# MITE infestation of germline accommodated by genome editing in *Blepharisma*

Brandon Kwee Boon Seah[1], Minakshi Singh[1], Christiane Emmerich[1], Aditi Singh[1], Christian Woehle[2], Bruno Huettel[2], Adam Byerly[3], Naomi Stover[4], Mayumi Sugiura[5], Terue Harumoto[5], Estienne Carl Swart[1,*]

[1] Max Planck Institute for Biology, Max-Planck-Ring 5, 72076 Tübingen, Germany

[2] Max Planck Genome Center Cologne, Max Planck Institute for Plant Breeding, Building B, Carl-von-Linné-Weg 10, 50829 Cologne, Germany

[3] Department of Computer Science and Information Systems, Bradley University, Peoria IL, USA

[4] Department of Biology, Bradley University, Peoria IL, USA

[5] Department of Chemistry, Biology, and Environmental Sciences, Faculty of Science, Nara Women's University, Nara 630-8506, Japan

* Correspondence: estienne.swart@tuebingen.mpg.de

## Summary

During a sophisticated developmental process, ciliates excise numerous internally eliminated sequences (IESs) from a germline genome copy, producing a functional somatic genome. Most IESs ultimately originate from transposons but homology is obscured by sequence decay. To obtain more representative perspectives on ciliate genome editing, we assembled forty thousand IESs of *Blepharisma stoltei*, from a much earlier-diverging lineage than existing models. Short IESs (< 115 bp) were largely non-repetitive, with a pronounced ~10 bp length periodicity, whereas longer IESs (max 7 kbp) were non-periodic and contained abundant interspersed repeats. Contrary to current models, the *Blepharisma* germline genome encodes few transposases. Instead, its most abundant repeat (8000 copies) was a Miniature Inverted-repeat Transposable Element (MITE), apparently a deletion derivative of a germline-limited Pogo-family transposon. We propose MITEs as an important and eventually self-limiting IES source. Rather than defending germline genomes against mobile elements, we argue that transposase domestication actually facilitates junk DNA accumulation.

## Keywords

micronucleus, macronucleus, DNA elimination, mobile element, selfish gene, nuclear dualism, sRNA, Ciliophora, protist

## Abbreviations

- IES - internally eliminated sequence
- LTR - long terminal repeat
- MAC - macronucleus
- MIC - micronucleus
- MITE - miniature inverted-repeat transposable element
- MITIES - miniature inverted-repeat transposable internally eliminated sequences
- TDR - terminal direct repeat
- TIR - terminal inverted repeat
- TSD - target site duplication

## Introduction

46

47 Ciliates are microbial eukaryotes that maintain separate germline and somatic genomes in

48 each cell, housed in two types of nuclei. During the sexual life cycle, germline micronuclei

49 (MICs) develop via a process of small RNA (sRNA)-assisted DNA elimination and DNA

50 amplification into new somatic macronuclei (MACs), which are the site of most gene

51 expression in vegetative cells. Germline-limited genome segments, called internally

52 eliminated sequences (IESs), are excised during development from MIC to MAC. The MAC

53 genome content is hence a subset of the germline MIC. Each of the few taxa studied so far

54 has its own peculiarities. For example, typical *Paramecium* IESs are short, have unique

55 sequence content, and are precisely excised, while *Tetrahymena* IESs are longer, more

56 repetitive, and imprecisely excised (Arnaiz et al., 2012; Feng et al., 2017; Hamilton et al.,

57 2016).

58 Ciliate IESs are thought to originate from cut-and-paste DNA transposons (Klobutcher and

59 Herrick, 1997) (Figure 1B), because: (i) 5'-TA-3' motifs at IES boundaries (*Euplotes*,

60 *Paramecium*) resemble the terminal direct repeats of Tc1/Mariner-superfamily transposons

61 (Klobutcher and Herrick, 1995); (ii) transposon-derived "domesticated" excisases are used to

62 remove IESs (Baudry et al., 2009; Cheng et al., 2010; Nowacki et al., 2009); and (iii) intact

63 transposons encoding transposases are mostly germline-limited (Arnaiz et al., 2012; Herrick

64 et al., 1985; Jahn et al., 1993; Le Mouël et al., 2003). Recently, IESs with non-autonomous

65 mobile elements that resemble miniature inverted-repeat transposable elements (MITEs)

66 have been reported in *Paramecium* spp. (Sellis et al., 2021). MITEs are deletion derivatives

67 of Tc1/Mariner transposons, generally short (<500 bp), lacking coding sequences, bounded

68 by terminal repeats, and are common in plants and animals (Feschotte et al., 2002).

69 However, the autonomous counterparts of most *Paramecium* putative MITEs, including the

70 most abundant ones with thousands of copies, have not been identified.

71 Developmental DNA elimination has been called "genome defense" because the process

72 removes IESs, which not only derive from selfish genetic elements (transposons), but are

73 often intragenic and hence deleterious if not removed (Yao et al., 2003). The "defense"

74 analogy was popularized due to parallels to other eukaryotes where small RNA-mediated

75 DNA heterochromatinization is thought to suppress mobile element proliferation (Coyne et

76 al., 2012; Grewal and Jia, 2007; Vogt and Mochizuki, 2013). Ciliates have been proposed to

77 use development-specific sRNAs to guide DNA elimination; in oligohymenophoreans, they

78 mark sequences for elimination (Mochizuki et al., 2002; Sandoval et al., 2014; Yao et al.,

79 2003), whereas spirotrich sRNAs mark sequences to be retained (Fang et al., 2012; Zahler

80 et al., 2012). Histone modifications are also required for elimination (Liu et al., 2007; Taverna

81  et al., 2002). sRNAs may not always be strictly necessary: in *Paramecium*, knockdown of

82  key sRNA biogenesis enzymes had a smaller effect on shorter IESs, and were only weakly

83  correlated with the more potent effects of knocking down the main IES excisase (Sandoval

84  et al., 2014; Swart et al., 2014).

85  Other phenomena during genome editing vary markedly between the few model species

86  studied in detail (reviews: (Chalker et al., 2013; Coyne et al., 2012; Rzeszutek et al., 2020)).

87  For example, in all species, germline chromosomes are fragmented into smaller somatic

88  ones to some degree, but spirotrichs produce extremely short somatic "nanochromosomes"

89  with only one or a few genes. "Unscrambling" of nonsequential MAC-destined sequences

90  into the correct order in the somatic genome occurs frequently in some spirotrichs, e.g.

91  *Oxytricha* and *Stylonychia* (Prescott and Greslin, 1992), infrequently in *Tetrahymena*

92  (Hamilton et al., 2016), and has not been reported in other ciliates (e.g. *Paramecium and*

93  *Euplotes*). Draft-quality germline genomes are available from only two out of eleven class-

94  level taxa (following taxonomy of Lynn, 2010): Oligohymenophorea (Arnaiz et al., 2012;

95  Guérin et al., 2017; Hamilton et al., 2016; Sellis et al., 2021) and Spirotrichea (Chen et al.,

96  2014) (Figure 1C).

97  Since it is not apparent which genome editing elements are common to all ciliates, we

98  targeted the heterotrich *Blepharisma stoltei* (class Heterotrichea), whose last common

99  ancestor with other ciliates with sequenced germline genomes is the last common ancestor

100  of all ciliates (Gao and Katz, 2014). *Blepharisma* has been a laboratory model for

101  photobiology (Giese, 1973) and mating factors (Kubota et al., 1973; Miyake and Beyer,

102  1974; Miyake et al., 1991; Sugiura and Harumoto, 2001), so cultivated strains and protocols

103  for inducing conjugation and development are available, and now too an accurate, highly

104  contiguous draft somatic genome (Singh et al., 2021). The somatic genome encodes a likely

105  IES excisase, *Blepharisma* PiggyMac (BPgm), most closely related to the main IES

106  excisases of *Paramecium* (PiggyMac) and *Tetrahymena* (Tpb2). Other somatic PiggyBac

107  paralogs are also present but lack a complete "catalytic triad", similar to the situation in

108  *Paramecium* (Bischerour et al., 2018). BPgm is upregulated during formation of the new

109  somatic MAC alongside other development-specific genes, including homologs of sRNA

110  biogenesis proteins implicated in genome editing (Singh et al., 2021).

111

112  In this study, we assembled a draft germline genome for *Blepharisma stoltei*. Through single

113  molecule long read sequencing and targeted assembly, we assembled IESs including many

114  with long, repetitive elements, which is not feasible with short read shotgun sequencing

115  alone. We found about ten thousand short (≤115 bp), precisely excised IESs with a periodic

4

116    length distribution like *Paramecium*'s. However most IESs (about thirty thousand) were

117    longer, up to several kbp, and, importantly, also include a Tc1/Mariner transposon whose

118    non-autonomous MITE was also the most abundant repeat in the genome. Complementing

119    the genomic analyses, we also identified small RNAs expressed during sexual development

120    with characteristics of scnRNAs that guide DNA elimination in other ciliates. These results

121    show that characteristics of germline-limited DNA in ciliates may be disjunct to phylogeny,

122    and also illustrate how MITEs could be an intermediate stage in the origin and proliferation of

123    IESs.

124 **Results**

125 *Detection and targeted assembly of ca. forty thousand germline-limited IESs*

126 To investigate the *Blepharisma* germline genome we enriched germline micronuclei from *B.*
127 *stoltei* strain ATCC 30299, and reconstructed 39799 IESs (13.2 Mbp total, average coverage
128 ~45x) scaffolded on the previously assembled 41 Mbp somatic genome (Singh et al., 2021)
129 using a mapping and targeted assembly approach for PacBio long reads (Seah and Swart,
130 2021). This MAC-scaffolded germline assembly is here referred to as the "MAC+IES"
131 assembly. About 70% of all predicted IESs were intragenic (within coding sequences or
132 introns), implying precise excision of IESs, as they would otherwise cause deleterious
133 translation frameshifts. However, genes occupied 77% of the somatic assembly (excluding
134 telomeres), so there was a small but statistically significant ($p = 3 \times 10^{-269}$) relative depletion
135 of intragenic IESs.

136 *A "hybrid" IES length distribution with periodic length peaks for short IESs*

137 Most IESs were short (median 255 bp, mean 331 bp), but the distribution was long-tailed
138 (90th percentile 603 bp, max 7251 bp). The length distribution was not unimodal, but had
139 multiple peaks at specific length values (Figure 1A, Table S1). It appeared to be a "hybrid"
140 distribution composed of two ranges: a "periodic" range, from ~65 to 115 bp (10778 IESs),
141 and a "non-periodic" range, >115 bp (29021).

142 The "periodic" IES size range contained sharp peaks every 10 to 11 bp, similar to the
143 periodicity of IESs in *Paramecium tetraurelia* (Arnaiz et al., 2012; Guérin et al., 2017). The
144 first peak in *B. stoltei* was centered at 65 bp, compared to 28 bp in *P. tetraurelia*, and there
145 was no "forbidden" peak. The most abundant "periodic" length peaks were at 72 bp and 110
146 bp. The "non-periodic" range (≥115 bp) contained isolated peaks at 153, 174, 228, and 389
147 bp, which has no obvious periodicity. Only 9701 IESs (total 1.36 Mbp) were contained within
148 the size classes represented by the above peaks (both periodic and non-periodic) (Table
149 S1), meaning that most IESs had lengths outside the peak values.

150 *IESs are bounded by heterogeneous direct and inverted terminal repeats*

151 In other ciliates, IES boundaries often have conserved terminal repeat motifs that could
152 reflect excisase cut site preferences or IES origins from specific classes of transposons
153 (Klobutcher and Herrick, 1997). We therefore searched for both direct and inverted terminal
154 repeats in *Blepharisma* IESs.

155    About three quarters of IESs (30212, 9.43 Mbp) were bounded by terminal direct repeats

156    (TDRs) that contained the subsequence TA ("TA-bound"). Other non-TA TDRs accounted for

157    another 6566 (2.85 Mbp); the remainder were not TDR-bound, though some may be

158    assembly errors (Figure 1A). Like most ciliates, *B. stoltei* genomes were AT-rich (somatic

159    33.5% GC, IESs 33.3% GC) but the number of TA- and TDR-bound sequences was unlikely

160    to be due to nucleotide composition alone (Figure 2A, 2B). The most common TDRs were

161    simple alternations of T and A (TA, TAT/ATA, TATA), especially in IESs up to 228 bp (Figure

162    2C), with the exception of TAA/TTA (see below).

163    Erroneous, low-frequency excision of MAC-destined sequences (MDSs) by the excision

164    machinery ("cryptic" IESs) was also detected in MAC DNA libraries, with a slight peak at 72

165    bp (Figure S1C). Of 10048 cryptic IESs, 56% were TA-bound; TAA/TTA-bound IESs were

166    also common, which suggests that the observed TDRs, including TAA/TTA, represented

167    intrinsic cut site preferences of the domesticated excisase(s) (Figure S1C to F).

168    Terminal inverted repeats (TIRs) at IES junctions were heterogeneous among IES size

169    classes (Figure 1D, Figure 2F), and no single TIR motif was generally conserved across all

170    *Blepharisma* IESs, unlike the common 5'-TAYNR-3' motif of *Paramecium* IESs. Considering

171    only TA-bound IESs, boundaries of "periodic" IESs had a weak consensus 5'-<u>TA</u>T rrn ttt t-3'

172    (weakly conserved bases in lowercase), whereas IES from "non-periodic" peaks had other

173    signatures, e.g. 5'-<u>TA</u>T Agn nnT TT-3' for both ~153 and ~174 bp IESs. Despite their

174    heterogeneity, TIRs were more common and longer than expected by chance, even with a

175    strict criterion of no gaps or mismatches (Figure 2D to F). Sequence clustering of long (≥10

176    bp) TIRs showed distinct TIRs associated with specific IES lengths. Additionally, 376

177    completely palindromic IESs were identified, of which 153 (40.7%) fell within the same ~228

178    bp length peak, despite comprising several apparently unrelated palindrome sequences

179    (Figure S2, Supplemental Information).

180    IESs in the ~389 bp size peak had distinctive TDRs and TIRs, suggesting they are a family

181    of "mobile IESs", i.e. homologous IESs inserted at nonhomologous genomic sites (Sellis et

182    al., 2021), described further below (see "Pogo/Tigger-family transposon with abundant

183    MITEs").

184    *Repeat elements are abundant in long, non-periodic IESs*

185    Mobile elements that have recently proliferated should appear as interspersed repeat

186    elements in the genome. As identified by RepeatModeler, a quarter of the assembly (12.7

187    Mbp, 23.3%) was composed of such interspersed repeats; like in other model ciliates (Chen

188    et al., 2014; Hamilton et al., 2016), they made up a greater proportion of germline-limited

7

189    IESs (71.0%) than the somatic genome (8.12%) (Figure 3A). The majority of sequence in

190    IESs ≥115 bp was annotated as repetitive, whereas the converse was true for shorter

191    "periodic" IESs (Figure 3C), paralleling short IESs in *Paramecium* which are mostly unique

192    sequences (Arnaiz et al., 2012).

193    Most interspersed repeats could not be classified to a known transposable element class by

194    RepeatClassifier (Figure 3B, Table S2). The most abundant classifiable type was

195    "DNA/TcMar-Tc2", all of which actually belonged to a single repeat family rnd-1_family-1,

196    followed by "LINE/RTE-X". The most abundant family, rnd-1_family-0, was unclassified and

197    made up 21.2% (2.69 Mbp) of total repeats. Families rnd-1_family-0 and rnd-1_family-1 were

198    related to each other and are discussed further below ("Pogo/Tigger-family transposon with

199    abundant MITEs").

200    Three non-periodic IES length peaks (153, 174, 389 bp) could be attributed to specific repeat

201    families, suggesting that they proliferated recently (Table S3, Figure 3C, S3B). This was

202    most pronounced for the ~389 bp peak, where 68.5% of the sequence content belonged to

203    rnd-1_family-0, whereas about a quarter of the ~153 and ~174 bp peaks was composed of

204    repeat families rnd-1_family-87 (palindromic) and rnd-1_family-82 respectively.

205    *Germline-limited repeats include few autonomous transposons but many MITEs*

206    Unlike *Tetrahymena* and *Oxytricha* where transposases are abundant in the germline-limited

207    IESs but rare in the somatic genome (Chen et al., 2014; Hamilton et al., 2016), *Blepharisma*

208    encoded only a few dozen identifiable transposase domains in either the germline-limited or

209    somatic genomes. Cut-and-paste DNA transposase domains of the DDE/D superfamily

210    identified in *Blepharisma* included DDE_1 and DDE_3 (Tc1/Mariner family), DDE_Tnp_1_7

211    (PiggyBac), DDE_Tnp_IS1595 (Merlin), and MULE (Mutator) (Figure 4E, Table S4). Not all

212    copies of DDE/D transposase domains in *Blepharisma* contained an intact catalytic triad,

213    suggesting that some may be inactive fragments or pseudogenes. Nonetheless, domains

214    with an intact triad were found in both germline-limited and somatic sequences. In general,

215    the expression level of somatic transposase genes was substantially higher than germline-

216    limited ones (Figure S4). This contrasts with the observations in *Oxytricha* of abundant

217    germline-limited transposase gene expression (Chen et al., 2014).

218    To identify intact transposon units, we examined the seven repeat families in the MAC+IES

219    assembly classified by RepeatClassifier (Figure 3B). Of these, only two were predominantly

220    germline-limited and represented by more than one full-length copy, namely rnd-1_family-1

221    and rnd-1_family-73 (Table S5). They contained distinct transposases from those found in

222    the MAC genome (Figure 4).

223 **Pogo/Tigger-family transposon with abundant MITEs**

224 Repeat elements of rnd-1_family-1 were bound by a ~30 bp terminal inverted repeat (TIR)

225 5'-CTC CCC CCC CCC CTC CGT GAG CGA ACA AAA-3' whose poly-C run length was

226 variable, possibly from assembly errors, and were flanked by a putative target site

227 duplication (TSD) 5'-TAA-3' (or its reverse complement 5'-TTA-3') (Figure 4B). All thirty

228 intact (≥95% of consensus length) copies of this family were found within IESs, and had high

229 sequence identity, with median 0.5% divergence from consensus.

230 The encoded transposase contained two domains characteristic of Pogo transposases from

231 the Tc1/Mariner superfamily: a DDE/D superfamily endonuclease domain (Pfam domain

232 DDE_1) and a helix-turn-helix domain (HTH_Tnp_Tc5) (Gao et al., 2020). The conserved

233 acidic residues ("catalytic triad") characteristic of DDE/D transposases (Yuan and Wessler,

234 2011) were also present, with the motif DD35D, i.e. all three residues were Asp, 35 a.a.

235 between the second and third conserved Asp. A phylogeny of the DDE_1 domain placed the

236 transposase in the Pogo/Tigger family, most closely related to the Tc2 subfamily and a

237 sequence from the oyster *Crassostrea*, all of which also had the DD35D motif (Figure 4A).

238 The transposase appeared to be germline-limited, with only ten partial Tblastn hits in the

239 somatic MAC genome (seven on "cruft" contigs) mostly overlapping the HTH_Tnp_Tc5

240 domain (17 to 84 a.a., E-values $2.3 \times 10^{-12}$ to $1.4 \times 10^{-6}$) and no matches to the DDE_1

241 domain. However, the TIR did not match previously characterized TIR signatures for the

242 Tc2, Fot, and Pogo subfamilies. A search of all *B. stoltei* IES sequences against HMMs for

243 known DNA transposon TIRs in the Dfam database found only three matches with E-value <

244 0.01, none from the above subfamilies.

245 The same TIR and TSD were also found in another repeat family rnd-1_family-0, which was

246 the most abundant repeat in the genome (Figure 4F), but these were short elements without

247 any predicted coding sequences. rnd-1_family-0 elements often constituted most of the ~389

248 bp IES size class (Figure 3C): the TSDs bounding the repeats (TAA/TTA) were the TDRs for

249 most of these IESs (Figure 2C), and the C-rich TIR motif corresponded to the C-rich IES

250 junctions (Figure 1D, Figure 2F). Copies of rnd-1_family-0 were also found nested in longer

251 IESs, suggesting recent proliferation (Figure S3C). Degenerated or partial copies were found

252 in shorter IESs (Figure 3C), with copies >5% divergence from consensus having median

253 length 308 bp, vs. 388 bp for copies <5% divergence (Figure 4D).

254 Therefore, we interpreted rnd-1_family-1 as a new Pogo/Tigger transposon, with a non-

255 autonomous derivative MITE, rnd-1_family-0. We propose the names Bogo for the

256 transposon and BogoMITE for its MITE, as well as the new term "MITIES" (miniature

257 inverted-repeat transposable internally eliminated sequences) to reflect their dual nature as

9

258    MITEs and IESs. Given their palindromic nature, sequences underlying rnd-1_family-87 and

259    rnd-1_family-160 repeats may also be MITIES.


260    **Tc1-family transposon with microsatellites**

261    Another IES-limited repeat family, rnd-1_family-73, also contained a DDE/D-type

262    transposase coding sequence . Twenty-two copies were >80% of the consensus length with

263    low sequence divergence (median 0.6% vs. consensus). A putative complete transposon

264    bounded by a TSD 5'-TATA-3' and a 38 bp TIR 5'-GTA CCC CCC CCC TCG TTT GTC GCA

265    TTT TCT AGT TTT TT-3' could be defined after manual curation of repeat boundaries

266    (Figure 4C). Nine of these were mobile IESs, with the TSDs corresponding to the IES

267    junctions. The remaining cases were nested in larger IESs alongside other repeat elements.

268    Ten repeats also contained a microsatellite with ~5 to 42 copies of its 10 bp repeat unit 5'-

269    GGG AAG GAC T-3' (Figure 4C) not found elsewhere in the genome. We propose the name

270    BstTc1 for this putative transposon.


271    The transposase encoded in full-length copies of BstTc1 contained a conserved DDE/D

272    superfamily domain DDE_3, phylogenetically affiliated to the Tc1 family although the exact

273    placement is unclear, grouping with only moderate support with Tc1 elements from

274    *Crassostrea* and *Hydra* (Figure 4A). Its catalytic triad motif DD34E differed from previously

275    reported motifs for the Tc1 family, DD41D, DD37D or DD36E (Dupeyron et al., 2020), so it

276    may be a novel subfamily.


277    *Non-LTR retrotransposon sequences in both the somatic and germline genomes*

278    Three retrotransposon repeat families in the MAC+IES assembly were classified by

279    RepeatClassifier, i.e. "LINE" or "LINE/RTE-X" (Table S5). Two of these were more closely

280    related with numerous very high identity sequences (>97%) (Figure 5A), suggesting recent

281    radiation of two related retrotransposon elements, while the third was more divergent (Figure

282    5B; Supplemental Information). Unlike the Bogo and BstTC1-derived elements, more

283    retrotransposon-derived sequences were detected in the *B. stoltei* somatic MAC genome

284    than in assembled IESs (Figure 4E, Table S5). However genes in IESs may be

285    undercounted because of (i) lower completeness of germline vs. somatic assembly; (ii)

286    indels caused by the lower accuracy of the uncorrected long reads used to assemble IESs

287    that prevent prediction; and (iii) shorter total length of IESs than somatic sequence.

288    Consistent with them being true somatic sequences, mappings of error-corrected long reads

289    from a MAC-enrichment library spanned well into flanking regions (Figure 5C; Figure S5A,

290    S5B). In each repeat family, some loci showed sharp dips in coverage suggesting partial

291  excision as IESs while other loci did not (Figure S5B). In MAC-enriched DNA, coverage of

292  such sequences is well above residual IES coverage (Figure S1B).

293  Twenty-nine genes in the main somatic assembly encoded full or partial copies of reverse

294  transcriptase domain RVT_1 (Singh et al., 2021). The four longest retrotransposon genes

295  also encoded an N-terminal apurinic/apyrimidinic endonuclease (Exo_endo_phos_2) domain

296  upstream of RVT_1. This domain pair is characteristic of some proteins from non-LTR

297  retrotransposons/LINE-like transposable elements, e.g. the BS element from *Drosophila*

298  *melanogaster* (UniProt Q95SX7) (Han, 2010; Udomkit et al., 1995). In contrast to the

299  development-specific upregulation of retrotransposon genes in *Tetrahymena* (Fillingham et

300  al., 2004) and *Oxytricha* (Chen et al., 2014), expression of *Blepharisma* genes encoding

301  proteins containing RVT_1 or Exo_endo_phos_2 domains was negligible in starved cells and

302  throughout a post-conjugation developmental time series, for both germline-limited and

303  somatic copies (Figure S4) (Singh et al., 2021). The only exception was a somatic APEX1

304  protein homolog (BSTOLATCC_MAC3189). APEX1 is involved in DNA repair (Fritz, 2000),

305  and Blastp best matches of this *Blepharisma* protein to GenBank's NR database are other

306  similarly annotated proteins.

307  Six retrotransposon-derived sequences from repeat family rnd-1_family-273 contained a

308  central IES that encoded almost half the amino acids of an Exo_endo_phos_2 endonuclease

309  domain (Figure 5D). Excision of the IES during development thus knocks out the

310  endonuclease domain in the somatic version of the gene. Furthermore, the repeat units as a

311  whole had >99% identity to each other over their ~4.1 kbp length, and were flanked by

312  dissimilar sequences (Figure 5D). The similar lengths of these IESs (173 to 182 bp), their

313  homologous location relative to the coding sequence, and their high sequence identity

314  (>96%) all point to a replication of an ancestral retrotransposon which coincidentally

315  contained a sequence recognized and excised as an IES. In two of these cases, the

316  endonuclease and reverse transcriptase domains can be linked into a single reading frame

317  when the IES is present (Figure 5D). None of *Blepharisma*'s putative domesticated

318  transposases are anywhere near as abundant as the retrotransposon repeats in the somatic

319  genome, let alone show signs of substantial recent replication.

320  *Development-specific 24 nt small RNAs are likely scnRNAs in Blepharisma stoltei*

321  Small RNA (sRNA) libraries were sequenced from a developmental time series, where two

322  complementary mating types of *B. stoltei* (strains ATCC 30299 and HT-IV) were separately

323  gamone-treated and mixed to initiate conjugation. Expression patterns of somatic genes

324  from mRNA-seq and the morphological staging have been reported in our sister report on

325  the MAC genome (Singh et al., 2021). Briefly: after mating types were mixed (0 h), cells

326  paired, produced gametic nuclei by meiosis and exchanged them (2 to 18 h), followed by

327  karyogamy (18 to 22 h) and development of the zygotic nuclei to new macronuclei (22 h

328  onwards). At 38 h, about a third of observed cells were exconjugants.

329  The most abundant sRNA length classes were 22 and 24 nt, comprising 32% and 30% of

330  the total reads respectively (Figure 6A). This is consistent with model ciliates, where Dicer-

331  generated, mRNA-derived siRNAs employed in gene silencing are typically 21 or 22 nt long,

332  whereas development-specific sRNAs are distinct and consistently ≥2 bp longer (Lepère et

333  al., 2009; Mochizuki et al., 2002).

334  Developmental dynamics of the 24 nt *Blepharisma* sRNAs resembled scnRNAs of other

335  species. Coverage of 24 nt sRNAs mapping to all feature types initially increased from 2 to 6

336  h and plateaued until 14 h. Coverage over IESs increased further from 14 h to 22 h, reaching

337  ~25 RPKM by the last time point (38 h), whereas coverage declined over coding sequences

338  (CDSs) and other genomic regions ("NON") after 14 h. The initial increase across all feature

339  types coincided with meiotic stages iv to viii of (Miyake et al., 1991) (Singh et al., 2021),

340  whereas the divergence between IESs and the rest of the genome corresponded to the

341  onset of karyogamy (Figure 6B). In contrast, 22 nt sRNAs were initially abundant (albeit with

342  high variance) at CDS and NON regions but low (<1 RPKM) at IESs, and declined sharply to

343  <5 RPKM in all features from 6 h onwards (Figure 6B).

344  *Blepharisma* 24 nt sRNAs had a strongly conserved 5'-U base preference, like scnRNAs in

345  other ciliates (Lepère et al., 2009; Mochizuki and Kurth, 2013; Zahler et al., 2012). For 24 nt

346  sRNAs mapping to IESs, all time points showed conserved 5'-U except for a slight decrease

347  at 6 h (Figure 6D, S6). 24 nt sRNAs mapping to CDSs only showed 5'-U bias after 6 h. We

348  interpret this to mean that 24 nt sRNAs mapping to IESs were predominantly scnRNAs at all

349  time points, whereas those mapping to CDSs initially comprised siRNAs and other types of

350  small RNAs, before being dominated by scnRNAs from 6 h onwards. In contrast, 22 nt

351  sRNAs mapping to CDSs showed no base biases at any time point, whereas 22 nt reads

352  mapping to IESs had a moderate 5'-U bias only from 6 h onwards. The latter may represent

353  true 22 nt scnRNAs, or fragments of originally 24 nt scnRNAs.

*Putative scnRNAs have lower coverage over periodic IESs and BogoMITE IESs*

355  Relative expression levels of putative scnRNAs differed between IES size classes. Based on

356  the IES length distribution and repeat content, we divided IESs into five groups: (1) short

357  "periodic" IESs (≤115 bp), (2) BogoMITEs, because that was the most abundant family, (3)

358  IESs with full-length Bogo transposons, (4) IESs with full-length BstTc1 transposons, and (5)

12

359    all other IESs ("non-periodic"). BogoMITEs and periodic IESs had lower scnRNA coverage

360    (max ~5 and 10 RPKM respectively) compared with nonperiodic IESs (~30 RPKM). The

361    former were comparable to or even lower than expression levels over non-IES features

362    (Figure 6C). Nonetheless, scnRNA coverage of BogoMITEs and periodic IESs showed an

363    initial increase then plateau, without the subsequent decline seen in non-IES regions. Bogo-

364    containing IESs had similar scnRNA coverage to other non-periodic IESs, but BstTc1-

365    containing IESs had higher coverage (Figure 6C).

366    Because of the repetitive sequence content in IESs and the short sRNA length, it is possible

367    that the expression levels calculated could be affected by mis-mapping. We reason that such

368    mismapping would not influence the results described above, because "periodic" IESs

369    (group 1) had low repetitive content, whereas the transposon-containing IESs (groups 2, 3,

370    4) each represented a single repeat family so any mismappings would be contained within

371    the same group and count towards the same RPKM value.

## 372 **Discussion**

373 Despite belonging to the earliest diverging lineage of ciliates sequenced to date, the

374 germline genome of *Blepharisma stoltei* has similarities to established model species,

375 especially the periodic lengths of short IESs like in *Paramecium*. It also provides fresh

376 observations, notably recent proliferation of non-autonomous MITEs that have autonomous

377 counterparts in the same genome, and retroelements in the somatic genome. Parallels

378 between *Paramecium* and *Blepharisma* suggest that ciliate germline characteristics may be

379 relatively plastic over evolutionary time and not strongly phylogenetically constrained.


380 *Comparison to IESs in other ciliates*

381 Most *Blepharisma* IESs are short, TA-bound, and intragenic, more similar to *Paramecium*

382 than *Tetrahymena* or spirotrichs. The most striking parallel is the sharply periodic length

383 distribution of short IESs with peaks every ~10 bp, coinciding with the DNA helical turn,

384 implying that the *Blepharisma* excisase complex has similar geometric constraints as

385 proposed for *Paramecium* (Arnaiz et al., 2012). *Blepharisma* "periodic" IESs are longer on

386 average and do not have a "forbidden" second peak, but the last peak (~110 bp; Figure 1A)

387 is still below the upper limit where such periodicity would be expected given the properties of

388 DNA (Figure 7 of (Arnaiz et al., 2012)). In contrast, *Tetrahymena thermophila* has a

389 continuous distribution (average length ~3 kbp) (Hamilton et al., 2016; Seah and Swart,

390 2021), while *Oxytricha trifallax* non-scrambled IESs (length ~20-100 bp) have weak

391 periodicity (Chen et al., 2014). Periodicity is consistent with a single primary IES excisase,

392 rather than multiple excisase families, which would smooth the length distribution.


393 Longer, nonperiodic IESs of *Blepharisma* contain more repeats, including whole

394 transposons, than short IESs. Unlike *Tetrahymena*, where 41.7% of high-confidence IESs

395 comprise putative autonomous transposons (Hamilton et al., 2016), some of which can be

396 grouped into families (Fillingham et al., 2004; Wuitschick et al., 2002), only a small fraction of

397 *Blepharisma*'s long IESs encode transposases, and their length distribution is not unimodal,

398 but long-tailed, with distinct peaks representing individual abundant families (Figure 3).

399 Germline-specific repeats and transposons across *Paramecium* spp. have recently been

400 surveyed (Sellis et al., 2021), but were likely underestimated because such repeats are

401 difficult to assemble from short-read data even with high coverage, as we saw with

402 *Blepharisma* BogoMITE elements, (Supplemental Information, Figure S1A).


403 The dynamics of *Blepharisma* 24 nt sRNAs are consistent with the scnRNA turnover model,

404 where RNA intermediates are produced from both IESs and MDSs (Malone et al., 2005;

405 Mochizuki and Gorovsky, 2005), but those from MDSs are selectively degraded, allowing the

406   remaining scnRNAs to mark IESs for excision. *Blepharisma* 24 nt sRNAs mapping to IESs

407   increase more than those mapping to CDSs during post-conjugation development (Figure

408   6B), complementing our finding that homologs of scnRNA biogenesis proteins, Dicer-like

409   (Dcl) and Piwi proteins, are highly upregulated during development (Singh et al., 2021).

410   Furthermore, higher coverage of *Blepharisma* scnRNAs in longer (presumably younger)

411   IESs than in short (~older) periodic IESs mirrors that of *Paramecium*, where younger IESs

412   are more likely to require scnRNAs for efficient excision (Lhuillier-Akakpo et al., 2014; Sellis

413   et al., 2021).

414   The longer an IES, the more likely it will contain a promoter by chance or contain one from a

415   transposase gene, thus giving rise to such sRNAs. This would explain the low 24 nt sRNA

416   level of BogoMITE IESs in contrast to their autonomous counterparts (Figure 6C), though

417   removal of both is essential. In contrast to the abundant Bogo transposon 24 nt sRNAs,

418   expression of these and other transposase genes in RNA-seq is negligible (Figure S4). This

419   raises the possibility that active, transcribed *Blepharisma* transposons are in fact silenced,

420   turning most of their transcripts into 24 nt sRNAs. This is contrary to the role of scnRNAs

421   proposed to target DNA for excision, but congruent with the role of sRNAs in transposon

422   silencing in other eukaryotes, from which the scnRNA biosynthesis enzymes originated

423   (Sandoval et al., 2014).

424   *Are MITEs a missing link in the IBAF model?*

425   The prevailing Invasion-Bloom-Abdication-Fade (IBAF) model for the evolution of IESs

426   hypothesizes that they originate from cut-and-paste DNA transposons that invade and

427   proliferate ("bloom") in the germline genome (Klobutcher and Herrick, 1997). Transposon

428   proliferation stops ("abdication") when its transposase is domesticated by a host promoter,

429   releasing the transposons from purifying selection, whereupon their sequences erode by drift

430   ("fade"). Depictions of the IBAF model usually show all the transposons expressing

431   transposases during "bloom", i.e. functioning as autonomous transposons (Feng and

432   Landweber, 2021; Klobutcher and Herrick, 1997). This is reasonable for *Tetrahymena* and

433   *Oxytricha*, which have hundreds of germline-encoded transposases that vastly outnumber

434   those in the somatic genome (Table S4). However, *Blepharisma* and *Paramecium* only have

435   a few dozen transposases, although germline-limited transposases may be underestimated,

436   especially for short-read assemblies.

437   This discrepancy can be resolved by taking MITIESs (MITE IESs) into account. In

438   *Blepharisma* this is best exemplified by the few autonomous Bogo transposon copies

439   compared to thousands of non-autonomous BogoMITEs. The narrow length distribution of

440    BogoMITEs, their high sequence identity, and occasional nested insertion inside unrelated

441    IESs are the clearest illustrations to date of recent MITE proliferation. Bogo is also the first

442    Pogo/Tigger transposon found in a ciliate germline genome; this subfamily is known to be

443    especially prone to MITE formation (Feschotte and Mouchès, 2000; Guermonprez et al.,

444    2008). The prevalence of IESs bound by terminal inverted repeats, including numerous

445    palindromic IESs (Figure 2D, S2), also suggest many more *Blepharisma* IESs are MITE

446    derivatives.

447    In *Paramecium* spp., MITEs of the Thon and Merou transposons have been identified but

448    only numbered about a dozen copies per genome, and their transposases belong to a

449    different transposase family than Bogo (Figure 4). The most abundant mobile IES family in

450    *Paramecium*, FAM_2183, is probably a MITE but its autonomous counterpart was not

451    reported (Sellis et al., 2021). MITEs as transposon/IES life cycle intermediates can hence

452    explain why *Blepharisma* and *Paramecium* have few MIC-encoded transposases compared

453    to *Oxytricha* and *Tetrahymena*, but nevertheless tens of thousands of IESs.

454    MITEs also provide a mechanism for transposon/IES proliferation self limitation (Figure 7A).

455    When MITEs outnumber the autonomous transposon, active transposase protein is more

456    likely to bind to target sites in MITEs than the full length transposon ("titration"), hindering the

457    replication of the autonomous version, giving time for loss-of-function mutations to inactivate

458    the transposases ("fade"). This "vertical inactivation" scenario (Hartl et al., 1997) was already

459    discussed in the original IBAF proposal (Klobutcher and Herrick, 1997), but no plausible

460    examples from ciliates were then known.

461    *Is "genome defense" a flawed analogy?*

462    The IBAF model also does not explain how ciliates can consistently and precisely excise

463    novel mobile elements from different transposon families that invade the germline genome.

464    The domesticated excisases of *Paramecium* (Baudry et al., 2009), *Tetrahymena* (Cheng et

465    al., 2010), and *Blepharisma* (Singh et al., 2021) belong to the PiggyBac family. Except for

466    *Tetrahymena* Tpb2, PiggyBacs are known to perform seamless excision, where the host

467    sequence after transposon excision is identical to that before insertion (Chen et al., 2020).

468    This would make them the ideal progenitor for IESs within coding sequences; indeed,

469    PiggyBac transposons are also known to produce MITEs (Mitra et al., 2013; Wang et al.,

470    2010). By extension, the first IESs probably originated from PiggyBac transposons. But what

471    about subsequent invasions by other transposons that leave behind "scars" upon excision?

472    Such imprecision would cause deleterious frameshift mutations in coding regions. How can

473    they invade the germline genome and yet avoid deleterious effects?

474    Part of the answer lies in the "hijacking" model proposed from *Paramecium* (Arnaiz et al.,

475    2012; Sellis et al., 2021), whereby the domestication of PiggyBac transposase changed the

476    dynamic for subsequent transposon invasions. New transposons would persist as IESs only

477    if they also encode a seamless excisase, or if they can also be recognized and cut by the

478    exapted PiggyBac transposase. The latter favors the invasion of transposons that produce a

479    TSD containing a submotif recognized as a cut site by PiggyBac (Figure 7B). The similarity

480    between IES and transposon boundaries would hence not be due to common origin or

481    sequence evolution after IES fixation in the germline (Klobutcher and Herrick, 1997), but

482    rather because of selection for transposons whose TSDs already match the excision site

483    preferences of domesticated PiggyBac. Analogous exaptation of TSDs for excision has been

484    demonstrated in another context: independent origin of introns from MITEs in at least two

485    different eukaryotes, where one of the TSDs produced upon MITE insertion was co-opted as

486    an intron splice site (Huff et al., 2016). Cross-talk between different (albeit related)

487    transposases for MITE transposition has also been documented (Feschotte et al., 2005).

488    We further argue that "genome defense" is a teleological expression that confuses cause

489    and effect. Domesticated excisases actually facilitate mobile element accumulation in the

490    germline, by shielding them from selection by effective exclusion from the somatic genome.

491    *Tetrahymena* is the exception that proves the rule: its domesticated excisase appears to be

492    imprecise; correspondingly, most of its IESs are intergenic, because intragenic IESs have

493    been efficiently removed by selection (Cheng et al., 2016; Feng et al., 2017). The origins of

494    gene silencing by DNA methylation in vertebrates have also been reinterpreted with similar

495    reasoning. Vertebrates have high levels of CpG methylation that inactivates transposons,

496    which was thus proposed to "compensate for" transposon proliferation in eukaryotic

497    genomes (Bestor, 1990). When seen from a non-teleological perspective, it is precisely

498    because CpG-mediated transposon inactivation is so effective, preventing exposure to

499    selection, that transposons persist, leading to larger genomes (Zhou et al., 2020).

500    *Why does the Blepharisma somatic genome contain retrotransposon sequences?*

501    Transposon-related sequences are typically germline-limited in other model ciliates, which

502    was formerly interpreted as successful "genome defense" keeping them out of the somatic

503    MAC genome (Chen et al., 2014; Fillingham et al., 2004; Guérin et al., 2017; Hamilton et al.,

504    2016; Swart et al., 2013). Counter to this, we found several retrotransposon-derived

505    sequences in the *Blepharisma* MAC genome (Figure 5; Table S2). Some show signs of

506    partial excision or possible absence of the locus in part of the population, but plenty have

507    uniform coverage typical of somatic sequences.

508 Recent retrotransposon proliferation in the soma and patchy distribution of different somatic

509 transposase classes across ciliates (Table S4) (Singh et al., 2021) suggest that "genome

510 defense" is at best leaky. We conjecture that if foreign DNA lacks suitable target sites

511 recognized by the excisase, it might still be marked by scnRNAs but fail to be excised or only

512 be partially excised (e.g. the IESs in Figure 5C). Such DNA would still be deleterious if

513 inserted intragenically.

514 Somatic MACs may be unable to repress mobile elements by heterochromatinization like

515 germline MICs and other eukaryotic nuclei. In *Tetrahymena*, most MAC DNA is not

516 associated with classical heterochromatin marks (Liu et al., 2007), while in *Paramecium*

517 MACs, H3K27me3 is not associated with transcription repression, despite being a classic

518 heterochromatin mark in multicellular eukaryotes (Drews et al., 2021). In such a permissive

519 expression environment, selection against mobile elements that are not already excised as

520 IESs may be especially effective, unless they are relatively transcriptionally inactive like the

521 *Blepharisma* retroelements. On the other hand, regular *Blepharisma* stock culture passaging

522 maintains a small effective population size, which would counteract selection against mobile

523 element accumulation in the soma.

524 The genome defense model may lead one to dismiss IES retention in the somatic genome

525 as excisase inefficiency or MIC contamination of the library, however, IES excision is not all-

526 or-nothing but a continuum. Experimental evolution experiments in *Paramecium* suggest IES

527 retention variability is itself a plastic and evolvable trait with consequences for somatic

528 genotypic diversity (Catania et al., 2021; Vitali et al., 2019). Assembly algorithms tend to

529 present an oversimplified, "pristine" view of somatic genomes, because they collapse

530 repetitive and lower-coverage regions, which are characteristic of mobile elements and

531 partially retained IESs. Accurate long read sequencing, haplotype-aware assemblers, and

532 sequence graphs will all play a role in building a more realistic picture of somatic genome

533 heterogeneity.

534 *Conclusion*

535 Why do we credit developmental DNA elimination with defending the genome, when natural

536 selection has been doing the hard work? Apart from technical biases during genome

537 assembly, there is also sampling bias by using lab strains. These are often clonal and largely

538 homozygous; if so, we would not observe accumulation of strongly deleterious foreign DNA

539 that actually needs defending against, but only IESs that have reached fixation and that are

540 already efficiently excised and non-deleterious. Purifying selection against deleterious IESs

541 has had to be indirectly observed, e.g. in the lack of intragenic IESs in *Tetrahymena*, where

18

542    excision is imprecise (Hamilton et al., 2016), and the statistical depletion of IES-like

543    sequences in the *Paramecium* somatic genome (Swart et al., 2014). Similar evolutionary

544    logic applies to prokaryotic CRISPR defense systems, where hidden fitness costs

545    (autoimmunity) have been underestimated because those individuals are removed by

546    selection (Stern et al., 2010), hence the phenomenon is easily misinterpreted as inheritance

547    of acquired traits (Weiss, 2015). Most studies on ciliate developmental DNA elimination to

548    date have focussed on the underlying molecular mechanisms, but to understand its origins

549    and evolution we should expand our view to diverse ciliates and their germline genomes

550    from natural populations.

## Figure Legends

**Figure 1. A "hybrid" IES length distribution with periodic length peaks for short IESs.**
(**A**) IES length histogram (0 to 500 bp (inset: full range), stacked bars for types of terminal direct repeats (TDRs) at IES boundaries. Peaks for IES size classes discussed are marked.
(**B**) Comparison of cut-and-paste DNA transposons (above) and ciliate genome editing (below), showing parallels between target site duplications (TSD) of transposons and terminal direct repeats (TDRs) bounding IESs, and effects of precise vs. imprecise excision.
(**C**) Diagrammatic tree of ciliates (following Lynn, 2008), branch lengths arbitrary. Genera with draft MIC genomes listed on right. (**D**) Sequence logos for MDS-IES junctions for TA-bound IESs of specific size classes, centered on the "TA". See also Figure S1.

**Figure 2. IESs are bounded by heterogeneous direct and inverted terminal repeats.** (**A**)
Numbers of terminal direct repeats (TDRs) per TDR length observed (blue) vs. number expected by random chance if bases were independently distributed (orange). (**B**) Ratio of observed to expected numbers of TDRs by length. (**C**) Length distributions of IESs containing TDRs of lengths 2, 3, 4, and 5 bp; the most abundant TDR sequences per TDR length are shown in color (sequences and their reverse complements are counted together, because TDRs could be encountered in either orientation, e.g. TAA/TTA), simple T/A alternations are in shades of blue. NB: plots in panel C have different vertical axis scales. (**D**)
Observed IESs per terminal inverted repeat (TIR) length vs. expected number by chance alone. (**E**) Same as panel D but for *P. tetraurelia*. (**F**) Lengths (scatter-overlaid boxplot) of IESs containing long TIRs (≥10 bp), grouped by their TIR sequence (rows). Each TIR-cluster is annotated with the median IES length (bp), cluster size (n), TDR consensus sequence, and TIR representative sequence. See also Figure S2.

**Figure 3. Repeat elements are abundant in long, non-periodic IESs.** (**A**) Total sequence length annotated as interspersed repeats vs. non-repetitive , in germline-limited vs. somatic parts of the genome. (**B**) Classification of repeat families by RepeatClassifier, and total annotated length per repeat class. (**C**) Total sequence length (vertical axis) per IES size class (horizontal axis), stacked plot of non-repetitive fraction vs. interspersed repeats , with the most abundant repeat families in the four non-periodic peaks overlaid in color. Inset: Distribution to 1000 bp. See also Figure S3.

**Figure 4. Germline-limited repeats include few autonomous transposons but many MITEs.** (**A**) Phylogenetic tree of DDE/D domains for Tc1/Mariner superfamily, including *B. stoltei* germline-limited (Bogo and BstTc1) and somatic transposases. (**B**) Diagram of features in Bogo and BogoMITE; TSD – target site duplications, TIR – terminal inverted

20

585  repeats, HTH_Tnp_Tc5, DDE_1 – conserved domains. (**C**) Diagram of features in BstTc1:

586  DDE_3 – conserved domain. (**D**) Histograms of sequence divergence from repeat family

587  consensus for copies of the Bogo and BogoMITE repeat families annotated by

588  RepeatMasker; for rnd-1_family-1, most low-divergence copies (<5% divergence) were short

589  fragments, but all full-length copies were low-divergence. (**E**) Counts of transposase-related

590  domains in different ciliates from six-frame translations of somatic vs. germline-limited

591  genome sequence. See also Figure S4. (**F**) Sequence logos for Bogo and BogoMITE repeat

592  boundaries, aligned on the terminal inverted repeats (TIRs) and terminal direct repeats

593  (TDRs). 3'-boundaries have been reverse complemented to show the TIRs. Sequence logos

594  were generated from alignments of full-length, intact Bogo elements (>1.8 kbp) and

595  BogoMITEs (between 385-395 bp), with columns comprising >90% gaps removed.

596  **Figure 5. Non-LTR retrotransposon sequences in both somatic and germline**

597  **genomes.** (**A**) Phylogeny of rnd-1_family-273 and rnd-1_family-276 retrotransposon

598  sequences. (**B**) Phylogeny of rnd-4_family-193 retrotransposon sequences. (**C**) Window of

599  mapped HiFi reads from sucrose gradient-purified MACs (grey) spanning a retrotransposon

600  gene with both an AP endonuclease domain and a reverse transcriptase domain (from rnd-

601  4_family-193). Only sequence columns with < 90% gaps are shown. (**D**) Multiple sequence

602  alignment of non-LTR retrotransposon copies from rnd-1_family-273. Schematic for

603  consequences of IES excision (Contig_45). Identity scale: green=100%; gold=30-99.9%;

604  red=0-29.9%. See also Figure S5.

605  **Figure 6. Development-specific 24 nt small RNAs are likely scnRNAs in *B. stoltei*.** (**A**)

606  Read length histogram for all sRNAs in the time series. (**B**) Relative expression (RPKM

607  units, vertical axis) of 22 and 24 nt sRNAs mapping to different feature types across time

608  series: blue - IES, orange - CDS, green - all other regions not annotated as IES or CDS

609  (including UTRs and intergenic regions which are difficult to delimit exactly with available

610  data). Timing of developmental stages inferred from morphology are labeled below (Singh et

611  al., 2021). (**C**) Relative expression of 22 and 24 nt sRNAs mapping to different categories of

612  IESs: containing full-length copies of BstTc1 and Bogo transposons, at least 90% covered by

613  BogoMITE elements, IESs in the periodic length range (< 115 bp), and all other IESs ("non-

614  periodic"). (**D**) Sequence logos for 22 and 24 sRNAs mapping to CDS and IES features in

615  controls and different time points (rows). See also Figure S6.

616  **Figure 7. Model for IES evolution in a ciliate genome with an existing domesticated**

617  **excisase.** (**A**) Graphs depict IES length distribution. (**1**) Invasion of germline genome by full

618  length transposon (green); existing IESs (blue) are excised by domesticated excisase. (**2**)

619  New transposon produces MITIES which are both MITES and IESs. (**3a**) If MITIES can be

21

620    excised by domesticated excisase, they proliferate and titrate the progenitor transposase. (**4**)

621    Proliferation of MITIES favors vertical inactivation of the full length transposon; loss of

622    function stops production of new MITIES, leading to eventual decay. (**3b**) If the MITE cannot

623    be excised by domesticated excisase (i.e. it is not an IES), it is more likely to cause

624    deleterious mutations upon insertion, and is therefore selected against and does not reach

625    fixation. (**B**) If a transposon TSD contains a submotif that can be recognized by the

626    domesticated excisase, it can theoretically be excised cleanly without leaving a "footprint",

627    avoiding potential frameshift mutations.

628    **Supplemental Figure Legends**

629    **Figure S1. Length distributions and retention scores for different IES assembly**

630    **methods, MAC library, and cryptic IESs.** (**A**) Comparison of IES reconstructions from

631    MIC-enrichment library sequenced with short reads by ParTIES (above) vs. from long reads

632    by BleTIES (below). Main panels: IES length histograms up to 500 bp, insets: IES retention

633    scores colored by TDR sequence type. Length peak at ~390 bp representing BogoMITE

634    element is present in BleTIES reconstruction but not ParTIES. (**B**) Conventional IESs:

635    retention scores computed from MAC-enriched library, sequenced with PacBio HiFi reads.

636    (**C**) "Cryptic" IESs from MAC read library: length histogram, colored by TDR sequence type.

637    (**D**) Retention scores of "cryptic IESs". (**E**) Length distribution of "cryptic" IESs that contain

638    "TTA" or "TAA" in their TDR, detail <500 bp, inset detail <150 bp. (**F**) Sequence logos of TA-

639    bound "cryptic" IES junctions centered on the TA motif, for all cryptic IESs (above) and the

640    subset in the ~72 bp size class (below). (**G**) Mapping pileup at IES with TA-containing TDR.

641    For aligned reads in panels E and F, dots: bases identical to reference, dashes: gaps

642    relative to reference, red bar: read clipping. (**H**) Mapping pileup at IES with non-TA-

643    containing TDR.

644    **Figure S2. Palindromic IESs clustering and length distribution.** Strip plots of IES lengths

645    for palindromic IESs (≥90% self-alignment identity), after they have been clustered by

646    sequence identity (rows represent clusters). Each cluster is annotated with the median IES

647    length and the cluster size. Insets: (**A**) Overall sequence length distribution histogram for all

648    palindromic IESs. The most common length of palindromic IESs is ~230 bp. (**B, C**)

649    Dendrogram of sequence distance and multiple sequence alignment of palindromic IESs

650    with ~230 bp length to illustrate that they comprise several distinct clusters of sequences.

651    **Figure S3. Most abundant repeat families in non-periodic IES size classes.** (**A**) Total

652    lengths (horizontal axis) of the top ten repeat families per IES size class (panel rows). (**B**)

653    Top repeat family (by sequence length) for each IES size class (panel rows); the total length

654    covered by that repeat family within IESs vs. the lengths of those IESs is shown in red,

655    superimposed on the total sequence vs. IES length distribution of IESs in general (grey).

656    Arrowheads mark centers of the size classes. (**C**) Examples of nested repeats within IESs.

657    Nested elements can be recognized when the two outer repeat elements belong to the same

658    family and align to consecutive parts of its family's consensus sequence, implying that the

659    inner element has likely been inserted into the middle of an existing element. Coordinates of

660    the split segments are relative to the repeat family consensus.

661    **Figure S4. Expression of genes with transposase domains.** Comparison of expression

662    levels for MAC- vs. MIC-limited transposase-related domains across developmental time

663    series; heatmap color scaled to log(transcripts per million). Domain architecture shown

664    diagrammatically.


665    **Figure S5. Non-LTR retrotransposon sequences in both somatic and germline**

666    **genomes.** (**A**) As in Figure 5A. (**B**) As in Figure 5A. Inset shows coverage across the entire

667    contig and position of the retrotransposon gene. (**C**) Alignment of MAC+IES and somatic

668    genomic sequences for Contig_44 retroelement genes from Figure 5A, showing how

669    excision of the central IES deletes part of the endonuclease domain and produces a

670    premature stop codon.


671    **Figure S6. Per-position base entropy of 22 nt and 24 nt sRNAs from developmental**

672    **time series.** Plots show conservation of 5'-U in 24 nt sRNAs. Each plot symbol represents

673    positional sequence entropy (symbol size) for a given nucleotide base (columns) and

674    position in the sRNA sequence (vertical axis) and time point (horizontal axis), in sRNAs

675    mapping to different feature types (rows).

## **Methods**

676

677   General reagents were analytical grade and purchased from Sigma-Aldrich or Merck unless

678   otherwise indicated.

679   *Ciliate strains origin and cultivation*

680   The strains used and their original isolation localities were: *Blepharisma stoltei* ATCC 30299,

681   Lake Federsee, Germany (Repak, 1968); *Blepharisma stoltei* HT-IV, Aichi prefecture, Japan

682   (Harumoto et al., 1998). Methods for cell cultivation and harvesting of material for

683   sequencing are described in our sister report (Singh et al., 2021).

684   *Enrichment of micronuclei, isolation and sequencing of genomic DNA*

685   *B. stoltei* ATCC 30299 cells were harvested and cleaned to yield 400 mL of cell suspension

686   (1600 cells/mL). This suspension was twice concentrated by centrifugation (100 g; 2 min;

687   room temperature) in pear-shaped flasks and in 50 mL tubes to ~8 mL. 10 mL chilled Qiagen

688   Buffer C1 (from the Qiagen Genomic DNA Buffer Set, Qiagen no. 19060) and 30 mL chilled,

689   autoclaved deionized water were added. The suspension was mixed by gently inverting the

690   tube until no clumps of cells were visible, and then centrifuged (1300 g; 15 min; 4°C). The

691   pellet was washed with chilled 2 mL Buffer C1 and 6 mL water, mixed by pipetting gently

692   with a wide-bore pipette tip, centrifuged (1300 g; 15 min; 4°C), and resuspended with chilled

693   2 mL Buffer C1 and 6 mL water by pipetting gently with a wide-bore pipette tip.

694   The nuclei suspension was layered over a discrete sucrose gradient of 20 mL 10% (w/v)

695   sucrose in TSC medium (0.1% (v/v) Triton X-100, 0.01% (w/v) spermidine trihydrochloride

696   and 5mM $CaCl_2$) on top of 40% (w/v) sucrose in TSC medium (Lauth et al., 1976). Gradients

697   were centrifuged (250 g; 10 min; 4°C). 10 to 12 mL fractions were collected by careful

698   pipetting from above, and the nuclei were pelleted by centrifugation (3000 g; 10 min; 4°C).

699   DNA was extracted from pelleted nuclei with the Qiagen Genomic tips 20/G and HMW DNA

700   extraction buffer set (Qiagen no. 19060) according to the manufacturer's instructions. DNA

701   concentration was measured by the Qubit dsDNA High-Sensitivity assay kit. Fragment size

702   distribution in each sample was assessed by a Femto Pulse analyzer.

703   *B. stoltei* ATCC 30299 DNA isolated from the MIC-enriched fraction on two separate

704   occasions was used to prepare two sets of DNA sequencing libraries. A low-input PacBio

705   SMRTbell library was prepared without shearing the DNA and was sequenced in the CLR-

706   (continuous long read) sequencing mode on a PacBio Sequel II instrument. Paired-end

707   short-read libraries were prepared for four sucrose gradient fractions (top (T), middle (M),

25

708     middle lower (ML), bottom (B)) and sequenced with 100 bp BGI-Seq paired-end reads on a

709     BGI-Seq instrument.


710     *IES prediction from PacBio subreads*

711     PacBio subreads (CLR reads) from a MIC-enriched sample (ENA accession ERR6548140)

712     were aligned to the somatic genome reference assembly (accession PRJEB40285) (Singh et

713     al., 2021) with minimap2 v2.17-r941 (Li, 2018), with options: -ax map-pb --secondary=no --

714     MD. Mapped reads were sorted and indexed with samtools v1.10 (Li et al., 2009), and then

715     used for predicting IESs with BleTIES MILRAA v0.1.9, with options: --type subreads --

716     junction_flank 5 --min_ies_length 15 --min_break_coverage 10 --

717     subreads_pos_max_cluster_dist 5. The BleTIES pipeline has been previously described

718     (Seah and Swart, 2021) and uses spoa v4.0.3 (Vaser et al., 2017) for assembly. After

719     inspecting the initial IES predictions, we removed IES predictions with length <50 bp and

720     retention score <0.075, which we judged to be more likely to be spurious or to have

721     insufficient coverage for an accurate assembly.

722     Terminal direct repeats (TDRs) at the boundary of a given IES were defined as a sequence

723     of any length that was exactly repeated on both ends of the IES, such that one copy lies

724     within the IES, and the other in the MAC-destined sequence. Because the sequence is

725     identical, it is not possible to determine from sequencing data alone where the physical

726     excision of the IES would occur; such ambiguous excision junctions have been termed

727     "floating IESs" (Sellis et al., 2021). Therefore, TDRs were always reported starting from the

728     left-most coordinate. If the TDR sequence contained 5'-TA-3', the corresponding IES was

729     also considered to be "TA-bound", even if the TDR was longer than the 2 bp 5'-TA-3'

730     sequence.

731     Reconstructed IES sequences were computationally inserted into the MAC assembly with

732     BleTIES Insert, to produce a hybrid MAC+IES assembly, which approximates the part of the

733     MIC genome that is collinear with the MAC.


734     *Identification and comparison of IES length classes*

735     Visual inspection of the length distribution of BleTIES-predicted IESs showed sharp peaks

736     every ~10 bp between ~65 and 115 bp. Peak calling on the graph of number of IESs (TA-

737     bound only) vs. length (bp) was performed with the function find_peaks from the Python

738     package scipy.signal v1.3.1 (Virtanen et al., 2020), with height cutoff 100. The ranges for

739     each IES size class were defined with the width at half peak height. In *Paramecium*

740     *tetraurelia*, where most IESs are TA-bound, the IES termini have a short, weakly conserved

741 inverted repeat (Arnaiz et al., 2012; Klobutcher and Herrick, 1995). To search for similar

742 motifs in *B. stoltei*, sequences flanking TA-bound IES junctions were extracted, with one

743 from each pair reverse-complemented so that the sequences were always in the orientation

744 5'-(MDS segment)-TA-(IES segment)-3'. Sequence logos of the junctions (10 bp MDS, 14 bp

745 within IES, not including the TA itself) were drawn for each IES length class with Weblogo

746 (Crooks et al., 2004). Only TA-bound IESs were used for the sequence logos because they

747 could be aligned relative to the 5'-TA-3' repeat, whereas for IESs bound by other types of

748 junctions there is no common reference point to align the boundaries of the IES.


749 *Probability of a pair of repeated sequences*

750 Under a null model where all bases in a sequence are independently and identically

751 distributed, the probability $P_n$ of having any possible sequence of length $n$ bounding a given

752 sequence feature (either a TDR or a TIR) is the sum of probabilities of all possible

753 sequences (each of which notated as $k$) of length $n$, squared: $P_n = \sum_{k \in K} p_k^2$, which can be

754 transformed to $P_n = (\sum_{b \in B} p_b^2)^n$, where *B* is the alphabet of bases and $p_b$ is the individual

755 probability of each base. The number of possible sequences $k$ of length $n$ is simply $|K| = |B|^n$.


756 The probability of having a repeat of length at least 2 is equal to the probability of having a

757 repeat of length 2, because all cases of repeat length > 2 implicitly have a repeat of length =

758 2. Therefore the probability of having a repeat of length exactly $n$, i.e. match in bases 1 to $n$,

759 and mismatch on base $n$+1 is $P_n \times Pr\,(mismatch) = P_n \times (1 - \sum_{b \in B} p_b^2)$. The expected

760 number of TDRs in *Blepharisma* were calculated by using the empirical base frequencies of

761 the MAC+IES genome assembly for $p_b$, and multiplying this probability by the number of

762 IESs.


763 *Identification of terminal inverted repeats (TIRs) and palindromes in IESs*

764 The BleTIES-assembled IES sequences for *Blepharisma* were used to identify exact,

765 ungapped terminal inverted repeats (TIRs). Starting from the ends of the IES sequence

766 immediately within the flanking TDRs, each base was compared to the reverse complement

767 of the corresponding base on the opposite end for a match, extending the TIR until a

768 mismatch was encountered, up to a maximum length of 25 bp. The same procedure was

769 used for *Paramecium tetraurelia* using IESs sequences downloaded from ParameciumDB

770 (https://paramecium.i2bc.paris-

771 saclay.fr/files/Paramecium/tetraurelia/51/annotations/ptetraurelia_mac_51_with_ies,

772 accessed 14 October 2021), except that the coordinates of TDRs were first renumbered and

773 extended beyond the "TA" motif if possible, following the BleTIES coordinate numbering

774    convention, in case there are potential TDRs that are longer than a simple TA. The expected

775    number of TIRs of given lengths under a null model was computed as described in

776    "Probability of a pair of sequences".

777    Long TIRs (≥10 bp) were clustered by sequence identity to look for IESs of potentially

778    related origin, using the cluster_fast algorithm (Edgar, 2010) implemented in Vsearch

779    v2.13.6 (Rognes et al., 2016) at 80% identity and the CD-HIT definition of sequence identity

780    (-iddef 0). For each resulting cluster of similar TIRs, the cluster centroid was used as the

781    representative sequence shown in Figure TIRS. TDRs associated with each cluster's IESs

782    were grouped by length, and for each TDR length a degenerate consensus was reported

783    with the degenerate_consensus function of the Bio.motifs module in Biopython v1.74.

784    Palindromic IESs were defined as IESs that align to their own reverse complement with a

785    sequence identity ≥90% (matching columns over sequence length); this definition was less

786    strict and permitted inexact matches unlike the TIR search, to allow for sequence divergence

787    and assembly errors. IES sequences were aligned with the PairwiseAligner function from

788    Bio.Align in BioPython v1.74, using global mode and parameter match_score = 1.0, with all

789    other scores set to zero.

790    Palindromic IESs were clustered with Vsearch cluster_fast as described above, except that

791    one sequence (BSTOLATCC_IES35757) was manually removed after inspection of results

792    because it appears to contain two different nested palindromic sequences. Cluster centroids

793    were aligned pairwise as above and used to calculate a matrix of edit distances (matching

794    columns / alignment length). The distance matrix was clustered with average linkage

795    clustering to produce a sequence distance dendrogram with the functions average and

796    dendrogram from scipy.cluster.hierarchy v1.3.1 (Virtanen et al., 2020).

797    *Comparison of intragenic:intergenic IES ratios*

798    Intragenic vs. intergenic IESs were defined by overlap of predicted IES annotations with

799    "gene" feature annotations on the MAC reference (ENA accession GCA_905310155), using

800    Bedtools v2.30.0 (Quinlan and Hall, 2010) and pybedtools v0.8.1 (Dale et al., 2011).

801    To test whether the underrepresentation of IESs within gene features was statistically

802    significant, compared to the null hypothesis of IESs and gene feature locations being

803    independently distributed, we assumed that the number of intragenic IESs would follow a

804    binomial distribution with individual probability equal to the fraction of the genome that is

805    covered by gene features. The p-value of the observed number of intragenic IESs would

806    then be equal to the cumulative probability density up to and including the observed value.

807 *Developmental time series small RNA-seq*

808 Complementary mating strains *B. stoltei* ATCC 30299 and HT-IV were pre-treated with

809 Gamone 2 and Gamone 1 respectively, and then mixed to initiate conjugation as described

810 previously; sRNA and mRNA were isolated from total RNA at the same time points

811 ("Conjugation time course", (Singh et al., 2021)). sRNA libraries were prepared with the

812 BGISeq-500 Small RNA Library protocol, which selects 18 to 30 nt sRNAs by polyacrylamide

813 gel electrophoresis, and sequenced on a BGISeq 500 instrument.

814 *Small RNA libraries mapping and comparison*

815 Small RNA libraries were mapped to the MAC+IES assembly with bowtie2 v2.4.2

816 (Langmead and Salzberg, 2012) using default parameters. Total reads mapping to CDS vs.

817 IES features were counted with featureCounts v2.0.1 (Liao et al., 2014). To account for

818 different total sequence lengths represented by CDSs, IESs, and intergenic regions, the read

819 counts were converted to relative expression values (reads per kbp transcript per million

820 reads mapped, RPKM (Mortazavi et al., 2008) ) using the total lengths of each feature type

821 in place of transcript length in the original definition of RPKM, with the following formula:

822 $10^9 \times$ (reads mapped to feature type) / (total reads mapped $\times$ total length of feature type).

823 Reads mapping to CDSs, IESs, or neither (but excluding tRNA and rRNA features) were

824 extracted with samtools view, with 22 and 24 nt reads extracted to separate files. Read

825 length distributions for each sequence length and feature type were summarized with

826 samtools stats.

827 *mRNA-seq read mapping*

828 To permit correct mapping of tiny introns RNA-seq data was mapped to the MAC genome

829 using a version of Hisat2 (Kim et al., 2019) with the static variable minIntronLen in hisat2.cpp

830 in the source code lowered to 9 from 20 (https://github.com/Swart-lab/hisat2/; commit hash

831 86527b9). Hisat2 was run with default parameters and parameters --min-intronlen 9 --max-

832 intronlen 30. It should be noted that spliced-reads do not span introns that are interrupted by

833 an IES due to the low maximum length, however such cases are not expected to occur

834 often.

835 *Gene prediction and domain annotation*

836 To predict protein-coding genes in IESs, non-IES nucleotides in the MAC+IES assembly

837 were first masked with 'N's. The Intronarrator pipeline (https://github.com/Swart-

29

838    lab/Intronarrator), a wrapper around Augustus (Stanke and Waack, 2003), was run with the

839    same parameters as for the *B. stoltei* MAC genome, i.e. a cut-off of 0.2 for the fraction of

840    spliced reads covering a potential intron, and ≥10 reads to call an intron (Singh et al., 2021).

841    Without masking, gene predictions around IESs were poor, with genuine MDS-limited genes

842    (with high RNA-seq coverage) frequently incorrectly extended into IES regions. The

843    possibility of genes spanning IES boundaries was not catered for.

844    Domain annotations for diagrams were generated with the InterproScan 5.44-79.0 pipeline

845    (Jones et al., 2014) incorporating HMMER (v3.3, Nov 2019, hmmscan) (Eddy, 2011).

846    For comparison of transposase-related domain content in MAC vs. MIC, reference

847    sequences were obtained from public databases for *Paramecium tetraurelia*

848    (https://paramecium.i2bc.paris-

849    saclay.fr/files/Paramecium/tetraurelia/51/annotations/ptetraurelia_mac_51_with_ies/),

850    *Tetrahymena thermophila* (http://www.ciliate.org/system/downloads/3-upd-cds-fasta-

851    2021.fasta), and *Oxytricha trifallax*

852    (https://oxy.ciliate.org/common/downloads/oxy/Oxy2020_CDS.fasta,

853    https://knot.math.usf.edu/mds_ies_db/data/gff/oxytri_mic_non_mds.gff). IES gene prediction

854    in *Blepharisma* was hampered by intermittent polynucleotide tract length errors, due to the

855    assembly of IESs from PacBio CLR reads. To mitigate this, a six-frame translation of the

856    MIC-limited genome regions was performed using a custom script, then scanned against the

857    Pfam-A database 32.0 (release 9) (Mistry et al., 2021) with hmmscan (HMMER), with i-E-

858    value cutoff ≤$10^{-6}$. Domains were annotated from the MAC genome with three different

859    methods: using published coding sequences ("cds" in Table S4), six-frame translations

860    ("6ft"), and six-frame translations split on stop codons ("6ft_split").

*Repeat annotation and clustering*

862    To evaluate the repetitive sequence content in IESs, we applied a repeat prediction and

863    annotation to the combined MAC+IES assembly, instead of clustering whole IESs by

864    sequence similarity. This was so that: (i) Repeats shared between the MDS and IES could

865    be identified. (ii) Complex structures such as nested repeats could be detected. (iii) Repeat

866    families were predicted *de novo*, permitting discovery of novel elements. (iv) Repeats did not

867    have to be strictly identical to be grouped into a family.

868    Interspersed repeat element families were predicted from the MAC+IES genome assembly

869    with RepeatModeler v2.0.1 (default settings, random number seed 12345) with the following

870    dependencies: rmblast v2.9.0+ (http://www.repeatmasker.org/RMBlast.html), TRF 4.09

871    (Benson, 1999), RECON (Bao and Eddy, 2002), RepeatScout 1.0.6 (Price et al., 2005),

872 RepeatMasker v4.1.1 (http://www.repeatmasker.org/RMDownload.html). Repeat families

873 were also classified in the pipeline by RepeatClassifier v2.0.1 through comparison against

874 RepeatMasker's repeat protein database and the Dfam database. Consensus sequences of

875 the predicted repeat families, produced by RepeatModeler, were then used to annotate

876 repeats in the MAC+IES assembly with RepeatMasker, using rmblast as the search engine.

877 The consensus sequences for rnd-1_family-0 and rnd-1_family-73 were manually curated for

878 downstream analyses. For rnd-1_family-0 (BogoMITE) the original consensus predicted by

879 RepeatModeler for rnd-1_family-0 was 784 bp long, but this was a spurious inverted

880 duplication of the basic ~390 bp unit; the duplication had been favored in the construction of

881 the consensus because RepeatModeler attempts to find the longest possible match to

882 represent each family. For family rnd-1_family-73 (containing BstTc1 transposon), the actual

883 repeat unit was longer than the boundaries predicted by RepeatModeler. In most IESs that

884 contain this repeat (19 of 22), it was flanked by and partially overlapping with short repeat

885 elements from families rnd-4_family-1308 and rnd-1_family-117, which are spurious

886 predictions. Repeat unit boundaries were manually defined by alignment of full length

887 repeats and their flanking regions.

888 Terminal inverted repeats of selected repeat element families were identified by aligning the

889 consensus sequence from RepeatModeler, and/or selected full-length elements, with their

890 respective reverse complements using MAFFT (Katoh and Standley, 2013) (plugin version

891 distributed with Geneious).

892 TIRs from the Dfam DNA transposon termini signatures database (v1.1,

893 https://www.dfam.org/releases/dna_termini_1.1/dna_termini_1.1.hmm.gz) (Storer et al.,

894 2021) were searched with hmmsearch (HMMer v3.2.1) against the IES sequences, to

895 identify matches to TIR signatures of major transposon subfamilies.

896 *Phylogenetic analysis of Tc1/Mariner-superfamily transposases*

897 Repeat family rnd-1_family-1 was initially classified as a "TcMar/Tc2" family transposable

898 element by RepeatClassifier. 30 full length copies (>95% of the consensus length) were

899 annotated by RepeatMasker, all of which fell within IESs and contained CDS predictions.

900 However, CDSs were of varying lengths because of frameshifts caused by indels, which may

901 be biological or due to assembly error; nonetheless, the nucleotide sequences had high

902 pairwise identity (about 98%, except for one outlier). We chose BSTOLATCC_MIC4025 as

903 the representative CDS sequence for phylogenetic analysis because it was one of the

904 longest predicted and both predicted Pfam domains (HTH_Tnp_Tc5 and DDE_1) appeared

905 to be intact.

906   For repeat family rnd-1_family-73, the initial classification was "DNA/TcMar-Tc1". As

907   described above, CDS predictions were of variable lengths, and the longest CDSs were not

908   necessarily the best versions of the sequence because of potential frameshift errors. For

909   phylogenetic analysis, we chose BSTOLATCC_MIC48344 as the representative copy,

910   because a complete *DDE_3* Pfam domain was predicted by HMMER that could align with

911   other DDE/D domains from reference alignments described below.

912   The representative CDSs of the rnd-1_family-1 and rnd-1_family-73 transposases were

913   aligned with MAFFT (E-INS-i mode) against a published DDE/D domain reference alignment

914   (Supporting Information Dataset_S01 of (Yuan and Wessler, 2011)) to identify the residues

915   at the conserved catalytic triad and the amino acid distance between the conserved

916   residues.

917   For the phylogenetic analysis of the DDE/D domains in the Tc1/Mariner superfamily, both

918   MAC- and MIC-limited genes containing DDE_1 and DDE_3 domains were separately

919   aligned for each Pfam domain with MAFFT v7.450 (algorithm: E-INS-i, scoring matrix:

920   BLOSUM62, Gap open penalty: 1.53) and trimmed to the DDE/D domain with Geneious and

921   incomplete domains were removed. As reference, 204 sequences from a published

922   alignment (Additional File 4 of (Dupeyron et al., 2020)) were selected to represent the 53

923   groups defined in that study, choosing only complete domains (with all three conserved

924   catalytic residues) and all *Oxytricha trifallax* TBE and *Euplotes crassus* Tec transposase

925   sequences. Thirteen *Paramecium* Tc1/Mariner DDE/D domain consensus sequences were

926   added (Additional File 4 of (Guérin et al., 2017)). Sequences were aligned with MAFFT (E-

927   INS-i mode) and trimmed to only the DDE/D domain boundaries with Geneious. Phylogeny

928   was inferred with FastTree2 v2.1.11 (Price et al., 2010) using the WAG substitution model.

929   The tree was visualized with Dendroscope v3.5.10 (Huson and Scornavacca, 2012), rooted

930   with bacterial IS630 sequences as outgroup

931   *Phylogenetic analysis of retrotransposon-derived sequences*

932   All the nucleotide sequences ≥500 bp for the repeat families identified by RepeatClassifier

933   as LINE or LINE/RTE-x: rnd-1_family-273, rnd-1_family-276 and rnd-4_family-193 were

934   aligned to one another with MAFFT v7.450 (automatic algorithm) (Katoh and Standley,

935   2013), with the option to automatically determine sequence direction (via the MAFFT plugin

936   for Geneious Prime (Kearse et al., 2012)). Since the alignment appeared to be poor between

937   the rnd-4-family-193 sequences and the rest, we generated separate alignments for this

938   family from the other two, also with MAFFT (E-INS-i mode). Maximum likelihood phylogenies

939  were generated by PhyML (Guindon et al., 2010) version 3.3.20180621 with the HKY85

940  substitution model.

## Data availability

942  Annotated draft MAC+IES genome for *Blepharisma stoltei* strain ATCC 30299 (European

943  Nucleotide Archive (ENA) Bioproject PRJEB46944 under accession GCA_914767885). IES

944  sequences and annotations, MAC gene predictions with intervening IESs, and gene

945  predictions within IESs (EDMOND, doi:10.17617/3.83; genome browser,

946  https://bleph.ciliate.org. Sequencing data for the MIC-enriched nuclear fractions (PacBio

947  CLR reads: ENA accession ERR6510520 and ERR6548140; BGI-seq reads: ENA

948  accessions ERR6474675, ERR6496962, ERR6497067, ERR6501836). Small RNA libraries

949  from developmental time series (ENA Bioproject PRJEB47200 under accessions

950  ERR6565537-ERR6565561). Repeat family predictions and annotations by RepeatModeler

951  and RepeatMasker (EDMOND, doi:10.17617/3.82). Alignment and phylogeny of Tc1/Mariner

952  superfamily transposase domains (EDMOND, doi:10.17617/3.JLWBFM)

## Acknowledgements

## Author contributions

959  Data curation: B.K.B.S., E.C.S., M.S.1. Formal analysis: B.K.B.S., M.S.1, E.C.S., C.W.

960  Funding acquisition: N.S., E.C.S. Investigation: B.K.B.S., M.S.1., E.C.S. Methodology:

961  B.K.B.S., M.S.1., E.C.S., M.S.2, T.H., A.S., C.E. Resources: M.S.2., T.H., C.W., B.H., A.B.,

962  N.S. Software: B.K.B.S., E.C.S., M.S.1, A.B., N.S. Supervision: M.S.2, T.H., E.C.S.

963  Visualization: B.K.B.S., M.S.1., E.C.S. Writing – original draft: B.K.B.S., M.S.1., E.C.S.

964  Writing – review & editing: B.K.B.S., M.S.1, E.C.S., M.S.2, T.H., A.S., C.W., N.S.

## Declaration of interests

966  The authors declare no competing interests.

# References

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

Bao, Z., and Eddy, S.R. (2002). Automated de novo identification of repeat sequence families in sequenced genomes. Genome Res. *12*, 1269–1276.

Baudry, C., Malinsky, S., Restituito, M., Kapusta, A., Rosa, S., Meyer, E., and Bétermier, M. (2009). PiggyMac, a domesticated piggyBac transposase involved in programmed genome rearrangements in the ciliate Paramecium tetraurelia. Genes Dev. *23*, 2478–2483.

Benson, G. (1999). Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids Res. *27*, 573–580.

Bestor, T.H. (1990). DNA methylation: evolution of a bacterial immune function into a regulator of gene expression and genome structure in higher eukaryotes. Philos. Trans. R. Soc. Lond. B Biol. Sci. *326*, 179–187.

Bischerour, J., Bhullar, S., Denby Wilkes, C., Régnier, V., Mathy, N., Dubois, E., Singh, A., Swart, E., Arnaiz, O., Sperling, L., et al. (2018). Six domesticated PiggyBac transposases together carry out programmed DNA elimination in Paramecium. ELife *7*.

Catania, F., Rothering, R., and Vitali, V. (2021). One cell, two gears: extensive somatic genome plasticity accompanies high germline genome stability in *Paramecium*. Genome Biol. Evol. *13*.

Chalker, D.L., Meyer, E., and Mochizuki, K. (2013). Epigenetics of ciliates. Cold Spring Harb. Perspect. Biol. *5*, a017764.

Cheng, C.-Y., Vogt, A., Mochizuki, K., and Yao, M.-C. (2010). A domesticated piggyBac transposase plays key roles in heterochromatin dynamics and DNA cleavage during programmed DNA deletion in Tetrahymena thermophila. Mol. Biol. Cell *21*, 1753–1762.

Cheng, C.-Y., Young, J.M., Lin, C.-Y.G., Chao, J.-L., Malik, H.S., and Yao, M.-C. (2016). The piggyBac transposon-derived genes TPB1 and TPB6 mediate essential transposon-like excision during the developmental rearrangement of key genes in *Tetrahymena thermophila*. Genes Dev. *30*, 2724–2736.

Chen, Q., Luo, W., Veach, R.A., Hickman, A.B., Wilson, M.H., and Dyda, F. (2020). Structural basis of seamless excision and specific targeting by piggyBac transposase. Nat. Commun. *11*, 3446.

34

1000    Chen, X., Bracht, J.R., Goldman, A.D., Dolzhenko, E., Clay, D.M., Swart, E.C., Perlman,
1001    D.H., Doak, T.G., Stuart, A., Amemiya, C.T., et al. (2014). The architecture of a scrambled
1002    genome reveals massive levels of genomic rearrangement during development. Cell *158*,
1003    1187–1198.

1004    Coyne, R.S., Lhuillier-Akakpo, M., and Duharcourt, S. (2012). RNA-guided DNA
1005    rearrangements in ciliates: is the best genome defence a good offence? Biol. Cell *104*, 309–
1006    325.

1007    Crooks, G.E., Hon, G., Chandonia, J.M., and Brenner, S.E. (2004). WebLogo: a sequence
1008    logo generator. Genome Res. *14*, 1188–1190.

1009    Dale, R.K., Pedersen, B.S., and Quinlan, A.R. (2011). Pybedtools: a flexible Python library
1010    for manipulating genomic datasets and annotations. Bioinformatics *27*, 3423–3424.

1011    Drews, F., Salhab, A., Karunanithi, S., Cheaib, M., Jung, M., Schulz, M.H., and Simon, M.
1012    (2021). Broad domains of histone marks in the highly compact *Paramecium* macronuclear
1013    genome. BioRxiv.

1014    Dupeyron, M., Baril, T., Bass, C., and Hayward, A. (2020). Phylogenetic analysis of the
1015    Tc1/mariner superfamily reveals the unexplored diversity of pogo-like elements. Mob. DNA
1016    *11*, 21.

1017    Eddy, S.R. (2011). Accelerated profile HMM searches. PLoS Comput. Biol. *7*, e1002195.

1018    Edgar, R.C. (2010). Search and clustering orders of magnitude faster than BLAST.
1019    Bioinformatics *26*, 2460–2461.

1020    Fang, W., Wang, X., Bracht, J.R., Nowacki, M., and Landweber, L.F. (2012). Piwi-interacting
1021    RNAs protect DNA against loss during *Oxytricha* genome rearrangement. Cell *151*, 1243–
1022    1255.

1023    Feng, Y., and Landweber, L.F. (2021). Transposon debris in ciliate genomes. PLoS Biol. *19*,
1024    e3001354.

1025    Feng, L., Wang, G., Hamilton, E.P., Xiong, J., Yan, G., Chen, K., Chen, X., Dui, W.,
1026    Plemens, A., Khadr, L., et al. (2017). A germline-limited piggyBac transposase gene is
1027    required for precise excision in *Tetrahymena* genome rearrangement. Nucleic Acids Res. *45*,
1028    9481–9502.

1029    Feschotte, C., and Mouchès, C. (2000). Evidence that a family of miniature inverted-repeat
1030    transposable elements (MITEs) from the *Arabidopsis thaliana* genome has arisen from a
1031    pogo-like DNA transposon. Mol. Biol. Evol. *17*, 730–737.

1032    Feschotte, C., Zhang, X., and Wessler, S.R. (2002). Miniature inverted-repeat trasnposable

35

1033   elements and their relationship to established DNA transposons. In Mobile DNA II, N.L.

1034   Craig, R. Craigie, M. Gellert, and A.M. Lambowitz, eds. (Washington, D.C.: ASM Press), pp.

1035   1147–1158.

1036   Feschotte, C., Osterlund, M.T., Peeler, R., and Wessler, S.R. (2005). DNA-binding specificity

1037   of rice mariner-like transposases and interactions with Stowaway MITEs. Nucleic Acids Res.

1038   *33*, 2153–2165.

1039   Fillingham, J.S., Thing, T.A., Vythilingum, N., Keuroghlian, A., Bruno, D., Golding, G.B., and

1040   Pearlman, R.E. (2004). A non-long terminal repeat retrotransposon family is restricted to the

1041   germ line micronucleus of the ciliated protozoan *Tetrahymena thermophila*. Eukaryotic Cell

1042   *3*, 157–169.

1043   Fritz, G. (2000). Human APE/Ref-1 protein. Int. J. Biochem. Cell Biol. *32*, 925–929.

1044   Gao, F., and Katz, L.A. (2014). Phylogenomic analyses support the bifurcation of ciliates into

1045   two major clades that differ in properties of nuclear division. Mol. Phylogenet. Evol. *70*, 240–

1046   243.

1047   Gao, B., Wang, Y., Diaby, M., Zong, W., Shen, D., Wang, S., Chen, C., Wang, X., and Song,

1048   C. (2020). Evolution of pogo, a separate superfamily of IS630-Tc1-mariner transposons,

1049   revealing recurrent domestication events in vertebrates. Mob. DNA *11*, 25.

1050   Giese, A.C. (1973). *Blepharisma*: The Biology of a Light-sensitive Protozoan (Stanford

1051   University Press).

1052   Grewal, S.I.S., and Jia, S. (2007). Heterochromatin revisited. Nat. Rev. Genet. *8*, 35–46.

1053   Guérin, F., Arnaiz, O., Boggetto, N., Denby Wilkes, C., Meyer, E., Sperling, L., and

1054   Duharcourt, S. (2017). Flow cytometry sorting of nuclei enables the first global

1055   characterization of *Paramecium* germline DNA and transposable elements. BMC Genomics

1056   *18*, 327.

1057   Guermonprez, H., Loot, C., and Casacuberta, J.M. (2008). Different strategies to persist: the

1058   pogo-like Lemi1 transposon produces miniature inverted-repeat transposable elements or

1059   typical defective elements in different plant genomes. Genetics *180*, 83–92.

1060   Guindon, S., Dufayard, J.-F., Lefort, V., Anisimova, M., Hordijk, W., and Gascuel, O. (2010).

1061   New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the

1062   performance of PhyML 3.0. Syst. Biol. *59*, 307–321.

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

1065   of *Tetrahymena thermophila* and relationship to the massively rearranged somatic genome.

1066    ELife *5*.

1067    Han, J.S. (2010). Non-long terminal repeat (non-LTR) retrotransposons: mechanisms, recent

1068    developments, and unanswered questions. Mob. DNA *1*, 15.

1069    Hartl, D.L., Lohe, A.R., and Lozovskaya, E.R. (1997). Modern thoughts on an ancyent

1070    marinere: function, evolution, regulation. Annu. Rev. Genet. *31*, 337–358.

1071    Harumoto, T., Miyake, A., Ishikawa, N., Sugibayashi, R., Zenfuku, K., and Iio, H. (1998).

1072    Chemical defense by means of pigmented extrusomes in the ciliate *Blepharisma japonicum*.

1073    Eur. J. Protistol. *34*, 458–470.

1074    Herrick, G., Cartinhour, S., Dawson, D., Ang, D., Sheets, R., Lee, A., and Williams, K.

1075    (1985). Mobile elements bounded by C4A4 telomeric repeats in *Oxytricha fallax*. Cell *43*,

1076    759–768.

1077    Huff, J.T., Zilberman, D., and Roy, S.W. (2016). Mechanism for DNA transposons to

1078    generate introns on genomic scales. Nature *538*, 533–536.

1079    Huson, D.H., and Scornavacca, C. (2012). Dendroscope 3: an interactive tool for rooted

1080    phylogenetic trees and networks. Syst. Biol. *61*, 1061–1067.

1081    Jahn, C.L., Doktor, S.Z., Frels, J.S., Jaraczewski, J.W., and Krikau, M.F. (1993). Structures

1082    of the *Euplotes crassus* Tec1 and Tec2 elements: identification of putative transposase

1083    coding regions. Gene *133*, 71–78.

1084    Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen,

1085    J., Mitchell, A., Nuka, G., et al. (2014). InterProScan 5: genome-scale protein function

1086    classification. Bioinformatics *30*, 1236–1240.

1087    Katoh, K., and Standley, D.M. (2013). MAFFT multiple sequence alignment software version

1088    7: improvements in performance and usability. Mol. Biol. Evol. *30*, 772–780.

1089    Kearse, M., Moir, R., Wilson, A., Stones-Havas, S., Cheung, M., Sturrock, S., Buxton, S.,

1090    Cooper, A., Markowitz, S., Duran, C., et al. (2012). Geneious Basic: an integrated and

1091    extendable desktop software platform for the organization and analysis of sequence data.

1092    Bioinformatics *28*, 1647–1649.

1093    Kim, D., Paggi, J.M., Park, C., Bennett, C., and Salzberg, S.L. (2019). Graph-based genome

1094    alignment and genotyping with HISAT2 and HISAT-genotype. Nat. Biotechnol. *37*, 907–915.

1095    Klobutcher, L.A., and Herrick, G. (1995). Consensus inverted terminal repeat sequence of

1096    *Paramecium* IESs: resemblance to termini of Tc1-related and *Euplotes* Tec transposons.

1097    Nucleic Acids Res. *23*, 2006–2013.

1098    Klobutcher, L.A., and Herrick, G. (1997). Developmental genome reorganization in ciliated

1099    protozoa: the transposon link. Prog. Nucleic Acid Res. Mol. Biol. *56*, 1–62.

1100    Kubota, T., Tokoroyama, T., Tsukuda, Y., Koyama, H., and Miyake, A. (1973). Isolation and

1101    structure determination of blepharismin, a conjugation initiating gamone in the ciliate

1102    blepharisma. Science *179*, 400–402.

1103    Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. Nat.

1104    Methods *9*, 357–359.

1105    Lauth, M.R., Spear, B.B., Heumann, J., and Prescott, D.M. (1976). DNA of ciliated protozoa:

1106    DNA sequence diminution during macronuclear development of *Oxytricha*. Cell *7*, 67–74.

1107    Le Mouël, A., Butler, A., Caron, F., and Meyer, E. (2003). Developmentally regulated

1108    chromosome fragmentation linked to imprecise elimination of repeated sequences in

1109    paramecia. Eukaryotic Cell *2*, 1076–1090.

1110    Lepère, G., Nowacki, M., Serrano, V., Gout, J.-F., Guglielmi, G., Duharcourt, S., and Meyer,

1111    E. (2009). Silencing-associated and meiosis-specific small RNA pathways in *Paramecium*

1112    *tetraurelia*. Nucleic Acids Res. *37*, 903–915.

1113    Lhuillier-Akakpo, M., Frapporti, A., Denby Wilkes, C., Matelot, M., Vervoort, M., Sperling, L.,

1114    and Duharcourt, S. (2014). Local effect of enhancer of zeste-like reveals cooperation of

1115    epigenetic and cis-acting determinants for zygotic genome rearrangements. PLoS Genet.

1116    *10*, e1004665.

1117    Liao, Y., Smyth, G.K., and Shi, W. (2014). featureCounts: an efficient general purpose

1118    program for assigning sequence reads to genomic features. Bioinformatics *30*, 923–930.

1119    Liu, Y., Taverna, S.D., Muratore, T.L., Shabanowitz, J., Hunt, D.F., and Allis, C.D. (2007).

1120    RNAi-dependent H3K27 methylation is required for heterochromatin formation and DNA

1121    elimination in Tetrahymena. Genes Dev. *21*, 1530–1545.

1122    Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics *34*,

1123    3094–3100.

1124    Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis,

1125    G., Durbin, R., and 1000 Genome Project Data Processing Subgroup (2009). The Sequence

1126    Alignment/Map format and SAMtools. Bioinformatics *25*, 2078–2079.

1127    Lynn, D.H. (2010). The Ciliated Protozoa (Dordrecht: Springer Netherlands).

1128    Malone, C.D., Anderson, A.M., Motl, J.A., Rexer, C.H., and Chalker, D.L. (2005). Germ line

1129    transcripts are processed by a Dicer-like protein that is essential for developmentally

1130    programmed genome rearrangements of *Tetrahymena thermophila*. Mol. Cell. Biol. *25*,

1131    9151–9164.

1132    Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G.A., Sonnhammer, E.L.L.,

1133    Tosatto, S.C.E., Paladin, L., Raj, S., Richardson, L.J., et al. (2021). Pfam: The protein

1134    families database in 2021. Nucleic Acids Res. *49*, D412–D419.

1135    Mitra, R., Li, X., Kapusta, A., Mayhew, D., Mitra, R.D., Feschotte, C., and Craig, N.L. (2013).

1136    Functional characterization of piggyBat from the bat *Myotis lucifugus* unveils an active

1137    mammalian DNA transposon. Proc Natl Acad Sci USA *110*, 234–239.

1138    Miyake, A., and Beyer, J. (1974). Blepharmone: a conjugation-inducing glycoprotein in the

1139    ciliate blepharisma. Science *185*, 621–623.

1140    Miyake, A., Rivola, V., and Harumoto, T. (1991). Double paths of macronucleus

1141    differentiation at conjugation in *Blepharisma japonicum*. Eur. J. Protistol. *27*, 178–200.

1142    Mochizuki, K., and Gorovsky, M.A. (2005). A Dicer-like protein in *Tetrahymena* has distinct

1143    functions in genome rearrangement, chromosome segregation, and meiotic prophase.

1144    Genes Dev. *19*, 77–89.

1145    Mochizuki, K., and Kurth, H.M. (2013). Loading and pre-loading processes generate a

1146    distinct siRNA population in *Tetrahymena*. Biochem. Biophys. Res. Commun. *436*, 497–502.

1147    Mochizuki, K., Fine, N.A., Fujisawa, T., and Gorovsky, M.A. (2002). Analysis of a piwi-related

1148    gene implicates small RNAs in genome rearrangement in tetrahymena. Cell *110*, 689–699.

1149    Mortazavi, A., Williams, B.A., McCue, K., Schaeffer, L., and Wold, B. (2008). Mapping and

1150    quantifying mammalian transcriptomes by RNA-Seq. Nat. Methods *5*, 621–628.

1151    Nowacki, M., Higgins, B.P., Maquilan, G.M., Swart, E.C., Doak, T.G., and Landweber, L.F.

1152    (2009). A functional role for transposases in a large eukaryotic genome. Science *324*, 935–

1153    938.

1154    Prescott, D.M., and Greslin, A.F. (1992). Scrambled actin I gene in the micronucleus of

1155    *Oxytricha nova*. Dev. Genet. *13*, 66–74.

1156    Price, A.L., Jones, N.C., and Pevzner, P.A. (2005). De novo identification of repeat families

1157    in large genomes. Bioinformatics *21 Suppl 1*, i351-8.

1158    Price, M.N., Dehal, P.S., and Arkin, A.P. (2010). FastTree 2 — approximately maximum-

1159    likelihood trees for large alignments. PLoS ONE *5*, e9490.

1160    Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing

1161    genomic features. Bioinformatics *26*, 841–842.

1162    Repak, A.J. (1968). Encystment and excystment of the heterotrichous ciliate *Blepharisma*

1163    *stoltei* Isquith. Journal of Protozoology *5*, 407–412.

1164    Rognes, T., Flouri, T., Nichols, B., Quince, C., and Mahé, F. (2016). VSEARCH: a versatile

1165    open source tool for metagenomics. PeerJ *4*, e2584.

1166    Rzeszutek, I., Maurer-Alcalá, X.X., and Nowacki, M. (2020). Programmed genome

1167    rearrangements in ciliates. Cell. Mol. Life Sci. *77*, 4615–4629.

1168    Sandoval, P.Y., Swart, E.C., Arambasic, M., and Nowacki, M. (2014). Functional

1169    diversification of Dicer-like proteins and small RNAs required for genome sculpting. Dev. Cell

1170    *28*, 174–188.

1171    Seah, B.K.B., and Swart, E.C. (2021). BleTIES: Annotation of natural genome editing in

1172    ciliates using long read sequencing. Bioinformatics *37*, 3929–3931.

1173    Sellis, D., Guérin, F., Arnaiz, O., Pett, W., Lerat, E., Boggetto, N., Krenek, S., Berendonk, T.,

1174    Couloux, A., Aury, J.-M., et al. (2021). Massive colonization of protein-coding exons by

1175    selfish genetic elements in *Paramecium* germline genomes. PLoS Biol. *19*, e3001309.

1176    Singh, M., Seah, B.K.B., Emmerich, C., Singh, A., Woehle, C., Huettel, B., Byerly, A., Stover,

1177    N.A., Sugiura, M., Harumoto, T., et al. (2021). The *Blepharisma stoltei* macronuclear

1178    genome: towards the origins of whole genome reorganization. BioRxiv.

1179    Stanke, M., and Waack, S. (2003). Gene prediction with a hidden Markov model and a new

1180    intron submodel. Bioinformatics *19 Suppl 2*, ii215-25.

1181    Stern, A., Keren, L., Wurtzel, O., Amitai, G., and Sorek, R. (2010). Self-targeting by CRISPR:

1182    gene regulation or autoimmunity? Trends Genet. *26*, 335–340.

1183    Storer, J., Hubley, R., Rosen, J., Wheeler, T.J., and Smit, A.F. (2021). The Dfam community

1184    resource of transposable element families, sequence models, and genome annotations.

1185    Mob. DNA *12*, 2.

1186    Sugiura, M., and Harumoto, T. (2001). Identification, characterization, and complete amino

1187    acid sequence of the conjugation-inducing glycoprotein (blepharmone) in the ciliate

1188    Blepharisma japonicum. Proc Natl Acad Sci USA *98*, 14446–14451.

1189    Swart, E.C., Bracht, J.R., Magrini, V., Minx, P., Chen, X., Zhou, Y., Khurana, J.S., Goldman,

1190    A.D., Nowacki, M., Schotanus, K., et al. (2013). The *Oxytricha trifallax* macronuclear

1191    genome: a complex eukaryotic genome with 16,000 tiny chromosomes. PLoS Biol. *11*,

1192    e1001473.

1193    Swart, E.C., Wilkes, C.D., Sandoval, P.Y., Arambasic, M., Sperling, L., and Nowacki, M.

1194    (2014). Genome-wide analysis of genetic and epigenetic control of programmed DNA

1195    deletion. Nucleic Acids Res. *42*, 8970–8983.

1196    Taverna, S.D., Coyne, R.S., and Allis, C.D. (2002). Methylation of histone h3 at lysine 9

1197    targets programmed DNA elimination in tetrahymena. Cell *110*, 701–711.

1198    Udomkit, A., Forbes, S., Dalgleish, G., and Finnegan, D.J. (1995). BS a novel LINE-like

1199    element in *Drosophila melanogaster*. Nucleic Acids Res. *23*, 1354–1358.

1200    Vaser, R., Sović, I., Nagarajan, N., and Šikić, M. (2017). Fast and accurate de novo genome

1201    assembly from long uncorrected reads. Genome Res. *27*, 737–746.

1202    Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D.,

1203    Burovski, E., Peterson, P., Weckesser, W., Bright, J., et al. (2020). SciPy 1.0: fundamental

1204    algorithms for scientific computing in Python. Nat. Methods *17*, 261–272.

1205    Vitali, V., Hagen, R., and Catania, F. (2019). Environmentally induced plasticity of

1206    programmed DNA elimination boosts somatic variability in *Paramecium tetraurelia*. Genome

1207    Res. *29*, 1693–1704.

1208    Vogt, A., and Mochizuki, K. (2013). A domesticated PiggyBac transposase interacts with

1209    heterochromatin and catalyzes reproducible DNA elimination in *Tetrahymena*. PLoS Genet.

1210    *9*, e1004032.

1211    Wang, S., Zhang, L., Meyer, E., and Matz, M.V. (2010). Characterization of a group of

1212    MITEs with unusual features from two coral genomes. PLoS ONE *5*, e10700.

1213    Weiss, A. (2015). Lamarckian Illusions. Trends Ecol. Evol. *30*, 566–568.

1214    Wuitschick, J.D., Gershan, J.A., Lochowicz, A.J., Li, S., and Karrer, K.M. (2002). A novel

1215    family of mobile genetic elements is limited to the germline genome in *Tetrahymena*

1216    *thermophila*. Nucleic Acids Res. *30*, 2524–2537.

1217    Yao, M.-C., Fuller, P., and Xi, X. (2003). Programmed DNA deletion as an RNA-guided

1218    system of genome defense. Science *300*, 1581–1584.

1219    Yuan, Y.-W., and Wessler, S.R. (2011). The catalytic domain of all eukaryotic cut-and-paste

1220    transposase superfamilies. Proc Natl Acad Sci USA *108*, 7884–7889.

1221    Zahler, A.M., Neeb, Z.T., Lin, A., and Katzman, S. (2012). Mating of the stichotrichous ciliate

1222    *Oxytricha trifallax* induces production of a class of 27 nt small RNAs derived from the

1223    parental macronucleus. PLoS ONE *7*, e42371.

1224    Zhou, W., Liang, G., Molloy, P.L., and Jones, P.A. (2020). DNA methylation enables

1225    transposable element-driven genome expansion. Proc Natl Acad Sci USA *117*, 19359–
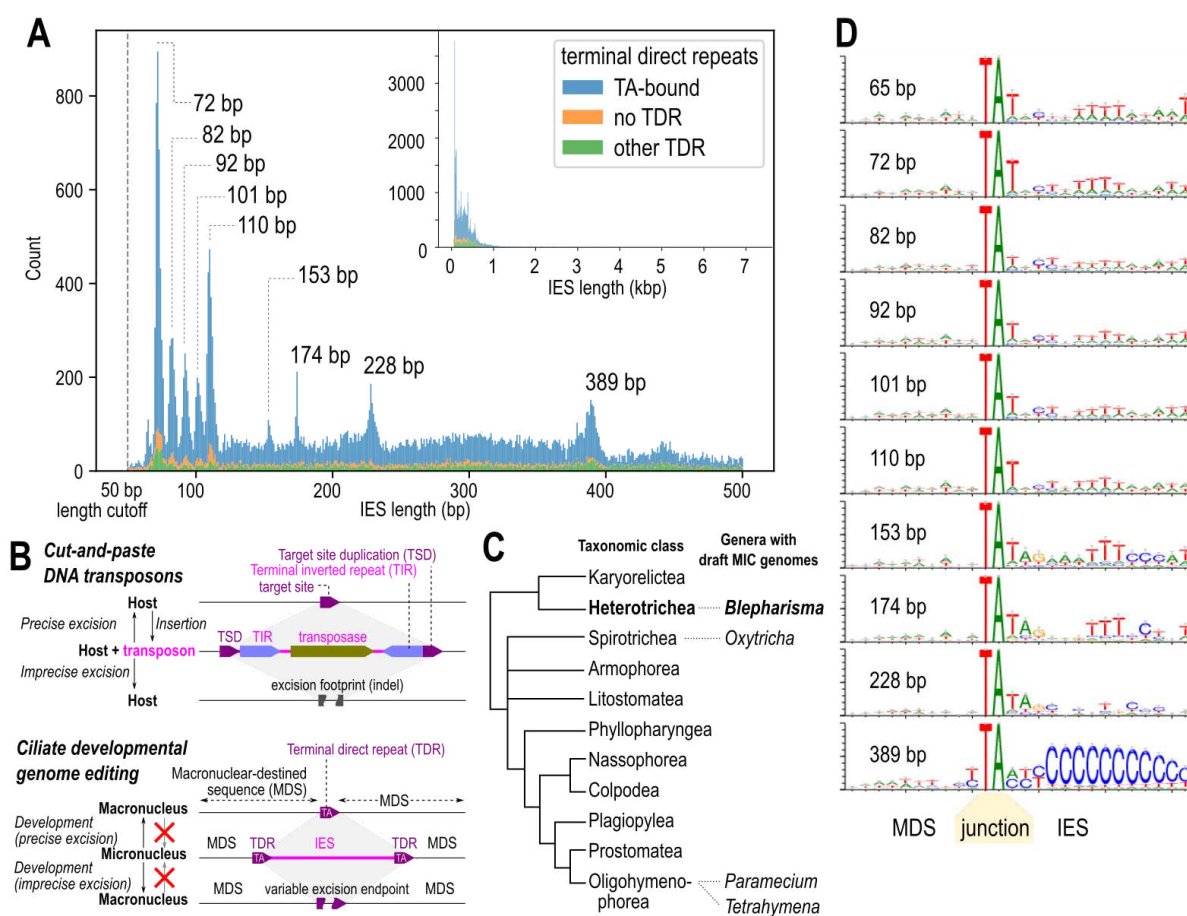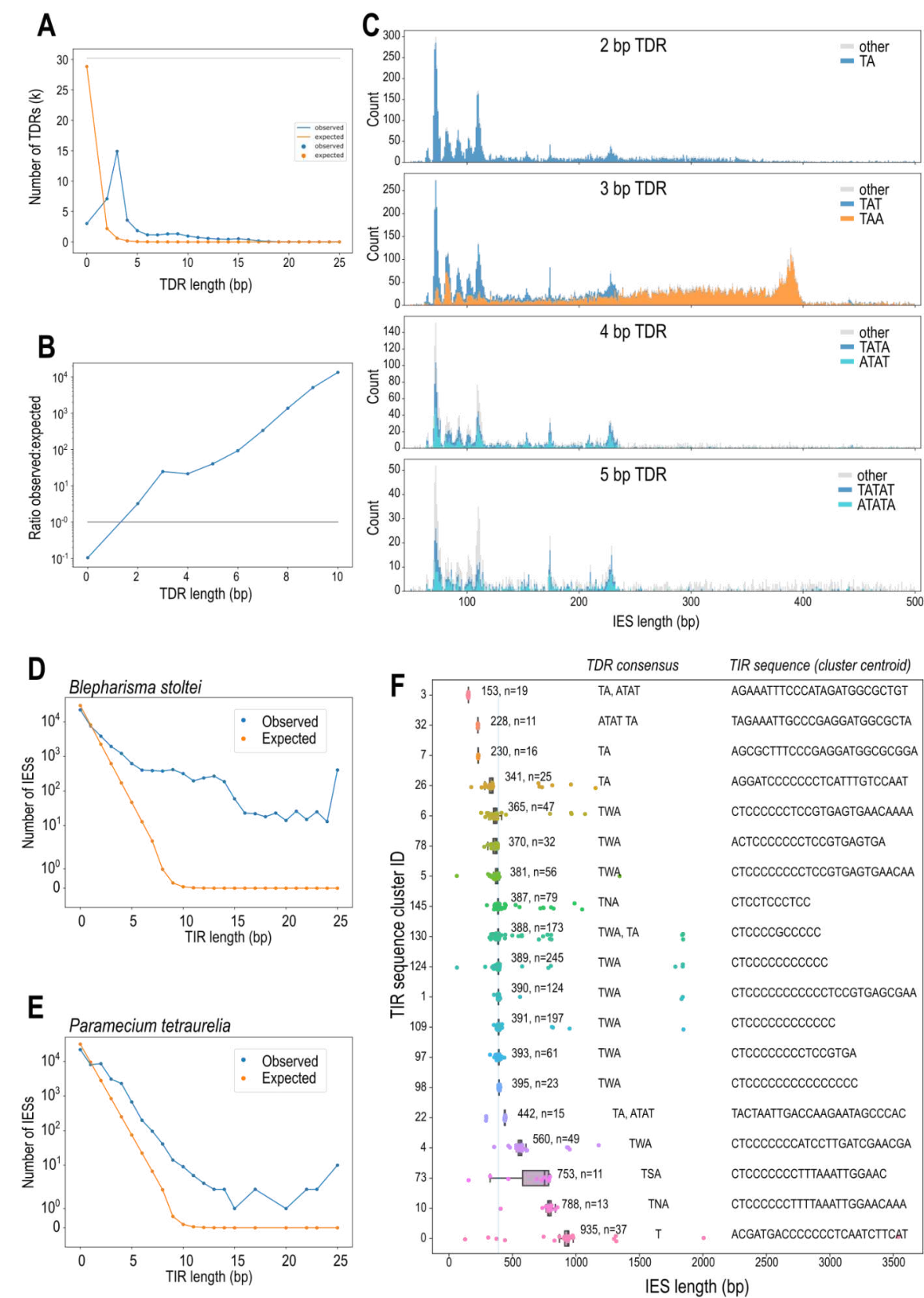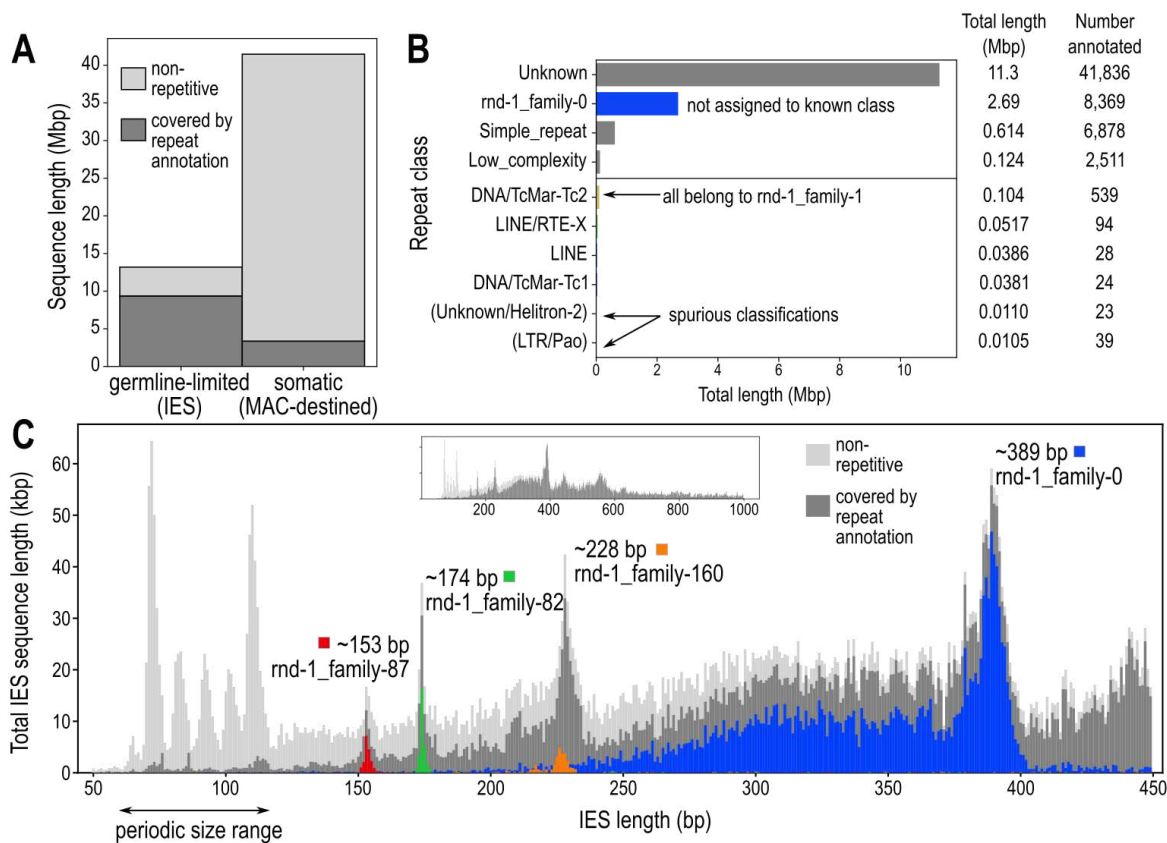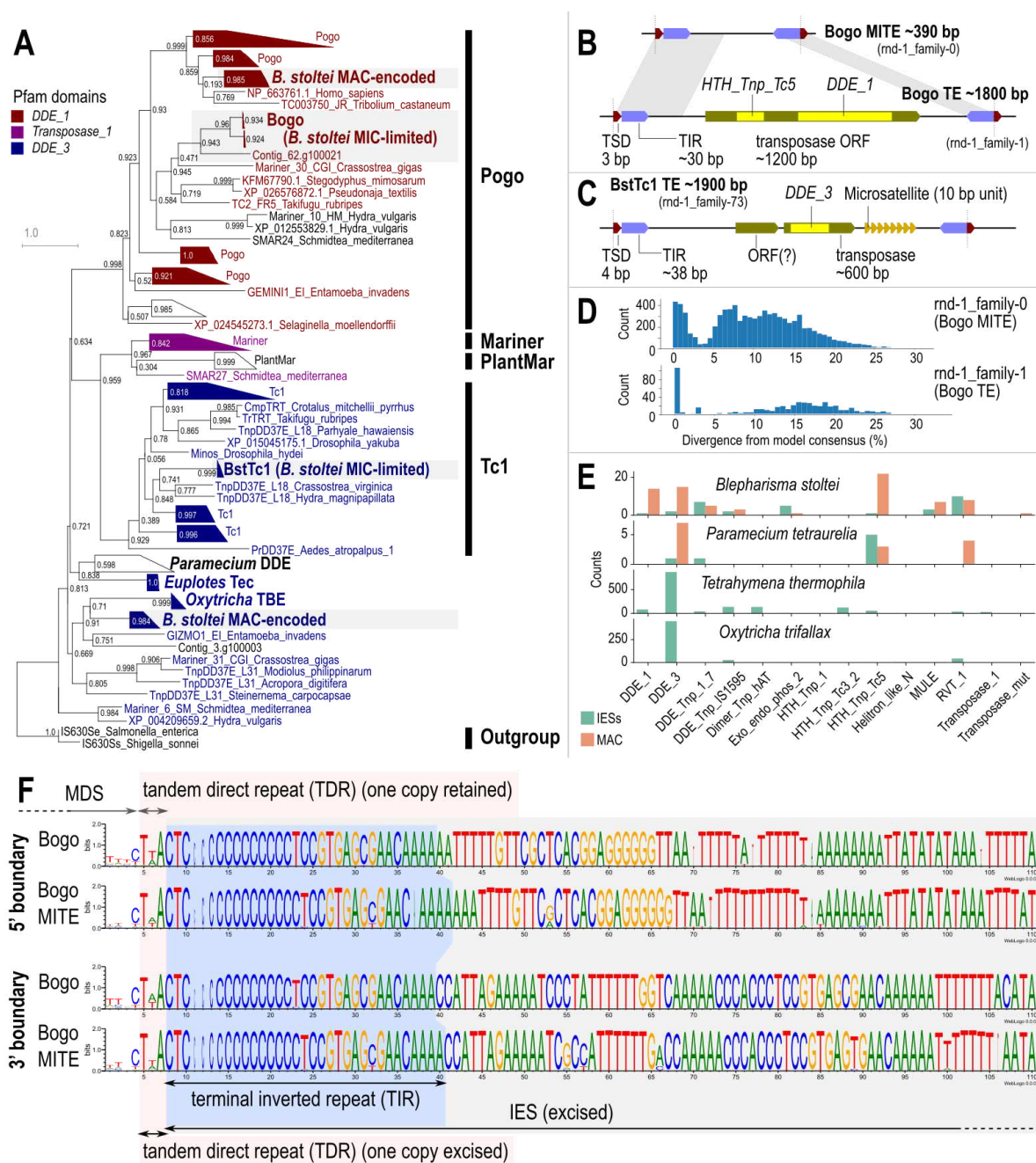
1226    19366.

**Figure 1.**

**Figure 2.**

**Figure 3.**

**Figure 4**.

**Figure 5.**

**Figure 6.**

**Figure 7.**

**Figure S1.**

**Figure S2.**

**Figure S3.**

**Figure S4**.

**Figure S5.**

**Figure S6**.