

1 Genome editing excisase origins illuminated by 2 somatic genome of *Blepharisma*

3 Minakshi Singh¹, Kwee Boon Brandon Seah¹, Christiane Emmerich¹, Aditi Singh¹, Christian
4 Woehle², Bruno Huettel², Adam Byerly³, Naomi Alexandra Stover⁴, Mayumi Sugiura⁵, Terue
5 Harumoto⁵, Estienne Carl Swart^{1,*}

- 6
- 7 1. Max Planck Institute for Biology, Tuebingen, Germany
 - 8 2. Max Planck Genome Center Cologne, Max Planck Institute for Plant Breeding, Cologne,
9 Germany
 - 10 3. Department of Computer Science and Information Systems, Bradley University, Peoria
11 IL, USA
 - 12 4. Department of Biology, Bradley University, Peoria, IL, USA
 - 13 5. Nara Women's University, Nara, Japan
 - 14 * Correspondence: estienne.swart@tuebingen.mpg.de

15 Summary

16 Massive DNA excision occurs regularly in ciliates, ubiquitous microbial eukaryotes with somatic
17 and germline nuclei in the same cell. Tens of thousands of internally eliminated sequences
18 (IESs) scattered throughout a copy of the ciliate germline genome are deleted during
19 development of the streamlined somatic genome. *Blepharisma* represents one of the two
20 earliest diverging ciliate classes, and, unusually, has dual pathways of somatic nuclear
21 development, making it ideal for investigating the functioning and evolution of these processes.
22 Here, we report the somatic genome assembly of *Blepharisma stoltei* strain ATCC 30299 (41
23 Mb), arranged as numerous alternative telomere-capped minichromosomes. This genome
24 encodes eight PiggyBac transposase homologs liberated from transposons. All are subject to
25 purifying selection, but just one, the putative IES excisase, has a complete catalytic triad. We
26 propose PiggyBac homologs were ancestral excisases that enabled evolution of extensive,
27 natural genome editing.

28
29

30 Keywords

31 Natural genome editing; transposase; transposon; somatic genome; minichromosome;
32 PiggyBac; PiggyMac; PiggyMic; DNA; sRNA.

33 Abbreviations

- 34 • MIC - micronucleus
- 35 • MAC - macronucleus
- 36 • IES - internally eliminated sequence
- 37 • MDS - macronuclear-destined sequence
- 38 • PacBio - Pacific Biosciences
- 39 • CLR - continuous long read (PacBio)
- 40 • CCS - circular consensus sequence (PacBio)
- 41 • HiFi - High-fidelity read (PacBio)
- 42 • ATAS - alternative telomere addition site
- 43 • PBLE - PiggyBac-like element
- 44 • PGBD - PiggyBac element-derived
- 45 • Pgm - PiggyMac
- 46 • PgmL - PiggyMac-like

47 Introduction

48 DNA excision in ciliates is a spectacular and widespread form of natural genome editing with
49 profound consequences for what germline and somatic genomes mean (Arnaiz et al., 2012;
50 Chen et al., 2014; Hamilton et al., 2016; Swart and Nowacki, 2015). Though the responsible
51 processes are under active study, much remains to be learnt from these master DNA
52 manipulators, including how and why this remarkable situation arose in them.

53
54 Knowledge of ciliate genome editing mechanisms is dominated by *Tetrahymena* and
55 *Paramecium* (class Oligohymenophorea), with additional input from *Oxytricha*, *Stylonychia* and
56 *Euplotes* (class Spirotrichea) (Chalker et al., 2013; Vogt et al., 2013). The remaining nine ciliate
57 classes await detailed characterization. To advance investigation of natural genome editing and
58 tackle questions about its origin we focused on the ciliate species *Blepharisma stoltei*. Together
59 with its sister-class, Karyorelictea, the class Heterotrichea, to which this ciliate species belongs,
60 represent the earliest branching ciliate lineages, more distantly related to current model ciliates
61 than those models are to each other (Lynn, 2010). Furthermore, the genus *Blepharisma* exhibits
62 distinctive alternative somatic nuclear developmental pathways, which have the potential to
63 disentangle genome editing processes from indirect influences of preceding pathways.

64
65 *Blepharisma* is a distinctive genus of single-celled ciliates known for the red, light-sensitive
66 pigment, blepharimin, in their sub-pellicular membranes (Giese, 1973), and unusual
67 nuclear/developmental biology (Figure 1) (Miyake et al., 1991). To date molecular investigations
68 and genomics of ciliates have predominantly focused on oligohymenophoreans and spirotrichs
69 (Figure 2, Table S1). In recent years, publication of a draft genome for the heterotrich ciliate,
70 *Stentor*, has facilitated revival of this genus for investigations of cellular regeneration
71 (Slabodnick and Marshall, 2014; Slabodnick et al., 2017; Zhang et al., 2021). However,
72 significant hurdles still need to be overcome to investigate genome editing in *Stentor coeruleus*
73 since requisite cell mating has not been observed in the reference somatic genome strain
74 (personal communication, Mark Slabodnick), and very high lethality has been reported for other
75 strains in which mating occurred (Rapport et al., 1976). We therefore focused on *Blepharisma*
76 which is amenable to such investigations, with controlled induction of mating, and, critically,
77 established procedures for investigating cellular and nuclear development from more than a
78 century of meticulous cytology (Friedl et al., 1983; Giese, 1973; Harumoto et al., 1998; Inaba,
79 1965; Kobayashi et al., 2015; Kumazawa, 1979; Miyake and Harumoto, 1990; Miyake et al.,

80 1979; Repak, 1968; Salvini et al., 1983; Sugiura et al., 2010; Terazima and Harumoto, 2004;
81 Young, 1937).

82
83 The *Blepharisma stoltei* strains used in the present study were originally isolated in Germany
84 (strain ATCC 30299) and Japan (strain HT-IV), with the former continuously cultured for over
85 fifty years, and the latter for over a decade. The cells are comparatively straightforward to
86 maintain, e.g., stable cultures can be established in a simple salt medium on a few grains of
87 rice. Due to their distinctive pigmentation and large size several *Blepharisma* species are
88 excellent subjects for introducing cell biology concepts to non-specialists, and are thus readily
89 available for educational purposes from commercial suppliers. They are ideal subjects for
90 behavioral and developmental investigations, e.g., as voracious predators of smaller ciliates and
91 other unicellular species, and also exhibit pronounced phenotypic plasticity, including forming
92 cysts and giant, cannibal cells under suitable conditions (Giese, 1973).

93
94 Like all ciliates (Prescott, 1994), *Blepharisma* cells have two types of nuclei: a macronucleus
95 (MAC) which is very large and transcriptionally active during vegetative growth, and a small,
96 generally transcriptionally inactive micronucleus (MIC), which serves as the germline (Figure 1A,
97 B). In vegetative propagation (asexual replication) of *Blepharisma*, cell fission results in half of
98 the MAC pinching off before distributing to each of the resulting daughter cells together with the
99 mitotically divided MICs. Upon starvation, *Blepharisma* cells, like other ciliates, are also capable
100 of sexual processes initiated by conjugation. Essential for developmental investigations, the
101 intricate ballet of nuclear movements and morphological changes occurring during *Blepharisma*
102 conjugation is well-documented (Miyake et al., 1991) (Figure 1C). During this process half of the
103 MICs in each of the cells undergo meiosis (meiotic MICs) and the rest do not (somatic MICs)
104 (Figure 1C). One of the meiotic MICs eventually gives rise to two haploid gametic nuclei. One
105 gametic MIC (the migratory nucleus) from each conjugating cell is exchanged with that of its
106 partner. In parallel in partnered cells, subsequent fusion of the migratory and stationary haploid
107 nuclei generates a zygotic nucleus (synkaryon), and after successive mitotic divisions gives rise
108 to both new MICs and new MACs (known as anlagen). The new MACs continue to mature,
109 eventually growing in size and DNA content (Miyake et al., 1991).

110
111 Conveniently for investigations of development and genome editing, *Blepharisma* is one of only
112 two ciliate genera, along with *Euplotes* (Katashima, 1959; Kimball, 1942; Luporini et al., 1983;
113 Vallesi et al., 1995), where conjugation has been shown to be mediated through pheromone-like

114 substances called gamones. *Blepharisma* has two mating types, distinguished by their gamone
115 production. Mating type I cells release gamone 1, a ~30 kDa glycoprotein (Miyake and Beyer,
116 1974; Sugiura and Harumoto, 2001); mating type II cells release gamone 2, calcium-3-(2'-
117 formylamino-5'-hydroxybenzoyl) lactate, a small-molecule effector (Kubota et al., 1973).
118 *Blepharisma* cells commit to conjugation when complementary mating types recognize each
119 other's gamones, with the cells remaining paired while meiosis and fertilization occur and
120 eventually new MACs begin to form.

121
122 As in model ciliates, we show in an accompanying paper that MIC-specific sequences are
123 removed to form a functional *Blepharisma* MAC genome (Seah, et al. 2022). Like other ciliates
124 the resulting MAC genome appears to have been freed of mobile elements and other forms of
125 junk DNA contained in the MIC genome (Klobutcher and Herrick, 1997). However, this situation
126 is an oversimplification of the actual MAC genome content (Seah, et al. 2022). In the best
127 studied ciliates, genome editing is thought to be coordinated or assisted by small RNAs (sRNAs)
128 (Chalker et al., 2013). Specific MIC-limited DNA segments — internally eliminated sequences
129 (IESs) — are excised by domesticated transposases (Arnaiz et al., 2012; Chalker et al., 2013;
130 Klobutcher and Herrick, 1995; Prescott, 1994). Large scale genome-wide DNA amplification
131 accompanies genome editing, producing thousands of copies in mature MACs of larger ciliate
132 species (Klobutcher and Herrick, 1997; Prescott, 1994).

133
134 We were motivated to investigate genome editing in *Blepharisma*, as, unlike model ciliates,
135 these cells can produce two kinds of anlagen, and because one of their two developmental
136 pathways skips the preceding series of mitoses, meioses, nuclear exchanges and fertilization
137 (Miyake et al., 1991) (Figure 1C). Primary anlagen mature in the conventional manner from
138 zygotic nuclei. Somatic MICs which have not undergone meiosis can give rise to secondary
139 anlagen, which can develop into mature macronuclei (Miyake et al., 1991). This occurs
140 frequently in strains with a high selfing frequency (conjugation among cells within a clonal
141 population), in preference to development of primary MAC anlagen (Miyake et al., 1991). This
142 alternative pathway of MAC development has also been observed experimentally after removal
143 of primary MAC anlagen by microsurgery (Miyake et al., 1991). As conjugation progresses, the
144 old (maternal) MACs are progressively degraded (Miyake et al., 1991). Since the *B. stoltei* MIC
145 genome has numerous gene-interrupting IESs (Seah et al. 2022), in principle, editing of DNA
146 needs to occur in both primary and secondary anlagen to produce functional MAC genomes.

147

148 Here we provide essential somatic genome and transcriptomic resources for *B. stoltei*. From
149 long-read sequencing, the *B. stoltei* MAC genome appears to be organized as numerous
150 alternative minichromosomes. Among *Blepharisma*'s MAC-encoded transposase genes we
151 identified were PiggyBac transposase homologs, which, thus far only reported in the distantly
152 related ciliates *Paramecium* and *Tetrahymena*. A few *Blepharisma* PiggyBac homologs are
153 substantially upregulated in MAC development, including the main candidate IES excisase.
154 Consistent with ancient origins of ciliate genome editing, *Blepharisma* shares pronounced
155 development-specific upregulation of homologs known to be involved in this process.
156 *Blepharisma* therefore represents an invaluable outgroup for investigations of genome editing
157 evolution.
158

159 Results

160 A compact somatic genome with a minichromosomal architecture

161 The draft *Blepharisma stoltei* ATCC 30299 MAC genome is compact (41 Mb) and AT rich (66%),
162 like most sequenced ciliate MAC genomes (Figure 2; Table S1, 2, Figure S1A). The genome is
163 gene-dense (25,711 predicted genes), with short intergenic regions, tiny, predominantly 15 and
164 16 bp introns (Figure S4; Supplemental information, "Tiny spliceosomal introns") and
165 untranslated regions (UTRs) (Figure 3A). *B. stoltei* uses an alternative nuclear genetic code with
166 UGA codons reassigned from stops to tryptophan (Figure S1B).

167
168 From joint variant calling of reads from strains ATCC 30299 and HT-IV, strain ATCC 30299
169 appears to be virtually homozygous, with only 1277 heterozygous single-nucleotide
170 polymorphisms (SNPs) compared to 193725 in strain HT-IV (i.e., individual heterozygosity of
171 3.08×10^{-5} vs. 4.67×10^{-3} respectively). Low SNP levels were likely beneficial for overall
172 genomic contiguity, since heterozygosity poses significant algorithmic challenges for assembly
173 software (Chin et al., 2016). For brevity's sake, we refer to this genome as the *Blepharisma*
174 MAC genome (and "*Blepharisma*" for the associated strain). Though the final assembly
175 comprises 64 telomere-to-telomere sequences, chromosomes and their ends are meaningless
176 given the extensive natural fragmentation of the *Blepharisma* MAC genome (characterized in the
177 next section), hence we simply refer to "contigs".

178
179 The basic telomere unit of *Blepharisma* is a permutation of CCCTAACA, like its heterotrich
180 relative *Stentor coeruleus* (Slabodnick et al., 2017) (Figure S2). Since a compelling candidate
181 for a telomerase ncRNA (TERC) could not be found in either *Blepharisma* or *Stentor* using
182 Infernal (Nawrocki et al., 2009) and RFAM models (RF00025 - ciliate TERC; RF00024 -
183 vertebrate TERC), it was not possible to delimit the repeat ends. Heterotrichs may use a
184 different or very divergent ncRNA. In contrast to the extremely short (20 bp) MAC telomeres of
185 spirotrichs like *Oxytricha* with extreme MAC genome fragmentation (Swart et al., 2013),
186 sequenced *Blepharisma* MAC telomeres are moderately long (Figure S2A), with a mode of 209
187 bp (~26 repeats of the 8 bp motif), extending to a few kilobases.

188
189 With a moderately strict definition of possessing at least three consecutive telomeric repeats,
190 one in eight reads in the *Blepharisma* HiFi library were telomere-bearing. Telomeric reads are
191 distributed across the entire genome (Figure 3B). Typically, a minority of mapped reads are
192 telomere-bearing at individual internal positions, and so we term them alternative telomere
193 addition sites (ATASs) (Figure 3B). We identified 46705 potential ATASs, the majority of which
194 (38686) were represented by only one mapped HiFi read.

195
196 The expected distance between telomeres, and hence the average MAC DNA molecule length,
197 is about 130 kb. This is consistent with the raw input MAC DNA lengths, which were mostly
198 longer than 10 kb and as long as 1.5 Mb (Figure S3A, B), and the small fraction (1.3%) of
199 *Blepharisma*'s HiFi reads bound by telomeres on both ends. Excluding the length of the
200 telomeres, telomere-bound reads may be as short as 4 kb (Figure S2B). Given the frequency of
201 telomere-bearing reads, we expect many additional two-telomere DNA molecules longer than 12
202 kb, the approximate maximum length of the HiFi reads (Figure S3A, B).

203
204 Since the lengths of the sequenced two-telomere DNA molecules on average imply that they
205 encode multiple genes, we propose classifying them as “minichromosomes”. This places them
206 between the “nanochromosomes” of ciliates like *Oxytricha* and *Stylonychia*, which typically
207 encode single genes and a few kilobases long (Aeschlimann et al., 2014; Swart et al., 2013),
208 and *Paramecium tetraurelia* and *Tetrahymena thermophila* MAC chromosomes which are
209 hundreds of kilobases to megabases long (Aury et al., 2006; Sheng et al., 2020; Zagulski et al.,
210 2004). The *Paramecium bursaria* MAC genome is considerably more fragmented than those of

211 other previously examined *Paramecium* species, and have thus also been classified as
212 minichromosomes (Cheng et al., 2020).

213 Key features of gene expression during new MAC development

214 To gain an overview of the molecular processes during *Blepharisma* genome editing, we
215 examined gene expression trends across development. Complementary *B. stoltei* strains were
216 treated with gamones of the opposite mating type, before mixing to initiate conjugation (Miyake
217 et al., 1991; Sugiura et al., 2012). Samples for morphological staging and RNA-seq were taken
218 at intervals from the time of mixing ("0 hour" time point) up to 38 hours.

219
220 During *Blepharisma* conjugation, meiosis begins around 2 h after conjugating cell pairs form and
221 continues up to 18 h, by when gametic nuclei generated by meiosis have been exchanged
222 (Figure 1C; Figure 4). This is followed by karyogamy and mitotic multiplication of the zygotic
223 nucleus (22 hours). At 26 h, new, developing primary MACs can be observed in the conjugating
224 pairs, as large, irregular bodies (Figure 4). These nuclei mature into the new MACs of the
225 exconjugant cell by 38 h, after which cell division generates two daughter cells. Smaller
226 secondary MACs, derived directly from MICs without all the intermediate nuclear stages, can
227 also be seen from 22 h, eventually disappearing, giving way to the primary MACs (Figure 4).

228
229 Examining gene expression at 26 h, when the majority of cells are forming a new MAC (Figure
230 4), we observe two broad trends: relatively stable constitutive gene expression (Table S5; Data
231 S3), e.g., an actin homolog (ENA accession: BSTOLATCC_MAC19444) and a bacteria-like
232 globin protein (BSTOLATCC_MAC21846), versus pronounced development-specific
233 upregulation (Table S6; Data S3), e.g., a histone (BSTOLATCC_MAC21995) an HMG box
234 protein (BSTOLATCC_MAC14030), and a translation initiation factor (eIF4E,
235 BSTOLATCC_MAC5291). We eschewed a shallow Gene Ontology (GO) enrichment analysis,
236 instead favoring close scrutiny of a smaller subset of genes strongly upregulated during new
237 MAC formation. For this, computational gene annotations in combination with BLASTP searches
238 and examination of literature associated with homologs was used. Ranking the relative gene
239 expression at 26 h vs. the average expression of starved, gamone treated, and 0 h cells, in
240 descending order, revealed numerous genes of interest, including homologs of proteins known
241 to be involved in genome editing in model ciliates (Table S6).

242

243 Among the top 100 genes ranked this way (69x to 825x upregulation) nine contain transposase
244 domains from PFAM: DDE_Tnp_1_7, DDE_3, MULE and DDE_Tnp_IS1595 (e.g.,
245 BSTOLATCC_MAC2188, BSTOLATCC_MAC14490, BSTOLATCC_MAC18054,
246 BSTOLATCC_MAC18052, respectively). We also observe small RNA (sRNA) biogenesis and
247 transport proteins, i.e., a Piwi protein (BSTOLATCC_MAC5406) and a Dicer-like protein
248 (BSTOLATCC_MAC1138; “Supplemental information”, “Homologs of small RNA-related proteins
249 involved in ciliate genome editing” and Figure S8), and a POT1 telomere-binding protein
250 homolog (POT1.4; BSTOLATCC_MAC1496; Supplemental information “Telomere-binding
251 protein paralogs”). Numerous homologs of genes involved in DNA repair and chromatin are also
252 present among these highly developmentally upregulated genes (“Supplemental information”,
253 “Development-specific upregulation of proteins associated with DNA repair and chromatin” and
254 “Development-specific histone variant upregulation”). The presence of proteins involved in either
255 transcription initiation or translation initiation among these highly upregulated genes suggests a
256 possible manner in which regulation of development-specific gene expression may be
257 coordinated (“Supplemental information”, “Development-specific upregulation of proteins
258 associated with initiation of transcription and translation”).

259 **A single *Blepharisma* PiggyBac homolog has a complete catalytic** 260 **triad**

261 In *Paramecium tetraurelia* and *Tetrahymena thermophila*, PiggyBac transposases are
262 responsible for IES excision during genome editing (Baudry et al., 2009; Cheng et al., 2010).
263 These transposases appear to have been domesticated, i.e., their genes are no longer
264 contained in transposons but are encoded in the somatic genome where they play an essential
265 genome development role (Baudry et al., 2009; Cheng et al., 2010). PiggyBac homologs
266 typically have a DDD catalytic triad, rather than the more common DDE triad of other DDE/D
267 transposases (Yuan and Wessler, 2011). The DDD catalytic motif is present in *Paramecium*
268 PiggyMac (Pgm) and *Tetrahymena* PiggyBac homologs Tpb1 and Tpb2 (Bischerour et al., 2018;
269 Cheng et al., 2010). Among ciliates, domesticated PiggyBac transposases have so far only been
270 reported in these model oligohymenophorean genera. Notably they have not been detected in
271 either the MAC or MIC genome of the spirotrich *Oxytricha trifallax* (Chen et al., 2014; Swart et
272 al., 2013).
273

274 We detected more transposase domains (9 distinct PFAM identifiers) in *Blepharisma* than any
275 other ciliate species we examined (Figure 5A). Using HMMER searches with the domain
276 characteristic of PiggyBac homologs, DDE_Tnp_1_7 (PF13843), we found eight homologs in *B.*
277 *stoltei* ATCC MAC genome and five additional ones within IESs, none of which were flanked by
278 terminal repeats (identified by RepeatModeler). We also found PiggyBac homologs in the MAC
279 genomes of *B. stoltei* HT-IV and *B. japonicum* R1072.

280
281 Reminiscent of *Paramecium tetraurelia*, which, among ten PiggyMac homologs, has just one
282 homolog with a complete catalytic triad (Bischerour et al., 2018), the DDD triad is preserved in
283 just a single *Blepharisma* PiggyBac homolog (Figure 5B; Contig_49.g1063,
284 BSTOLATCC_MAC17466). This gene is strongly upregulated during development from 22 to 38
285 h, when new MACs develop and IES excision is required (Figure 5B). In a multiple sequence
286 alignment the canonical catalytic triad second aspartate of a lower-expressed, MIC-limited
287 PiggyBac is offset by one amino acid (Data S5).

288
289 There are significant similarities in the basic properties of *Blepharisma* and *Paramecium* IESs,
290 detailed in the *Blepharisma* MIC genome report (Seah et al. 2022). Consequently, adopting the
291 *Paramecium* nomenclature, we refer to the primary candidate IES excisase as *Blepharisma*
292 PiggyMac (BPgm) and the other somatic homologs as BPgm-Likes (BPgmLs). By extension, we
293 refer to their close relatives which are germline-limited as PiggyMics (Figure 5B).

294
295 Other than the PFAM DDE_Tnp_1_7 domain, three *Blepharisma* MAC genome-encoded
296 PiggyBac homologs also possess a short, characteristic cysteine-rich domain (CRD) (Figure
297 5C), which is absent from the other BPgmLs and PiggyMics. PiggyBac CRDs have been
298 classified into three different groups and are essential for *Paramecium* IES excision (Guérineau
299 et al., 2021). In *Blepharisma*, the CRD consists of five cysteine residues arranged as CxxC-
300 CxxCxxxxH-Cxxx(Y)H (where C, H, Y and x respectively denote cysteine, histidine, tyrosine and
301 any other residue). Two *Blepharisma* homologs possess this CRD without the penultimate
302 tyrosine residue, while the third contains a tyrosine residue before the final histidine. This -YH
303 feature towards the end of the CxxC-CxxCxxxxH-Cxxx(Y)H CRD is shared by all the PiggyBac
304 homologs we found in *Condyllostoma*, the bat PiggyBac-like element (PBLE) and human
305 PiggyBac element-derived (PGBD) proteins PGBD2 and PGBD3. In contrast, PiggyBac
306 homologs from *Paramecium* and *Tetrahymena* have a CRD with six cysteine residues arranged

307 in the variants of the motif CxxC-CxxC-Cx{2-7}Cx{3,4}H, and group together with human
308 PGDB4 and *Spodoptera frugiperda* PBLE (Figure 5C).

309 PiggyBac transposases are subject to purifying selection and 310 originated early in ciliate evolution

311 Previous experiments involving individual or paired gene knockdowns of most of the ten
312 *Paramecium tetraurelia* PiggyMac(-like) paralogs led to substantial IES retention, even though
313 only one PiggyMac gene (Pgm) has the complete catalytic triad, indicating that all these proteins
314 are functional (Bischerour et al., 2018). To examine functional constraints on *Paramecium*
315 PiggyMac homologs we examined non-synonymous (d_N) to synonymous substitution rates (d_S),
316 i.e. $\omega = d_N/d_S$, for pairwise codon sequence alignments using two closely related *Paramecium*
317 species (*P. tetraurelia* and *P. octaurelia*). All d_N/d_S values for pairwise comparisons of each of
318 the catalytically incomplete *P. tetraurelia* PgmLs versus the complete Pgm, were less than 1,
319 ranging from 0.01 to 0.25 (Table S7). All d_N/d_S values for pairwise comparisons between *P.*
320 *tetraurelia* and *P. octaurelia* PiggyBac orthologs were also substantially less than 1, ranging
321 from 0.02 to 0.11 (Table S8). Since $d_N/d_S = 1$ indicates genes evolving neutrally (Yang and
322 Nielsen, 2000), none of these genes are likely pseudogenes, and all appear subject to similar
323 purifying selection.

324
325 Only one of *Blepharisma*'s eight MAC and five MIC PiggyBac homologs has the complete,
326 characteristic DDD triad necessary for catalysis. In pairwise comparisons of each of the MAC
327 homologs with incomplete/missing triads versus the complete one d_N/d_S ranges from 0.0076 to
328 0.1351 (Table S9). The pairwise non-synonymous to synonymous substitution rates of the
329 PiggyMics in comparison to the BPgm were also much less than 1 (range 0.007 to 0.2),
330 indicating they are also subject to purifying selection.

331
332 We detected PiggyBac homologs in two other heterotrichs, but not the oligohymenophorean
333 *Ichthyophthirius multifiliis* ("Supplemental information"). To determine whether the *Blepharisma*
334 PiggyBac homologs share a common ciliate ancestor with the oligohymenophorean PiggyBacs,
335 or whether they arose from independent acquisitions in major ciliate groups, we created a large
336 phylogeny of PiggyBac homologs representative of putative domesticated transposases from
337 *Blepharisma stoltei* ATCC 30299, *Condyllostoma magnum*, *Paramecium* spp., *Tetrahymena*
338 *thermophila*, as well as PiggyBac-like elements (PBLEs (Bouallègue et al., 2017)) from diverse

339 eukaryotes (Figure 6; Data S1). All the heterotrichous ciliates PiggyBac homologs, ie. BPgm,
340 BPgmLs 1-7 and PiggyMics grouped together with the *Condyllostoma* Pgms. The ciliate Pgms
341 and PgmLs largely cluster as a single clade, with the exception of PiggyMic 5, which appears as
342 a low-support outgroup to opisthokont, archaeplastid and stramenopile PiggyBac-like elements.
343 PiggyMic 5 has the shortest detected DDE_Tnp_1_7 domain (26 a.a.), and appeared poorly
344 aligned relative to the other homologs.

345 *Blepharisma*'s MAC genome encodes additional domesticated 346 transposases

347 Three *Blepharisma* MAC genome-encoded proteins possess PFAM domain DDE_1 (PF03184;
348 Figure 7). The most common domain combinations for this domain, aside from proteins with it
349 alone (5898 sequences; PFAM version 35), are with an N-terminal PFAM domain
350 HTH_Tnp_Tc5 (PF03221) alone (2240 sequences), and both an N-terminal CENP-B_N domain
351 (PF04218) and central HTH_Tnp_Tc5 domain (1255 sequences). The CENP-B_N domain is
352 characteristic of numerous transposases, notably the Tigger and PogoR families (Gao et al.,
353 2020). Though pairwise sequence identity is low amongst the *Blepharisma* DDE_1 proteins
354 (avg. 28.3%) in their multiple sequence alignment, the CENP-B_N domain in one of them
355 appears to align reasonably well to corresponding regions in the two proteins lacking this
356 domain, suggesting it deteriorated beyond the recognition capabilities of HMMER3 and the
357 given PFAM domain model. BLASTp matches for all three proteins in GenBank are annotated
358 either as Jerky or Tigger homologs (Jerky transposases belong to the Tigger transposase family
359 (Gao et al., 2020)). Given that none of the *Blepharisma* MAC DDE_1 domain proteins appears
360 to have a complete catalytic triad, it is unlikely they are involved in transposition or IES excision.
361

362 Six MAC-encoded transposases containing the DDE_3 domain (PF13358) are present in
363 *Blepharisma*, all of which are substantially upregulated in MAC development and five of which
364 possess the complete DDE catalytic triad (Figure 7B). The DDE_3 domain is characteristic of
365 DDE transposases encoded by the Telomere-Bearing Element transposons (TBEs) of *Oxytricha*
366 *trifallax* (Williams et al., 1993; Witherspoon et al., 1997), which, despite being MIC genome-
367 limited, are proposed to be involved in IES excision (Nowacki et al., 2009). DDE_3-containing
368 transposons, called Tec elements, are found in another spirotrichous ciliate, *Euplotes crassus*,
369 but no role in genome editing has been established for these (Jahn et al., 1993). TBEs and Tec
370 elements do not share obvious features, other than both possessing an encoded protein

371 belonging to the IS630-Tc1 transposase (super)-family (Doak et al., 1994). All six *Blepharisma*
372 DDE_3 genes have at least 150x HiFi read coverage, consistent with their presence in *bona fide*
373 MAC DNA.

374
375 As judged by BLASTP searches in which most of the top hundred best matches are classified
376 are “IS630 family” transposases, *Blepharisma* MAC-encoded DDE_3 domain transposases are
377 more closely related to the IS630 transposase family than to *Oxytricha* TBE transposases and
378 *Euplotes* Tec transposases. One of the BLAST top hits is a MIC genome-encoded protein in
379 *Oxytricha trifallax* with a DDE_3 domain which is not a TBE transposase (GenBank accession:
380 KEJ83017.1). IS630 transposases diverge considerably from Tc1-Mariner transposases, and
381 hence are considered an outgroup to them (Dupeyron et al., 2020). IS630-related transposases
382 encoded by Anchois transposons have also been detected in the *Paramecium tetraurelia* MIC
383 genome (Arnaiz et al., 2012). Given that all but one of the *B. stoltei* paralogs appear to possess
384 a complete catalytic triad, there is a possibility that they may be involved in some IES excision.

385
386 Among other ciliates with draft MAC genomes we examined, the IS1595- and MULE
387 transposase-like domains (PFAM PF12762 and PF10551) have so far only been observed in the
388 spirotrichs *Oxytricha* and *Stylonychia* (Aeschlimann et al., 2014; Swart et al., 2013).
389 DDE_Tnp_IS1595 domains are characteristic of the Merlin transposon superfamily and MULE is
390 part of the Mutator transposon superfamily (Yuan and Wessler, 2011). Currently no particular
391 functions have been demonstrated for these proteins in these ciliates, but their genes were
392 substantially upregulated during their development (Chen et al., 2014; Swart et al., 2013). Both
393 transposase-like domains are found in MAC-encoded proteins in *Blepharisma* and their
394 underlying genes are upregulated during MAC development (Figure 7C, Figure S7). Consistent
395 with the notion of transposase domestication, the genes encoding DDE_Tnp_IS1595 and MULE
396 proteins appear to lack flanking transposon terminal inverted repeats. Members of both IS1595
397 and MULE transposases also appear to have complete catalytic triads.

398
399 In addition to cut-and-paste transposases, we detected a family (> 30 copies) of APE-type non-
400 LTR retrotransposase genes encoding proteins with two characteristic domains, i.e., an APE
401 endonuclease domain (PFAM “exo_endo_phos_2”; PF14529) and a reverse transcriptase
402 domain (PFAM “RVT_1”; PF00078) present on adjacent genes. Unlike the conventional
403 transposase-derived genes in *B. stoltei*, the expression of all these genes throughout the
404 conditions we examined is negligible, and some also appear to be truncated pseudogenes (Data

405 S3; workbook “RVT1 + exo_endo_phos_2”). Since it is necessary to understand the relationship
406 of these sequences with respect to IESs, and that they are not due to residual MIC DNA
407 contamination, their analysis is reported in the context of the *Blepharisma stoltei* MIC genome
408 (Seah et al. 2022).

409 Discussion

410 The genus *Blepharisma* represents one of the earliest diverging ciliate lineages, the
411 heterotrichs, forming an outgroup to the best-studied and deeply divergent oligohymenophorean
412 and spirotrich ciliates (Lynn, 2010). *Blepharisma* species thus provide a vantage point to
413 compare unique processes that have accompanied the evolution of nuclear and genomic
414 dimorphism in ciliates, particularly the extensive genomic editing occurring during MAC
415 development. The annotated draft *B. stoltei* ATCC 30299 MAC genome and associated
416 transcriptomic data provide the basis for comparative studies of genome editing.

417 *Blepharisma* PiggyMac is the primary candidate IES excisase

418 A considerable body of evidence implicates PiggyBac homologs in IES excision of the
419 oligohymenophorean ciliates *Tetrahymena* and *Paramecium* (Arnaiz et al., 2012; Baudry et al.,
420 2009; Bischerour et al., 2018; Cheng et al., 2010; Feng et al., 2017). The responsible IES
421 excisases in the less-studied spirotrichs, *Oxytricha*, *Stylonychia* and *Euplotes*, are not as
422 evident. *Oxytricha*'s TBE transposases are considered to be involved in IES excision, but are
423 encoded by full-length germline-limited transposons and are absent from the MAC (Nowacki et
424 al., 2009), unlike the primary, MAC genome-encoded IES excisase (Tpb2) in *Tetrahymena* and
425 the *Paramecium* PiggyMacs and PiggyMac-likes. The pronounced developmental upregulation
426 of numerous additional MAC- and MIC-encoded transposases in *Oxytricha* raises the possibility
427 that transposases other than those of TBEs could also be involved in IES excision (Chen et al.,
428 2014; Swart et al., 2013). Knowledge of IESs in other ciliates is sparse, primarily confined to the
429 phyllopharyngean *Chilodonella uncinata* (Zufall and Katz, 2007; Zufall et al., 2012). As far as we
430 are aware, no specific IES excisases have been proposed for them.

431
432 In current models of IES excision, MIC-limited sequence demarcation by deposition of
433 methylation marks on histones occurs in an sRNA-dependent process (Chalker et al., 2013).
434 These sequences are recognized by domesticated transposases whose excision is supported

435 by additional proteins that somehow recognize these marks (Chalker et al., 2013). Together with
436 MIC sequencing we observed abundant, development-specific sRNA production in *Blepharisma*
437 resembling other model ciliates (Seah et al. 2022). Homologs of proteins implicated in ciliate
438 genome editing were present among the genes most highly differentially upregulated during new
439 MAC development, notably including Dicer-like and Piwi proteins which are candidate genes
440 responsible for development-specific sRNA biogenesis (Figure S8).

441
442 Since the oligohymenophorean PiggyBac homologs are clear IES excisases, we sought and
443 found eight homologs of these genes in the *Blepharisma* MAC genome and five in the IESs.
444 *Blepharisma* is the first ciliate genus aside from *Tetrahymena* and *Paramecium* in which such
445 proteins have been reported, and distantly related to both. Additional searches revealed clear
446 PiggyBac homologs in *Condylostoma magnum*, and a weaker pair of matches in *Stentor*
447 *coeruleus*, suggesting that these are a common feature of heterotrich ciliates. Reminiscent of
448 *Paramecium tetraurelia*, in which just one of the nine PiggyBac homologs, PiggyMac, has a
449 complete DDD catalytic triad (Bischerour et al., 2018), a single *Blepharisma* PiggyBac homolog
450 has a complete canonical DDD catalytic triad, and its gene is highly upregulated during MAC
451 development. As is characteristic of PiggyBac homologs, each of these three proteins also has a
452 C-terminal, cysteine-rich, zinc finger domain. The organization of the heterotrich PiggyBac
453 homolog zinc finger domains is more similar to comparable domains of *Homo sapiens* PGBD2
454 and PGBD3 homologs than the zinc finger domains in *Paramecium* and *Tetrahymena* PiggyBac
455 homologs.

456
457 Since the discovery of multiple PiggyBac homologs (PiggyMac-likes) in *Paramecium* there have
458 been questions about their role. Aside from PiggyMac, all PiggyMac-likes have incomplete
459 catalytic triads, and are thus likely catalytically inactive, but nevertheless their gene knockdowns
460 lead to pronounced IES retention (Bischerour et al., 2018). It has therefore been proposed that
461 the PiggyMac-likes may function as heteromeric multi-subunit complexes in conjunction with
462 PiggyMac during DNA excision (Bischerour et al., 2018). On the other hand, cryo-EM structures
463 available for moth PiggyBac transposase support a model in which these proteins function as a
464 homodimeric complex *in vitro* (Chen et al., 2020). Furthermore, the primary *Tetrahymena*
465 PiggyBac, Tpb2, is able to perform cleavage *in vitro* alone (Cheng et al., 2010). In other
466 eukaryotes, domesticated PiggyBacs without complete catalytic triads are thought to be retained
467 due to co-option of their DNA-binding domains (Sarkar et al., 2003). One possibility for such
468 purely DNA-binding transposase-derived proteins in ciliates could be in competitively regulating

469 (taming) the excision of DNA by the catalytically active transposases. Future experimental
470 analyses of the BPgm and the BPgm-likes could aid in resolving the conundrums and
471 understanding of possible interactions between catalytically active and inactive transposases.

472 *Blepharisma* has additional domesticated transposases whose 473 roles await determination

474 In addition to the PiggyBac homologs, we found MAC genome-encoded transposases with the
475 PFAM domains “DDE_1”, “DDE_3”, “DDE_Tnp_IS1595” and “MULE” in *Blepharisma*. All the
476 genes encoding these proteins lack flanking terminal repeats characteristic of active
477 transposons, suggesting they are further classes of domesticated transposases. In *Blepharisma*
478 and numerous other organisms, the DDE_1 domains co-occur with CENPB domains. Two such
479 proteins represent totally different proposed exaptations in mammals (centromere-binding
480 protein) and fission yeast (regulatory protein) (Casola et al., 2008; Hohmann, 1993; Mojzita and
481 Hohmann, 2006). Given the great evolutionary distances involved, there is no reason to expect
482 that the *Blepharisma* homologs have either function. None of the three proteins with co-
483 occurring DDE_1 and CENPB domains have a complete catalytic triad, making it unlikely that
484 these are active transposases or IES excisases, though all three are noticeably upregulated
485 during MAC development. Six proteins with the PFAM domain DDE_3 are also encoded by
486 *Blepharisma* MAC genes, of which five possess a complete catalytic triad. DDE_3 domains are
487 also characteristic of TBE transposases in *Oxytricha* and Tec transposases in *Euplotes*. All the
488 “DDE_3” protein genes are upregulated during conjugation in *B. stoltei*, peaking during new
489 MAC development. A number of DDE_Tnp_IS1595 and MULE domain-containing proteins have
490 complete catalytic triads and also show pronounced upregulation during *Blepharisma* MAC
491 development.

492
493 All ciliate species have MAC genome-encoded transposase families (Figure 5A). Though
494 upregulation of some of these homologs in model ciliates has been noted (Chen et al., 2014;
495 Swart et al., 2013; Vogt and Mochizuki, 2013), their roles remain to be determined. Aside from
496 the timing of IES excisase expression to coincide with new MAC genome formation, the manner
497 in which the excisases perform excision is also crucial. Upon excision, classical cut-and-paste
498 transposases in eukaryotes typically leave behind additional bases, notably including the target-
499 site duplication arising when they were inserted, forming a “footprint” (van Luenen et al., 1994).
500 PiggyBac homologs are unique in performing precise, “seamless” excision in eukaryotes (Elick

501 et al., 1996), conserving the number of bases at the site of transposon insertion after excision, a
502 property that makes them popular for genetic engineering (Chen et al., 2020). *Tetrahymena*
503 Tpb2 is the one exception among PiggyBac homologs associated with imprecise excision in this
504 eukaryote (Cheng et al., 2010). Since intragenic IESs are abundant in *Blepharisma*, like
505 *Paramecium* and unlike *Tetrahymena*, it is essential that these are excised precisely.

506
507 Though there are clearly numerous additional domesticated transposases with complete
508 catalytic triads and whose genes are substantially upregulated during *Blepharisma*
509 development, whether they are capable of excision, and if this is precise, needs to be
510 established. *Tetrahymena* has distinct domesticated transposases that excise different subsets
511 of IESs, namely those that are predominant, imprecisely excised and intergenic (by Tpb2)
512 (Cheng et al., 2010), versus those that are rare, precisely excised and intragenic (by Tpb1 and
513 Tpb6) (Cheng et al., 2016; Feng et al., 2017). We could envisage if the additional *Blepharisma*
514 domesticated transposases are still capable of excision, but not a precise form, an involvement
515 in excision of a subset of the numerous intergenic IESs.

516 A single origin of PiggyBac homologs within ciliates is the most 517 parsimonious scenario

518 Though phylogenetic analyses indicate *Tetrahymena*, *Paramecium* and *Blepharisma* PiggyBac
519 homologs form a monophyletic clade the lack of PiggyBac homologs in some ciliate classes and
520 potentially the oligohymenophorean *Ichthyophthirius multifiliis* raises the question whether
521 PiggyBac IES excisases were lost or replaced in these lineages, or rather gained independently
522 from the same source by heterotrichs and a subset of oligohymenophoreans. We think the
523 former is more likely, and consistent with a long-standing hypothesis for ancestral IES excisase
524 substitution in particular ciliate lineages (Klobutcher and Herrick, 1997). However, the alternative
525 cannot be dismissed, because non-model ciliates, where the genome assembly quality allows
526 reliable gene and domain annotations, have only been sparsely sampled.

527 Future directions

528 The *B. stoltei* ATCC 30299 MAC genome together with the corresponding MIC genome (Seah et
529 al., 2022) pave the way for future investigations of genome editing in the context of a peculiar,
530 direct pathway to new MAC genome development skipping the upstream complexity of the

531 standard pathway (Miyake et al., 1991). The pair of *B. stoltei* strains used are both now low
532 frequency selfers, in which the conventional, indirect MAC development pathway dominates.
533 Comparisons with fresh, high frequency *Blepharisma* selfers collected from the wild will facilitate
534 comparative gene expression analyses with the direct MAC development pathway, which will
535 assist in distinguishing expression upregulation due to meiotic and fertilization processes
536 preceding indirect new MAC development.

537 Methods

538 Strains and localities

539 The strains used and their original isolation localities were: *Blepharisma stoltei* ATCC 30299,
540 Lake Federsee, Germany (Repak, 1968); *Blepharisma stoltei* HT-IV, Aichi prefecture, Japan;
541 *Blepharisma japonicum* R1072, from an isolate from Bangalore, India (Harumoto et al., 1998).

542 Cell cultivation, harvesting and cleanup

543 For genomic DNA isolation *B. stoltei* ATCC 30299 and HT-IV cells were cultured in Synthetic
544 Medium for *Blepharisma* (SMB) (Miyake and Beyer, 1973) at 27°C. *Blepharisma*s were fed
545 *Chlorogonium elongatum* grown in Tris-acetate phosphate (TAP) medium (Andersen, 2004) at
546 room temperature. *Chlorogonium* cells were pelleted at 1500 g at room temperature for 3
547 minutes to remove most of the TAP medium, and resuspended in 50 mL SMB. 50 ml of dense
548 *Chlorogonium* was used to feed 1 litre of *Blepharisma* culture once every three days.

549
550 *Blepharisma stoltei* ATCC 30299 and HT-IV cells used for RNA extraction were cultured in
551 Lettuce medium inoculated with *Enterbacter aerogenes* and maintained at 25°C (Miyake et al.,
552 1990).

553
554 *Blepharisma* cultures were concentrated by centrifugation in pear-shaped flasks at 100 g for 2
555 minutes using a Hettich Rotanta 460 centrifuge with swing out buckets. Pelleted cells were
556 washed with SMB and centrifuged again at 100 g for 2 minutes. The washed pellet was then
557 transferred to a cylindrical tube capped with a 100 µm-pore nylon membrane at the base and
558 immersed in SMB to filter residual algal debris from the washed cells. The cells were allowed to
559 diffuse through the membrane overnight into the surrounding medium. The next day, the
560 cylinder with the membrane was carefully removed while attempting to minimize dislodging any
561 debris collected on the membrane. Cell density after harvesting was determined by cell counting
562 under the microscope.

563 DNA isolation, library preparation and sequencing

564 *B. stoltei* macronuclei were isolated by sucrose gradient centrifugation (Lauth et al., 1976). DNA
565 was isolated with a Qiagen 20/G genomic-tip kit according to the manufacturer's instructions.
566 Purified DNA from the isolated MACs was fragmented, size selected and used to prepare
567 libraries according to standard PacBio HiFi SMRTbell protocols. The libraries were sequenced in
568 circular consensus mode to generate HiFi reads.

569
570 Total genomic DNA from *B. stoltei* HT-IV and *B. stoltei* ATCC 30299 was isolated with the
571 SigmaAldrich GenElute Mammalian genomic DNA kit. A sequencing library was prepared with a
572 NEBnext FS DNA Library Prep Kit for Illumina and sequenced on an Illumina HiSeq 3000
573 sequencer, generating 150 bp paired-end reads.

574
575 Total genomic DNA from *B. japonicum* was isolated with the Qiagen MagAttract HMW DNA kit.
576 A long-read PacBio sequencing library was prepared using the SMRTbell® Express Template
577 Preparation Kit 2.0 according to the manufacturers' instructions and sequenced on an PacBio
578 Sequel platform with 1 SMRT cell. Independently, total genomic DNA from *B. japonicum* was
579 isolated with the SigmaAldrich GenElute Mammalian genomic DNA kit and an sequencing
580 library was prepared with the TruSeq Nano DNA Library Prep Kit (Illumina) and sequenced on
581 an Illumina NovaSeq6000 to generate 150 bp paired-end reads.

582 Gamone 1/ Cell-Free Fluid (CFF) isolation and conjugation activity 583 assay

584 *B. stoltei* ATCC 30299 cells were cultured and harvested and concentrated to a density of 2000
585 cells/mL according to the procedure described in "Cell cultivation, Harvesting and Cleanup". This
586 concentrated cell culture was incubated overnight at 27°C. The next day, the cells were
587 harvested, and the supernatant collected and preserved at 4°C at all times after extraction. The
588 supernatant was then filtered through a 0.22 µm-pore filter. BSA (10 mg/mL) was added to
589 produce the final CFF at a final BSA concentration of 0.01%.

590
591 To assess the activity of the CFF, serial dilutions of the CFF were made to obtain the gamone
592 activity in terms of units (U) (Miyake, 1981). The activity of the isolated CFF was 2¹⁰ U.

593 Conjugation time course and RNA isolation for high-throughput 594 sequencing

595 *B. stoltei* cells for the complementary strains, ATCC 30299 and HT-IV, were cultivated and
596 harvested by gentle centrifugation to achieve a final cell concentration of 2000 cells/ml for each
597 strain. Non-gamone treated ATCC 30299 (A1) and HT-IV cells (H1) were collected (time point: -
598 3 hours). Strain ATCC 30299 cells were then treated with synthetic gamone 2 (final
599 concentration 1.5 µg/mL) and strain HT-IV cells were treated with cell-free fluid with a gamone 1
600 activity of $\sim 2^{10}$ U/ml for three hours (Figure S6).

601
602 Homotypic pair formation in both cultures was checked after three hours. More than 75% of the
603 cells in both cultures formed homotypic pairs. At this point the samples A2 (ATCC 30299) and
604 H2 (HT-IV) were independently isolated for RNA extraction as gamone-treated control cells just
605 before mixing. For the rest of the culture, homotypic pairs in both cultures were separated by
606 pipetting them gently with a wide-bore pipette tip. Once all pairs had been separated, the two
607 cultures were mixed together. This constitutes the experiment's 0-h time point. The conjugating
608 culture was observed and samples collected for RNA isolation or cell fixation at 2 h, 6 h, 14 h,
609 18 h, 22 h, 26 h, 30 h and 38 h (Figure S6). Further details of the sample staging approach are
610 described in (Miyake et al., 1991) and (Sugiura et al., 2012). At each time point including
611 samples A1, H1, A2 and H2, 7 mL of culture was harvested for RNA-extraction using Trizol. The
612 total RNA obtained was then separated into a small RNA fraction < 200 nt and a fraction with
613 RNA fragments > 200 nt using the Zymo RNA Clean and Concentrator-5 kit according to the
614 manufacturer's instructions. RNA-seq libraries were prepared by BGI according to their standard
615 protocols and sequenced on a BGISEq 500 instrument.

616
617 Separate 2 mL aliquots of cells at each time point for which RNA was extracted were
618 concentrated by centrifuging gently at 100 rcf. 50 µL of the concentrated cells were fixed with
619 Carnoy's fixative (ethanol:acetic acid, 6:1), stained with DAPI and imaged to determine the state
620 of nuclear development (Miyake et al., 1991).

621

622 Cell fixation and imaging

623 *B. stoltei* cells were harvested as above (“Cell cultivation”), and fixed with an equal volume of
624 “ZFAE” fixative, containing zinc sulfate (0.25 M, Sigma Aldrich), formalin, glacial acetic acid and
625 ethanol (Carl Roth), freshly prepared by mixing in a ratio of 10:2:2:5. Fixed cells were pelleted
626 (1000 g; 1 min), resuspended in 1% TritonX-100 in PHEM buffer to permeabilize (5 min; room
627 temperature), pelleted and resuspended in 2% (w/v) formaldehyde in PHEM buffer to fix further
628 (10 min; room temp.), then pelleted and washed twice with 3% (w/v) BSA in TBSTEM buffer
629 (~10 min; room temp.). For indirect immunofluorescence, washed cells were incubated with
630 primary antibody rat anti-alpha tubulin (Abcam, ab6161; 1:100 dilution in 3% w/v BSA/TBSTEM;
631 60 min; room temp.) then secondary antibody goat anti-rat IgG H&L labeled with AlexaFluor 488
632 (Abcam, ab150157, 1:500 dilution in 3% w/v BSA/TBSTEM; 20 min; room temp.). Nuclei were
633 counterstained with DAPI (1 µg/mL) in 3% (w/v) BSA/TBSTEM. A z-stack of images was
634 acquired using a confocal laser scanning microscope (Leica TCS SP8), equipped with a HC PL
635 APO 40x 1.30 Oil CS2 objective and a 1 photomultiplier tube and 3 HyD detectors, for DAPI
636 (405 nm excitation, 420-470 nm emission) and Alexa Fluor 488 (488 nm excitation, 510-530 nm
637 emission). Scanning was performed in sequential exposure mode. Spatial sampling was
638 achieved according to Nyquist criteria. ImageJ (Fiji) (Schindelin et al., 2012) was used to adjust
639 image contrast and brightness and overlay the DAPI and AlexaFluor 488 channels. The z-stack
640 was temporally color-coded.

641
642 For a nuclear 3D reconstruction (Figure 1B), cells were fixed in 1% (w/v) formaldehyde and
643 0.25% (w/v) glutaraldehyde. Nuclei were stained with Hoechst 33342 (Invitrogen) (5 µM in the
644 culture media), and imaged with a confocal laser scanning microscope (Zeiss, LSM780)
645 equipped with an LD C-Apochromat 40x/1,1 W Korr objective and a 32 channel GaAsP array
646 detector, with 405 nm excitation and 420-470 nm emission. Spatial sampling was achieved
647 according to Nyquist criteria. The IMARIS (Bitplane) software v8.0.2 was used for three-
648 dimensional reconstructions and contrast adjustments.

649 Genome assembly

650 Two MAC genome assemblies for *B. stoltei* ATCC 30299 (70x and 76x coverage) were
651 produced with Flye (version 2.7-b1585) (Kolmogorov et al., 2019) for the two separate PacBio
652 Sequel II libraries (independent replicates) using default parameters and the switches: --pacbio-

653 hifi -g 45m. The approximate genome assembly size was chosen based on preliminary Illumina
654 genome assemblies of approximately 40 Mb. Additional assemblies using the combined
655 coverage (145x) of the two libraries were produced using either Flye version 2.7-b1585 or 2.8.1-
656 b1676, and the same parameters. Two rounds of extension and merging were then used, first
657 comparing the 70x and 76x assemblies to each other, then comparing the 145x assembly to the
658 former merged assembly. Assembly graphs were all relatively simple, with few tangles to be
659 resolved (Figure S5B). Minimap2 (Li, 2018) was used for pairwise comparison of the assemblies
660 using the parameters: -x asm5 --frag=yes --secondary=no, and the resultant aligned sequences
661 were visually inspected and manually merged or extended where possible using Geneious
662 (version 2020.1.2) (Kearse et al., 2012).

663
664 Visual inspection of read mapping to the combined assembly was then used to trim off contig
665 ends where there was little correspondence between the assembly consensus and the mapped
666 reads - which we classify as "cruft". Read mapping to cruft regions was often lower or uneven,
667 suggestive of repeats. Alternatively, these features could be due to trace MIC sequences, or
668 sites of alternative chromosome breakage during development which lead to sequences that are
669 neither purely MAC nor MIC. A few contigs with similar dubious mapping of reads at internal
670 locations, which were also clear sites of chromosome fragmentation (evident by abundant
671 telomere-bearing reads in the vicinity) were split apart and trimmed back as for the contig ends.
672 Telomere-bearing reads mapped to the non-trimmed region nearest to the trimmed site were
673 then used to define contig ends, adding representative telomeric repeats from one of the
674 underlying sequences mapped to each of the ends. The main genome assembly with gene
675 predictions can be obtained from the European Nucleotide Archive (ENA) (PRJEB40285;
676 accession GCA_905310155). "Cruft" sequences are also available from the same accession.

677
678 Two separate assemblies were generated for *Blepharisma japonicum*. A genome assembly for
679 *Blepharisma japonicum* strain R1072 was generated from Illumina reads, using SPAdes
680 genome assembler (v3.14.0) (Prijbelski et al., 2020). An assembly with PacBio Sequel long
681 reads was produced with Ra (v0.2.1) (Vaser and Sikic, 2019), which uses the Overlap-Layout-
682 Consensus paradigm. The assembly produced with Ra was more contiguous, with 268 contigs,
683 in comparison to 1510 contigs in the SPAdes assembly, and was chosen as the reference
684 assembly for *Blepharisma japonicum* (ENA accession: ERR6474383).

685

686 *Condylostoma magnum* genomic reads (study accession PRJEB9019) from a previous study
687 (Swart et al., 2016) were reassembled to improve contiguity and remove bacterial
688 contamination. Reads were trimmed with `bbduk.sh` from the BBmap package v38.22
689 (<https://sourceforge.net/projects/bbmap/>), using minimum PHRED quality score 2 (both ends)
690 and k-mer trimming for Illumina adapters and Phi-X phage sequence (right end), retaining only
691 reads ≥ 25 bp. Trimmed reads were error-corrected and reassembled with SPAdes v3.13.0
692 (Prijbelski et al., 2020) using k-mer values 21, 33, 55, 77, 99. To identify potential contaminants,
693 the unassembled reads were screened with phyloFlash v3.3b1 (Gruber-Vodicka et al., 2020)
694 against SILVA v132 (Quast et al., 2013); the coding density under the standard genetic code
695 and prokaryotic gene model were also estimated using Prodigal v2.6.3 (Hyatt et al., 2010).
696 Plotting the coverage vs. GC% of the initial assembly showed that most of the likely bacterial
697 contigs (high prokaryotic coding density, lower coverage, presence of bacterial SSU rRNA
698 sequences) had $\geq 40\%$ GC, so we retained only contigs with $< 40\%$ GC as the final *C. magnum*
699 genome bin. The final assembly is available from the ENA bioproject PRJEB48875 (accession
700 GCA_920105805).

701
702 All assemblies were inspected with the quality assessment tool QUAST (Gurevich et al., 2013).

703 Variant calling

704 Illumina total genomic DNA-seq libraries for *B. stoltei* strains ATCC 30299 (ENA accession:
705 ERR6061285) and HT-IV (ERR6064674) were mapped to the ATCC 30299 reference assembly
706 with `bowtie2` v2.4.2 (Langmead and Salzberg, 2012). Alignments were tagged with the MC tag
707 (CIGAR string for mate/next segment) using `samtools` (Danecek et al., 2021) `fixmate`. The BAM
708 file was sorted and indexed, read groups were added with `bamaddrg` (commit 9baba65,
709 <https://github.com/ekg/bamaddrg>), and duplicate reads were removed with Picard
710 `MarkDuplicates` v2.25.1 (<http://broadinstitute.github.io/picard/>). Variants were called from the
711 combined BAM file with `freebayes` v1.3.2 (Garrison and Marth, 2012) in diploid mode, with
712 maximum coverage 1000 (option `-g`). The resultant VCF file was combined and indexed with
713 `bcftools` v1.12 (Danecek et al., 2021), then filtered to retain only SNPs with quality score > 20 ,
714 and at least one alternate allele.

715 Annotation of alternative telomere addition sites

716 Alternative telomere addition sites (ATASs) were annotated by mapping PacBio HiFi reads to
717 the curated reference MAC assembly described above, using minimap2 and the following flags:
718 -x asm20 --secondary=no --MD. We expect reads representing alternative telomere additions to
719 have one portion mapping to the assembly (excluding telomeric regions), with the other portion
720 containing telomeric repeats being soft-clipped in the BAM record. For each mapped read with a
721 soft-clipped segment, we extracted the clipped sequence, and the coordinates and orientation of
722 the clip relative to the reference. We searched for ≥ 24 bp tandem direct repeats of the telomere
723 unit (i.e., ≥ 3 repeats of the 8 bp unit) in the clipped segment with NCRF v1.01.02 (Harris et al.,
724 2019), which can detect tandem repeats in the presence of noise, e.g., from sequencing error.
725 The orientation of the telomere sequence, the distance from the end of the telomeric repeat to
726 the clip junction ('gap'), and the number of telomere-bearing reads vs. total mapped reads at
727 each junction were also recorded. Junctions with zero gap between telomere repeat and clip
728 junction were annotated as ATASs. The above procedure was implemented in the MILTEL
729 module of the software package BleTIES v0.1.3 (Seah and Swart, 2021).

730
731 MILTEL output was processed with Python scripts depending on Biopython (Cock et al., 2009),
732 pybedtools (Dale et al., 2011), Bedtools (Quinlan and Hall, 2010), and Matplotlib (Hunter, 2007),
733 to summarize statistics of junction sequences and telomere permutations at ATAS junctions,
734 and to extract genomic sequences flanking ATASs for sequence logos. Logos were drawn with
735 Weblogo v3.7.5 (Crooks et al., 2004), with sequences oriented such that the telomere would be
736 added on the 5' end of the ATAS junctions.

737
738 To calculate the expected minichromosome length, we assumed that ATASs were independent
739 and identically distributed in the genome following a Poisson distribution. About 47×10^3 ATASs
740 were annotated, supported on average by a single read. Given a genome of 42 Mbp at $145 \times$
741 coverage, the expected rate of encountering an ATAS is $47 \times 10^3 / (145 \times 42 \text{ Mbp})$, so the
742 distance between ATASs (i.e., the minichromosome length) is exponentially distributed with
743 expectation $(145 \times 42 \text{ Mbp}) / 47 \times 10^3 = 130 \text{ kbp}$.

744 RNA-seq read mapping

745 To permit correct mapping of tiny introns RNA-seq data was mapped to the *B. stoltei* ATCC
746 30299 MAC genome using a version of HISAT2 (Kim et al., 2019) with modified source code,
747 with the static variable minIntronLen in hisat2.cpp lowered to 9 from 20 (change available in the
748 HISAT2 github fork: <https://github.com/Swart-lab/hisat2/>; commit hash 86527b9). HISAT2 was
749 run with default parameters and parameters --min-intronlen 9 --max-intronlen 500. It should be
750 noted that RNA-seq from timepoints in which *B. stoltei* ATCC 30299 and *B. stoltei* HT-IV cells
751 were mixed together were only mapped to the former genome assembly, and so reads for up to
752 three alleles may map to each of the genes in this assembly.

753 Genetic code prediction

754 We used the program PORC (Prediction Of Reassigned Codons; available from
755 <https://github.com/Swart-lab/PORC>) previously written to predict genetic codes in protist
756 transcriptomes (Swart et al., 2016) to predict the *B. stoltei* genetic code. This program was used
757 to translate the draft *B. stoltei* ATCC 30299 genome assembly in all six frames (with the
758 standard genetic code). Like the program FACIL (Dutilh et al., 2011) that inspired PORC, the
759 frequencies of amino acids in PFAM (version 34.0) protein domain profiles aligned to the six
760 frame translation by HMMER 3.1b2 (Eddy, 2011) (default search parameters; domains used for
761 prediction with conditional E-values < 1e-20), and correspondingly also to the underlying codon,
762 are used to infer the most likely amino acid encoded by each codon (Figure S1B).

763 Gene prediction

764 We created a wrapper program, Intronarrator, to predict genes in *Blepharisma* and other
765 heterotrichs, accommodating their tiny introns. Intronarrator can be downloaded and installed
766 together with dependencies via Conda from GitHub (<https://github.com/Swart-lab/Intronarrator>).
767 Intronarrator directly infers introns from spliced RNA-seq reads mapped by HISAT2 from the
768 entire developmental time course we generated. RNA-seq reads densely cover almost the entire
769 *Blepharisma* MAC genome, aside from intergenic regions, and most potential protein-coding
770 genes (Figure 4B). After predicting the introns and removing them to create an intron-minus
771 genome, Intronarrator runs AUGUSTUS (version 3.3.3) using its intronless model. It then adds
772 back the introns to the intronless gene predictions to produce the final gene predictions.

773

774 Introns are inferred from “CIGAR” string annotations in mapped RNA-seq BAM files, using the
775 regular expression “[0-9]+M([0-9][0-9])N[0-9]+M” to select spliced reads. For intron inference we
776 only used primary alignments with: MAPQ >= 10; just a single “N”, indicating one potential
777 intron, per read; and at least 6 mapped bases flanking both the 5’ and 3’ intron boundaries (to
778 limit spurious chance matches of a few bases that might otherwise lead to incorrect intron
779 prediction). The most important parameters for Intronarrator are a cut-off of 0.2 for the fraction of
780 spliced reads covering a potential intron, and a minimum of 10 or more spliced reads to call an
781 intron. The splicing fraction cut-off was chosen based on the overall distribution of splicing
782 (Figure S4A-C). From our visual examination of mapped RNA-seq reads and gene predictions,
783 values less than this were typically “cryptic” excision events (Saudemont et al., 2017) which
784 remove potentially essential protein-coding sequences, rather than genuine introns. Intronarrator
785 classifies an intron as sense (7389 in total, excluding alternative splicing), when the majority of
786 reads (irrespective of splicing) mapping to the intron are the same strand, and antisense (554 in
787 total) when they are not. The most frequently spliced intron was chosen in rare cases of
788 overlapping alternative intron splicing.

789
790 To eliminate spurious prediction of protein-coding genes overlapping ncRNA genes, we also
791 incorporated ncRNA prediction in Intronarrator. Infernal (Nawrocki et al., 2009) (default
792 parameters; e-value < 1e-6) was used to predict a restricted set of conserved ncRNAs models
793 (i.e., tRNAs, rRNAs, SRP, and spliceosomal RNAs) from RFAM 14.0 (Kalvari et al., 2018).
794 These ncRNAs were hard-masked (with “N” characters) before AUGUSTUS gene prediction.
795 Both Infernal ncRNA predictions (excluding tRNAs) and tRNA-scan SE 2.0 (Chan et al., 2019)
796 (default parameters) tRNA predictions are annotated in the *B. stoltei* ATCC 30299 assembly
797 deposited in the European Nucleotide Archive.

798
799 Since we found that *Blepharisma stoltei*, like *Blepharisma japonicum* (Swart et al., 2016), uses a
800 non-standard genetic code, with UGA codon translated as tryptophan, gene predictions use the
801 “The Mold, Protozoan, and Coelenterate Mitochondrial Code and the Mycoplasma/Spiroplasma
802 Code (transl_table=4)” from the NCBI genetic codes. The default AUGUSTUS gene prediction
803 parameters override alternative (mitochondrial) start codons permitted by NCBI genetic code 4,
804 other than ATG. So, all predicted *B. stoltei* gene coding sequences begin with ATG.

805
806 RNA-seq read mapping relative to gene predictions of Contig_1 of *B. stoltei* ATCC30299 was
807 visualized with PyGenomeTracks (Lopez-Delisle et al., 2021).

808 Assessment of genome completeness

809 A BUSCO (version 4.0.2) (Waterhouse et al., 2018) analysis of the assembled MAC genomes of
810 *B. stoltei* and *B. japonicum* was performed on the set of predicted proteins (BUSCO mode -prot)
811 using the BUSCO Alveolata database. The completeness of the *Blepharisma* genomes was
812 compared to the protein-level BUSCO analysis of the published genome assemblies of ciliates
813 *T. thermophila*, *P. tetraurelia*, *S. coeruleus* and *I. multifiliis* (Figure S1).

814 Gene annotation

815 Pannzer2 (Törönen et al., 2018) (default parameters) and EggNog (version 2.0.1) (Huerta-
816 Cepas et al., 2019) were used for gene annotation. Annotations were combined and are
817 available from the Max Planck Society's Open Research Repository, Edmond
818 (<https://dx.doi.org/10.17617/3.8c>). Protein domain annotations were performed using hmmscan
819 from HMMER3 (version 3.3, Nov 2019) (Eddy, 2011) vs. the PFAM database (Pfam-A.full, 33.0,
820 retrieved on June 23, 2020) with default parameters.

821 Gene expression analysis

822 Features from RNA-seq reads mapped to the *B. stoltei* ATCC 30299 MAC and MAC+IES
823 genomes over the developmental time-course were extracted using featureCounts from the
824 Subread package (Liao et al., 2014). Further analysis was performed using the R software
825 environment. Genes with a total read count of less than 50, across all timepoints, were filtered
826 out of the dataset. The remaining genes were passed as a DGElist object to edgeR (Robinson
827 et al., 2010). Each time point, representing one library, was normalized for library size using the
828 edgeR function calcNormFactors. The normalized read counts were transformed into TPM
829 (transcripts per million) values (Li et al., 2010; Wagner et al., 2012). The TPM-values for
830 different genes were compared across timepoints to examine changes in gene expression.
831 Heatmaps showing log₂(TPM) changes across timepoints were plotted using the tidyverse
832 collection of R packages (<https://www.tidyverse.org/>) and RColorBrewer
833 (<https://rdr.io/cran/RColorBrewer/>). Tabulated gene expression estimates together with protein
834 annotations are available from Edmond (<https://dx.doi.org/10.17617/3.8c>).

835 Sequence visualization and analysis

836 Nucleotide and amino acid sequences were visualized using Geneious Prime (Biomatters Ltd.)
837 (Kearse et al., 2012). Multiple sequence alignments were performed with MAFFT version 7.450
838 (Kato and Standley, 2013; Kato et al., 2002). Phylogenetic trees were constructed with
839 PhyML version 3.3.20180621 (Guindon et al., 2010).

840 Orthogroup inference and analysis of orthogroup clusters

841 OrthoFinder version 2.5.2 with default parameters (i.e., using Diamond for searching, MAFFT for
842 multiple alignment and FastTree for phylogenies) was used to define orthogroups, i.e., sets of
843 genes descended from the last common ancestor of the chosen species. Proteomes for the
844 following ciliate species were used: *Tetrahymena thermophila*, *Oxytricha trifallax*, *Stentor*
845 *coeruleus* (data from ciliate.org (Stover et al., 2012)); *Euplotes octocarinatus* (EOGD (Wang et
846 al., 2018)); *Paramecium tetraurelia*, *Paramecium caudatum* (data from ParameciumDB (Arnaiz
847 et al., 2020)); plus *Perkinsus marinus* ATCC 50983 (GenBank accession: AAXJ000000000) as a
848 non-ciliate outgroup. Orthogroup clusters are available as Data S2, or from Edmond
849 (<https://dx.doi.org/10.17617/3.8c>).

850 Identification and correction of MIC-encoded PiggyBac homologs

851 We sought coding regions present within *Blepharisma* IESs to gauge the expression and type of
852 MIC-limited genes (IES assembly and gene prediction described in Seah et al. 2022). After gene
853 prediction within IESs with Intronator, predicted protein domains were annotated by HMMER
854 (v3.3) (Eddy, 2011). Several transposase families were represented in protein domains
855 identified with coding regions of IESs. However, gene prediction within IESs was hampered by
856 the presence of intermittent A-residues in the consensus sequence which occur due to the
857 inaccuracy inherent in long-reads, from which the IES regions were assembled. These errors
858 cause IES gene-prediction to falter by generating inaccurate ORFs. To circumvent this, a six-
859 frame translation of the MIC-limited genome regions was performed using a custom script,
860 which was then used to detect PFAM domains, using HMMER and the Pfam-A database 32.0
861 (release 9) (Mistry et al., 2021). Domain annotations for diagrams were generated with the
862 InterproScan 5.44-79.0 pipeline (Jones et al., 2014)

863 Four instances of the Pfam domain DDE_Tnp_1_7, characteristic of PiggyBac transposases,
864 were detected in an initial gene prediction within *Blepharisma* IESs. The four genes
865 corresponding to the DDE_Tnp_1_7 domain had high RNA-seq coverage of combined reads
866 from all timepoints across development. The IESs with the PiggyBac domains on Contig 17 and
867 Contig 39 each had two ORFs with a partial DDE_1_7 domain, separated by a few hundred bp.
868 Alignment of short-read MIC-enriched DNA reads mapped to the IES regions containing the
869 putative PiggyBac homologs indicated that several A-nucleotides in the assembled IESs were
870 insertion errors in the IES assembly, which were corrected with the short-read alignment. Open
871 reading frames of predicted genes in these corrected regions were adjusted accordingly. The
872 prefix “cORF” (corrected ORFs) was used to indicate the short-read corrected sequences of the
873 PiggyMics.

874
875 Short-read MIC-enriched DNA sequences were aligned to the IES regions containing putative
876 PiggyBac homologs with Hisat2 (2.0.0-beta) with modified source code (described above). Indel
877 errors in the IES assembly were corrected manually, then used to predict coding regions. Pfam
878 domains were annotated on MIC PiggyBac homologs with corrected ORFs using the
879 InterproScan (v. 1.1.4) (Quevillon et al., 2005) plugin in Geneious v11.1.5 (Biomatter Ltd.).
880 DDE_Tnp_1_7 domains were detected in the corrected ORFs, which in some cases spanned
881 IES regions lacking predicted genic regions before correction. A multiple sequence alignment of
882 the correct MIC PiggyBac homologs with other ciliate PiggyBac-derived proteins (PGBDs) and
883 eukaryotic PiggyBac-like elements (PBLEs) that contain the PiggyBac transposase domain
884 DDE_Tnp_1_7 (PF13843) was performed with MAFFT (v4.1) via the Geneious plugin (algorithm
885 L-INS-i, BLOSUM62 scoring matrix, gap open penalty 1.53, offset value 0.123). A phylogenetic
886 tree was constructed using the FastTree (v 2.1.11) plugin for Geneious (Whelan-Goldman
887 model).

888 d_N/d_S estimation

889 We generated pairwise coding sequence alignments of PiggyMac paralog nucleotide sequences
890 from *P. tetraurelia* and *P. octaurelia* using MAFFT version 7.450 (Kato and Standley, 2013)
891 (Kato et al., 2002) (algorithm: “auto”, scoring matrix: 200PAM/k=2, gap open penalty 1.53,
892 offset value 0.123) using the “translation align” panel of Geneious Prime (version 2020.1.2)
893 (Kearse et al., 2012). PAML version 4.9 (Yang, 2007) was used to estimate d_N/d_S values in
894 pairwise mode (runmode = -2, seqtype = 1, CodonFreq = 2). For *Blepharisma stoltei*, we

895 generated pairwise coding sequence alignments of the *Blepharisma* PiggyMac homolog, BPgm
896 (Contig_49.g1063; BSTOLATCC_MAC17466), with the *Blepharisma* Pgm-likes (BPgmLs) using
897 Translation Align panel of Geneious v11.1.5 (Genetic code: *Blepharisma*, Protein alignment
898 options: MAFFT alignment (v7.450) (Katoch and Standley, 2013), scoring matrix: BLOSUM62,
899 Gap open penalty: 1.53, offset value: 0.1). PAML version 4.9 was used to estimate dN/dS values
900 in pairwise mode (runmode = -2, seqtype = 1, CodonFreq = 2).

901 Phylogenetic analysis

902 Protein sequences of PBLEs were obtained from Bouallègue et al (Bouallègue et al., 2017).
903 Protein sequences of *Paramecium* and *Tetrahymena* Pgms and PgmLs were obtained from
904 ParameciumDB (Arnaiz et al., 2020) (PGM, PGMLs1-5) and ciliate.org (Stover et al., 2012)
905 (Tpb1, Tpb2, Tpb7, LIA5), respectively. *Condyllostoma* and *Blepharisma* Pgms and PgmLs were
906 obtained from genome assemblies (accessions GCA_920105805 and GCA_905310155,
907 respectively). Sequence manipulation was done using Geneious (Biomatters Ltd.). The
908 Geneious plug-in for InterProScan (Jones et al., 2014) was used to identify DDE_Tnp_1_7
909 domains using the PFAM-A database (Mistry et al., 2021). The DDE_Tnp_1_7 domain and
910 regions adjacent to it were extracted and aligned using the MAFFT plug-in (v7.450) for
911 Geneious (Katoch and Standley, 2013) (Algorithm: L-INS-i, Scoring matrix: BLOSUM62, Gap
912 open penalty: 1.53, Offset value: 0.123). Phylogenetic trees using this alignment were generated
913 with the FastTree2 (v2.2.11) Geneious plug-in using the Whelan-Goldman model. The
914 phylogenetic trees were visualized with FigTree (v1.4.4) (Andrew Rambaut,
915 <http://tree.bio.ed.ac.uk/>).

916 Repeat annotation

917 Interspersed repeat element families were predicted with RepeatModeler v2.0.1 (default
918 settings, random number seed 12345) with the following dependencies: rmbblast v2.9.0+
919 (<http://www.repeatmasker.org/RMBlast.html>), TRF 4.09 (Benson, 1999), RECON (Bao and
920 Eddy, 2002), RepeatScout 1.0.6 (Price et al., 2005), RepeatMasker v4.1.1
921 (<http://www.repeatmasker.org/RMDownload.html>). Repeat families were also classified in the
922 pipeline by RepeatClassifier v2.0.1 through comparison against RepeatMasker's repeat protein
923 database and the Dfam database. Consensus sequences of the predicted repeat families,
924 produced by RepeatModeler, were then used to annotate repeats with RepeatMasker, using
925 rmbblast as the search engine.

926
927 Terminal inverted repeats (TIRs) of selected repeat element families were identified by aligning
928 the consensus sequence from RepeatModeler, and/or selected full-length elements, with their
929 respective reverse complements using MAFFT (Kato and Standley, 2013) (plugin version
930 distributed with Geneious). TIRs from the Dfam DNA transposon termini signatures database
931 (v1.1, https://www.dfam.org/releases/dna_termini_1.1/dna_termini_1.1.hmm.gz) (Storer et al.,
932 2021) were searched with hmmsearch (HMMer v3.2.1) against the IES sequences, to identify
933 matches to TIR signatures of major transposon subfamilies.
934

935 Data and code availability

936 The draft *Blepharisma stoltei* ATCC 30299 MAC genome assembly is accessible from
937 bleph.ciliate.org and from the European Nucleotide Archive (ENA) bioproject PRJEB40285
938 under the accession GCA_905310155. PacBio CCS reads (ERR5873783 and ERR5873334)
939 and subreads (ERR5962314) used to assemble the genome are also available from ENA.
940 Illumina DNA-seq data for the *B. stoltei* ATCC 30299 and HT-IV strains is available from
941 accessions ERR6061285 and ERR6064674, respectively. The RNA-seq developmental time
942 course is available from the bioproject PRJEB45374 (accessions ERR6049461-ERR6049485).
943
944 Illumina and PacBio Sequel sequencing data for *Blepharisma japonicum* strain R1702 is
945 available from the ENA bioproject PRJEB46921 (Illumina accessions: ERR6473251,
946 ERR6474356; PacBio accession: ERR6474383).
947
948 Code availability for software we generated or modified is indicated in place in Methods.
949

950 Acknowledgements

951 This paper is dedicated in memory of Akio Miyake and his decades of inspirational *Blepharisma*
952 research. We thank Federico Buonanno for the provision of *B. stoltei* ATCC 30299 cells and
953 culturing advice, Christa Lanz and the MPI for Biology's genome center, Sebastien Colin and the
954 MPI for Biology's Light Microscopy Facility for the 3D nuclear reconstruction, and Adrian Streit

955 for discussion. Research reported in this publication was supported by the National Institutes of
956 Health (award No. P40OD010964) to N.A.S and the Max Planck Society.
957

958 Author contributions

959 Conceptualization, M.S.1., K.B.B.S., E.C.S.; Methodology, M.S.1., K.B.B.S., C.E., A.S., C.W.,
960 B.H., E.C.S.; Software, M.S.1., K.B.B.S., E.C.S.; Investigation, M.S.1., K.B.B.S., C.W., B.H.,
961 E.C.S.; Writing – Original Draft, M.S.1., K.B.B.S., A.S., C.W., B.H., E.C.S.; Writing – Review &
962 Editing, M.S.1, K.B.B.S., A.S., E.C.S.; Funding Acquisition, E.C.S.; Resources, A.B. and N.A.S.;
963 Supervision, M.S.2., T.H., E.C.S.
964

965 Declaration of interests

966 The authors declare no competing interests.
967
968

969 Bibliography

970 Aeschlimann, S.H., Jönsson, F., Postberg, J., Stover, N.A., Petera, R.L., Lipps, H.-J., Nowacki,
971 M., and Swart, E.C. (2014). The draft assembly of the radically organized *Stylonychia lemnae*
972 macronuclear genome. *Genome Biol. Evol.* 6, 1707–1723.

973 Andersen, R.A. (2004). *Algal Culturing Techniques*.

974 Arnaiz, O., Mathy, N., Baudry, C., Malinsky, S., Aury, J.-M., Denby Wilkes, C., Garnier, O.,
975 Labadie, K., Lauderdale, B.E., Le Mouël, A., et al. (2012). The *Paramecium* germline genome
976 provides a niche for intragenic parasitic DNA: evolutionary dynamics of internal eliminated
977 sequences. *PLoS Genet.* 8, e1002984.

978 Arnaiz, O., Meyer, E., and Sperling, L. (2020). ParameciumDB 2019: integrating genomic data
979 across the genus for functional and evolutionary biology. *Nucleic Acids Res.* 48, D599–D605.

980 Aury, J.-M., Jaillon, O., Duret, L., Noel, B., Jubin, C., Porcel, B.M., Ségurens, B., Daubin, V.,

- 981 Anthouard, V., Aiach, N., et al. (2006). Global trends of whole-genome duplications revealed by
982 the ciliate *Paramecium tetraurelia*. *Nature* **444**, 171–178.
- 983 Baudry, C., Malinsky, S., Restituto, M., Kapusta, A., Rosa, S., Meyer, E., and Bétermier, M.
984 (2009). PiggyMac, a domesticated piggyBac transposase involved in programmed genome
985 rearrangements in the ciliate *Paramecium tetraurelia*. *Genes Dev.* **23**, 2478–2483.
- 986 Bischerour, J., Bhullar, S., Denby Wilkes, C., Régnier, V., Mathy, N., Dubois, E., Singh, A.,
987 Swart, E., Arnaiz, O., Sperling, L., et al. (2018). Six domesticated PiggyBac transposases
988 together carry out programmed DNA elimination in *Paramecium*. *ELife* **7**.
- 989 Bouallègue, M., Rouault, J.-D., Hua-Van, A., Makni, M., and Capy, P. (2017). Molecular
990 Evolution of piggyBac Superfamily: From Selfishness to Domestication. *Genome Biol. Evol.* **9**,
991 323–339.
- 992 Casola, C., Hucks, D., and Feschotte, C. (2008). Convergent domestication of pogo-like
993 transposases into centromere-binding proteins in fission yeast and mammals. *Mol. Biol. Evol.*
994 **25**, 29–41.
- 995 Chalker, D.L., Meyer, E., and Mochizuki, K. (2013). Epigenetics of ciliates. *Cold Spring Harb.*
996 *Perspect. Biol.* **5**, a017764.
- 997 Chan, P.P., Lin, B.Y., Mak, A.J., and Lowe, T.M. (2019). tRNAscan-SE 2.0: Improved Detection
998 and Functional Classification of Transfer RNA Genes. *BioRxiv*.
- 999 Cheng, C.-Y., Vogt, A., Mochizuki, K., and Yao, M.-C. (2010). A domesticated piggyBac
1000 transposase plays key roles in heterochromatin dynamics and DNA cleavage during
1001 programmed DNA deletion in *Tetrahymena thermophila*. *Mol. Biol. Cell* **21**, 1753–1762.
- 1002 Cheng, C.-Y., Young, J.M., Lin, C.-Y.G., Chao, J.-L., Malik, H.S., and Yao, M.-C. (2016). The
1003 piggyBac transposon-derived genes TPB1 and TPB6 mediate essential transposon-like excision
1004 during the developmental rearrangement of key genes in *Tetrahymena thermophila*. *Genes Dev.*
1005 **30**, 2724–2736.
- 1006 Cheng, Y.-H., Liu, C.-F.J., Yu, Y.-H., Jhou, Y.-T., Fujishima, M., Tsai, I.J., and Leu, J.-Y. (2020).
1007 Genome plasticity in *Paramecium bursaria* revealed by population genomics. *BMC Biol.* **18**, 180.
- 1008 Chen, Q., Luo, W., Veach, R.A., Hickman, A.B., Wilson, M.H., and Dyda, F. (2020). Structural
1009 basis of seamless excision and specific targeting by piggyBac transposase. *Nat. Commun.* **11**,
1010 3446.
- 1011 Chen, X., Bracht, J.R., Goldman, A.D., Dolzhenko, E., Clay, D.M., Swart, E.C., Perlman, D.H.,

- 1012 Doak, T.G., Stuart, A., Amemiya, C.T., et al. (2014). The architecture of a scrambled genome
1013 reveals massive levels of genomic rearrangement during development. *Cell* *158*, 1187–1198.
- 1014 Chin, C.-S., Peluso, P., Sedlazeck, F.J., Nattestad, M., Concepcion, G.T., Clum, A., Dunn, C.,
1015 O'Malley, R., Figueroa-Balderas, R., Morales-Cruz, A., et al. (2016). Phased diploid genome
1016 assembly with single-molecule real-time sequencing. *Nat. Methods* *13*, 1050–1054.
- 1017 Cock, P.J.A., Antao, T., Chang, J.T., Chapman, B.A., Cox, C.J., Dalke, A., Friedberg, I.,
1018 Hamelryck, T., Kauff, F., Wilczynski, B., et al. (2009). Biopython: freely available Python tools for
1019 computational molecular biology and bioinformatics. *Bioinformatics* *25*, 1422–1423.
- 1020 Crooks, G.E., Hon, G., Chandonia, J.M., and Brenner, S.E. (2004). WebLogo: a sequence logo
1021 generator. *Genome Res.* *14*, 1188–1190.
- 1022 Dale, R.K., Pedersen, B.S., and Quinlan, A.R. (2011). Pybedtools: a flexible Python library for
1023 manipulating genomic datasets and annotations. *Bioinformatics* *27*, 3423–3424.
- 1024 Danecek, P., Bonfield, J.K., Liddle, J., Marshall, J., Ohan, V., Pollard, M.O., Whitwham, A.,
1025 Keane, T., McCarthy, S.A., Davies, R.M., et al. (2021). Twelve years of SAMtools and BCFtools.
1026 *Gigascience* *10*.
- 1027 Doak, T.G., Doerder, F.P., Jahn, C.L., and Herrick, G. (1994). A proposed superfamily of
1028 transposase genes: transposon-like elements in ciliated protozoa and a common “D35E” motif.
1029 *Proc Natl Acad Sci USA* *91*, 942–946.
- 1030 Dupeyron, M., Baril, T., Bass, C., and Hayward, A. (2020). Phylogenetic analysis of the
1031 Tc1/mariner superfamily reveals the unexplored diversity of pogo-like elements. *Mob. DNA* *11*,
1032 21.
- 1033 Dutilh, B.E., Jurgelenaite, R., Szklarczyk, R., van Hijum, S.A.F.T., Harhangi, H.R., Schmid, M.,
1034 de Wild, B., François, K.-J., Stunnenberg, H.G., Strous, M., et al. (2011). FACIL: Fast and
1035 Accurate Genetic Code Inference and Logo. *Bioinformatics* *27*, 1929–1933.
- 1036 Eddy, S.R. (2011). Accelerated profile HMM searches. *PLoS Comput. Biol.* *7*, e1002195.
- 1037 Elick, T.A., Bauser, C.A., and Fraser, M.J. (1996). Excision of the piggyBac transposable
1038 element in vitro is a precise event that is enhanced by the expression of its encoded
1039 transposase. *Genetica* *98*, 33–41.
- 1040 Feng, L., Wang, G., Hamilton, E.P., Xiong, J., Yan, G., Chen, K., Chen, X., Dui, W., Plemens,
1041 A., Khadr, L., et al. (2017). A germline-limited piggyBac transposase gene is required for precise
1042 excision in *Tetrahymena* genome rearrangement. *Nucleic Acids Res.* *45*, 9481–9502.

- 1043 Friedl, E., Miyake, A., and Heckmann, K. (1983). Requirement of successive protein syntheses
1044 for the progress of meiosis in *Blepharisma*. *Exp. Cell Res.* *145*, 105–113.
- 1045 Gao, B., Wang, Y., Diaby, M., Zong, W., Shen, D., Wang, S., Chen, C., Wang, X., and Song, C.
1046 (2020). Evolution of pogo, a separate superfamily of IS630-Tc1-mariner transposons, revealing
1047 recurrent domestication events in vertebrates. *Mob. DNA* *11*, 25.
- 1048 Garrison, E., and Marth, G. (2012). Haplotype-based variant detection from short-read
1049 sequencing (arXiv).
- 1050 Giese, A.C. (1973). *Blepharisma: The Biology of a Light-sensitive Protozoan* (Stanford
1051 University Press).
- 1052 Gruber-Vodicka, H.R., Seah, B.K.B., and Pruesse, E. (2020). phyloFlash: Rapid Small-Subunit
1053 rRNA Profiling and Targeted Assembly from Metagenomes. *MSystems* *5*.
- 1054 Guérineau, M., Bessa, L., Moriau, S., Lescop, E., Bontems, F., Mathy, N., Guittet, E.,
1055 Bischerour, J., Bétermier, M., and Morellet, N. (2021). The unusual structure of the PiggyMac
1056 cysteine-rich domain reveals zinc finger diversity in PiggyBac-related transposases. *Mob. DNA*
1057 *12*, 12.
- 1058 Guindon, S., Dufayard, J.-F., Lefort, V., Anisimova, M., Hordijk, W., and Gascuel, O. (2010).
1059 New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the
1060 performance of PhyML 3.0. *Syst. Biol.* *59*, 307–321.
- 1061 Gurevich, A., Saveliev, V., Vyahhi, N., and Tesler, G. (2013). QUASt: quality assessment tool
1062 for genome assemblies. *Bioinformatics* *29*, 1072–1075.
- 1063 Hamilton, E.P., Kapusta, A., Huvos, P.E., Bidwell, S.L., Zafar, N., Tang, H., Hadjithomas, M.,
1064 Krishnakumar, V., Badger, J.H., Caler, E.V., et al. (2016). Structure of the germline genome of
1065 *Tetrahymena thermophila* and relationship to the massively rearranged somatic genome. *ELife*
1066 *5*.
- 1067 Harris, R.S., Cechova, M., and Makova, K.D. (2019). Noise-cancelling repeat finder: uncovering
1068 tandem repeats in error-prone long-read sequencing data. *Bioinformatics* *35*, 4809–4811.
- 1069 Harumoto, T., Miyake, A., Ishikawa, N., Sugibayashi, R., Zenfuku, K., and Iio, H. (1998).
1070 Chemical defense by means of pigmented extrusomes in the ciliate *Blepharisma japonicum*.
1071 *Eur. J. Protistol.* *34*, 458–470.
- 1072 Hohmann, S. (1993). Characterisation of PDC2, a gene necessary for high level expression of
1073 pyruvate decarboxylase structural genes in *Saccharomyces cerevisiae*. *Mol. Gen. Genet.* *241*,

- 1074 657–666.
- 1075 Huerta-Cepas, J., Szklarczyk, D., Heller, D., Hernández-Plaza, A., Forslund, S.K., Cook, H.,
1076 Mende, D.R., Letunic, I., Rattei, T., Jensen, L.J., et al. (2019). eggNOG 5.0: a hierarchical,
1077 functionally and phylogenetically annotated orthology resource based on 5090 organisms and
1078 2502 viruses. *Nucleic Acids Res.* 47, D309–D314.
- 1079 Hunter, J.D. (2007). Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* 9, 90–95.
- 1080 Hyatt, D., Chen, G.-L., Locascio, P.F., Land, M.L., Larimer, F.W., and Hauser, L.J. (2010).
1081 Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC*
1082 *Bioinformatics* 11, 119.
- 1083 Inaba, F. (1965). Conjugation between two strains of *Blepharisma*. *J. Protozool.* 12, 146–151.
- 1084 Jahn, C.L., Doktor, S.Z., Frels, J.S., Jaraczewski, J.W., and Krikau, M.F. (1993). Structures of
1085 the *Euplotes crassus* Tec1 and Tec2 elements: identification of putative transposase coding
1086 regions. *Gene* 133, 71–78.
- 1087 Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen, J.,
1088 Mitchell, A., Nuka, G., et al. (2014). InterProScan 5: genome-scale protein function
1089 classification. *Bioinformatics* 30, 1236–1240.
- 1090 Kalvari, I., Argasinska, J., Quinones-Olvera, N., Nawrocki, E.P., Rivas, E., Eddy, S.R., Bateman,
1091 A., Finn, R.D., and Petrov, A.I. (2018). Rfam 13.0: shifting to a genome-centric resource for non-
1092 coding RNA families. *Nucleic Acids Res.* 46, D335–D342.
- 1093 Katashima, R.Y.O. (1959). Mating Types in *Euplotes eurystomus*. *J. Protozool.* 6, 75–83.
- 1094 Katoh, K., and Standley, D.M. (2013). MAFFT multiple sequence alignment software version 7:
1095 improvements in performance and usability. *Mol. Biol. Evol.* 30, 772–780.
- 1096 Katoh, K., Misawa, K., Kuma, K., and Miyata, T. (2002). MAFFT: a novel method for rapid
1097 multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* 30, 3059–
1098 3066.
- 1099 Kearse, M., Moir, R., Wilson, A., Stones-Havas, S., Cheung, M., Sturrock, S., Buxton, S.,
1100 Cooper, A., Markowitz, S., Duran, C., et al. (2012). Geneious Basic: an integrated and
1101 extendable desktop software platform for the organization and analysis of sequence data.
1102 *Bioinformatics* 28, 1647–1649.
- 1103 Kimball, R.F. (1942). The nature and inheritance of mating types in *euplotes patella*. *Genetics*
1104 27, 269–285.

- 1105 Kim, D., Paggi, J.M., Park, C., Bennett, C., and Salzberg, S.L. (2019). Graph-based genome
1106 alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* 37, 907–915.
- 1107 Klobutcher, L.A., and Herrick, G. (1995). Consensus inverted terminal repeat sequence of
1108 *Paramecium* IESs: resemblance to termini of Tc1-related and *Euplotes* Tec transposons.
1109 *Nucleic Acids Res.* 23, 2006–2013.
- 1110 Klobutcher, L.A., and Herrick, G. (1997). Developmental genome reorganization in ciliated
1111 protozoa: the transposon link. *Prog. Nucleic Acid Res. Mol. Biol.* 56, 1–62.
- 1112 Kobayashi, M., Miura, M., Takusagawa, M., Sugiura, M., and Harumoto, T. (2015). Two possible
1113 barriers blocking conjugation between different megakaryotypes of *Blepharisma*. *Zool. Sci.* 32,
1114 53–61.
- 1115 Kolmogorov, M., Yuan, J., Lin, Y., and Pevzner, P.A. (2019). Assembly of long, error-prone
1116 reads using repeat graphs. *Nat. Biotechnol.* 37, 540–546.
- 1117 Kriventseva, E.V., Kuznetsov, D., Tegenfeldt, F., Manni, M., Dias, R., Simão, F.A., and
1118 Zdobnov, E.M. (2019). OrthoDB v10: sampling the diversity of animal, plant, fungal, protist,
1119 bacterial and viral genomes for evolutionary and functional annotations of orthologs. *Nucleic
1120 Acids Res.* 47, D807–D811.
- 1121 Kubota, T., Tokoroyama, T., Tsukuda, Y., Koyama, H., and Miyake, A. (1973). Isolation and
1122 structure determination of blepharimin, a conjugation initiating gamone in the ciliate
1123 *blepharisma*. *Science* 179, 400–402.
- 1124 Kumazawa, H. (1979). Homopolar Grafting in *Blepharisma japonicum*. *Journal of Experimental
1125 Zoology* 207, 1–16.
- 1126 Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat.
1127 Methods* 9, 357–359.
- 1128 Lauth, M.R., Spear, B.B., Heumann, J., and Prescott, D.M. (1976). DNA of ciliated protozoa:
1129 DNA sequence diminution during macronuclear development of *Oxytricha*. *Cell* 7, 67–74.
- 1130 Liao, Y., Smyth, G.K., and Shi, W. (2014). featureCounts: an efficient general purpose program
1131 for assigning sequence reads to genomic features. *Bioinformatics* 30, 923–930.
- 1132 Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34, 3094–
1133 3100.
- 1134 Li, B., Ruotti, V., Stewart, R.M., Thomson, J.A., and Dewey, C.N. (2010). RNA-Seq gene
1135 expression estimation with read mapping uncertainty. *Bioinformatics* 26, 493–500.

- 1136 Lopez-Delisle, L., Rabbani, L., Wolff, J., Bhardwaj, V., Backofen, R., Grüning, B., Ramírez, F.,
1137 and Manke, T. (2021). pyGenomeTracks: reproducible plots for multivariate genomic datasets.
1138 *Bioinformatics* 37, 422–423.
- 1139 van Luenen, H.G., Colloms, S.D., and Plasterk, R.H. (1994). The mechanism of transposition of
1140 Tc3 in *C. elegans*. *Cell* 79, 293–301.
- 1141 Luporini, P., Miceli, C., and Ortenzi, C. (1983). Evidence that the ciliate *Euplotes raikovi*
1142 releases mating-inducing factors (gamones). *J. Exp. Zool.* 226, 1–9.
- 1143 Lynn, D.H. (2010). *The Ciliated Protozoa* (Dordrecht: Springer Netherlands).
- 1144 Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G.A., Sonnhammer, E.L.L.,
1145 Tosatto, S.C.E., Paladin, L., Raj, S., Richardson, L.J., et al. (2021). Pfam: The protein families
1146 database in 2021. *Nucleic Acids Res.* 49, D412–D419.
- 1147 Miyake, A. (1981). Cell interaction by gamones in *blepharisma*. In *Sexual Interactions in*
1148 *Eukaryotic Microbes*, (Elsevier), pp. 95–129.
- 1149 Miyake, A., and Beyer, J. (1973). Cell interaction by means of soluble factors (gamones) in
1150 conjugation of *Blepharisma intermedium*. *Exp. Cell Res.* 76, 15–24.
- 1151 Miyake, A., and Beyer, J. (1974). Blepharhormone: a conjugation-inducing glycoprotein in the ciliate
1152 *blepharisma*. *Science* 185, 621–623.
- 1153 Miyake, A., and Harumoto, T. (1990). Asymmetrical cell division in *Blepharisma japonicum*:
1154 difference between daughter cells in mating-type expression. *Exp. Cell Res.* 190, 65–68.
- 1155 Miyake, A., Tulli, M., and Nobili, R. (1979). Requirement of protein synthesis in the initiation of
1156 meiosis and other nuclear changes in conjugation of *Blepharisma*. *Exp. Cell Res.* 120, 87–93.
- 1157 Miyake, A., Harumoto, T., Salvi, B., and Rivola, V. (1990). Defensive function of pigment
1158 granules in *Blepharisma japonicum*. *Eur. J. Protistol.* 25, 310–315.
- 1159 Miyake, A., Rivola, V., and Harumoto, T. (1991). Double paths of macronucleus differentiation at
1160 conjugation in *Blepharisma japonicum*. *Eur. J. Protistol.* 27, 178–200.
- 1161 Mojzita, D., and Hohmann, S. (2006). Pdc2 coordinates expression of the THI regulon in the
1162 yeast *Saccharomyces cerevisiae*. *Mol. Genet. Genomics* 276, 147–161.
- 1163 Nawrocki, E.P., Kolbe, D.L., and Eddy, S.R. (2009). Infernal 1.0: inference of RNA alignments.
1164 *Bioinformatics* 25, 1335–1337.
- 1165 Nowacki, M., Higgins, B.P., Maquilan, G.M., Swart, E.C., Doak, T.G., and Landweber, L.F.

- 1166 (2009). A functional role for transposases in a large eukaryotic genome. *Science* 324, 935–938.
- 1167 Prescott, D.M. (1994). The DNA of ciliated protozoa. *Microbiol. Rev.* 58, 233–267.
- 1168 Prjibelski, A., Antipov, D., Meleshko, D., Lapidus, A., and Korobeynikov, A. (2020). Using
1169 SPAdes de novo assembler. *Curr. Protoc. Bioinformatics* 70, e102.
- 1170 Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies, J., and
1171 Glöckner, F.O. (2013). The SILVA ribosomal RNA gene database project: improved data
1172 processing and web-based tools. *Nucleic Acids Res.* 41, D590-6.
- 1173 Quevillon, E., Silventoinen, V., Pillai, S., Harte, N., Mulder, N., Apweiler, R., and Lopez, R.
1174 (2005). InterProScan: protein domains identifier. *Nucleic Acids Res.* 33, W116-20.
- 1175 Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing genomic
1176 features. *Bioinformatics* 26, 841–842.
- 1177 Rapport, E.W., Rapport, D.J., Berger, J., and Kupers, V. (1976). Induction of conjugation in
1178 *Stentor coeruleus*. *Trans. Am. Microsc. Soc.* 95, 220–224.
- 1179 Repak, A.J. (1968). Encystment and excystment of the heterotrichous ciliate *Blepharisma stoltei*
1180 Isquith. *Journal of Protozoology* 5, 407–412.
- 1181 Robinson, M.D., McCarthy, D.J., and Smyth, G.K. (2010). edgeR: a Bioconductor package for
1182 differential expression analysis of digital gene expression data. *Bioinformatics* 26, 139–140.
- 1183 Salvini, M., Durante, M., and Nobili, R. (1983). Characterization of macronuclear DNA in
1184 *Blepharisma japonicum*. *Protoplasma* 117, 82–88.
- 1185 Sarkar, A., Sim, C., Hong, Y.S., Hogan, J.R., Fraser, M.J., Robertson, H.M., and Collins, F.H.
1186 (2003). Molecular evolutionary analysis of the widespread piggyBac transposon family and
1187 related “domesticated” sequences. *Mol. Genet. Genomics* 270, 173–180.
- 1188 Saudemont, B., Popa, A., Parmley, J.L., Rocher, V., Blugeon, C., Necșulea, A., Meyer, E., and
1189 Duret, L. (2017). The fitness cost of mis-splicing is the main determinant of alternative splicing
1190 patterns. *Genome Biol.* 18, 208.
- 1191 Schindelin, J., Arganda-Carreras, I., Frise, E., Kaynig, V., Longair, M., Pietzsch, T., Preibisch,
1192 S., Rueden, C., Saalfeld, S., Schmid, B., et al. (2012). Fiji: an open-source platform for
1193 biological-image analysis. *Nat. Methods* 9, 676–682.
- 1194 Seah, B.K.B., and Swart, E.C. (2021). BleTIES: Annotation of natural genome editing in ciliates
1195 using long read sequencing. *Bioinformatics* 37, 3929–3931.

- 1196 Sheng, Y., Duan, L., Cheng, T., Qiao, Y., Stover, N.A., and Gao, S. (2020). The completed
1197 macronuclear genome of a model ciliate *Tetrahymena thermophila* and its application in genome
1198 scrambling and copy number analyses. *Sci. China Life Sci.* **63**, 1534–1542.
- 1199 Slabodnick, M.M., and Marshall, W.F. (2014). *Stentor coeruleus*. *Curr. Biol.* **24**, R783-4.
- 1200 Slabodnick, M.M., Ruby, J.G., Reiff, S.B., Swart, E.C., Gosai, S., Prabakaran, S., Witkowska, E.,
1201 Larue, G.E., Fisher, S., Freeman, R.M., et al. (2017). The Macronuclear Genome of *Stentor*
1202 *coeruleus* Reveals Tiny Introns in a Giant Cell. *Curr. Biol.* **27**, 569–575.
- 1203 Stover, N.A., Punia, R.S., Bowen, M.S., Dolins, S.B., and Clark, T.G. (2012). *Tetrahymena*
1204 Genome Database Wiki: a community-maintained model organism database. *Database (Oxford)*
1205 **2012**, bas007.
- 1206 Sugiura, M., and Harumoto, T. (2001). Identification, characterization, and complete amino acid
1207 sequence of the conjugation-inducing glycoprotein (blepharhormone) in the ciliate *Blepharisma*
1208 *japonicum*. *Proc Natl Acad Sci USA* **98**, 14446–14451.
- 1209 Sugiura, M., Shiotani, H., Suzaki, T., and Harumoto, T. (2010). Behavioural changes induced by
1210 the conjugation-inducing pheromones, gamone 1 and 2, in the ciliate *Blepharisma japonicum*.
1211 *Eur. J. Protistol.* **46**, 143–149.
- 1212 Sugiura, M., Tanaka, Y., Suzaki, T., and Harumoto, T. (2012). Alternative gene expression in
1213 type I and type II cells may enable further nuclear changes during conjugation of *Blepharisma*
1214 *japonicum*. *Protist* **163**, 204–216.
- 1215 Suzuki, S. (1957). Parthenogenetic conjugation in *Blepharisma undulans japonicus* Suzuki. *Bull.*
1216 *Yamagata Univ. Natural Sci.* **4**, 69–84.
- 1217 Swart, E.C., and Nowacki, M. (2015). The eukaryotic way to defend and edit genomes by sRNA-
1218 targeted DNA deletion. *Ann. N. Y. Acad. Sci.* **1341**, 106–114.
- 1219 Swart, E.C., Bracht, J.R., Magrini, V., Minx, P., Chen, X., Zhou, Y., Khurana, J.S., Goldman,
1220 A.D., Nowacki, M., Schotanus, K., et al. (2013). The *Oxytricha trifallax* macronuclear genome: a
1221 complex eukaryotic genome with 16,000 tiny chromosomes. *PLoS Biol.* **11**, e1001473.
- 1222 Swart, E.C., Serra, V., Petroni, G., and Nowacki, M. (2016). Genetic Codes with No Dedicated
1223 Stop Codon: Context-Dependent Translation Termination. *Cell* **166**, 691–702.
- 1224 Terazima, M.N., and Harumoto, T. (2004). Defense function of pigment granules in the ciliate
1225 *Blepharisma japonicum* against two predatory protists, *Amoeba proteus* (Rhizopodea) and
1226 *Climacostomum virens* (Ciliata). *Zool. Sci.* **21**, 823–828.

- 1227 Törönen, P., Medlar, A., and Holm, L. (2018). PANNZER2: a rapid functional annotation web
1228 server. *Nucleic Acids Res.* *46*, W84–W88.
- 1229 Vallesi, A., Giuli, G., Bradshaw, R.A., and Luporini, P. (1995). Autocrine mitogenic activity of
1230 pheromones produced by the protozoan ciliate *Euplotes raikovi*. *Nature* *376*, 522–524.
- 1231 Vaser, R., and Sikic, M. (2019). Yet another de novo genome assembler. *BioRxiv*.
- 1232 Vogt, A., and Mochizuki, K. (2013). A domesticated PiggyBac transposase interacts with
1233 heterochromatin and catalyzes reproducible DNA elimination in *Tetrahymena*. *PLoS Genet.* *9*,
1234 e1004032.
- 1235 Vogt, A., Goldman, A.D., Mochizuki, K., and Landweber, L.F. (2013). Transposon domestication
1236 versus mutualism in ciliate genome rearrangements. *PLoS Genet.* *9*, e1003659.
- 1237 Wagner, G.P., Kin, K., and Lynch, V.J. (2012). Measurement of mRNA abundance using RNA-
1238 seq data: RPKM measure is inconsistent among samples. *Theory Biosci.* *131*, 281–285.
- 1239 Wang, R.-L., Miao, W., Wang, W., Xiong, J., and Liang, A.-H. (2018). EOGD: the *Euplotes*
1240 *octocarinatus* genome database. *BMC Genomics* *19*, 63.
- 1241 Waterhouse, R.M., Seppey, M., Simão, F.A., Manni, M., Ioannidis, P., Klioutchnikov, G.,
1242 Kriventseva, E.V., and Zdobnov, E.M. (2018). BUSCO applications from quality assessments to
1243 gene prediction and phylogenomics. *Mol. Biol. Evol.* *35*, 543–548.
- 1244 Wick, R.R., Schultz, M.B., Zobel, J., and Holt, K.E. (2015). Bandage: interactive visualization of
1245 de novo genome assemblies. *Bioinformatics* *31*, 3350–3352.
- 1246 Williams, K., Doak, T.G., and Herrick, G. (1993). Developmental precise excision of *Oxytricha*
1247 *trifallax* telomere-bearing elements and formation of circles closed by a copy of the flanking
1248 target duplication. *EMBO J.* *12*, 4593–4601.
- 1249 Witherspoon, D.J., Doak, T.G., Williams, K.R., Seegmiller, A., Seger, J., and Herrick, G. (1997).
1250 Selection on the protein-coding genes of the TBE1 family of transposable elements in the
1251 ciliates *Oxytricha fallax* and *O. trifallax*. *Mol. Biol. Evol.* *14*, 696–706.
- 1252 Yang, Z. (2007). PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* *24*,
1253 1586–1591.
- 1254 Yang, Z., and Nielsen, R. (2000). Estimating synonymous and nonsynonymous substitution
1255 rates under realistic evolutionary models. *Mol. Biol. Evol.* *17*, 32–43.
- 1256 Young, D. (1937). Macronuclear reorganization in *Blepharisma undulans*. Doctoral dissertation.

- 1257 Yuan, Y.-W., and Wessler, S.R. (2011). The catalytic domain of all eukaryotic cut-and-paste
1258 transposase superfamilies. *Proc Natl Acad Sci USA* *108*, 7884–7889.
- 1259 Zagulski, M., Nowak, J.K., Le Mouël, A., Nowacki, M., Migdalski, A., Gromadka, R., Noël, B.,
1260 Blanc, I., Dessen, P., Wincker, P., et al. (2004). High coding density on the largest *Paramecium*
1261 *tetraurelia* somatic chromosome. *Curr. Biol.* *14*, 1397–1404.
- 1262 Zhang, K.S., Blauch, L.R., Huang, W., Marshall, W.F., and Tang, S.K.Y. (2021). Microfluidic
1263 guillotine reveals multiple timescales and mechanical modes of wound response in *Stentor*
1264 *coeruleus*. *BMC Biol.* *19*, 63.
- 1265 Zufall, R.A., and Katz, L.A. (2007). Micronuclear and macronuclear forms of beta-tubulin genes
1266 in the ciliate *Chilodonella uncinata* reveal insights into genome processing and protein evolution.
1267 *J. Eukaryot. Microbiol.* *54*, 275–282.
- 1268 Zufall, R.A., Sturm, M., and Mahon, B.C. (2012). Evolution of germline-limited sequences in two
1269 populations of the ciliate *Chilodonella uncinata*. *J. Mol. Evol.* *74*, 140–146.
- 1270
- 1271
- 1272

1273 Figure captions

1274 **Figure 1. *Blepharisma* nuclei and nuclear development during conjugation. A.** *B. stoltei*
1275 ATCC 30299 cell stained with anti-alpha tubulin-Alexa488 (depth-color coded red to yellow) and
1276 DAPI (cyan). **B.** Snapshot of a 3D reconstruction (Imaris, Bitplane) from CLSM images of
1277 Hoechst 33342 (dsDNA dye, Invitrogen™) fluorescence (Ex405 nm / Em420-470 nm). **C.**
1278 Schematic of nuclear processes occurring during conjugation (classified according to, and
1279 modified from (Miyake et al., 1991)). Nuclear events occurring before and up to, but not
1280 including fusion of the gametic nuclei (syngamy) are classified into sixteen pre-gamic stages
1281 where the MICs undergo meiosis and the haploid products of meiotic MICs are exchanged
1282 between the conjugating cells, followed by karyogamy. After karyogamy, cells are classified into
1283 10 stages S (synkaryon), D1 (1st mitosis), I1 (1st interphase), D2 (2nd mitosis), I2 (2nd interphase),
1284 D3 (3rd mitosis), I3 (3rd interphase), D4 (4th mitosis), E1 (1st embryonic stage), E2 (2nd embryonic
1285 stage). After E2, the exconjugants divide further and are classified into 6 stages of cell division
1286 (CD1-6) which we did not follow here. See also Figure 4.

1287
1288 **Figure 2. Basic properties of ciliate MAC genomes.** In cell diagrams MACs are green and
1289 MICs are small black dots in close proximity to MACs. Citations for genome properties are in
1290 Data S1. See also Figure S1.

1291
1292 **Figure 3. A gene-dense somatic genome with a minichromosomal architecture. A.** HiFi
1293 (DNA) and RNA-seq coverage across a representative *B. stoltei* ATCC30299 MAC genome
1294 contig (Contig_1). Y scale is linear for HiFi reads and logarithmic (base 10) for RNA-seq. Plus
1295 strand (relative to the contig) RNA-seq coverage is green; minus strand RNA-seq coverage is
1296 blue. Between the RNA-seq coverage graphs each arrow represents a predicted gene. Two
1297 orthogroups classified by OrthoFinder are shown. **B.** Mapping of a subset telomere-containing
1298 HiFi reads to a *B. stoltei* MAC genome contig region, with alternative telomere addition sites
1299 (ATASs) shown by blue (5') or mauve (3') arrows. Pink bars at read ends indicate soft-masking,
1300 typically of telomeric repeats. See also Figure S2-5.

1301
1302 **Figure 4. Developmental staging of *B. stoltei* for RNA-seq.** Classification of nuclear
1303 morphology into stages is according to previous descriptions (Miyake et al., 1991). Nuclear
1304 events occurring before and up to, but not including fusion of the gametic nuclei (syngamy) are
1305 classified into sixteen stages indicated by roman numerals. These are the pre-gamic stages of

1306 conjugation where the MICs undergo meiosis and the haploid products of meiotic MICs are
1307 exchanged between the conjugating cells. Stages after syngamy are classified into 10 stages as
1308 in Figure 1. Illustration of various cell stages adapted from (Suzuki, 1957)). Stacked bars show
1309 the proportion of cells at each time point at different stages of development, preceded by the
1310 number of cells inspected (n). See also Figure S6.

1311
1312 **Figure 5. MAC genome-encoded transposases in ciliates and properties of a putative**
1313 ***Blepharisma* IES excisase. A.** Presence/absence matrix of PFAM transposase domains
1314 detected in predicted MAC genome-encoded ciliate proteins. Ciliate classes are indicated before
1315 the binomial species names. **B.** DDE_Tnp_1_7 domain phylogeny with PFAM domain
1316 architecture and gene expression heatmap for *Blepharisma*. “Mixing” indicates when cells of the
1317 two complementary mating types were mixed. Outgroup: PiggyBac element from *Trichoplusia ni*.
1318 Catalytic residues: D- aspartate, D'- aspartate residue with 1 aa translocation. **C.** Cysteine-rich
1319 domains of PiggyBac homologs. PBLE transposases: Ago (*Aphis gossypii*); Bmo (*Bombyx*
1320 *morí*); Cag (*Ctenoplusia agnata*); Har (*Helicoverpa armigera*); Hvi (*Heliothis virescens*); PB-Tni
1321 (*Trichoplusia ni*); Mlu (PiggyBat from *Myotis lucifugus*); PLE-wu (*Spodoptera frugiperda*).
1322 Domesticated PGBD transposases: Oni (*Oreochromis niloticus*); Pny (*Pundamilia nyererei*);
1323 Lia5, Tpb1, Tpb2, Tpb6 and Tpb7 (*Tetrahymena thermophila*); Pgm, PgmL1, PgmL2,
1324 PgmL3a/b/c, PgmL4a/b, PgmL5a/b (*Paramecium tetraurelia*); Tru (*Takifugu rubripes*); Pgbd2,
1325 Pgbd3 and Pgbd4 (*Homo sapiens*).

1326
1327
1328 **Figure 6. Phylogeny of ciliate PiggyBac homologs, eukaryotic PBLEs and PGBD5**
1329 **homologs.** Highlighted clade contains all PiggyBac homologs found in Heterotrichea, containing
1330 MAC and MIC-limited homologs of PiggyMac from *Blepharisma* and PiggyMac homologs of
1331 *Condylostoma magnum*. The tree is rooted at the PiggyBac-like element of *Entamoeba*
1332 *invadens*.

1333
1334 **Figure 7. DDE_1, DDE_3 and DDE_Tnp_IS1595 domain-containing proteins in**
1335 ***Blepharisma*.** **A.** DDE_1 domain phylogeny with PFAM domain architecture and gene
1336 expression heatmap for *Blepharisma*. **B.** DDE_3 domain phylogeny with PFAM domain
1337 architecture and gene expression heatmap for *Blepharisma*. **C.** DDE_Tnp_IS1595 domain
1338 phylogeny with PFAM domain architecture and gene expression heatmap for *Blepharisma*. See
1339 also Figure S7.

1340

1341 Supplemental figure captions

1342 **Figure S1. Analysis of assembly completeness and genetic code. A.** Completeness of the
1343 *B. stoltei* ATCC 30299 MAC assembly was estimated by the percentage of BUSCOs found in
1344 the assembly with reference to the OrthoDB v10 alveolate database (Kriventseva et al., 2019).
1345 The nature of the ortholog-matches is indicated by characters followed by counts: C (complete
1346 orthologs) - light blue, D (duplicated orthologs) - dark blue, F (fragmented orthologs) - yellow
1347 and M (missing orthologs) - red. **B.** Prediction for *B. stoltei* ATCC 30299 MAC genome by
1348 PORC; codons that are stops in the standard genetic code are highlighted in orange.

1349

1350 **Figure S2. Properties of minichromosomes, telomeres, and alternative telomere addition**
1351 **sites. A.** Length distribution of telomeres of telomere-bearing HiFi reads. **B.** Length distribution
1352 of HiFi reads delimited by telomeres. **C.** Diagram of a telomere-bearing read mapped onto
1353 genome reference at an ATAS. Sequence which is ambiguously chromosomal or telomeric is
1354 “junction sequence”; junction coordinate which maximizes telomere repeat length on the read is
1355 the “first identifiable breakpoint”; the coordinate maximizing alignment length to reference is the
1356 “last identifiable breakpoint”. The last telomeric unit permutation at the last identifiable
1357 breakpoint is underlined (length 8 bp). **D.** Mean base frequencies in +/- 1 kbp flanking ATAS
1358 junctions. **E.** Sequence logos of chromosomal sequence at ATAS junctions, sorted by which
1359 permutation of the telomeric repeat is present (plot labels). Logos are aligned to the “last
1360 identifiable breakpoint” between positions 20 and 21; telomeric repeats on telomere-bearing
1361 reads begin to the left of the breakpoint. **F.** Frequencies of 2-mers in whole genome (blue), in
1362 telomeres (green), and at ATAS junctions (chromosomal side after last identifiable breakpoint,
1363 orange). **G.** Histogram of junction sequence lengths for ATASs in *B. stoltei*. **H.** Counts of each
1364 telomere repeat permutation at ATAS junctions (last identifiable breakpoint).

1365

1366 **Figure S3. Femto Pulse analyses of *B. stoltei* MAC DNA and POT1 phylogeny. A.** Mapping
1367 of PacBio CLR reads with 3 consecutive telomeric repeats to a representative *T. thermophila*
1368 MAC chromosome (Chr_001 from ciliate.org). **B.** Length distribution of input MAC DNA sizes
1369 prior to fragmentation and library preparation (Femto Pulse; LM = lower maker) - replicate 1.
1370 RFU=relative fluorescent units. **C.** Length distribution of input MAC DNA sizes prior to
1371 fragmentation and library preparation (Femto Pulse; LM = lower maker) - replicate 2. **C.** POT1

1372 paralog phylogeny, PFAM domain architecture, and gene expression in *Blepharisma*. Diagram
1373 elements as described in Figure 5B.

1374
1375 **Figure S4. Intron splicing.** **A.** Distribution of intron splicing fraction of candidate sense introns
1376 in the *B. stoltei* MAC genome. **B.** Distribution of intron splicing fractions of introns according to
1377 intron lengths. **C.** Distribution of intron splicing fraction of candidate antisense introns. **D.**
1378 Distribution of intron lengths from predicted genes. **E.** Sequence logos for 15 bp introns (splicing
1379 frequency > 0.5). **F.** Sequence logos for all predicted 16 nt introns, and 16 nt introns with “A” at
1380 either position -7 or -6 (counting from the 3’ end). The number of introns underlying the logos
1381 are indicated to the right. **G.** Distribution of intron splicing fractions of introns according to intron
1382 lengths. **H.** Sample of RNA-seq reads mapped to a GT-GG intron from gene
1383 BSTOLATCC_MAC21551 (Contig_57.g761). Translation in alternative reading frames
1384 downstream of the predicted intron leads to premature stop codons soon after the intron.

1385
1386 **Figure S5. *B. stoltei* ATCC30299 MAC genome orthogroups and assembly graph.** **A.**
1387 Clustered orthogroups (Data S2) in the *B. stoltei* MAC genome. **B.** Bandage (Wick et al., 2015)
1388 representation of Flye 2.8.1 assembly graph. Edges corresponding to contigs are colored by
1389 coverage (brightest pink = 160x, black=0x).

1390
1391 **Figure S6. Experimental approach for conjugation RNA-seq time series.** Complementary
1392 mating type strains of *Blepharisma stoltei* were harvested and cleaned by starving overnight.
1393 The cleaned cultures were treated in a time-staggered format, with gamones of the
1394 complementary mating type, where gamone 2 was a solution of the synthetic gamone 2 calcium
1395 salt and gamone 1 was provided as the cell-free fluid (CFF) harvested from mating-type I cells.
1396 Two sets of time-staggered gamone-treated cultures were used for the time series. Set I,
1397 indicated by the solid line, was mixed and used to observe and collect samples at 0 hours, 2
1398 hours, 6 hours, 26 hours and 30 hours after mixing. Set II, indicated by the dashed lines, was
1399 mixed and used to observe and collect samples at 14 hours, 18 hours, 22 hours and 38 hours
1400 after mixing. Test tubes indicate Trizol samples prepared for RNA-extraction which were stored
1401 at -80 °C before processing. Cells collected for imaging were obtained shortly before the
1402 remainder were transferred into Trizol.

1403
1404 **Figure S7. MULE domain transposases in *Blepharisma*.** MULE domain phylogeny with
1405 PFAM domain architecture and gene expression heatmap for *Blepharisma*.

1406
1407 **Figure S8. Small RNA-related proteins in *Blepharisma*.** A. ResIII, Helicase_c and
1408 Ribonuclease_3 domain phylogeny with PFAM domain architecture and gene expression
1409 heatmap for *Blepharisma*. B. PIWI domain phylogeny with PFAM domain architecture and gene
1410 expression heatmap for *B. stoltei*.
1411
1412 **Figure S9. Histones and histone-domain-containing proteins in *Blepharisma*.** Gene
1413 expression heatmaps are shown as in previous figures, are clustered according to major histone
1414 type as classified using HistoneDB domain models. Domains from PFAM and HistoneDB are
1415 shown to the right.

Figure 1

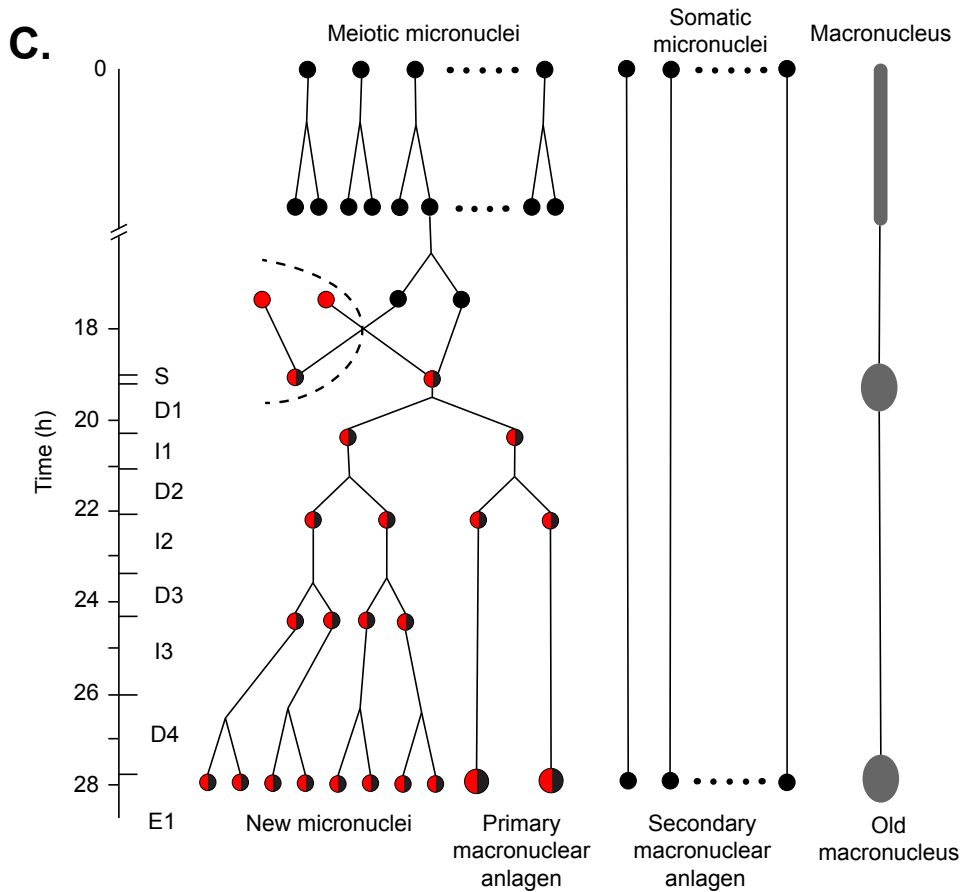
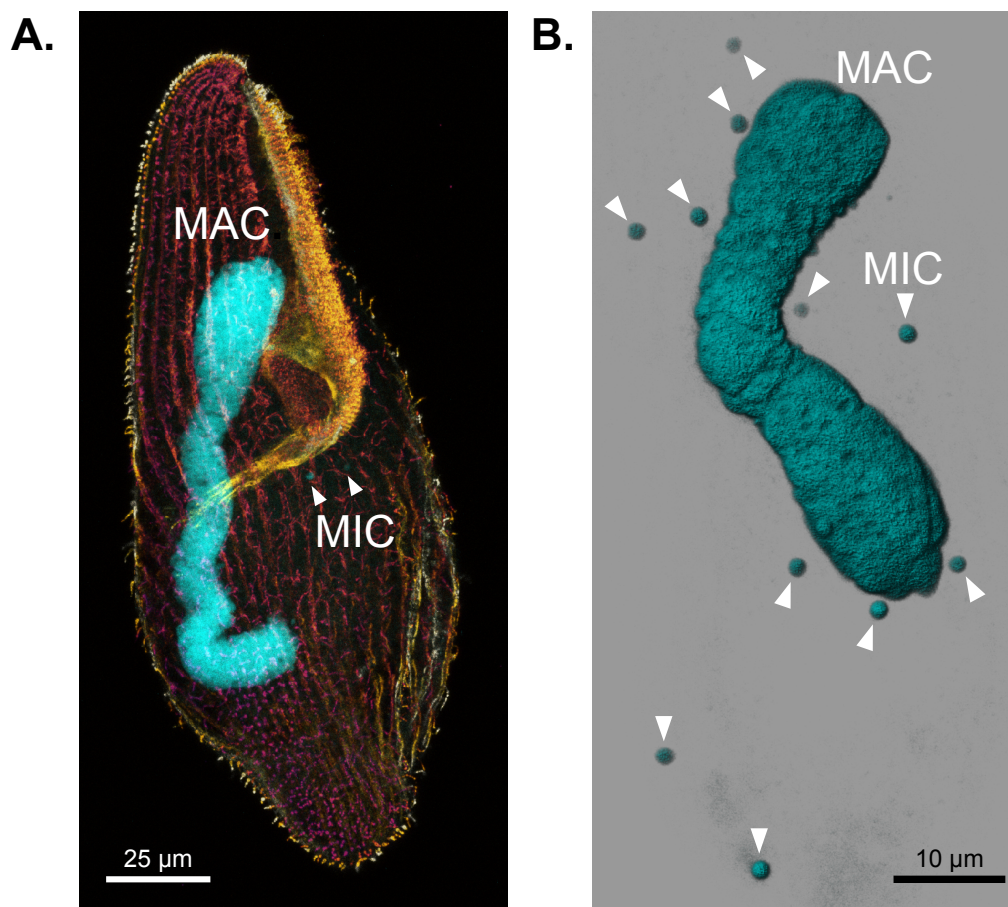
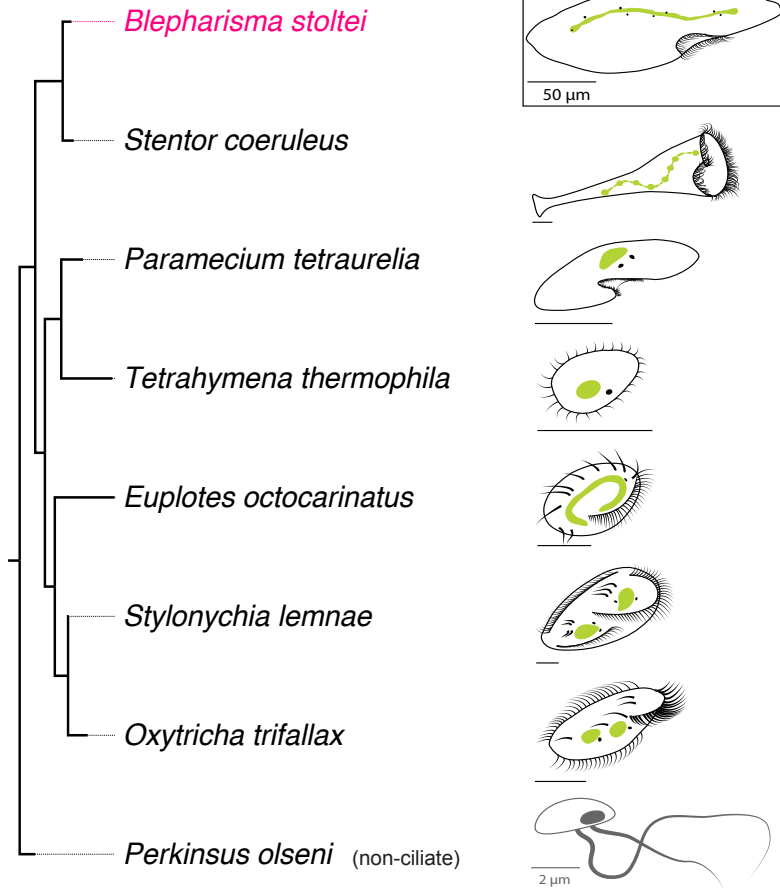


Figure 2.



Species	Genome size (~Mb)	Genome architecture	Genes (zygosity)	Codon reassignments		
				UAA	UAG	UGA
<i>Blepharisma stoltei</i> (ATCC 302099)	41	Mini-chromosomes	25726 (n)	*	*	W
<i>Stentor coeruleus</i>	77	?	31426 (?)	*	*	*
<i>Paramecium tetraurelia</i>	72	Chromosomes	39642 (n)	Q	Q	*
<i>Tetrahymena thermophila</i>	103	Chromosomes	26258 (n)	Q	Q	*
<i>Euplotes octocarinatus</i>	88	Nano-chromosomes	29076 (n)	*	*	C
<i>Stylonychia lemnae</i>	52	Nano-chromosomes	15102 (n)	Q	Q	*
<i>Oxytricha trifallax</i>	50	Nano-chromosomes	18400 (n)	Q	Q	*
<i>Perkinsus olseni</i>	63	Chromosomes	17342 (4n)	*	*	*

Figure 3.

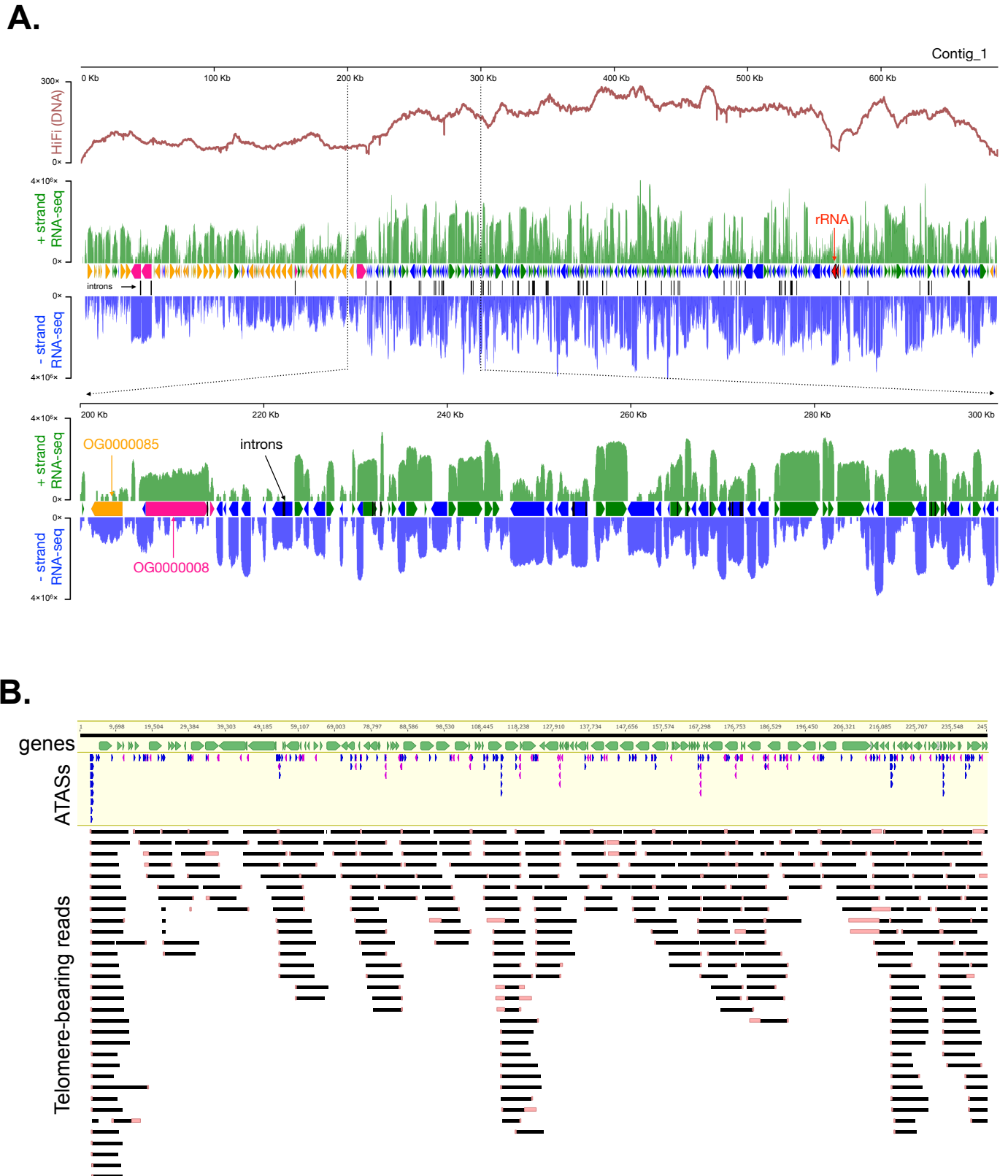


Figure 4

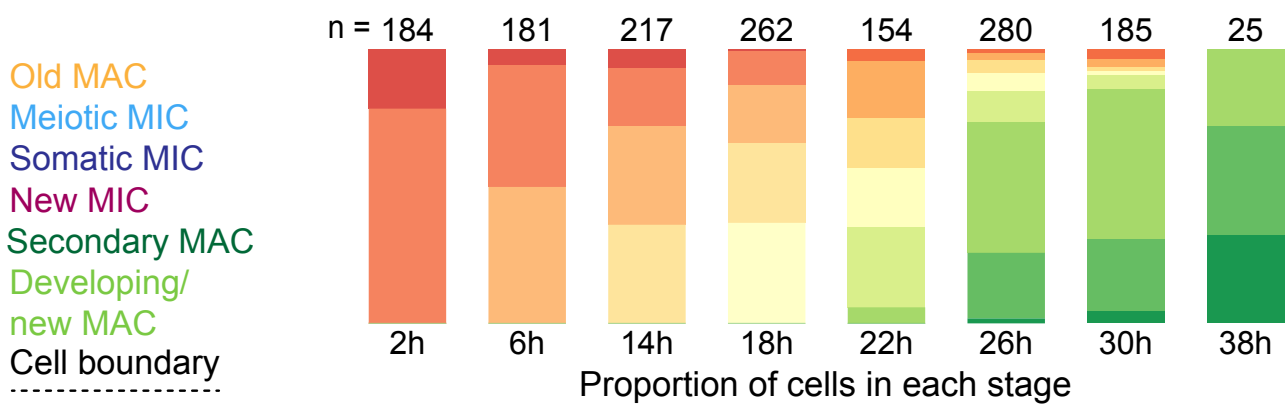
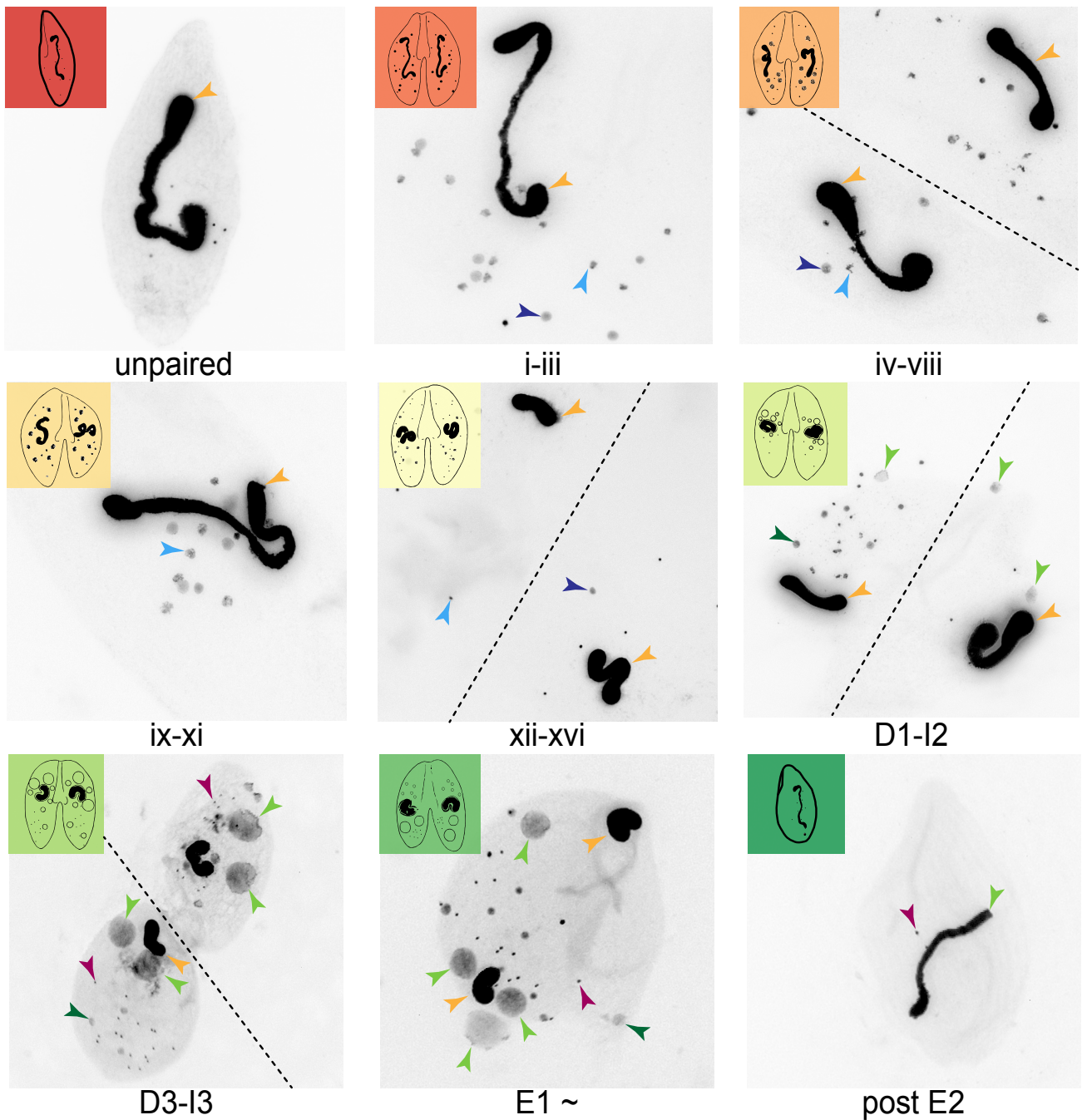


Figure 5.

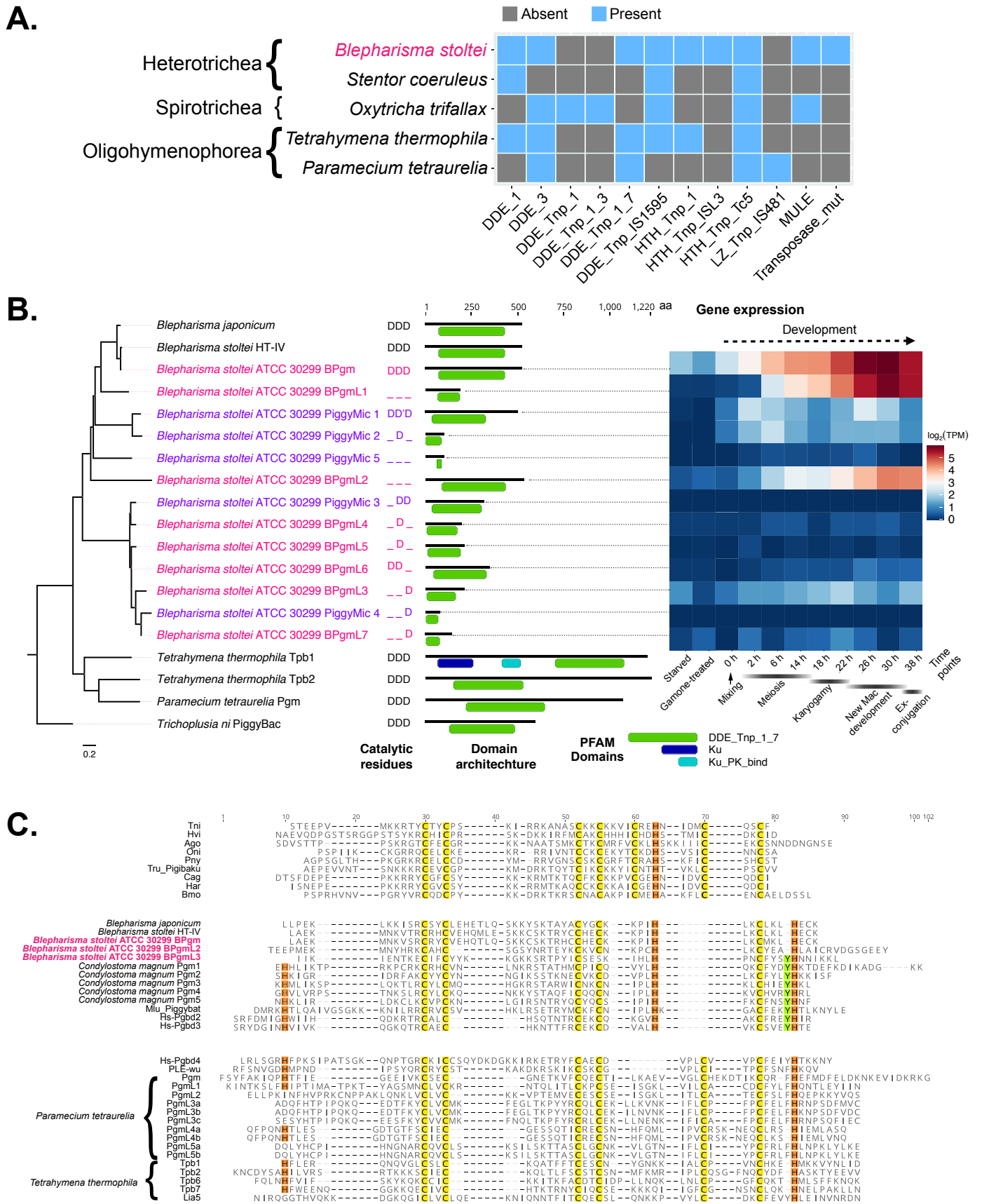
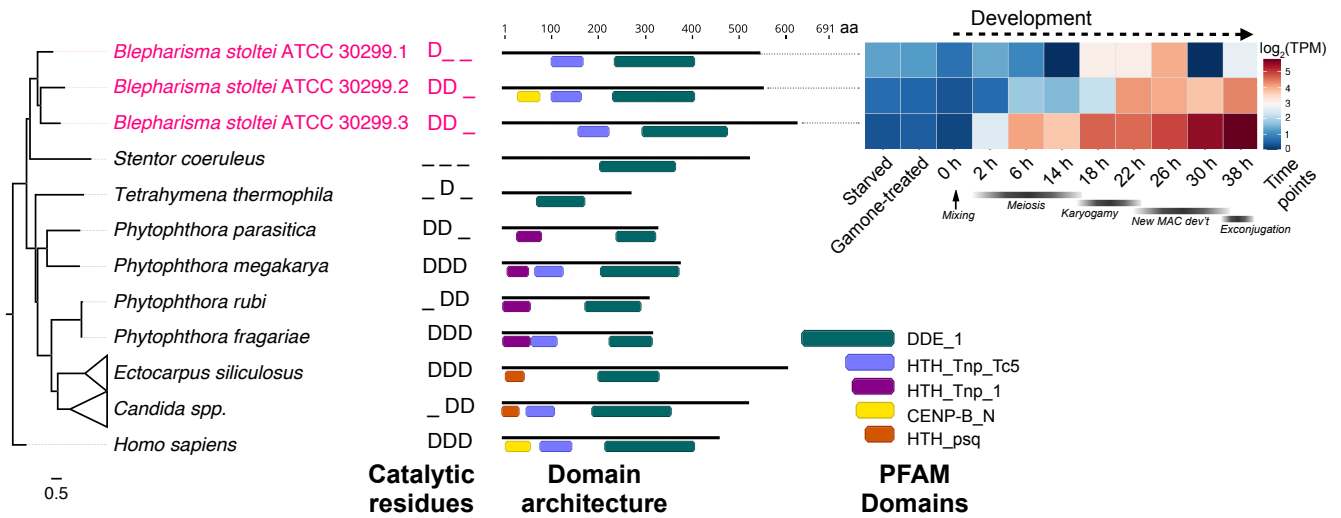
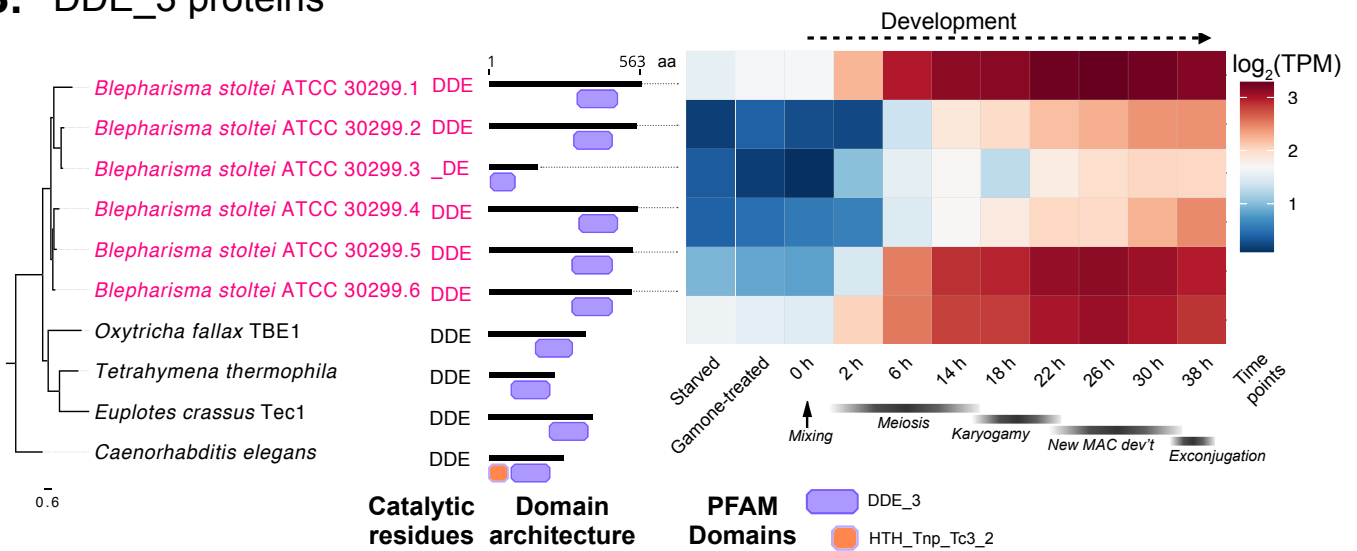


Figure 7

A. DDE_1 proteins



B. DDE_3 proteins



C. DDE_Tnp_IS1595 proteins

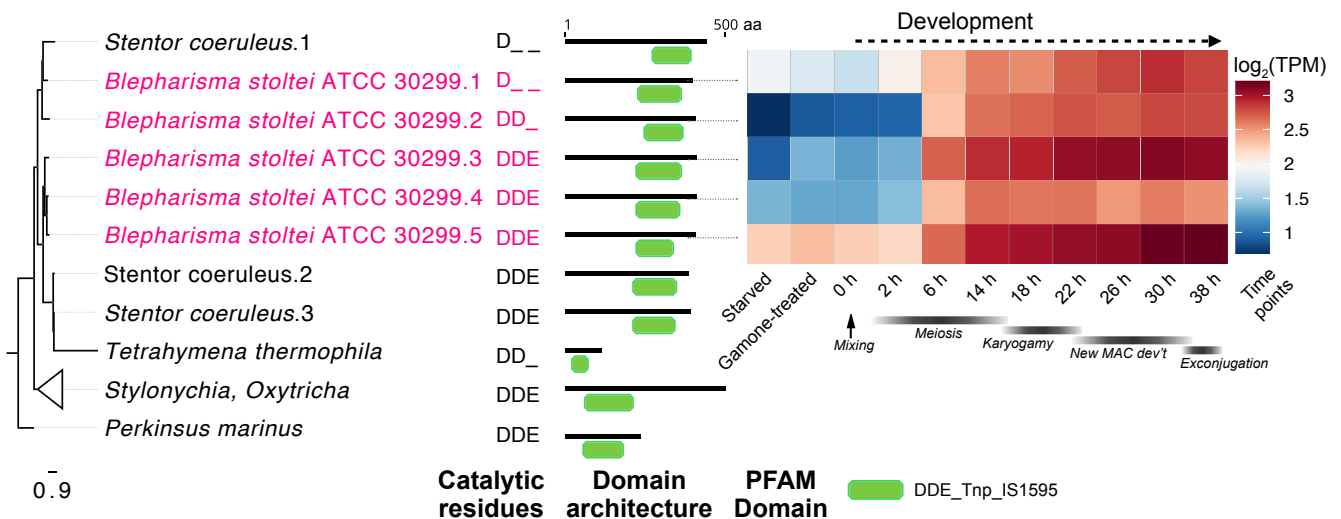
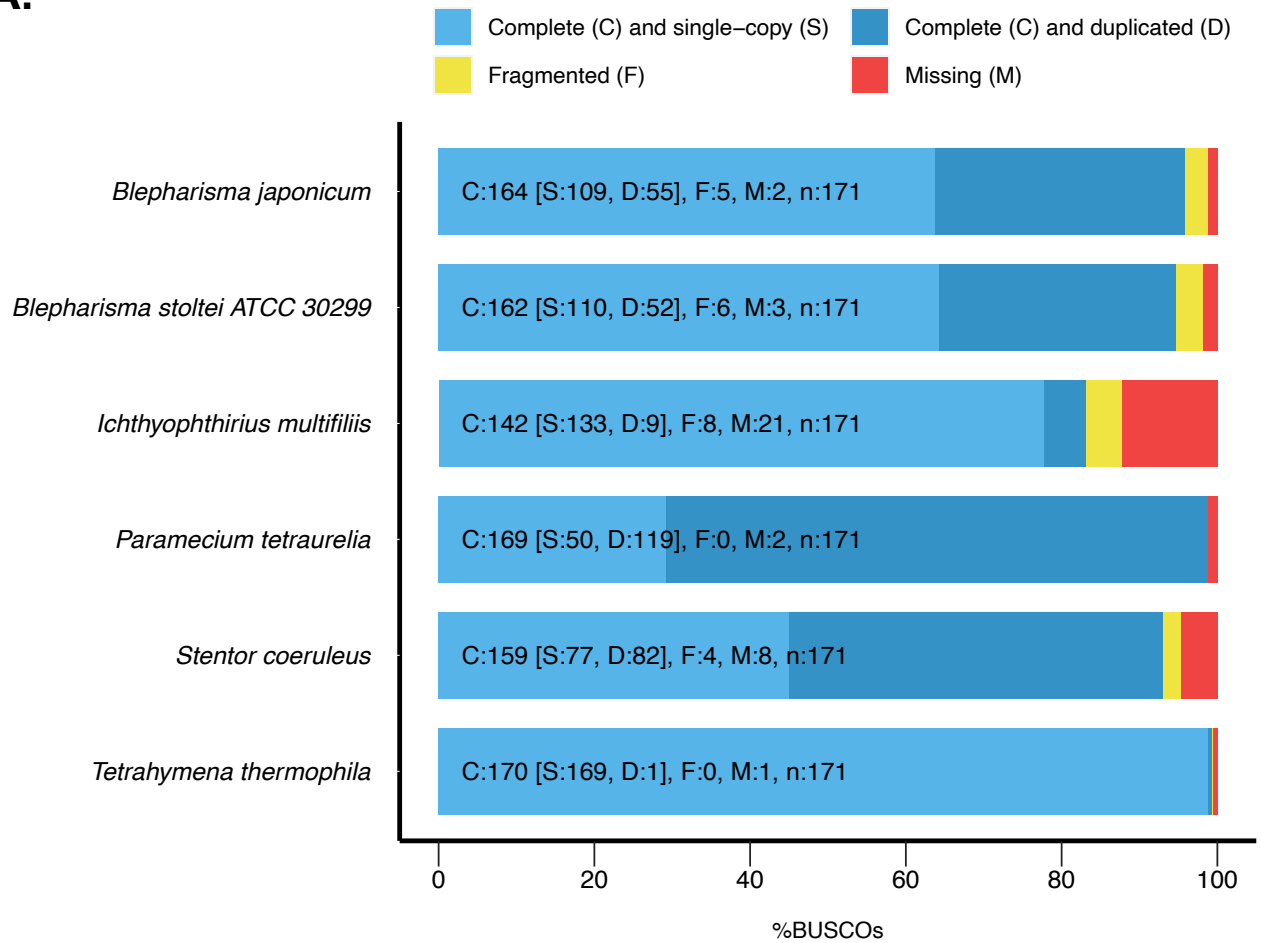


Figure S1

A.



B.

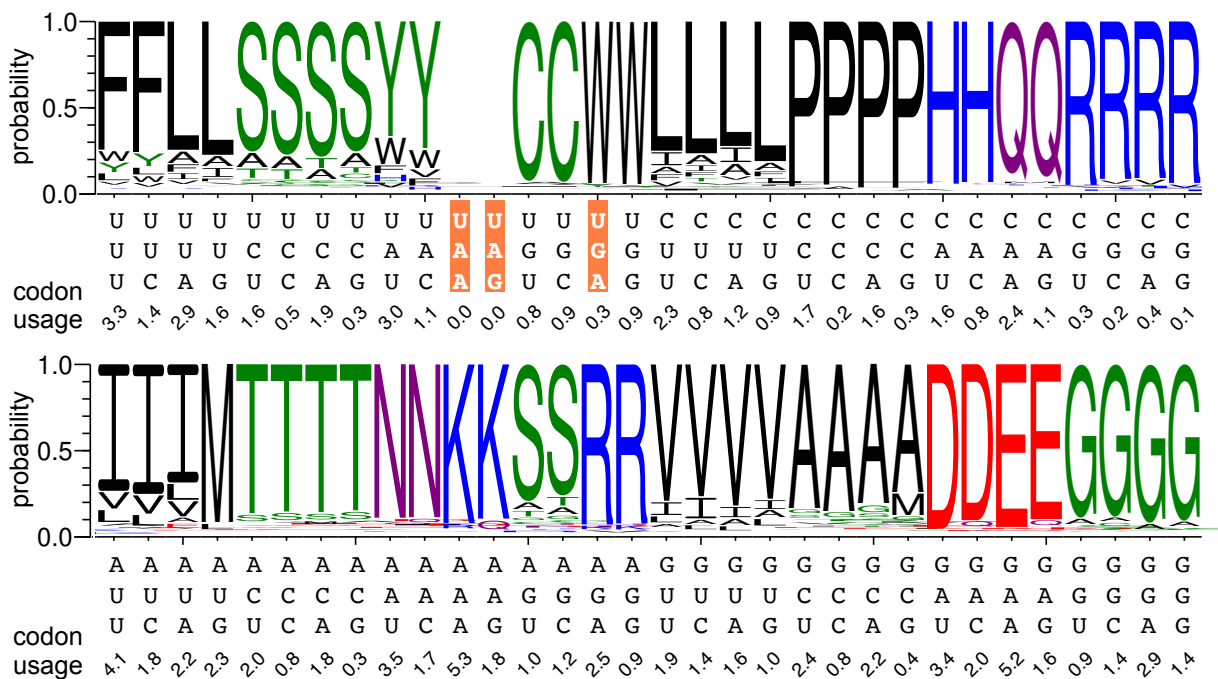


Figure S2

bioRxiv preprint doi: <https://doi.org/10.1101/2021.12.14.471807>; this version posted April 27, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

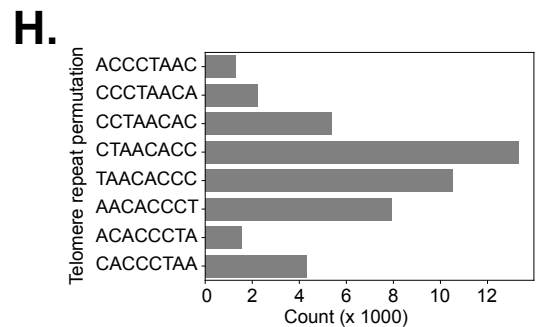
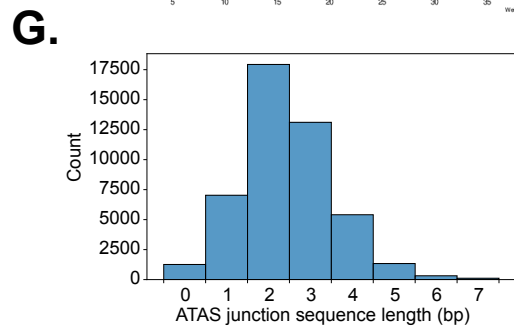
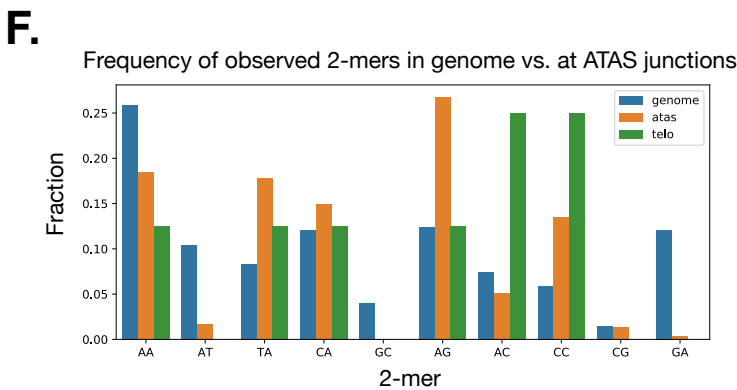
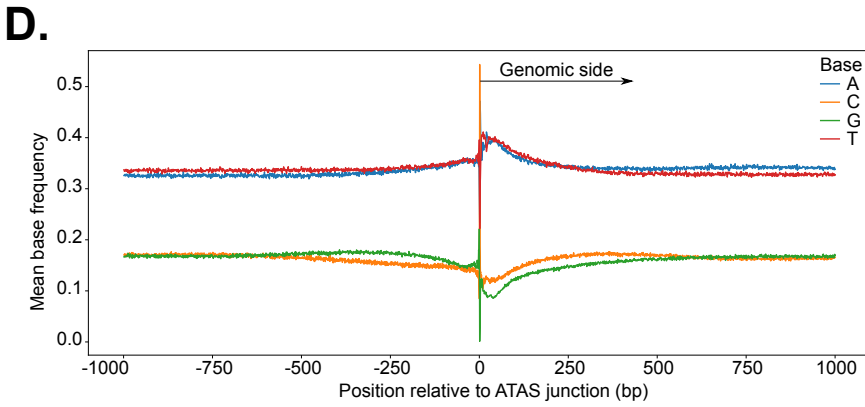
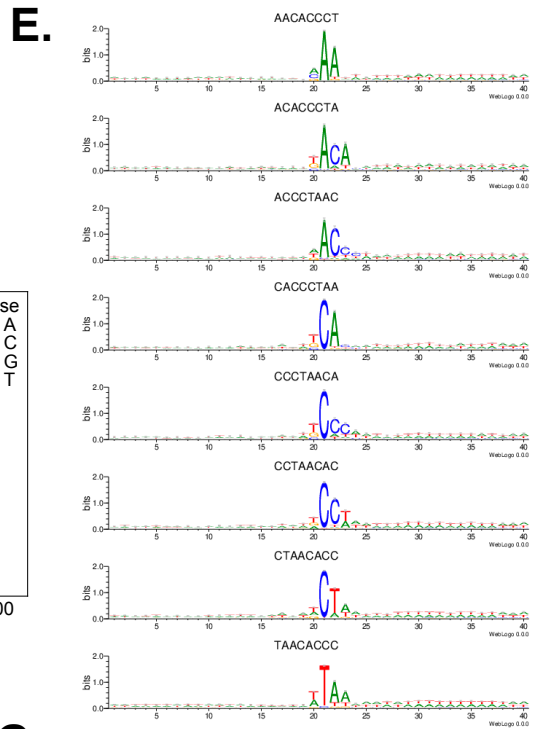
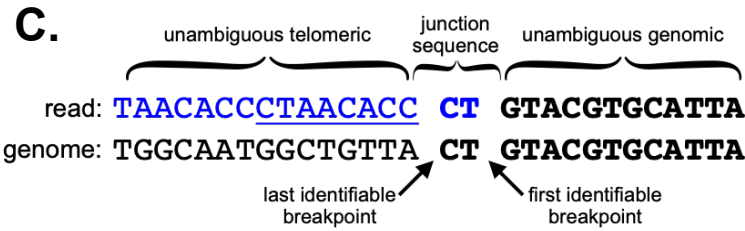
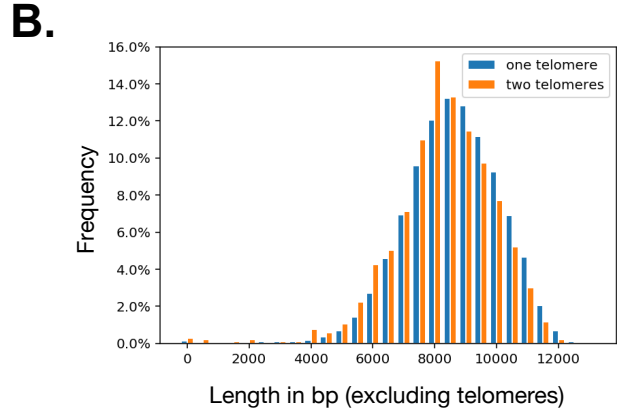
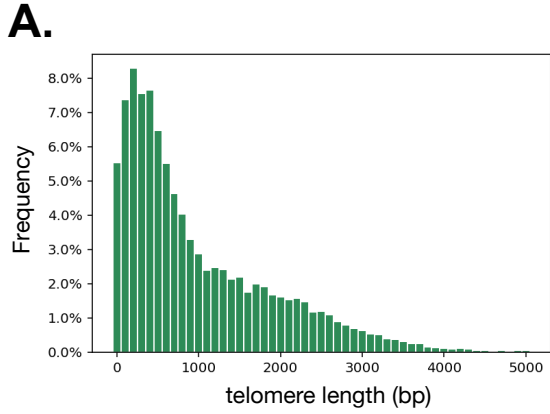


Figure S3

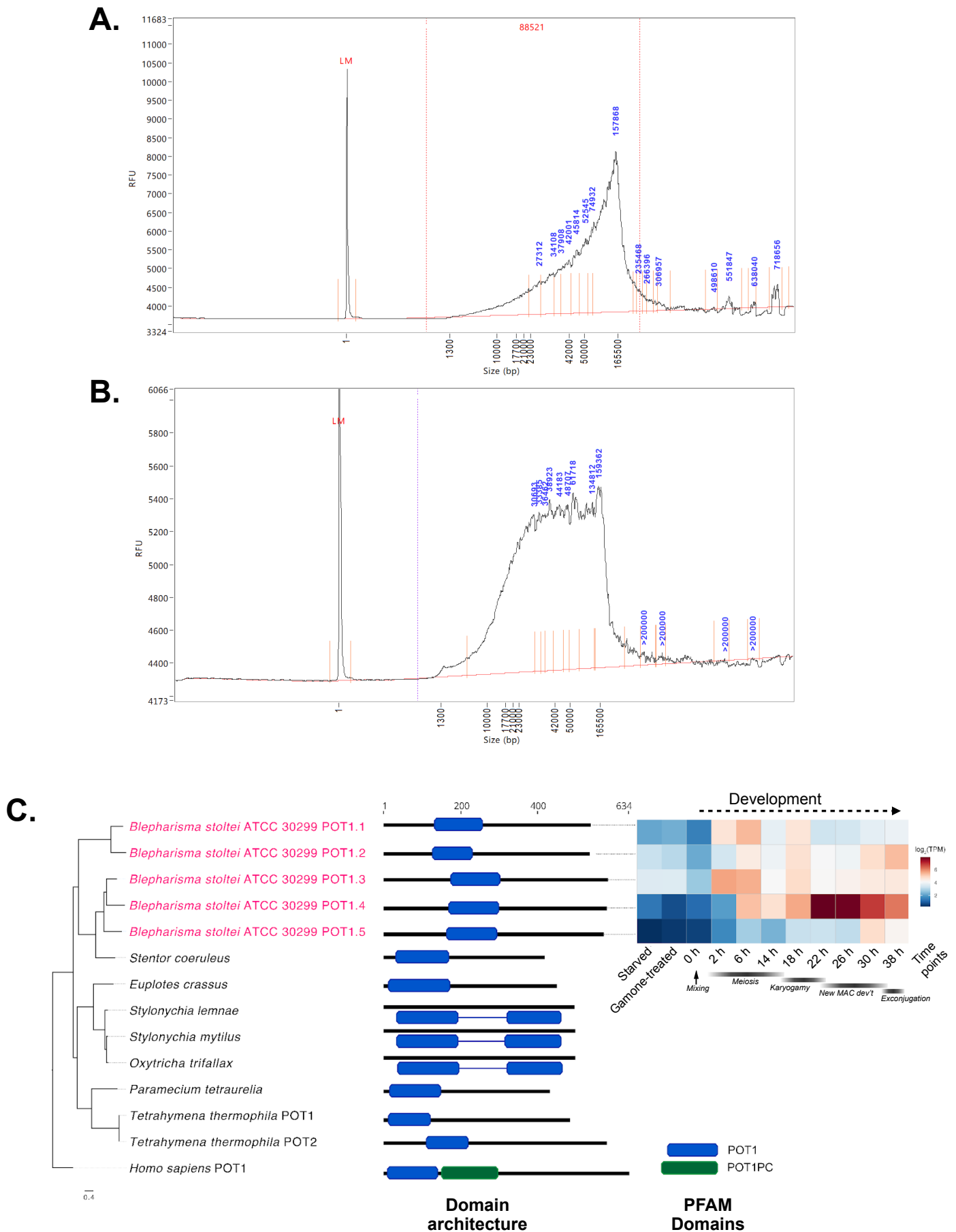


Figure S4

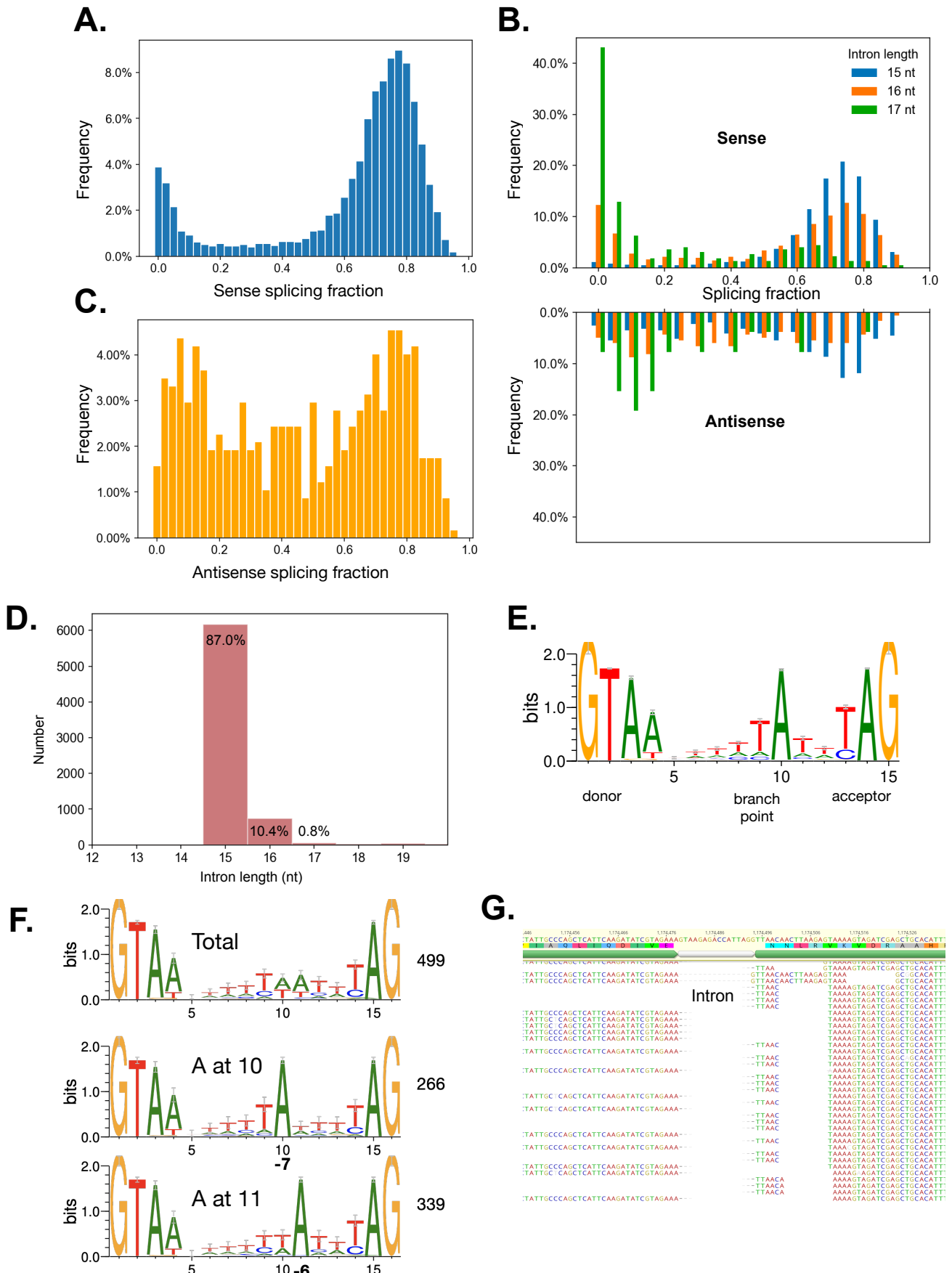
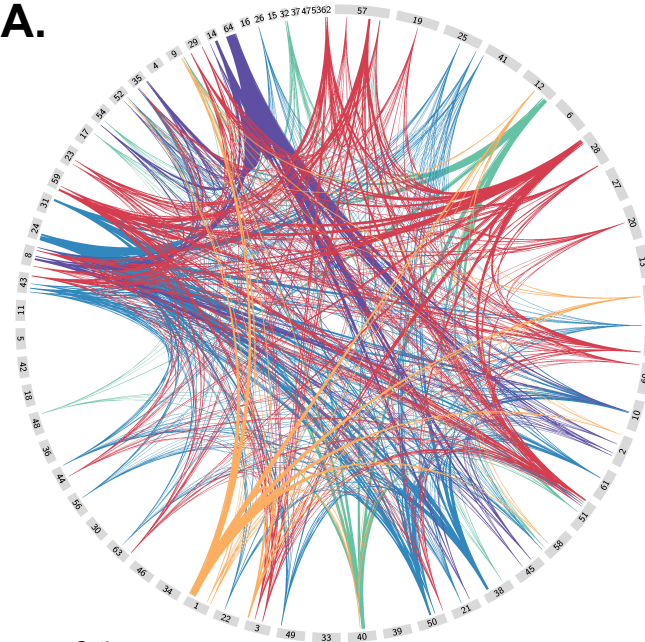


Figure S5

A.



Orthogroups:

OG0000085 OG0000018 OG0000019
OG0000014 OG0000052

B.

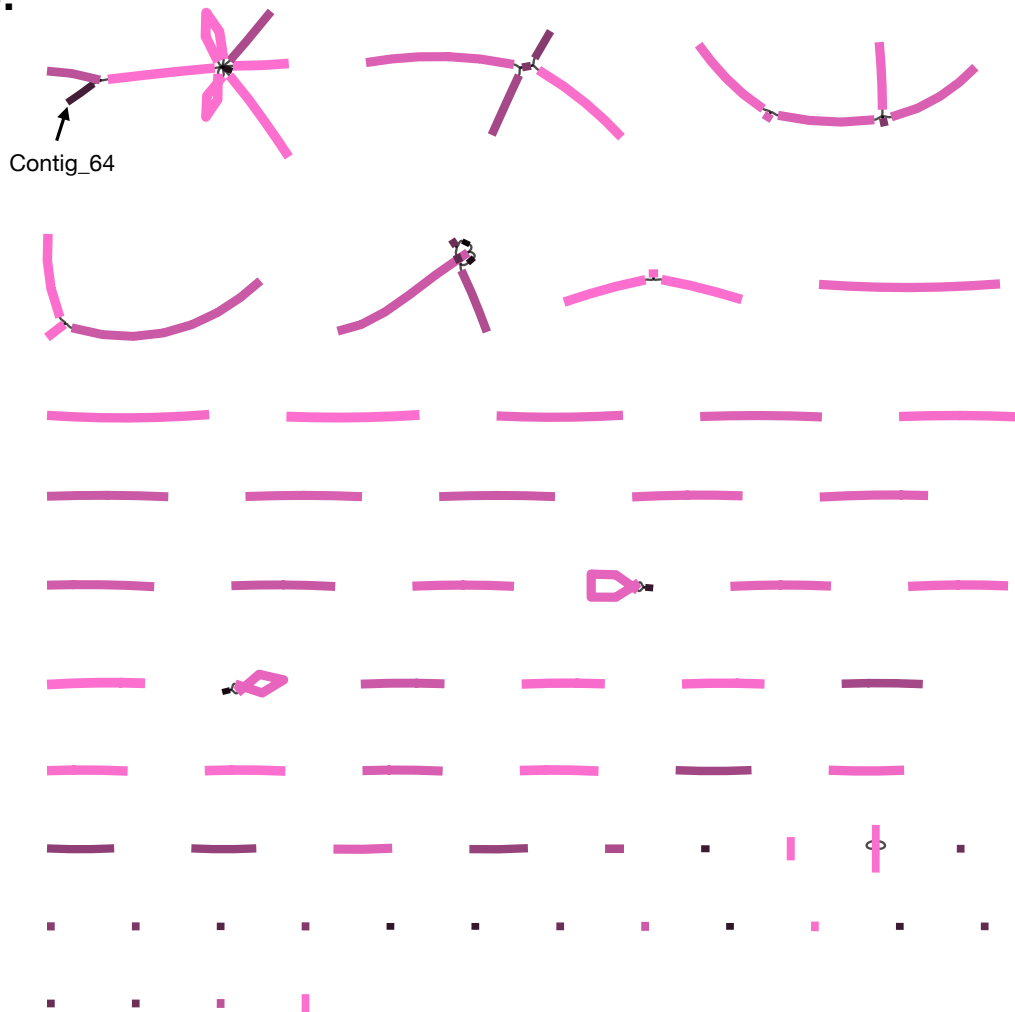


Figure S6

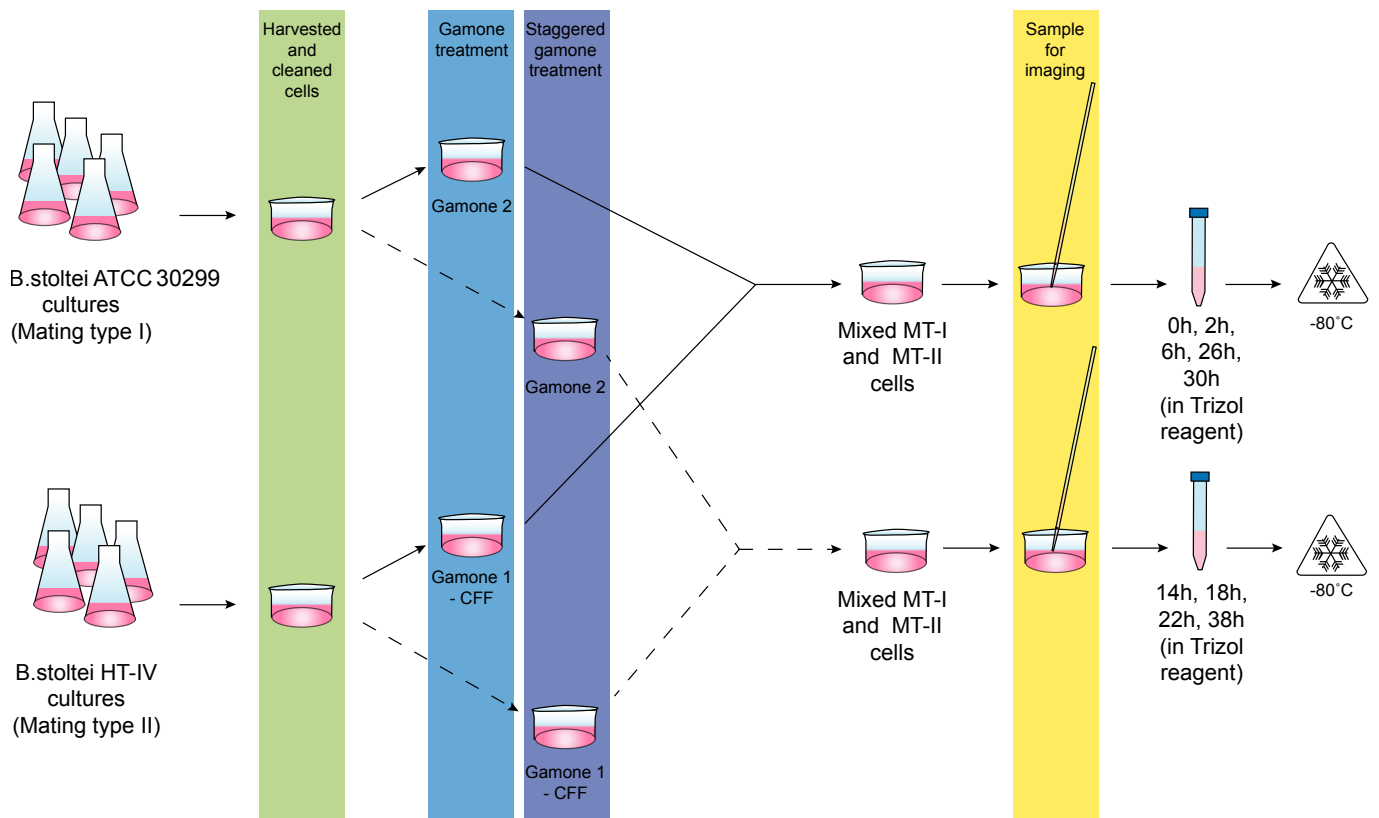


Figure S7

MULE domain-containing proteins

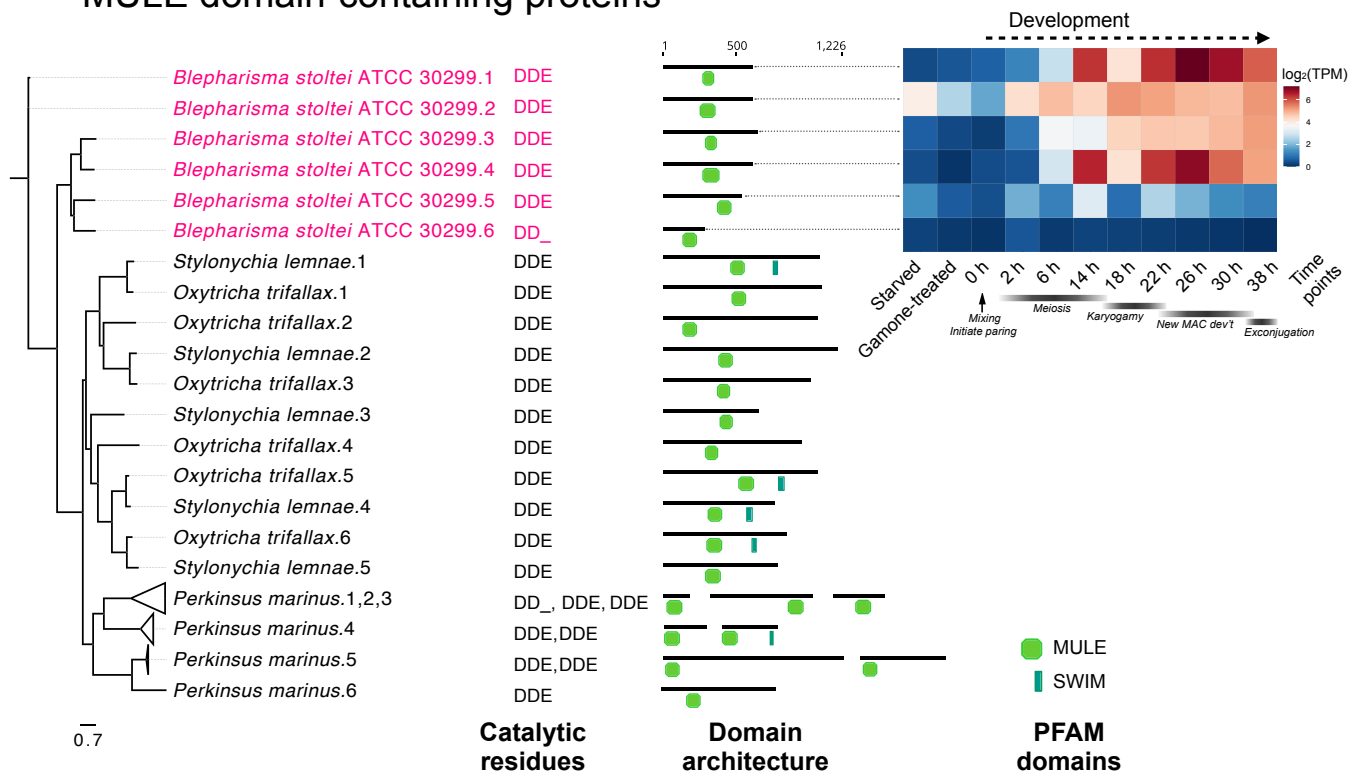
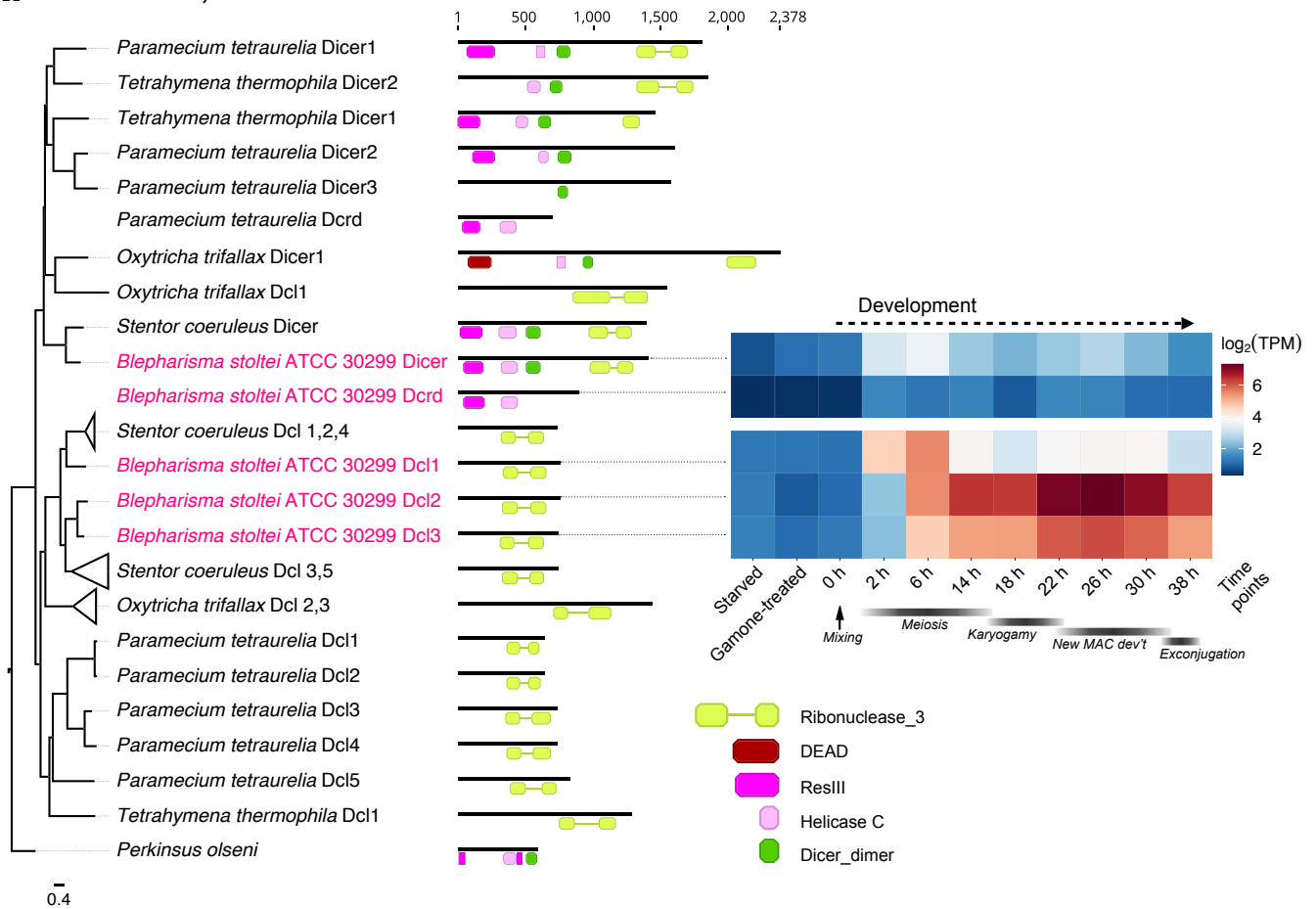


Figure S8

A. Dicers, Dicer-likes and Dicer-derivatives



B. Piwi proteins

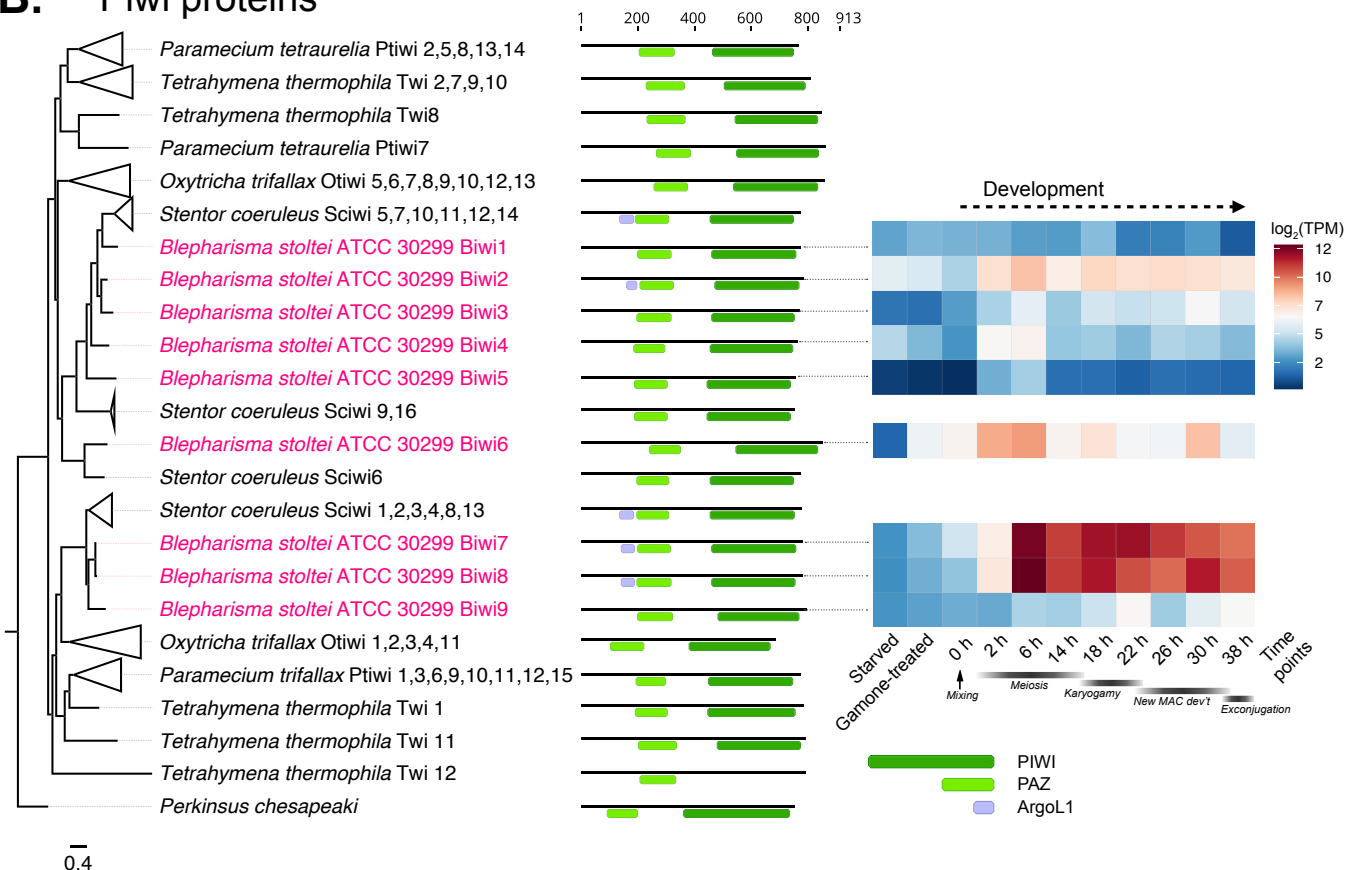


Figure S9

