

1 **Genomic and transcriptomic analyses of the**  
2 **subterranean termite *Reticulitermes speratus*: gene**  
3 **duplication facilitates social evolution**

4

5 Shuji Shigenobu<sup>\*+1,2</sup>, Yoshinobu Hayashi<sup>+3</sup>, Dai Watanabe<sup>4,5</sup>, Gaku Tokuda<sup>6</sup>,  
6 Masaru Y Hojo<sup>6,19</sup>, Kouhei Toga<sup>5,7</sup>, Ryota Saiki<sup>5</sup>, Hajime Yaguchi<sup>5,8</sup>, Yudai  
7 Masuoka<sup>5,9</sup>, Ryutaro Suzuki<sup>5,10</sup>, Shogo Suzuki<sup>5</sup>, Moe Kimura<sup>11</sup>, Masatoshi  
8 Matsunami<sup>4,12</sup>, Yasuhiro Sugime<sup>4</sup>, Kohei Oguchi<sup>4,10,13</sup>, Teruyuki Niimi<sup>2,14</sup>, Hiroki  
9 Gotoh<sup>15,20</sup>, Masaru K Hojo<sup>8</sup>, Satoshi Miyazaki<sup>16</sup>, Atsushi Toyoda<sup>17</sup>, Toru Miura<sup>\*4,13</sup>,  
10 Kiyoto Maekawa<sup>\*18</sup>

11

12 1) NIBB Research Core Facilities, National Institute for Basic Biology, Okazaki,  
13 444-8585 Japan

14 2) Department of Basic Biology, School of Life Science, The Graduate University for  
15 Advanced Studies, SOKENDAI, Nishigonaka 38, Myodaiji, Okazaki, Aichi 444-  
16 8585, Japan

17 3) Department of Biology, Keio University, Hiyoshi, Yokohama, 223-8521, Japan

18 4) Faculty of Environmental Earth Science, Hokkaido University, Sapporo,  
19 Hokkaido, 060-0810, Japan

20 5) Graduate School of Science and Engineering, University of Toyama, Toyama,  
21 930-8555, Japan

22 6) Tropical Biosphere Research Center, COMB, University of the Ryukyus,  
23 Nishihara, Okinawa 903-0213, Japan

24 7) Department of Biosciences, College of Humanities and Sciences, Nihon  
25 University, Setagaya-ku, Tokyo, 156-8550, Japan

26 8) Department of Bioscience, School of Science and Technology, Kwansai Gakuin  
27 University, Sanda, Hyogo, 669-1337, Japan

28 9) Institute of Agrobiological Sciences, NARO (National Agriculture and Food  
29 Research Organization), Tsukuba, Japan

30 10) Bioproduction Research Institute, National Institute of Advanced Industrial  
31 Science and Technology (AIST), Tsukuba, Japan

32 11) School of Science, University of Toyama, Toyama, 930-8555, Japan

33 12) Graduate School of Medicine, University of the Ryukyus, Nishihara, Okinawa  
34 903-0125, Japan

35 13) Misaki Marine Biological Station, School of Science, The University of Tokyo,  
36 Miura, Kanagawa, 238-0225, Japan

- 37 14) Division of Evolutionary Developmental Biology, National Institute for Basic  
38 Biology, Okazaki, 444-8585 Japan
- 39 15) Ecological Genetics Laboratory, Department of Genomics and Evolutionary  
40 Biology, National Institute of Genetics, Mishima, Shizuoka, 411-8540, Japan
- 41 16) Graduate School of Agriculture, Tamagawa University, Machida, Tokyo 194-  
42 8610, Japan
- 43 17) Advanced Genomics Center, National Institute of Genetics, Mishima, Shizuoka,  
44 411-8540 Japan
- 45 18) Faculty of Science, Academic Assembly, University of Toyama, Gofuku,  
46 Toyama, 930-8555, Japan
- 47 19) Global Education Institute, University of the Ryukyus, Nishihara, Okinawa 903-  
48 0213, Japan
- 49 20) Department of Biological Sciences, Faculty of Science, Shizuoka University,  
50 Shizuoka, Shizuoka 422-8529, Japan
- 51
- 52
- 53 \* corresponding authors. Email: [shige@nibb.ac.jp](mailto:shige@nibb.ac.jp), [miu@mmbs.s.u-tokyo.ac.jp](mailto:miu@mmbs.s.u-tokyo.ac.jp),  
54 [kmaekawa@sci.u-toyama.ac.jp](mailto:kmaekawa@sci.u-toyama.ac.jp)
- 55
- 56 + contributed equally
- 57

## 58 Summary

59 Termites are model social organisms characterized by a polyphenic caste system.  
60 Subterranean termites (Rhinotermitidae) are ecologically and economically  
61 important species, including acting as destructive pests. Rhinotermitidae occupies  
62 an important evolutionary position within the clade representing an intermediate  
63 taxon between the higher (Termitidae) and lower (other families) termites. Here, we  
64 report the genome, transcriptome and methylome of the Japanese subterranean  
65 termite *Reticulitermes speratus*. The analyses highlight the significance of gene  
66 duplication in social evolution in this termite. Gene duplication associated with  
67 caste-biased gene expression is prevalent in the *R. speratus* genome. Such  
68 duplicated genes encompass diverse categories related to social functions,  
69 including lipocalins (chemical communication), cellulases (wood digestion and  
70 social interaction), lysozymes (social immunity), geranylgeranyl diphosphate  
71 synthase (social defense) and a novel class of termite lineage-specific genes with  
72 unknown functions. Paralogous genes were often observed in tandem in the  
73 genome, but the expression patterns were highly variable, exhibiting caste biases.  
74 Some duplicated genes assayed were expressed in caste-specific organs, such as  
75 the accessory glands of the queen ovary and frontal glands in soldier heads. We  
76 propose that gene duplication facilitates social evolution through regulatory  
77 diversification leading to caste-biased expression and subfunctionalization and/or  
78 neofunctionalization that confers caste-specialized functions.  
79

## 80 Significance Statement

81 Termites are model social organisms characterized by a sophisticated caste  
82 system, where distinct castes arise from the same genome. Our genomics data of  
83 Japanese subterranean termite provides insights into the evolution of the social  
84 system, highlighting the significance of gene duplication. Gene duplication  
85 associated with caste-biased gene expression is prevalent in the termite genome.  
86 Many of the duplicated genes were related to social functions, such as chemical  
87 communication, social immunity and defense, and they often expressed in caste-  
88 specific organs. We propose that gene duplication facilitates social evolution

89 through regulatory diversification leading to caste-biased expression and functional  
90 specialization. In addition, since subterranean termites are ecologically and  
91 economically important species including destructive pests in the world, our  
92 genomics data serves as a foundation for these studies.

## 93 Introduction

94           The evolution of eusociality, i.e., animal societies defined by the  
95 reproductive division of labor, cooperative brood care and multiple overlapping  
96 generations, represents one of the major transitions in evolution, having increased  
97 the level of biological complexity (1). Eusocial insects such as bees, wasps, ants  
98 and termites show sophisticated systems based on the division of labor among  
99 castes, which is one of the pinnacles of eusocial evolution (2). Recent advances in  
100 molecular biological technologies and omics studies have revealed many molecular  
101 mechanisms underlying eusociality and have led to the establishment of a new field  
102 of study known as “sociogenomics” (3). The genomes of major eusocial  
103 hymenopteran lineages, i.e., ants, bees and wasps, have been sequenced, and the  
104 differences in gene expression, DNA methylation (4–6)(7) (8–11) and histone  
105 modification (12) (13) among castes have been explored. These sociogenomics  
106 studies in hymenopterans revealed some genetic bases of social evolution,  
107 including the co-option of genetic toolkits of conserved genes, changes in protein-  
108 coding genes, cis-regulatory evolution leading to genetic network reconstruction,  
109 epigenetic modifications and taxonomically restricted genes (TRG) (14, 15).

110

111           Isoptera (termites) is another representative insect lineage exhibiting highly  
112 sophisticated eusociality and a wide range of social complexities (16). Termites are  
113 hemimetabolous and diploid insects that are phylogenetically distant from  
114 hymenopterans with holometaboly and haplodiploidy. Termite societies are  
115 characterized by reproductives of both sexes, workers and soldiers. In the termite  
116 lineage, eusociality is thought to have evolved once, although the levels of social  
117 complexity features, such as colony size, feeding habitat, symbiosis with  
118 microorganisms and caste developmental pathways, diverged among termite  
119 species. These characteristics are especially different between the two major  
120 termite sublineages, i.e., the early-branching families (called “lower” termites) and  
121 the most apical family Termitidae (“higher” termites) (Fig. 1a). To date, based on  
122 the whole-genome sequences of a few termite species, the commonality and  
123 diversity of genetic repertoires between Isoptera and Hymenoptera or between  
124 termite lineages and their solitary outgroup (i.e., cockroach) have been investigated  
125 (17, 18). Additionally, in *Zootermopsis nevadensis*, clear differences in gene

126 expression levels among castes (17) and in DNA methylation between alates and  
127 workers (19) were detected.

128

129         Among the more than 2900 extant species of termites (Isoptera) (20),  
130 subterranean termites (Rhinotermitidae), especially two genera, *Reticulitermes* and  
131 *Coptotermes*, occupy an important evolutionary position (Fig 1a). Recent  
132 phylogenetic studies showed that Rhinotermitidae is paraphyletic, and a clade  
133 including *Reticulitermes*, *Coptotermes* and *Heterotermes* was shown to be sister to  
134 Termitidae (21, 22). In particular, *Reticulitermes* exhibits intermediate social  
135 complexity between those of higher (Termitidae) and lower (all the other families)  
136 termites (23), for example, this genus displays primitive feeding ecology and gut  
137 symbiont features, a relatively complex colony structure and a caste development  
138 mode termed the bifurcated pathway (Fig. 1b). Moreover, *Reticulitermes* is the most  
139 common termite group in palearctic (24) and nearctic (25) regions and a major pest  
140 causing serious damage to human-made wooden structures (26). For these  
141 reasons, members of this genus are probably among the most studied termites  
142 (16). Nevertheless, despite their evolutionary, ecological and economic relevance,  
143 subterranean termites remain an understudied group in terms of both genetics and  
144 genomics.

145

146         In this study, we targeted the Japanese subterranean termite *Reticulitermes*  
147 *speratus*. We conducted whole-genome sequencing, caste-specific RNA-seq  
148 analysis and whole-genome bisulfite sequencing of *R. speratus* to understand the  
149 genomic, transcriptomic and epigenetic bases of the social life of this termite  
150 species. *R. speratus* nymphoids are almost exclusively produced  
151 parthenogenetically by automixis with terminal fusion in primary queens, such that  
152 the genome should be homozygous at most loci (27), which provides an advantage  
153 in genome sequencing. We also compared the omics data of *R. speratus* with those  
154 of sequenced higher and lower termites (Fig. 1a). Our integrative analyses of the  
155 genome and transcriptome of *R. speratus* and other termites revealed that gene  
156 duplications are often associated with caste-biased gene expression and caste-  
157 specific functions, which highlights the significant role of gene duplication in  
158 eusocial evolution in the termite lineage.

159

## 160 Results and Discussion

### 161 Genomic features of *Reticulitermes speratus*

162 Genome sequencing of *R. speratus* was performed with genomic DNA isolated  
163 from female secondary reproductives (nymphoids) [Fig. 1b]. *R. speratus* nymphoids  
164 are almost exclusively produced parthenogenetically by automixis with terminal  
165 fusion in primary queens, such that the genome should be homozygous at most loci  
166 (27) and thus ease de Bruijn-graph-based genome assembly. We generated a total  
167 of 86 Gb of Illumina HiSeq sequence data and assembled them *de novo* into 5817  
168 scaffolds with an N50 of 1.97 Mb and total size of 881 Mb [Table 1], covering 88%  
169 of the genome based on the genome size (1.0 Gb) estimated by flow cytometry  
170 (28). The assembled *R. speratus* genome has high coverage of coding regions,  
171 capturing 99.2% (98.5% complete; 0.7% fragmented) of 1367 Insecta  
172 benchmarking universal single-copy orthologs (BUSCOs) (29) [Table 1]. The *R.*  
173 *speratus* genome is rich in repetitive elements, which make up 40.4% of the  
174 genome. A total of 15,591 protein-coding genes were predicted by combining the  
175 reference-guided assembly of RNA-seq reads (36 libraries derived from different  
176 castes, sexes and body parts; see below for details) and homology-based gene  
177 prediction followed by manual curation of gene families of interest [Fig. 1c]. Whole-  
178 genome bisulfite sequencing revealed extensive gene body methylation of the *R.*  
179 *speratus* genome, amounting to 8.8% of methylated cytosines in the CG context  
180 [Supplementary Fig. 1]. The genome-wide DNA methylation landscape was similar  
181 to that of a dampwood termite *Z. nevadensis* (12%) (19). These omics data and a  
182 genome browser are available at <http://www.termite.nibb.info/retsp/>.

183 We compared the *R. speratus* gene repertoire with those of 88 other  
184 arthropods, including the two termites *Z. nevadensis* and *Macrotermes natalensis*  
185 (17, 30). Ortholog analysis showed that 12,032 (82.9%) of the 15,591 genes in *R.*  
186 *speratus* genes were shared with other arthropods, and 1773 were taxonomically  
187 restricted (TRGs) to Isoptera, among which 430 were shared with the other two  
188 termites and 1343 were unique to *R. speratus* [Fig. 1c]. Whole-genome comparison  
189 with two sequenced termites, *M. natalensis* and *Z. nevadensis*, showed a high  
190 degree of synteny conservation [Fig. 1d]. We identified 2799 syntenic blocks (N50:  
191 858.4 kb) shared with *M. natalensis* that covered 95.1% of the *R. speratus* genome

192 where 560.4 Mb of nucleotides was aligned, while 3650 syntenic blocks (N50: 591.1  
193 kb) shared with *Z. nevadensis* covered 72.1% of the *R. speratus* genome where  
194 116.7 Mb was aligned. Only a few cases of large genomic rearrangements were  
195 found between termite genomes, at least, at the contiguity level of the current  
196 assemblies, suggesting overall conservation of genome architecture in the lineage  
197 of termites over 135 MY (Fig. 1a). Interestingly, despite such high conservation of  
198 macrosynteny, interruptions or breaks in local synteny were observed and often  
199 associated with tandem gene duplications. For example, when we examined  
200 regions containing large tandem gene duplications (> 5-gene tandem duplications),  
201 synteny between the *R. speratus* and *M. natalensis* genomes was interrupted in 10  
202 of 21 regions (examples shown in Supplementary Fig 3).

203

## 204 Transcriptome differentiation among castes

205 Distinct castes arise from the same genome, a phenomenon called caste  
206 polyphenism which is a distinctive hallmark in social insects (31, 32). To elucidate  
207 caste-biased gene expression in order to understand the mechanism underlying the  
208 caste-specific phenotypes, we compared the transcriptomes of three castes  
209 (primary reproductives, workers and soldiers) in *R. speratus* [Fig. 1b;  
210 Supplementary Table 2]. We sequenced 36 RNA-seq libraries, representing three  
211 biological replicates of both sexes and two body parts (“head” and “thorax +  
212 abdomen”) for each of the three castes.

213 The results clearly showed that termite castes were distinctively  
214 differentiated at the gene expression level. The multidimensional scaling (MDS) plot  
215 depicted the three castes as clearly distinct transcriptomic clusters for both the  
216 head and thorax + abdomen transcriptomes [Fig. 2a]. However, little sexual  
217 difference was detected within each caste, although reproductives showed  
218 substantial transcriptomic differences in thorax + abdomen samples between  
219 queens and kings, probably due to the difference in the reproductive organs [Fig.  
220 2a]. Using a generalized linear model (GLM) with caste and sex as explanatory  
221 variables, we identified 1579 and 2076 genes differentially expressed among castes  
222 in head and thorax + abdomen samples, respectively, with the criteria of a false  
223 discovery rate (FDR)-corrected  $P < 0.01$ , while we identified only 6 and 79 genes



224 that were differentially expressed between sexes in the head and thorax +  
225 abdomen samples, respectively, with the same criteria. We focused on the genes  
226 that were differentially expressed among castes (caste-DEGs) and further classified  
227 them into three categories of caste-biased genes (i.e., reproductive-, worker-, and  
228 soldier-biased genes), with a criterion of >2-fold higher expression relative to that in  
229 the other two castes. These caste-biased genes should account for the specialized  
230 functions of each caste. Soldier samples exhibited the highest number of caste-  
231 specific genes, suggesting the highly specialized functions of the soldier caste. This  
232 is consistent with the finding of a previous RNA-seq analysis of the Eastern  
233 subterranean termite *Reticulitermes flavipes*, reporting that a majority of DEGs were  
234 soldier-specific (33); 73 of the 93 DEGs identified were up- or downregulated  
235 specifically in the soldier caste. In addition to these soldier-specific *R. flavipes*  
236 genes (e.g., *troponin C* and fatty *acyl-CoA reductase*), the caste-biased genes  
237 identified in our transcriptome analysis of *R. speratus* included genes previously  
238 reported as caste-biased genes in other termites (34–37), e.g., *vitellogenin*  
239 (reproductives), *geranylgeranyl pyrophosphate synthase* (soldiers), and *beta-*  
240 *glucosidase* (probably associated with cellulase; workers). This consistency  
241 between transcriptome analyses of different termite species indicates that the RNA-  
242 seq analysis in this study is reliable and that the regulation and perhaps the  
243 functions of these caste-biased genes are conserved across the termite lineage.

244

245 The caste-DEGs were enriched for Gene Ontology terms related to a wide  
246 array of functions [Supplementary Table 4], such as hormone metabolism, chitin  
247 metabolism, hydrolase activity, oxidoreductase activity, lipid metabolism, signaling  
248 and lysozyme activity. Protein motifs enriched in the caste-DEGs were also  
249 identified, including cytochrome P450, lipocalin, lysozyme, glycosyl hydrolase  
250 family, TGF-beta and trypsin [Supplementary Table 5]. Among the 1773  
251 taxonomically restricted genes (TRGs) that were restricted to Isoptera (see above),  
252 the termite-shared TRGs showed strong enrichment for caste-DEG (Fisher's exact  
253 test,  $P < 1.0e-7$  for head samples,  $P < 1.0e-10$  for thorax+abdomen samples), while  
254 the TRGs found only in *R. speratus* (orphan genes) did not ( $P = 0.99$  and  $P =$   
255  $0.97$ , respectively).

256

257 To investigate the relationship between caste-biased gene expression and  
258 DNA methylation, we analyzed differential methylation levels among three castes  
259 (reproductives, workers, and soldiers). The BS-seq data showed that the global  
260 CpG methylation patterns were very similar among the castes [Supplementary Fig.  
261 2ab], in contrast to the methylation pattern of *Z. nevadensis*, in which DNA  
262 methylation differed strongly between castes (winged adults vs. final-instar larvae)  
263 and was strongly linked to caste-specific splicing (19). Instead, gene body DNA  
264 methylation of *R. speratus* seems to be important for the expression of  
265 housekeeping genes, as reported in the drywood termite *Cryptotermes secundus*  
266 (18). Housekeeping genes exhibited a high degree of gene body methylation in all  
267 castes of *R. speratus*, while caste-biased genes showed a significantly lower level  
268 of DNA methylation [Supplementary Fig. 2cd].

## 269 Gene duplication and caste-biased gene expression

270 Evolutionary novelties are often brought about by gene duplications (38) (reviewed  
271 in (39)), and the transition to eusociality in Hymenoptera has been associated with  
272 gene family expansion (18, 40, 41). Our ortholog analysis comparing the *R.*  
273 *speratus* gene repertoire with those of 88 other arthropods identified 1396  
274 multigene families duplicated in the *R. speratus* genome. Interestingly, compared to  
275 the genome as a whole, the set of caste-DEGs identified above was significantly  
276 enriched for genes in multigene families (X-squared = 218.62, df = 1, p-value <  
277 2.2e-16). We also calculated the tau score as a proxy of caste specificity of gene  
278 expression for all genes and found that duplicated genes were significantly more  
279 caste-specific than single-copy genes (p < 2.2e-16, Wilcoxon rank sum test) in both  
280 transcriptome data sets (head and thorax+abdomen) [Fig. 2b]. Additionally, gene  
281 set tests showed that sets of duplicated genes were differentially expressed in all  
282 pairwise comparisons between castes [Fig. 2c]. These data highlight the important  
283 roles of gene duplication in the caste evolution of termites.  
284

285 Multigene families related to caste-specific traits in *R.*  
286 *speratus*

287 Caste-biased multigene families were associated with diverse functional categories,  
288 some of which were strongly related to caste-specific behaviors and tasks. Here,  
289 we highlight five families, namely, lipocalins (protein transporters for social  
290 communication and physiological signaling), cellulases (carbohydrate-active  
291 enzymes for worker wood digestion), lysozymes (immune-related genes for social  
292 immunity), geranylgeranyl diphosphate synthases (metabolic enzymes for the  
293 production of soldier defensive chemicals), and a novel termite-specific gene family  
294 with unknown functions, as examples of multigene families relevant to termite  
295 sociality. Molecular evolution studies have shown that the redundancy caused by  
296 gene duplication may allow one paralog to acquire a new function  
297 (neofunctionalization) or divide the ancestral function among paralogs  
298 (subfunctionalization) (38, 39). We are particularly interested in the evolutionary  
299 impact of gene duplication on caste specialization through neo/subfunctionalization.

300 Lipocalins

301 Lipocalins belong to a family of proteins, with molecular recognition properties such  
302 as the ability to bind a range of small hydrophobic molecules (e.g. pheromones)  
303 and specific cell surface receptors, and to form complexes with soluble  
304 macromolecules (42). A previous study identified a gene of the lipocalin family,  
305 SOL1, that is exclusively expressed in the mandibular glands of mature soldiers of  
306 the rotten-wood termite *Hodotermopsis sjostedti* (43). SOL1 is thought to function  
307 as a signaling molecule for defensive social interactions among termite colony  
308 members (31). Moreover, RNA-seq analysis showed that a lipocalin gene, *Neural*  
309 *Lazarillo homolog 1* (*ZnNLaz1*), was specifically expressed in soldier-destined larvae  
310 in an incipient colony of *Z. nevadensis* (44). Gene function and protein localization  
311 analyses suggested that ZnNLaz1 was a crucial regulator of soldier differentiation  
312 through the regulation of trophallactic interactions with a queen. Thus, it was of  
313 interest that the lipocalin-related motif (Pfam PF00061; lipocalin/cytosolic fatty-acid  
314 binding protein family) was significantly enriched in the list of caste-DEGs (FDR <  
315 0.05; Supplementary Table 5).

316 We identified 18 lipocalin family genes in the *R. speratus* genome [Fig. 3a-c;  
317 Supplementary Table 7]. The number of lipocalin genes was larger than those in  
318 other insects [Fig. 3c]. Phylogenetic analysis of lipocalin family genes identified in  
319 arthropods, including three termite species, revealed a highly dynamic evolutionary  
320 history of this protein family [Fig. 3a]. A couple of subfamilies, namely clades A and  
321 B, had experienced extensive expansion in the termite lineage. The most drastic  
322 expansion was found in clade A, which includes *H. sjostedti* SOL1. In this clade, 7,  
323 9 and 5 genes were identified in *R. speratus*, *M. natalensis* and *Z. nevadensis*,  
324 respectively, and extensive and independent gene expansions occurred in each  
325 species. Clade B was also composed of genes with a termite lineage specific  
326 expansion. In many cases, these lipocalin genes were found in tandem arrays in  
327 the *R. speratus* genome [Fig. 3b]. The inferred phylogenetic tree indicated that  
328 duplications in each clade occurred after the divergence of termites from a common  
329 ancestor.

330 A comparison of the transcriptome among castes revealed that most  
331 lipocalin genes (15 of 18) showed caste-biased gene expression [Fig. 2, Fig. 3ab].  
332 The caste specificity, however, varied among genes, regardless of sequence  
333 similarity and positional proximity on the genome. In particular, the expression  
334 levels of genes in clades A and B drastically changed among castes. For example,  
335 *RS008823* and *RS008824* displayed soldier-specific expression, the expression of  
336 *RS013912* was biased toward workers, and *RS013913* was downregulated in  
337 soldiers. *RS008881* and *RS008884* were exclusively expressed in queen bodies  
338 (thorax + abdomen), while *RS008882*, a gene next to the aforementioned two  
339 genes, showed quite different expression patterns and high expression levels in  
340 heads, especially those of workers. These results indicate that termite lipocalin  
341 genes underwent dynamic expansion in terms of gene repertoire, and regulatory  
342 diversification of caste-biased expression. This gene expansion and regulatory  
343 diversification of lipocalins may facilitate the evolution of the molecules involved in  
344 signaling during caste development and among individuals through social  
345 interactions.

346 To address the caste-specific function of the lipocalin paralogs, the  
347 expression patterns of several selected caste-biased lipocalin genes were  
348 examined by *in situ* hybridization [Fig. 3d-h, Supplementary Fig. 6]. *RS008881*, a  
349 queen-biased lipocalin gene, was found to be expressed exclusively in the  
350 accessory glands of the ovary [Fig. 3d]. The next gene on the same scaffold,

351 *RS008882*, was shown to be specifically expressed in worker antennae and  
352 maxillary/labial palps [Fig. 3f-h]. *RS008823*, a soldier-biased gene, was expressed  
353 exclusively in the frontal gland cells of the soldier heads [Fig. 3e]. Note that ovaries  
354 and frontal glands develop during postembryogenesis in a caste-specific manner  
355 (i.e., ovaries in queens and frontal glands in soldiers) in the pathway of caste  
356 differentiation in *R. speratus*. Antennae and maxillary/labial palps are not caste-  
357 specific but crucial sensory organs, especially for blind termite immatures, such as  
358 workers. Given that animal lipocalins generally work as carrier proteins (45), there  
359 is a possibility that focal termite lipocalins bind and convey some molecules to the  
360 targets from caste-specific organs (e.g., egg-recognition pheromone and soldier  
361 defensive and/or inhibiting substances; (46–48)), or participate in sensory  
362 reception, such as the role of odorant-binding proteins (49).

### 363 Cellulases

364 Lignocellulose degradation in termites is achieved by a diverse array of  
365 carbohydrate-active enzymes (CAZymes) produced by the host and their intestinal  
366 symbionts. The repertoire of CAZyme families in the genome of *R. speratus* did not  
367 show considerable differences from those of other nonxylophagous insects, such as  
368 a honeybee and a fruit fly (Supplementary Fig. 4). However, we found gene family  
369 expansion and expressional diversification for glycoside hydrolase family (GH) 1  
370 and GH9 members. The majority of GH1 and GH9 members are  $\beta$ -glucosidase  
371 (BGs; EC 3.2.1.21) and endo- $\beta$ -1,4-glucanases (EGs; EC 3.2.1.4), respectively,  
372 which are essential for cellulose digestion in termites (50).

373 We identified 16 GH1 paralogs [Supplementary Table 8]. Such gene expansion  
374 of GH1 was also observed in the genome of other termites, but the reason for the  
375 gene expansion remains elusive (51). Although the phylogenetic tree divided these  
376 GH1 paralogs into four distinct groups (clades A to D in Fig. 4a), most of them were  
377 tandemly located in the genome of *R. speratus* (Fig. 4bc). The predominantly  
378 expressed BG gene was *RS004136*, while the expression of this gene was clearly  
379 biased toward the bodies (thorax + abdomen) of reproductives and workers (Fig.  
380 4b). This gene formed a rigid clade with *bona fide* BGs reported from the salivary  
381 glands or midgut of termites (clade A in Fig. 4a) (52), suggesting that this gene is  
382 involved in cellulose digestion in *R. speratus*. Indeed, *in situ* hybridization analysis  
383 showed that *RS004136* was specifically expressed in the salivary glands of workers

384 (Fig. 4de, Supplementary Fig. 7a). Other GH1 members showed a wide variety of  
385 expression patterns across castes and body parts [Fig. 4b]. Some of them might  
386 have diversified their functions, other than wood digestion, related to termite  
387 sociality, such as egg-recognition pheromones (53). A typical example displaying  
388 such diversification was *RS004624*, which was expressed specifically in the  
389 abdomens of queens [Fig. 4b]. The peptide sequence of this gene showed a  
390 monophyletic relationship with that of Neofem2 of *Cryptotermes secundus* (clade D  
391 in Fig. 4a), which is a queen recognition pheromone probably functioning in the  
392 suppression of reproductive emergence (54). *In situ* hybridization showed that  
393 *RS004624* was specifically expressed in the accessory glands of queen ovaries  
394 (Fig. 4fg, Supplementary Fig. 7b), suggesting that *RS004624* is involved in  
395 enzymatic activities in queen-specific glands. Together with the results for a queen-  
396 biased lipocalin (*RS008884*), this finding indicates that the queen accessory glands  
397 may produce some queen-specific pheromones. Like lipocalins, GH1 paralogs are  
398 also typical examples of multigene family members participating in caste-specific  
399 tasks, which may be acquired by gene duplication resulting in neo- or  
400 subfunctionalization.

401 We found four paralogs of GH9 in *R. speratus* [Supplementary Fig. 5,  
402 Supplementary Table 8]. Although several insect GH9 EGs have acquired the  
403 ability to hydrolyze hemicellulose (55), neo- or subfunctionalization of termite EGs  
404 has yet to be clarified. Intriguingly, we found that the GH9 member *RS006396* was  
405 weakly but uniformly expressed across all termite body parts and castes. This result  
406 suggests that some GH9 members also perform a function other than that of  
407 cellulase, as is the case for GH1.

## 408 Lysozymes

409 The immune system of termites is of particular interest, because the group living of  
410 termites with nonsclerotized and nonpigmented epidermis and microbe-rich habitat  
411 puts them at high risk for pathogenic infections (56). Thus, defense against  
412 pathogenic microbes is important for termites. In the *R. speratus* genome we  
413 identified 251 immune-related genes [Supplementary Information 1.9,  
414 Supplementary Table 26]. The repertoire and number of immune-related genes of  
415 *R. speratus* showed no large differences compared to those of other insect species,  
416 but a notable exception was found for lysozymes [Supplementary Fig. 8].

417 Lysozymes are involved in bacteriolysis through hydrolysis of  $\beta$ -1,4-linkages in the  
418 peptidoglycans present in bacterial cell walls, and three distinct types of lysozymes,  
419 chicken- or conventional-type (c-type), goose-type (g-type), and invertebrate-type (i-  
420 type) lysozymes, have been found in animals (57). We identified 13 and 3 genes  
421 encoding c-type and i-type lysozymes, respectively, and the number of lysozyme  
422 genes was larger than those in other insects [Fig. 5, Supplementary Table 9].  
423 Phylogenetic analysis revealed that c-type lysozymes underwent extensive gene  
424 duplications in the sublineage leading to *R. speratus* [Fig. 5a]. Seven c-type  
425 lysozymes formed a tandem array on scaffold\_859 [Fig. 5b], probably generated by  
426 repeated tandem gene duplication events. Interestingly, most of the c-type  
427 lysozyme genes showed caste-biased expression. Three genes (*RS014698*,  
428 *RS100022*, and *RS100023*) exhibited high expression levels compared to those of  
429 other lysozyme genes and were expressed in a soldier-specific manner, while  
430 *RS100026* was expressed in a worker-specific manner and *RS100024* and  
431 *RS100025* were highly expressed in both workers and soldiers [Fig. 5b]. The  
432 differential expression patterns of the lysozyme genes in *R. speratus* may represent  
433 division of labor among castes in terms of colony-level immunity.

434 It is also possible that duplicated lysozymes may have functions other than  
435 immunity. A previous study indicated that the salivary glands of *R. speratus* secrete  
436 c-type lysozymes to digest bacteria ingested by termites through social feeding  
437 behavior (58). The same lysozyme genes are also expressed in the queen ovaries  
438 and eggs and play a role in egg recognition as proteinaceous pheromones in *R.*  
439 *speratus* (48, 53). We could not find identical sequences of these lysozyme genes  
440 in our gene models, but these sequences were most closely related to *RS002400*  
441 with 88% nucleotide identity, which occupied the basal position of the lineage-  
442 specific gene expansion (Fig. 5a).

## 443 GGPP synthase

444 Whole-genome comparison of *R. speratus* with *Z. nevadensis* and *M. natalensis*  
445 revealed a 270 kb *R. speratus*-specific fragment in scaffold\_31, while the rest of this  
446 scaffold showed very high syntenic conservation among the three termites [Fig 6a].  
447 We found that the *R. speratus*-specific region was encompassed by a tandemly  
448 duplicated gene cluster composed of 13 genes encoding geranylgeranyl  
449 diphosphate (GGPP) synthase [Fig. 6b, Supplementary Table 10]. GGPP synthase

450 catalyzes the consecutive condensation of an allylic diphosphate with three  
451 molecules of isopentenyl diphosphate to produce GGPP, an essential precursor for  
452 the biosynthesis of diterpenes, carotenoids and retinoids (59–61). The extensive  
453 duplication of GGPP synthase paralogs observed in *R. speratus* is unusual  
454 because the genomes of other insects surveyed have only a single copy of GGPP  
455 synthase gene. The phylogenetic analysis of GGPP synthase homologs revealed  
456 two clusters, a possibly ancestral group (including *RS007484*) and an apical group  
457 (including other paralogs identified) [Fig. 6c]. The latter cluster also contained some  
458 GGPP synthase paralogs obtained from the termitid *Nasutitermes takasagoensis*  
459 (34)

460 Transcriptome data indicated that all of the GGPP synthase genes, except  
461 *RS007484* which was a member of the ancestral group in the phylogenetic tree,  
462 showed caste-biased expression, and caste specificity varied across the paralogs  
463 [Fig 6b]. Specifically, *RS100010*, *RS007480*, *RS100012*, *RS100015*, *RS100016*,  
464 *RS100017* and *RS007483* showed soldier-specific expression, while *RS007481*,  
465 *RS007482* and *RS100013* showed reproductive-specific expression [Fig 6b].  
466 Several GGPP synthase genes have been identified in some termite species and  
467 are known to function in a caste-specific manner; for example, the soldiers of *N.*  
468 *takasagoensis* synthesize defensive polycyclic diterpenes by high expression of the  
469 GGPP synthase gene in the frontal gland to use chemical defense (62). It has been  
470 reported that the soldiers of *Reticulitermes* have a frontal gland in which diterpenes  
471 are synthesized, although the biological role is not fully understood (63–65).  
472 Consequently, it is possible that the soldier-specific GGPP synthases identified to  
473 date are involved in chemical defense. Indeed, *in situ* hybridization revealed that  
474 the soldier-specific GGPP synthase *RS100016* was expressed exclusively in the  
475 soldier frontal gland, as shown in a previous study (66) [Fig. 6d, Supplementary Fig.  
476 7c]. It is also possible that reproductive-specific GGPP synthases are involved in  
477 the metabolism of other diterpenes, such as pheromone synthesis, especially  
478 *RS007481*, which shows strong queen-specific expression in the thorax and  
479 abdomen and may play a role in the synthesis of queen substances.

480 Under the branch-site (BS) model of codon substitutions (67), significant  
481 positive selection was detected on five branches of *R. speratus* GGPP synthase  
482 family tree [Fig. 6e]: ancestral branches #1 and #2, and the branches leading to  
483 *RS100017* (branch #3), *RS100012* (branch #4) and *RS007483* (branch #5). These  
484 results suggest that all GGPP synthase paralogs of *R. speratus* except the



485 ancestral type *RS007484* have experienced positive selection and finally acquired  
486 novel roles for the production of defensive and/or pheromonal substances.  
487

## 488 The TY family, a novel gene family restricted to termites

489 Numerous studies have shown that novel genes (e.g., TRG) play important roles in  
490 the evolution of novel social phenotypes in hymenopteran social insects (8, 68, 69).  
491 We found that termite-shared TRGs showed strong enrichment for caste-DEGs  
492 (see above). A striking example of caste-biased TRGs is a tandem array of three  
493 novel genes [Fig. 7a; Supplementary Table 11], *RS001196*, *RS001197* and  
494 *RS001198*, that have no significant homologs in any organisms outside termite  
495 clades. These three genes were expressed at extremely high levels (up to 250,000  
496 RPKM), which constituted approximately 30% of the worker head transcriptome,  
497 and strongly biased across the three castes [Fig. 7a]. Each gene was composed of  
498 a single exon encoding a short peptide ~60 aa in length that contained a secretion  
499 signal peptide in the N-terminal region followed by a middle part rich in charged  
500 amino acid residues and C-terminal part rich in polar amino acids with unusually  
501 high number of tyrosine residues [Fig. 7b]. Here, we named this novel class of  
502 peptides the termite-specific tyrosine-rich peptide family (TY family). The three TY  
503 genes showed modest sequence similarity with each other, suggesting that they are  
504 paralogs derived by tandem duplication. TY family orthologs were also found in the  
505 genomes of *Z. nevadensis* and *M. natalensis* [Fig. 7b]. We estimated pairwise  
506 evolutionary rates (the ratio of nonsynonymous to synonymous substitutions, i.e.,  
507 dN/dS) between *R. speratus* and *Z. nevadensis* for these three peptides. The  
508 dN/dS for each gene ranged from 0.03 to 0.16 (Fig. 7c), indicating that they evolved  
509 under strong purifying selection and suggesting a conserved function in the termite  
510 lineage. Indeed, *Z. nevadensis* orthologs were also expressed at a high level in the  
511 soldier and worker castes in a pattern similar to that in *R. speratus*.

## 512 Facilitation of caste specification by gene duplication

513 Recent advances in sociogenomics in different social insects are promoting our  
514 understanding of the genetic bases of social evolution, which include the co-option

515 of genetic toolkits of conserved genes, changes in protein-coding genes, cis-  
516 regulatory evolution leading to genetic network reconstruction, epigenetic  
517 modifications and TRGs (15, 70). In addition to these components, our genomic  
518 and transcriptomic analyses in *R. speratus* highlighted the significance of gene  
519 duplication for caste specialization. Gene duplication is, in general, a key source of  
520 genetic innovation that plays a role in the evolution of phenotypic complexity; gene  
521 duplication allows for subsequent divergent evolution of the resultant gene copies,  
522 enabling evolutionary innovations in protein functions and/or expression patterns  
523 (71–73). Regarding eusocial evolution in insects, Gadagkar (74) first pointed out the  
524 importance of gene duplication; 'genetic release followed by diversifying evolution'  
525 made possible the appearance of multiple caste phenotypes in social insects. Many  
526 decades later, genomic analyses revealed gene family expansion, especially in  
527 relation to chemical communication, in both ants (odorant receptors; (75–77)) and  
528 termites (ionotropic receptors; (18)). Based on Godagkar's hypothesis, duplicated  
529 genes can be released from the constraints of original selection, leading to new  
530 directional evolution, i.e., for caste-specific functions (e.g., queen- or worker-trait  
531 genes). However, the detailed roles and significance of gene duplication in social  
532 evolution have been elusive.

533         This study revealed that gene duplication associated with caste-biased gene  
534 expression is prevalent in the *R. speratus* genome. The list of duplicated genes  
535 encompasses a wide array of functional categories related to the social behaviors in  
536 termites as exemplified by transporters such as lipocalins (communication and  
537 physiological signaling; cf. (31, 44)), digestive enzymes such as carbohydrate-  
538 active enzymes, immune-related genes such as lysozymes (social immunity), and  
539 metabolic enzymes such as GGPP synthase (social defense). This study  
540 demonstrated that caste-specific expression patterns differed among in-paralogs.  
541 Although such paralogous genes were often observed in tandem in the genome,  
542 the expression patterns were often independent from one another, showing  
543 differential caste biases in many cases. Additionally, discordant caste biases in  
544 transcriptional expression were observed among closely related paralogs with  
545 similar coding sequences, as represented by little correlation between phylogenetic  
546 position and caste specificity (Fig. 3a). Although the regulatory and evolutionary  
547 mechanisms underlying caste-biased expression patterns are elusive, these  
548 examples strongly suggest that gene duplications have facilitated caste  
549 specialization, leading to social evolution in termites.

550           After the gain of caste-biased gene regulation, subfunctionalization and/or  
551 neofunctionalization seems to have occurred, leading to caste-specific expression  
552 and caste-specialized functions. For example, in the case of lipocalin family,  
553 lipocalin paralogs were generated by lineage-specific functional expansion in caste-  
554 specific organs or tissues: a queen-specific lipocalin (*RS008881*) was expressed  
555 specifically in the ovarian accessory glands, while a soldier-biased lipocalin  
556 (*RS008823*), was expressed exclusively in the frontal glands in soldier heads [Fig.  
557 3de]. Taken together, we hypothesize that, in termites, caste specification through  
558 gene duplication proceeds by the following three steps: 1) gene family expansion by  
559 tandem gene duplication, 2) regulatory diversification leading to an expression  
560 pattern restricted to a certain caste, and 3) subfunctionalization and/or  
561 neofunctionalization of the gene products conferring caste-specific functions. As an  
562 exaptation of these steps, the case in which one (or some) of the multiple functions  
563 of pleiotropic genes are allocated and specialized to a duplicated gene copy might  
564 have led to caste-specific subfunctionalization (38, 39).

565           Recently, it was suggested that the evolution of phenotypic differences  
566 among castes in the honey bee was associated with the gene duplication, by  
567 showing that duplicated genes had higher levels of caste-biased expression  
568 compared to singleton genes (78). It was also shown that the level of gene  
569 duplication was correlated with social complexity in bees (superfamily Apoidea)  
570 (78). Given the independent origin of eusociality in termites and honeybees, gene  
571 duplications might be a shared mechanism facilitating the evolution of caste  
572 systems in social insects.

573

574

## 575   Materials and Methods

### 576   Insects

577   All mature colonies of *Reticulitermes speratus* used for genome, RNA, and Bisulfite  
578 sequencing (BS-seq), were collected in Furudo, Toyama Prefecture, Japan  
579 [Supplementary Table 1]. Detailed sample information is described in *SI Appendix*,  
580 *Supplementary Methodology*.

## 581 Sample collection, genome sequencing and assembly

582 All colonies of *Reticulitermes speratus* used for genome, RNA, and bisulfite  
583 sequencing were collected at Furudo, Toyama Prefecture, Japan. Detailed sample  
584 information is described in Supplementary Table 1 and *SI Appendix*,  
585 *Supplementary Methodology*.

586

587 We used female secondary reproductives (nymphoids I and II) for genome  
588 sequencing. We excluded the gut and ovaries of nymphoids to avoid contamination  
589 by DNAs from the king or other microorganisms. Genomic DNA was isolated from  
590 each individual using a Genomic-tip 20/G (Qiagen). We used 5 microsatellite loci  
591 (Rf6-1, Rf21-1, Rf24-2, Rs02, and Rs03) to confirm whether they were homozygous  
592 at these loci and shared the same genotype. The purified genomic DNA purified  
593 was fragmented with a Covaris S2 sonicator (Covaris), size-selected with  
594 BluePippin (Sage Science), and then used to create two pair-end libraries using a  
595 TruSeq DNA Sample Preparation Kit (Illumina) with insert sizes of ~250 and ~800  
596 bp [Supplementary Table 3]. Four Mate-pair libraries with peaks at ~3 kb, ~5 kb, ~8  
597 kb and ~10 kb, respectively, were also created using a Nextera Mate Pair Sample  
598 Preparation Kit (Illumina) [Supplementary Table 3]. These libraries were sequenced  
599 using an Illumina HiSeq system with 2 × 151 bp paired-end sequencing protocol.  
600 Reads of the pair-end and mate-pair libraries were assembled using ALLPATHS-  
601 LG (build# 47878), with default parameters. BUSCO  
602 v4.0.6(29)(<https://busco.ezlab.org/>) was used in quantitative measuring for the  
603 assessment of genome assembly using *insecta\_odb10* as the lineage input. A  
604 genome browser was built using JBrowse (<https://jbrowse.org/>)

## 605 Gene prediction

606 A protein-coding gene reference set was generated with two main sources of  
607 evidence, aligned *R. speratus* transcripts and aligned homologous proteins of other  
608 insects, and a set of *ab initio* gene predictions. RNA-seq reads were assembled *de*  
609 *novo* using Trinity, and then mapped to the genome using Exonerate. We  
610 processed homology evidence at the protein level using the reference proteomes of  
611 7 sequenced insects including *Z. nevadensis* and Blattodea protein sequences  
612 predicted from RNA-seq of *Periplaneta americana* and *Nasutitermes*  
613 *takasagoensis*. These proteins were split-mapped to the *R. speratus* genome with

614 Exonerate. These models were merged using the EvidenceModeler (EVM), which  
615 yielded 15584 gene models. Seventy-four genes were manually inspected and  
616 corrected. In particular, tandemly duplicated genes were liable to be incorrect gene  
617 prediction with erroneous exon–exon connections across homologs. The final set of  
618 15591 genes was designated as Rspe OGS1.0 [Supplementary Data 2  
619 (DOI:10.6084/m9.figshare.14267381)]. The quality of theOGS1.0 was evaluated by  
620 assessing two types of evidence, homology and expression. Among 15591 genes,  
621 12996 (83.3%) showed any hits in the NCBI nr database, 10440 (70.0%) included  
622 known protein motifs defined in the Pfam database, and 14302 (91.7%) showed  
623 evidence of expression with a threshold of RPKM = 1.0 in any sample of caste-  
624 specific RNA-seq data. In sum, 15577 (99.9%) had evidence for the presence of  
625 homologs and/or expression.

## 626 Orthology inference and gene duplication analysis

627 Orthology determination among three termites: Orthologous genes among the  
628 proteomes of three termite species, *R. speratus*, *Z. nevadensis*, and *M. natalensis*  
629 (gene models RspeOGS1.0, ZnevOGSv2.229, and MnatOGS3, respectively), were  
630 determined by pairwise comparisons with InParanoid v4.1 followed by three-  
631 species comparison with MultiParanoid. *M. natalensis* gene set, MnatOGS3, was  
632 built in this study using a similar pipeline as used for *R. speratus*.  
633 Ortholog analysis with arthropod proteomes: Orthology relationships of *R. speratus*  
634 genes (OGS1.0) with other arthropod genes were analyzed by referring to the  
635 OrthoDB gene orthology database ver.8 (87 arthropod species)  
636 (<https://www.orthodb.org/>). We grouped *R. speratus* genes with the OrthoDB  
637 ortholog group using a two-step clustering procedure. For each *R. speratus* protein,  
638 BLASTP was used to find similar proteins among the arthropod proteins, and the  
639 ortholog group of the top hit was provisionally assigned to the query *R. speratus*  
640 gene. Then, the ortholog grouping was evaluated by comparing the similarity level  
641 (BLAST bit score) among members within the focal ortholog group. We keep the  
642 grouping if the BLAST bit score between the query *R. speratus* gene and top  
643 arthropod gene was higher than the minimal score within the original cluster  
644 members. Among 15591 *R. speratus* OGS1.0 genes, 12434 were clustered into  
645 9033 OrthoDB Arthropod ortholog groups. Gene duplication was assessed based  
646 on this clustering. If two or more members of one species were included in a single  
647 ortholog group, they were regarded as a multigene family.

## 648 RNA-seq

649 W4–5 workers (old workers) and soldiers were collected from each colony. To  
650 collect primary reproductives, dealated adults were chosen randomly from each  
651 colony in accordance with the method of the previous literature (79), and female–  
652 male pairs were mated (Supplementary Table 1). Kings and queens were sampled  
653 after 4 months. Each individual was divided into head and body parts (thorax +  
654 abdomen). We prepared RNA-seq libraries for 12 categories based on castes  
655 (reproductives, workers and soldiers), sexes (males and females) and body parts  
656 (head, and thorax + abdomen). Three biological replications of the 12 categories  
657 were made with three different field colonies totaling 36 RNA-seq libraries  
658 [Supplementary Table 2]. All Illumina libraries prepared using a TruSeq Stranded  
659 mRNA Library Prep kit were subjected to a single-end sequencing of 101 bp  
660 fragments on HiSeq 2500. The cleaned reads were mapped onto the genome with  
661 TopHat v2.1.0 guided by the OGS1.0 gene models. Transcript abundances were  
662 estimated using featureCounts and normalized with the trimmed mean of M-values  
663 (TMM) algorithm in edgeR. Differentially expressed genes among castes and  
664 between sexes were detected in each body part (head / thorax and abdomen) using  
665 a generalized linear model with two factors, namely, caste and sex using edgeR  
666 with the conditions set as false discovery rate (FDR) < 0.01 and the log2 fold  
667 change of the expression level > 1.

## 668 Data Availability

669 Data from whole-genome sequencing, transcriptome sequencing, and methylome  
670 sequencing have been deposited in the DDBJ database under BioProject  
671 accessions PRJDB2984, PRJDB5589 and PRJDB11323, respectively. The  
672 analyzed data including genome assembly, gene prediction, annotation, and gene  
673 expression are available through FigShare  
674 (<https://doi.org/10.6084/m9.figshare.c.5483235>). The *R. speratus* genome browser  
675 is available at <http://www.termite.nibb.info/retsp/>.

## 676 Code availability

677 Custom R and Ruby scripts were deposited into Github  
678 ([https://github.com/termiteg/retsp\\_genome\\_paper](https://github.com/termiteg/retsp_genome_paper)).

679

## 680 **Acknowledgments**

681 We thank R. H. Suzuki and A. Karasawa for experimental support, T. Nishiyama  
682 and M. Hasebe for discussion on genome analyses, N. Kanasaki and K. Kai for  
683 rearing insects, T. Shibata, S. Ohi, T. Aizu, H. Ishizaki, H. Asao for next-generation  
684 sequencing (NGS), and K. Yamaguchi for NGS data management. Computations  
685 were partially performed on the supercomputers at the Data Integration and  
686 Analysis Facility, National Institute for Basic Biology. This study was funded by the  
687 JSPS/MEXT KAKENHI Grant Numbers 25128705, 24570022, 16K07511,  
688 JP19H03273, 22128008, 19K22294, 221S0002 and NIBB Collaborative Research  
689 Programs (20-323).

690

691

## 692 **Author contributions**

693 S.S., Y.H., T.M., and K.M. designed and managed the project. D.W., K.T., R.S.,  
694 H.Y., Y.M., R.S., and K.M. collected samples. D.W., R.S., Y.M., and R.S. performed  
695 the DNA extraction. S.S., and A.T. performed the library construction and genome  
696 sequencing. Y.H., D.W., K.T., R.S., H.Y., and Y.M. generated the RNA-Seq data.  
697 S.S., Y.H., and R.S. generated the BS-Seq data. S.S., Y.H., D.W., G.T., M.Y.H.,  
698 K.T., M.M., Y.S., K.O., T.N., H.G., M.K.H., and S.M. contributed to the genome  
699 assembly and annotation. S.Su, and M.K. performed histological analyses. S.S.,  
700 Y.H., G.T., T.M., and K.M. drafted the manuscript. All authors contributed to the  
701 final version of the manuscript.

702

703 Tables

704 **Table 1. Summary of *Reticulitermes speratus* genome assembly, annotation**  
705 **and methylome**

Genome	No. scaffolds	5,817
	No. contigs	63,310
	Total length	881 Mb
	Scaffold N50	1,967 kb
	Contig N50	37.5 kb
	Longest scaffold	14.3 Mb
	GC%	39.70%
	No. Ns	63 Mb
	Completeness (BUSCO insecta_odb10)	C:98.5% [S:98.1%, D:0.4%]
Annotation	No. genes (coding)	15,591
	Repeat content	40.4%
Methylome	%methylated CpG	8.79%

706



## 707 Figure legends

708 **Figure 1: Phylogenetic position of *Reticulitermes speratus* in Blattodea, its**  
709 **developmental pathway, and evolution of the gene repertoire and genome**  
710 **structure.**

711 **(a)** Phylogenetic tree of termites and cockroaches. Estimated divergence dates  
712 (mya: million years ago) are based on Bucek et al. (80). *R. speratus* is marked in  
713 red, and two termites mainly compared in this study are marked with bold  
714 characters. **(b)** Developmental pathway of *R. speratus*. There are 2 larval stages  
715 before the molt into a nymph (with wing buds) or worker (no wing buds). There are  
716 6 imaginal stages, and the 6th-stage nymphs molt into alates, which are primary  
717 reproductives (queen and king). Secondary reproductives (neotenics) differentiate  
718 from the 3rd- to 6th-stage nymphs. In the apterous line, there are at least 5 stages  
719 of workers. Some workers in the colony molt into presoldiers and soldiers. Female  
720 neotenics used for genome sequencing and 3 castes used for RNA-seq are marked  
721 with asterisks. **(c)** Gene repertoire of *R. speratus* categorized by orthology. *R.*  
722 *speratus* genes were compared to those of 88 arthropods and grouped into three  
723 classes: orthologs shared with other arthropods (labeled 'arthropod'), orthologs  
724 shared with other termites (*Z. nevadensis* and/or *M. natalensis*) but with no  
725 orthologs in other arthropods (labeled 'termite'), and orphan genes unique to *R.*  
726 *speratus* (labeled 'no hit'). **(d)** High conservation of synteny between termite  
727 genomes revealed by dot plots generated by comparing *R. speratus* with *Z.*  
728 *nevadensis* and *M. natalensis*. Scaffolds longer than 2.0 Mb in the *R. speratus*  
729 assembly are used for plotting. Forward alignments are plotted in red and reverse  
730 alignments are plotted in blue.

731

732 **Figure 2: Caste-specific transcriptome analysis and the enrichment of**  
733 **duplicate genes for caste-biased genes.**

734 **(a)** Multidimensional scaling (MDS) plot of RNA-seq data showing relatedness  
735 between the expression profiles of different castes (reproductive, soldier and  
736 worker) and sexes (male and female). The left panel plots RNA-seq data from head  
737 samples, and the right panel plots data from thorax + abdomen samples. Three  
738 biological replicates were analyzed for each condition and plotted individually. **(b)**  
739 Numbers of caste-biased genes with >2-fold higher expression levels than the other

740 two castes. Colours in each bar indicate the differences of RNA-seq data obtained.  
741 **(c)** Violin plots showing the distribution of tau indexes of duplicate genes and single  
742 genes. Tau values range between 0 and 1, with low values indicating invariable and  
743 constitutive expression between castes and higher values supporting caste  
744 specificity. In both body part samples, the tau values of duplicate genes were  
745 significantly greater than those of single genes ( $p < 2.2e-16$ , Wilcoxon rank sum  
746 test). **(d)** Enrichment of caste-DEGs (differentially expressed genes among castes)  
747 for duplicate genes. In each comparison between castes (soldier vs worker,  
748 reproductive vs worker, and reproductive vs soldier), all genes are ranked and  
749 ordered by log-fold-change value along the horizontal axis. Black bars mark the  
750 positions of genes. Genes of sociality-related functions highlighted in the text are  
751 selected and plotted in the lower panels. Curved lines in the upper panel show  
752 relative enrichment of the duplicate genes (blue line) or single genes (green line)  
753 relative to uniform ordering.

754

755 **Figure 3: Lipocalin genes in *R. speratus*.**

756 **(a)** Maximum likelihood (ML) tree of lipocalin homologs based on the amino acid  
757 sequences obtained with a log gamma (LG) model. Branches leading to clade A  
758 and clade B, which show gene family expansion specific to termite sublineages, are  
759 marked in yellow and green, respectively. Color gradients in the outer tracks show  
760 the expression levels as averaged  $\log(\text{RPKM}+1)$  values in three castes  
761 (reproductive, soldier, and worker). Expression levels of head samples and thorax +  
762 abdomen samples are shown in purple and green, respectively. Caste-DEGs  
763 (differentially expressed genes among castes) are marked as R, S, or W beside the  
764 color gradients, indicating biases toward the reproductive, soldier, or worker caste,  
765 respectively. **(b)** Lipocalin multigene clusters in the *R. speratus* genome and their  
766 relative expression levels among castes. The heatmap shows the Z-scores of the  
767  $\log(\text{RPKM}+1)$  values in the caste-specific transcriptome. **(c)** Comparison of the  
768 number of lipocalin subclasses among representative arthropods. Note clades A  
769 and B are specific to termites. **(d)** Vertical cryosection of the queen abdomen  
770 subjected to *in situ* hybridization with an antisense DIG-labeled *RS008881* mRNA  
771 probe. The accessory gland cell layer is stained dark (arrowhead), in contrast to the  
772 other ovarian tissues, including the spermatika (asterisk). Bar = 0.2 mm. **(e)**  
773 Photographs of *in situ* hybridization for *RS008823* mRNA in the soldier head. The  
774 front of the head is on the left side. The gland cell layer surrounding the frontal

775 gland reservoir (R) is stained dark (arrowhead). The asterisk indicates the brain.  
776 Bar = 0.1 mm. **(f, g, h)** Vertical cryosection of the worker antenna (f) and horizontal  
777 cryosections of the worker labial palp (g, right palp) and maxillary palp (h, the last  
778 segment of the left (upper) and right (lower) palp) subjected to *in situ* hybridization  
779 for *RS008882* mRNA. Tissues around some sensilla are stained dark (arrowhead).  
780 Bar = 0.1 mm. Photographs of cryosections hybridized with sense probes (negative  
781 controls) are shown in Supplementary Fig. 6a-c.

782

783 **Figure 4: Glycoside hydrolase family (GH) 1 in the *R. speratus* genome.**

784 **(a)** ML tree of GH1 genes based on the amino acid sequences obtained with a  
785 LG+G+I model. Fourteen of 16 GH1 genes in *R. speratus* were used; two genes  
786 (*RS004146* and *RS100005*) were removed from the analysis due to incomplete  
787 retrieval of the coding sequences from gapped scaffolds. GH1 subclasses are  
788 colored and labeled A, B, C, and D. **(b)** GH1 multigene clusters in the *R. speratus*  
789 genome and their expression levels. Letters A-D on the gene structures represent  
790 GH1 subclasses categorized in the phylogenetic tree in (a). The heatmap shows  
791 the Z-scores of the  $\log(\text{RPKM}+1)$  values in the caste-specific transcriptome. **(c)**  
792 Synteny comparison around the GH1 multigene cluster region (orange rectangle)  
793 between *R. speratus* and *M. natalensis* genomes. **(d)** Vertical cryosection of the  
794 worker thorax subjected to *in situ* hybridization with an antisense DIG-labeled  
795 *RS004136* mRNA probe. The head part is on the right side. Bar = 0.2 mm. **(e)**  
796 Magnified view of the worker thorax. The salivary gland cells are specifically stained  
797 dark (arrowhead). Bar = 0.1 mm. **(f)** Vertical cryosection of the queen abdomen  
798 subjected to *in situ* hybridization for *RS004624* mRNA. Bar = 0.2 mm. **(g)** Magnified  
799 view of the queen ovary. The accessory gland cell layer is stained dark  
800 (arrowhead), in contrast to the other ovarian tissues, including ovarioles with two  
801 oocytes (asterisks). Bar = 0.1 mm. See Supplementary Fig. 7a-b for negative  
802 controls of the *in situ* hybridization experiments (d-g).

803

804 **Figure 5: Lysozyme family in the *R. speratus* genome.**

805 **(a)** ML tree of lysozyme genes with a GTR+G model. The red curve indicates a  
806 lineage-specific gene expansion observed in the *R. speratus* genome for a c-type  
807 lysozyme. **(b)** Lysozyme multigene clusters in the *R. speratus* genome and their  
808 relative expression levels among castes. The heatmap shows the Z-scores of the  
809  $\log(\text{RPKM}+1)$  values in the caste-specific transcriptome.

810

811 **Figure 6: Geranylgeranyl diphosphate (GGPP) synthase homologs in the *R.***  
812 ***speratus* genome.**

813 **(a)** Synteny comparison around GGPP synthase loci among three termites, *R.*  
814 *speratus*, *M. natalensis* and *Z. nevadensis*. *R. speratus*-specific insertions were  
815 found, where GGPP synthase paralogs were tandemly duplicated in the *R. speratus*  
816 genome. **(b)** Genomic location and gene expression of *R. speratus* GGPPS  
817 homologs. The heatmap shows the expression level calculated by mean-centered  
818  $\log(\text{RPKM}+1)$ . Yellow indicates high expression, while blue denotes low expression.  
819 Black represents the mean level of expression among the castes. Note that the  
820 heatmap of *RS007484* is almost entirely black for all samples, which indicates that  
821 expression was invariable among castes, while most of the rest of the paralogs  
822 showed caste-biased expression. **(c)** ML tree of GGPP synthase homologs with a  
823 LG+G model. *R. speratus* genes are marked with blue circles. **(d)** Vertical  
824 cryosection of the soldier head subjected to *in situ* hybridization for *RS100016*  
825 mRNA. The front of the head is on the left side. The gland cell layer surrounding the  
826 frontal gland reservoir (R) is stained dark (arrowhead). The asterisk indicates the  
827 brain. The frontal pore (P) discharging frontal gland secretion is also observed. Bar  
828 = 0.1 mm. See Supplementary Fig. 7c for the negative control experiment. **(e)**  
829 Molecular evolutionary analysis of *R. speratus* GGPP synthase homologs by the  
830 PAML branch-site test. Detected positive selection is marked with a single asterisk \*  
831 ( $p < 0.05$ ) or double asterisks \*\* ( $p < 0.01$ ) next to the corresponding branches.

832

833 **Figure 7: The TY family, a novel secretion gene family identified from termite**  
834 **taxonomically restricted genes.**

835 **(a)** Genomic locations and caste-biased expression patterns of TY family genes. **(b)**  
836 Multiple alignment of TY homologs of *R. speratus* and *Z. nevadensis*. Protein motifs  
837 and structural characteristics are represented. **(c)** Orthology of TY homologs in  
838 three termites and the results of the Ka/Ks analysis.

839

840

## 841 References

- 842 1. J. M. Smith, E. Szathmary, *The Major Transitions in Evolution* (OUP Oxford,  
843 1997).
- 844 2. E. O. Wilson, Others, The insect societies. *The insect societies*. (1971).
- 845 3. G. E. Robinson, C. M. Grozinger, C. W. Whitfield, Sociogenomics: social life in  
846 molecular terms. *Nat. Rev. Genet.* **6**, 257–270 (2005).
- 847 4. C. M. Grozinger, Y. Fan, S. E. R. Hoover, M. L. Winston, Genome-wide  
848 analysis reveals differences in brain gene expression patterns associated with  
849 caste and reproductive status in honey bees (*Apis mellifera*). *Molecular*  
850 *Ecology* **16**, 4837–4848 (2007).
- 851 5. R. Bonasio, *et al.*, Genomic comparison of the ants *Camponotus floridanus*  
852 and *Harpegnathos saltator*. *Science* **329**, 1068–1071 (2010).
- 853 6. P. G. Ferreira, *et al.*, Transcriptome analyses of primitively eusocial wasps  
854 reveal novel insights into the evolution of sociality and the origin of alternative  
855 phenotypes. *Genome Biol.* **14**, R20 (2013).
- 856 7. R. Bonasio, *et al.*, Genome-wide and caste-specific DNA methylomes of the  
857 ants *Camponotus floridanus* and *Harpegnathos saltator*. *Curr. Biol.* **22**, 1755–  
858 1764 (2012).
- 859 8. B. Feldmeyer, D. Elsner, S. Foitzik, Gene expression patterns associated with  
860 caste and reproductive status in ants: worker-specific genes are more derived  
861 than queen-specific ones. *Mol. Ecol.* **23**, 151–161 (2014).
- 862 9. S. Patalano, *et al.*, Molecular signatures of plastic phenotypes in two eusocial  
863 insect species with simple societies. *Proceedings of the National Academy of*  
864 *Sciences* **112**, 201515937–201513975 (2015).
- 865 10. R. Libbrecht, P. R. Oxley, L. Keller, D. J. C. Kronauer, Robust DNA  
866 Methylation in the Clonal Raider Ant Brain. *Curr. Biol.* **26**, 391–395 (2016).
- 867 11. D. S. Standage, *et al.*, Genome, transcriptome and methylome sequencing of  
868 a primitively eusocial wasp reveal a greatly reduced DNA methylation system  
869 in a social insect. *Mol. Ecol.* **25**, 1769–1784 (2016).
- 870 12. D. F. Simola, *et al.*, Epigenetic (re)programming of caste-specific behavior in  
871 the ant *Camponotus floridanus*. *Science* **351**, aac6633 (2016).
- 872 13. K. M. Glastad, B. G. Hunt, M. A. D. Goodisman, DNA methylation and  
873 chromatin organization in insects: insights from the Ant *Camponotus*  
874 *floridanus*. *Genome Biol. Evol.* **7**, 931–942 (2015).
- 875 14. S. M. Rehan, A. L. Toth, Climbing the social ladder: the molecular evolution of  
876 sociality. *Trends Ecol. Evol.* **30**, 426–433 (2015).

- 877 15. A. L. Toth, S. M. Rehan, Molecular Evolution of Insect Sociality: An Eco-Evo-  
878 Devo Perspective. *Annu. Rev. Entomol.* **62**, 419–442 (2017).
- 879 16. Y. Roisin, J. Korb, “Social Organisation and the Status of Workers in Termites”  
880 in *Biology of Termites: A Modern Synthesis*, D. E. Bignell, Y. Roisin, N. Lo,  
881 Eds. (Springer Netherlands, 2011), pp. 133–164.
- 882 17. N. Terrapon, *et al.*, Molecular traces of alternative social organization in a  
883 termite genome. *Nat. Commun.* **5**, 3636 (2014).
- 884 18. M. C. Harrison, *et al.*, Hemimetabolous genomes reveal molecular basis of  
885 termite eusociality. *Nat Ecol Evol* **2**, 557–566 (2018).
- 886 19. K. M. Glastad, K. Gokhale, J. Liebig, M. A. D. Goodisman, The caste- and sex-  
887 specific DNA methylome of the termite *Zootermopsis nevadensis*. *Sci. Rep.* **6**,  
888 37110 (2016).
- 889 20. K. Krishna, D. A. Grimaldi, V. Krishna, M. S. Engel, Treatise on the Isoptera of  
890 the world.(Bulletin of the American Museum of Natural History, no. 377)  
891 (2013).
- 892 21. D. J. G. Inward, A. P. Vogler, P. Eggleton, A comprehensive phylogenetic  
893 analysis of termites (Isoptera) illuminates key aspects of their evolutionary  
894 biology. *Mol. Phylogenet. Evol.* **44**, 953–967 (2007).
- 895 22. T. Bourguignon, *et al.*, The evolutionary history of termites as inferred from 66  
896 mitochondrial genomes. *Mol. Biol. Evol.* **32**, 406–421 (2014).
- 897 23. E. L. Vargo, C. Husseneder, Biology of subterranean termites: insights from  
898 molecular studies of Reticulitermes and Coptotermes. *Annu. Rev. Entomol.* **54**,  
899 379–403 (2009).
- 900 24. HARRIS, WV, Termites of the Palearctic Region. *Biology of Termites*, 295–313  
901 (1970).
- 902 25. WEESNER, F. M, Termites of the nearctic region. *The Biology of Termites*,  
903 477–522 (1970).
- 904 26. S. Govorushko, Economic and ecological importance of termites: A global  
905 review. *Entomol. Sci.* **22**, 21–35 (2019).
- 906 27. K. Matsuura, *et al.*, Queen Succession Through Asexual Reproduction in  
907 Termites. *Science* **323**, 1687–1687 (2009).
- 908 28. S. Koshikawa, S. Miyazaki, R. Cornette, T. Matsumoto, T. Miura, Genome size  
909 of termites (Insecta, Dictyoptera, Isoptera) and wood roaches (Insecta,  
910 Dictyoptera, Cryptocercidae). *Naturwissenschaften* **95**, 859–867 (2008).
- 911 29. M. Seppey, M. Manni, E. M. Zdobnov, BUSCO: Assessing Genome Assembly  
912 and Annotation Completeness. *Methods Mol. Biol.* **1962**, 227–245 (2019).
- 913 30. M. Poulsen, *et al.*, Complementary symbiont contributions to plant

- 914 decomposition in a fungus-farming termite. *Proc. Natl. Acad. Sci. U. S. A.* **111**,  
915 14500–14505 (2014).
- 916 31. T. Miura, Developmental regulation of caste-specific characters in social-insect  
917 polyphenism. *Evol. Dev.* **7**, 122–129 (2005).
- 918 32. H. F. Nijhout, Development and evolution of adaptive polyphenisms. *Evol. Dev.*  
919 **5**, 9–18 (2003).
- 920 33. T. Wu, G. K. Dhimi, G. J. Thompson, Soldier-biased gene expression in a  
921 subterranean termite implies functional specialization of the defensive caste.  
922 *Evol. Dev.* **20**, 3–16 (2018).
- 923 34. M. Hojo, S. Shigenobu, K. Maekawa, T. Miura, G. Tokuda, Duplication and  
924 soldier-specific expression of geranylgeranyl diphosphate synthase genes in a  
925 nasute termite *Nasutitermes takasagoensis*. *Insect Biochem. Mol. Biol.* **111**,  
926 103177 (2019).
- 927 35. M. E. Scharf, D. Wu-Scharf, B. R. Pittendrigh, G. W. Bennett, Caste- and  
928 development-associated gene expression in a lower termite. *Genome Biol.* **4**,  
929 R62 (2003).
- 930 36. M. M. Steller, S. Kambhampati, D. Caragea, Comparative analysis of  
931 expressed sequence tags from three castes and two life stages of the termite  
932 *Reticulitermes flavipes*. *BMC Genomics* **11**, 463 (2010).
- 933 37. T. Weil, M. Rehli, J. Korb, Molecular basis for the reproductive division of  
934 labour in a lower termite. *BMC Genomics* **8**, 198 (2007).
- 935 38. S. Ohno, *Evolution by Gene Duplication* (Springer, Berlin, Heidelberg, 1970).
- 936 39. H. Innan, F. Kondrashov, The evolution of gene duplications: classifying and  
937 distinguishing between models. *Nat. Rev. Genet.* **11**, 97–108 (2010).
- 938 40. D. F. Simola, *et al.*, A chromatin link to caste identity in the carpenter ant  
939 *Camponotus floridanus*. *Genome Res.* **23**, 486–496 (2013).
- 940 41. K. M. Kapheim, *et al.*, Social evolution. Genomic signatures of evolutionary  
941 transitions from solitary to group living. *Science* **348**, 1139–1143 (2015).
- 942 42. D. R. Flower, The lipocalin protein family: structure and function. *Biochem. J*  
943 **318 ( Pt 1)**, 1–14 (1996).
- 944 43. T. Miura, *et al.*, Soldier caste-specific gene expression in the mandibular  
945 glands of *Hodotermopsis japonica* (Isoptera: termopsidae). *Proc. Natl. Acad.*  
946 *Sci. U. S. A.* **96**, 13874–13879 (1999).
- 947 44. H. Yaguchi, *et al.*, A lipocalin protein, Neural Lazarillo, is key to social  
948 interactions that promote termite soldier differentiation. *Proc. Biol. Sci.* **285**  
949 (2018).
- 950 45. M. Ruiz, D. Sanchez, C. Correnti, R. K. Strong, M. D. Ganfornina, Lipid-binding

- 951 properties of human ApoD and Lazarillo-related lipocalins: functional  
952 implications for cell differentiation. *FEBS J.* **280**, 3928–3943 (2013).
- 953 46. Y. Mitaka, *et al.*, Caste-Specific and Sex-Specific Expression of  
954 Chemoreceptor Genes in a Termite. *PLoS One* **11**, e0146125 (2016).
- 955 47. D. Watanabe, H. Gotoh, T. Miura, K. Maekawa, Social interactions affecting  
956 caste development through physiological actions in termites. *Front. Physiol.* **5**,  
957 127 (2014).
- 958 48. K. Matsuura, T. Tamura, N. Kobayashi, T. Yashiro, S. Tatsumi, The  
959 antibacterial protein lysozyme identified as the termite egg recognition  
960 pheromone. *PLoS One* **2**, e813 (2007).
- 961 49. P. Pelosi, J.-J. Zhou, L. P. Ban, M. Calvello, Soluble proteins in insect  
962 chemical communication. *Cell. Mol. Life Sci.* **63**, 1658–1676 (2006).
- 963 50. H. Watanabe, G. Tokuda, Cellulolytic Systems in Insects. *Annu. Rev. Entomol.*  
964 **55**, 609–632 (2009).
- 965 51. G. Tokuda, Plant cell wall degradation in insects: Recent progress on  
966 endogenous enzymes revealed by multi-omics technologies. *Adv. In Insect*  
967 *Phys.* **57**, 97 (2019).
- 968 52. J. Ni, G. Tokuda, Lignocellulose-degrading enzymes from termites and their  
969 symbiotic microbiota. *Biotechnol. Adv.* **31**, 838–850 (2013).
- 970 53. K. Matsuura, T. Yashiro, K. Shimizu, S. Tatsumi, T. Tamura, Cuckoo Fungus  
971 Mimics Termite Eggs by Producing the Cellulose-Digesting Enzyme  $\beta$ -  
972 Glucosidase. *Curr. Biol.* **19**, 30–36 (2009).
- 973 54. J. Korb, T. Weil, K. Hoffmann, K. R. Foster, M. Rehli, A gene necessary for  
974 reproductive suppression in termites. *Science* **324**, 758 (2009).
- 975 55. M. Shelomi, B. Wipfler, X. Zhou, Y. Pauchet, Multifunctional cellulase enzymes  
976 are ancestral in Polyneoptera. *Insect Mol. Biol.* **29**, 124–135 (2020).
- 977 56. Q. Gao, G. J. Thompson, Social context affects immune gene expression in a  
978 subterranean termite. *Insectes Soc.* **62**, 167–170 (2015).
- 979 57. L. Callewaert, C. W. Michiels, Lysozymes in the animal kingdom. *J. Biosci.* **35**,  
980 127–160 (2010).
- 981 58. A. Fujita, I. Shimizu, T. Abe, Distribution of lysozyme and protease, and amino  
982 acid concentration in the guts of a wood-feeding termite, *Reticulitermes*  
983 *speratus* (Kolbe): possible digestion of symbiont bacteria transferred by  
984 trophallaxis. *Physiological Entomology* **26**, 116–123 (2001).
- 985 59. WEST, C. A, Biosynthesis of diterpenes. *Biosynthesis of Isoprenoid*  
986 *Compounds* **1**, 375–412 (1981).
- 987 60. K. Ogura, T. Koyama, Enzymatic Aspects of Isoprenoid Chain Elongation.



- 988            *Chem. Rev.* **98**, 1263–1276 (1998).
- 989    61. K. C. Wang, S. Ohnuma, Isoprenyl diphosphate synthases. *Biochim. Biophys.*  
990            *Acta* **1529**, 33–48 (2000).
- 991    62. M. Hojo, T. Matsumoto, T. Miura, Cloning and expression of a geranylgeranyl  
992            diphosphate synthase gene: insights into the synthesis of termite defence  
993            secretion. *Insect Mol. Biol.* **16**, 121–131 (2007).
- 994    63. G. D. Prestwich, The chemicals of termite societies (Isoptera). *Sociobiology*  
995            **14**, 175–191 (1988).
- 996    64. L. J. Nelson, L. G. Cool, B. T. Forschler, M. I. Haverty, Correspondence of  
997            Soldier Defense Secretion Mixtures with Cuticular Hydrocarbon Phenotypes  
998            for Chemotaxonomy of the Termite Genus *Reticulitermes* in North America. *J.*  
999            *Chem. Ecol.* **27**, 1449–1479 (2001).
- 1000    65. A. Quintana, *et al.*, Interspecific variation in terpenoid composition of defensive  
1001            secretions of European *Reticulitermes* termites. *J. Chem. Ecol.* **29**, 639–652  
1002            (2003).
- 1003    66. M. Hojo, K. Toga, D. Watanabe, T. Yamamoto, K. Maekawa, High-level  
1004            expression of the Geranylgeranyl diphosphate synthase gene in the frontal  
1005            gland of soldiers in *Reticulitermes speratus* (Isoptera: Rhinotermitidae). *Arch.*  
1006            *Insect Biochem. Physiol.* **77**, 17–31 (2011).
- 1007    67. Z. Yang, R. Nielsen, Codon-substitution models for detecting molecular  
1008            adaptation at individual sites along specific lineages. *Mol. Biol. Evol.* **19**, 908–  
1009            917 (2002).
- 1010    68. B. R. Johnson, N. D. Tsutsui, Taxonomically restricted genes are associated  
1011            with the evolution of sociality in the honey bee. *BMC Genomics* **12**, 164  
1012            (2011).
- 1013    69. S. Sumner, The importance of genomic novelty in social evolution. *Mol. Ecol.*  
1014            **23**, 26–28 (2014).
- 1015    70. C. R. Smith, *et al.*, Genetic and genomic analyses of the division of labour in  
1016            insect societies. *Nat. Rev. Genet.* **9**, 735–748 (2008).
- 1017    71. M. Lynch, J. S. Conery, The evolutionary fate and consequences of duplicate  
1018            genes. *Science* **290**, 1151–1155 (2000).
- 1019    72. M. Long, E. Betrán, K. Thornton, W. Wang, The origin of new genes: glimpses  
1020            from the young and old. *Nat. Rev. Genet.* **4**, 865–875 (2003).
- 1021    73. G. C. Conant, K. H. Wolfe, Probabilistic cross-species inference of orthologous  
1022            genomic regions created by whole-genome duplication in yeast. *Genetics* **179**,  
1023            1681–1692 (2008).
- 1024    74. R. Gadagkar, The evolution of caste polymorphism in social insects: Genetic  
1025            release followed by diversifying evolution. *Journal of Genetics* **76**, 167–179

- 1026 (1997).
- 1027 75. S. K. McKenzie, I. Fetter-Pruneda, V. Ruta, D. J. C. Kronauer, Transcriptomics  
1028 and neuroanatomy of the clonal raider ant implicate an expanded clade of  
1029 odorant receptors in chemical communication. *Proceedings of the National  
1030 Academy of Sciences* **113**, 14091–14096 (2016).
- 1031 76. X. Zhou, *et al.*, Phylogenetic and transcriptomic analysis of chemosensory  
1032 receptors in a pair of divergent ant species reveals sex-specific signatures of  
1033 odor coding. *PLoS Genet.* **8**, e1002930 (2012).
- 1034 77. X. Zhou, *et al.*, Chemoreceptor Evolution in Hymenoptera and Its Implications  
1035 for the Evolution of Eusociality. *Genome Biol. Evol.* **7**, 2407–2416 (2015).
- 1036 78. L. M. Chau, M. A. D. Goodisman, Gene duplication and the evolution of  
1037 phenotypic diversity in insect societies. *Evolution* **71**, 2871–2884 (2017).
- 1038 79. K. Maekawa, K. Ishitani, H. Gotoh, R. Cornette, T. Miura, Juvenile Hormone  
1039 titre and vitellogenin gene expression related to ovarian development in  
1040 primary reproductives compared with nymphs and nymphoid reproductives of  
1041 the termite *Reticulitermes speratus*. *Physiological Entomology* **35**, 52–58  
1042 (2010).
- 1043 80. A. Bucek, *et al.*, Evolution of Termite Symbiosis Informed by Transcriptome-  
1044 Based Phylogenies. *Curr. Biol.* **29**, 3728–3734.e4 (2019).

Fig. 1

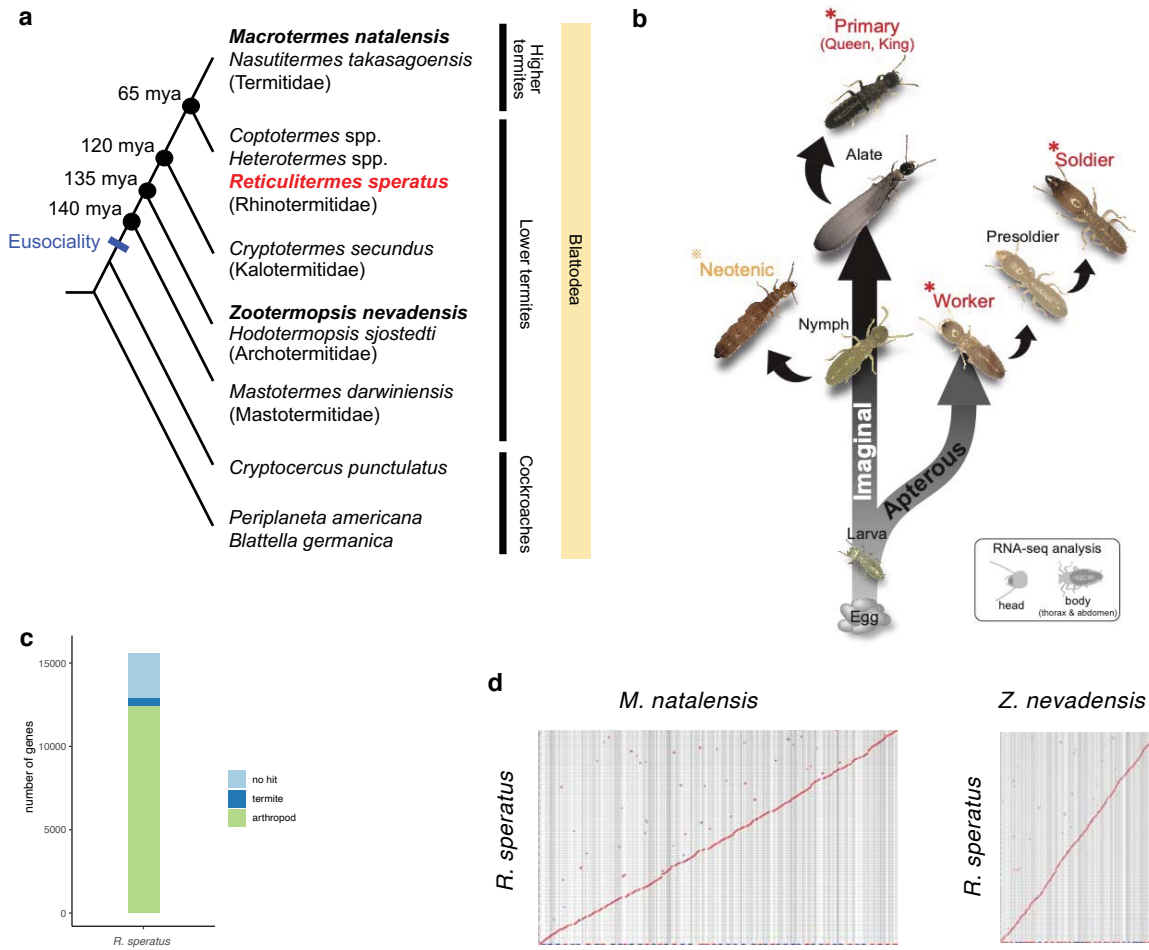


Fig. 2

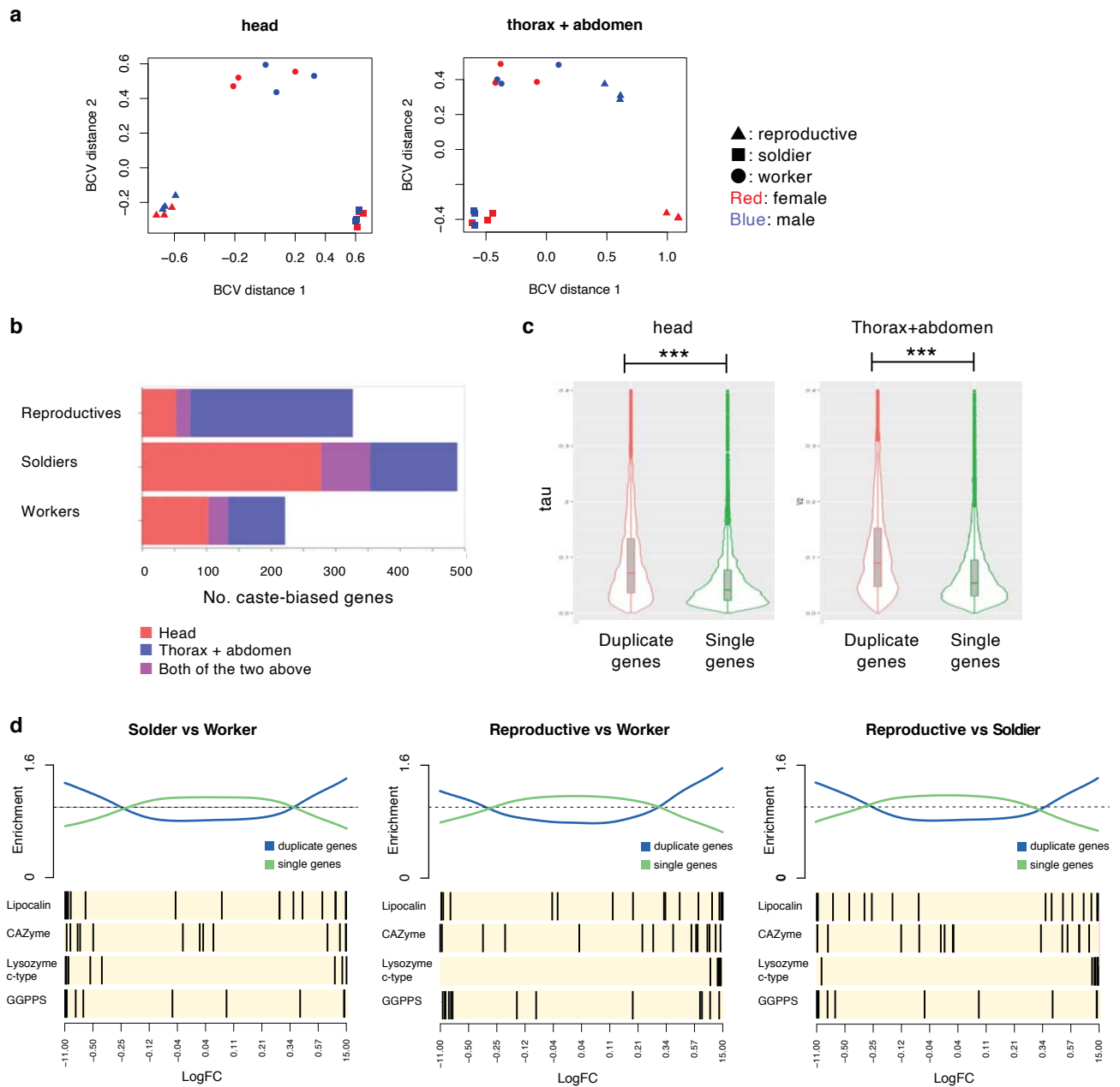


Fig. 3

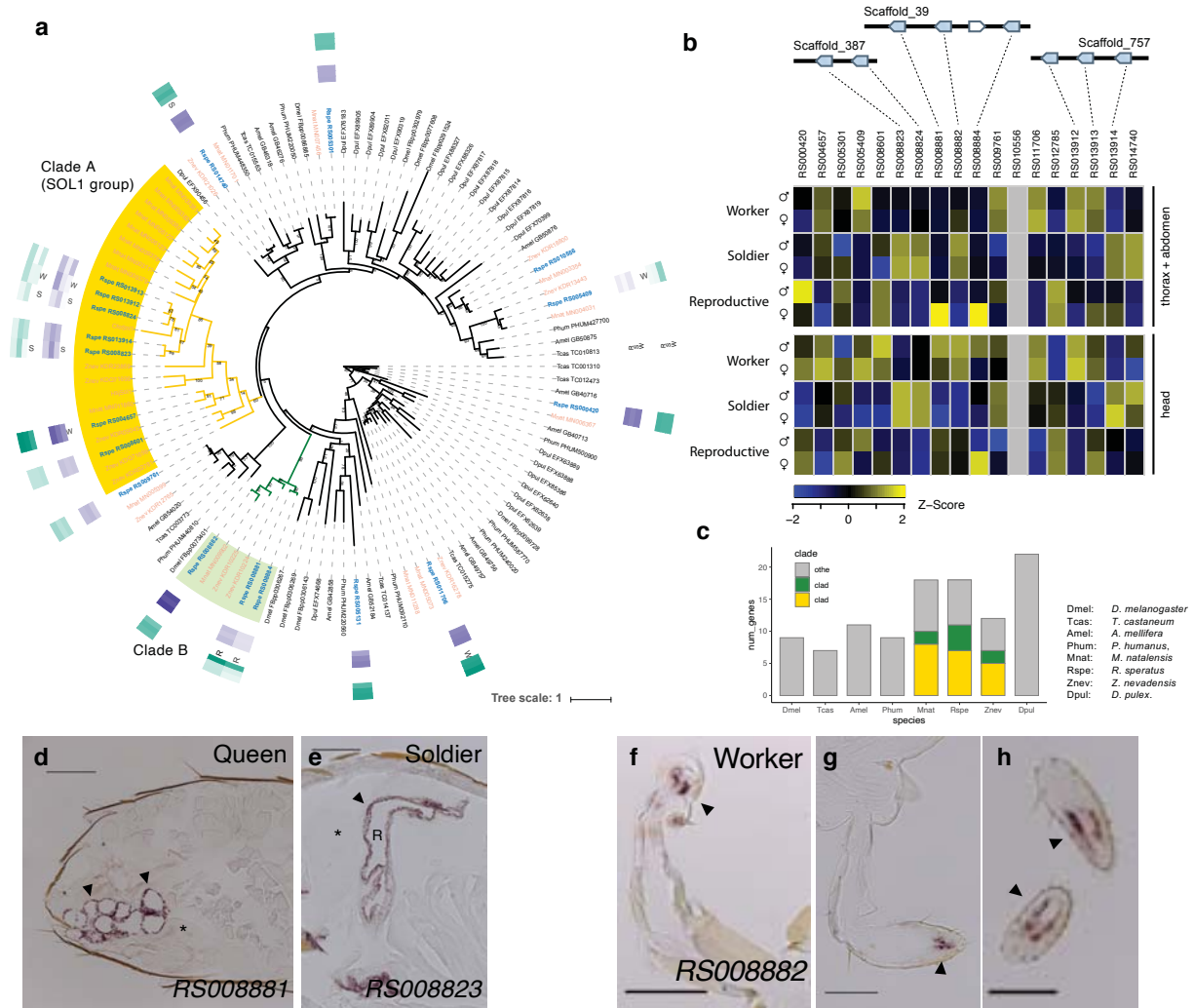


Fig. 4

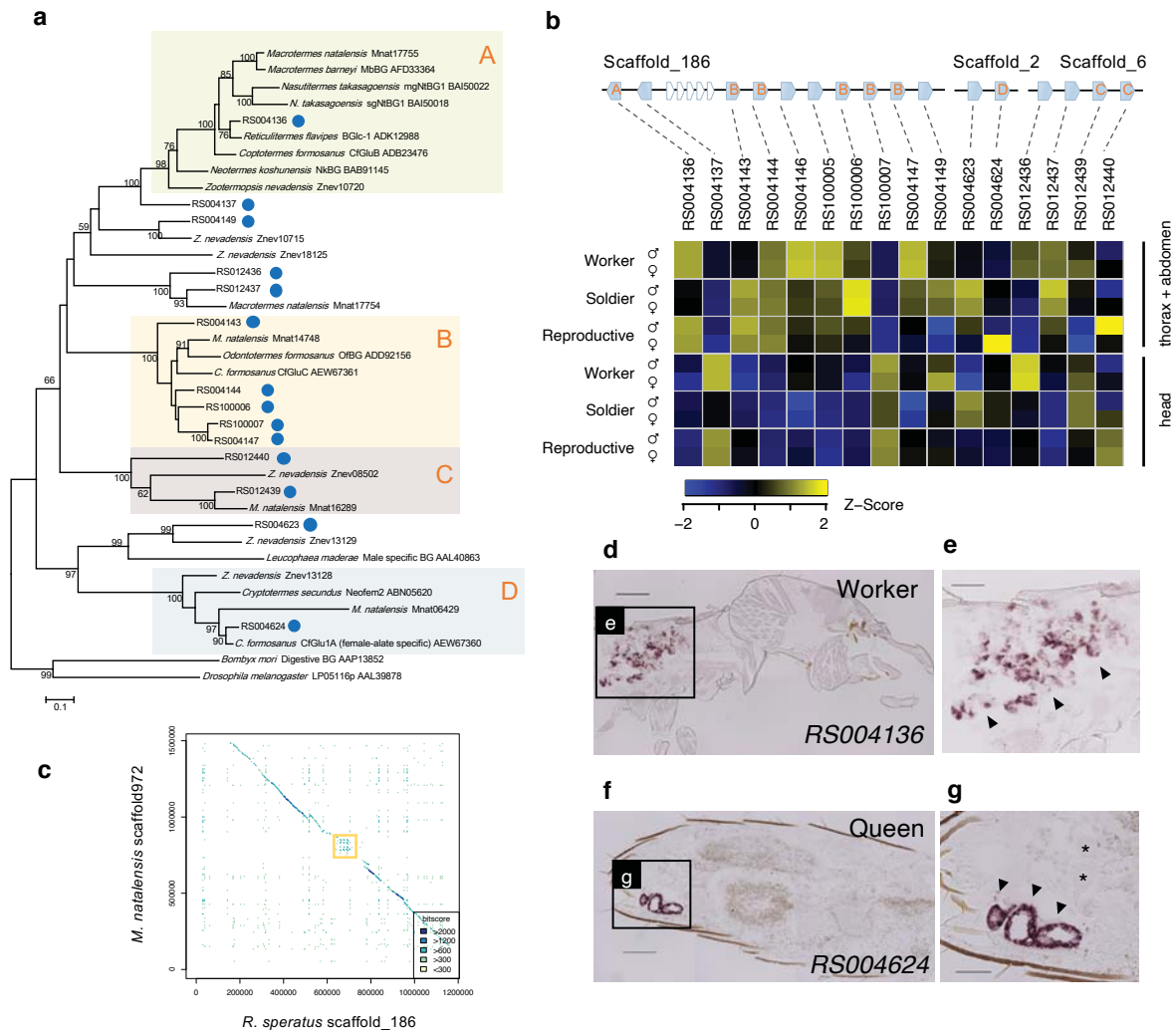


Fig. 5

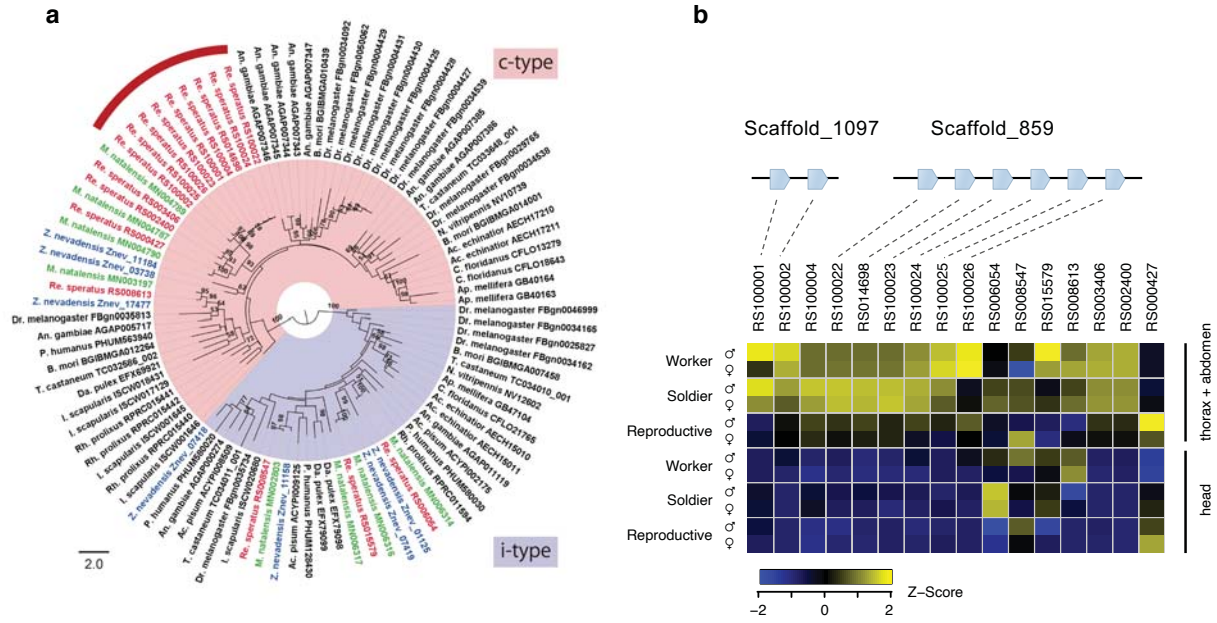


Fig. 6

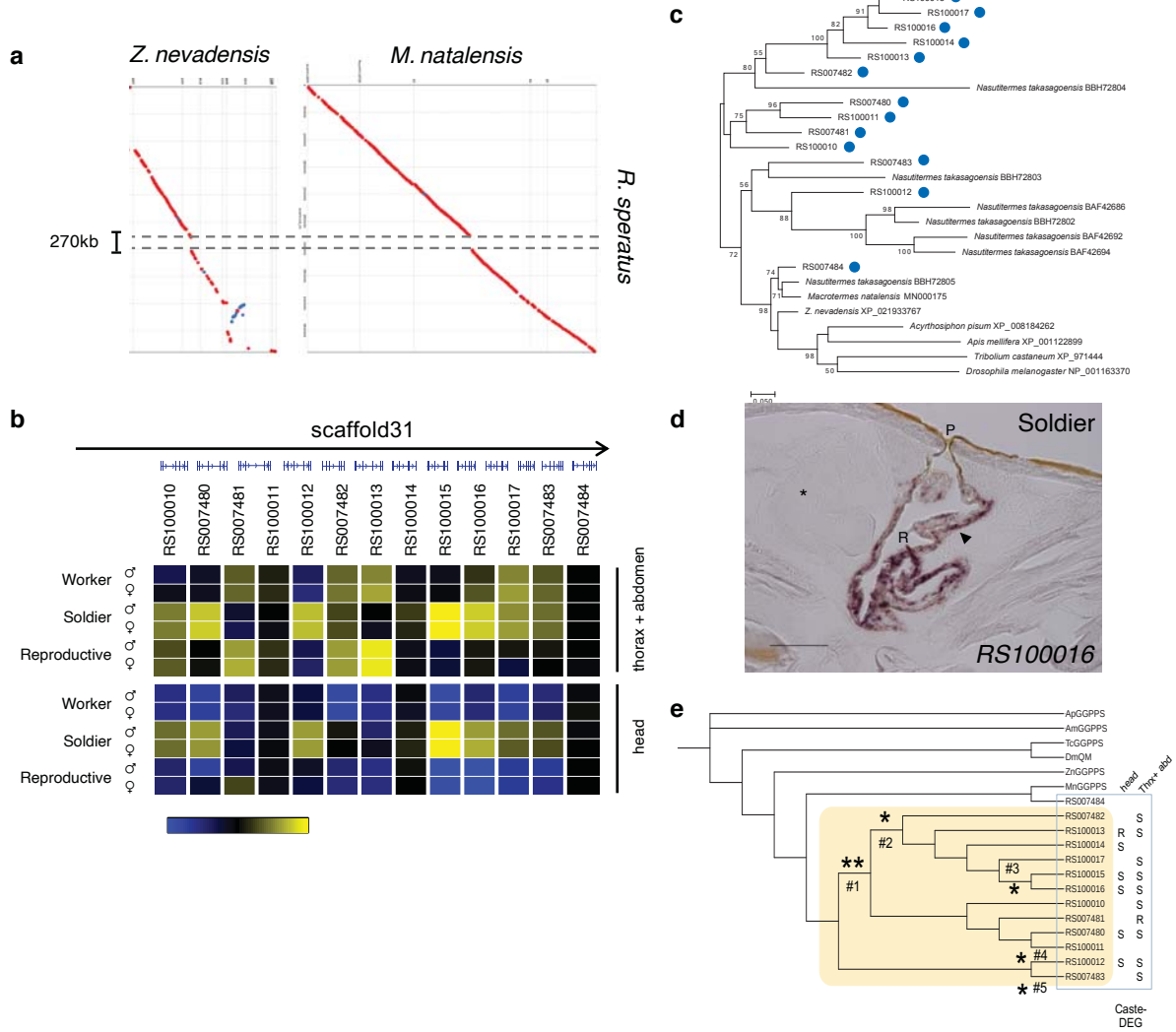




Fig. 7

