroup  
78 species and *D. melanogaster* population genomics data from 205 inbred lines from  
79 Drosophila Genetic Reference Panel, Freeze 2.0, DGRP2 (Huang et al., 2014). We applied a  
80 maximum likelihood method (Yang, 2007) to compute dN/dS ratio ( $\omega$ ) using the protein-coding  
81 sequences of five closely related melanogaster subgroup species (*D. melanogaster*, *D.*  
82 *simulans*, *D. sechellia*, *D. yakuba* and *D. erecta*). We estimated the proportions of adaptive  
83 changes ( $\alpha$ ) in each gene by applying an extension of MK test named asymptotic MK (Messer  
84 and Petrov, 2013; Uricchio et al., 2019) using *D. simulans* as outgroup. We then calculated the  
85 rate of adaptive changes ( $\omega_a$ ) of each gene by multiplying  $\omega$  to  $\alpha$  ( $\omega_a = \alpha\omega$ ) (Moutinho et al.,  
86 2019) using *D. yakuba* as the outgroup species (See methods). The rate of nonadaptive  
87 changes can be further calculated by  $\omega_{na} = \omega - \omega_a$ . Finally, we successfully assigned  $\omega$  to 12118  
88 protein coding genes and  $\omega_a$  and  $\omega_{na}$  to 7192 genes.

89 For each of *D. melanogaster* genes subjecting the same pipeline of analysis, we further  
90 obtained 17 different structural or functional properties (see Methods), which can be further  
91 divided into two categories: structure-related properties and function-related properties.  
92 Specifically, structure-related properties include ratio of secondary structures (helix ratio, sheet  
93 ratio, helix+sheet ratio, coil ratio), intrinsic structural disorder (ISD), relative solvent accessibility  
94 (RSA); while function-related properties include gene pseudo-age, protein length, number of  
95 protein-protein interactions (PPI numbers), ratio of protein-binding sites (PPI-site ratio), ratio of  
96 DNA-binding sites (DNA-site ratio) and gene expression patterns such as male expression level,  
97 female expression level, mean expression level, male specificity and tissue specificity. The  
98 properties along with gene-specific protein evolution ( $\omega$ ,  $\omega_a$  and  $\omega_{na}$ ) are available in  
99 supplementary file S1.

### 100 ***Molecular interactions contribute to the variations of protein sequence evolution***

101 ***and adaptation.*** In order to identify the determinants that drive protein evolution ( $\omega$ ,  $\omega_a$  and

102  $\omega_{na}$ ), we calculated the Pearson's correlations of  $\omega$ ,  $\omega_a$  and  $\omega_{na}$  with all the structure- and  
103 function-related properties. The correlation coefficient ( $r$ ) and corresponding p-values ( $p$ ) of  
104 each of the properties were listed in Table 1. Interestingly, we observed that for structure-related  
105 properties (secondary structure ratios, ISD, and RSA), variation of  $\omega$  is dominated by  
106 nonadaptive changes ( $\omega_{na}$ ) (Figure S1). Taking RSA as an example, we observed that RSA  
107 strongly correlates with both  $\omega$  ( $r=0.16$ ,  $p=1e-73$ ) and  $\omega_{na}$  ( $r=0.15$ ,  $p=3e-35$ ), while weakly  
108 correlates with  $\omega_a$  ( $r=0.06$ ,  $p=1e-6$ ). These correlations suggest that, under the constraints of  
109 structure-related properties, relaxation of purifying selection may play a more important role in  
110 determine protein evolution. These are in line with previous studies that proteins with less  
111 structural constraints, i.e. those harboring more disordered, exposed sites display faster  
112 evolutionary and nonadaptive evolutionary rate (Afanasyeva et al., 2018; Moutinho et al., 2019)

113 However, for function-related properties (gene pseudo-age, protein length, PPI number,  
114 PPI-site ratio, DNA-site ratio and gene expression patterns), the importance of  $\omega_a$  in shaping  
115 protein evolution begin to emerge (Figure S1). For example, when considering tissue specificity,  
116 the correlation efficient ( $r$ ) of  $\omega$  is 0.30 ( $p=2e-205$ ), while  $r$  of  $\omega_a$  and  $\omega_{na}$  are 0.16 ( $p=3e-35$ ) and  
117 0.17 ( $p=2e-42$ ), respectively. In such cases, the correlations of  $\omega_a$  and  $\omega_{na}$  almost contributed  
118 equally to the variation of protein sequence evolutionary rates,  $\omega$ . Interestingly, among the  
119 function-related properties, we found that molecular interactions, i.e., protein interactions,  
120 strongly positively correlates with  $\omega$ ,  $\omega_a$  and  $\omega_{na}$  (Table 1). We also noticed that for molecular  
121 interactions, compared to other function-related properties, variations of  $\omega_a$  contributes slightly  
122 to variations of  $\omega$ . This could be a result of intercorrelations of molecular interactions and ISD or  
123 RSA (Table S1), since disordered regions and exposed regions are often responsible for  
124 interacting with other molecules (Keskin et al., 2008; Van Der Lee et al., 2014). These results  
125 highlight the non-neglected contributions of functional constraints, including molecular  
126 interactions, on the adaptive evolution of protein-coding sequence.

127 **Complex correlations of protein length and male expression level with protein**  
128 **evolutionary rates.** To better clarify and visualize the correlations of  $\omega$ ,  $\omega_a$ , and  $\omega_{na}$  with gene  
129 properties in a refined fashion, we divided *D. melanogaster* genes into 15 groups according to  
130 the ascending orders of  $\omega$  values and compared these properties of different gene groups, while  
131 ensuring that each gene group contains the same total number of amino acids (Figure S2).  
132 Overall, for most of the properties being investigated, we observed similar correlations as shown  
133 in Table 1. For example, fast evolving genes are relatively young, short, lowly expressed, male  
134 or tissue specific, abundant of disordered, exposed residues, excluded in protein-protein  
135 network center hubs, and abundant of protein and DNA binding sites.

136 In contrast to previous observations, we found complex (nonlinear) correlations of  $\omega$   
137 gene groups with protein length and gene expression levels (Figure 1). For protein length, our  
138 Pearson correlation analysis (Table 1) and a number of previous studies have suggested a  
139 strong negative correlations with  $\omega$  (Lipman et al., 2002; Moutinho et al., 2019). However, we  
140 observed that some proteins with the slowest evolutionary rates, i.e. with the smallest  $\omega$  values,  
141 are significantly shorter than other gene groups with intermediate evolutionary rates (Fig. 1A).  
142 These include highly conserved genes such as eIF1A ( $\omega=0.0001$ , 148 a.a), rala ( $\omega=0.0001$ , 201  
143 a.a.), ctp ( $\omega=0.0001$ , 89 a.a.), and Mlc-c ( $\omega=0.0001$ , 153 a.a.).

144 Similar complex correlations were also observed in male expression level and mean  
145 expression level (Fig. 1BC). We found that, when checking male expression level and mean  
146 expression level, the gene group that shows the largest mean  $\omega$  has higher expression than  
147 those with intermediate  $\omega$ . Such U-shape correlations were not observed in female expression  
148 levels. Although protein length and mean expression levels of genes are known to be strongly  
149 correlated with protein evolutionary rates as listed in Table 1 and also in other references  
150 (Drummond et al., 2005; Lipman et al., 2002; Zhang and Yang, 2015), fast evolving genes can  
151 also be moderately or highly expressed, especially in male *D. melanogaster*. For example,  
152 many seminal fluid proteins show high  $\omega$  values and are highly expressed, such as Sfp60F  
153 ( $\omega=0.77$ , 82 a.a.), EbpII ( $\omega=0.68$ , 66 a.a), Acp36DE ( $\omega=0.68$ , 912 a.a.), and Dup99B ( $\omega=0.63$ ,  
154 54 a.a.). These proteins evolve at very fast rates (Begun and Lindfors, 2005; Swanson et al.,  
155 2001), contain various range of amino acids (54 in Dup99B to 912 in Acp36DE), and are  
156 moderately or highly expressed in male *D. melanogaster* (TPM ranging from 440 for Acp36DE  
157 to 3189 for Sfp60F), lowly expressed in female (TPM all around 1, presumably in spermatheca).  
158 We listed all the genes and protein length and expression levels in each  $\omega$  gene group, which  
159 can be found in supplementary file S2.

160 Since tissue specificity and male specificity both strongly correlates with  $\omega$ ,  $\omega_a$ , and  $\omega_{na}$   
161 (Table 1), we asked whether male specificity would be a redundant property compared to tissue  
162 specificity to indicate protein evolution due to the complex correlations of male expression  
163 levels. To answer this question, we classified *D. melanogaster* genes into 15 groups according  
164 to ascending values of male specificity. We then did similar classification to classify all the  
165 genes into 15 groups according to ascending values of tissue specificity (Figure 1). As  
166 expected, we found that tissue specificity positively correlates with  $\omega$ ,  $\omega_a$  and  $\omega_{na}$  (Fig. S3).  
167 However, we observed complex correlations for male specificity gene groups. Specifically, gene  
168 group with the lowest male specificity show significantly higher  $\omega$ ,  $\omega_a$  and  $\omega_{na}$  than its following

169 gene group (Fig. S3). This could be a result of fast evolving female-biased genes (Yang et al.,  
170 2016) included in this gene group.

171

### 172 **Putative molecular interaction sites are hotspots for protein adaptive evolution**

173 Having established that molecular interactions positively correlates with the adaptation of  
174 protein sequence, we next investigate whether residues involved in molecular interactions are  
175 targets for adaptive evolution. To tackle this question, we predicted protein-protein interaction  
176 sites (PPI-sites) and DNA binding sites (DNA-sites) for each of *D. melanogaster* protein  
177 sequence (see Methods). In addition, we characterized allosteric residues as surface and  
178 interior critical residues with STRESS model (Clarke et al., 2016) for all the structural models.  
179 We also extracted putative binding sites from STRESS Monte Carlo (MC) simulations. We  
180 calculated  $\omega$ ,  $\omega_a$  and  $\omega_{na}$  for residues in each of the putative molecular interaction category.  
181 Strikingly, we observed that residues involved in protein-protein interactions, DNA binding and  
182 ligand binding exhibited higher rates of adaptive evolution compared to their corresponding null  
183 sites (Fig. 2A-C). In addition, allosteric residues at protein surface showed higher adaptation  
184 rates than allosteric residues at protein interior or residues that are not involved in ligand binding  
185 from STRESS simulations (Fig. 2C).

186 Since we observed significant positive intercorrelations between PPI and DNA binding  
187 with ISD and RSA (Table S1), we next asked whether the increase of  $\omega_a$  in protein-protein  
188 interactions sites or DNA binding sites was caused by the increase of disorder or site exposure.  
189 We calculated and compared  $\omega$ ,  $\omega_a$  and  $\omega_{na}$  for putative PPI and DNA binding sites with  
190 different levels of ISD or RSA. Remarkably, we found that  $\omega_a$  of these putative binding sites  
191 remains similar among different levels of ISD or RSA (Fig. S4, left column). The results suggest  
192 that putative PPI or DNA binding events in proteins can result in elevated adaptation rates  
193 regardless their structural disorder or site exposure. While for residues that are not associated  
194 with putative PPI or DNA binding, we also observed increase in  $\omega_a$  when increasing ISD or RSA  
195 (Fig. S4, right column), which could be the result of some other yet unknown underlying  
196 mechanisms or inaccuracy of putative binding sites predictions.

197 In order to gain better understanding of adaptation in molecular interaction sites, we  
198 further visualized positive selections that are associated with molecular interactions. We first  
199 investigated whether adaptive evolution is associated with particular protein structures or protein  
200 families. To do this, we looked into fast-adaptive proteins with the largest ~15% rates of  
201 adaptation ( $\omega_a > 0.15$ ) that are linked to high quality structural models. Interestingly, among  
202 these proteins, we found 45 enriched as trypsin-like cysteine/serine peptidase domain and 17



203 7TM chemoreceptors, suggesting widespread adaptive evolution acting on these protein  
204 families or protein domains in *D. melanogaster* (Table S2). Many of the 7TM chemoreceptors  
205 are olfactory and gustatory genes, which shows adaptive evolution in various species such as  
206 *Drosophila* and mosquito (Hill et al., 2002; Lawniczak and Begun, 2007; McBride, 2007; Wu et  
207 al., 2009). In addition to these two protein families, recurrent positive selections acting on some  
208 other fast-adaptive proteins were identified in previous studies in *Drosophila* and mammals, and  
209 the possible adaptive evolution mechanisms have been linked to exogenous ligand binding, for  
210 example, serine protease inhibitors (serpin), Toll-like receptor 4 (TLR-4), and cytochrome P450  
211 (Jiggins and Kim, 2007; Slodkowicz and Goldman, 2020).

212 In order to visualize the link between adaptive evolution and molecular interactions in the  
213 two protein families with frequent adaptive evolution, we showed significant positive selections  
214 and molecular interactions in two representatives: CG10232 and Or67a, each for trypsin-like  
215 cysteine/serine peptidase domain and 7TM chemoreceptors, respectively. We observed that in  
216 both cases, positively selected sites highly overlapped with predicted or inferred binding pockets  
217 (Fig. 2D-E). Specifically, in CG10232, we found clusters of positive selected sites around NAG  
218 binding sites that are inferred from a crystal structure of serine protease (PDB code: 2XXL) (Fig.  
219 2D), while in Or67a, positively selected sites expand around the putative odorant binding  
220 channel formed by helices S1-S6 in extracellular regions (Butterwick et al., 2018) (Fig. 2E).

221 Except for these examples that are associated with exogenous ligand or exogenous  
222 peptide binding, we also identified two previously not described examples where adaptive  
223 evolution might be linked to endogenous protein binding: Spatzle (spz, Fig. 2F) and Cul6 (Fig.  
224 2G). Spatzle can bind to Toll-like receptors (TLR) and trigger humoral innate immune response.  
225 We built the missing loop in Spatzle in the crystal structure of Toll/Spatzle complex (PDB code  
226 4BV4) according to the dimeric crystal structure of Spatzle (PDB code 3E07). In this complex  
227 structural model, we observed several positively selected sites in Toll-4/Spatzle interfaces (Fig.  
228 2F). Cul6, another example, is a protein in cullins family in *D. melanogaster*. The cullins protein  
229 family are known as scaffold proteins that assemble multi-subunit Cullin-RING E3 ubiquitin  
230 ligase by forming SCF complex with F box and RING-box (Rbx) proteins (Zheng et al., 2002).  
231 We constructed the putative Cul6 contained SCF complex by superimposition to the crystal  
232 structure of the Cul1-Rbx1-Skp1-F box<sup>Skp2</sup> SCF ubiquitin ligase complex (Zheng et al., 2002). In  
233 the structural model, we observed positive selected sites in Cul6 clustered around the binding  
234 sites of RING-box protein, Rbx1, and F-box protein, Skp1 (Fig. 2G).

235

236 **Frequent adaptive evolution and co-adaptative evolution in genes involved in**  
237 **reproduction, immune system and environmental information processing**

238 To find out whether specific biological functions were associated with fast-adaptive genes, we  
239 applied DAVID Go analysis with genes that have largest ~15% rates of adaptation ( $\omega_a > 0.15$ ).  
240 The significant Go terms are frequently linked to serine-type endopeptidase activity,  
241 reproduction, protein lysis, chemosensory and other related biological functions (Table S3). As  
242 these fast-adaptive genes tend to be enriched in similar biological functions, we asked whether  
243 these genes are evolved co-adaptively, i.e., whether these proteins are interacting with each  
244 other frequently. To test this possibility, we obtained PPI of *D. melanogaster* from STRING  
245 database (Szklarczyk et al., 2019) and analyzed protein-protein interactions among fast-  
246 adaptive proteins. We found that fast-adaptive proteins tend to interact with each other more  
247 frequently than expected (PPI enrichment p-value  $< 1.0e-16$ ). In the PPI network of fast-  
248 adaptive proteins, we observed 7 strongly connected sub-clusters with at least 5 members (Fig.  
249 3A, Table S4). Proteins in these sub-clusters are enriched in biological processes such as  
250 reproduction, immune response, defense response to bacterium and virus, RNA interference,  
251 chitin metabolic, etc., which are in line with the Go analysis of fast-adaptive genes (Table S5-  
252 S10).

253 We next asked whether co-adaptation plays a role in the adaptive evolution of interacting  
254 proteins to a broader extend, including both fast- and slow-adaptive proteins. To address this  
255 question, we analyzed and compared adaptation rates of all PPIs available in STRING database  
256 with high confidence in *D. melanogaster* and we found that protein partners of fast-adaptive  
257 proteins ( $\omega_a > 0.15$ ) have significantly larger maximum/average  $\omega_a$  compared to slow-adaptive  
258 proteins (Figure 4). We further analyzed and visualized adaptive evolutionary rates of proteins in  
259 PPI networks of 9 different biological pathways extracted from KEGG pathways, including  
260 immune system, xenobiotics biodegradation, response to environment, aging and development,  
261 genetic information processing, sensory system, transport and catabolism, cell growth and  
262 death and metabolism. We observed that, in these PPI networks, proteins with relatively large  
263  $\omega_a$  tend to interact with each other (Figure S5A, S5B). We also noticed that, for pathways that  
264 are previously known as adaptation-hotspots, e.g., immune system, fast-adaptive proteins can  
265 act as central nodes and are co-adaptively evolved with other fast-adaptive proteins (Figure  
266 S5C). While in pathways such as transport and catabolism, fast-adaptive proteins are mainly at  
267 PPI periphery. In line with these findings, we found that  $\omega_a$  are larger in pathways that harbor  
268 fast-adaptive proteins as central nodes than other pathways (Figure S6).



269 **Physical interactions contribute to co-adaptation of fast-adaptive genes.** Having  
270 established that molecular interactions contribute to adaptive evolution of protein sequence, we  
271 then investigated whether these physical molecular interactions could drive protein-protein co-  
272 adaptation. To do this, we looked into interacting fast-adaptive protein pairs that are associated  
273 known or inferred complex structural models. For inferred complex structural models, we  
274 superimposed the structural models of the pair of proteins onto their high resolution homologous  
275 complex structures. Here we observed and illustrated co-adaptation at PPI interface in two  
276 examples: Toll-4/Spatzle and Spn28Db/CG18563 (Figure 3).

277 **Toll-4/Spatzle.** Toll-4 is a member of toll-like receptors. Previous studies have shown strong  
278 evidence of adaptive evolution of Toll-4 in *Drosophila* and mammals (Levin and Malik, 2017;  
279 Slodkowitz and Goldman, 2020). Toll-4 can bind to Spatzle and trigger further innate immune  
280 responses with high confidence (inferred from STRING database). In the previous section, we  
281 showed that several positively selected sites in Spatzle overlap with Toll-Spatzle interfaces.  
282 Here, we further showed that, in Toll-4, considerable number of significant positively selected  
283 sites were located at interface for Spatzle (Fig. 3B), which is in line with a previous study of Toll-  
284 4 in *D. willistoni* (Levin and Malik, 2017).

285 **Spn28Db/CG18563.** Spn28Db is one of the serine protease inhibitors in *D. melanogaster* that  
286 are expressed in male accessory glands, while CG18563 belongs to the protein family of  
287 trypsin-like cysteine/serine peptidase domain. The interactions between the two proteins were  
288 predicted with high confidence from STRING database, and the molecular interactions can be  
289 inferred from existing crystal structure of serpin and bacteria protease complex (PDB code  
290 1EZK). We observed many positive selected sites at the molecular interface between the two  
291 proteins (Fig. 3C), suggesting that physical interactions might play a role in the co-adaptation of  
292 the two proteins.

293

### 294 **Most clinally differentiated SNPs in protein-coding genes are adaptive**

295 To find out the relations between short-term adaptation to local environments and long-term  
296 adaptive evolution, we extracted residues with significant  $F_{ST}$  SNPs from clinal variations  
297 (Svetec et al., 2016). We then computed evolutionary rates ( $\omega$ ), adaptation rates ( $\omega_a$ ) and non-  
298 adaptation rates ( $\omega_{na}$ ) of these residues as in previous section. We observed that these  
299 residues have much higher ratio of adaptation rates over non-adaptation rates than genome-  
300 wide random expectations (Figure 5), suggesting that these residues have higher proportions of  
301 adaptive changes, and that they can be hotspots for adaptive evolution. To further characterize  
302 structural and functional properties of short-term genetic variations, we mapped significant

303 nonsynonymous  $F_{ST}$  residues to different structural and functional characteristics, such as ISD,  
304 RSA, PPI-sites, DNA-sites and ligand-binding sites. We found that these nonsynonymous SNPs  
305 follow the patterns of adaptive changes. For example, they were enriched disordered regions  
306 and protein surfaces, and were significantly more likely to be involved in protein-protein  
307 interactions and ligand-binding than expectation (Table S11-S15).

308

## 309 **Discussion**

310 In this study, we systematically studied the impact of structure- and function-related gene  
311 properties on protein sequence evolution and adaptation in *D. melanogaster* genome. We found  
312 that, compared to protein structure-related properties, such as intrinsic structural disorder (ISD)  
313 and relative solvent accessibility (RSA), function-related properties, such as tissue specificity  
314 and male specificity, contribute more extensively to protein sequence adaptive evolution.  
315 Especially, we noticed that molecular interactions in proteins contribute to the variation of  
316 protein sequence adaptive evolution. In line with this result, we detected that molecular  
317 interaction sites are hotspots for adaptative evolution. We confirmed that proteins that are fast  
318 adaptive are enriched in GO terms that are associated reproduction, immunity and  
319 environmental information processing. Furthermore, we revealed that fast-adaptive proteins  
320 tend to interact with each other frequently and protein partners of these fast-adaptive proteins  
321 tend to have higher adaptation rates, suggesting that co-adaptive evolution might be common in  
322 *D. melanogaster*. By looking at interacting fast-adaptive proteins, we further demonstrated that  
323 physical interactions may contribute to the mechanisms of co-adaptative evolution of fast-  
324 adaptive proteins.

325 Extensive studies have been conducted to uncover the main drivers that govern protein  
326 sequence evolutionary rate (Zhang and Yang, 2015). Gene expression level was proved to be a  
327 major determinant (Zhang and Yang, 2015) through mechanisms such as the pressure for  
328 translational robustness, i.e., robustness to translational missense errors (Drummond et al.,  
329 2005). Here, we showed that caveat exists when we looked at gene expression levels in male  
330 *D. melanogaster*. Previous studies have revealed that male biased or female biased genes can  
331 be fast evolving (Yang et al., 2016). On the other hand, many male biased genes can be highly  
332 expressed in testis, which results in a complex correlation between protein sequence  
333 evolutionary rate and male expression level or even mean expression level of *D. melanogaster*.  
334 The unique evolutionary property of these male biased or specific genes could be caused by the  
335 unique transcriptional scanning mechanism in testis (Xia et al., 2020). We propose that tissue  
336 specificity might be a better quantity when considering the impact of gene expression profile on

337 protein sequence evolution in *D. melanogaster*. In addition to male expression level, a similar  
338 complex correlation was observed for protein length. It has been the notion that short proteins  
339 tend to evolve faster than long proteins, which may be biologically relevant or byproduct of other  
340 factors such as selection on buried and exposed sites (Moutinho et al., 2019). Here, we  
341 demonstrated that, in *D. melanogaster*, although protein length is strongly negatively correlated  
342 with protein sequence evolutionary rate, genes that have the slowest evolutionary rates tend to  
343 be relatively short. This could be caused by the fact that under essential functional constraint,  
344 genes can undergo strong purifying selections, while essential genes such as secreted proteins  
345 are constrained to be smaller, and that essential genes could be shorter than other genes (Chen  
346 et al., 2020).

347 It has been recognized that protein surface and intrinsic disorder regions are frequent  
348 targets for adaptive evolution and contribute to the variations of protein sequence adaptive  
349 evolution (Afanasyeva et al., 2018; Moutinho et al., 2019). However, the detailed mechanisms  
350 underlying these observations remains unclear. One possible explanation would be that these  
351 regions are frequently linked to intermolecular interactions (Afanasyeva et al., 2018; Moutinho et  
352 al., 2019). For example, Moutinho et al hypothesized that molecular interactions involved in  
353 host-pathogen coevolution were the major driver of protein adaptation (Moutinho et al., 2019).  
354 Here, we further identified that proportions of possible molecular interaction sites inside proteins  
355 contribute to the variations of protein sequence adaptive evolution and that these molecular  
356 interaction sites or regulatory sites at protein surface can be hotspots of protein adaptation.  
357 Indeed, some specific molecular interactions have been linked to adaptive evolution in several  
358 case studies (Bachtrog, 2008; Hughes and Nei, 1988; Levin and Malik, 2017; Schott et al.,  
359 2014) and large-scale studies based on proteins with high quality structural models (Slodkowicz  
360 and Goldman, 2020). In the latter study, the authors showed that positive selections in  
361 mammals tend to cluster closer to binding sites of exogenous ligands than expected by chance  
362 (Slodkowicz and Goldman, 2020), suggesting an important role of function important regions in  
363 adaptive evolution. Here, we extend the conclusion to *D. melanogaster* genome, including  
364 proteins with or without high resolution structural models. We also showed that except for  
365 exogenous ligands, endogenous ligands might also contribution to adaptive evolution, while the  
366 latter might explain why interacting proteins tend to evolve co-adaptively.

367 Notably, previous studies have revealed that multi-interface proteins tend to be evolving  
368 more slowly than single-interface proteins (Kim et al., 2006), which seems to be contradictory to  
369 our results that proteins with more interaction sites evolve faster and have faster adaptation  
370 rates. Here, we argue that, in our study, we used sequence profile to predict molecular

371 interaction sites in proteins at a genomic scale, rather than only looking into proteins with high  
372 resolution structures. In this way, we may capture many weak or transient interactions, which  
373 are thought to be evolving faster than obligate and conserved interactions (Mintseris and Weng,  
374 2005). Meanwhile, we did not exclude intrinsic disordered regions (IDR) or intrinsic disordered  
375 proteins (IDP) in our study, which are widespread in *D. melanogaster* genome. It has been  
376 suggested that IDR/IDP tend to evolve fast due to lack of structural restraints (Echave et al.,  
377 2016). In the functional aspect, IDR/IDP are thought to be promiscuous binders through many  
378 multiple binding mechanisms, including forming static, semi-static, and fuzzy or dynamic  
379 complexes (Uversky, 2019), suggesting that the evolution of IDR/IDP cannot be explained  
380 merely by the lack of structural restraints. Actually, IDP and IDR in human genome were found  
381 to be undergoing extensive adaptive evolution (Afanasyeva et al., 2018). At last, it has been  
382 recognized that, except for allosteric regulations, encounter complexes (Gabdouline and Wade,  
383 1999) might also play an important role in mediating intermolecular interactions, such as  
384 protein-protein association (Tang et al., 2006) and protein-ligand binding (Re et al., 2019). Since  
385 encounter residues that are responsible for encounter complexes do not reside in conserved  
386 binding interfaces, these residues could be under relaxed purifying selections or even positive  
387 selections, which could be another yet-to-identify mechanism that contribute to protein  
388 sequence adaptive evolution.

389 In consistent with previous studies in *D. melanogaster* (Begun and Lindfors, 2005;  
390 Begun and Whitley, 2000; Lazzaro et al., 2004), we showed that fast-adaptive proteins are  
391 enriched in molecular functions such as reproduction, immunity and environmental information  
392 processing. We further demonstrated that fast-adaptive proteins tend to interact with each other  
393 more frequently than random expectations, suggesting co-adaptation might be common among  
394 fast-adaptive proteins. Mechanisms that contribute to the co-adaptation could be: (1) interacting  
395 fast-adaptive proteins are often enriched in similar molecular functions and under similar  
396 selective pressure; (2) interacting fast-adaptive undergo co-evolution through physical  
397 interactions. In this study we showed two examples that adaptive evolution could occur at  
398 protein-protein interface, which suggest that physical interactions could contribute to the co-  
399 adaptation of fast-adaptive proteins in *D. melanogaster*. Moreover, we showed that co-  
400 adaptation might exist to a broader extend rather than only among fast-adaptive proteins.  
401 Specifically, proteins that interact with fast-adaptive proteins tend to have higher adaptation  
402 rates. Since molecular interactions contribute to adaptive evolution, it is reasonable to  
403 hypothesize that co-adaptation at this broader extend could be regulated by these interactions.  
404 Actually, it has been suggested that interacting proteins tend to have similar evolutionary rates

405 and the possible mechanism would be the co-evolution of physical interactions (Pazos and  
406 Valencia, 2008).

407 It has been suggested that populations in different local environments can have genetic  
408 variances that result in local adaptations. In this study, we found that loci with significant genetic  
409 variance among populations harbor higher proportions of long-term adaptive changes and these  
410 loci follow similar patterns as adaptive changes, i.e. they are enriched in disordered regions,  
411 protein surfaces, and functionally important regions. These results suggest that population  
412 differentiation of protein-coding genes can be an important basis for long-term adaptive  
413 evolution. Importantly, our results indicate that most of the clinal amino-acid changes are  
414 adaptive, suggesting that non-selective forces play a non-essential role in the SNPs that show  
415 strong geographic differences. Our results also support a large effect of spatially varying  
416 selection on protein sequence and structures (Storz and Kelly, 2008).

417 It should be noted that studies at the genomic scale that aim to uncover the function- or  
418 structure-related constraints imposed on protein sequence evolution and adaptation share  
419 similar limitations that for most of the proteins or residues, structural or functional information  
420 would be incomplete or even missing. Thus, in this study, we used highly accurate neural-  
421 network based tools to predict molecular interactions, secondary structures, intrinsic structural  
422 disorder, relative solvent accessibility for each of the protein in *D. melanogaster* genome. In this  
423 way we were able to identify key factors that impact protein sequence evolution and adaptation  
424 in a less accurate but rather systematic fashion. We hope that with the availability of more and  
425 more curated structural, functional information and complex structural models of proteins in the  
426 near future, we will be able to uncover the precise role of molecular interactions in protein  
427 sequence adaptive evolution.

428

## 429 **Material and Methods**

430  **$d_N/d_S$  ratio ( $\omega$ ).** We used a maximum likelihood method to infer  $d_N/d_S$  ratio ( $\omega$ ) of *D.*  
431 *melanogaster* protein-coding genes using the genome sequences of five species in  
432 *melanogaster* subgroup (*D. melanogaster*, *D. simulans*, *D. sechellia*, *D. yakuba*, and *D. erecta*).  
433 The protein-coding sequences were extracted from the alignments of 26 insects, which were  
434 obtained from UCSC Genome Browser (<http://hgdownload.soe.ucsc.edu/downloads.html>). The  
435 sequences were further processed by GeneWise (Birney et al., 2004) to remove possible  
436 insertions and deletions using the longest isoforms of the corresponding *D. melanogaster*  
437 protein sequences as references (FlyBase version r6.15) (Thurmond et al., 2019). The  
438 processed sequences were then realigned by PRANK -codon function (Löytynoja, 2014). We

439 used codeml in PAML (Yang, 2007) to compute gene-specific  $\omega$  using M0 model. We removed  
440 sequences that have more than 15% of their nucleotides not aligned (gaps) to *D. melanogaster*  
441 genes in more than 2 species. To further avoid numeric errors and ensure reasonable  
442 estimations, we only retained relatively divergent sequences that are: (1) divergent with dS  
443 larger than 0.3, (2) less divergent with dS larger than 0.1 and dN smaller than 0.001 ( $dS \gg dN$ ).  
444 At last, there were 12118 genes in total passed all the criteria and were assigned gene specific  
445  $\omega$ , containing 6,538,872 amino acids. We also calculated site-specific  $\omega$  by using likelihood ratio  
446 tests (LRT) comparing M7 model against M8 model (Yang et al., 2005).

447 **Rate of adaptive and nonadaptive changes.** We recalled all SNPs of 205 inbred lines from  
448 the Drosophila Genetic Reference Panel (DGRP), Freeze 2.0 (Huang et al., 2014)  
449 (<http://dgrp2.gnets.ncsu.edu>). We then generated 410 alternative genomes using all monoallelic  
450 and bi-allelic SNP data sets. We extracted the coding sequences of *D. melanogaster* genes  
451 from the generated alternative genomes, removed all possible insertions and deletions using  
452 GeneWise (Birney et al., 2004) as described above. We then align all the coding sequences to  
453 their corresponding aligned CDS sequences using PRANK -codon function (Löytynoja, 2014).  
454 We removed polymorphisms segregating at frequencies smaller than 5% to reduce possible  
455 slightly deleterious mutations (Charlesworth and Eyre-Walker, 2008). In order to avoid possible  
456 effects of low divergence between *D. simulans* and *D. melanogaster* (Keightley and Eyre-  
457 Walker, 2012), we used *D. yakuba* as outgroup to estimate nonsynonymous polymorphisms  
458 (Pn), synonymous polymorphisms (Ps), nonsynonymous substitutions (Dn) and synonymous  
459 substitutions (Ds) by MK.pl (Begun et al., 2007; Langley et al., 2012). Similar as Begun et al.  
460 (Begun et al., 2007), we only analyzed genes with at least six variants for each of substitutions,  
461 polymorphisms, nonsynonymous changes and synonymous changes. We used an extension of  
462 MK test, asymptotic MK (Messer and Petrov, 2013; Uricchio et al., 2019), to estimate the  
463 proportions of adaptive changes ( $\alpha$ ). The rate of adaptive changes ( $\omega_a$ ) was then calculated as  
464  $\omega_a = \omega\alpha$  and the rate of non-adaptive changes as  $\omega_{na} = \omega - \omega_a$ . Details of the asymptotic MK  
465 test were as following:

466 (1) Classical McDonald–Kreitman test. According to Smith and Eyre-Walker (Smith and Eyre-  
467 Walker, 2002), the proportions of adaptive changes for protein-coding genes can be calculated  
468 as following:

469 
$$\alpha = 1 - \frac{DsPn}{DnPs}$$

470 According to this equation, we could estimate the proportion of adaptive changes and carried  
471 out classical MK test by applying Fisher's exact test.



472 (2) Asymptotic estimation of  $\alpha$ . A known problem of the classical estimation of  $\alpha$  above is the  
473 accumulation of slightly deleterious mutations at low frequencies. We therefore used an  
474 extension of MK test, asymptotic MK test approach (Messer and Petrov, 2013) to estimate the  
475 proportions of adaptive changes. As in original aMK, we defined  $\alpha(x)$  as a function of derived  
476 allele frequency ( $x$ ):

$$477 \quad \alpha(x) = 1 - \frac{DsPn(x)}{DnPs(x)}$$

478 where  $Pn(x)$  and  $Ps(x)$  are number of non-synonymous and synonymous polymorphisms at  
479 frequency  $x$ , respectively. However, the original approach may suffer from numeric errors when  
480 there were very few polymorphic sites, which is quite common in many of *D. melanogaster*  
481 genes. To make the estimations more robust while preserving the same asymptote, we further  
482 define  $Pn(x)$  and  $Ps(x)$  as total number of  $Pn$  and  $Ps$  above frequency  $x$  as described in  
483 Uricchio et al (Uricchio et al., 2019). We fitted  $\alpha(x)$  to an exponential curve of  $\alpha(x) \approx \exp(-bx)+c$   
484 using Imfit (Newville et al., 2014) and determined the asymptotic value of  $\alpha$  at the limit of  $x$ , 1.0.  
485 We then estimate the rate of adaptive changes ( $\omega_a$ ) as

$$486 \quad \omega_a = \frac{N_a/L_N}{dS} = \frac{dN_a}{dS} = \frac{dN_a}{dN} \cdot \frac{dN}{dS} = \alpha\omega$$

487 where  $N_a$  is the number of adaptive changes and  $dN_a=N_a/L_N$  is the number of adaptive changes  
488 per nonsynonymous site. Finally, we calculated the rate of nonadaptive changes ( $\omega_{na}$ ) as  
489  $\omega_{na}=\omega-\omega_a$ . The final dataset contains 7192 protein-coding genes, with smallest  $\omega_a$  being 0.00  
490 and largest being 1.29.

491 **Structure-/function- related properties of *D. melanogaster* proteins.** We obtained function-  
492 related properties mentioned in main text as following. We derived *D. melanogaster* gene ages  
493 (Kondo et al., 2017; Zhang et al., 2010) for genes that are specific to *Drosophila*, and from  
494 GenTree (Shao et al., 2019) for genes that are beyond *Drosophila* clade. We then assigned a  
495 pseudo-age to each of the genes. Specifically, there are 11 age groups from “cellular  
496 organisms”, assigning to a pseudo age value of 0, to “melanogaster”, assigning a pseudo age  
497 value of 10. We downloaded *D. melanogaster* protein-protein interaction (PPI) from STRING  
498 database (Szklarczyk et al., 2019). A cut-off of combined score larger than 0.7 was used to  
499 retain high confident PPI for further analysis. We then used BSpred (Mukherjee and Zhang,  
500 2011) to predict protein-protein interaction (PPI) sites and DRNApred (Yan and Kurgan, 2017)  
501 to predict DNA binding sites. For each protein, we calculated ratios of protein interaction  
502 residues (PPI-site ratio) and ratios of DNA binding residues (DNA-site ratio) by dividing total  
503 predicted protein interaction sites and DNA binding sites over protein length, respectively. For  
504 structure-related properties, we used DeepCNF (Wang et al., 2016) to predict these properties

505 for each gene, including three-state secondary structures (helix, sheet, and coil), structural  
506 disorder, relative solvent accessibility (RSA). Further, we calculated the ratios of helix, sheet,  
507 helix+sheet, and coil residues of each gene from predicted secondary structures. For each  
508 gene, we computed intrinsic structural disorder (ISD) and relative solvent accessibility (RSA), as  
509 protein-length normalized summations of the probabilities of each residue being disorder and  
510 exposed, respectively.

511 **Gene expression patterns.** We downloaded gene expression profile from FlyAtlas2 (Leader et  
512 al., 2018). We converted FPKM to TPM by normalizing FPKM against the summation of all  
513 FPKMs as following:

$$514 \quad \text{TPM}_i = \frac{\text{FPKM}_i}{\sum \text{FPKM}_j} \times 10^6$$

515 After TPM conversion, we only retained genes with expression level larger than 0.1 TPM for  
516 further analysis. We treated male and female whole-body TPM as male and female expression  
517 levels. We calculated mean expression level by averaging male and female TPM. We used  
518 following Z-score to describe male specificities of *D. melanogaster* genes:

$$519 \quad zscore = \frac{\text{TPM}(\text{male expression}) - \text{TPM}(\text{female expression})}{\sqrt{sd^2(\text{male expression}) + sd^2(\text{female expression})}}$$

520 We calculated tissue specificities of genes using tau values (Yanai et al., 2005) based on the  
521 expression profiles of 27 different tissues.

522 **High quality 3D structures of *D. melanogaster* proteins.** We downloaded high-quality  
523 structures or structural models of *D. melanogaster* proteins from protein data bank (PDB)  
524 (Burley et al., 2019), SWISS-MODEL Repository (Bienert et al., 2017), and MODBASE (Pieper  
525 et al., 2011), with descending priorities. For example, if there were 3D structures of a same  
526 protein or protein region in multiple databases, we first considered high-resolution structures  
527 from PDB; if no structures were found in PDB, we then considered SWISS-MODEL Repository;  
528 and at last from MODBASE. In addition, we used blastp (Camacho et al., 2009) to search  
529 homologs of each *D. melanogaster* protein against all PDB sequences with E-value threshold of  
530 0.001. We further carried out comparative structural modeling using RosettaCM (Song et al.,  
531 2013) to model high-quality structural models of proteins or protein regions that were not  
532 available in PDB, SWISS-MODEL Repository and MODBASE. For each RosettaCM simulation,  
533 we used no more than 5 most significant hits from blastp search. For proteins that are in  
534 complex forms, we only extracted monomers for further analysis. At last, we obtained 14543

535 high quality structural models, corresponding to 11284 genes. These structural models contain  
536 2,691,913 unique amino acids, 41.2% of all the residues in genes that were assigned  $\omega$ .  
537 **Evolutionary rates of different structural/functional sites.** We classified amino acids into  
538 different classes of structural/functional properties. Specifically, we classified three classes for  
539 both ISD and RSA according the probability of residues being disordered or exposed: ordered  
540 or buried (0.00 to 0.33), medium (0.33 to 0.67), disordered or exposed (0.67 to 1.00). For both  
541 PPI and DNA binding, we classified two classes: PPI-site or DNA-site (binding sites), None-PPI  
542 or None-DNA (corresponding null sites for PPI or DNA binding). For residues that have 3D  
543 structures, we used STRESS (Clarke et al., 2016) to predict putative ligand binding sites and  
544 allosteric sites from all the high-quality structures or structural models. The allosteric sites were  
545 further classified as surface critical or interior critical according to their locations. We then  
546 classified these residues into four groups: LIG (ligand binding sites), Surf. Crit. (surface critical  
547 sites), Interior Crit. (interior critical sites) and Others (other sites). For each of the site classes,  
548 we randomly sampled 1,00 sequences, each containing 10,000 amino acids. We computed  $\omega$ ,  
549  $\omega_a$ , and  $\omega_{na}$  for the randomly sampled sequences similar as the steps described in the above  
550 sections.

551

## 552 Reference

- 553 Afanasyeva, A., Bockwoldt, M., Cooney, C.R., Heiland, I., and Gossmann, T.I. (2018). Human  
554 long intrinsically disordered protein regions are frequent targets of positive selection. *Genome*  
555 *Res.* 28, 975–982.
- 556 Bachtrog, D. (2008). Positive selection at the binding sites of the male-specific lethal complex  
557 involved in dosage compensation in *Drosophila*. *Genetics* 180, 1123–1129.
- 558 Begun, D.J., and Lindfors, H.A. (2005). Rapid evolution of genomic Acp complement in the  
559 melanogaster subgroup of *Drosophila*. *Mol. Biol. Evol.* 22, 2010–2021.
- 560 Begun, D.J., and Whitley, P. (2000). Adaptive evolution of relish, a *Drosophila* NF-  
561 kappaB/IkappaB protein. *Genetics* 154, 1231–1238.
- 562 Begun, D.J., Holloway, A.K., Stevens, K., Hillier, L.W., Poh, Y.-P., Hahn, M.W., Nista, P.M.,  
563 Jones, C.D., Kern, A.D., Dewey, C.N., et al. (2007). Population Genomics: Whole-Genome  
564 Analysis of Polymorphism and Divergence in *Drosophila simulans*. *PLoS Biol.* 5, e310.
- 565 Bienert, S., Waterhouse, A., De Beer, T.A.P., Tauriello, G., Studer, G., Bordoli, L., and  
566 Schwede, T. (2017). The SWISS-MODEL Repository-new features and functionality. *Nucleic*  
567 *Acids Res.* 45, D313–D319.
- 568 Birney, E., Clamp, M., and Durbin, R. (2004). GeneWise and Genomewise. *Genome Res.* 14,

569 988–995.

570 Burley, S.K., Berman, H.M., Bhikadiya, C., Bi, C., Chen, L., Di Costanzo, L., Christie, C.,  
571 Dalenberg, K., Duarte, J.M., Dutta, S., et al. (2019). RCSB Protein Data Bank: Biological  
572 macromolecular structures enabling research and education in fundamental biology,  
573 biomedicine, biotechnology and energy. *Nucleic Acids Res.* *47*, D464–D474.

574 Butterwick, J.A., del Mármol, J., Kim, K.H., Kahlson, M.A., Rogow, J.A., Walz, T., and Ruta, V.  
575 (2018). Cryo-EM structure of the insect olfactory receptor Orco. *Nature* *560*, 447–452.

576 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden,  
577 T.L. (2009). BLAST+: architecture and applications. *BMC Bioinformatics* *10*, 421.

578 Charlesworth, J., and Eyre-Walker, A. (2008). The McDonald-Kreitman test and slightly  
579 deleterious mutations. *Mol. Biol. Evol.* *25*, 1007–1015.

580 Chen, H., Zhang, Z., Jiang, S., Li, R., Li, W., Zhao, C., Hong, H., Huang, X., Li, H., and Bo, X.  
581 (2020). New insights on human essential genes based on integrated analysis and the  
582 construction of the HEGIAP web-based platform. *Brief. Bioinform.* *21*, 1397–1410.

583 Clarke, D., Sethi, A., Li, S., Kumar, S., Chang, R.W.F., Chen, J., and Gerstein, M. (2016).  
584 Identifying Allosteric Hotspots with Dynamics: Application to Inter- and Intra-species  
585 Conservation. *Structure* *24*, 826–837.

586 Drummond, D.A., Bloom, J.D., Adami, C., Wilke, C.O., and Arnold, F.H. (2005). Why highly  
587 expressed proteins evolve slowly. *Proc. Natl. Acad. Sci. U. S. A.* *102*, 14338–14343.

588 Echave, J., Spielman, S.J., and Wilke, C.O. (2016). Causes of evolutionary rate variation among  
589 protein sites. *Nat. Rev. Genet.* *17*, 109–121.

590 Enard, D., Cai, L., Gwennap, C., and Petrov, D.A. (2016). Viruses are a dominant driver of  
591 protein adaptation in mammals. *Elife* *5*, e12469.

592 Gabdoulline, R.R., and Wade, R.C. (1999). On the protein-protein diffusional encounter  
593 complex. *J. Mol. Recognit.* *12*, 226–234.

594 Goldman, N., Thorne, J.L., and Jones, D.T. (1998). Assessing the impact of secondary structure  
595 and solvent accessibility on protein evolution. *Genetics* *149*, 445–458.

596 Hill, C.A., Fox, A.N., Pitts, R.J., Kent, L.B., Tan, P.L., Chrystal, M.A., Cravchik, A., Collins, F.H.,  
597 Robertson, H.M., and Zwiebel, L.J. (2002). G Protein-Coupled Receptors in *Anopheles*  
598 *gambiae*. *Science* (80-. ). *298*, 176–178.

599 Huang, W., Massouras, A., Inoue, Y., Peiffer, J., Ràmia, M., Tarone, A.M., Turlapati, L., Zichner,  
600 T., Zhu, D., Lyman, R.F., et al. (2014). Natural variation in genome architecture among 205  
601 *Drosophila melanogaster* Genetic Reference Panel lines. *Genome Res.* *24*, 1193–1208.

602 Hughes, A.L., and Nei, M. (1988). Pattern of nucleotide substitution at major histocompatibility

603 complex class I loci reveals overdominant selection. *Nature* 335, 167–170.

604 Jiggins, F.M., and Kim, K.W. (2007). A screen for immunity genes evolving under positive  
605 selection in *Drosophila*. *J. Evol. Biol.* 20, 965–970.

606 Keightley, P.D., and Eyre-Walker, A. (2012). Estimating the rate of adaptive molecular evolution  
607 when the evolutionary divergence between species is small. *J. Mol. Evol.* 74, 61–68.

608 Keskin, O., Gursoy, A., Ma, B., and Nussinov, R. (2008). Principles of protein-protein  
609 interactions: What are the preferred ways for proteins to interact? *Chem. Rev.* 108, 1225–1244.

610 Kim, P.M., Lu, L.J., Xia, Y., and Gerstein, M.B. (2006). Relating three-dimensional structures to  
611 protein networks provides evolutionary insights. *Science* (80-. ). 314, 1938–1941.

612 Kondo, S., Vedanayagam, J., Mohammed, J., Eizadshenass, S., Kan, L., Pang, N., Aradhya, R.,  
613 Siepel, A., Steinhauer, J., and Lai, E.C. (2017). New genes often acquire male- specific  
614 functions but rarely become essential in *Drosophila*. 1841–1846.

615 Kosiol, C., Vinař, T., da Fonseca, R.R., Hubisz, M.J., Bustamante, C.D., Nielsen, R., and Siepel,  
616 A. (2008). Patterns of Positive Selection in Six Mammalian Genomes. *PLoS Genet.* 4,  
617 e1000144.

618 Langley, C.H., Stevens, K., Cardeno, C., Lee, Y.C.G., Schridder, D.R., Pool, J.E., Langley, S.A.,  
619 Suarez, C., Corbett-Detig, R.B., Kolaczkowski, B., et al. (2012). Genomic variation in natural  
620 populations of *Drosophila melanogaster*. *Genetics* 192, 533–598.

621 Lawniczak, M.K.N., and Begun, D.J. (2007). Molecular population genetics of female-expressed  
622 mating-induced serine proteases in *Drosophila melanogaster*. *Mol. Biol. Evol.* 24, 1944–1951.

623 Lazzaro, B.P., Scurman, B.K., and Clark, A.G. (2004). Genetic basis of natural variation in *D.*  
624 *melanogaster* antibacterial immunity. *Science* 303, 1873–1876.

625 Leader, D.P., Krause, S.A., Pandit, A., Davies, S.A., and Dow, J.A.T. (2018). FlyAtlas 2: A new  
626 version of the *Drosophila melanogaster* expression atlas with RNA-Seq, miRNA-Seq and sex-  
627 specific data. *Nucleic Acids Res.* 46, D809–D815.

628 Van Der Lee, R., Buljan, M., Lang, B., Weatheritt, R.J., Daughdrill, G.W., Dunker, A.K.,  
629 Fuxreiter, M., Gough, J., Gsponer, J., Jones, D.T., et al. (2014). Classification of intrinsically  
630 disordered regions and proteins. *Chem. Rev.* 114, 6589–6631.

631 Levin, T.C., and Malik, H.S. (2017). Rapidly evolving Toll-3/4 genes encode male-specific Toll-  
632 like receptors in *drosophila*. *Mol. Biol. Evol.* 34, 2307–2323.

633 Lin, Y.S., Hsu, W.L., Hwang, J.K., and Li, W.H. (2007). Proportion of solvent-exposed amino  
634 acids in a protein and rate of protein evolution. *Mol. Biol. Evol.* 24, 1005–1011.

635 Lindblad-Toh, K., Garber, M., Zuk, O., Lin, M.F., Parker, B.J., Washietl, S., Kheradpour, P.,  
636 Ernst, J., Jordan, G., Mauceli, E., et al. (2011). A high-resolution map of human evolutionary

- 637 constraint using 29 mammals. *Nature* 478, 476–482.
- 638 Lipman, D.J., Souvorov, A., Koonin, E. V., Panchenko, A.R., and Tatusova, T.A. (2002). The  
639 relationship of protein conservation and sequence length. *BMC Evol. Biol.* 2, 20.
- 640 Löytynoja, A. (2014). Phylogeny-aware alignment with PRANK. *Methods Mol. Biol.* 1079, 155–  
641 170.
- 642 McBride, C.S. (2007). Rapid evolution of smell and taste receptor genes during host  
643 specialization in *Drosophila sechellia*. *Proc. Natl. Acad. Sci. U. S. A.*
- 644 Messer, P.W., and Petrov, D.A. (2013). Frequent adaptation and the McDonald-Kreitman test.  
645 *Proc. Natl. Acad. Sci. U. S. A.* 110, 8615–8620.
- 646 Mintseris, J., and Weng, Z. (2005). Structure, function, and evolution of transient and obligate  
647 protein-protein interactions. *Proc. Natl. Acad. Sci. U. S. A.* 102, 10930–10935.
- 648 Moutinho, A.F., Trancoso, F.F., Dutheil, J.Y., and Zhang, J. (2019). The Impact of Protein  
649 Architecture on Adaptive Evolution. *Mol. Biol. Evol.* 36, 2013–2028.
- 650 Mukherjee, S., and Zhang, Y. (2011). Protein-protein complex structure predictions by  
651 multimeric threading and template recombination. *Structure* 19, 955–966.
- 652 Newville, M., Ingargiola, A., Stensitzki, T., and Allen, D.B. (2014). LMFIT: Non-Linear Least-  
653 Square Minimization and Curve-Fitting for Python. Zenodo.
- 654 Nielsen, R., Bustamante, C., Clark, A.G., Glanowski, S., Sackton, T.B., Hubisz, M.J., Fledel-  
655 Alon, A., Tanenbaum, D.M., Civello, D., White, T.J., et al. (2005). A scan for positively selected  
656 genes in the genomes of humans and chimpanzees. *PLoS Biol.* 3, e170.
- 657 Obbard, D.J., Welch, J.J., Kim, K.-W., and Jiggins, F.M. (2009). Quantifying Adaptive Evolution  
658 in the *Drosophila* Immune System. *PLoS Genet.* 5, e1000698.
- 659 Pál, C., Papp, B., and Hurst, L.D. (2001). Highly expressed genes in yeast evolve slowly.  
660 *Genetics* 158, 927–931.
- 661 Palmer, W.H., Hadfield, J.D., and Obbard, D.J. (2018). RNA-interference pathways display high  
662 rates of adaptive protein evolution in multiple invertebrates. *Genetics* 208, 1585–1599.
- 663 Pazos, F., and Valencia, A. (2008). Protein co-evolution, co-adaptation and interactions. *EMBO*  
664 *J.* 27, 2648–2655.
- 665 Pieper, U., Webb, B.M., Barkan, D.T., Schneidman-Duhovny, D., Schlessinger, A., Braberg, H.,  
666 Yang, Z., Meng, E.C., Pettersen, E.F., Huang, C.C., et al. (2011). ModBase, a database of  
667 annotated comparative protein structure models, and associated resources. *Nucleic Acids Res.*  
668 39, D465–D474.
- 669 Ramsey, D.C., Scherrer, M.P., Zhou, T., and Wilke, C.O. (2011). The relationship between  
670 relative solvent accessibility and evolutionary rate in protein evolution. *Genetics* 188, 479–488.



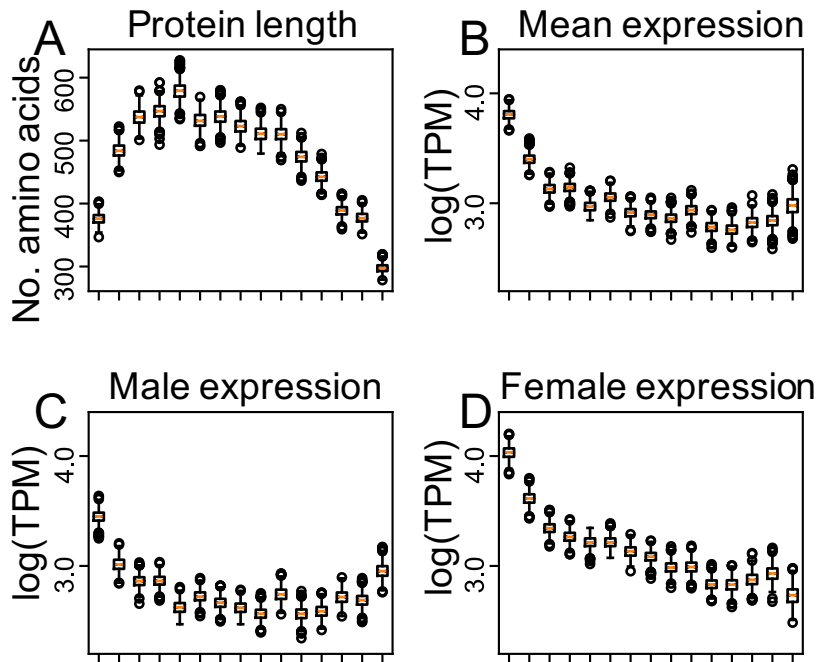
- 671 Re, S., Oshima, H., Kasahara, K., Kamiya, M., and Sugita, Y. (2019). Encounter complexes and  
672 hidden poses of kinaseinhibitor binding on the free-energy landscape. *Proc. Natl. Acad. Sci. U.*  
673 *S. A.* *116*, 18404–18409.
- 674 Sackton, T.B., Lazzaro, B.P., Schlenke, T.A., Evans, J.D., Hultmark, D., and Clark, A.G. (2007).  
675 Dynamic evolution of the innate immune system in *Drosophila*. *Nat. Genet.* *39*, 1461–1468.
- 676 Schott, R.K., Refvik, S.P., Hauser, F.E., López-Fernández, H., and Chang, B.S.W. (2014).  
677 Divergent positive selection in rhodopsin from lake and riverine cichlid fishes. *Mol. Biol. Evol.*  
678 *31*, 1149–1165.
- 679 Shao, Y., Chen, C., Shen, H., He, B.Z., Yu, D., Jiang, S., Zhao, S., Gao, Z., Zhu, Z., Chen, X.,  
680 et al. (2019). GenTree, an integrated resource for analyzing the evolution and function of  
681 primate-specific coding genes. *Genome Res.* *29*, 682–696.
- 682 Sironi, M., Cagliani, R., Forni, D., and Clerici, M. (2015). Evolutionary insights into host-  
683 pathogen interactions from mammalian sequence data. *Nat. Rev. Genet.* *16*, 224–236.
- 684 Slodkowitz, G., and Goldman, N. (2020). Integrated structural and evolutionary analysis reveals  
685 common mechanisms underlying adaptive evolution in mammals. *Proc. Natl. Acad. Sci.* *117*,  
686 5977–5986.
- 687 Smith, N.G.C., and Eyre-Walker, A. (2002). Adaptive protein evolution in *Drosophila*. *Nature.*
- 688 Song, Y., Dimaio, F., Wang, R.Y.R., Kim, D., Miles, C., Brunette, T., Thompson, J., and Baker,  
689 D. (2013). High-resolution comparative modeling with RosettaCM. *Structure* *21*, 1735–1742.
- 690 Storz, J.F., and Kelly, J.K. (2008). Effects of Spatially Varying Selection on Nucleotide Diversity  
691 and Linkage Disequilibrium: Insights From Deer Mouse Globin Genes. *Genetics* *180*, 367–379.
- 692 Svetec, N., Cridland, J.M., Zhao, L., and Begun, D.J. (2016). The Adaptive Significance of  
693 Natural Genetic Variation in the DNA Damage Response of *Drosophila melanogaster*. *PLoS*  
694 *Genet.* *12*, e1005869.
- 695 Swanson, W.J., Clark, a G., Waldrip-Dail, H.M., Wolfner, M.F., and Aquadro, C.F. (2001).  
696 Evolutionary EST analysis identifies rapidly evolving male reproductive proteins in *Drosophila*.  
697 *Proc. Natl. Acad. Sci. U. S. A.* *98*, 7375–7379.
- 698 Szklarczyk, D., Gable, A.L., Lyon, D., Junge, A., Wyder, S., Huerta-Cepas, J., Simonovic, M.,  
699 Doncheva, N.T., Morris, J.H., Bork, P., et al. (2019). STRING v11: Protein-protein association  
700 networks with increased coverage, supporting functional discovery in genome-wide  
701 experimental datasets. *Nucleic Acids Res.* *47*, D607–D613.
- 702 Tang, C., Iwahara, J., and Clore, G.M. (2006). Visualization of transient encounter complexes in  
703 protein-protein association. *Nature* *444*, 383–386.
- 704 Thurmond, J., Goodman, J.L., Strelets, V.B., Attrill, H., Gramates, L.S., Marygold, S.J.,

- 705 Matthews, B.B., Millburn, G., Antonazzo, G., Trovisco, V., et al. (2019). FlyBase 2.0: The next  
706 generation. *Nucleic Acids Res.* 47, D759–D765.
- 707 Uricchio, L.H., Petrov, D.A., and Enard, D. (2019). Exploiting selection at linked sites to infer the  
708 rate and strength of adaptation. *Nat. Ecol. Evol.* 3, 977–984.
- 709 Uversky, V.N. (2019). Intrinsically Disordered Proteins and Their “Mysterious” (Meta)Physics.  
710 *Front. Phys.* 7.
- 711 Wang, S., Li, W., Liu, S., and Xu, J. (2016). RaptorX-Property: a web server for protein structure  
712 property prediction. *Nucleic Acids Res.* 44, W430–W435.
- 713 Wu, D.D., Wang, G.D., Irwin, D.M., and Zhang, Y.P. (2009). A profound role for the expansion  
714 of trypsin-like serine protease family in the evolution of hematophagy in mosquito. *Mol. Biol.*  
715 *Evol.* 26, 2333–2341.
- 716 Xia, B., Yan, Y., Baron, M., Wagner, F., Barkley, D., Chiodin, M., Kim, S.Y., Keefe, D.L., Alukal,  
717 J.P., Boeke, J.D., et al. (2020). Widespread Transcriptional Scanning in the Testis Modulates  
718 Gene Evolution Rates. *Cell* 180, 248–262.e21.
- 719 Yan, J., and Kurgan, L. (2017). DRNApred, fast sequence-based method that accurately  
720 predicts and discriminates DNA- and RNA-binding residues. *Nucleic Acids Res.* 45, gkx059.
- 721 Yanai, I., Benjamin, H., Shmoish, M., Chalifa-Caspi, V., Shklar, M., Ophir, R., Bar-Even, A.,  
722 Horn-Saban, S., Safran, M., Domany, E., et al. (2005). Genome-wide midrange transcription  
723 profiles reveal expression level relationships in human tissue specification. *Bioinformatics* 21,  
724 650–659.
- 725 Yang, Z. (2007). PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* 24,  
726 1586–1591.
- 727 Yang, L., Zhang, Z., and He, S. (2016). Both Male-Biased and Female-Biased Genes Evolve  
728 Faster in Fish Genomes. *Genome Biol. Evol.* 8, 3433–3445.
- 729 Yang, Z., Wong, W.S.W., and Nielsen, R. (2005). Bayes empirical Bayes inference of amino  
730 acid sites under positive selection. *Mol. Biol. Evol.* 22, 1107–1118.
- 731 Zhang, J., and He, X. (2005). Significant impact of protein dispensability on the instantaneous  
732 rate of protein evolution. *Mol. Biol. Evol.* 22, 1147–1155.
- 733 Zhang, J., and Yang, J.R. (2015). Determinants of the rate of protein sequence evolution. *Nat.*  
734 *Rev. Genet.* 16, 409–420.
- 735 Zhang, Y.E., Vibranovski, M.D., Krinsky, B.H., and Long, M. (2010). Age-dependent  
736 chromosomal distribution of male-biased genes in *Drosophila*. *Genome Res.* 20, 1526–1533.
- 737 Zheng, N., Schulman, B.A., Song, L., Miller, J.J., Jeffrey, P.D., Wang, P., Chu, C., Koeppe, D.M.,  
738 Elledge, S.J., Pagano, M., et al. (2002). Structure of the Cul1-Rbx1-Skp1-F boxSkp2 SCF

739 ubiquitin ligase complex. Nature 416, 703–709.

740

741



742

743

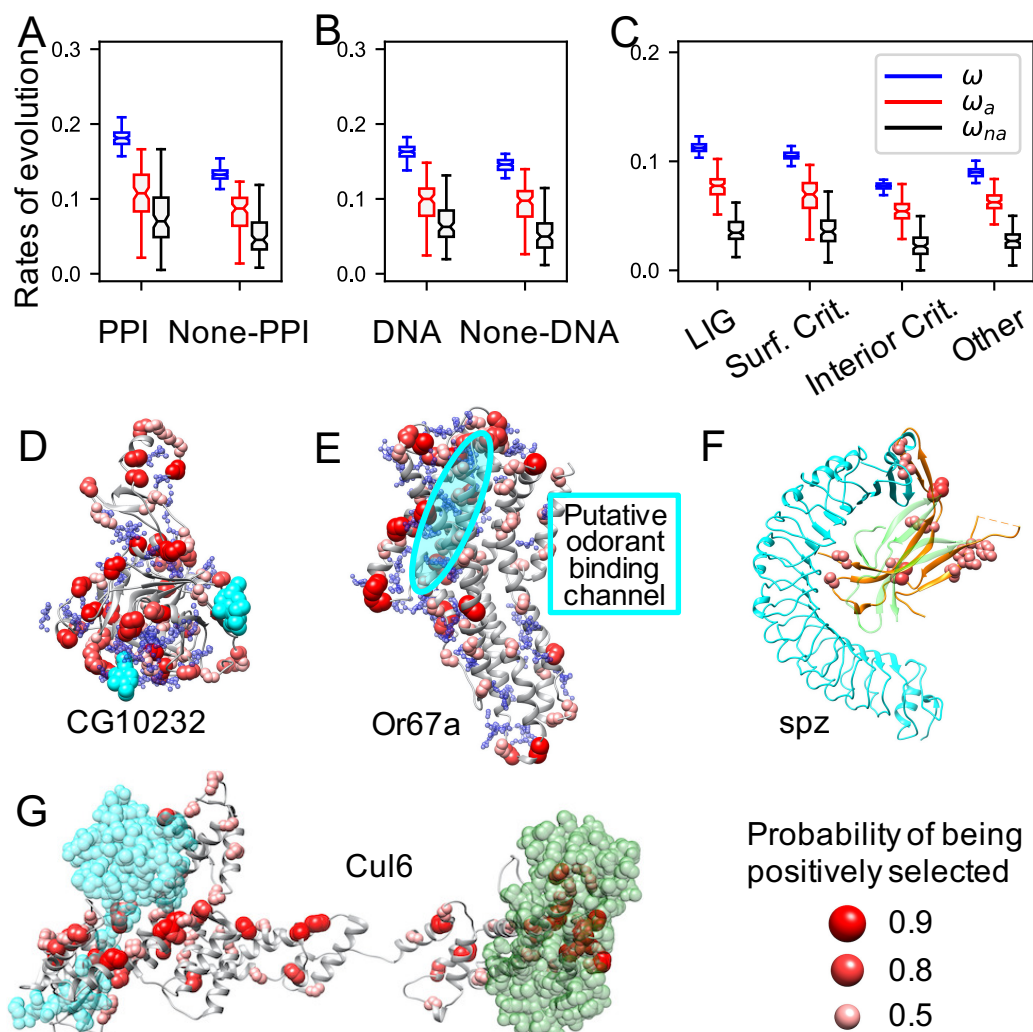
744

745

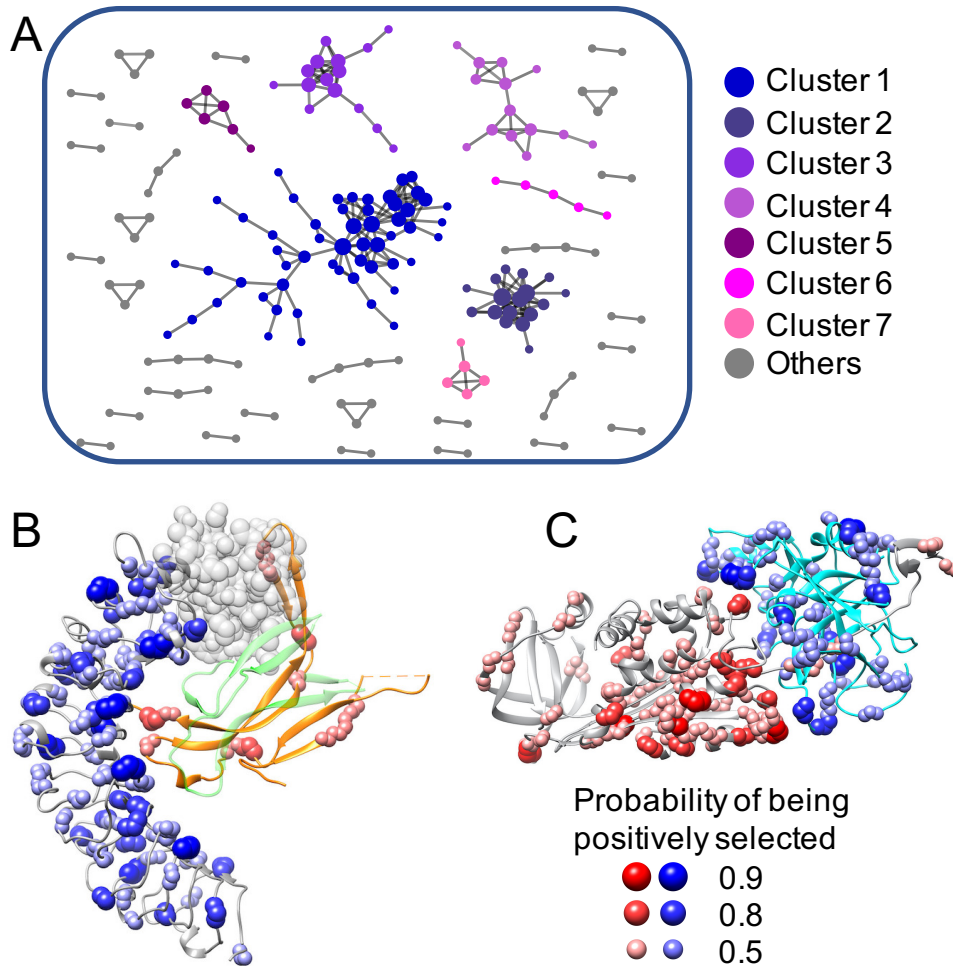
746

747

**Figure 1.** Protein length (A), mean (B), male (C) and female expression (D) levels in each gene groups divided by ascending  $\omega$  values. Values for each gene group and each gene property were computed through 1,000 bootstrapping steps. Obvious complex U-shaped correlations with  $\omega$  were observed for protein length (A) and male expression level (C).

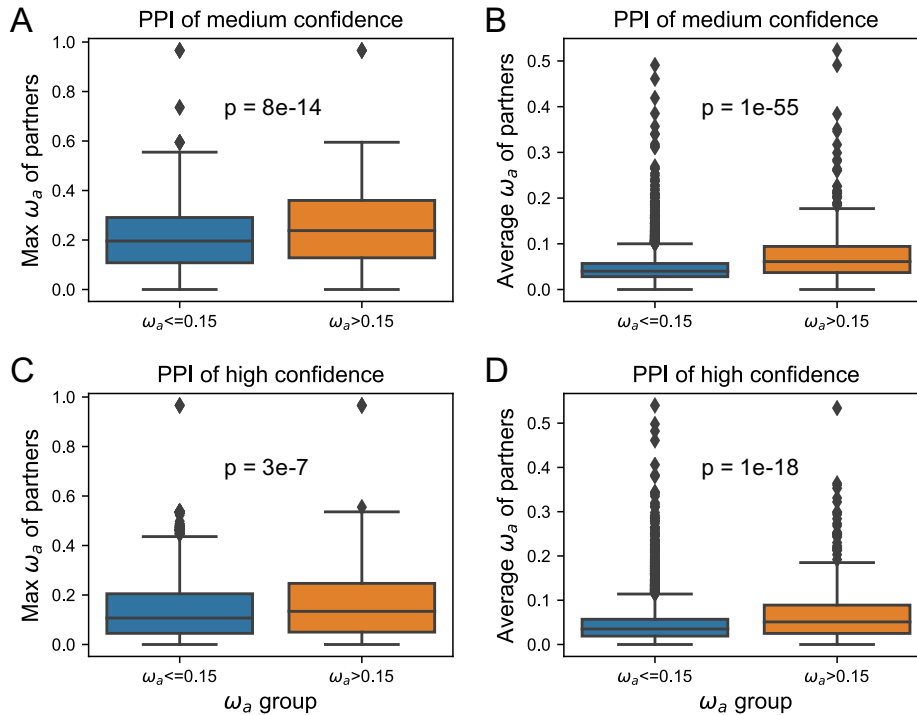


748  
 749 Figure 2. Adaptive evolution in molecular interaction sites. Protein-protein interaction sites (A),  
 750 DNA binding sites (B) and putative ligand binding sites (C) show higher adaptation rates than  
 751 none binding sites. Examples of positive selection around molecular interaction sites in high  
 752 quality structural models of CG10232 (D), Or67a (E), spz (F), and Cul6 (G). Except for spz  
 753 (PDB code 3e07), the other proteins are obtained from SWISS model repository. Putative ligand  
 754 binding pockets of CG10232 (D) and Or67a (E) are shown in blue spheres. Ligands including  
 755 interacting proteins are shown in cyan or green: NAG of CG10232 in cyan (D), Toll receptor of  
 756 spz in cyan (F), RING-box protein in cyan and F-box protein in green for Cul6 (G). The putative  
 757 odorant binding channel of Or67a is highlighted in cyan circle (E). The ligand poses in (D, F and  
 758 G) are obtained by superimposition from structure 2XXL, 4BV4 and 1LDK, respectively.  
 759  
 760



761  
762  
763  
764  
765  
766  
767  
768  
769  
770

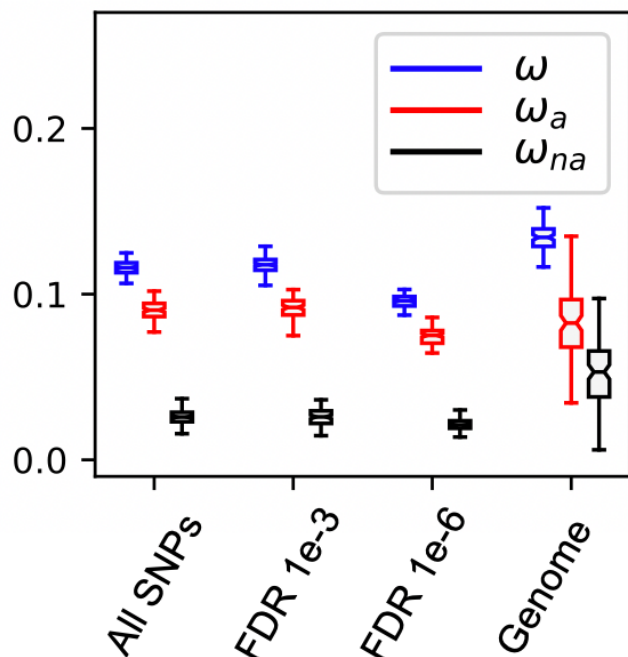
Figure 3. Co-adaptation of fast-adaptive proteins. (A) Sub-clusters of PPI networks of fast-adaptive proteins. Only proteins with at least one partner were shown. Examples of molecular interactions that might regulate co-adaptation in fast-adaptive proteins: (B) Toll-4 (gray) and spz (orange, with green representing the other spz monomer), (C) Spn28Db (gray, serine protease inhibitor 28Db) and CG18563 (cyan, with Go term “serine-type endopeptidase activity”). A putative N-terminus (transparent beads) of Toll-4 were built by superimposition from 4LXR, since the N-terminus were missing in the structural model. Complex structural model of Spn28Db and CG18563 was inferred from 1EZK.



771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781

Figure 4. Co-adaptation of PPIs in *D. melanogaster*. For fast-adaptive proteins, adaptation rates of their partners (orange box plot) are significantly larger compared to slow adaptive proteins (blue box plot). Max  $\omega_a$  of protein partners are shown in (A and C) and averaged  $\omega_a$ , of protein partners are shown in (B and D). PPI from STRING with median confidence (combined score larger than 0.4) are shown in (A and B), and PPI with high confidence (combined score larger than 0.7) are shown in (C and D).





782  
783 Figure 5. Adaptive evolution in FST SNPs. The significant SNPs at different FDR cutoffs all  
784 show much higher proportions of adaptation than genome-wide expectation.  
785

786  
787

**Table 1.** Pearson correlation coefficients between  $\omega$ ,  $\omega_a$  and  $\omega_{na}$  and gene properties <sup>a</sup>

Categories	Properties	$\omega$	$\omega_a$	$\omega_{na}$
Function-related properties	Gene age	0.55 (0)	0.34 (7e-154)	0.41 (2e-224)
	Protein length <sup>b</sup>	-0.22 (2e-136)	-0.13 (1e-25)	-0.30 (2e-142)
	Mean expression <sup>b</sup>	-0.21 (2e-109)	-0.11 (2e-18)	-0.11 (1e-18)
	Male expression <sup>b</sup>	-0.06 (5e-10)	0.00 (8e-1)	-0.03 (6e-2)
	Female expression <sup>b</sup>	-0.29 (2e-205)	-0.17 (2e-42)	-0.16 (3e-35)
	Male specificity	0.21 (2e-104)	0.12 (3e-22)	0.10 (3e-14)
	Tissue Specificity	0.30 (4e-226)	0.18 (1e-45)	0.19 (3e-48)
	PPI number <sup>b</sup>	-0.28 (1e-217)	-0.14 (4e-29)	-0.19 (4e-58)
	PPI-site ratio	0.14 (1e-50)	0.05 (7e-6)	0.10 (2e-16)
	DNA-site ratio	0.25 (8e-164)	0.12 (3e-23)	0.23 (3e-79)
Structure-related properties	Helix ratio	-0.05 (4e-7)	-0.01 (3e-1)	-0.05 (4e-5)
	Sheet ratio	-0.04 (2e-5)	0.00 (9e-1)	-0.01 (3e-1)
	Helix+sheet ratio	-0.09 (3e-25)	-0.02 (1e-1)	-0.07 (1e-9)
	Coil ratio	0.10 (8e-27)	0.01 (2e-1)	0.08 (8e-11)
	ISD	0.17 (8e-82)	0.04 (1e-3)	0.12 (2e-24)
	RSA	0.16 (7e-87)	0.06 (1e-6)	0.15 (3e-35)
Protein evolution	$\omega$	1.00 (0)	0.65 (0)	0.78 (0)
	$\omega_a$	0.65 (0)	1.00 (0)	0.03 (7e-3)
	$\omega_{na}$	0.78 (0)	0.03 (7e-3)	1.00 (0)

788 <sup>a</sup> Pearson correlation coefficient R were listed along with corresponding P-values in  
789 parentheses.

790 <sup>b</sup> To better estimate the correlations for sequence length, expression levels and PPI numbers,  
791 we used logarithmic scales rather than absolute values, which could vary dramatically from near  
792 zero to thousands.  
793 Abbreviations in this table: ISD, intrinsic structural disorder; RSA, relative solvent accessibility;  
794 PPI number, protein-protein interaction number; PPI-site ratio, ratio of protein-protein interaction  
795 sites; DNA-site ratio, ratio of DNA-binding sites.  
796