

Hemimetabolous genomes reveal molecular basis of termite eusociality

Mark C Harrison,^{1*} Evelien Jongepier,^{1*} Hugh M. Robertson,^{2*} Nicolas Arning,¹
Tristan Bitard-Feildel,¹ Hsu Chao,³ Christopher P. Childers,⁴ Huyen Dinh,³
Harshavardhan Doddapaneni,³ Shannon Dugan,³ Johannes Gowin,^{5,6} Carolin
Greiner,^{5,6} Yi Han,³ Haofu Hu,⁷ Daniel S.T. Hughes,³ Ann-Kathrin Huylmans,⁸
Carsten Kemena,¹ Lukas P.M. Kremer,¹ Sandra L. Lee,³ Alberto Lopez-Ezquerro,¹
Ludovic Mallet,¹ Jose M. Monroy-Kuhn,⁵ Annabell Moser,⁵ Shwetha C. Murali,³
Donna M. Muzny,³ Saria Otani,⁷ Maria-Dolors Piulachs,⁹ Monica Poelchau,⁴
Jiaxin Qu,³ Florentine Schaub,⁵ Ayako Wada-Katsumata,¹⁰ Kim C. Worley,³
Qiaolin Xie,¹¹ Guillem Ylla,⁹ Michael Poulsen,⁷ Richard A. Gibbs,³ Coby Schal,¹⁰
Stephen Richards,³ Xavier Belles,^{9†} Judith Korb,^{5,6†} Erich Bornberg-Bauer^{1†}

¹Institute for Evolution and Biodiversity, University of Münster, Münster, Germany.

²Department of Entomology, University of Illinois at Urbana-Champaign, Urbana IL, USA.

³Human Genome Sequencing Center, Department of Human and Molecular Genetics,
Baylor College of Medicine, Houston, TX, USA.

⁴USDA-ARS, National Agricultural Library, Beltsville, MD, USA.

⁵Evolutionary Biology & Ecology, University of Freiburg, Freiburg, Germany.

⁶Behavioral Biology, University of Osnabrück, Osnabrück, Germany.

⁷Ecology and Evolution, University of Copenhagen, Universitetsparken 15, 2100 Copenhagen, Denmark.

⁸Institute of Science and Technology Austria, Klosterneuburg, Austria.

⁹Institut de Biologia Evolutiva, CSIC-University Pompeu Fabra, Barcelona, Spain.

¹⁰Department of Entomology and Plant Pathology, North Carolina State University, Raleigh, NC, USA.

¹¹China National GeneBank, Beijing Genomics Institute(BGI)-Shenzhen,Shenzhen,518083,China

†Corresponding authors. E-mail: xavier.belles@ibe.upf-csic.es (XB);
judith.korb@biologie.uni-freiburg.de (JK); ebb@uni-muenster.de (EBB)

*These authors contributed equally to this work.

Around 150 million years ago, eusocial termites evolved from within the cockroaches, 50 million years before eusocial Hymenoptera, such as bees and ants, appeared. Here, we report the first, 2GB genome of a cockroach, *Blattella germanica*, and the 1.3GB genome of the drywood termite, *Cryptotermes secundus*. We show evolutionary signatures of termite eusociality by comparing the genomes and transcriptomes of three termites and the cockroach against the background of 16 other eusocial and non-eusocial insects. Dramatic adaptive changes in genes underlying the production and perception of pheromones confirm the importance of chemical communication in the termites. These are accompanied by major changes in gene regulation and the molecular evolution of caste determination. Many of these results parallel molecular mechanisms of eusocial evolution in Hymenoptera. However, the specific solutions are remarkably different, thus revealing a striking case of convergence in one of the major evolutionary transitions in biological complexity.

1 Eusociality, the reproductive division of labour with overlapping generations and cooperative brood care,
2 is one of the major evolutionary transitions in biology¹. Although rare, eusociality has been observed
3 in a diverse range of organisms, including shrimps, mole-rats and several insect lineages^{2,3,4}. A particu-
4 larly striking case of convergent evolution occurred within the holometabolous Hymenoptera and in the
5 hemimetabolous termites (Isoptera), which are separated by ca. 400 my of evolution⁵. Termites evolved
6 within the cockroaches around 150 mya, towards the end of the Jurassic^{6,7}, about 50 my before the first
7 bees and ants appeared⁵. Therefore, identifying the molecular mechanisms common to both origins of
8 eusociality is crucial to understanding the fundamental signatures of these rare evolutionary transitions.
9 While the availability of genomes from many eusocial and non-eusocial hymenopteran species⁸ has allowed
10 extensive research into the origins of eusociality within ants and bees^{9,10,11}, a paucity of genomic data
11 from cockroaches and termites has precluded large-scale investigations into the evolution of eusociality
12 in this hemimetabolous clade.

13 The conditions under which eusociality arose from within the cockroaches differ greatly from
14 those present in the non-eusocial ancestors of eusocial Hymenoptera. Termites and cockroaches are
15 hemimetabolous and so show a direct development, while holometabolous hymenopterans complete the
16 adult body plan during metamorphosis. In termites, workers are immatures and only reproductive castes
17 are adults¹², while in Hymenoptera, adult workers and queens represent the primary division of labour.
18 Moreover, termites are diploid and their colonies consist of both male and female workers, and usually
19 a queen and king dominate reproduction. This is in contrast to the haplodiploid system found in Hy-
20 menoptera, in which all workers and dominant reproductives are female. It is therefore intriguing that
21 strong similarities have evolved convergently within the termites and the hymenopterans, such as differ-
22 entiated castes and a nest life with reproductive division of labour. The termites can be subdivided into
23 wood-dwelling and foraging termites. The former belong to the lower termites and produce simple, small

24 colonies with totipotent workers that can become reproductives. Foraging termites (some lower and all
25 higher termites) form large, complex societies, in which worker castes can be irreversible¹². Similarly,
26 within ants, bees and wasps, varying levels of eusociality exist.

27 Here we provide insights into the genomic signatures of eusociality within the termites. We analysed
28 the genomes of three termite species with differing levels of social complexity and compared them to the
29 first cockroach genome, as a closely-related non-eusocial outgroup. Furthermore, differences in expression
30 between nymphs and adults of the cockroach were compared to differences in expression between workers
31 and reproductives of the three termites, in order to gain insights into how expression patterns changed
32 along with the evolution of castes. Using fifteen additional insect genomes to infer background gene
33 family turnover rates, we analysed the evolution of gene families along the transition from non-social
34 cockroaches to eusociality in the termites. In this study we concentrated particularly on two hallmarks
35 of insect eusociality, as previously described for Hymenoptera, with the expectation that similar patterns
36 occurred along with the emergence of termites. These are the evolution of a sophisticated chemical
37 communication, which is essential for the functioning of a eusocial insect colony^{3,13,14} and major changes
38 in gene regulation along with the evolution of castes^{9,10}. Additionally, we tested the hypothesis that the
39 high levels of transposable elements present in cockroach and termite genomes allowed the evolution of
40 gene families which were essential to the transition to eusociality.

41 **Evolution of genomes and proteomes**

42 We sequenced and assembled the genome of the German cockroach, *Blattella germanica* (Ectobiidae),
43 and of the lower, drywood termite, *Cryptotermes secundus* (Kalotermitidae; for assembly statistics see
44 supplementary table S1). The cockroach genome (2.0 Gb) is considerably larger than all three termite
45 genomes. The genome size of *C. secundus* (1.30 Gb) is comparable to the higher, fungus-growing termite,
46 *Macrotermes natalensis*, (1.31 Gb, Termitidae)¹⁵ but more than twice as large as the lower, dampwood
47 termite, *Zootermopsis nevadensis* (562 Mb, Termopsidae)¹⁶. The smaller genomes of termites compared
48 to the cockroach are in line with previous size estimations based on C-values¹⁷. The proteome of *B. ger-*
49 *manica* (29,216 proteins) is also much larger than in the termites, where we find the proteome size in
50 *C. secundus* (18,162) to be similar to the other two termites (*M. natalensis*: 16,140; *Z. nevadensis*:
51 15,459; Fig. 1). In fact, the *B. germanica* proteome was the largest among all 21 arthropod species anal-
52 ysed here (20 insects and the centipede *Strigamia maritima*; Fig. 1). Strong evidential support for over
53 80% of these proteins in *B. germanica* (see supporting material) and large expansions in many manually
54 annotated gene families offer high confidence in the accuracy of this proteome size. We compared gene
55 expression between nymphs (5th and 6th instars) and female reproductive adults in *B. germanica*, and

56 between workers, queens and kings in each of the three termites. Gene expression differed significantly
57 ($p < 0.05$) between female reproductives and nymphs in 2457 genes for *B. germanica*. In the termites
58 3369 (*C. secundus*) to 6756 (*Z. nevadensis*) genes differed significantly between queens and workers,
59 which are arguably analogous to female adults and nymphs in the cockroach (Fig. 2).

60 The transitions to eusociality in ants¹⁰ and bees⁹ have been linked to major changes in gene family
61 sizes. Similarly, we detected significant gene family changes on the branch leading to the termites (7
62 expansions and 10 contractions; Fig. S1, table S2). The numbers of species-specific, significant expan-
63 sions and contractions of gene families varied within termites (*Z. nevadensis*: 15/5; *C. secundus*: 27/3;
64 *M. natalensis*: 24/20; tables S3-S5). Interestingly, in *B. germanica* we measured 93 significant gene fam-
65 ily expansions but no contractions (table S6), which contributed to the large proteome. The *C. secundus*
66 and *B. germanica* genomes contain similar proportions of repetitive content (both 55%; table S7), which
67 is higher than in both *Z. nevadensis* (28%) and the higher termite, *M. natalensis* (46%)¹⁸. This is in
68 contrast to the reported negative correlation between repetitive content and the level of eusociality in
69 bees⁹. As also found in *Z. nevadensis* and *M. natalensis*¹⁸, LINEs and especially the subfamily BovB
70 were the most abundant transposable elements (TEs) in the *B. germanica* and *C. secundus* genomes,
71 indicating that a proliferation of LINEs may have occurred in the ancestors of Blattodea (cockroaches
72 and termites). We hypothesised that these high levels of TEs may be driving the high turnover in gene
73 family sizes within the termites and *B. germanica*¹⁹. Expanded gene families indeed had more repetitive
74 content within 10 kb flanking regions in all three termites ($p < 1.3 \times 10^{-8}$; Wald *t*-test; table S8-S9), in
75 particular in the higher termite *M. natalensis*. In contrast, gene family expansions were not correlated
76 with TE content in flanking regions for *B. germanica*. These results suggest a major expansion of LINEs
77 at the root of the Blattodea clade contributed to the evolution of gene families within termites, likely via
78 unequal crossing-over¹⁹; however, the expansions in *B. germanica* were not facilitated by TEs. It can
79 therefore be concluded that the large expansion of LINEs within Blattodea allowed the evolution of gene
80 families which ultimately facilitated the transition to eusociality.

81 Out of 729 non-saturated (synonymous substitution rate: $dS < 3$) 1-to-1 protein orthologs between
82 the termites and the two closest related, available non-eusocial species, *B. germanica* and the orthopteran
83 *Locusta migratoria*, we found 165 (22.6%) to be evolving significantly faster (ratio of nonsynonymous to
84 synonymous nucleotide substitution rates: dN/dS or ω) among the termites. These genes were enriched in
85 functions related to carbohydrate metabolism (table S10), which was also over-represented in genes with
86 higher ω values in eusocial compared to non-eusocial bees¹¹. Functions related to oxidation-reduction
87 processes, including a number of mitochondrial genes, were also enriched among genes with a higher ω
88 within termites. This is consistent with the finding that mitochondrial genes were found to be evolving
89 under positive selection during the evolution of ants²⁰. One hundred (60.6%) of the genes with a signifi-

90 cantly higher ω within the termites were evolving even faster on the branch leading to the higher termite,
91 *M. natalensis*. The ten most significant of these genes have functions related to signaling, cell transport,
92 glycogen metabolism, transcription regulation, proteolysis and morphogenesis (table S11). These findings
93 support the notion that changes in gene regulation, diet and developmental pathways have facilitated
94 the transition to higher eusociality and a change from simple wood-dwelling colonies to large, complex,
95 foraging societies.

96 CHC production

97 Despite their different ancestry, both termites and eusocial hymenopterans are characterised by the pro-
98 duction of caste-specific cuticular hydrocarbons (CHCs)^{21,22,23}, which are often crucial for regulating
99 reproductive division of labour and chemical communication. Accordingly, we find changes in the ter-
100 mites in three groups of proteins involved in the synthesis of CHCs: desaturases (introduction of double
101 bonds²⁴), elongases (extension of C-chain length²⁵) and CYP4G1 (last step of CHC biosynthesis²⁶).

102 Desaturases are thought to be important for division of labour and social communication in ants²⁷. As
103 previously described for ants²⁷, Desat B genes are the most abundant desaturase family in the termites
104 and the cockroach (table S12), especially in *M. natalensis* where we found ten gene copies (significant
105 expansion; $p = 0.0024$; table S5; Fig. S7). As in ants, especially the First Desaturases (Desat A - Desat E)
106 vary greatly in their expression between castes and species in the three termites (Fig. 2; table S13)²⁷.
107 Both in *Z. nevadensis* and *M. natalensis*, most desaturases are more highly expressed in worker castes
108 than in queens, while these genes are generally more evenly expressed between castes in *C. secundus*.
109 In *B. germanica* 4 out of 7 Desat B genes are over-expressed in nymphs compared to female adults
110 and only one is more highly expressed in female adults (table S13). This pattern has been maintained
111 in *Z. nevadensis* (1 queen, 2 worker genes) and *M. natalensis* (5 worker genes), in which most Desat B
112 genes are worker-specific. In contrast to ants, where these genes are under strong purifying selection²⁷, we
113 found significant positive selection within the Desat B genes for the highly eusocial termite, *M. natalensis*,
114 (codeml site models 7 & 8; $p = 1.1 \times 10^{-16}$), indicating a diversification in function, possibly related to their
115 greater diversification of worker castes (major and minor workers, major and minor soldiers). Although
116 desaturases are often discussed in the context of CHC production and chemical communication, their
117 biochemical roles are quite diverse²⁷, and the positive selection we observe for *M. natalensis* may, at least
118 in part, be related to their rather different ecology of foraging and fungus farming rather than nest mate
119 recognition. Future experimental verification of the function of these genes will help better understand
120 these observed genomic and transcriptomic patterns.

121 Underlining an increased importance of CHC communication in termites, the expression patterns

122 of elongases (extension of C-chain length) differ considerably in the termites compared to the cockroach
123 (Fig. 2; table S14). In contrast to *B. germanica*, in which elongases are both nymph- (6 genes) and adult-
124 biased (5 genes), only one or two elongase genes in each termite are queen-biased in their expression,
125 while many are worker-biased. As with the desaturases, a group of *M. natalensis* elongases also reveal
126 significant signals of positive selection (codeml branch-site test; $p = 4 \times 10^{-4}$), further indicating a greater
127 diversification of CHC production in this higher termite.

128 The last step of CHC biosynthesis, the production of hydrocarbons from long-chain fatty aldehydes,
129 is catalyzed by a P450 gene, CYP4G1, in *Drosophila melanogaster*²⁶. We found one copy of CYP4G1 in
130 *B. germanica*, *Z. nevadensis* and *C. secundus*, but three copies in *M. natalensis*, reinforcing the greater
131 importance of CHC synthesis in this higher termite. Such P450 genes have experimentally been shown to
132 be crucial for maintaining reproductive division in the termite *C. secundus*²⁸. Corroborating the known
133 importance of maternal CHCs in *B. germanica*²⁹, CYP4G1 is over-expressed in female adults compared
134 to nymphs (Fig. 2; table S15). In each of the termites, however, CYP4G1 is more highly expressed in
135 workers (or kings in *C. secundus*) compared to queens (Fig. 2; table S15), adding support that, compared
136 to cockroach nymphs, a change in the dynamics and turnover of CHCs in termite workers has taken place.

137 Perception of chemical cues

138 Insects perceive chemical cues from toxins, pathogens, food and pheromones with three major families
139 of chemoreceptors, the Odorant (ORs), Gustatory (GRs) and Ionotropic (IRs) Receptors³⁰. Especially
140 ORs have been linked to colony communication in eusocial Hymenoptera, where they abound^{13,14}. In-
141 terestingly, as previously detected for *Z. nevadensis*¹⁶, the OR repertoire is substantially smaller in
142 *B. germanica* and all three termites compared to hymenopterans. IRs, on the other hand, which are
143 less frequent in hymenopterans, are strongly expanded in the cockroach and termite genomes (Fig. 3 &
144 Fig. S6). Intronless IRs, which are known to be particularly divergent³¹, show the greatest cockroach-
145 and Blattodea-specific expansions (Fig. 3a, Blattodea-, Cockroach- and Group D-IRs). By far the most
146 IRs among all investigated species were found in *B. germanica* (455 complete gene models), underlining
147 that the capacity for detecting many different kinds of chemosensory cues is crucial for this generalist
148 that thrives in challenging, human environments. In line with a specialisation in diet and habitat, the
149 total number of IRs is lower within the termites (*Z. nevadensis*: 141; *C. secundus*: 135; *M. natalensis*:
150 75). Nevertheless, IRs are more numerous in termites than in all other analysed species (except *Nasonia*
151 *vitripennis*: 111). This is strikingly similar to the pattern for ORs in Hymenoptera, which are also highly
152 numerous in non-eusocial outgroups as well as in eusocial species^{13,32}.

153 We scanned each IR group for signs of species-specific positive selection. Within the Blattodea-specific

154 intronless IRs, we found several codon positions under significant positive selection for the higher termite,
155 *M. natalensis* (codeml site models 7 & 8; $p < 1.7 \times 10^{-10}$). The positively evolving codons are situated
156 within the two ligand-binding lobes of the receptors (Fig. 3c), showing a diversification of ligand specificity
157 has occurred along with the transition to higher eusociality and a change from wood-feeding to fungus-
158 farming in this higher termite. In total, only two IRs were differentially expressed between nymphs and
159 adult females in *B. germanica*. Underlining a change in expression along with the evolution of castes, we
160 found 35 IRs to be differentially expressed between workers and queens in *Z. nevadensis*, 12 in *C. secundus*
161 and 11 in *M. natalensis* (Fig. 3, table S16). The possible role of IRs in pheromonal communication has
162 been highlighted both in the cockroach *Periplaneta americana*³³ and in *D. melanogaster*³⁴, where several
163 IRs show sex-biased expression.

164 One group of ORs (orange clade in Fig. 3b) is evolving under significant positive selection at codon
165 positions within the second transmembrane domain in *M. natalensis* (codeml site model; $p = 1.1 \times 10^{-11}$)
166 and *C. secundus* ($p = 5.6 \times 10^{-16}$; Fig. 3d). Such a variation in the transmembrane domain can be related
167 to ligand binding specificity, as has been shown for a polymorphism in the third transmembrane domain
168 for an OR in *D. melanogaster*^{35,36}, adding further support for an adaptive evolution of chemoreceptors,
169 in line with the greater need for a sophisticated colony communication in the termites. Similar to IRs, a
170 higher proportion of ORs were differentially expressed between workers and queens in the three termites
171 than between nymphs and adults in the cockroach (Fig. 2; table S17), highlighting a change in expression
172 and function along with the transition to eusociality. The evolution of chemoreceptors along with the
173 emergence of the termites can also be related to adaptation processes and changes in diet compared to
174 the cockroach. Experimental verification will help pinpoint which receptors are particularly important
175 for communication.

176 Changes in gene regulation in termites

177 The development of distinct castes underlying division of labour is achieved via differential gene expres-
178 sion. Major changes in gene regulation have been reported as being central to the transition to eusociality
179 in bees⁹ and ants¹⁰. Accordingly, we found major changes in DNA methylation patterns (levels per 1-to-1
180 ortholog) among the termites compared to four other hemimetabolous insect species (Fig. 4a). This is
181 revealed by CpG depletion patterns ($\text{CpG}_{o/e}$), a reliable predictor of DNA methylation^{37,38}, correlating
182 more strongly between the termites than among any of the other analysed hemimetabolous insects (Fig.
183 4). In other words, within orthologous genes, DNA methylation levels differ greatly between termites and
184 other hemimetabolous species but remain conserved among termite species. Furthermore, a higher pro-
185 portion of genes were putatively DNA methylated ($\text{CpG}_{o/e} < 0.5$) within the termites (40.7% to 50.6%)

186 compared to other hemimetabolous species (11.5% to 34.0%), as also described for eusocial compared to
187 solitary bees⁹.

188 Levels of DNA methylation correlated negatively with caste-specificity of expression for each of the
189 termites. This is confirmed by a positive correlation between CpG_{o/e} (negative association with level of
190 DNA methylation) and log₂-fold change of expression between queens and workers (Pearson's $r = 0.32$
191 to 0.36 ; $p < 2.2 \times 10^{-16}$). The caste-specific expression of unmethylated genes in termites is reflected
192 in the enrichment of GO terms related to sensory perception, regulation of transcription, signalling and
193 development, whereas methylated genes are mainly related to general metabolic processes (Fig. 4b, tables
194 S18). These results show strong parallels to findings for eusocial Hymenoptera^{39,40,41,42}. This is in stark
195 contrast to the non-eusocial cockroach, *B. germanica*, where there was only a very weak relationship
196 between CpG_{o/e} and differential expression between nymphs and adult females ($r = 0.14$), nor were any
197 large differences apparent in enriched GO terms between methylated and non-methylated genes (Fig.
198 4b).

199 Our results argue in favour of a diminished role of DNA methylation in caste-specific expression
200 within eusocial insects, as recently shown^{37,43}. In fact, DNA methylation appears to be important for
201 the regulation of house-keeping genes because methylated genes are related to general biological processes
202 (further supported by lower CpG_{o/e} within 1-to-1 orthologs than in non-conserved genes)⁴⁴, while caste-
203 specific genes are 'released' from this type of gene regulation. However, a recent study linked caste-specific
204 DNA methylation to alternative splicing in *Z. nevadensis*⁴⁵.

205 Major biological transitions are often accompanied by expansions of transcription factor (TF) families,
206 such as genes containing zinc-finger (ZF) domains⁴⁶. We also observed large differences in ZF families
207 within the termites compared to *B. germanica*. Many ZF families were reduced or absent in termites,
208 while different, unrelated ZF gene families were significantly expanded (tables S2-S5). Queen-biased
209 genes were significantly over-represented among ZF genes for termites ($p < 2 \times 10^{-10}$; χ^2 test; table
210 S19), indicating an important role of ZF genes in the regulation of genes related to caste-specific tasks
211 and colony organisation in the termites. This is in contrast to the significant under-representation of
212 differentially expressed ZF genes within *B. germanica* ($p = 1.42 \times 10^{-5}$; χ^2 -test). Interestingly, two other
213 important TF families (bHLH and bZIP)⁴⁶, which were not expanded in the termites, showed no caste-
214 specific expression pattern ($p > 0.05$). These major upheavals in ZF gene families and their caste-specific
215 expression show that major changes in TFs accompanied the evolution of termites, strikingly similar to
216 the evolution of ants¹⁰.

217 Endocrine regulation

218 Hemimetabolous eusociality is characterised by differentiated castes, which represent different develop-
219 mental stages. This is in contrast to eusocial Hymenoptera, in which workers and reproductives are adults.
220 While cockroaches develop directly through several nymphal stages before becoming reproductive adults,
221 termite development is more phenotypically plastic, and workers are essentially immatures (Fig. 2).
222 In wood-dwelling termites, such as *C. secundus* and *Z. nevadensis*, worker castes are non-reproductive
223 immatures that are totipotent to develop into other castes, while in the higher termite, *M. natalensis*,
224 workers can be irreversibly defined instars. It is therefore clear that a major change during the evolution
225 of termites occurred within developmental pathways. Accordingly, we found changes in expression and
226 gene family size of several genes related both to molting and metamorphosis.

227 In the synthesis of the molting hormone, 20-hydroxyecdysone, the six Halloween genes (5 Cytochrome
228 P450s and a Rieske-domain oxygenase) play a key role^{47,48}. Only one Halloween gene, Shade (Shd;
229 CYP314A1), which mediates the final step of 20-hydroxyecdysone synthesis, is differentially expressed
230 between the final nymphal stages and adults females in *B. germanica* (Fig. 2; table S20), consistent with
231 its role in the nymphal or imaginal molt. In the three termites, the Halloween genes show varying caste-
232 specific expression (Fig. 2; table S20), showing that ecdysone plays a significant role in the regulation
233 of caste differences. Ecdysteroid kinase genes (EcK), which convert the insect molting hormone into its
234 inactive state, ecdysone 22-phosphate, for storage⁴⁹, are only over-expressed in female adults compared
235 to nymphs in *B. germanica* (16/51 genes, Fig.2, table S21). In termites, however, where the gene copy
236 number is reduced (18 to 20 per species), these important molting genes appear to have evolved worker-
237 specific functions (Fig. 2; table S21).

238 Whereas 20-hydroxyecdysone promotes molting, juvenile hormone (JH) represses imaginal develop-
239 ment in pre-adult instars⁵⁰. JH is important in caste differentiation in eusocial insects, including ter-
240 mites^{12,51}. Hemolymph juvenile hormone binding proteins (JHBP), which transport JH to its target
241 tissues⁵², are reduced within the termites (21 to 33 genes) but significantly expanded in *B. germanica*
242 (51 copies). Thirteen of the JHBP genes are over-expressed in adult females and only 8 in nymphs in
243 *B. germanica*. In both *Z. nevadensis* (15 worker-specific and 1 queen-specific) and *M. natalensis* (11
244 worker-specific and 4 queen-specific), on the other hand, JHBPs are significantly more worker-biased (p
245 < 0.01 , χ^2 test; table S22; Fig. 2). In *C. secundus*, expression is more varied, with 5 worker-biased, 8
246 king-biased and 3 queen-biased genes (Fig. 2; table S22).

247 These changes in copy number and caste-specific expression of genes involved in metamorphosis and
248 molting within termites compared to the German cockroach demonstrate that changes occurred in the
249 control of the developmental pathway along with the evolution of castes. However, this interpretation

250 needs to be experimentally verified.

251 **Conclusions**

252 These results, considered alongside many studies on eusociality in Hymenoptera^{9,10,13,27}, provide evi-
253 dence that major changes in gene regulation and the evolution of sophisticated chemical communication
254 are fundamental to the transition to eusociality in insects. Strong changes in DNA methylation patterns
255 correlated with broad-scale modifications of expression patterns. Many of these modified expression
256 patterns remained consistent among the three studied termite species and occurred within protein path-
257 ways essential for eusocial life, such as CHC production, chemoperception, ecdysteroid synthesis and JH
258 transport. Many of the mechanisms implicated in the evolution of eusociality in the termites occurred
259 convergently around 50 my later in the phylogenetically distant Hymenoptera. However, several details
260 are unique due to the distinct conditions within which eusociality arose. One important difference is the
261 higher TE content within cockroaches and termites, which likely facilitated changes in gene family sizes,
262 supporting the transition to eusociality. However, the most striking difference is the apparent importance
263 of IRs for chemical communication in the termites, compared to ORs in Hymenoptera. According to our
264 results, the non-eusocial ancestors of termites possessed a broad repertoire of IRs, which favoured the
265 evolution of important functions for colony communication in these chemoreceptors within the termites,
266 whereas in the solitary ancestors of eusocial hymenopterans ORs were most abundant^{13,32}. The parallel
267 expansions of different chemoreceptor families in these two independent origins of eusociality indicate that
268 convergent selection pressures existed during the evolution of colony communication in both lineages.

References

- 269 1. Szathmáry, E. & Maynard Smith, J. The major evolutionary transitions. *Nature* **374**, 227–232
270 (1995).
- 271
- 272 2. Andersson, M. The Evolution of Eusociality. *Annual Review of Ecology and Systematics* **15**, 165–189
273 (1984).
- 274 3. Wilson, E. O. *The insect societies* (Harvard University Press, Cambridge, MA, 1971).
- 275 4. Rubenstein, D. R. & Abbot, P. The evolution of social evolution. In *Comparative Social Evolution*
276 (Cambridge University Press, Cambridge, 2017).
- 277 5. Misof, B. *et al.* Phylogenomics resolves the timing and pattern of insect evolution. *Science* **346**,
278 763–767 (2014).
- 279 6. Legendre, F. *et al.* Phylogeny of Dictyoptera: Dating the Origin of Cockroaches, Praying Mantises
280 and Termites with Molecular Data and Controlled Fossil Evidence. *PLOS ONE* **10**, e0130127 (2015).
- 281 7. Bourguignon, T. *et al.* The Evolutionary History of Termites as Inferred from 66 Mitochondrial
282 Genomes. *Molecular Biology and Evolution* **32**, 406–421 (2015).
- 283 8. Elsner, D., Kremer, L. P., Arning, N. & Bornberg-Bauer, E. Comparative genomic approaches to
284 investigate molecular traits specific to social insects. *Current Opinion in Insect Science* **16**, 87–94
285 (2016).
- 286 9. Kapheim, K. M. *et al.* Genomic signatures of evolutionary transitions from solitary to group living.
287 *Science* **348**, 1139–1143 (2015).
- 288 10. Simola, D. F. *et al.* Social insect genomes exhibit dramatic evolution in gene composition and
289 regulation while preserving regulatory features linked to sociality. *Genome Research* **23**, 1235–1247
290 (2013).
- 291 11. Woodard, S. H. *et al.* Genes involved in convergent evolution of eusociality in bees. *Proceedings of*
292 *the National Academy of Sciences* **108**, 7472–7477 (2011).
- 293 12. Korb, J. & Hartfelder, K. Life history and development - a framework for understanding develop-
294 mental plasticity in lower termites. *Biological Reviews* **83**, 295–313 (2008).
- 295 13. Zhou, X. *et al.* Chemoreceptor Evolution in Hymenoptera and Its Implications for the Evolution of
296 Eusociality. *Genome Biology and Evolution* **7**, 2407–2416 (2015).

- 297 14. Tribble, W. *et al.* Orco mutagenesis causes loss of antennal lobe glomeruli and impaired social behavior
298 in ants. *bioRxiv* 112532 (2017).
- 299 15. Poulsen, M. *et al.* Complementary symbiont contributions to plant decomposition in a fungus-farming
300 termite. *Proceedings of the National Academy of Sciences* **111**, 14500–14505 (2014).
- 301 16. Terrapon, N. *et al.* Molecular traces of alternative social organization in a termite genome. *Nature*
302 *Communications* **5**, 3636 (2014).
- 303 17. Gregory, T. R. Animal Genome Size Database. <http://www.genomesize.com/> (2017).
- 304 18. Korb, J. *et al.* A genomic comparison of two termites with different social complexity. *Frontiers in*
305 *Genetics* **6** (2015).
- 306 19. Kazazian, H. H. Mobile Elements: Drivers of Genome Evolution. *Science* **303**, 1626–1632 (2004).
- 307 20. Roux, J. *et al.* Patterns of Positive Selection in Seven Ant Genomes. *Molecular Biology and Evolution*
308 **31**, 1661–1685 (2014).
- 309 21. Oystaeyen, A. V. *et al.* Conserved Class of Queen Pheromones Stops Social Insect Workers from
310 Reproducing. *Science* **343**, 287–290 (2014).
- 311 22. Weil, T., Hoffmann, K., Kroiss, J., Strohm, E. & Korb, J. Scent of a queen—cuticular hydrocarbons
312 specific for female reproductives in lower termites. *Naturwissenschaften* **96**, 315–319 (2009).
- 313 23. Dietemann, V., Peeters, C., Liebig, J., Thivet, V. & Hölldobler, B. Cuticular hydrocarbons mediate
314 discrimination of reproductives and nonreproductives in the ant *Myrmecia gulosa*. *Proceedings of the*
315 *National Academy of Sciences* **100**, 10341–10346 (2003).
- 316 24. Dallerac, R. *et al.* A $\delta 9$ desaturase gene with a different substrate specificity is responsible for the
317 cuticular diene hydrocarbon polymorphism in *Drosophila melanogaster*. *Proceedings of the National*
318 *Academy of Sciences* **97**, 9449–9454 (2000).
- 319 25. Finck, J., Berdan, E. L., Mayer, F., Ronacher, B. & Geiselhardt, S. Divergence of cuticular hydro-
320 carbons in two sympatric grasshopper species and the evolution of fatty acid synthases and elongases
321 across insects. *Scientific Reports* **6**, srep33695 (2016).
- 322 26. Qiu, Y. *et al.* An insect-specific P450 oxidative decarbonylase for cuticular hydrocarbon biosynthesis.
323 *Proceedings of the National Academy of Sciences* **109**, 14858–14863 (2012).
- 324 27. Helmkampf, M., Cash, E. & Gadau, J. Evolution of the insect desaturase gene family with an
325 emphasis on social Hymenoptera. *Molecular Biology and Evolution* 456–471 (2015).

- 326 28. Hoffmann, K., Gowin, J., Hartfelder, K. & Korb, J. The scent of royalty: a p450 gene signals
327 reproductive status in a social insect. *Molecular Biology and Evolution* **31**, 2689–2696 (2014).
- 328 29. Fan, Y., Eliyahu, D. & Schal, C. Cuticular hydrocarbons as maternal provisions in embryos and
329 nymphs of the cockroach *Blattella germanica*. *Journal of Experimental Biology* **211**, 548–554 (2008).
- 330 30. Joseph, R. M. & Carlson, J. R. Drosophila Chemoreceptors: A Molecular Interface Between the
331 Chemical World and the Brain. *Trends in Genetics* **31**, 683–695 (2015).
- 332 31. Croset, V. *et al.* Ancient Protostome Origin of Chemosensory Ionotropic Glutamate Receptors and
333 the Evolution of Insect Taste and Olfaction. *PLOS Genetics* **6**, e1001064 (2010).
- 334 32. Robertson, H. M., Gadau, J. & Wanner, K. W. The insect chemoreceptor superfamily of the parasitoid
335 jewel wasp *Nasonia vitripennis*. *Insect Molecular Biology* **19**, 121–136 (2010).
- 336 33. Chen, Y., He, M., Li, Z.-Q., Zhang, Y.-N. & He, P. Identification and tissue expression profile of
337 genes from three chemoreceptor families in an urban pest, *Periplaneta americana*. *Scientific Reports*
338 **6** (2016).
- 339 34. Koh, T.-W. *et al.* The Drosophila IR20a Clade of Ionotropic Receptors Are Candidate Taste and
340 Pheromone Receptors. *Neuron* **83**, 850–865 (2014).
- 341 35. Pellegrino, M., Steinbach, N., Stensmyr, M. C., Hansson, B. S. & Vosshall, L. B. A natural poly-
342 morphism alters odour and DEET sensitivity in an insect odorant receptor. *Nature* **478**, 511–514
343 (2011).
- 344 36. Nichols, A. S. & Luetje, C. W. Transmembrane Segment 3 of Drosophila melanogaster Odorant
345 Receptor Subunit 85b Contributes to Ligand-Receptor Interactions. *Journal of Biological Chemistry*
346 **285**, 11854–11862 (2010).
- 347 37. Bewick, A. J., Vogel, K. J., Moore, A. J. & Schmitz, R. J. Evolution of DNA methylation across
348 insects. *Molecular Biology and Evolution* 654–655 (2017).
- 349 38. Park, J. *et al.* Comparative Analyses of DNA Methylation and Sequence Evolution Using *Nasonia*
350 Genomes. *Molecular Biology and Evolution* **28**, 3345–3354 (2011).
- 351 39. Elango, N., Hunt, B. G., Goodisman, M. A. D. & Yi, S. V. DNA methylation is widespread and
352 associated with differential gene expression in castes of the honeybee, *Apis mellifera*. *Proceedings of*
353 *the National Academy of Sciences* **106**, 11206–11211 (2009).

- 354 40. Standage, D. S. *et al.* Genome, transcriptome and methylome sequencing of a primitively eusocial
355 wasp reveal a greatly reduced DNA methylation system in a social insect. *Molecular Ecology* **25**,
356 1769–1784 (2016).
- 357 41. Patalano, S. *et al.* Molecular signatures of plastic phenotypes in two eusocial insect species with
358 simple societies. *Proceedings of the National Academy of Sciences* **112**, 13970–13975 (2015).
- 359 42. Rehan, S. M., Glastad, K. M., Lawson, S. P. & Hunt, B. G. The Genome and Methylome of a
360 Subsocial Small Carpenter Bee, *Ceratina calcarata*. *Genome Biology and Evolution* **8**, 1401–1410
361 (2016).
- 362 43. Libbrecht, R., Oxley, P. R., Keller, L. & Kronauer, D. J. C. Robust DNA Methylation in the Clonal
363 Raider Ant Brain. *Current Biology* **26**, 391–395 (2016).
- 364 44. Foret, S., Kucharski, R., Pittelkow, Y., Lockett, G. A. & Maleszka, R. Epigenetic regulation of
365 the honey bee transcriptome: unravelling the nature of methylated genes. *BMC Genomics* **10**, 472
366 (2009).
- 367 45. Glastad, K. M., Gokhale, K., Liebig, J. & Goodisman, M. A. D. The caste- and sex-specific DNA
368 methylome of the termite *Zootermopsis nevadensis*. *Scientific Reports* **6**, 37110 (2016).
- 369 46. Schmitz, J. F., Zimmer, F. & Bornberg-Bauer, E. Mechanisms of transcription factor evolution in
370 Metazoa. *Nucleic Acids Research* **44**, 6287–6297 (2016).
- 371 47. Rewitz, K. F., Rybczynski, R., Warren, J. T. & Gilbert, L. I. The Halloween genes code for cy-
372 tochrome P450 enzymes mediating synthesis of the insect moulting hormone. *Biochemical Society*
373 *Transactions* **34**, 1256–1260 (2006).
- 374 48. Lang, M. *et al.* Mutations in the neverland Gene Turned *Drosophila pachea* into an Obligate Specialist
375 Species. *Science* **337**, 1658–1661 (2012).
- 376 49. Sonobe, H. *et al.* Purification, Kinetic Characterization, and Molecular Cloning of a Novel Enzyme,
377 Ecdysteroid 22-Kinase. *Journal of Biological Chemistry* **281**, 29513–29524 (2006).
- 378 50. Jindra, M., Belles, X. & Shinoda, T. Molecular basis of juvenile hormone signaling. *Current Opinion*
379 *in Insect Science* **11**, 39–46 (2015).
- 380 51. Korb, J. Juvenile Hormone: A Central Regulator of Termite Caste Polyphenism. In Kent, A. Z. a.
381 C. F. (ed.) *Advances in Insect Physiology*, vol. 48 of *Genomics, Physiology and Behaviour of Social*
382 *Insects*, 131–161 (Academic Press, 2015). DOI: 10.1016/bs.aiip.2014.12.004.

- 383 52. Kolodziejczyk, R. *et al.* Insect Juvenile Hormone Binding Protein Shows Ancestral Fold Present in
384 Human Lipid-Binding Proteins. *Journal of Molecular Biology* **377**, 870–881 (2008).

385 **Acknowledgements:** We thank three anonymous referees for their helpful, constructive feedback. We
386 also thank Oliver Niehuis for allowing use of the unpublished *E. danica* genome, Jürgen Gadau and Chris
387 Smith for comments and advice on the manuscript, Jonathan Schmitz for assistance with analyses and
388 proof-reading the manuscript. JK thanks Charles Darwin University (Australia), especially Prof. Stephen
389 Garnett and the Horticulture and Aquaculture team for providing logistic support to collect *C. secundus*.
390 The Parks and Wildlife Commission, Northern Territory, the Department of the Environment, Water,
391 Heritage and the Arts gave permission to collect (Permit number 36401) and export (Permit WT2010-
392 6997) the termites. USDA is an equal opportunity provider and employer. MCH and EJ supported by
393 DFG grant BO2544/11-1 to EBB. JK by University of Osnabrück and DFG grant KO1895/16-1. XB and
394 MDP supported by Spanish Ministerio de Economía y Competitividad (CGL2012-36251 and CGL2015-
395 64727-P to XB, and CGL2016-76011-R to MDP), including FEDER funds, and by Catalan Government
396 (2014 SGR 619). CS: grants from US Department of Housing and Urban Development (NCHHU-0017-
397 13), National Science Foundation (IOS-1557864), Alfred P. Sloan Foundation (2013-5-35 MBE), National
398 Institute of Environmental Health Sciences (P30ES025128) to Center for Human Health and the Envi-
399 ronment, and Blanton J. Whitmire Endowment. MP is supported by a Villum Kann Rasmussen Young
400 Investigator Fellowship (VKR10101).

401

402 **Author contributions:** E.B-B. conceived, managed and coordinated the project; M.C.H., E.J. and
403 H.M.R. are joint first authors. J.K. conceived and managed *C. secundus* sequencing project, coordi-
404 nated termite-related analyses; S.R. conceived and managed *B. germanica* sequencing project; S.R.,
405 S.D., S.L.L., H.C., H.V.D., H.D., Y.H., J.Q., S.C.M., D.S.T.H., K.C.W., D.M.M. and R.A.G. carried out
406 *B. germanica* library construction, genome sequencing and assembly; C.S., A.W.K. provided biological
407 material through full-sib mating for *B. germanica*; X.B. and C.S. co-managed the *B. germanica* analy-
408 sis; M.P. and C.P.C. implemented Web Apollo data traces; S.O. and M.P. provided biological material
409 for *M. natalensis*; C.G., J.G., J.M.M.-K., A.M., F.S., H.H. & J.K. coordinated and carried out DNA
410 and RNA sequencing for *C. secundus*; M-D.P., X.B. and G.Y. coordinated transcriptome sequencing
411 of *B. germanica*; L.M. performed automated gene prediction on *C. secundus*; E.J. and N.A. improved
412 assembly and annotation for *B. germanica* & *C. secundus*, compared and analysed genome sizes and qual-
413 ity. E.J., N.A. and L.P.M.K. analysed TEs; M.C.H. analysed CpG patterns and signatures of selection;
414 T.B-F., E.J., C.K., L.P.M.K. and A.L-E. performed orthology and phylogenetic analyses; L.P.M.K., E.J.,
415 H.M.R. and M.C.H. analysed gene family evolution; A.L-E., E.J. and M.C.H. analysed transcriptomes
416 and performed DE analyses; T.B-F. and A.L-E. carried out orthoMCL clustering; H.M.R. corrected
417 gene models for chemoreceptors; C.K. and E.J. for desaturases and elongases; A-K.H. and M.C.H. of
418 Cytochrome p450s; E.B-B and M.C.H drafted and wrote the manuscript; X.B., M-D.P., J.K. contributed
419 to sections of the manuscript; E.J., L.P.M.K., A.L-E., C.K., M.C.H. wrote and organized Supplementary
420 Materials; L.P.M.K., N.A., A.L-E., M.C.H. and E.B-B. prepared figures for the manuscript. All authors
421 read, corrected and commented on the manuscript.

422 MATERIALS AND METHODS

423 Genome sequencing and assembly

424 Genomic DNA from a single *Blattella germanica* male from an inbred line (strain: American Cyanamid
425 = Orlando Normal) was used to construct two paired-end (180 bp and 500 bp inserts) and one of the two
426 mate pair libraries (2 kb inserts). An 8kb mate pair library was constructed from a single female. The
427 libraries were sequenced on an Illumina HiSeq2000 sequencing platform. The 413 Gb of raw sequence
428 data were assembled with Allpaths LG¹, then scaffolded and gap-filled using the in-house tools Atlas-Link
429 v.1.0 (<https://www.hgsc.bcm.edu/software/atlas-link>) and Atlas gap-fill v.2.2. For *Cryptotermes*
430 *secundus*, three paired-end libraries (250 bp, 500 bp and 800 bp inserts) and three mate pair libraries
431 (2 kb, 5 kb and 10 kb inserts) were constructed from genomic DNA that was extracted from the head and
432 thorax of 1 000 individuals, originating from a single, inbred field colony. The libraries were sequenced on
433 an Illumina HiSeq2000 sequencing platform. The *C. secundus* genome was assembled using SOAPdenovo
434 (v.2.04)² with optimised parameters, followed by gapcloser (v1.10, released with SOAPdenovo) and kgf
435 (v1.18, released with SOAPdenovo).

436 Transcriptome sequencing and assembly

437 For annotation purposes, twenty-two whole body RNAseq samples from various developmental stages
438 were obtained for *B. germanica*. For *C. secundus* RNAseq libraries were obtained for three workers,
439 four queens and four kings, based on degutted, whole body extracts. In addition, we sequenced 10
440 *M. natalensis* RNAseq libraries from three queens, one king and six pools of workers. All libraries were
441 constructed using the Illumina (TruSeq) RNA-Seq kit.

442 For protein coding gene annotation, *B. germanica* reads were assembled with *de novo* Trinity (version
443 r2014-04-13)³. The *C. secundus* reads were assembled using i) Cufflinks on reads mapped with TopHat
444 (version2.2.1)^{4,5}, ii) *de novo* Trinity³; and iii) genome-guided Trinity on reads mapped with TopHat.

445 Repeat annotation

446 A custom *C. secundus* and *B. germanica* repeat library was constructed using a combination of homology-
447 based and *de novo* approaches, including RepeatModeler/RepeatClassifier (<http://www.repeatmasker.org/RepeatModeler.html>), LTRharvest/LTRdigest⁶ and TransposonPSI (<http://transposonpsi.sourceforge.net/>). The *ab initio* repeat library was complemented with the RepBase (update 29-

450 08-2016)⁷ and SINE repeat databases, filtered for redundancy with CD-hit and classified with Repeat-
451 Classifier. RepeatMasker (version open-4.0.6, <http://www.repeatmasker.org>) was used to mask the
452 *C. secundus* and *B. germanica* genome. Repeat content for the other studied species (Fig. 1) was
453 obtained from the literature^{8,9,10,11,12,13,14}.

454 Protein-coding gene annotation

455 The *B. germanica* genome was annotated with Maker (version 2.31.8)¹⁵, using (i) the species-specific
456 repeat library, (ii) *B. germanica* transcriptome data (22 whole body RNAseq samples), and (iii) the swis-
457 sprot/uniprot database (last accessed: 21-01-2016) plus the *C. secundus* and *Zootermopsis nevadensis*
458 protein sequences for evidence-based gene model predictions. AUGUSTUS (version 3.2)¹⁶, GeneMark-
459 ES Suite (version 4.21)¹⁷ and SNAP¹⁸ were used for *ab initio* predictions. *Cryptotermes secundus*
460 protein-coding genes were predicted using homology-based, *ab initio* and expression-based methods,
461 and integrated into a final gene set (see supplementary information). Gene structures were predicted
462 by GeneWise¹⁹. The *ab initio* annotations were predicted with AUGUSTUS²⁰ and SNAP¹⁸, retained
463 if supported by both methods and integrated with the homology-based predictions using GLEAN²¹.
464 Transcriptome-based gene models were merged with PASA²² and tested for coding potential with CPC²³
465 and OrfPredictor²⁴. PASA gene models were merged with the homology-based and *ab initio* gene set, re-
466 taining the PASA models in case of overlap. Desaturases, elongases, chemosensory receptors, Cytochrome
467 P450's and genes involved in the juvenile hormone pathway were manually curated in Blattodea.

468 Differential gene expression

469 The *C. secundus* and *M. natalensis* RNAseq libraries, were complemented with nine published *Z. nevaden-*
470 *sis* libraries, yielding 2 to 6 libraries from workers, queens and kings for each termite. These were com-
471 pared to six of the *B. germanica* libraries: two from 5th instar nymphs, two from 6th instar nymphs and
472 two from adult females. Reads were mapped to the genome using HiSat2²⁵. Read counts per gene where
473 obtained using htseq-count and DESeq2²⁶ was used for differential expression analysis.

474 Protein orthology

475 In addition to *B. germanica*, *C. secundus*, *Z. nevadensis* and *M. natalensis*, 18 other insect proteomes
476 were included in our analyses; *L. migratoria*, *R. prolixus*, *E. danica*, *D. melanogaster*, *A. aegypti*, *T. cas-*
477 *taneum*, *N. vitripennis*, *P. canadensis*, *A. mellifera*, *H. saltator*, *L. humile*, *C. floridanus*, *P. barbatus*,

478 *S. invicta*, *A. echinator* and *A. cephalotes*; as well as for the centipede, *S. maritima*, as an outgroup
479 (for sources see Table S23). These proteomes were grouped in to orthologous clusters with OrthoMCL²⁷,
480 with a granularity of 1.5.

481 IR and OR identification, phylogeny and structure

482 Ionotropic receptors (IRs) were identified using two custom Hidden Markov Models (HMMs) obtained
483 with `hmmbuild` and `hmmcompress` of the HMMER suite²⁸. The first HMM comprises the IR's ion channel and
484 ligand-binding domain based on a MAFFT²⁹ protein alignment of 76 IRs from 15 species (Table S24).
485 The second HMM was built to distinguish IRs from iGluRs, IR8a and IR25a, which have an additional
486 amino-terminal domain (ATD)³⁰. For this we built an HMM from 48 protein sequences (Table S24). The
487 proteomes were scanned with `pfam_scan` and the two custom HMMs, where proteins that matched the
488 IR HMM, but not the ATD HMM were annotated as IRs. ORs were identified based on the Pfam domain
489 PF02949 (7tm Odorant receptor).

490 Multiple sequence alignments of IRs and ORs were obtained with `hmmalign`²⁸, using the Pfam OR
491 HMM PF02949 and custom IR HMM to guide the alignment. Gene trees were computed with FastTree³¹
492 (options: `-pseudo -spr 4 -mlacc 2 -slowlni`) and visualised with iTOL v3³². Putative IR ligand-
493 binding residues and structural regions were identified based on the alignments with *D. melanogaster* IRs
494 and iGluRs of known structure³³.

495 Gene family expansions and contractions

496 For the analyses of gene family expansions and contractions, the hierarchical clustering algorithm
497 MC-UPGMA³⁴ was used, with a ProtoLevel cutoff of 80³⁵. Protein families were further divided into
498 sub-families if they contained more than 100 proteins in a single species, or more than an average of 35
499 proteins per species. Proteins were blasted against the RepeatMasker TE database (E-value < 10⁻⁵) and
500 clusters where > 50% of the proteins were identified as transposable elements were discarded. Clade- and
501 species-specific protein family expansions and contractions, were identified with CAFE v3.0³⁶ using the
502 same protocol as^{37,38} (see also Supplementary material).

503 TE-facilitated expansions

504 The repeat content in the 10 kb flanking regions of *B. germanica*, *C. secundus*, *Z. nevadensis* and
505 *M. natalensis* genes was calculated using bedtools³⁹. CDS' from neighbouring genes were removed and

506 the repeat content was analysed using Generalized Linear Mixed Models (glmmPQL implemented in the
507 R⁴⁰ package MASS⁴¹) with binomial error distribution. Fixed predictors included gene family expansion,
508 species ID and their interaction. Cluster ID was fitted as random factor to avoid pseudo-replication.
509 Significance was assessed based on the Wald-*t* test (R package aod⁴²) at $\alpha < 0.05$. Main and interaction
510 effects for each of the genomic regions are listed in table S8. Model parameters are listed in table S9.

511 Evolutionary rates

512 The rate of protein evolution (ω ; ratio of non-synonymous to synonymous substitutions) was estimated
513 for the OrthoMCL 1-to-1 orthologs in *L. migratoria*, *B. germanica*, *Z. nevadensis*, *C. secundus* and
514 *M. natalensis*. Protein sequences were aligned with t-coffee⁴³. CDS alignments were obtained with
515 pal2nal.pl⁴⁴ and trimmed with Gblocks⁴⁵. To identify genes with different rates of protein evolution
516 within the termites compared to outgroups, a set of codeml branch models was used (model = 2; NSsites
517 = 0; PAML suite⁴⁶). Specifically, we compared the null model (H_0 : one ω across all branches) to i) H_{A1} :
518 allowing for termite-specific ω ((Lmig,Bger,(Znev#1,(Csec#1, Mnat#1)#1)); and ii) H_{A2} : allowing
519 for different ω for different levels of eusocial complexity (Lmig,Bger,(Znev#1,(Csec#1, Mnat#2)#1)).
520 LR-tests were performed on unsaturated models (dS < 3) and p-values were Bonferroni-corrected. Gene
521 ontology enrichment of genes with significantly higher rates of protein evolution in termites was performed
522 with the TopGo⁴⁷ package in R.

523 To test for positive selection within gene families of interest, i) site model tests 7 and 8 were performed
524 (model = 0; NSsites = 7 8) on species-specific CDS alignments or ii) branch-site test (model = 2; NSsites
525 = 2; fix.omega = 1 for null model and 0 for alternative model) on multi-species alignments. Protein
526 sequences were aligned using MAFFT²⁹ with the E-INS-i strategy, and CDS alignments were created
527 using pal2nal.pl⁴⁴. Phylogenetic trees were created with FastTree³¹. Alignments were trimmed using
528 Gblocks (settings: -b2 = 21; -b3 = 20; -b4 = 5; -b5 = a). Models were compared using LR test and
529 where $p < 0.05$, Bayes Empirical Bayes (BEB) results were consulted for codon positions under positive
530 selection.

531 CpG depletion patterns and GO enrichment

532 To estimate DNA methylation we compared observed to expected CpG counts within CDS sequences^{48, 49}.
533 A low $CpG_{o/e}$ indicates a high level of DNA methylation, as the cytosine of methylated CpGs often
534 mutate to thymines. Expected CpG counts were calculated by dividing the product of cytosine and
535 guanine counts by the sequence length. The PCA in figure 3 was created using the R function precomp

536 on log transformed CpG_{o/e} values for all 1-to-1 orthologs for the seven hemimetabolous species. These
537 orthologs were extracted from the OrthoMCL results. The 3D plot was created with the plot3d command
538 from the R package rgl.

539 CpG depleted (first quartile) and enriched genes (fourth quartile) were tested for enrichment of Gene
540 Ontology terms. Pfam protein domains were obtained for *B. germanica*, *Z. nevadensis*, *C. secundus*
541 and *M. natalensis* protein sequences using PfamScan⁵⁰. Corresponding GO terms were obtained with
542 Pfam2GO. GO-term over-representation was assessed using TopGO⁴⁷ package in R. Enrichment analysis
543 was performed using the weight algorithm selecting nodesize=10 to remove terms with less than 10
544 annotated GO terms. After that GO terms classified as significant (topGOFisher;0.01) were visualized
545 using R package tagcloud (<https://cran.r-project.org/web/packages/tagcloud/>).

546 Data availability

547 The data reported in this study are archived at the following databases: NCBI (genomes sequences), SRA
548 (genomic and transcriptomic reads), i5k Workspace@NAL & Dryad (annotations). Detailed accession
549 information is tabulated in the Supplementary Materials (table S26).

550 Scripts and output files are available on request to E.B.B.

551 References (Materials and Methods)

- 552 1. Gnerre, S. *et al.* High-quality draft assemblies of mammalian genomes from massively parallel se-
553 quence data. *Proceedings of the National Academy of Sciences* **108**, 1513–1518 (2011).
- 554 2. Li, Y., Hu, Y., Bolund, L. & Wang, J. State of the art de novo assembly of human genomes from
555 massively parallel sequencing data. *Human Genomics* **4**, 271 (2010).
- 556 3. Grabherr, M. G. *et al.* Full-length transcriptome assembly from RNA-Seq data without a reference
557 genome. *Nature Biotechnology* **29**, 644–652 (2011).
- 558 4. Kim, D. *et al.* TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions
559 and gene fusions. *Genome Biology* **14**, R36 (2013).
- 560 5. Roberts, A., Trapnell, C., Donaghey, J., Rinn, J. L. & Pachter, L. Improving RNA-Seq expression
561 estimates by correcting for fragment bias. *Genome Biology* **12**, R22 (2011).
- 562 6. Ellinghaus, D., Kurtz, S. & Willhoeft, U. LTRharvest, an efficient and flexible software for de novo
563 detection of LTR retrotransposons. *BMC Bioinformatics* **9**, 18 (2008).
- 564 7. Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive elements in
565 eukaryotic genomes. *Mobile DNA* **6**, 11 (2015).
- 566 8. Chipman, A. D. *et al.* The First Myriapod Genome Sequence Reveals Conservative Arthropod Gene
567 Content and Genome Organisation in the Centipede *Strigamia maritima*. *PLOS Biology* **12**, e1002005
568 (2014).
- 569 9. Mesquita, R. D. *et al.* Genome of *Rhodnius prolixus*, an insect vector of Chagas disease, reveals
570 unique adaptations to hematophagy and parasite infection. *Proceedings of the National Academy of
571 Sciences* **112**, 14936–14941 (2015).
- 572 10. Nene, V. *et al.* Genome Sequence of *Aedes aegypti*, a Major Arbovirus Vector. *Science* **316**, 1718–
573 1723 (2007).
- 574 11. Leadership, O. p. *et al.* Insights into social insects from the genome of the honeybee *Apis mellifera*.
575 *Nature* **443**, 931–949 (2006).
- 576 12. Gadau, J. *et al.* The genomic impact of 100 million years of social evolution in seven ant species.
577 *Trends in Genetics* **28**, 14–21 (2012).

- 578 13. Richards, S. *et al.* The genome of the model beetle and pest *Tribolium castaneum*. *Nature* **452**,
579 949–955 (2008).
- 580 14. Wang, X. *et al.* The locust genome provides insight into swarm formation and long-distance flight.
581 *Nature Communications* **5**, 2957 (2014).
- 582 15. Holt, C. & Yandell, M. MAKER2: an annotation pipeline and genome-database management tool
583 for second-generation genome projects. *BMC Bioinformatics* **12**, 491 (2011).
- 584 16. Keller, O., Kollmar, M., Stanke, M. & Waack, S. A novel hybrid gene prediction method employing
585 protein multiple sequence alignments. *Bioinformatics* **27**, 757–763 (2011).
- 586 17. Borodovsky, M., Mills, R., Besemer, J. & Lomsadze, A. Prokaryotic Gene Prediction Using GeneMark
587 and GeneMark.hmm. In *Current Protocols in Bioinformatics* (John Wiley & Sons, Inc., 2002). DOI:
588 10.1002/0471250953.bi0405s01.
- 589 18. Korf, I. Gene finding in novel genomes. *BMC Bioinformatics* **5**, 59 (2004).
- 590 19. Birney, E., Clamp, M. & Durbin, R. GeneWise and Genomewise. *Genome Research* **14**, 988–995
591 (2004).
- 592 20. Stanke, M. *et al.* AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Research*
593 **34**, W435–W439 (2006).
- 594 21. Elsik, C. G. *et al.* Creating a honey bee consensus gene set. *Genome Biology* **8**, R13 (2007).
- 595 22. Campbell, M. A., Haas, B. J., Hamilton, J. P., Mount, S. M. & Buell, C. R. Comprehensive analysis
596 of alternative splicing in rice and comparative analyses with Arabidopsis. *BMC Genomics* **7**, 327
597 (2006).
- 598 23. Kong, L. *et al.* CPC: assess the protein-coding potential of transcripts using sequence features and
599 support vector machine. *Nucleic Acids Research* **35**, W345–W349 (2007).
- 600 24. Min, X. J., Butler, G., Storms, R. & Tsang, A. OrfPredictor: predicting protein-coding regions in
601 EST-derived sequences. *Nucleic Acids Research* **33**, W677–W680 (2005).
- 602 25. Pertea, M., Kim, D., Pertea, G. M., Leek, J. T. & Salzberg, S. L. Transcript-level expression analysis
603 of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nature Protocols* **11**, 1650–1667
604 (2016).
- 605 26. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq
606 data with DESeq2. *Genome Biology* **15**, 550 (2014).

- 607 27. Li, L., Stoeckert, C. J. & Roos, D. S. OrthoMCL: Identification of Ortholog Groups for Eukaryotic
608 Genomes. *Genome Research* **13**, 2178–2189 (2003).
- 609 28. Eddy, S. R. Profile hidden Markov models. *Bioinformatics* **14**, 755–763 (1998).
- 610 29. Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7: Improve-
611 ments in Performance and Usability. *Molecular Biology and Evolution* **30**, 772–780 (2013).
- 612 30. Croset, V. *et al.* Ancient Protostome Origin of Chemosensory Ionotropic Glutamate Receptors and
613 the Evolution of Insect Taste and Olfaction. *PLOS Genetics* **6**, e1001064 (2010).
- 614 31. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree: Computing Large Minimum Evolution Trees
615 with Profiles instead of a Distance Matrix. *Molecular Biology and Evolution* **26**, 1641–1650 (2009).
- 616 32. Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and annotation
617 of phylogenetic and other trees. *Nucleic Acids Research* **44**, W242–245 (2016).
- 618 33. Benton, R., Vannice, K. S., Gomez-Diaz, C. & Vosshall, L. B. Variant Ionotropic Glutamate Receptors
619 as Chemosensory Receptors in *Drosophila*. *Cell* **136**, 149–162 (2009).
- 620 34. Loewenstein, Y., Portugaly, E., Fromer, M. & Linial, M. Efficient algorithms for accurate hierarchical
621 clustering of huge datasets: tackling the entire protein space. *Bioinformatics* **24**, i41–i49 (2008).
- 622 35. Rappoport, N., Linial, N. & Linial, M. ProtoNet: charting the expanding universe of protein se-
623 quences. *Nature Biotechnology* **31**, 290–292 (2013).
- 624 36. Han, M. V., Thomas, G. W. C., Lugo-Martinez, J. & Hahn, M. W. Estimating Gene Gain and
625 Loss Rates in the Presence of Error in Genome Assembly and Annotation Using CAFE 3. *Molecular
626 Biology and Evolution* **30**, 1987–1997 (2013).
- 627 37. Simola, D. F. *et al.* Social insect genomes exhibit dramatic evolution in gene composition and
628 regulation while preserving regulatory features linked to sociality. *Genome Research* **23**, 1235–1247
629 (2013).
- 630 38. Kapheim, K. M. *et al.* Genomic signatures of evolutionary transitions from solitary to group living.
631 *Science* **348**, 1139–1143 (2015).
- 632 39. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.
633 *Bioinformatics* **26**, 841–842 (2010).
- 634 40. Team, R. C. R: A language and environment for statistical computing (2012).

- 635 41. Venables, W. & Ripley, B. *Modern Applied Statistics with S* (Springer, New York, 2002), fourth edn.
- 636 42. Lesnoff, M., Lancelot & R. *aod: Analysis of Overdispersed Data* (2012). R package version 1.3.
- 637 43. Notredame, C., Higgins, D. G. & Heringa, J. T-Coffee: A novel method for fast and accurate multiple
638 sequence alignment. *Journal of Molecular Biology* **302**, 205–217 (2000).
- 639 44. Suyama, M., Torrents, D. & Bork, P. PAL2nal: robust conversion of protein sequence alignments
640 into the corresponding codon alignments. *Nucleic Acids Research* **34**, W609–W612 (2006).
- 641 45. Castresana, J. Selection of conserved blocks from multiple alignments for their use in phylogenetic
642 analysis. *Molecular Biology and Evolution* **17**, 540–552 (2000).
- 643 46. Yang, Z. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Molecular Biology and Evolution*
644 **24**, 1586–1591 (2007).
- 645 47. Alexa, A. & Rahnenfuhrer, J. topGO: Enrichment analysis for Gene Ontology (2010).
- 646 48. Bewick, A. J., Vogel, K. J., Moore, A. J. & Schmitz, R. J. Evolution of DNA methylation across
647 insects. *Molecular Biology and Evolution* 654–655 (2017).
- 648 49. Park, J. *et al.* Comparative Analyses of DNA Methylation and Sequence Evolution Using Nasonia
649 Genomes. *Molecular Biology and Evolution* **28**, 3345–3354 (2011).
- 650 50. Finn, R. D. *et al.* The Pfam protein families database: towards a more sustainable future. *Nucleic*
651 *Acids Research* **44**, D279–D285 (2016).
- 652 51. Bell, W. J., Roth, L. M. & Nalepa, C. A. *Cockroaches: ecology, behavior, and natural history* (JHU
653 Press, Baltimore, Maryland, 2007).

Figures

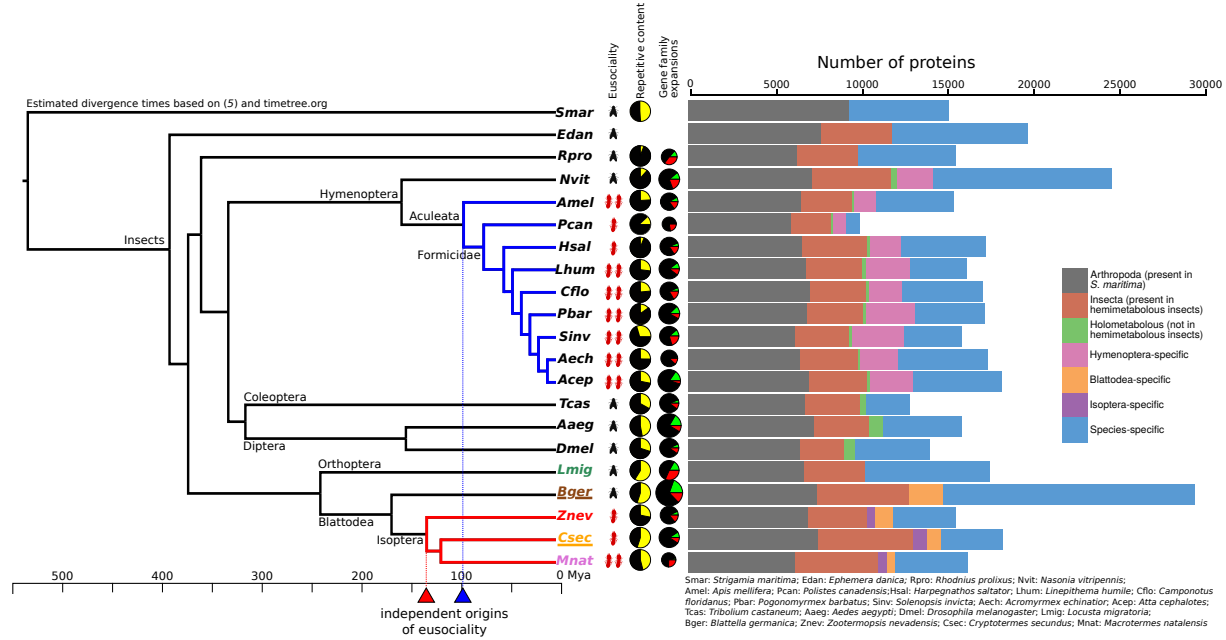


Figure 1: **Phylogenetic, genomic and proteomic comparisons of 20 insect species.** From left to right: Phylogenetic tree of 20 insect species with *Strigamia maritima* (centipede) as outgroup; level of eusociality (one red insect: simple eusociality; two red insects: advanced eusociality; black fly: non-eusocial); fractions of repetitive content (yellow) within genomes of selected species (for sources see supplementary material); proportions of species-specific gene family expansions (green), contractions (red) and stable gene families (black), size of pies represents relative size of gene family change (based on total numbers). Bar chart showing protein orthology across taxonomic groups within each genome.

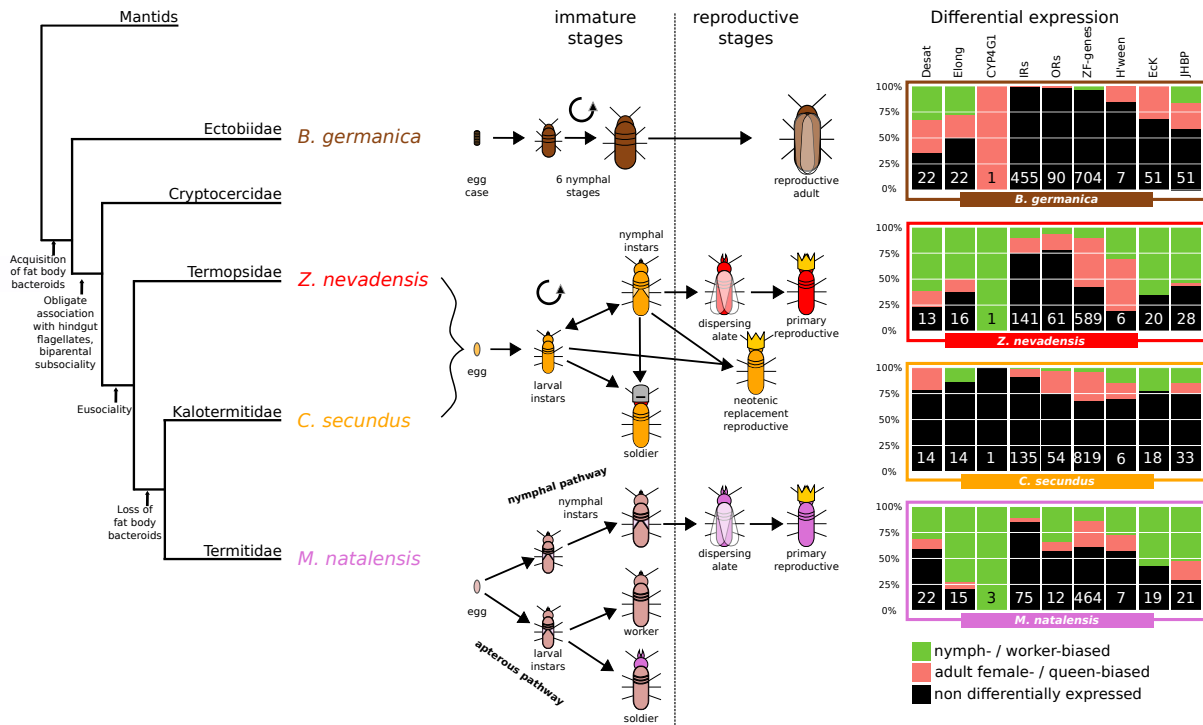


Figure 2: Comparison of developmental pathways between *B. germanica*, the lower termites, *Z. nevadensis* and *C. secundus*, and the higher termite, *M. natalensis*. Shown from left to right are: a simple phylogeny⁵¹ describing important novelties along the evolutionary trajectory to termites; life cycles; differential expression between workers and queens (between nymphs and adult females in *B. germanica*) of selected gene families (Desat = desaturases, Elong = elongases, H'ween = Halloween genes; numbers denote total numbers of genes in each gene family).

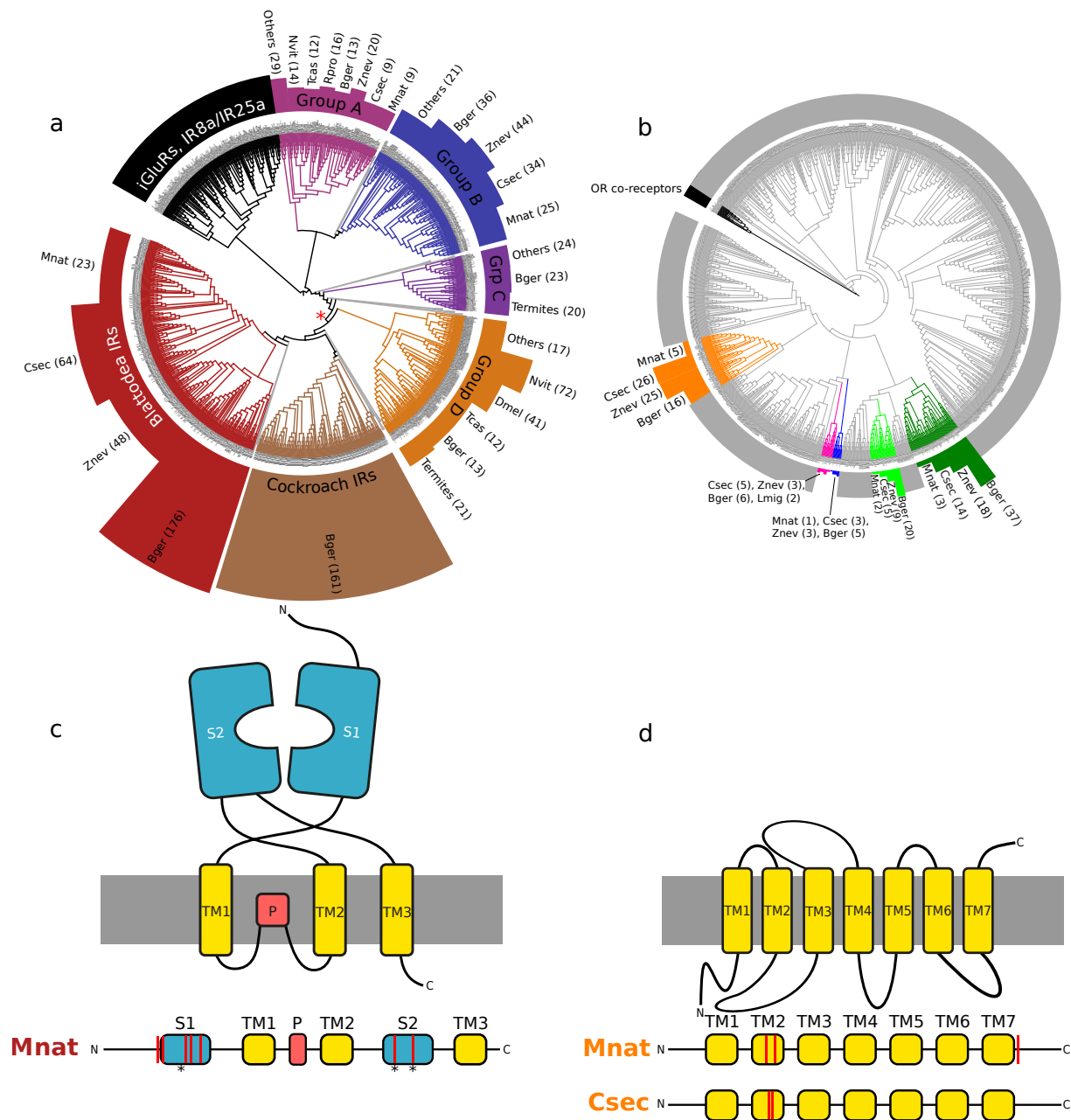


Figure 3: Expansions, contractions and positive selection within IRs and ORs in termites. IR (a) and OR (b) gene trees of 13 insect species. Only well supported clades (support values > 85) that include *B. germanica* or termite genes are highlighted within the gene trees. Lengths of coloured bars represent number of genes per species within each of these clades. Red asterisk in (a) denotes putative root of intronless IRs. 2D structure and sites under positive selection (red bars; codeml site models 7 & 8) for Blattodea-IR genes in *M. natalensis* ($p < 1.7 \times 10^{-10}$) (c) (asterisks denote putative ligand binding sites³³) and orange OR genes in *M. natalensis* ($p = 1.1 \times 10^{-11}$) and *C. secundus* ($p = 5.6 \times 10^{-16}$) (d).

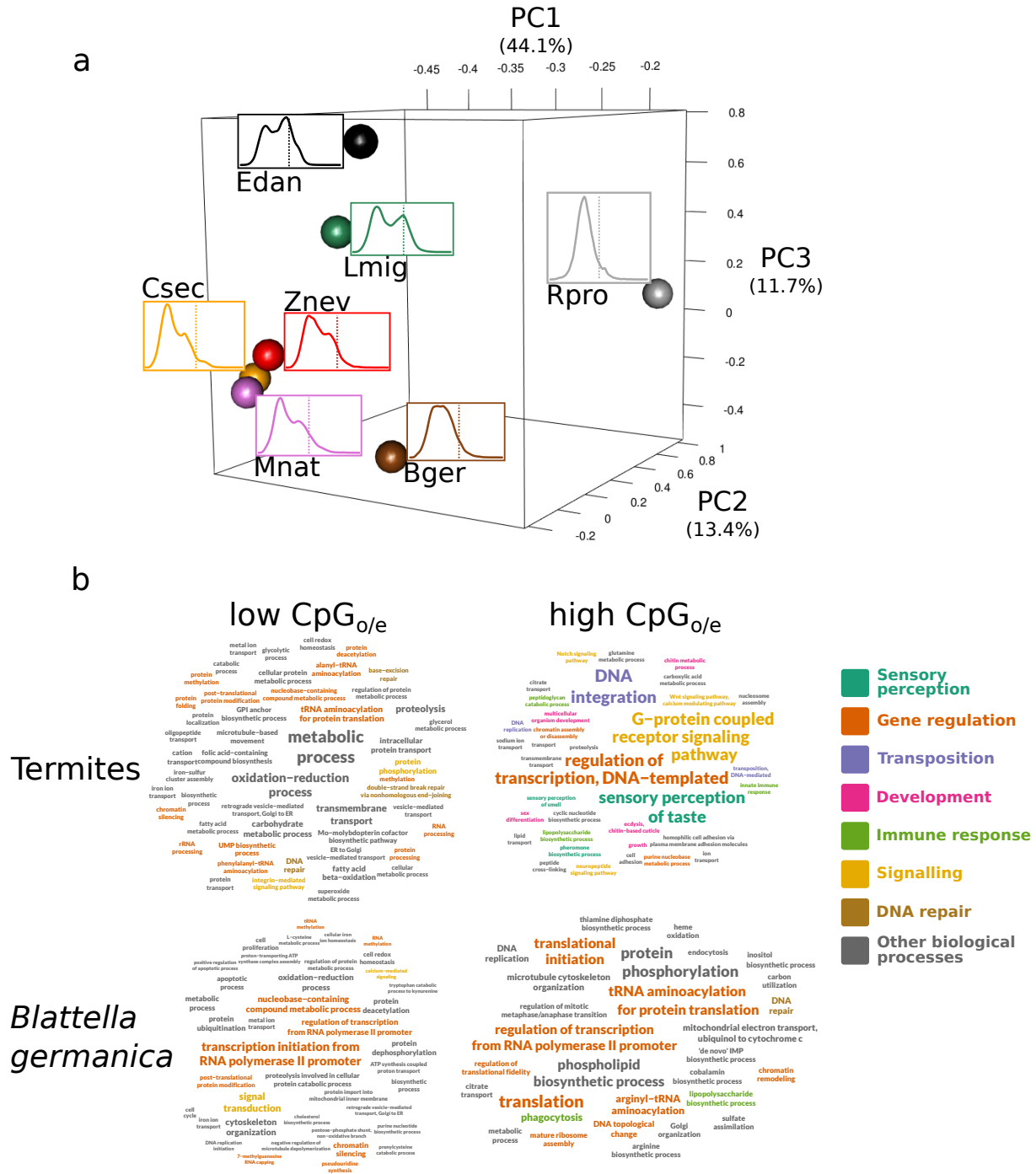


Figure 4: $CpG_{o/e}$ of seven hemimetabolous insects. a) PCA of DNA methylation patterns among 2664 1-to-1 orthologs, estimated via $CpG_{o/e}$. Spheres represent positions of species within 3D PCA; curves are distribution of $CpG_{o/e}$ with dotted line showing $CpG_{o/e} = 1$. b) Tag clouds of enriched ($p < 0.05$) GO terms (biological processes) among lower (left) and higher quartile (right) of $CpG_{o/e}$ within termites (top) and *B. germanica* (bottom). For termites, genes were merged from all three species for analysing GO term enrichment.

High $CpG_{o/e}$ indicates low level of DNA methylation and vice versa.