

[Click here to view linked References](#)

Jiang Page 1

1 **Comparative genomic analysis of *SET*-domain family reveals the**
2 **origin, expansion, and putative function of the arthropod-specific**
3 ***SmydA* genes as histone modifier in insects**

4 **Feng Jiang^{1,*}, Qing Liu^{1,2,*}, Yanli Wang^{2,3}, Jie Zhang¹, Huimin Wang¹, Tianqi**
5 **Song³, Meiling Yang², Xianhui Wang^{2,#}, Le Kang^{1,2,#}**

6 ¹ Beijing Institutes of Life Science, Chinese Academy of Sciences, Beijing, China

7 ² State Key Laboratory of Integrated Management of Pest Insects and Rodents,
8 Institute of Zoology, Chinese Academy of Sciences, Beijing, China

9 ³ Institute of Applied Biology, Shanxi University, Taiyuan, Shanxi, China

10 *These authors contributed equally to this study.

11 Corresponding authors:

12 Le Kang, Ph.D. and Professor

13 Institute of Zoology, Chinese Academy of Sciences

14 Beijing 100101, China

15 Tel: 86-10-6480-7219

16 Fax: 86-10-6480-7099

17 E-mail: lkang@ioz.ac.cn

18 OR

19 Xianhui, Wang Ph.D. and Professor

20 Institute of Zoology, Chinese Academy of Sciences

21 Beijing 100101, China

22 Tel: 86-10-64807220

23 Fax: 86-10-6480-7099

24 E-mail: wangxh@ioz.ac.cn

Evolution of *SET* Genes in Insects

25 Abstract

1
2
3 26 The *SET* domain is an evolutionarily conserved motif present in histone lysine
4
5 27 methyltransferases, which are important in the regulation of chromatin and gene
6
7
8 28 expression in animals. In this study, we searched for *SET* domain-containing genes
9
10
11 29 (*SET* genes) in all of the 147 arthropod genomes sequenced so far to understand the
12
13 30 evolutionary history by which *SET* domain have evolved in insects. Phylogenetic and
14
15 31 ancestral state reconstruction analysis revealed a arthropod-specific *SET* gene family,
16
17 32 named *SmydA*, which is ancestral to arthropod animals and specifically diversified
18
19
20 33 during insect evolution. Considering that pseudogenization is the most probable fate
21
22 34 of the new emerging gene copies, we provided experimental and evolutionary
23
24
25 35 evidence to demonstrate their essential functions. Fluorescence *in situ* hybridization
26
27 36 analysis and *in vitro* methyltransferase activity assays showed that the *SmydA-2* gene
28
29
30 37 was transcriptionally active and retained the original histone methylation activity.
31
32 38 Expression knockdown by RNA interference significantly increased mortality,
33
34
35 39 implying that the *SmydA* genes may be essential for insect survival. We further
36
37 40 showed predominantly strong purifying selection on the *SmydA* gene family and a
38
39
40 41 potential association between the regulation of gene expression and insect phenotypic
41
42 42 plasticity by transcriptome analysis. Overall, these data suggest that the *SmydA* gene
43
44 43 family retains essential functions that may possibly define novel regulatory pathways
45
46
47 44 in insects. This work provides insights into the roles of lineage-specific domain
48
49
50 45 duplication in insect evolution.

51
52 46 *Key words:* insects, domain, gene duplication, histone modification.
53
54
55
56
57

58 Evolution of *SET* Genes in Insects
59
60
61
62
63
64
65

47 **Background**

1
2
3 48 Protein domains are functional and structural units that are evolutionary well
4
5 49 conserved across species [1]. Specific protein domains are often linked to discrete
6
7
8 50 biological function; therefore, the frequent duplication, gain, and loss of protein
9
10
11 51 domains play substantial roles in functional novelty [2]. Domain duplication can be
12
13 52 achieved via frequent domain-containing gene family expansion. Thus, the member
14
15 53 number of a gene family that contains domains can be expanded, representing a
16
17
18 54 common method by which divergence to domain sequences can lead to the
19
20 55 evolutionary novelty of domain-containing genes [3]. In taxonomically related
21
22 56 species, the expansion of conserved gene families through gene duplication is
23
24
25 57 widespread in metazoan genomes [4]. Gene duplication may increase species fitness
26
27 58 by subfunctionalization or neofunctionalization [5, 6]. Subfunctionalization results in
28
29
30 59 the symmetric division of the functional capability of the original gene among the
31
32 60 duplicated genes [7]. Neofunctionalization allows the original copy to maintain its
33
34
35 61 function and permits the new copy to diverge under relaxed selective constraints or
36
37 62 positive selection for a novel function. Rapid domain diversification followed by gene
38
39
40 63 duplications in particular lineages is important for the adaptation of lineage-specific
41
42 64 ecological specializations [8].

43
44 65 Histones are highly alkaline proteins in cell nuclei that package and order the
45
46 66 nuclear DNA into nucleosomes, which are the main components of chromatin.
47
48
49 67 Histone modifications are a major epigenetic regulatory mechanism for phenotypic
50
51
52 68 plasticity in insects. Inhibition of histone deacetylation affects developmental
53
54 69 plasticity both in ants (*Camponotus floridanus*) and honeybees (*Apis mellifera*) [9, 10].
55
56
57 70 Genome-wide profiling of histone modifications revealed an important role of histone
58
59 Evolution of *SET* Genes in Insects
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

71 H3 lysine 27 acetylation in the caste differentiation of ants [11]. Methylations of
72 histone H3 lysine 27 and histone H3 lysine 36 are more abundant in queen ovaries
73 than in larvae, implying that histone methylation plays a specific role in honey bees
74 [12]. In recent years an increasing number of publications have established histone
75 lysine methylation as a central epigenetic modification in regulation of chromatin and
76 transcription. The *SET* domain, which is observed in many histone lysine
77 methyltransferases, is widely and probably universally distributed in metazoan
78 species. This protein family typically comprises an approximately 130 amino
79 acid-long *SET* domain, which was identified in the strongest PEV suppressor gene
80 *Su(var)3-9*, in the Pc-G gene Enhancer of zeste [E(z)] and in the activating *trx-G* gene
81 *Trithorax* of *Drosophila* [13]. The *SET* domain possesses a catalytic activity that
82 transfers a methyl group to the amino group of lysine residues of nuclear histones
83 from S-adenosyl-L-methionine. Based on their biochemical characteristics, *SET*
84 domain is capable of catalyzing mono-, di- or tri-methylation of their lysine
85 substrates. *SET* domain-dependent methylation has been identified in a wide range of
86 lysine residues in different histones: K4 (K is the abbreviation for lysine), K9, K27,
87 K36, and K79 in histone H3; K20 in histone H4; K59 in the globular domain of
88 histone H4; and K26 in histone H1B [14]. Methylation of lysine residues in histone
89 proteins is an important post-translational epigenetic event that regulates gene
90 expression by serving as an epigenetic marker for the recruitment of complexes that
91 participate in the organization of chromatin structure [15]. The importance of
92 *SET*-domain containing genes is strongly supported by the involvement of this protein
93 family in diverse biological mechanisms, such as transcriptional activation,
94 transcriptional repression, enhancer function, mRNA splicing and DNA replication

Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

95 [16]. Therefore, expectedly, the regulation of various *SET*-domain containing genes
96 are increasing correlated with diverse epigenetic phenomena which, for example,
97 include epigenetic control in plants, centromeric gene silencing in yeasts,
98 repeat-induced point mutations in fungi, DNA elimination in *Tetrahymena*, germline
99 chromatin silencing in worms and heterochromatin formation in flies [17].

100 Insects constitute a remarkably diverse group of organisms that make up a vast
101 majority of known species with their importance including biodiversity, agricultural,
102 and human health concerns. The insect lineage comprises species that are both
103 cosmopolitan distributed and geographically restricted, showing a broad range of
104 adaptation diversity. The evolutionary history of gene families is not confounded by
105 whole-genome duplication, and the major topology of insect species is well resolved
106 [18]. Therefore, the insect lineage offers an excellent model to study domain/gene
107 evolution in the context of gene family dynamics [19-23]. Insect *SET*
108 domain-containing genes (*SET* genes) have been identified in a limited number of
109 representative insect species without complicated analysis [24-26]. The *Smyd*
110 subfamilies of *SET* genes have expanded in a few insects from Diptera and
111 Hymenoptera, and several members of the *Smyd* subfamilies show significant changes
112 in gene expression in response to phenotypic plasticity in ants [27, 28]. However, the
113 evolutionary history of insect *SET* genes remains largely unknown because the *SET*
114 genes from a broad range of insect species have not been combined in a single
115 evolutionary framework. Therefore, a comprehensive study of the origin and
116 diversification of the *SET* gene family in insects is required. Accurate classification of
117 *SET*-domain containing genes can pave the fundamental way to further understanding
118 the epigenetic basis of gene regulation in insects.

Evolution of *SET* Genes in Insects

119 In the present study, we aimed to ascertain the origin and diversification of *SET*
120 genes in insects. We searched for *SET* genes in the 130 insect genomes and the 17
121 arthropod genomes as outgroups. These 130 insect species include both
122 hemimetabolous and holometabolous insects and cover all the insect species for
123 which genome data have been fully available and annotated so far. Our phylogenetic
124 analysis revealed that an important diversification of arthropod-specific *SET* genes,
125 *SmydA*, occurred during insect evolution. Experimental evidence of the important
126 functions of *SmydA* genes in insects was obtained through fluorescence *in situ*
127 hybridization, *in vitro* methyltransferase activity assay, and survival assay after
128 expression knockdown. Furthermore, we compared the gene expression patterns and
129 examined the selection signatures of *SmydA* genes in the four representative insects
130 exhibiting phenotypic plasticity. These results provides insights into the regulatory
131 roles of lineage-specific domain duplication in insect evolution.

132

133 **Results**

134 **Identification and phylogenetic classification of *SET* genes**

135 We comprehensively searched for *SET* genes in a wide range of sequenced insect
136 species, which included 130 insect species from 14 insect orders (Supplementary
137 Table S1). The *SET* genes were defined by the presence of the *SET* domain as
138 predicted by the HMMER search, and their gene models were manually improved.
139 Seventeen non-insect arthropods were also included to achieve ancestral status along
140 with insect evolution. In total, 4,498 *SET* genes were identified in the 147 arthropod
141 genomes (Supplementary Table S2). The genes showing potential pseudogene signals
142 were removed in these identified *SET* genes. A database webserver

Evolution of *SET* Genes in Insects

143 (<http://159.226.67.242:8080/>) has been constructed to select, retrieve, and analyze the
144 data in this study. In insects, the number of *SET* genes found per species ranges from
145 16 in the scuttle fly *Megaselia scalaris* to 81 in the mosquito *Culex quinquefasciatus*
146 (Table 1 and see Supplementary Table S3 for the full list of summary of *SET* genes in
147 the 147 arthropod genomes). This observation suggests that the size of *SET* genes
148 varies significantly among different insect lineages. Although the genome size of the
149 migratory locust *Locusta migratoria* is approximately 30-fold that of the fruit fly
150 *Drosophila melanogaster* [29], the number of *SET* genes in locusts is comparable
151 with that of flies. Thus, the genome sizes and number of *SET* genes are not linearly
152 correlated. The specificity of certain substrates is reflected by the classification of
153 *SET* genes, and *SET* genes can be classified into seven major conserved groups,
154 namely: Suv, Ash, Trx, E(z), PRDM, SMYD, and SETD [24]. We performed
155 phylogenetic analysis of the *SET* genes for representative species to obtain insights
156 into the evolution of insect *SET* genes. Multiple sequence alignments of complete
157 proteins could not accurately determine the homologous sites of *SET* genes because of
158 the considerably different sequence lengths and domain architectures of these genes.
159 Thus, alignment-based methods using Bayesian inferences for *SET* domain sequences
160 and alignment-free methods based on feature frequency profiles for complete protein
161 sequences were conducted to infer phylogenetic relationships. The overall tree
162 topologies (Figure 1) inferred using the two methods were generally consistent. Based
163 on the previous nomenclature system [24], the phylogenetic tree topology allows the
164 grouping of insect *SET* genes into seven major conserved groups, generally showing
165 slight fluctuation in the member sizes in each conserved group. The protein domains
166 for each *SET* gene were annotated using the InterProScan package. In general, the

Evolution of *SET* Genes in Insects

167 *SET* genes in the same conserved group exhibited a similar domain composition,
168 suggesting that the domain architectures support the conserved group classification
169 inferred through the phylogenetic analysis. In addition to the *SET* genes in the
170 conserved groups, a large number of *SET* genes could not be classified into known
171 conserved groups on the basis of the phylogenetic analysis. These unclassified genes
172 act as potential “arthropod-specific” genes. The lineage-specificity was further
173 verified through reciprocal BLAST search against known *SET* genes of nematodes
174 and humans.

175

176 **Ancestral states of the *SET* gene family in insects**

177 A character matrix that represents the present/absent states for each *SET* homologous
178 group was constructed to infer the ancestral states of interior nodes along with the
179 species tree using the Mesquite program. The ancestral states at different nodes could
180 infer the appearances/losses of the *SET* homologous group that occurred at and above
181 the level of orders (Figure 2A). The grouping of *SET* homologous genes for each
182 species was inferred using the OrthoMCL program with the corresponding
183 orthologous *SET* gene in *D. melanogaster*, and the grouping reliability was supported
184 by the phylogenetic analysis (Supplementary Figure S1–S5). The putative ancestral
185 state was composed of 19 *SET* homologous groups present in the last common
186 ancestor (LCA) of the studied arthropod species. Generally the insect species
187 possessed more *SET* homologous groups than the chelicerata species studied,
188 suggesting that *SET* homologous groups considerably expanded during insect
189 evolution. At the interior clades, novel *SET* homologous groups emerged several
190 times. Only few losses of *SET* homologous groups, such as the loss of *SmydA-3*, were

Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
191 observed at the interior clades. The large fluctuation of *SET* homologous groups in
192 each species indicates that these groups experienced rapid lineage-specific
193 expansion/contraction within insect orders. For example, in Hymenoptera, the number
194 of *SET* homologous groups ranged from 18 (covering 23 *SET* genes) in the jumping
195 ant *Harpegnathos saltator* to 30 (covering 52 *SET* genes) in the parasitoid wasp
196 *Nasonia vitripennis*. In Diptera, 13 *SET* homologous groups (covering 14 *SET* genes)
197 were found in *M. scalaris*, and the oriental fruit fly *Bactrocera dorsalis* possessed
198 only 31 *SET* homologous groups (covering 45 *SET* genes). A large number of
199 arthropod specific *SET* homologous groups cannot be classified into the seven major
200 conserved groups, which revealed their origin after the emergence of main arthropod
201 lineages. Nevertheless, at least six of these groups were present among insect species
202 belonging to different orders, indicating their broad conservation in insects (Figure
203 2B).

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
204 *SET* domains do not just function as an independent unit, as in many proteins it
205 co-occurs with multiple other protein domains to regulate their target specificity and
206 catalysis [16]. We surveyed the gene ontology (GO) classification of proteins by
207 integrating biological knowledge into three hierarchies, namely, biological process,
208 molecular function, and cellular component, to assess the function innovation of
209 domain acquisition globally. The common GO categories included histone lysine
210 methylation (GO:0034968), regulation of transcription (GO:0006355), protein
211 binding (GO:0005515), nucleic acid binding (GO:0003676), and metal ion binding
212 (GO:0046872) (Figure 3A). Partitioning of *SET* gene families between the conserved
213 and arthropod specific groups revealed that GO categories could be shared between
214 the two groups or be assigned exclusively to one group. The GO categories, which

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

215 were only exclusive in the arthropod specific groups, included RNA
216 methyltransferase activity (GO:0008173), metallocarboxypeptidase activity
217 (GO:0004181), lysozyme activity (GO:0003796), homophilic cell adhesion
218 (GO:0007156), sulfotransferase activity (GO:0008146) and so on.

219

220 **Emergence of arthropod lineage-specific *SET* gene families**

221 Pairwise BLAST search against all the *SET* genes indicated that the arthropod specific
222 *SET* genes showed considerable amino acid similarity to the SMYD groups, which
223 contain a conserved core consisting of a *SET* domain and a MYND (Myeloid
224 translocation protein, Nervy, Deaf) zinc finger domain [30]. The arthropod specific
225 *SET* genes also contain the *SET* and MYND domains and were named *SmydA* [28].
226 We performed the phylogenetic analysis of the SMYD genes through Bayesian
227 inferences. The majority of the SMYD genes could be classified into 11 monophyletic
228 clades, which exhibited similar high Bayesian posterior probability values (Figure
229 3B). In a global view, these SMYD genes fell into two distinct branches, which
230 correspond with the conserved SMYD and *SmydA* groups. These results could
231 exclude the possibility that the *SmydA* groups have raised from multiple independent
232 gain events by duplications from deeply diverged SMYD genes of insects. Indeed,
233 *SmydA* genes were absent from in all Chelicerata species investigated but present in
234 the genomes of crustacean species and insect species, suggesting that *SmydA* genes
235 may have originated prior to the divergence of Crustacea and Insecta. *SmydA-1*,
236 *SmydA-2*, *SmydA-3*, and *SmydA-6* were already present before the split of Crustacea
237 with other insects, showing clues for their ancient duplication events. The strong
238 support for distinct individual lineages of paralogous genes implied that multiple
Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

239 duplications occurred within the order level; the most notable case was the detection
240 of three copies of *SmydA-3* in the red flour beetle *Tribolium castaneum*
241 (Supplementary Table 2). *SmydA-1/SmydA-4* and *SmydA-6* were subjected to
242 additional rounds of duplication in Lepidoptera and Orthoptera, respectively. The
243 genes annotated as *SmydA-8* and *SmydA-9* in *D. melanogaster* previously formed a
244 single clade alone with a high Bayesian posterior probability value (0.99), suggesting
245 a specific duplication event in *Drosophila*. Therefore, the *SmydA* groups differed
246 considerably in the number of genes in each insect order, implying the complexity of
247 their evolutionary histories.

248 To shed light into the evolutionary history of *SmydA* genes, we determined the
249 location and gene order of *SmydA* genes in the four holometabolous species with
250 available chromosome-level genome assemblies or genome-scale genetic linkage
251 maps (Figure 3C). In Diptera, the syntenic gene orders could be inferred from the four
252 ancient *SmydA* genes, namely, *SmydA-1*, *SmydA-2*, *SmydA-3*, and *SmydA-6*, all of
253 which may have been present in the ancestor of insects and crustaceans. An
254 insect-specific *SmydA-9* could be observed in the majority of insect orders, including
255 both hemimetabolous and holometabolous insects. *SmydA-9* showed syntenic
256 conservation with the four ancient genes. This gene order was also conserved when
257 *SmydA* genes in insects distantly related from other insect orders were examined.
258 Almost all of the five synteny-anchoring genes were maintained in both the
259 coleopteran species *T. castaneum* and hymenoptera species *A. mellifera*, with an
260 exception of *SmydA-2* that was missed in *A. mellifera*. In contrast to those in *T.*
261 *castaneum* and *A. mellifera*, the reversed order of *SmydA-3* and *SmydA-6* in Dipteran
262 species implies that an intrachromosome transfer event of genomic segments occurred

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

263 before the emergence of Diptera. Duplication events could also occur in the early
264 diversification of arthropod species. No orthologous *SmydA-4* gene was detected the
265 chelicerata species, indicating that duplication event contributes to the emergence of
266 *SmydA-4* gene in Pancrustacea species. *SmydA-4* was present all the hemimetabolous
267 insect orders studied, as well as in the holometabolous insect orders Lepidoptera,
268 Coleoptera, and Diptera. The absence of *SmydA-4* in all the 32 hymenopteran species
269 suggested that subsequent loss of *SmydA-4* could be traced back to the ancestor of the
270 hymenopteran lineage before the divergence of wasp, ants, and bees. In the SMYD
271 phylogenetic tree, the Bayesian inferences supported the grouping of *SmydA-3*,
272 *SmydA-4*, and *SmydA-6*. Three of the four species exhibited a accordant location of
273 *SmydA-3/SmydA-4/SmydA-6* in the syntenic regions. In addition to the old duplication
274 events that categorized the divergent duplicates into distinct *SmydA* subfamilies (e.g.,
275 *SmydA-3* and *SmydA-4*), recent duplications within an insect order were also
276 observed. The three copies of *SmydA-3* in *T. castaneum*, which spanned within a 4.2
277 kb genomic region, were observed in tandem array between the two syntenic genes
278 *SmydA-1* and *SmydA-6*. The closeness in protein sequence and genomic location
279 implies an evolutionary origin of these three copies of *SmydA-3* via local duplication.
280 Overall, our data suggest that the order of *SmydA* genes was conserved over a
281 remarkable wide range of holometabolous insect orders.

282 283 **Selective pressures acting on *SmydA* genes**

284 Functional differentiations or mutations leading to pseudogene formation were the
285 two major causes for sequence divergence between new duplicates and their
286 orthologous counterpart. Synonymous substitutions are assumed to accumulate at a

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

287 constant rate; hence, the ratios of nonsynonymous substitution per nonsynonymous
288 site (d_N) to synonymous substitution per synonymous site (d_S) are deemed to be an
289 indicator to measure the relative rates of evolution for protein sequences. The four
290 genes (*ACYPI26757* and *ACYPI55839* in *Acyrtosiphon pisum*; Px015362.1 and
291 Px001029.1 in *Plutella xylostella*) showing signals of recombination were removed
292 from the further selection analysis. We estimated a global d_N/d_S ratio (one ratio, model
293 M0) for these *SET* genes to determine whether the *SmydA* genes have been under
294 different selection pressures than the other conserved *SET* genes. The d_N/d_S ratios (ω
295 = d_N/d_S ratio) of *SET* genes varied from low (0.0007, Ez, CG6502) to high (0.1627,
296 *Smyd4-1*, CG1868), indicating a variance in the rates of protein evolution on different
297 *SET* genes (Table 2). The ω values among the conserved *SET* genes (excluding the
298 SMYD genes) ranged from 0.0007 to 0.0624 (mean ω = 0.0185). The conserved
299 SMYD and *SmydA* groups showed ω values in the ranges of 0.055–0.1627 (mean ω =
300 0.1020) and 0.0052–0.1623 (mean ω = 0.0884), respectively. Overall, both the
301 conserved SMYD and *SmydA* (P = 0.0003 and P = 0.0178, Wilcoxon signed-rank
302 tests with Bonferroni correction, respectively) groups exhibited significantly higher ω
303 values than the conserved *SET* genes (Figure 3D). However, the distributions of ω
304 values of the conserved SMYD and *SmydA* groups were statistically indistinguishable
305 (P = 1.0000, Wilcoxon signed-rank tests with Bonferroni correction).

307 **Function approval of *SmydA* genes**

308 We attempted to determine whether the *SmydA* genes retained histone methylase
309 activities to approve the non-pseudogenization process of these genes. We expressed
310 *SmydA-2* as a randomly selected representative and performed *in vitro* histone
Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
311 methylation activity assays using histones as substrates in the migratory locust. As
312 shown in Figure 4A, Western blot analysis detected increased lysine methylation on
313 histone H3 compared with the controls, indicating that *SmydA-2* possesses
314 methyltransferase activity on histones. Similar to that of the other conserved SMYD
315 genes, the methyltransferase activity of *SmydA-2* was also dependent on S-adenosyl
316 methionine. Fluorescence *in situ* hybridization analysis provided further tissue
317 expression evidence to support the reliability of the *SmydA-2* gene function. Obvious
318 fluorescence signals were observed in the brain and epidermal cells of cuticle in the
319 locusts (Figure 4B). These cells did not show any hybridization signal for the negative
320 controls. The origin and evolution of new emerging genes undergo an increased
321 expression breadth of new duplicated genes over evolutionary time [31, 32]. Thus, we
322 determined the expression levels of the *SmydA-2* gene using quantitative real-time
323 polymerase chain reaction (qPCR) analysis in the different tissues. qPCR data showed
324 that the *SmydA-2* gene was expressed in a broad range of tissues, including brains,
325 testes, ovaries, cuticles, and legs (Figure 4C). The broad expression pattern suggests
326 that the *SmydA-2* gene is less tissue specific and may serve as a functional gene in
327 multiple tissues [32].

41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
328 To determine whether the *SmydA-2* gene plays an essential role during
329 development [33], we knocked its expression down by using RNA interferences in the
330 locusts. Compared with the controls, the relative mRNA level of the *SmydA-2* gene
331 decreased by approximately 70% after injecting double-strand RNAs (Supplementary
332 Figure S6). After injection of *dsSmydA-2*, we observed large numbers of dead locusts,
333 which did not display obvious defect phenotype. As shown in Figure 4D,
334 Kaplan–Meier survival estimates indicate that injection of locusts with *dsSmydA-2*

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

335 significantly increased mortality when compared with the controls ($\chi^2 = 6.260$, $df = 1$,
336 $P = 0.0123$, Chi-square tests).

337

338 **Expression and selection analysis of *SmydA* genes in response to phenotypic plasticity**

339 Epigenetic reprogramming that modifies chromatin structure through histone
340 modifiers contributes to orchestrate the generation and maintenance of phenotypic
341 plasticity, which is a key trait for the success of insects. Therefore, we compared the
342 expression patterns of histone-modifier *SET* genes in four representative insects
343 exhibiting phenotypic plasticity, namely, locust density-dependent behavior, aphid
344 seasonal morphs, dietary-mediated interactions of bees and ants. Specially, we
345 performed differential expression analysis between gregarious and solitary locusts,
346 between asexual and sexual morphs in *A. pisum*, between queens and workers in *A.*
347 *mellifera*, and between large workers and queens in *Acromyrmex echinator*. In all the
348 four species, a number of differentially expressed genes (DEGs) were detected
349 between the two alternative phenotypes using the criteria of a false discovery rate
350 (FDR)-corrected $P < 0.05$. In terms of DEG number, a large portion of *SET* genes
351 showed significant changes in gene expression (16 in 31, 52%, in *A. mellifera*; 25 in
352 59, 50%, in *A. pisum*; 13 in 29, 45%, in *L. migratoria*; and 11 in 27, 41%, in *A.*
353 *echinator*). Compared with those DEGs observed at the genome-wide level (DEGs in
354 total), the sensitivities of DEG number in *SET* genes in the four insects were even
355 prominent, emphasizing the important regulatory role of *SET* genes in phenotypic
356 transition ($P_s < 0.02$, Chi-square tests). Overlapping of the differentially expressed
357 *SET* genes derived from the same ortholog could provide a clue of their convergent
358 function in phenotypic transition. We found three *SET* genes, namely, *Hmt4-20*, *Set2*,

Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

359 and *SmydA-5*, showed significant changes in gene expression simultaneously in three
360 of the four insect species studied.

361 Assuming that a non-pseudogene gene should not be randomly expressed, we
362 compared the expression pattern of the duplication-derived *SmydA* genes to their
363 derived ancestral SMYD genes in response to environment-dependent phenotypic
364 plasticity (Figure 5). The majority of *SET* genes from the conserved SMYD (30 in 30
365 in total, 100%) and *SmydA* (13 in 17 in total, 76%) groups were expressed in at least
366 one insect. No significant differences ($P = 0.710$, Chi-Square tests) in the number of
367 expressed genes were observed between the two groups. A number of DEGs were
368 detected in both the conserved SMYD and *SmydA* groups in the four insect species.
369 All the two *SmydA* genes in *A. pisum* and all the four *SmydA* genes in *A. echinatio*
370 were also differentially expressed. We also obtained significant results in three of the
371 six *SmydA* genes of *L. migratoria* and in two of the five *SmydA* genes of *A. mellifera*
372 between the two alternative phenotypes. The DEG number in the *SmydA* groups did
373 not show significant deviation from those in the conserved SMYD group in the four
374 insects ($P_s > 0.2$, Fisher's exact tests). This result suggests that the *SmydA* genes
375 might not be randomly expressed and that they did not represent as pseudogenes or
376 transcriptional byproducts. Thus, the *SmydA* genes may preserve a regulatory role,
377 indicating the function similarity to their ancestral SMYD genes.

378 The free ratio model of *SmydA* genes fitted the data significantly better than the
379 one model (model M0) using likelihood ratio tests ($P_s < 0.001$), indicating
380 heterogeneous rates of sequence evolution along the gene tree of *SmydA* genes.
381 Therefore, we tested whether the differentially expressed *SmydA* genes between
382 alternative phenotypes (foreground branches) evolved under different selective

Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

383 pressures than those in the remaining branches (background branch) (Supplementary
384 Figure S7). The branch model was much better supported by the data than the model
385 M0 for *SmydA-5* in *A. mellifera* and *SmydA-1* in *L. migratoria* (Table 3). Fixing $\omega = 1$
386 for the foreground branch did not result in an improved fit over the branch model with
387 the unconstrained foreground branch (the null neutral model and the alternative
388 model). This result suggests that the ω values in the external branch were smaller than
389 1 for *SmydA-3* and *SmydA-5* in *A. mellifera*, *SmydA-1* in *L. migratoria*, and *SmydA-3*
390 in *A. echinatio*. Only *SmydA-1* in *L. migratoria* exhibited elevated ω values, and a
391 branch-site model allowing heterogeneous ω values across sequences and branches
392 identified four sites (5M, 11K, 93P, and 105C) under positive selection.

393

394 Discussion

395 In this study, the phylogenetic analyses allowed the subdivision of the insect *SET*
396 genes into seven major conserved groups and one arthropod-specific *SmydA* group.
397 We inferred many *SmydA* gene duplication events along insect evolution, suggesting
398 an important diversification of the *SmydA* genes occurred during insect evolutionary
399 processes. With the *SmydA-2* genes in locusts as representatives, the maintenance of
400 essential gene function was confirmed from the experimental evidence of *in vitro*
401 methyltransferase activity, *in situ* mRNA expression, and phenotypes after expression
402 knockdown. Based on the examination of distribution pattern and selection signatures
403 across insects, our data indicated that extensive pseudogenization unlikely occurred
404 for the *SmydA* genes. Finally, the transcriptome analyses of the four insects showed
405 that several *SmydA* genes are involved in insect phenotype plasticity, suggesting that

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

406 *SmydA* genes contributed novelties for insect adaptive evolution. This data suggest a
407 role of diverged regulatory functions after their duplication in insects.

408 *SmydA* genes represent a class of arthropod-specific genes that are only present in
409 the LCA of insect species and crustacean species, suggesting their origin after the split
410 of chelicerates from Pancrustacea species. Conservation of five ancient *SmydA* genes
411 in a wide range of species suggests they probably originated from duplication events
412 of conserved SMYD genes predating the diversification of insects. Although a few
413 cases of whole-genome duplication have been documented in chelicerates, evidence
414 that whole-genome duplication occurs widely in arthropod evolution remains lacking
415 [34]. Therefore, gene duplication rather than whole-genome duplication possibly
416 leads to the emergence of multiple copies of ancient *SmydA* genes in the LCA of
417 Pancrustacea species. The clear split of conserved SMYD and *SmydA* genes excluded
418 the possibility that multiple independent duplication events from conserved SMYD
419 genes resulted in the current repertoire of *SmydA* genes in insects. This result suggests
420 that the five ancient *SmydA* genes were first produced from a single ancestral gene,
421 which was derived from conserved SMYD genes. The five ancient *SmydA* genes were
422 thus the source from which insect-specific *SmydA* duplications were subsequently
423 produced in insects. Determining the location and order of multiple gene members at
424 the genomic scale sheds light on the evolutionary history of gene family. The closely
425 linked manner in genomic location suggests that homologous recombination and
426 functional differentiation may be a major force to shape the evolution of *SmydA* genes
427 in insects. For instance, in dipteran and lepidopteran insects, homologous
428 recombination may give rise to *SmydA-6* via the duplication events of *SmydA-3*
429 because *SmydA-3* and *SmydA-6* were in close proximity to each other in both genomic

Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

430 location and phylogenetic trees. The tandem organization of three *SmydA-3* copies in
431 *T. castaneum* may also result from species-specific duplications via homologous
432 recombination. Retrotransposition events may represent another contributing force for
433 generating unlinked *SmydA* genes; these events can also generate intronless
434 retroposed gene copies [35]. However, the retrotransposition events could not be
435 inferred from the presence of signature of intron–exon structure because of the
436 subsequent insertion in deeply diverged duplicates, such as *SmydA-5*. Conserved gene
437 orders between species from Lepidoptera, Coleoptera, and Diptera revealed a high
438 degree of macrosyntenic gene order of the five ancient *SmydA* genes during
439 approximately 348 million years of evolutions splitting these insects [36]. This
440 observation implies strong constraints for preserving the conserved gene order of
441 *SmydA* genes in insects. Currently, whether this macro-syntenic gene order is
442 preserved outside holometabolous insects cannot be determined because
443 chromosome-level genome assemblies or genome-scale genetic linkage maps are not
444 available in hemimetabolous insects. This issue would be addressed when the genome
445 assembly is considerably improved in the future.

446 Selective pressures were significantly weaker for the SMYD genes than for the
447 six conserved groups (Suv, Ash, Trx, E(z), PRDM, and SETD). Compared with the
448 six conserved groups, SMYD genes were the least conserved gene group and,
449 concordantly, the least constrained one. Nevertheless, the ω values of SMYD genes
450 ranged from 0.0052 for *SmydA-2* to 0.1627 for *Smyd4-1*. $\omega \ll 1$ was consistent with
451 their broad conservation across insects, implying their essential functional roles. This
452 observation suggests that purifying selection is the main force governing the evolution
453 of SMYD genes. The distributions of ω values of the conserved SMYD and *SmydA*

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

454 genes were statistically indistinguishable, indicating a symmetrical rate of sequence
455 evolution. Thus, purifying selection is subject to the conserved SMYD and *SmydA*
456 genes, but their intensity may be relaxed compared with other *SET* genes. Both the
457 GO analysis and the *in vitro* methyltransferase activity assay suggest that *SmydA*
458 genes, similar to their conserved SMYD ancestors, are sufficient to perform the
459 original function relating to histone methylation [37]. GO ontology analysis implied
460 that the *SmydA* genes have developed to acquire novel functions. These functions
461 were absent in the conserved SMYD genes, indicating that the *SmydA* genes may have
462 undergone functional differentiation. Gene duplications that occurred in specific
463 lineages are important in contributing to lineage-specific adaptive processes [38].
464 After gene duplication, purifying selection is expected in both gene copies if
465 duplication can confer a selective advantage [7]. By contrast, one of the two copies
466 can evolve either under relaxed purifying selection when no immediate advantage is
467 shown from gene duplication or under positive selection when a new function is
468 acquired via advantageous mutations [39]. Overall, these data suggest that the *SmydA*
469 genes may not represent redundant gene copies that are under pseudogenization.

470 Several members of the SMYD family of histone methyltransferases have
471 undergone a dramatic expansion in the insect lineage [27]. These SMYD genes were
472 identified as caste-specific genes in ants (*Harpegnathos saltator*), suggesting that
473 these histone modifiers play dedicated regulatory roles in insect phenotypic plasticity.
474 However, the biological significance of the differential expressions of these genes
475 remains unknown [40]. Our study further verified the presence of the differential
476 expression patterns of the SMYD genes in the four other insects that also possessed
477 adaptive phenotypic plasticity. Consequently, the understanding of convergent

Evolution of *SET* Genes in Insects

1 478 regulatory roles of the SMYD genes in insect phenotypic plasticity was extended.
2 479 Histone lysine methyltransferase catalyzes methyl group transfer to the amino group
3
4 480 of lysine residues of histones by means of the *SET* domain, a domain presented within
5
6 481 many proteins that regulate diverse development processes [41]. Histone lysine
7
8 482 methylation on specific residues is associated with distinct signatures of gene
9
10 483 expression, thereby serving as a chromatin modulator for epigenetic regulation [42].
11
12 484 Future studies should understand how the expanded SMYD gene family can quickly
13
14 485 become essential and identify the roles of the duplicated SMYD genes in insects,
15
16 486 despite the expectation of redundant functionality at the beginning of new duplicated
17
18 487 gene evolution [33].
19
20
21
22
23
24
25
26

27 489 **Materials and Methods**

28 490 **Identification of insect *SET* genes**

29 491 Genome assemblies and official gene sets of 130 insect species, including 62 dipteran
30
31 492 insects, 33 hymenopteran insects, 10 hemipteran insects, 7 coleopteran species, 9
32
33 493 lepidopteran insects, and representatives from Orthoptera, Phthiraptera,
34
35 494 Phasmatoptera, Trichoptera, Thysanoptera, Isoptera, Blattodea, Ephemeroptera and
36
37 495 Odonata, were obtained from their respective genome databases (Supplementary
38
39 496 Table S1). Among the basal arthropod species, we included 17 arthropod genomes
40
41 497 from 10 chelicerate species, five crustacean species and two non-insect hexapod
42
43 498 species.
44
45
46
47
48
49

50
51 499 The hidden Markov model-based HMMER program was used to identify the *SET*
52
53 500 domain containing proteins using PF00856 in the Pfam database [43, 44]. The
54
55 501 resulting genes with stop codons or frameshift mutations were subsequently manually
56
57
58
59

60
61
62
63
64
65
Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

502 checked. The obvious incorrect gene models were improved with transcriptome data
503 through the GeneWise program [45]. The PSILC program was used to identify the
504 potential pseudogenes [46]. Gene Ontology (GO) categories were determined via
505 scanning protein sequences against Interpro member databases using various
506 profile-based and hidden Markov models in the InterProScan package [47].

507

508 **Phylogenetic analysis, ancestral state reconstructions, and tests for selection**

509 Alignment-based methods using Bayesian inferences for *SET* domain sequences and
510 alignment-free methods based on feature frequency profiles for complete protein
511 sequences were used to infer phylogenetic relationships of *SET* genes across insects.
512 Multiple alignments were generated using the MAFFT alignment software [48].
513 According to the Akaike information criterion, the model of molecular evolution with
514 the best fit to the data was determined by using the ProtTest software [49]. Bayesian
515 reconstruction of phylogeny was conducted using the MrBayes 3.2.1 software for
516 10,000,000 generations [50]. The first 25% of the trees were discarded as burn-in. The
517 alignment-free and distance-based methods for phylogenetic tree building were
518 implemented by means of the feature frequency profile method
519 (<http://sourceforge.net/projects/ffp-phylogeny/>).

520 We constructed a character matrix that represents present/absent states for each
521 *SET* gene family to reconstruct the ancestral states of interior clades. The grouping of
522 the *SET* gene family was inferred from the OrthoMCL software with the
523 corresponding orthologous *SET* gene in *D. melanogaster*. Ancestral state
524 reconstruction was implemented in the Mesquite program using maximum likelihood
525 approaches (<http://mesquiteproject.org/>). This process requires a phylogeny tree of all

Evolution of *SET* Genes in Insects

1 526 the species studied. Single-copy orthologous gene families were inferred from the
2 527 core eukaryote CEGMA gene sets from each species using the OrthoMCL software
3
4 528 [51]. The resulting 455 single-copy orthologous gene families were used to construct
5
6
7 529 the species tree, which is consistent with the phylogenomic tree recently inferred from
8
9
10 530 transcriptome data [18].

11 531

15 532 **Expression of SMYD family genes in response to phenotypic plasticity**

17 533 The transcriptome data for gregarious and solitary locusts in *L. migratoria*, asexual
18
19 534 and sexual morphs in *A. pisum*, queens and workers in *A. mellifera*, and minor and
20
21
22 535 major workers in *A. echinator* were retrieved from the NCBI database under
23
24
25 536 accession numbers PRJNA79681, GSE56830, GSE61253, and GSE51576,
26
27
28 537 respectively. The raw reads were preprocessed to remove adapters and low-quality
29
30
31 538 bases using the Trimmomatic software; these reads were then mapped to the genome
32
33 539 assembly using the Tophat2 software [52, 53]. Raw counts of each gene were
34
35 540 calculated and annotated using the HT-seq package in Python, and the trimmed mean
36
37
38 541 of M value normalization method was used to normalize raw counts [54]. Differential
39
40 542 expression analysis was performed using the edgeR package at an FDR cut-off of 0.05
41
42 543 [55].

43 544

48 545 **Function approval of *SmydA-2* genes via experimental evidence**

49
50
51 546 A fluorescence *in situ* analysis of *SmydA-2* was performed on the brains and
52
53 547 integuments of locust nymphs. Biotin-labeled antisense and sense probes of *SmydA-2*
54
55 548 were produced from pGEM-T Easy plasmids (Promega) by using the T7/SP6 RNA

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
549 transcription system (Roche) following the manufacturer's protocol. The brains and
550 integuments were fixed in 4% paraformaldehyde overnight. The paraffin-embedded
551 slides (5 μm thick) were deparaffinized in xylene, rehydrated with an ethanol gradient,
552 digested with 20 $\mu\text{g}/\text{mL}$ proteinase K (Roche) at 37 $^{\circ}\text{C}$ for 15 min, and then incubated
553 with *SmydA-2* probe at 60 $^{\circ}\text{C}$ for 5 min. The slides were hybridized for 7–15 h at
554 37 $^{\circ}\text{C}$ and washed in 0.2 \times SSC and 2% BSA at 4 $^{\circ}\text{C}$ for 5 min. The biotin-labeled
555 probes of *SmydA-2* were detected with a streptavidin horseradish peroxidase
556 conjugate and fluorescein tyramide substrate using a TSA kit (Perkin Elmer). Images
557 for fluorescence signals were acquired using an LSM 710 confocal fluorescence
558 microscope (Zeiss).

24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
559 The recombinant proteins for *SmydA-2* and the negative controls of translation
560 system were produced using the TNT protein expression system (Promega). For *in*
561 *vitro* methyltransferase assay, 2 mg of unmodified histone H3 peptides (Sino
562 Biological) were incubated with 1 mg of recombinant protein and 0.1 mM
563 S-adenosyl-methionine (SAM, NEB) in a reaction buffer containing 50 mM Tris-HCl
564 (pH 8.0), 10% glycerol, 20 mM KCl, 5 mM MgCl_2 , 1 mM DTT, and 1 mM PMSF at
565 30 $^{\circ}\text{C}$ for 2 h. The reaction mixtures were subjected to electrophoresis on SDS-PAGE,
566 and the methylation activities were detected in Western blotting using anti-pan methyl
567 lysine antibody (Abcam). Anti-histone H3 (Abcam) was used as endogenous control
568 for protein samples.

48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
569 Locusts were reared in large, well-ventilated cages (40 cm \times 40 cm \times 40 cm) at a
570 density of 500–1000 insects per container. These colonies were reared under a 14:10
571 light/dark photo regime at 30 $^{\circ}\text{C}$ and were fed fresh wheat seedlings and bran.
572 Double-stranded RNAs of *SmydA-2* and green fluorescent protein (GFP) were

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

573 prepared using the T7 RiboMAX Express RNAi system (Promega) in accordance
574 with the manufacturer's protocols. Second-instar locusts were injected with
575 double-stranded RNAs in the second ventral segment of the abdomen. Total RNAs
576 were isolated using TRIzol reagent (Thermo Fisher Scientific) and then
577 reverse-transcribed into cDNA using M-MLV reverse transcriptase (Promega). The
578 mRNA levels were quantified using the SYBR Green expression assays on a
579 LightCycler 480 instrument (Roche). Survival data were analyzed using the
580 Kaplan–Meier method, and survival curves were compared using log-rank testing for
581 the *dsSmydA-2* and *dsGFP* curves.

582

583 **Signature of selection detected through likelihood ratio tests**

584 Protein sequences of *SET* genes were aligned with the MAFFT alignment software
585 [48] and the back-translated into corresponding nucleotide sequences. Gene
586 conversion was detected using the recombination detection program GENECONV
587 version 1.81a. To assess the contribution of natural selection during the diversification
588 of the *SET* gene family in insects, the ratios of nonsynonymous substitution per
589 nonsynonymous site (d_N) to synonymous substitution per synonymous site (d_S) across
590 the phylogenetic tree of the species were calculated using the software package
591 PAML version 4.48a [56]. The basic model M0 (null model) assumes the ratio $\omega =$
592 d_N/d_S is invariable (one-ratio model) among all branches examined, whereas the
593 alternative model allows the ω ratio to vary in different tree branches in the
594 phylogenetic tree [57, 58]. Likelihood ratio tests were applied to compare the null and
595 alternative models, which estimated ω ratio separately for different branches,
596 assuming a priori and the background branches. A significantly higher likelihood of
Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

597 the alternative model than the null model indicates a better fit to the data, indicating a
598 variation of selective pressures in different tree branches [57, 58].
599

600 **Declarations**

601 **List of abbreviations**

602 *SET* genes, *SET* domain-containing genes; E(z), Enhancer of zeste; LCA, last
603 common ancestor; GO, gene ontology; MYND, Myeloid translocation protein; qPCR,
604 quantitative real-time polymerase chain reaction; DEGs, differentially expressed
605 genes; FDR, false discovery rate; SAM, S-adenosyl-methionine; GFP, green
606 fluorescent protein; PP, posterior probability

607 **Ethics approval and consent to participate**

608 All animal procedures were licensed under the Institutional Animal Care and Use
609 Committee of the Institute of Zoology, Chinese Academy of Sciences.

610 **Consent for publication**

611 Not applicable

612 **Competing interests**

613 The authors declare they have no competing interests.

614 **Funding**

615 This research was supported by the following grants: National Natural Science
616 Foundation of China Grants (Nos. 31301915, 31472051, and 31672353) and Strategic
617 Priority Research Program of the Chinese Academy of Sciences (No. XDB11010000).

618 **Authors' contributions**

619 F.J., X.W., and L.K conceived and designed the experiments. F.J. and Q. L analyzed
620 and interpreted the data. F.J., Q. L., Y.W., J.Z., H.W., T.S., and M.Y. performed the
621 experiments. F.J., Q.L., and L.K wrote the paper.

622 **Acknowledgements**

623 The computational resources were provided by the Research Network of
624 Computational Biology and the Supercomputing Center at Beijing Institutes of Life
625 Science, Chinese Academy of Sciences.

626 **Availability of supporting data and materials**

627 The dataset supporting the conclusions of this article is available in
628 <http://159.226.67.242:8080/>.

630 **References**

- 631 1. Elofsson A, Sonnhammer EL: **A comparison of sequence and structure**
632 **protein domain families as a basis for structural genomics.** *Bioinformatics*
633 1999, **15**(6):480-500.
- 634 2. Itoh M, Nacher JC, Kuma K, Goto S, Kanehisa M: **Evolutionary history and**
635 **functional implications of protein domains and their combinations in**
636 **eukaryotes.** *Genome Biol* 2007, **8**(6):R121.
- 637 3. Sakarya O, Conaco C, Egecioglu O, Solla SA, Oakley TH, Kosik KS:
638 **Evolutionary expansion and specialization of the PDZ domains.** *Mol Biol*
639 *Evol* 2010, **27**(5):1058-1069.

58 Evolution of *SET* Genes in Insects

- 640 4. Johnson BR, Tsutsui ND: **Taxonomically restricted genes are associated**
641 **with the evolution of sociality in the honey bee.** *BMC Genomics* 2011,
642 **12**:164.
- 643 5. Qian W, Zhang J: **Genomic evidence for adaptation by gene duplication.**
644 *Genome Res* 2014, **24**(8):1356-1362.
- 645 6. Lynch M, Conery JS: **The evolutionary fate and consequences of duplicate**
646 **genes.** *Science* 2000, **290**(5494):1151-1155.
- 647 7. Innan H, Kondrashov F: **The evolution of gene duplications: classifying and**
648 **distinguishing between models.** *Nat Rev Genet* 2010, **11**(2):97-108.
- 649 8. Moore AD, Bornberg-Bauer E: **The dynamics and evolutionary potential of**
650 **domain loss and emergence.** *Mol Biol Evol* 2012, **29**(2):787-796.
- 651 9. Simola DF, Graham RJ, Brady CM, Enzmann BL, Desplan C, Ray A, Zwiebel
652 LJ, Bonasio R, Reinberg D, Liebig J *et al*: **Epigenetic (re)programming of**
653 **caste-specific behavior in the ant *Camponotus floridanus*.** *Science* 2016,
654 **351**(6268):aac6633.
- 655 10. Spannhoff A, Kim YK, Raynal NJ, Gharibyan V, Su MB, Zhou YY, Li J,
656 Castellano S, Sbardella G, Issa JP *et al*: **Histone deacetylase inhibitor**
657 **activity in royal jelly might facilitate caste switching in bees.** *EMBO Rep*
658 2011, **12**(3):238-243.
- 659 11. Simola DF, Ye C, Mutti NS, Dolezal K, Bonasio R, Liebig J, Reinberg D,
660 Berger SL: **A chromatin link to caste identity in the carpenter ant**
661 ***Camponotus floridanus*.** *Genome Res* 2013, **23**(3):486-496.
- 662 12. Dickman MJ, Kucharski R, Maleszka R, Hurd PJ: **Extensive histone**
663 **post-translational modification in honey bees.** *Insect Biochem Mol Biol*
664 2013, **43**(2):125-137.
- 665 13. Jenuwein T, Laible G, Dorn R, Reuter G: **SET domain proteins modulate**
666 **chromatin domains in eu- and heterochromatin.** *Cell Mol Life Sci* 1998,
667 **54**(1):80-93.
- 668 14. Dillon SC, Zhang X, Trievel RC, Cheng X: **The SET-domain protein**
669 **superfamily: protein lysine methyltransferases.** *Genome Biol* 2005,
670 **6**(8):227.
- 671 15. Boros IM: **Histone modification in *Drosophila*.** *Brief Funct Genomics* 2012,
672 **11**(4):319-331.
- 673 16. Herz HM, Garruss A, Shilatifard A: **SET for life: biochemical activities and**
674 **biological functions of SET domain-containing proteins.** *Trends Biochem*
675 *Sci* 2013, **38**(12):621-639.
- 676 17. Jenuwein T: **The epigenetic magic of histone lysine methylation.** *FEBS J*
677 2006, **273**(14):3121-3135.
- 678 18. Misof B, Liu S, Meusemann K, Peters RS, Donath A, Mayer C, Frandsen PB,
679 Ware J, Flouri T, Beutel RG *et al*: **Phylogenomics resolves the timing and**
680 **pattern of insect evolution.** *Science* 2014, **346**(6210):763-767.
- 681 19. Ferguson LC, Green J, SurrIDGE A, Jiggins CD: **Evolution of the insect**
682 **yellow gene family.** *Mol Biol Evol* 2011, **28**(1):257-272.
- 683 20. Helmkampf M, Cash E, Gadau J: **Evolution of the insect desaturase gene**
684 **family with an emphasis on social Hymenoptera.** *Mol Biol Evol* 2015,
685 **32**(2):456-471.
- 686 21. Tanaka K, Diekmann Y, Hazbun A, Hijazi A, Vreede B, Roch F, Sucena E:
687 **Multispecies Analysis of Expression Pattern Diversification in the**

Evolution of *SET* Genes in Insects

- 688 **Recently Expanded Insect Ly6 Gene Family.** *Mol Biol Evol* 2015,
689 **32(7):1730-1747.**
- 690 22. Urena E, Pirone L, Chafino S, Perez C, Sutherland JD, Lang V, Rodriguez
691 MS, Lopitz-Otsoa F, Blanco FJ, Barrio R *et al*: **Evolution of SUMO**
692 **Function and Chain Formation in Insects.** *Mol Biol Evol* 2016,
693 **33(2):568-584.**
- 694 23. Benton R: **Multigene Family Evolution: Perspectives from Insect**
695 **Chemoreceptors.** *Trends Ecol Evol* 2015, **30(10):590-600.**
- 696 24. Zhang L, Ma H: **Complex evolutionary history and diverse domain**
697 **organization of SET proteins suggest divergent regulatory interactions.**
698 *The New phytologist* 2012, **195(1):248-263.**
- 699 25. Vidal NM, Graziotin AL, Iyer LM, Aravind L, Venancio TM: **Transcription**
700 **factors, chromatin proteins and the diversification of Hemiptera.** *Insect*
701 *Biochem Mol Biol* 2016, **69:1-13.**
- 702 26. Rider SD, Jr., Srinivasan DG, Hilgarth RS: **Chromatin-remodelling proteins**
703 **of the pea aphid, Acyrthosiphon pisum (Harris).** *Insect Mol Biol* 2010, **19**
704 **Suppl 2:201-214.**
- 705 27. Bonasio R, Zhang G, Ye C, Mutti NS, Fang X, Qin N, Donahue G, Yang P, Li
706 Q, Li C *et al*: **Genomic comparison of the ants Camponotus floridanus and**
707 **Harpegnathos saltator.** *Science* 2010, **329(5995):1068-1071.**
- 708 28. Calpena E, Palau F, Espinos C, Galindo MI: **Evolutionary History of the**
709 **Smyd Gene Family in Metazoans: A Framework to Identify the Orthologs**
710 **of Human Smyd Genes in Drosophila and Other Animal Species.** *PLoS*
711 *One* 2015, **10(7):e0134106.**
- 712 29. Wang X, Fang X, Yang P, Jiang X, Jiang F, Zhao D, Li B, Cui F, Wei J, Ma C
713 *et al*: **The locust genome provides insight into swarm formation and**
714 **long-distance flight.** *Nat Commun* 2014, **5:2957.**
- 715 30. Thompson EC, Travers AA: **A Drosophila Smyd4 homologue is a**
716 **muscle-specific transcriptional modulator involved in development.** *PLoS*
717 *One* 2008, **3(8):e3008.**
- 718 31. Kaessmann H: **Origins, evolution, and phenotypic impact of new genes.**
719 *Genome Res* 2010, **20(10):1313-1326.**
- 720 32. Assis R, Bachtrog D: **Neofunctionalization of young duplicate genes in**
721 **Drosophila.** *Proc Natl Acad Sci U S A* 2013, **110(43):17409-17414.**
- 722 33. Kemkemer C, Long M: **New genes important for development.** *EMBO Rep*
723 2014, **15(5):460-461.**
- 724 34. Kenny NJ, Chan KW, Nong W, Qu Z, Maeso I, Yip HY, Chan TF, Kwan HS,
725 Holland PW, Chu KH *et al*: **Ancestral whole-genome duplication in the**
726 **marine chelicerate horseshoe crabs.** *Heredity (Edinb)* 2016, **116(2):190-199.**
- 727 35. Kaessmann H, Vinckenbosch N, Long M: **RNA-based gene duplication:**
728 **mechanistic and evolutionary insights.** *Nat Rev Genet* 2009, **10(1):19-31.**
- 729 36. Hedges SB, Dudley J, Kumar S: **TimeTree: a public knowledge-base of**
730 **divergence times among organisms.** *Bioinformatics* 2006,
731 **22(23):2971-2972.**
- 732 37. Tan X, Rotllant J, Li H, De Deyne P, Du SJ: **SmyD1, a histone**
733 **methyltransferase, is required for myofibril organization and muscle**
734 **contraction in zebrafish embryos.** *Proc Natl Acad Sci U S A* 2006,
735 **103(8):2713-2718.**

Evolution of *SET* Genes in Insects

- 736 38. Domazet-Lošo T, Tautz D: **An evolutionary analysis of orphan genes in**
737 ***Drosophila***. *Genome Res* 2003, **13**(10):2213-2219.
- 738 39. Pegueroles C, Laurie S, Alba MM: **Accelerated evolution after gene**
739 **duplication: a time-dependent process affecting just one copy**. *Mol Biol*
740 *Evol* 2013, **30**(8):1830-1842.
- 741 40. Bonasio R: **The role of chromatin and epigenetics in the polyphenisms of**
742 **ant castes**. *Brief Funct Genomics* 2014, **13**(3):235-245.
- 743 41. Leinhart K, Brown M: **SET/MYND Lysine Methyltransferases Regulate**
744 **Gene Transcription and Protein Activity**. *Genes (Basel)* 2011,
745 **2**(1):210-218.
- 746 42. Sims RJ, 3rd, Nishioka K, Reinberg D: **Histone lysine methylation: a**
747 **signature for chromatin function**. *Trends Genet* 2003, **19**(11):629-639.
- 748 43. Finn RD, Mistry J, Tate J, Coghill P, Heger A, Pollington JE, Gavin OL,
749 Gunasekaran P, Ceric G, Forslund K *et al*: **The Pfam protein families**
750 **database**. *Nucleic Acids Res* 2010, **38**(Database issue):D211-222.
- 751 44. Finn RD, Clements J, Eddy SR: **HMMER web server: interactive sequence**
752 **similarity searching**. *Nucleic Acids Res* 2011, **39**(Web Server issue):W29-37.
- 753 45. Birney E, Clamp M, Durbin R: **GeneWise and Genomewise**. *Genome Res*
754 2004, **14**(5):988-995.
- 755 46. Coin L, Durbin R: **Improved techniques for the identification of**
756 **pseudogenes**. *Bioinformatics* 2004, **20** Suppl 1:i94-100.
- 757 47. Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, McWilliam H,
758 Maslen J, Mitchell A, Nuka G *et al*: **InterProScan 5: genome-scale protein**
759 **function classification**. *Bioinformatics* 2014, **30**(9):1236-1240.
- 760 48. Katoh K, Asimenos G, Toh H: **Multiple alignment of DNA sequences with**
761 **MAFFT**. *Methods Mol Biol* 2009, **537**:39-64.
- 762 49. Darriba D, Taboada GL, Doallo R, Posada D: **ProtTest 3: fast selection of**
763 **best-fit models of protein evolution**. *Bioinformatics* 2011, **27**(8):1164-1165.
- 764 50. Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, Höhna S,
765 Larget B, Liu L, Suchard MA, Huelsenbeck JP: **MrBayes 3.2: efficient**
766 **Bayesian phylogenetic inference and model choice across a large model**
767 **space**. *Syst Biol* 2012, **61**(3):539-542.
- 768 51. Li L, Stoeckert CJ, Jr., Roos DS: **OrthoMCL: identification of ortholog**
769 **groups for eukaryotic genomes**. *Genome Res* 2003, **13**(9):2178-2189.
- 770 52. Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for**
771 **Illumina sequence data**. *Bioinformatics* 2014, **30**(15):2114-2120.
- 772 53. Trapnell C, Pachter L, Salzberg SL: **TopHat: discovering splice junctions**
773 **with RNA-Seq**. *Bioinformatics* 2009, **25**(9):1105-1111.
- 774 54. Robinson MD, Oshlack A: **A scaling normalization method for differential**
775 **expression analysis of RNA-seq data**. *Genome Biol* 2010, **11**(3):R25.
- 776 55. Robinson MD, McCarthy DJ, Smyth GK: **edgeR: a Bioconductor package**
777 **for differential expression analysis of digital gene expression data**.
778 *Bioinformatics* 2010, **26**(1):139-140.
- 779 56. Yang Z: **PAML 4: phylogenetic analysis by maximum likelihood**. *Mol Biol*
780 *Evol* 2007, **24**(8):1586-1591.
- 781 57. Yang Z: **Likelihood ratio tests for detecting positive selection and**
782 **application to primate lysozyme evolution**. *Mol Biol Evol* 1998,
783 **15**(5):568-573.

Evolution of *SET* Genes in Insects

784 58. Yang Z, Nielsen R: **Synonymous and nonsynonymous rate variation in**
785 **nuclear genes of mammals.** *J Mol Evol* 1998, **46**(4):409-418.

786

787

788 **Figures**

789 **Figure 1. Phylogenetic analysis of *SET* genes in insects.** A phylogeny using
790 Bayesian inference is generated from the domain protein sequence of *SET* genes. One
791 representative is elected for each order. The protein domains, which are labeled with
792 different colors based on the domain type, are shown in the exterior circle of the
793 phylogenetic tree.

794

795 **Figure 2. Diversification of arthropod-specific *SET* genes.** (A) Inference of
796 ancestral sets of *SET* homologous groups along the evolution of insects. The gains and
797 losses of *SET* homologous groups are indicated in the internal nodes of the
798 phylogenetic tree. The number in parentheses indicates the number of species in each
799 order. (B) Distribution pattern of *SET* genes in arthropod orders. One representative is
800 elected for each order. Red color indicates presence of *SET* genes, and blue color
801 indicates absence of *SET* genes.

802

803 **Figure 3. Evolution of *SmydA* genes in insects.** (A) Gene ontology categories of the
804 conserved and arthropod-specific groups of *SET* genes. The gene ontology categories,
805 which are only present in the arthropod-specific group, are highlighted in red. (B)
806 Phylogenetic tree of the SMYD gene family of the representative species selected
807 from each order. The phylogenetic tree is constructed using the Bayesian inference
808 method. The Bayesian posterior probability (PP) values are indicated only for the
Evolution of *SET* Genes in Insects

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

809 internal nodes to improve clarity; consequently, the *SET* genes are grouped into
 810 different monophyletic clades (SMYD subfamilies). Red and orange circles indicate
 811 $PP > 90\%$ and $PP > 70\%$, respectively. (C) Conserved syntenies for *SmydA* genes in
 812 four holometabolous species. (D) Distributions of ω ($\omega = d_N/d_S$ ratio) values of the
 813 conserved SMYD and *SmydA* groups of *SET* genes.

814

815 **Figure 4. Function approval of *SmydA-2* genes through experimental evidence.**

816 (A) *In vitro* methyltransferase assay of histone H3 of *SmydA-2* in locusts. Anti-pan
 817 methyl lysine antibody recognizes histone H3 *in vitro* methylated with *SmydA-2*.
 818 Anti-histone H3 serves as endogenous control for protein samples. The analyses were
 819 carried out in three replicates. $**P < 0.01$. (B) Expression evidence of *SmydA-2* in the
 820 brain and cuticle of locusts via fluorescence *in situ* hybridization analysis. Green
 821 signals indicate the expression of *SmydA-2* /control, and blue signals indicate nuclear
 822 staining with Hoechst. (C) Relative gene expression of *SmydA-2* in the different
 823 tissues. mRNA levels are quantified using the SYBR Green expression assays on a
 824 LightCycler 480 instrument. The qPCR data are shown as the mean \pm SEM ($n = 6$).
 825 (D) Survival analysis of the locusts after *SmydA-2* double-strand RNA injection. Data
 826 are analyzed through the Kaplan–Meier survival curve comparison of the *dsSmydA-2*
 827 and *dsGFP* groups for three replicates.

828

829 **Figure 5. Differential expression analysis in insects showing phenotype plasticity.**

830 Alternative phenotype includes gregarious and solitary phases in *Locusta migratoria*
 831 (LOCMI), asexual and sexual morphs in *Acyrtosiphon pisum* (ACYPI), queens and

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

832 workers in *Apis mellifera* (APIME), and minor workers and major workers in
833 *Acromyrmex echinator* (ACREC).

834 **Tables**

835 **Table 1. Summary of *SET* genes in insect genomes.**

836 **Table 2. Tests of rate heterogeneity acting on *SET* genes in insects.**

837 **Table 3. Signatures of selection acting on differentially expressed *SET* genes in**
838 **response to phenotypic plasticity.**

839

840 **Supplementary Data**

841 **Supplementary Table S1. The arthropod genome data involved in this study.**

842 **Supplementary Table S2. *SET* genes in the 147 arthropod genomes.**

843 **Supplementary Table S3. Summary of *SET* genes in the 147 arthropod genomes.**

844 **Supplementary Figure S1. Phylogenetic analysis of the *SET* genes in Lepidoptera**
845 **using Bayesian inferences.** The *SET* gene families labeled with different colors are
846 shown in the exterior circle of the phylogenetic tree. The insect species involved are
847 represented with different colors of the external branch.

848 **Supplementary Figure S2. Phylogenetic analysis of the *SET* genes in Diptera**
849 **using Bayesian inferences.** The *SET* gene families labeled with different colors are
850 shown in the exterior circle of the phylogenetic tree. The insect species involved are
851 represented with different colors of the external branch. The representative species are
852 selected to improve clarity.

853 **Supplementary Figure S3. Phylogenetic analysis of the *SET* genes in Hemiptera**
854 **using Bayesian inferences.** The *SET* gene families labeled with different colors are

Evolution of *SET* Genes in Insects

1 855 shown in the exterior circle of the phylogenetic tree. The insect species involved are
2 856 represented with different colors of the external branch.

3
4 857 **Supplementary Figure S4. Phylogenetic analysis of the *SET* genes in**
5
6
7 858 **Hymenoptera using Bayesian inferences.** The *SET* gene families labeled with
8
9 859 different colors are shown in the exterior circle of the phylogenetic tree. The insect
10
11 860 species involved are represented with different colors of the external branch. The
12
13 861 representative species are selected to improve clarity.

14
15
16 862 **Supplementary Figure S5. Phylogenetic analysis of the *SET* genes in Coleopteran**
17
18 863 **using Bayesian inferences.** The *SET* gene families labeled with different colors are
19
20 864 shown in the exterior circle of the phylogenetic tree. The insect species involved are
21
22 865 represented with different colors of the external branch.

23
24
25 866 **Supplementary Figure S6. Effects of RNA interference of the mRNA expression**
26
27 867 **levels of *SmydA-2* in locust brains.** The locusts are injected with double-stranded
28
29 868 RNAs into the second ventral segment of the abdomen. Due to the systemic RNA
30
31 869 interference in locusts, the brain, which is spatially distant from the abdomen, is used
32
33 870 in qPCR assays to guarantee effective expression knockdown. qPCR data are shown
34
35 871 as the mean \pm SEM (n = 6). **P < 0.01.

36
37 872 **Supplementary Figure S7. Tree topology and branch labeling for tests of**
38
39 873 **selection on *SET* genes.** APIME, *Apis mellifera*; ACREC, *Acromyrmex echinator*;
40
41 874 LOCMI, *Locusta migratoria*. Supplementary Table S1 presents the abbreviation of
42
43 875 insect species.

44
45
46
47
48
49
50
51 876

Table 1. Summary of *SET* genes in insect genomes.

Order	Genus	SMYD	SETD	PRDM	Ash	Suv	Trx	Ez	AS	Total
Coleoptera	<i>Agrilus</i> (1)	4	1	2	3	3	3	1	9	26
Coleoptera	<i>Anoplophora</i> (1)	7	1	2	3	3	3	2	7	28
Coleoptera	<i>Dendroctonus</i> (1)	5	1	1	3	3	3	1	12	29
Coleoptera	<i>Leptinotarsa</i> (1)	10	1	1	2	5	3	1	9	32
Coleoptera	<i>Onthophagus</i> (1)	4	1	1	3	4	3	1	10	27
Coleoptera	<i>Oryctes</i> (1)	6	1	1	3	3	1	1	9	25
Coleoptera	<i>Tribolium</i> (1)	6	2	1	3	3	3	1	15	34
Phthiraptera	<i>Pediculus</i> (1)	6	1	1	3	4	3	1	9	28
Blattodea	<i>Blattella</i> (1)	4	2	2	4	3	2	1	7	25
Diptera	<i>Aedes</i> (2)	11-12	1	2	3-4	2-3	3-4	1-2	11-12	34-38
Diptera	<i>Anopheles</i> (19)	6-19	1	1-2	1-3	2-3	2-3	1	4-11	20-37
Diptera	<i>Bactrocera</i> (2)	4-5	1	1-2	3-4	4	3-6	1-2	13-22	31-45
Diptera	<i>Ceratina</i> (1)	5	1	1	2	4	3	1	11	28
Diptera	<i>Ceratitis</i> (1)	5	1	1	3	3	3	1	14	31
Diptera	<i>Culex</i> (1)	40	1	1	13	2	9	1	14	81
Diptera	<i>Drosophila</i> (22)	4-5	1	1	3-4	3-5	2-4	1	7-14	24-31
Diptera	<i>Glossina</i> (6)	4-5	1	1	3-4	2-5	3-4	1	12-15	29-34
Diptera	<i>Lucilia</i> (1)	5	1	1	3	3	3	1	12	29
Diptera	<i>Lutzomyia</i> (1)	6	1	1	3	3	2	1	10	27
Diptera	<i>Mayetiola</i> (1)	13	1	1	9	6	4	1	25	60
Diptera	<i>Megaselia</i> (1)	2	1	1	3	2	1	1	5	16
Diptera	<i>Musca</i> (1)	5	1	1	3	3	3	1	20	37
Diptera	<i>Phlebotomus</i> (1)	5	1	1	4	3	3	1	6	24
Diptera	<i>Belgica</i> (1)	27	2	1	3	5	4	1	12	55
Diptera	<i>Stomoxys</i> (1)	5	1	1	3	2	3	1	16	32
Ephemeroptera	<i>Ephemer</i> (1)	18	1	1	3	3	2	1	12	42
Hemiptera	<i>Acyrtosiphon</i> (1)	14	1	0	2	10	4	1	31	63
Hemiptera	<i>Cimex</i> (1)	4	1	2	3	5	3	1	5	24
Hemiptera	<i>Diaphorina</i> (1)	3	1	1	4	4	3	2	11	29
Hemiptera	<i>Gerris</i> (1)	6	1	1	3	3	3	1	8	26
Hemiptera	<i>Halyomorpha</i> (1)	5	1	1	2	5	3	1	8	26
Hemiptera	<i>Homalodisca</i> (1)	5	2	2	2	5	4	1	8	29
Hemiptera	<i>Nilaparvata</i> (1)	4	1	6	2	4	4	1	7	29
Hemiptera	<i>Oncopeltus</i> (1)	6	1	1	2	5	4	1	7	27
Hemiptera	<i>Pachypsylla</i> (1)	1	1	2	2	3	1	1	9	20
Hemiptera	<i>Rhodnius</i> (1)	6	1	1	2	2	2	1	6	21
Hymenoptera	<i>Acromyrmex</i> (1)	7	2	1	3	3	3	1	7	27
Hymenoptera	<i>Apis</i> (3)	6-7	1	1	3	3-4	1-3	1	7-9	22-29
Hymenoptera	<i>Athalia</i> (1)	7	1	2	2	3	2	1	8	26
Hymenoptera	<i>Atta</i> (1)	8	1	1	3	4	3	1	7	28
Hymenoptera	<i>Bombus</i> (2)	7-8	1	1	3	4	3	1	8-10	29-30
Hymenoptera	<i>Camponotus</i> (1)	8	2	1	2	3	2	1	8	27
Hymenoptera	<i>Cardiocondyla</i> (1)	7	2	1	3	4	3	1	10	31
Hymenoptera	<i>Cephus</i> (1)	6	1	1	2	3	2	1	6	22
Hymenoptera	<i>Cerapachys</i> (1)	5	1	1	2	3	3	1	6	22
Hymenoptera	<i>Ceratosolen</i> (1)	8	1	1	3	3	2	1	9	28
Hymenoptera	<i>Copidosoma</i> (1)	17	1	1	3	4	2	1	16	45
Hymenoptera	<i>Dufourea</i> (1)	7	2	1	3	4	3	1	7	28

Hymenoptera	<i>Eufriesea</i> (1)	6	2	1	3	4	3	1	8	28
Hymenoptera	<i>Fopius</i> (1)	9	1	1	3	4	1	1	9	29
Hymenoptera	<i>Habropoda</i> (1)	8	2	1	3	4	3	1	8	30
Hymenoptera	<i>Harpegnathos</i> (1)	8	2	0	1	2	1	1	8	23
Hymenoptera	<i>Linepithema</i> (1)	7	2	1	3	4	3	1	8	29
Hymenoptera	<i>Megachile</i> (1)	7	2	1	3	3	3	1	8	28
Hymenoptera	<i>Melipona</i> (1)	7	2	1	3	4	3	1	8	29
Hymenoptera	<i>Microplitis</i> (1)	18	1	1	3	4	3	2	8	40
Hymenoptera	<i>Monomorium</i> (1)	6	1	1	2	3	2	1	5	21
Hymenoptera	<i>Nasonia</i> (1)	17	1	1	3	4	2	1	23	52
Hymenoptera	<i>Orussus</i> (1)	11	2	1	2	3	3	1	7	30
Hymenoptera	<i>Pogonomyrmex</i> (1)	5	2	1	2	4	3	1	8	26
Hymenoptera	<i>Polistes</i> (1)	6	1	1	1	4	2	1	6	22
Hymenoptera	<i>Solenopsis</i> (1)	2	1	1	3	3	3	1	7	21
Hymenoptera	<i>Trichogramma</i> (1)	15	1	1	3	4	1	1	26	52
Hymenoptera	<i>Vollenhovia</i> (1)	6	1	1	3	4	2	1	3	21
Hymenoptera	<i>Hypothenemus</i> (1)	9	1	1	3	3	3	1	8	29
Hymenoptera	<i>Wasmannia</i> (1)	7	1	1	3	3	3	1	6	25
Isoptera	<i>Zootermopsis</i> (2)	6	1	2	2	4	3	1	10	29
Lepidoptera	<i>Bombyx</i> (1)	4	2	1	3	4	3	1	8	26
Lepidoptera	<i>Danaus</i> (1)	5	1	1	3	5	3	1	10	29
Lepidoptera	<i>Heliconius</i> (1)	5	1	1	2	4	3	1	6	23
Lepidoptera	<i>Papilio</i> (2)	6	1	1	3	2-4	2	1	9-11	26-27
Lepidoptera	<i>Lerema</i> (1)	4	1	2	3	3	3	1	10	27
Lepidoptera	<i>Melitaea</i> (1)	5	1	1	3	1	3	1	8	23
Lepidoptera	<i>Manduca</i> (1)	6	2	7	7	5	5	2	29	63
Lepidoptera	<i>Plutella</i> (1)	5	4	1	4	5	6	0	13	38
Odonata	<i>Ladona</i> (1)	3	2	2	3	4	3	1	9	27
Orthoptera	<i>Locusta</i> (1)	9	1	1	3	4	3	1	7	29
Phasmatoptera	<i>Timema</i> (1)	3	1	1	3	5	3	1	6	23
Thysanoptera	<i>Frankliniella</i> (1)	6	2	8	3	5	3	1	21	49
Trichoptera	<i>Limnephilus</i> (1)	3	1	0	2	3	2	1	6	18

Note: The numbers in parenthesis indicates the number of the species in each genus. The dash is used to represent the range of *SET* gene number in each genus. AS, arthropod-specific.

Table 2. Tests of rate heterogeneity acting on *SET* genes in insects.

	Gene	One Ratio Likelihood	One Ratio ω	Free Ratio Likelihood	df	<i>P</i>
	<i>Smyd3</i>	-4833.870633	0.055	-4833.870633	16	<0.001
SMYD	<i>Smyd4-1</i>	-17270.85481	0.1627	-17140.2931	58	<0.001
	<i>Smyd4-2</i>	-13187.36796	0.1125	-13112.10598	44	<0.001
	<i>Smyd4-3</i>	-20488.96316	0.1069	-20364.99139	66	<0.001
	<i>Smyd4-4</i>	-15552.36608	0.1112	-15475.97917	44	<0.001
	<i>Smyd5</i>	-21495.43548	0.0633	-21329.01303	64	<0.001
	<i>upSET(MLL5)</i>	-7286.598116	0.0103	-7247.800191	62	0.087
	<i>Set8</i>	-6450.096636	0.0321	-6386.997507	60	<0.001
	<i>Hmt4-20</i>	-3523.660744	0.0079	-3478.339497	56	<0.001
SETD	<i>SETD</i>	-9030.115692	0.033	-9009.972504	34	0.212
PRDM	<i>Blimp-1</i>	-2679.981724	0.0051	-2664.129882	52	0.988
	<i>Mes-4</i>	-5530.425067	0.0163	-5504.225668	56	0.612
Ash	<i>ash1</i>	-4995.315864	0.0122	-4947.987993	60	<0.001
	<i>Set2</i>	-5636.021533	0.0118	-5570.266003	60	<0.001
	<i>Su(var)3-9</i>	-4351.473377	0.0212	-4308.872564	32	<0.001
Suv	<i>egg</i>	-15308.27271	0.0624	-15214.54477	54	<0.001
	<i>CG4565</i>	-7168.675146	0.056	-7114.254055	46	<0.001
	<i>G9a</i>	-4641.585219	0.0091	-4604.810574	54	0.040
	<i>trx</i>	-3897.22035	0.0031	-3877.624919	58	0.972
Trx	<i>Set1</i>	-3733.003015	0.0026	-3700.07484	60	0.281
	<i>trr</i>	-4549.712	0.0114	-4471.116449	60	<0.001
E(z)	<i>Ez</i>	-3368.302419	0.0007	-3355.922925	61	1.000
	<i>SmydA-1</i>	-10066.85883	0.0904	-9995.276076	34	<0.001
	<i>SmydA-2</i>	-11858.79656	0.0052	-11812.61641	30	<0.001
	<i>SmydA-3</i>	-13902.68842	0.0817	-13842.81154	56	<0.001
SMYDA	<i>SmydA-4</i>	-9602.742487	0.0254	-9583.599425	26	0.057
	<i>SmydA-5</i>	-13748.76916	0.1179	-13656.26849	50	<0.001
	<i>SmydA-6</i>	-12142.19779	0.1623	-12043.99319	42	<0.001
	<i>SmydA-9</i>	-13258.40628	0.1357	-13193.53611	52	<0.001

Note: Accounting for the unequal genome sequencing efforts between different insect families, we selected one species within each genus to be representative of the genus.

Table 3. Signatures of selection acting on differential expressed *SET* genes in response to phenotypic plasticity.

Model-Parameters	APIME		LOCMI	ACREC		
	<i>SmydA-3</i>	<i>SmydA-5</i>	<i>SmydA-1</i>	<i>SmydA-3</i>	<i>SmydA-5</i>	<i>SmydA-9</i>
Basic models						
M0: ω	0.082	0.118	0.090	0.082	0.118	0.136
Branch models						
B0: lnL	-13914.741	-13749.007	-10088.904	-13905.140	-13749.047	-13259.370
B0: ω_0 ($\omega_1 = 1$)	0.077	0.113	0.090	0.081	0.117	0.135
BA: lnL	-13901.138	-13745.405	-10056.182	-13901.922	-13748.719	-13258.338
BA: ω_0, ω_1	0.080, 0.142	0.115, 0.313	0.095, 0.003	0.081, 0.177	0.118, 0.181	0.135, 0.186
Branch-site models						
A0: p_{2a} ($\omega_2 = 1$)	0.078	0.059	0.111	0.082	0.155	0.096
AA: p_{2a}, ω_2	0.078, 1.000	0.025, 3.102	0.109, 8.895	0.082, 1.000	0.155, 1.000	0.011, 19.742
Positively selected sites (BEB)	5 M 11 K 93 P 105 C					
LRT, P						
M0 versus BA	0.078	0.009	<0.001	0.216	0.752	0.712
BA versus B0	<0.001	0.007	<0.001	0.011	0.418	0.151
A0 versus AA	1.000	0.802	0.022	1.000	1.000	0.082

ω , the ratios of nonsynonymous substitution per nonsynonymous site to synonymous substitution per synonymous site; ω_0, ω_1 , background and foreground ω values, respectively; APIME, *Apis mellifera*; ACREC, *Acromyrmex echinator*; LOCMI, *Locusta migratoria*.

Figure 1

[Click here to download Figure figure1_20161006.pdf](#)

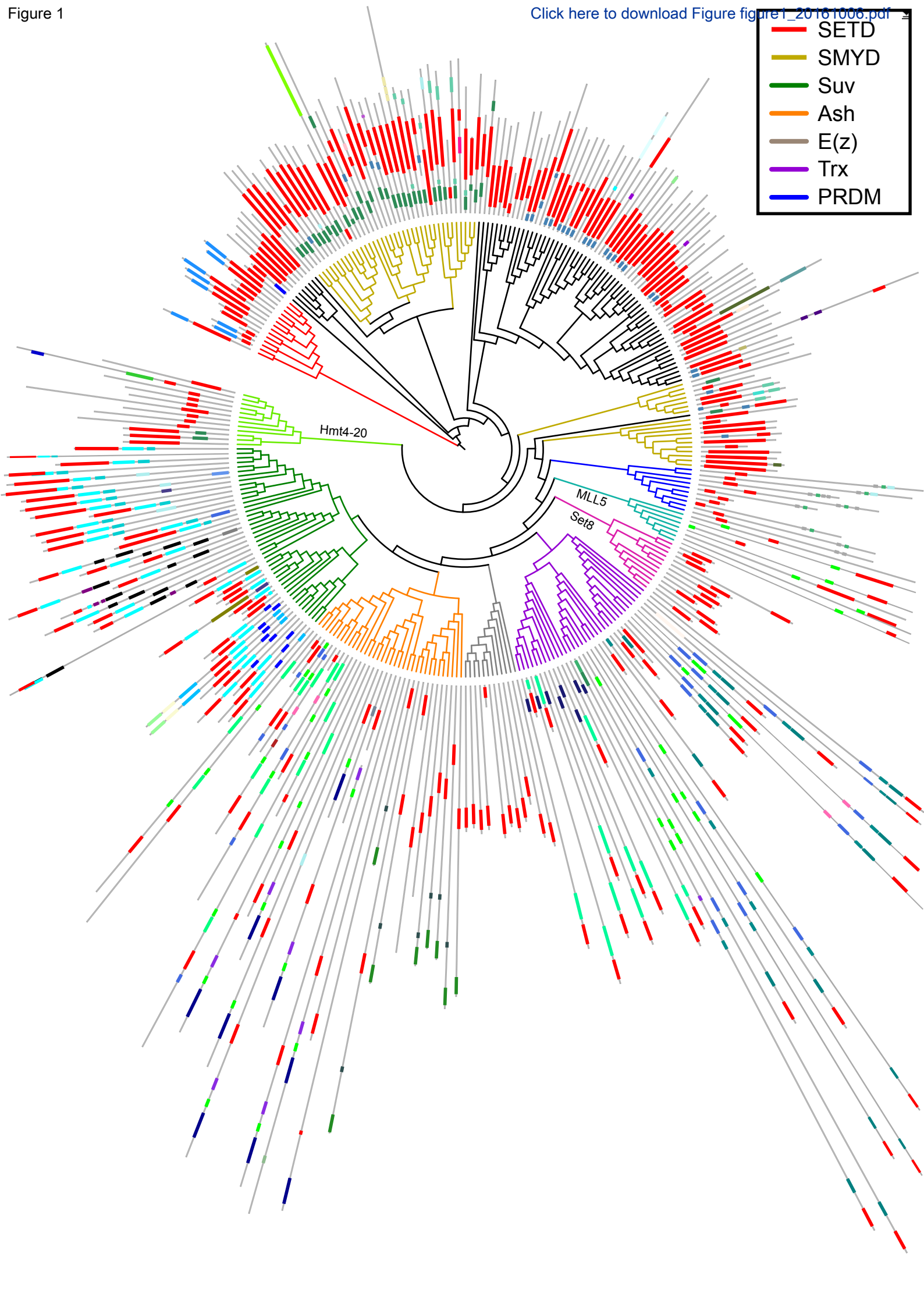
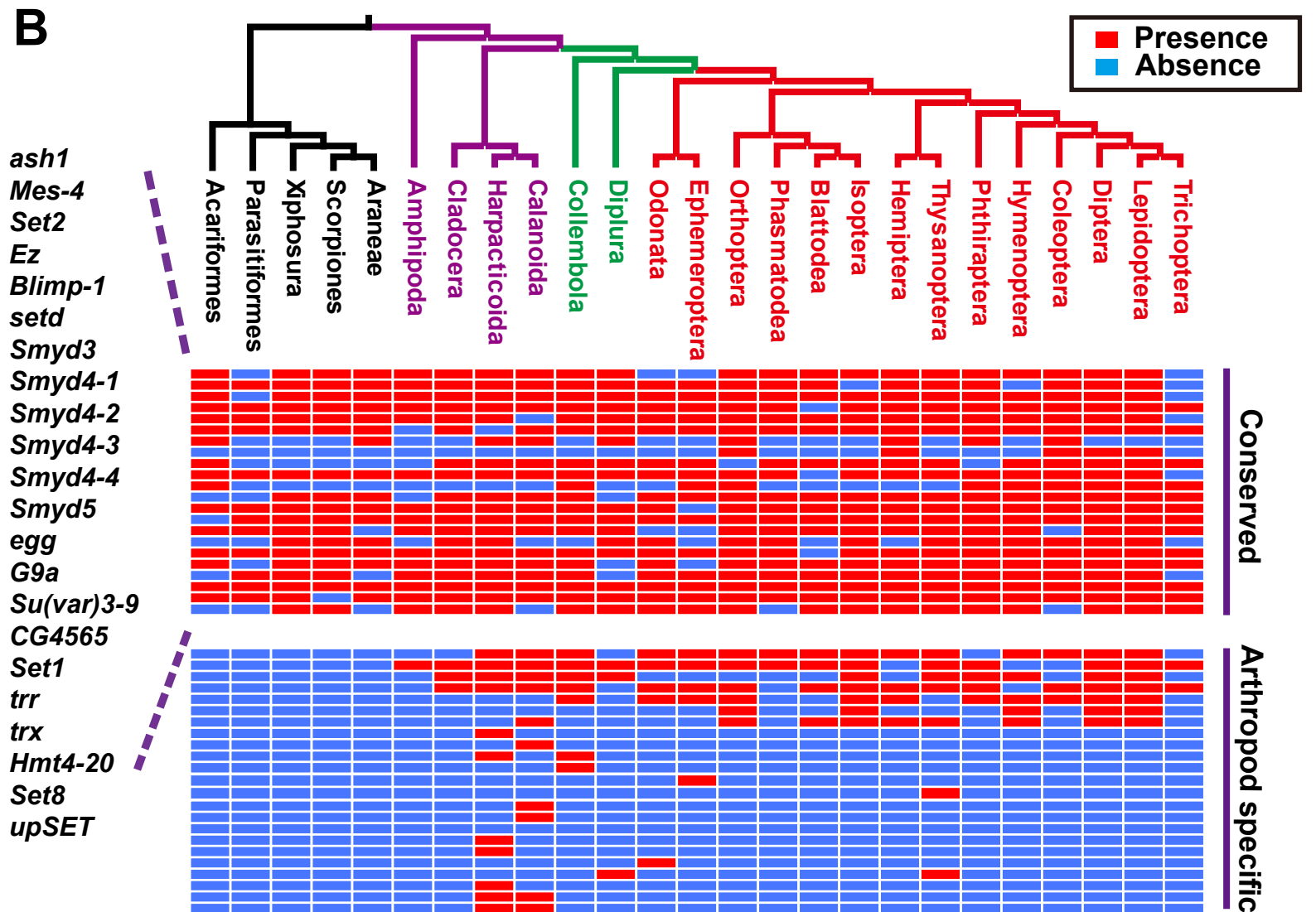
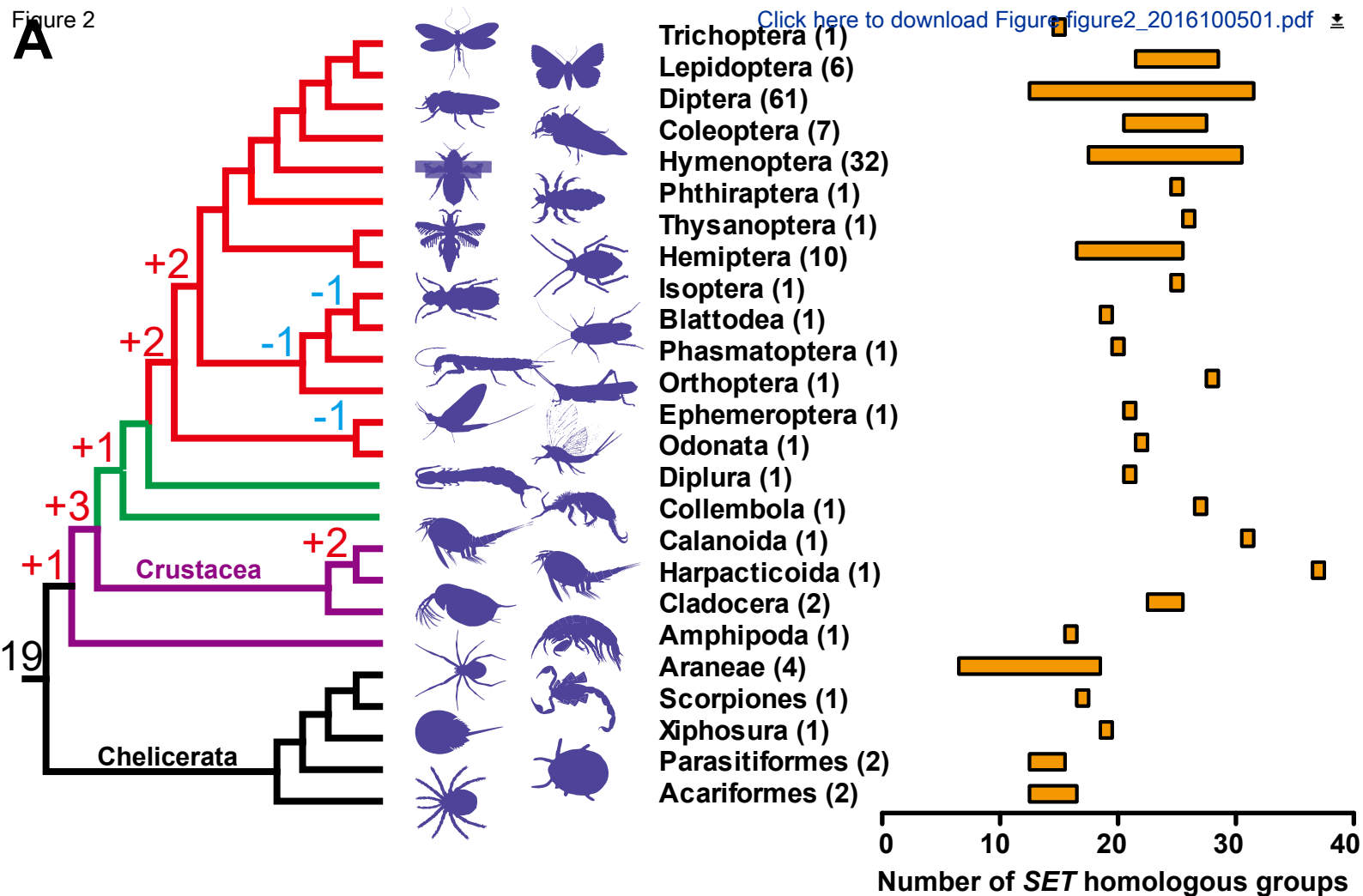
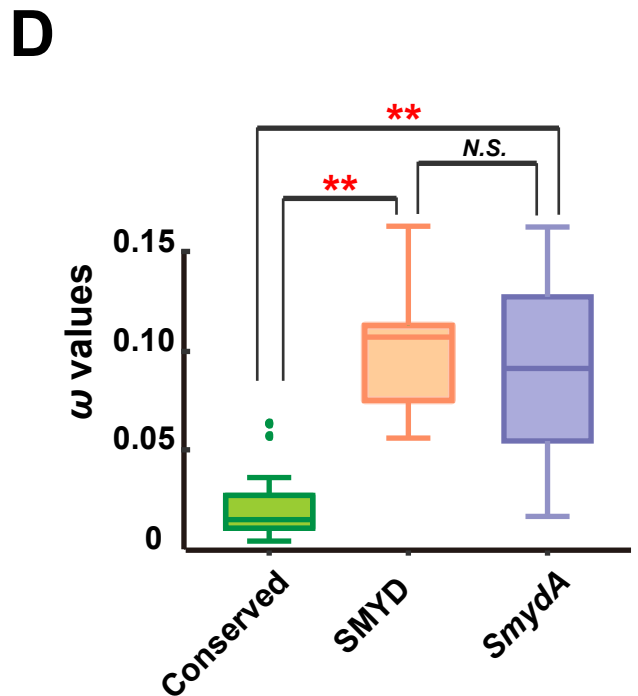
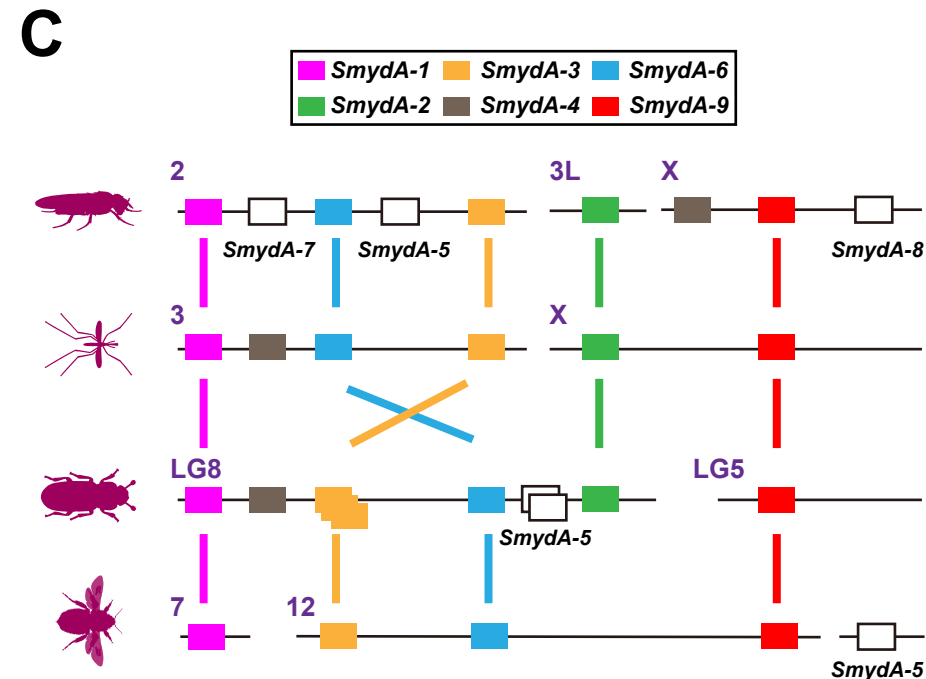
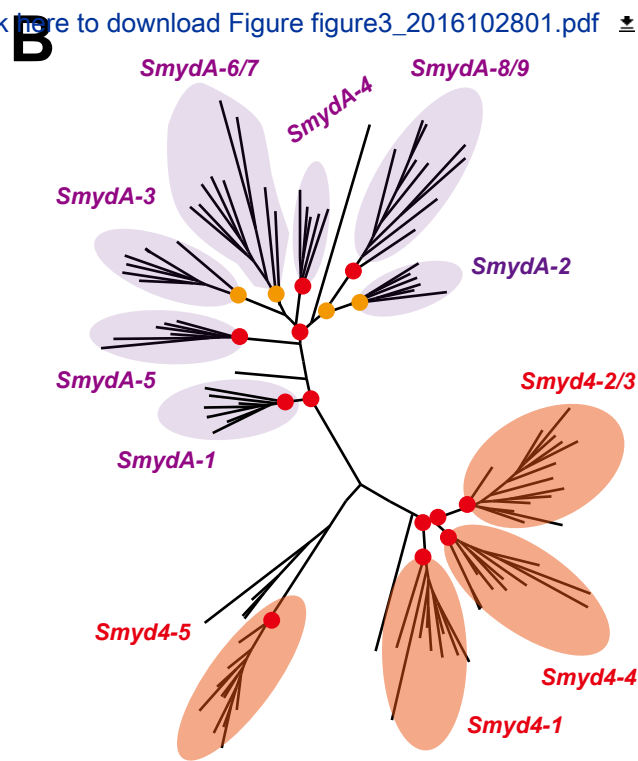
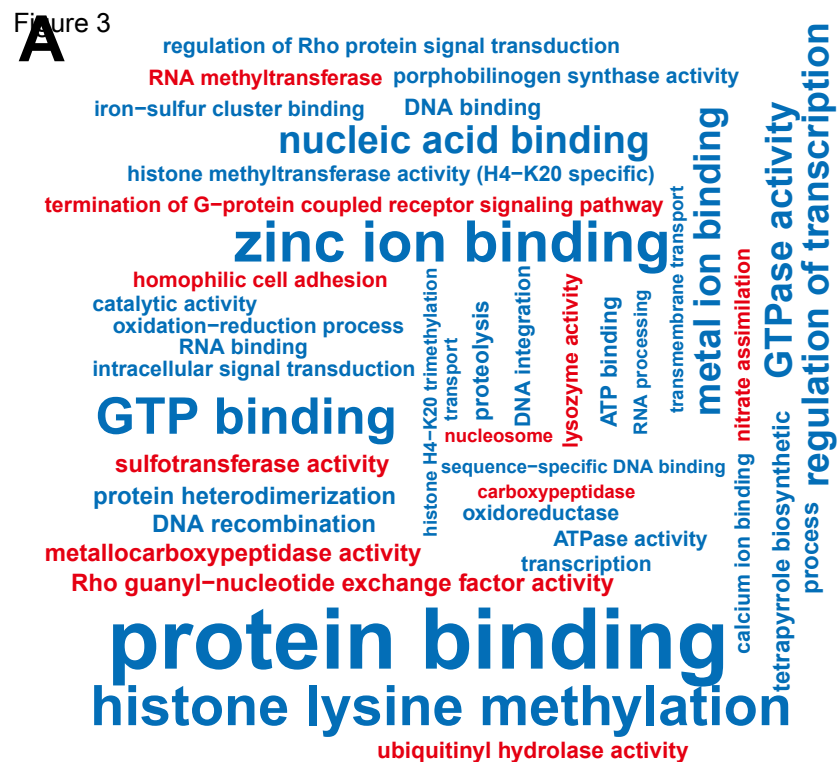


Figure 2





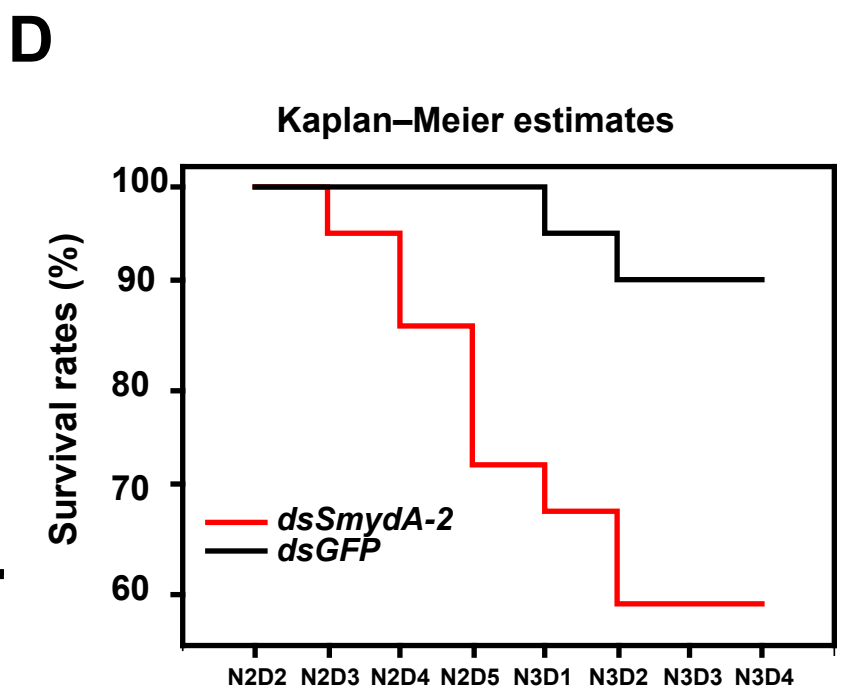
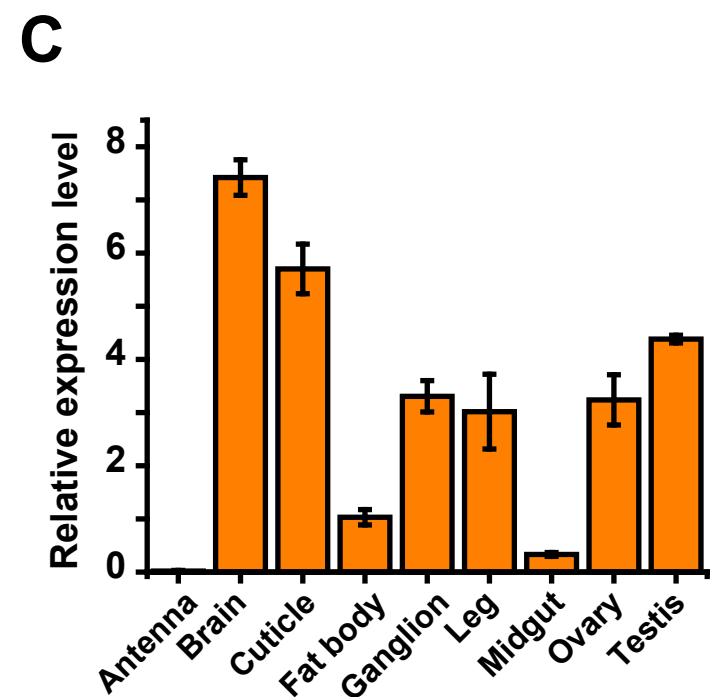
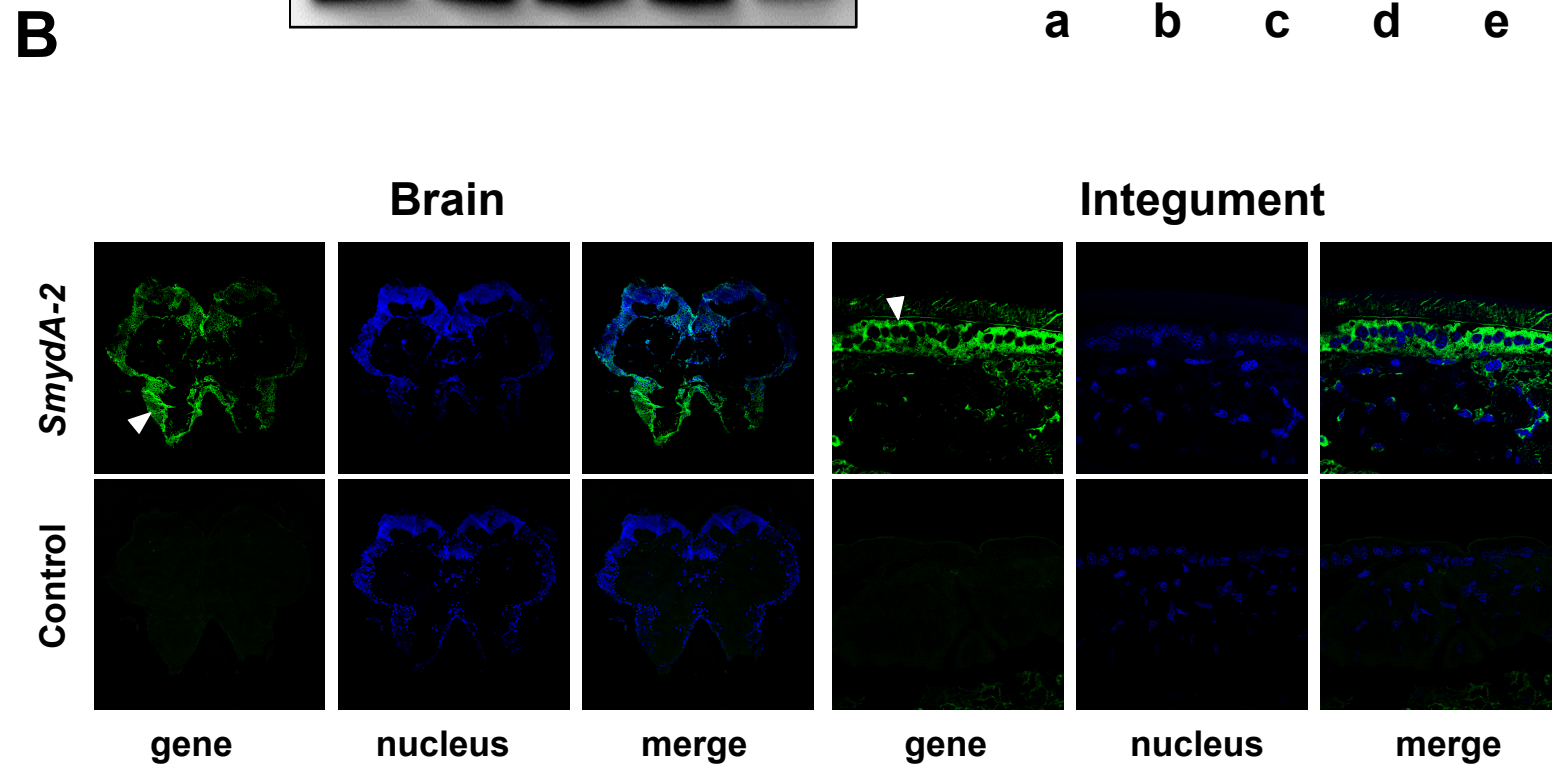
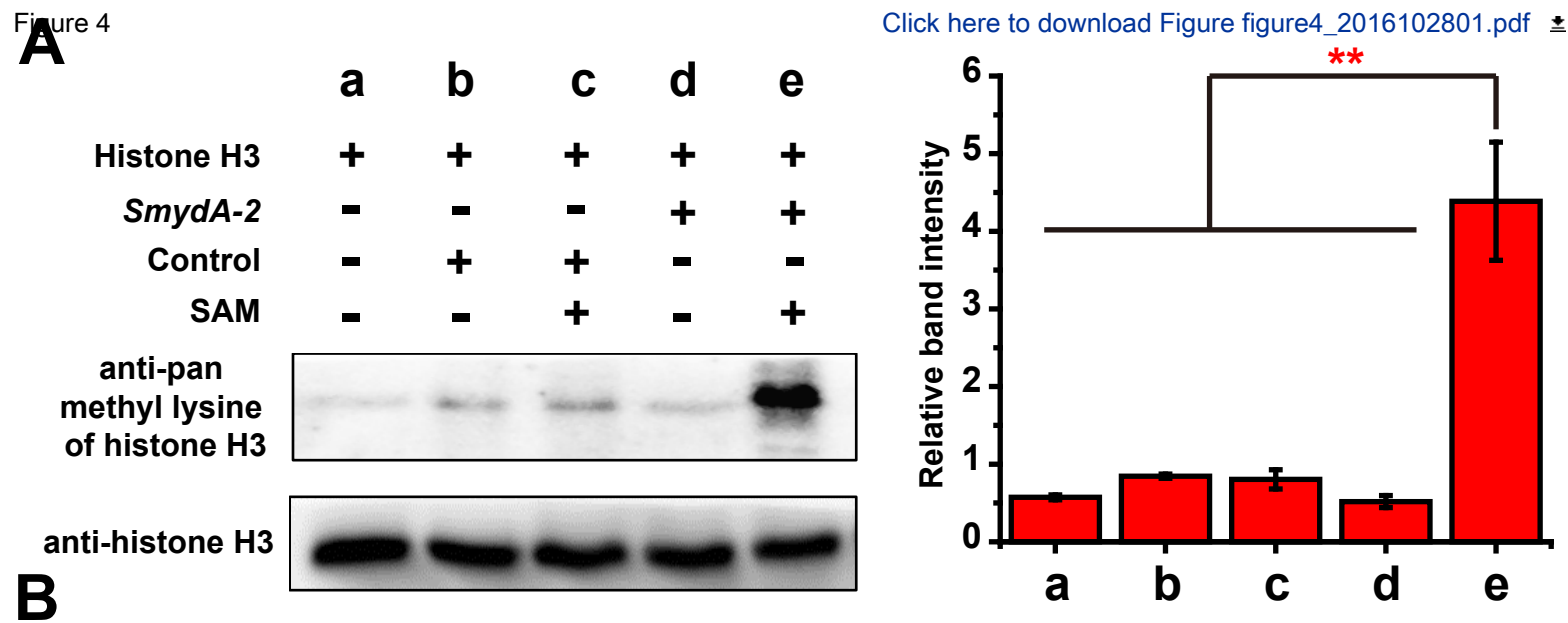
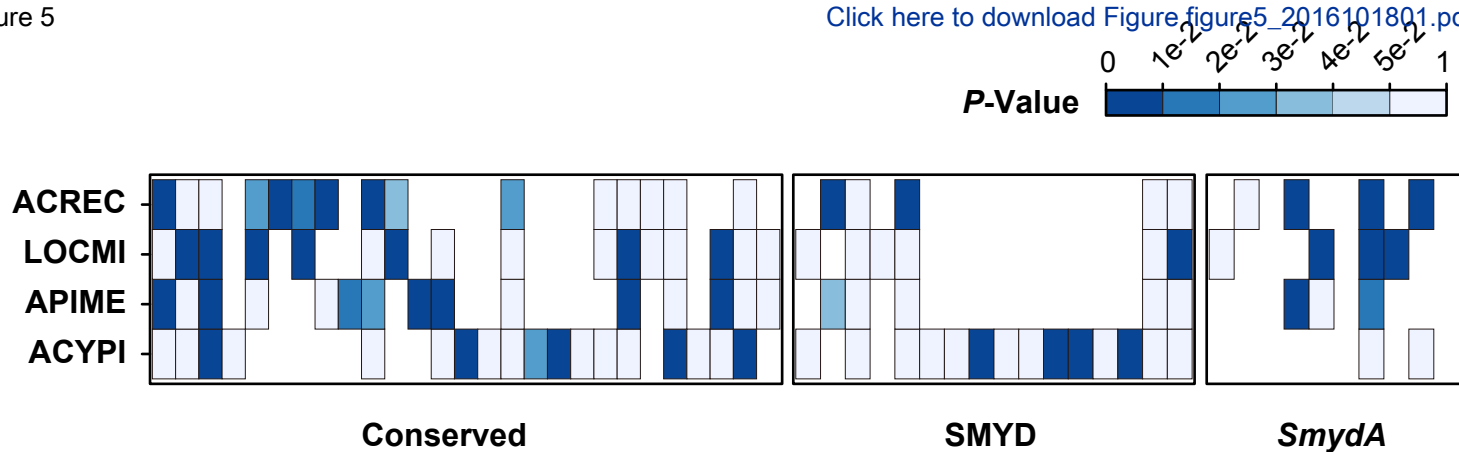


Figure 5

[Click here to download Figure figure5_2016101801.pdf](#)





[Click here to access/download](#)

Supplementary Material

[insectsSetDomain_supply_16102701_giga.pdf](#)

