# Transcribed germline-limited coding sequences in *Oxytricha trifallax*

**Richard V. Miller[1,2], Rafik Neme[1,3], Derek M. Clay[1,2], Jananan S. Pathmanathan[1,4],**

**Michael W. Lu[1,6], V. Talya Yerlici[1,5], Jaspreet S. Khurana[1,6], and Laura F. Landweber[1,6,*]**

**[1] Department of Biochemistry and Molecular Biophysics, Columbia University, New York,**

**NY 10032, USA**

**[2] Department of Molecular Biology, Princeton University, Princeton, NJ 08544, USA**

**[3] Current address: Department of Chemistry and Biology, Universidad del Norte,**

**Barranquilla, Colombia**

**[4] Current address: School of Environmental and Biological Sciences, Rutgers University,**

**New Brunswick, NJ 08901, USA**

**[5] Current address: Department of Laboratory Medicine and Pathobiology, Faculty of**

**Medicine, University of Toronto, Toronto, ON M5G 1M1, Canada**

**[6] Current address: Strand Therapeutics, Cambridge, MA 02139, USA**

**[6] Department of Biological Sciences, Columbia University, New York, NY 10027, USA**

Running Head: *Oxytricha* transcribed germline-limited ORFs

Keywords: germline; genome rearrangement; DNA elimination; noncoding RNA; ciliate;

micronucleus

*To whom correspondence should be addressed:

**Laura F. Landweber**

Department of Biochemistry and Molecular Biophysics, Columbia University, New York, NY

10032, USA

Phone: +1 212 305-3898

Email: Laura.Landweber@columbia.edu

25

**Abstract**

27    **The germline-soma divide is a fundamental distinction in developmental biology,**

28    **and different genes are expressed in germline and somatic cells throughout metazoan life**

29    **cycles. Ciliates, a group of microbial eukaryotes, exhibit germline-somatic nuclear**

30    **dimorphism within a single cell with two different genomes. The ciliate *Oxytricha trifallax***

31    **undergoes massive RNA-guided DNA elimination and genome rearrangement to produce a**

32    **new somatic macronucleus (MAC) from a copy of the germline micronucleus (MIC). This**

33    **process eliminates noncoding DNA sequences that interrupt genes and also deletes**

34    **hundreds of germline-limited open reading frames (ORFs) that are transcribed during**

35    **genome rearrangement. Here, we update the set of transcribed germline-limited ORFs**

36    **(TGLOs) in *O. trifallax*. We show that TGLOs tend to be expressed during nuclear**

37    **development and then are absent from the somatic MAC. We also demonstrate that**

38    **exposure to synthetic RNA can reprogram TGLO retention in the somatic MAC and that**

39    **TGLO retention leads to transcription outside the normal developmental program. These**

40    **data suggest that TGLOs represent a group of developmentally regulated protein coding**

41    **sequences whose gene expression is terminated by DNA elimination.**

42

43 **Introduction**

44       Ciliates are a lineage of microbial eukaryotes characterized by functional nuclear

45 differentiation. Each ciliate cell has one or more somatic macronuclei (MAC) and one or more

46 germline micronuclei (MIC). The somatic MAC contains the somatic genome, consisting of over

47 17,000 gene-sized nanochromosomes that are transcribed throughout the organism's life cycles

48 (Swart et al. 2013; Lindblad et al. 2019). The germline genome is a fragmented and scrambled

49 version of the somatic genome that undergoes a complex process of DNA deletion and

50 rearrangement during sexual reproduction (Chen et al. 2014).

51       Previous studies have shown that *Oxytricha*'s sexual rearrangement cycle is guided by

52 several noncoding RNA pathways. In the early stages of the sexual life cycle, bidirectional

53 transcription across the length of nanochromosomes produce thousands of long template RNAs

54 from the parental MAC (Lindblad et al. 2017). These transcripts guide the rearrangement of

55 macronuclear destined sequences (MDSs) during development, and previous experiments

56 showed that injection of synthetic template RNAs could program aberrant rearrangements

57 (Nowacki et al. 2008; Bracht et al. 2017; Nowacki et al. 2011). Millions of 27-nucleotide long

58 PIWI-associated small RNAs (piRNAs) are abundant during early *Oxytricha* rearrangement and

59 interact with the *Oxytricha* PIWI ortholog Otiwi-1. These piRNAs also derive from the parental

60 MAC. Their role is to protect the sequences they target against DNA deletion during

61 development of the zygotic MAC. Injection of synthetic piRNA sequences that target internal

62 eliminated sequences (IESs) that interrupt MDSs in the MIC can prevent their deletion during

63 rearrangement and program their retention in the MAC (Fang et al. 2012). Programmed IES

64 retention is now used as a genetic tool to create somatic knockout strains in *Oxytricha* (Khurana

65 et al. 2018; Beh et al. 2019).

66      Besides IESs and transposons that are eliminated during development, *Oxytricha* has

67      other classes of germline-specific MIC DNA sequences (Chen et al. 2016). Analysis of the

68      germline MIC genome together with transcriptome-guided gene prediction previously uncovered

69      810 germline-limited protein coding genes encoded in the MIC genome (Chen et al. 2014).

70      These germline-limited genes are specifically transcribed during rearrangement, and 26% of

71      them had demonstrated translation of peptides present in a survey of one developmental time

72      point.

73      Other lineages also have germline-limited protein coding sequences, including the ciliate

74      *Tetrahymena thermophila* (Hamilton et al. 2016; Lin et al. 2016; Feng et al. 2017), the parasitic

75      roundworm *Ascaris suum* (Wang et al. 2012; Wang et al. 2017), and the sea lamprey *Petromyzon*

76      *marinus* (Bryant et al. 2016; Smith et al. 2009; Smith et al. 2012; Timoshevskiy et al. 2016;

77      Timoshevskiy et al. 2017). Protein coding sequences are discarded in all these cases, and genes

78      eliminated from somatic lineage cells are typically predicted to have functions in the germline

79      and embryogenesis (Smith et al. 2012; Bryant et al. 2016). The songbird *Taeniopygia guttata* has

80      a germline-limited chromosome that is deleted from somatic lineage cells (Pigozzi and Solari

81      1998; Pigozzi and Solari 2005; Itoh et al. 2009; Biederman et al. 2018; Kinsella et al. 2019;

82      Torgasheva et al. 2019).

83      Here, we update and expand the set of transcribed germline-limited ORFs (TGLOs) in

84      *Oxytricha* and provide functional experiments that reprogram the somatic retention of a small

85      number of TGLOs to test the hypothesis that developmental deletion is the main mechanism to

86      repress their gene expression during asexual growth. Like the previous set of germline-limited

87      genes, we show that TGLOs contain several predicted functions and conserved domains that

88      could be involved in Oxytricha development. This work also identified a locus, g111288, that is

89    retained in the somatic MAC of a subset of progeny cells, revealing an example of a strain-

90    specific macronuclear chromosome.

91 **Materials and methods**

92 **Illumina library preparation and sequencing**

93      Genomic DNA was collected from mated *O. trifallax* cells at various developmental

94 time-points using the Nucleospin genomic DNA spin column column (Machery-Nagle). Illumina

95 DNA sequencing libraries were prepared using the NEBNext Ultra II library preparation kit

96 (New England Biolabs). 2 x 250 bp paired end sequencing reads were obtained using an Illumina

97 HiSeq 2500, and remaining adapter sequences were trimmed using Trim Galore! software in the

98 Galaxy cloud computing environment.

99      Total RNA was extracted from mated *O. trifallax* cells at various developmental time-

100 points using Trizol reagent (Thermo Fisher, Waltham, MA, USA). Contaminating DNA was

101 removed using a Turbo DNase kit (Thermo Fisher, Waltham, MA, USA). Poly-adenylated

102 transcripts were enriched using the NEBNext Poly(A) mRNA Magnetic Isolation Module (New

103 England Biolabs, Ipswich, MA, USA). RNA sequencing libraries were prepared using the

104 ScriptSeq version 2 kit (Illumina, San Diego, CA, USA). 2 x 75 bp paired end sequencing reads

105 were obtained using an Illumina HiSeq 2500, and remaining adapter sequences were trimmed

106 using Trim Galore! software in the Galaxy cloud computing environment.

107 **TGLO computational prediction**

108      We predicted TGLOs using a previously published pipeline for germline-limited gene

109 prediction with some modifications (Chen et al. 2014). We predicted coding sequences with

110 AUGUSTUS (version 3.3.0) (Stanke et al. 2006) using a gene prediction model trained on *O.*

111 *trifallax* somatic MAC genes and transcripts as hints. We generated hint files for the gene

112 prediction software by mapping RNA-seq data from cells collected at various time points to the

113 germline MIC genome using HISAT2 (version 2.0.5). We ran AUGUSTUS with the options --

114     UTR=on and --alternatives-from-evidence=true. We filtered AUGUSTUS gene predictions to

115     keep only models supported by hints including at least four supporting RNA-seq reads and

116     greater than 80% of the coding sequence covered by RNA-seq reads to obtain the high

117     transcription dataset. We kept only models supported by hints including at least two supporting

118     RNA-seq reads and required greater than 20% of the coding sequence be covered by RNA-seq

119     reads to obtain the low transcription dataset. We also removed candidate sequences with more

120     than a minimal number of whole cell genomic DNA reads mapped from asexually growing

121     cultures of either parental genotype or a pool of F1 cells to ensure that MAC encoded candidates

122     were removed while accounting for the fact that some MIC encoded sequences will be present in

123     whole cell sequencing reads.

124     **DNA sequencing analysis**

125     Genomic DNA sequencing reads were aligned to the *O. trifallax* MIC genome assembly

126     using BWA-MEM (version 0.7.17) with the -M option to mark short split alignments as

127     supplementary alignments. Alignment files were processed using the Samtools software package

128     (version 0.1.20) (Li et al. 2009). FeatureCounts software (version 2.0.0) (Liao et al. 2014) was

129     used to assess the raw number of reads mapping to *O. trifallax* genome features (Burns et al.

130     2016). Relative DNA copy number changes for each genome feature were normalized using the

131     R/Bioconductor package DESeq2 (version 1.26.0) (Love et al. 2014). Heat maps showing

132     normalized DNA copy number during the developmental life cycle were generated using the

133     log2 normalized copy number values and the pheatmap R package (version 1.0.12).

134     **Transcriptome sequencing analysis**

135     Poly(A)-selected RNA sequencing reads were aligned to the *O. trifallax* MAC genome

136     assembly and MIC genome assembly using HISAT2 (version 2.0.4) and Bowtie2 in the local

7

137    alignment mode, respectively. Relative DNA copy number changes were normalized using the

138    R/Bioconductor package DESeq2. Alignment files were processed using the Samtools software

139    package (version 0.1.20) (Li et al. 2009). FeatureCounts software (version 2.0.0) (Liao et al.

140    2014) was used to assess the raw number of reads mapping to *O. trifallax* genome features

141    (Burns et al. 2016). Relative RNA expression changes for each genome feature were normalized

142    using the R/Bioconductor package DESeq2 (version 1.26.0) (Love et al. 2014). Heat maps

143    showing normalized RNA expression during the developmental life cycle were generated using

144    the log2 normalized copy number values and the pheatmap R package (version 1.0.12). Two

145    timepoints of triplicate RNA-seq reads (12 hr and 36 hr) from the late time-course were

146    previously uploaded to the European Nucleotide Archive under the project number

147    PRJEB32087.

148    **Small RNA sequencing analysis**

149         Previously sequenced Otiwi-1-dependent piRNAs (Fang et al. 2012) were aligned to the

150    *O. trifallax* MIC genome assembly using Bowtie2 (version 2.3.4.1) in the local alignment mode.

151    Alignment files were processed using the Samtools software package (version 0.1.20) (Li et al.

152    2009), and alignments were viewed in the context of the MIC genome using the Integrative

153    Genomics Viewer (version 2.7.2) (Robinson et al. 2011).

154    **Mass spectrometry analysis**

155         Raw data were analyzed using MaxQuant (version 1.6.3.4) to search against a combined

156    database containing previously published macronuclear-encoded and MIC-limited genes in

157    addition to either highly-transcribed or lowly-transcribed TGLOs (Chen et al. 2014). Searches

158    were performed using Trypsin/P as the enzyme with a maximum of two missed cleavages,

159    methionine oxidation and protein N-terminal acetylation as variable modifications, cysteine

8

160    carbamidomethylation as a fixed modification, precursor mass tolerances of 20 ppm for the first

161    search and 4.5 ppm for the main search, and a maximum FDR of 1% for both peptides and

162    proteins.

**Cell culture**

164        *Oxytricha trifallax* cells were cultured in Petri dishes or large Pyrex dishes containing

165    Pringsheim medium (0.11 mM $Na_2HPO_4$, 0.08mM $MgSO_4$, 0.85 mM $Ca(NO_3)_2$, 0.35 mM KCl,

166    pH 7.0) and fed *Chlamydomonas reinhardtii* and *Klebsiella pneumoniae* according to previously

167    published methods (Khurana et al. 2014). Matings were performed by starving the compatible

168    parental mating types 310 and 510, mixing the mating types, and diluting to a concentration of

169    about 5,000 cells per milliliter in Pringsheim medium and plating the cells in 10 cm plastic Petri

170    dishes. Matings were assessed several hours after mixing mating types by calculating the

171    percentage of paired cells per total cells.

**Reverse transcription PCR (RT-PCR)**

173        Cell cultures or mating time-courses were concentrated by centrifugation and total RNA

174    was extracted using Trizol. Turbo DNase (Thermo Fisher, Waltham, MA, USA) was used to

175    digest DNA before extracting RNA again. Eluted DNA-free total RNA was reverse transcribed

176    using oligo (dT) and AMV reverse transcriptase (New England Biolabs, Ipswich, MA, USA).

177    PCR was performed using cDNA template and Phusion High Fidelity DNA polymerase (New

178    England Biolabs, Ipswich, MA, USA).

**Nanochromosome assembly**

180        Pooled F1 cells were sequenced using Illumina sequencing. Short reads were mapped to

181    the germline MIC genome. Reads mapping to g111288 were isolated. Next, we searched for the

9

182    5' and 3' end of an arbitrary read mapping to g111288 in the other reads. We iterated the process

183    of searching for the 5' or 3' end of each read in the remaining reads until we found a read

184    terminating with a telomere repeat ($C_4A_4$). We manually assembled the sequences of the reads

185    into an g111288 nanochromosome.

186    **In vitro transcription**

187    To prepare long single-stranded RNA (ssRNA) transcripts for microinjection, PCR

188    primers were first designed to use Phusion High-Fidelity DNA polymerase (New England

189    Biolabs, Ipswich, MA, USA) to amplify the coding sequence of the desired TGLO and add a T7

190    promoter to the gene. The T7-flanked product was cloned using the TOPO TA cloning kit

191    (Thermo Fisher, Waltham, MA, USA) and Sanger sequenced (Genewiz, South Plainfield, NJ,

192    USA) to verify the insert. In vitro transcription was done using the HiScribe T7 High Yield RNA

193    Synthesis Kit according to the manufacturer's instructions (New England Biolabs, Ipswich, MA,

194    USA).

195    **RNA injection**

196    In vitro transcribed RNA was extracted using Trizol and resuspended to a concentration

197    of 3 micrograms per microliter. ssRNA was microinjected into mating cells at 12 hours post-

198    mixing according to previously published protocols (Fang et al. 2012). Post-injected cells were

199    allowed to recover in Volvic water for two days before picking single cells and plating them in

200    Volvic to establish clonal lines.

201    **5' rapid amplification of cDNA ends (5' RACE)**

202    We used a published 5' RACE protocol (Scotto-Lavino et al. 2006) with minor changes.

203    Briefly, total RNA was extracted in Trizol (Thermo Fisher, Waltham, MA, USA) and treated

204     with Turbo DNase (Ambion). One microgram of DNase-treated total RNA was reverse

205     transcribed using AMV reverse transcriptase (New England Biolabs, Ipswich, MA, USA) and a

206     gene-specific primer for either the germline-limited gene or actin II control. cDNA was poly(A)

207     tailed using terminal transferase (New England Biolabs, Ipswich, MA, USA). The A-tailed

208     cDNA was amplified using two rounds of PCR amplification using Phusion High-Fidelity DNA

209     Polymerase (New England Biolabs, Ipswich, MA, USA). The first round of amplification was

210     done over 15 cycles, the first round product was diluted 1:1000, the diluted first round product

211     was amplified over 35 cycles, and the products from the second round of amplification were

212     resolved on an agarose gel and stained with ethidium bromide (Bio-Rad, Hercules, CA).

213     **RT-qPCR**

214          As we did previously, we reverse transcribed total RNA from two different times during

215     the organism's life cycle using random hexamer primers. This cDNA was used as template in a

216     series of RT-qPCR experiments to detect the expression of either germline-limited ORF

217     candidate or actin. We used Power Sybr Green qPCR Master Mix (Thermo Fisher, Waltham,

218     MA, USA) and custom qPCR primers (Integrated DNA Technologies, Coralville, IA, USA) and

219     performed the reaction using a CFX384 Touch Real-Time PCR Detection System (Bio-Rad,

220     Hercules, CA, USA). We analyzed the Cq values using a standard curve method and compared

221     the number of transcripts in each sample to the number of small subunit mitochondrial rRNA.

11

**Southern hybridization**

222

223      1 μg of genomic DNA was resolved on a 1% agarose gel, and ethidium bromide was used

224    for visualization. MAC DNA was purified according to previously published methods (Swart et

225    al. 2013). Dilute PCR products were used as a control to approximate the expected copy number

226    in the genomic DNA lanes. The 1 Kb Plus DNA ladder (Thermo Fisher, Waltham, MA, USA)

227    was used as a size standard. After gel electrophoresis, DNA was blotted onto a nylon membrane,

228    detected using a digoxigenin-labeled DNA probe, and detected using chemiluminescence

229    according to a previously published protocol (Yerlici et al. 2019)

**Primers**

230

231    The following primers were synthesized by Integrated DNA Technologies (Coralville, IA, USA)

232    for use in this study.

233    g104149 retention fwd: 5'-CGATGATGATGCAGAGCAGTGGAGGCTTAG-3'

234    g104149 retention rev: 5'-CATATCGTGTTCATTCATGTAAGATAACTACTGCTTG-3'

235    g67186 retention fwd: 5'-CAATTCACATAATCCTCTATTTCTGCAACTTTTTCTAGAC-3'

236    g67186 retention rev: 5'-

237    GAATTATTTGTAAATACTTGACTGACTCATTGTTGATAAAATGATTTAC-3'

238    QT RACE: 5'-CCAGTGAGCAGAGTGACGAGGACTCGAGCTCAAGC-3' (Scotto-Lavino

239    2006)

240    QO RACE: 5'-CCAGTGAGCAGAGTGACG-3' (Scotto-Lavino 2006)

241    QI RACE: 5'-GAGGACTCGAGCTCAAGC-3' (Scotto-Lavino 2006)

242    Actin II RT: 5'-GTGGTGAAGTTATATCCTCTCTTGGCCAATAATG-3'

243    Actin II GSP 1: 5'-TGGCATGAGGAATTGCGTAACCTTCATAGA-3'

244    Actin II GSP 2: 5'-TCCATCTCCAGAGTCAAGCACAACACC-3'

245     g104149 RT: 5'-TTGGGTAAATTCTGGCCAACTCCCTTG-3'

246     g104149 GSP 1: 5'-CCAAGCTTCTCTGCACCTCATCCGTGAACA-3'

247     g104149 GSP 2: 5'-GTCTGCCCATCCACGATTTCACTGACC-3'

248     g67186 RT: 5'-AGCCTTGGTCCCTTCTGAGGCAG-3'

249     g67186 GSP 1: 5'-CCTGGCAAGAGCAACTTGACAGCAC-3'

250     g67186 GSP 2: 5'-GAGAGGCCAGAGGCTTCATTGCATACC-3'

251     g104149 gene qPCR fwd: 5'-CCAAGCTTCTCTGCACCTCATCCGTGAACA-3'

252     g104149 gene qPCR rev: 5'-AAGGTCAGTGAAATCGTGGATGGGCAGACT-3'

253     g67186 gene qPCR fwd: 5'-TGCAATGAAGCCTCTGGCCTCTCA-3'

254     g67186 gene qPCR rev: 5'-CCTGGCAAGAGCAACTTGACAGCAC-3'

255     g67186 upstream qPCR fwd: 5'-

256     CAATTCAATAGCACCGAATAGAAAGCTTATTTTATACAAGGATTAG-3'

257     g67186 upstream qPCR fwd: 5'-

258     CTAGATTTAATTAAAACTTGAAATGTCTACAGCCCATTAATAATTCG-3'

259     Actin II qPCR fwd: 5'-GGTGTTGTGCTTGACTCTGGAGATGGA-3'

260     Actin II qPCR rev: 5'-TGGCATGAGGAATTGCGTAACCTTCATAGA-3'

261     Mitochondrial 23S rDNA qPCR fwd: 5'-GATAGGGACCGAACTGTCTCACG-3' (Nowacki et

262     al. 2009)

263     Mitochondrial 23S rDNA qPCR rev: 5'-CATATCCTGGTTGTGAATAATCTTCCAAGGG-3'

264     (Nowacki et al. 2009)

265     Telomere primer 1: 5'-

266     ACTATAGGGCACGCGTGGTCGACGGCCCGGGCTGGTCCCCAAAACCCCAAAACCCC

267     AAAA -3' (Nowacki et al. 2008)

268     Telomere primer 2: 5'-ACTATAGGGCACGCGTGGT-3' (Nowacki et al. 2008)

13

269     g43073 TSP 1: 5'-GCCAGGTAGTTGCAAGCGCTCTCGAGAG-3'

270     g43073 TSP 2: 5'-GCTCAAAGTTTTAACTACTTGATTGAAGTGTAGATTTGGCAATC-3'

271     g104149 TSP 1: 5'-GTAAATTCTGGCCAACTCCCTTGAGTTCCAAGCTTC-3'

272     g104149 TSP 2: 5'-CAAAGTCTGCCCATCCACGATTTCACTGACCTTTG-3'

273     g93797 TSP 1: 5'-GCCCAATTCATATGCTGCTTCTTTGAGCCACTTG-3'

274     g93797 TSP 2: 5'-GATCTGGTTTTCACAGTTGAGGTAGTAGTAGTAG-3'

275     g111288 fwd PCR: 5'-CTCTACTCTCTTAGGTCTCCCTCTGCCATT-3'

276     g111288 rev PCR: 5'-AGCGGCCTGAAACTTTGTAAGGAGTAAGAT-3'

277     Actin II fwd PCR: 5'-GACTCAAATTATGTTTGAAGTCTTCAATGTACCTTGCC-3'

278     Actin II rev PCR: 5'-GTGGTGAAGTTATATCCTCTCTTGGCCAATAATG-3'

279     g111288 nanochromosome gene fwd qPCR: 5'-CAGGCCGCTTTAACTGCAACCATAGTTG-

280     3'

281     g111288 nanochromosome gene rev qPCR: 5'-

282     GGAAATTGAGCCAACTTTACAGTTAGAGCC-3'

283     g111288 nanochromosome MDS2 fwd qPCR: 5'-

284     CTTTCCTACAAATCCCCTTAAATTTCCAGTCTTGTAC-3'

285     g111288 nanochromosome MDS2 rev qPCR: 5'-

286     GTACCATGCTAGGATGTTATTGAAATCATAGAAGAC-3'

287     g111288 nanochromosome MDS4 fwd qPCR: 5'-

288     CGTCAAATTCAGTAACTAGCTCAGGTACGTC-3'

289     g111288 nanochromosome MDS4 rev qPCR: 5'-CTACCCTCCCGAGGAAAATACCTGG-3'

290     g111288 nanochromosome MDS7 fwd qPCR: 5'-

291     CTGAAATGGCTGTATCTATGGTTATTATAAAGAATTAGTG-3'

14

292   g111288 nanochromosome MDS7 rev qPCR: 5'-CAATCATCACTCTCCCTAACCGTACCTC-

293   3'

294   g111288 nanochromosome IES6 fwd qPCR: 5'-

295   GGGAAGTTATTTTATTATGAGTTTAGGTTGCATTCATTC-3'

296   g111288 nanochromosome IES6 rev qPCR: 5'-

297   GAATGAAAATGAGTGAATTAAGAATTTTAATGAAGTATGATATAACATTC-3'

298   **Bioinformatic analyses**

299        Short read DNA sequences were locally aligned to reference sequences using Bowtie 2

300   (Langmead and Salzberg 2012) or BWA-MEM. Short read RNA sequences were aligned to

301   reference sequences using HISAT2 (Kim et al. 2019). Sanger sequencing DNA reads were

302   aligned to reference sequences using the Geneious aligner in the Geneious software package

303   (version 5.9) (Biomatters, Ltd., Auckland, New Zealand) with default parameters.

304   **Data availability**

305        All cell stocks are available upon request. Illumina sequencing datasets were uploaded to

306   the NCBI Short Read Archive under the BioProject PRJNA665991. The authors affirm that all

307   data necessary for confirming the conclusions of the article are present within the manuscript and

308   figures.

15

309    **Results**

310    **Thousands of transcribed germline-limited open reading frames (TGLOs) are expressed**

311    **during development.**

312    We examined potential germline-limited coding sequences in the *Oxytricha trifallax* MIC

313    genome by searching for transcribed germline-limited open reading frames, which we refer to as

314    TGLOs. We adapted a computational pipeline originally used to identify 810 germline-limited

315    protein coding genes expressed during *Oxytricha trifallax* development (Figure 1A, left) (Chen

316    et al. 2014). First, we used Augustus gene prediction (Stanke et al. 2006) and RNA sequencing

317    hints from throughout the organism's life cycle to predict 217,805 potential coding sequences in

318    the germline genome. To exclude potential coding sequences that are present in the somatic

319    MAC genome or are not transcribed at significant levels, we restrict TGLOs to computationally

320    predicted ORFs with virtually no DNA sequencing coverage in the MAC genome of both

321    parental strains. Another requirement is that they have RNA expression in at least one timepoint

322    during the organism's life cycle. To set read mapping thresholds appropriate for the variable

323    sequencing depth of individual RNA and DNA libraries, we used a Monte Carlo approach in

324    which the predicted 217,805 candidate loci were randomly shuffled 100 times throughout the

325    germline-limited portion of the MIC genome, while recording the distribution of the number of

326    DNA and RNA reads mapped to the random loci. The distributions of DNA or RNA reads

327    mapped to randomly shuffled TGLO loci were treated as the background germline-limited

328    coverage. We required that TGLOs have a number of DNA sequencing reads mapping to them

329    from either parent or the F1 progeny that is no greater than the fifth percentile from the

330    background germline-limited coverage simulation (i.e. no reads mapped per TGLO). On the

331    other hand, highly expressed TGLOs should have RNA sequencing coverage equal to at least the

332    95th percentile from the random distribution (i.e. four reads mapped per TGLO). We also used a

16

333 lower RNA sequencing threshold (i.e. a minimum of two reads mapped per TGLO) because at

334 least one experimentally confirmed TGLO was not present in the high transcription TGLO

335 dataset. CD-HIT (Fu et al. 2012) and RepeatMasker (Smit et al. 2013) were used to cluster

336 similar sequences and to remove sequences associated with repetitive elements. The final

337 mutually exclusive datasets contained 4342 and 6296 TGLOs with high and low transcription

338 levels, respectively (Figure 1A, center). Like the previously reported germline-limited gene

339 dataset, TGLOs tend to be intron-poor, with 8.8% and 6.4% of high and low transcription

340 TGLOs, respectively, containing introns compared to 64.9% of MAC encoded genes. These

341 datasets update our previous estimates and contain 279 (213 and 66, resp.) of the 810 germline-

342 limited genes predicted in Chen et al. (2014) (Figure 1A, right) (Chen et al. 2014), with some of

343 the reduction attributed to strain-specific differences described below.

344  The previous set of 810 germline-limited genes included functional predictions (Chen et

345 al. 2014). We investigated conserved domains and putative gene functions using the functional

346 annotation tool eggNOG mapper (version 2) (Huerta-Cepas et al. 2017). 111 high transcription

347 TGLOs and 245 low transcription TGLOs mapped to conserved eggNOG orthology clusters

348 (version 5.0) (Figure 1B) (Huerta-Cepas et al. 2019). 54 TGLOs with functional predictions were

349 previously-predicted germline-limited genes (42 and 12 in high and low transcription TGLOs,

350 respectively). Predicted functions and conserved domains included several potentially involved

351 in DNA rearrangement and epigenetic regulation, including MT-A70, miRNA methylation,

352 DNA helicase, PHD zinc finger, and high mobility group.

353  Protein expression of TGLOs could also suggest a function role for a subset of predicted

354 coding sequences. One quarter (26%) of the original 810 germline-limited genes had peptides

355 identified in a nuclear proteome extracted from mid-rearrangement cells at a single timepoint

356 (Chen et al. 2014), and we queried the new TGLO datasets against this previously published

17

357     peptide dataset. 144 high and 48 low transcription TGLOs (101 and 42 newly discovered,

358     respectively) were present in this limited 40 hour proteomic survey. Several peptides from the

359     developmental survey were also mapped to TGLOs with eggNOG functional predictions (Figure

360     1B, blue text).

361         The previously published set of germline-limited genes was limited to developmental

362     gene expression, with most germline-limited genes transcribed beginning 40 hours after mixing

363     of parental cells (Chen et al. 2014). We assessed the transcription profiles of TGLOs throughout

364     the organism's developmental life cycle using a deeply sequenced set of developmental RNA

365     sequencing libraries. Two partially overlapping triplicate RNA sequencing time-courses across

366     post-zygotic development showed that RNA expression from both high (Figure 1C) and low

367     transcription TGLOs also clustered toward the later stages of rearrangement. Conversely, a

368     random sample of one thousand somatic MAC-encoded genes had a diverse set of RNA

369     expression profiles during the same time-course, suggesting that TGLOs are enriched in

370     developmental expression.

371     **TGLO genes are eliminated after gene expression.**

372         By definition, TGLO DNA sequences are restricted to the germline MIC. Since the

373     germline genome is diploid, TGLOs are present at a copy number equal to twice the number of

374     micronuclei per cell. Since DNA copy number changes significantly throughout MAC

375     development (Spear and Lauth 1976), we studied DNA copy number changes and elimination of

376     TGLOs during development. A preliminary copy-number study indicated that most TGLOs are

377     eliminated by the end of the developmental life cycle, but the DNA copy number profiles of

378     TGLOs are heterogeneous, with some showing very little copy number variation throughout

379     development, leaving it unclear whether the loci are eliminated from the developing somatic

380     MAC by the end of the sexual life cycle (Figure 2A).

18

381    Since we previously reported that telomeres are permissive to transcription in *O. trifallax*,

382    unlike in other lineages (Beh et al. 2019), we amplified several TGLO loci via telomere

383    suppression PCR (Chang et al. 2004) to determine whether telomeres are added upstream of

384    these loci before DNA elimination. We found that three out of six sampled TGLOs—

385    representing both high and low transcription TGLOs—had telomeres added near the ORF during

386    mid to late development and before their elimination from the developing somatic MAC (Figure

387    2B), consistent with their transcriptional pattern.


388    **Strain-specific germline-limited ORFs**

389    Our studies uncovered one case of a germline-encoded ORF that was also present at a

390    low copy level in the somatic MAC of one parent. The protein coding locus, OXYTRIMIC_220

391    ("g111288"),  was included in the previously reported set of 810 MIC-limited genes, but it does

392    not encode any conserved functional domains nor was it detected in a developmental mass

393    spectrometry survey (Chen et al. 2014). The initial Augustus gene prediction identified this ORF.

394    However, it was later excluded from the pipeline after incorporating new DNA sequencing

395    libraries from the parent strains and F1 progeny, which suggested that g111288 is present in the

396    somatic MAC of at least one parental strain.

397    We used PCR to amplify g111288 from parental genomic DNA to test whether the locus

398    is present in the somatic genome of either parent strain. We found that the coding sequence was

399    abundant in strain JRB510, which was not the reference strain used for genome sequencing

400    (Swart et al. 2013; Chen et al. 2014). In addition, we found that several cell lines derived from

401    either single F1 progeny or genetically manipulated F1 lines also contained g111288 at

402    detectable DNA copy levels (Figure 3A). In addition, the g111288 locus varied in DNA copy

403    level in individual F1 lines derived from different parental crosses.

404     Since g111288 appeared to be present in the MAC genome of only  parental strain,

405     JRB510, and germline limited in the reference strain JRB310, we investigated the nature of the

406     putative g111288 somatic MAC nanochromosome. Next generation sequencing reads from a

407     pool of F1 progeny cells were mapped to the germline MIC genome. This allowed assembly of

408     an entire g111288 nanochromosome with telomeres at both ends and indicated that it derives

409     from seven MDSs with the g111288 open reading frame entirely contained within the first MDS

410     (Figure 3B). RNA sequencing from developmental time-points confirmed that g111288 is

411     transcribed from 40 to 60 hr after mixing of both parental strains (Figure 3C). In addition,

412     alignment of RNA-seq reads to the other six MDSs on the g111288 nanochromosome suggested

413     the possibility that the other six MDSs of the g111288 nanochromosome could have coding

414     potential. To assess the nanochromosome's relative copy number in different cell lines, we

415     performed qPCR to target different amplicons across the g111288 nanochromosome using

416     template genomic DNA from parental cells and F1 progeny lines. A two order of magnitude

417     copy number increase was consistently observed in the JRB510 parent line relative to the

418     reference JRB310 strain (Figure 3D). Moreover, three F1 lines displayed copy levels somewhat

419     higher than the JRB510 parental strain, and the other two F1 lines appeared to have few to no

420     copies of the nanochromosome, like strain JRB310. Southern hybridization with a probe

421     targeting a MAC-specific MDS-MDS junction region confirmed the presence of the

422     nanochromosome in MAC DNA from parental strain JRB510 as well as two F1 cell lines

423     (SLC89 and SLC92; Seegmiller et al. 1996) (Figure 3E).

424     Since g111288 is present in the somatic genome of several F1 lines and at a low level in

425     one parent, we assessed whether the coding sequence is transcribed during asexual (vegetative)

426     growth. However, we did not detect any transcripts from this locus outside the middle and late

427     stages of developmental, corresponding to approximately 48 hours after mixing of mating-

20

428    compatible cells (Figure 3F). Swart et al. (2013) previously reported that many other MAC

429    nanochromosomes have developmental-specific expression (Swart et al. 2013), suggesting that

430    g111288 is a strain-specific nanochromosome, retained only in the MAC genome of JRB510 and

431    passed on to its F1 progeny.


432    **Few ncRNAs map to TGLO loci**

433        *Oxytricha*'s genome rearrangements and DNA deletion are regulated by noncoding

434    RNAs (ncRNAs). For example, Otiwi-1-bound piRNAs map to retained MDSs but not germline-

435    limited regions or IESs (Fang et al. 2012), and long template RNAs map to nanochromosomes in

436    the MAC genome (Lindblad et al. 2017). Hence, we mapped template RNAs and Otiwi-1-

437    associated piRNAs to the MIC genome and assessed their coverage in TGLO loci and the

438    g111288 locus. We found that Otiwi-1 piRNAs map to MDSs more heavily than TGLOs (Figure

439    4A). Otiwi-1 piRNAs aligned to g111288, which is retained at a low somatic copy level in one

440    parent (Figure 4B), but piRNAs are present at a reduced level compared to neighboring MDSs.

441    Template RNA coverage was also significantly higher in MDSs compared to TGLOs (Figure

442    4C), although the strain-specific TGLO g111288 lacked any template RNAs despite being

443    encoded by the JRB510 MAC (Figure 4D).


444    **Synthetic RNA injection can protect TGLO loci from genomic deletion**

445        g111288 presents an example of a potential coding sequence that is present in the somatic

446    MAC of one strain while eliminated as a TGLO in another strain. We decided to test whether

447    exposure to artificial RNAs during development could reprogram the germline-limited status of

448    TGLOs, thereby retaining them on MAC nanochromosomes. Given our previous observations

449    that exposure to non-coding RNAs can reprogram IES retention in the MAC (Fang et al. 2012;

450    Khurana et al. 2018 RNA; Beh et al. 2019) we used RNA injection to test whether exposure to

21

451    targeting RNA could reprogram the retention of two TGLO loci during development (Figure

452    5A). We targeted two TGLO loci that are encoded in the IESs of other MAC loci. One of the two

453    candidates, g67186, was previously predicted to encode a histone 2B gene (Chen et al. 2014),

454    while the other, g104149, did not contain any predicted conserved domains. The two candidates

455    are also among the highest expressed TGLOs that mapped within IESs, facilitating our strategy

456    (Figure 5A). Importantly, we also observed that our candidate TGLOs lacked Otiwi-1 piRNAs

457    and template RNAs during the sexual life cycle (Figure 5B), suggesting that the cell does not

458    endogenously encode their somatic retention during the sexual life cycle.

459        PCR from cell cultures derived from single injected cells, followed by Sanger sequencing

460    indicated that RNA injection did reprogram IES+TGLO retention in some progeny, with varying

461    levels of retention based on differences in PCR band sizes. Some products contained small

462    deletions in the retained sequence relative to the reference MIC locus, but no deletions affected

463    the ORF (Figure 5C and 5D). No F1 lines from uninjected WT parental cells contained the

464    TGLO sequences, suggesting that RNA injection specifically programs the somatic DNA

465    retention (Figure 5E right and Figure 5F right).

466        RNA programmed IES retention was previously shown to be heritable after subsequent

467    sexual cycles, so we also tested whether the IES+TGLO insertions were retained after

468    backcrossing to a parental strain. PCR amplification from genomic DNA of backcrossed pools of

469    cells indicated that the retained TGLO g104149 was partially heritable for at least two more

470    generations (Figure 5E, left). The other retained TGLO, g67186, was partially heritable for one

471    backcrossed generation (Figure 5F, left). A second band corresponding to the wild-type product

472    was present in both backcrosses, consistent with the presence of WT nanochromosomes in the

473    backcrosses to the wild-type parental strain.

**Retained TGLOs are transcribed outside usual developmental program**

474

475        Our engineered strains that retain TGLO loci are unique in their ability to encode

476    previously eliminated germline sequences in their macronucleus. Genome-wide transcription

477    start site profiling in asexually growing *O. trifallax* cells showed that transcription initiation

478    typically occurs in the subtelomeric sequence of somatic nanochromosomes that encode a single

479    gene, and this is usually within approximately one hundred bases of the transcribed coding

480    sequence (Beh et al. 2019). Since the retained TGLO reading frames are nested within the

481    protein coding sequences of a flanking gene, but also retain their own putative upstream and

482    downstream regulatory sequences, we assessed the expression of retained TGLOs. We collected

483    total RNA from asexually growing cells with the retained TGLO, as well as WT parental lines,

484    and a WT developmental time-course when TGLOs are normally transcribed, and amplified

485    cDNA ends using 5' RACE (Figure 6A). We found that retained TGLO loci were now

486    transcribed during both the asexual life cycle as well as at their normal developmental pattern

487    (Figure 6B bottom and Figure 6C bottom). The size of the RACE products were similar for the

488    retained lines as well as during normal developmental expression, suggesting that the

489    endogenous TSS was used for gene expression

490        Given the structural differences between the somatic MAC nanochromosome in asexually

491    growing cells and the differentiating MAC during the sexual life cycle, the transcriptional

492    environment of the two nuclei could differ greatly. We used qRT-PCR to compare the

493    transcription levels of retained TGLO loci during the asexual life cycle vs. WT TGLO

494    expression during development, finding that the transcription level of retained TGLOs is

495    approximately an order of magnitude higher during the WT developmental timepoint compared

496    to artificial expression during the asexual life cycle in retained lines (Figures 6D and 6E).

23

497 **Discussion**

498      Here, we introduce the definition of TGLO as a transcribed germline-limited DNA

499 sequence with the ability to encode a putative protein. We show that the *O. trifallax* germline

500 MIC genome contains abundant TGLOs that are transcribed to varying degrees in WT cells

501 during development, and are then eliminated from the somatic MAC. This suggests that TGLO

502 gene expression may be regulated by DNA elimination. The conserved domains and predicted

503 functions found in TGLO datasets also support this hypothesis. Moreover, as ciliates have

504 heterochromatic MIC genomes that are not amenable to transcription (Gorovsky and Woodard

505 1969), and previous observations demonstrated that *Oxytricha*'s germline MIC lacks RNA

506 polymerase II expression (Khurana et al. 2014). Therefore, it is an attractive hypothesis that this

507 lineage may have evolved mechanisms of shutting down gene transcription by programmed

508 DNA elimination after activating gene expression during development.

509      The earlier report of 810 germline-limited genes in *O. trifallax* assumed that germline-

510 limited coding sequences would be deleted before the cell returned to the asexual life cycle

511 (Chen et al. 2014). Here we present evidence instead that the timing of DNA elimination of

512 TGLOs is heterogeneous during the sexual life cycle. Furthermore, we note the transient addition

513 of *de novo* telomeres in unexpected locations accompanying TGLO transcription, a step that

514 might activate them for transcription. Conceptually similar, in a related ciliate *Euplotes crassus*,

515 DNA processing during the sexual life cycle is responsible for modulating the transcription of

516 one of three telomerase catalytic subunit genes (Karamysheva et al. 2003). Finally, our DNA

517 sequencing results suggest that most TGLOs are indeed eliminated from the somatic MAC by the

518 end of the sexual life cycle. However, we cannot exclude the possibility that a subset of TGLOs

519 persist longer, as further research into later developmental time-points could reveal.

24

520     We also observed that at least one germline-encoded ORF, g111288, is actually present at

521     a low somatic copy level in one parental cell line. Unlike TGLOs, g111288 is variably retained

522     as a high copy nanochromosome in some F1 progeny. Presumably, the presence of ncRNAs

523     derived from one parent can program retention of the chromosome in F1 cells, but the

524     incomplete penetrance of somatic g111288 heritability correlates with its low somatic copy

525     number in the JRB510 cell line. Curiously, g111288 does not appear to be transcribed from the

526     somatic MAC in either the parent nor F1 progeny. This is unexpected because the entire coding

527     sequence is present on its own nanochromsome along with its putative upstream and downstream

528     regulatory sequences. However, it is possible that its gene expression requires other *trans*-acting

529     regulatory factors specific to the developmental life cycle.

530     The case of g111288 is also noteworthy because it appears capable of being either

531     germline-restricted or somatic-encoded. At the level of smaller MDS or IES regions, flexibility

532     between being retained vs. deleted has been observed before but on an evolutionary timescale

533     (Mollenbeck et al. 2006) rather than an intraspecies difference (Vitali et al. 2019). This feature

534     itself could contribute to the birth of new genes, since new coding sequences can sometimes arise

535     from retained noncoding sequences if transcribed and functional (Neme and Tautz 2016; Neme

536     et al. 2017). A previous study in *Tetrahymena* reported that a set of developmentally transcribed

537     somatic minichromosomes are gradually eliminated from the MAC after genome rearrangement

538     (Lin et al. 2016). Moreover, a specific minichromosome in one *Tetrahymena* species might be

539     germline-limited in another species. This *Tetrahymena* example and our functional experiments

540     that reprogram somatic TGLO retention in *O. trifallax* suggest that TGLOs might be a reservoir

541     of sequences with somatic coding potential. We can envision an evolutionary model by which

542     germline-encoded sequences can gain access to the somatic genome where they would be

543    expressed. A deeper intraspecies survey of MAC and MIC genomes, together with

544    developmental RNAseq to survey expression, would be needed to test this hypothesis.

545           Our ability to program the somatic retention of specific TGLOs via ncRNA injection is a

546    unique feature of the present study. This had the ability to unmask gene expression of targeted

547    TGLOs outside their normal developmental program. *Tetrahymena thermophila* also has non-

548    maintained chromosomes that are lost soon after expression during development and can be

549    fused to adjacent regions to program their retention in the somatic MAC (Feng et al. 2017). Here

550    we have extended this general phenomenon to *Oxytricha* and showed that somatic retention

551    subverts the cell's endogenous transcription of the gene locus. This supports the hypothesis that

552    TGLO elimination represses their gene expression. In our example the misexpression of a single

553    TGLO locus did not affect cell viability, but the ensemble of loci may need to be silenced during

554    asexual growth.

555         .

## Acknowledgements

## References

563

564     Allen, S. E., Hug, I., Pabian, S., Rzeszutek, I., Hoehener, C., and Nowacki, M. (2017). Circular

565         concatemers of ultra-short DNA segments produce regulatory RNAs. 168: 990-999.

566     Allen, S. E. and Nowacki, M. (2020). Roles of noncoding RNAs in ciliate genome architecture.

567         J. Mol. Biol. S0022-2836: 30026-30027.

568     Beh, L. Y., Debelouchina, G. T., Clay, D. M., Thompson, R. E., Lindblad, K. A., Hutton, E. R.,

569         Bracht, J. R., Sebra, R. P., Muir, T. W., and Landweber, L. F. (2019). Identification of a

570         DNA N6-adenine methyltransferase complex and its impact on chromatin organization.

571         Cell 177: 1781-1796.

572     Bryant, S. A., Herdy, J. R., Amemiya, C. T., and Smith, J. J. (2016). Characterization of

573         somatically-eliminated genes during development of the sea lamprey (*Petromyzon*

574         *marinus*). Mol. Biol. and Evol. 33: 2337-2344.

575     Burns, J., Kukushkin, D., Lindblad, K., Chen, X., Jonoska, N., and Landweber, L. F. (2016). : a

576         database of ciliate genome rearrangements. Nucleic Acids Res. 44: D703-D709.

577     Chang, W. J., Stover, N. A., Addis, V. M., and Landweber, L. F. (2004). A micronuclear locus

578         containing three protein-coding genes remains linked during macronuclear development

579         in the spirotrichous ciliate Holosticha. Protist, 155: 245-255.

580     Chen, X., Bracht, J. R., Goldman, A. D., Dolzhenko, E., Clay, D. M., Swart, E. C., Perlman, D.

581         H., Doak, T. G., Stuart, A., Amemiya, C. T., Sebra, R. P., and Landweber, L. F. (2014).

582         The architecture of a scrambled genome reveals massive levels of genomic

583         rearrangement during development. Cell 158: 1187-1198.

584    Clay, D. M., Yerlici, V. T., Villano, D. J., and Landweber, L. F. (2019). Programmed

585        Chromosome Deletion in the Ciliate Oxytricha trifallax. G3 9: 3105-3118.

586    Curtis, E. A. and Landweber, L. F. (2006). Evolution of gene scrambling in ciliate micronuclear

587        genes. Ann. N. Y. Acad. Sci. 870: 349-350.

588    Fang, W., Wang, X., Bracht, J. R., Nowacki, M., and Landweber, L. F. (2012). Piwi-interacting

589        RNAs protect DNA against loss during *Oxytricha* genome rearrangement. Cell 151:

590        1243–1255.

591    Feng, L., Wang, G., Hamilton, E. P., Xiong, J., Yan, G., Chen, K., Chen, X., Dui, W., Plemens,

592        A., Khadr, L., et al. (2017). A germline-limited piggyBac transposase gene is required for

593        precise excision in *Tetrahymena* genome rearrangement. Nucleic Acids Res. 45: 9481–

594        9502.

595    Fu, L., Niu, B., Zhu, Z., Wu, S., and Li, W. (2012). CD-HIT: accelerated for clustering the next-

596        generation sequencing data. Bioinformatics, 28: 3150-3152.

597    Gorovsky, M. A. and Woodard, J. (1969). Studies on nuclear structure and function in

598        *Tetrahymena pyriformis*: I. RNA synthesis in macro-and micronuclei. The Journal of cell

599        biology, 42: 673-682.

600    Hamilton, E. P., Kapusta, A., Huvos, P. E., Bidwell, S. L., Zafar, N., Tang, H., Hadjithomas, M.,

601        Krishnakumar, V., Badger, J. H., Caler, E. V., et al. (2016). Structure of the germline

602        genome of *Tetrahymena thermophila* and relationship to the massively rearranged

603        somatic genome. eLife 5: e19090.

604    Huerta-Cepas, J., Forslund, K., Coelho, L. P., Szklarczyk, D., Jensen, L. J., von Mering, C., and

605        Bork, P. (2017). Fast genome-wide functional annotation through orthology assignment

606        by eggNOG-mapper. Mol. Biol. Evol. 34: 2115–2122.

607     Huerta-Cepas, J., Szklarczyk, D., Heller, D., Hernández-Plaza, A., Forslund, S. K., Cook, H.,

608         Mende, D. R., Letunic, I., Rattei, T., Jensen, L. J., et al. (2019). eggNOG 5.0: a

609         hierarchical, functionally and phylogenetically annotated orthology resource based on

610         5090 organisms and 2502 viruses. Nucleic Acids Res. 47: D309–D314.

611     Itoh, Y., Kampf, K., Pigozzi, M. I., and Arnold, A. P. (2009). Molecular cloning and

612         characterization of the germline-restricted chromosome sequence in the zebra finch.

613         Chromosoma 118: 527-536.

614     Karamysheva, Z., Wang, L., Shrode, T., Bednenko, J., Hurley, L. A., and Shippen, D. E. (2003).

615         Developmentally programmed gene elimination in *Euplotes crassus* facilitates a switch in

616         the telomerase catalytic subunit. Cell, 113: 565-576.

617     Khurana, J. S., Clay, D. M., Moreira, S., Wang, X., and Landweber, L. F. (2018). Small RNA-

618         mediated regulation of DNA dosage in the ciliate Oxytricha. RNA 24: 18-29.

619     Khurana, J. S., Wang, X., Chen, X, Perlman, D. H., and Landweber, L. F. (2014). Transcription-

620         independent functions of an RNA polymerase II subunit, Rpb2, during genome

621         rearrangement in the ciliate, *Oxytricha trifallax*. Genetics 197: 839-849.

622     Kim, D., Paggi, J. M., Park, C., Bennett, C., and Salzberg, S. L. (2019). Graph-based genome

623         alignment and genotyping with HISAT2 and HISAT-genotype. Nat. Biotechnol. 37: 907–

624         915.

625     Kinsella, C. M., Ruiz-Ruano, F. J., Dion-Côté, A. M., Charles, A. J., Gossmann, T. I., Cabrero,

626         J., Kappei, D., Hemmings, N., Simons, M. J. P., Camacho, J. P. M., et al. (2019).

627         Programmed DNA elimination of germline development genes in songbirds. Nat.

628         Commun. 10: 5468.

629    Langmead, B. and Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. Nat.

630        Methods 9: 357-359.

631    Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.,

632        Durbin, R., and 1000 Genome Project Data Processing Subgroup. (2009). The sequence

633        alignment/map format and SAMtools. Bioinformatics, 25: 2078-2079.

634    Liao, Y., Smyth, G. K., and Shi, W. (2014). featureCounts: an efficient general purpose program

635        for assigning sequence reads to genomic features. Bioinformatics, 30: 923-930.

636    Lin, C. G., Lin, I. T., and Yao, M. C. (2016). Programmed minichromosome elimination as a

637        mechanism for somatic genome reduction in *Tetrahymena thermophila*. PLoS Genet. 12:

638        e1006403.

639    Lindblad, K. A., Bracht, J. R., Williams, A. E., and Landweber, L. F. (2017). Thousands of

640        RNA-cached copies of whole chromosomes are present in the ciliate *Oxytricha* during

641        development. RNA 23: 1200-1208.

642    Lindblad, K. A., Pathmanathan, J. S., Moreira, S., Bracht, J. R., Sebra, R. P., Hutton, E. R., and

643        Landweber, L. F. (2019). Capture of complete ciliate chromosomes in single sequencing

644        reads reveals widespread chromosome isoforms. BMC Genom. 20: 1037.

645    Love, M. I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and

646        dispersion for RNA-seq data with DESeq2. Genome Biol. 15: 550.

647    Meyer, G. F. and Lipps, H. J. (1981). The formation of polytene chromosomes during

648        macronuclear development of the hypotrichous ciliate Stylonychia mytilus. Chromosoma,

649        82: 309-314.

650　　　Möllenbeck, M., Cavalcanti, A. R., Jönsson, F., Lipps, H. J., and Landweber, L. F. (2006).

651　　　　　　Interconversion of germline-limited and somatic DNA in a scrambled gene. J. Mol. Evol.,

652　　　　　　63: 69-73.

653　　　Neme, R., Amador, C., Yildirim, B., McConnell, E., and Tautz D. (2017). Random sequences are

654　　　　　　an abundant source of bioactive RNAs or peptides. Nat. Ecol. Evol. 1: 0217.

655　　　Neme, R. and Tautz, D. (2016). Fast turnover of genome transcription across evolutionary time

656　　　　　　exposes entire non-coding DNA to de novo gene emergence. eLife 5: e09977.

657　　　Nowacki, M., Higgins, B. P., Maquilan, G. M., Swart, E. C., Doak, T. G., and Landweber, L. F.

658　　　　　　(2009). A functional role for transposases in a large eukaryotic genome. Science 324:

659　　　　　　925-928.

660　　　Nowacki, M., Vijayan, V., Zhou, Y., Schotanus, K., Doak, T. G., and Landweber, L. F. (2008).

661　　　　　　RNA-mediated epigenetic programming of a genome-rearrangement pathway. Nature

662　　　　　　451: 153–158.

663　　　Pigozzi, M. I. and Solari, A. J. (1998). Germ cell restriction and regular transmission of an

664　　　　　　accessory chromosome that mimics a sex body in the zebra finch, *Taeniopygia guttata*.

665　　　　　　Chromosome Res. 6: 105-113.

666　　　Pigozzi, M. I. and Solari, A. J. (2005). The germ-line-restricted chromosome in the zebra finch:

667　　　　　　recombination in females and elimination in males. Chromosoma 114: 403-409.

668　　　Prescott, D. M. (1994). The DNA of ciliated protozoa. Microbiol. Rev. 58: 233–267.

669　　　Scotto-Lavino, E., Du, G., and Frohman, M. A. (2006). 5' end cDNA amplification using classic

670　　　　　　RACE. Nat. Protoc. 1: 2555-2562.

671    Seegmiller, A., Williams, K. R., Hammersmith, R. L., Doak, T. G., Witherspoon, D., Messick,

672        T., Storjohann, L. L., and Herrick, G. (1996). Internal eliminated sequences interrupting

673        the Oxytricha 81 locus: allelic divergence, conservation, conversions, and possible

674        transposon origins. Mol. Biol. Evol. 13: 1351-1362.

675    Smit, A. F. A., Hubley, R., and Green, P. RepeatMasker Open-4.0. (2013)

676        <http://www.repeatmasker.org>.

677    Smith, J. J., Antonacci, F., Eichler, E. E., and Amemiya, C. T. (2009). Programmed loss of

678        millions of base pairs from a vertebrate genome. Proc. Natl. Acad. Sci. U.S.A. 106:

679        11212-11217.

680    Smith, J. J., Baker, C., Eichler, E. E., and Amemiya, C. T. (2012). Genetic consequences of

681        programmed genome rearrangement. Curr. Biol. 22: 1524–1529.

682    Spear, B. B. and Lauth, M. R. (1976). Polytene chromosomes of Oxytricha: biochemical and

683        morphological changes during macronuclear development in a ciliated protozoan.

684        Chromosoma 54: 1-13.

685    Stanke, M., Keller, O., Gunduz, I., Hayes, A., Waack, S., and Morgenstern, B. (2006).

686        AUGUSTUS: ab initio prediction of alternative transcripts. Nucleic Acids Res. 34:

687        W435-W439.

688    Swart, E. C., Bracht, J. R., Magrini, V., Minx, P., Chen, X., Zhou, Y., Khurana, J. S., Goldman,

689        A. D., Nowacki, M., Schotanus, K., et al. (2013). The *Oxytricha trifallax* macronuclear

690        genome: a complex eukaryotic genome with 16,000 tiny chromosomes. PLoS Biol. 11:

691        e1001473.

692    Timoshevskiy, V. A., Herdy, J. R., Keinath, M. C., and Smith, J. J. (2016). Cellular and

693         molecular features of developmentally programmed genome rearrangement in a

694         vertebrate (sea lamprey: *Petromyzon marinus*). PLoS Genet. 12: e1006103.

695    Timoshevskiy, V. A., Lampman, R. T., Hess, J. E., Porter, L. L., Smith, J. J. (2017). Deep

696         ancestry of programmed genome rearrangement in lampreys. Dev. Biol. 429: 31-34.

697    Torgasheva, A. A., Malinovskaya, L. P., Zadesenets, K. S., Karamysheva, T. V., Kizilova, E. A.,

698         Akberdina, E. A., Pristyazhnyuk, I. E., Shnaider, E. P., Volodkina, V. A., Saifitdinova, A.

699         F. et al. (2019). Germline-restricted chromosome (GRC) is widespread among songbirds.

700         Proc. Natl. Acad. Sci. U.S.A. 116: 11845-11850.

701    Vitali, V., Hagen, R., and Catania, F. (2019). Environmentally induced plasticity of programmed

702         DNA elimination boosts somatic variability in *Paramecium tetraurelia*. Genome Res.,

703         29: 1693-1704.

704    Wang, J., Gao, S., Mostovoy, Y., Kang, Y., Zagoskin, M., Sun, Y., Zhang, B., White, L. K.,

705         Easton, A., Nutman, T. B. et al. (2017). Comparative genome analysis of programmed

706         DNA elimination in nematodes. Genome Res. 27: 2001-2014.

707    Wang, J., Mitreva, M., Berriman, M., Thorne, A., Magrini, V., Koutsovoulos, G., Kumar, S.,

708         Blaxter, M. L., and Davis, R. E. (2012). Silencing of germline-expressed genes by DNA

709         elimination in somatic cells. Dev. Cell 23: 1072–1080.

710    Yerlici, V. T., Lu, M. W., Hoge, C. R., Miller, R. V., Neme, R., Khurana, J. S., Bracht, J. R., and

711         Landweber, L. F. (2019). Programmed genome rearrangements in Oxytricha produce

712         transcriptionally active extrachromosomal circular DNA. Nucleic Acids Res. 47: 9741–

713         9760.

714    **Figure legends**

715    **Figure 1: Germline-limited ORFs are expressed during *Oxytricha trifallax* genome**

716    **rearrangement.**

717    (**A**) Left: Pipeline for predicting TGLOs in *Oxytricha trifallax* germline MIC genome. Center:

718    Total number of computationally predicted candidates remaining after each pipeline step. Right:

719    Total number of MIC-limited genes (Chen et al. 2014) remaining after each pipeline step. (**B**)

720    EggNOG mapper-predicted functions and conserved domains in TGLOs. Blue text indicates that

721    peptides from the associated TGLOs were present in a single nuclear proteome surveyed during

722    rearrangement (Chen et al. 2014). (**C**) $Log_2$-normalized RNA-seq read counts of High and Low

723    transcription TGLOs and one thousand randomly selected somatic MAC-encoded genes across

724    the *Oxytricha trifallax* developmental life cycle (hours labeled post mixing of compatible mating

725    types). Color scale refers to the $log_2$-normalized RNA expression.

726    **Figure 2: TGLOs are eliminated from the developing MAC.**

727    (**A**) $Log_2$-normalized DNA copy number of High and Low transcription TGLOs across the

728    *Oxytricha trifallax* developmental life cycle. Color scale refers to the $log_2$-normalized DNA copy

729    number. (**B**) Telomere suppression PCR targeting the upstream telomere addition site of selected

730    TGLOs in genomic DNA samples collected throughout the *Oxytricha trifallax* developmental

731    life cycle.

732    **Figure 3: Parental cells can carry a strain-specific germline-limited ORF.**

733    (**A**) Top: PCR targeting g111288 or Actin II using genomic DNA from F1 lines, parent lines, and

734    other mutant F1 lines used in this study. Bottom: Genome track showing the approximate

735    location of g111288 PCR primers. Yellow: g111288, light blue: flanking MDSs. (**B**) The

736    germline genome locus containing g111288 with mapped F1 reads from a pool of asexually

737    growing F1 cells. Yellow: g111288, light blue: MDSs, dark blue: assembled g111288 MDSs

738    from pooled F1 reads, gray triangles: observed telomere addition sites. **(C)** The germline genome

739    locus (bottom) containing g111288 (yellow) and strain-specific MDSs (dark blue) with mapped

740    RNA-seq coverage (black) from several time-points during asexual growth (starved or encysted

741    cells) and hours post mixing of mating types during the sexual life cycle. **(D)** Top: Copy number

742    relative to mitochondrial rDNA based on qPCR targeting several amplicons on the g111288

743    nanochromosome, an IES within the corresponding germline locus, and two unrelated somatic

744    loci. Bottom: The germline genome locus containing g111288 with qPCR primer locations

745    indicated. Yellow: g111288, light blue: MDSs, dark blue: assembled g111288 MDSs from

746    pooled F1 reads, black arrows: qPCR primers. **(E)** Top: Southern blot of parental and F1 MAC

747    DNA detected using an MDS-MDS junction spanning DNA probe. Bottom: MIC genome track

748    showing the portions of MDSs 1 and 2 detected. **(F)** Top: RT-PCR targeting g111288 or Actin II

749    using RNA from the same cell lines as in (A). Bottom: Genome track showing the approximate

750    location of g111288 RT-PCR primers. Yellow: g111288, light blue: MDSs.


751    **Figure 4: TGLO loci have few Otiwi-1 piRNAs and template RNAs.**

752    **(A)** Distribution of normalized mapping quality-filtered Otiwi-1 piRNA read counts (Fang et al.

753    2012) mapped to High and Low transcription TGLOs compared to MDSs. Read counts were

754    normalized to reads per kilobase million (RPKM). Brackets and asterisks indicate statistically

755    significant differences between distributions. Statistical significance was assessed using the non-

756    parametric Kolmogorov–Smirnov (KS) test, and $P<0.05$ was considered statistically significant.

757    **(B)** The germline genome locus containing the strain-specific TGLO g111288 (yellow), MDSs

758    (blue), and mapped Otiwi-1-associated piRNA coverage (gray) from several time-points during

759    rearrangement. **(C)** Distribution of normalized mapping quality-filtered template RNA read

760    counts (Lindblad et al. 2017) mapped to High and Low transcription TGLOs compared to MDSs.

761    Read counts were normalized to RPKM. Brackets and asterisks indicate statistically significant

762    differences between distributions. Statistical significance was assessed using the non-parametric

763    KS test, and P<0.05 was considered statistically significant. **(D)** The germline genome locus

764    containing the strain-specific TGLO g111288 (yellow), MDSs (blue), and mapped template

765    RNA coverage (gray) from several time-points during rearrangement.


766    **Figure 5: RNA injection programs heritable TGLO retention.**

767    **(A)** Synthetic RNA injection scheme to program the retention of a TGLO (yellow) in an IES

768    between two MDSs (blue). Possible products can include telomere-capped (black)

769    nanochromosomes with the entire IES plus TGLO flanked by the MDSs of the wild-type

770    flanking locus. **(B)** The germline genome loci containing the programmed retention candidate

771    TGLOs g104149 and g67186 (yellow), MDSs (blue), and mapped piRNA or template RNA

772    coverage (gray) from several time-points during rearrangement.**(C)** Top: Cell culture PCR

773    targeting the IES containing g104149 from cell lines derived from single RNA injected mating

774    pairs. Middle: The expected retention product containing g104149 with PCR primer locations.

775    Yellow: g104149, light blue: MDSs, black arrows: PCR primers. Bottom: Sanger sequencing

776    chromatograms from PCR reactions in (B) aligned to the expected retention product containing

777    g104149 (yellow). **(D)** Top: Cell culture PCR targeting the IES containing the predicted histone

778    2B TGLO g67186 from cell lines derived from single RNA injected mating pairs. Middle: The

779    expected retention product containing g67186 with PCR primer locations. Yellow: g67186, light

780    blue: MDSs, black arrows: PCR primers. Bottom: Sanger sequencing chromatograms from PCR

781    reactions aligned to the expected retention product containing g67186 (yellow). **(E)** Top: PCR

782    targeting the IES containing g104149 using genomic DNA from parental cells, F1 retention cells,

783    F1 retention cells backcrossed to parental cells, and unmanipulated F1 lines. Bottom: The

37

784    expected retention product containing g104149 with PCR primers. Yellow: g104149, light blue:

785    MDSs, black arrows: PCR primers. **(F)** Top: PCR targeting the IES containing the predicted

786    histone 2B TGLO g67186 using genomic DNA from parental cells, F1 retention cells, F1

787    retention cells backcrossed to parental cells, and unmanipulated F1 lines. Bottom: The expected

788    retention product containing g67186 with PCR primers. Yellow: g67186, light blue: MDSs,

789    black arrows: PCR primers.


790    **Figure 6: Retained TGLOs are misexpressed during asexual life cycle.**

791    **(A)** Possible transcription start sites (black arrows) on a hypothetical rearranged somatic

792    nanochromosome after RNA injection to retain TGLOs (yellow). Green: target transcript

793    deriving from TGLO's putative upstream regulatory sequence. **(B)** Germline genome locus

794    containing g104149 (yellow) and gene-specific 5' RACE primers used to amplify transcription

795    start site. **(C)** 5' RACE products targeting the g104149 or Actin II transcription start site in RNA

796    from F1 retention cells, parental cells, and mid-rearrangement mated cells. TdT: terminal

797    transferase. **(D)** Germline genome locus containing g67186 (yellow) and gene-specific 5' RACE

798    primers used to amplify transcription start sites. **(E)** 5' RACE products targeting the g67186 or

799    Actin II transcription start site in RNA from F1 retention cells, parental cells, and mid-

800    rearrangement mated cells. TdT: terminal transferase. **(F)** g104149 or Actin II RNA transcript

801    levels based on qRT-PCR relative to mitochondrial rRNA. Error bars: standard deviation of three

802    biological replicates. **(G)** g67186 or Actin II RNA transcript levels based on qRT-PCR relative

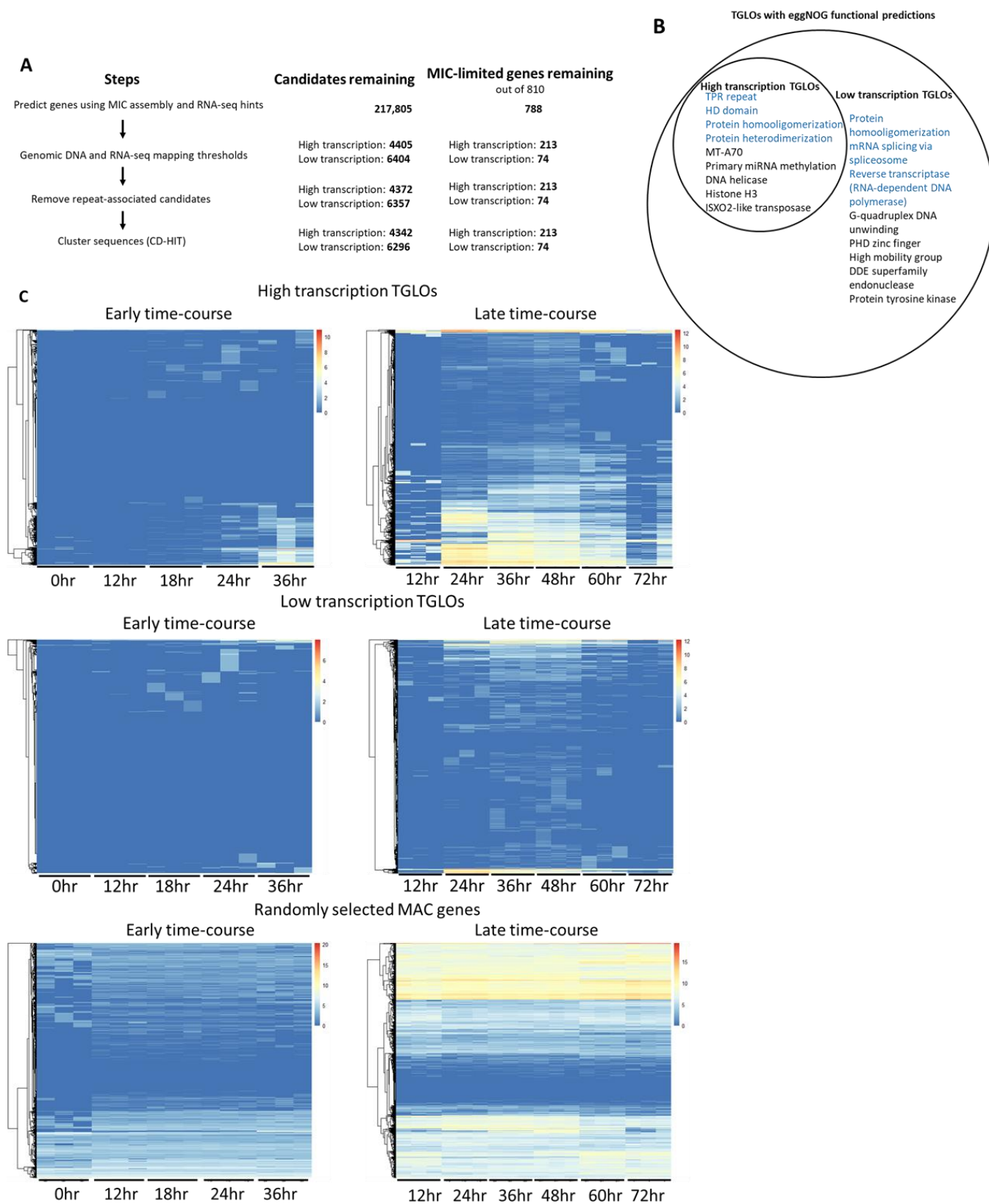803    to mitochondrial rRNA. Error bars: standard deviation of three biological replicates.

Figure 1

**A**

## High transcription TGLOs

### Early time-course



0hr  12hr  18hr  24hr  36hr

### Late time-course



24hr  36hr  48hr  60hr  72hr

## Low transcription TGLOs

### Early time-course



0hr  12hr  18hr  24hr  36hr

### Late time-course



24hr  36hr  48hr  60hr  72hr
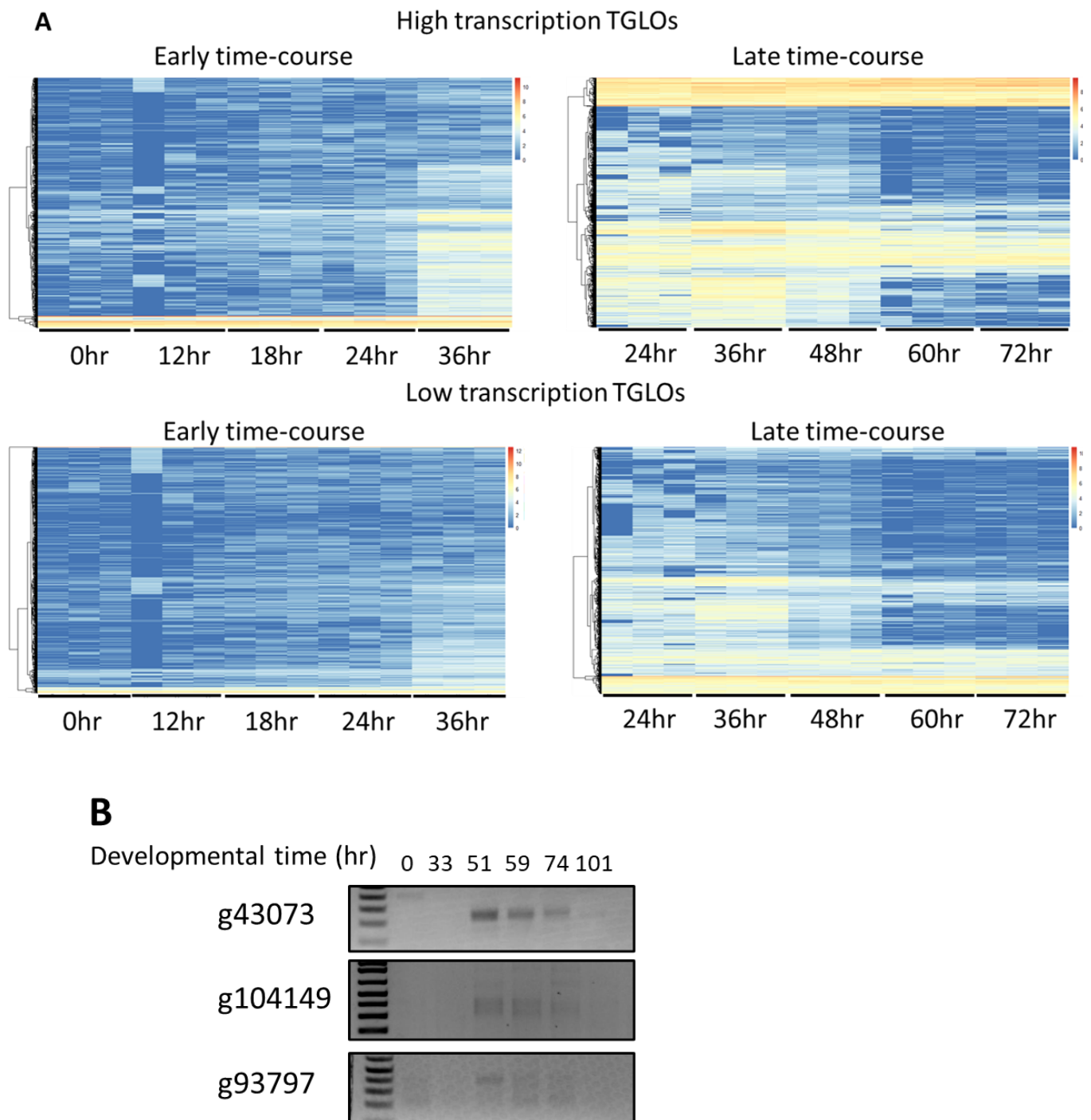
**B**

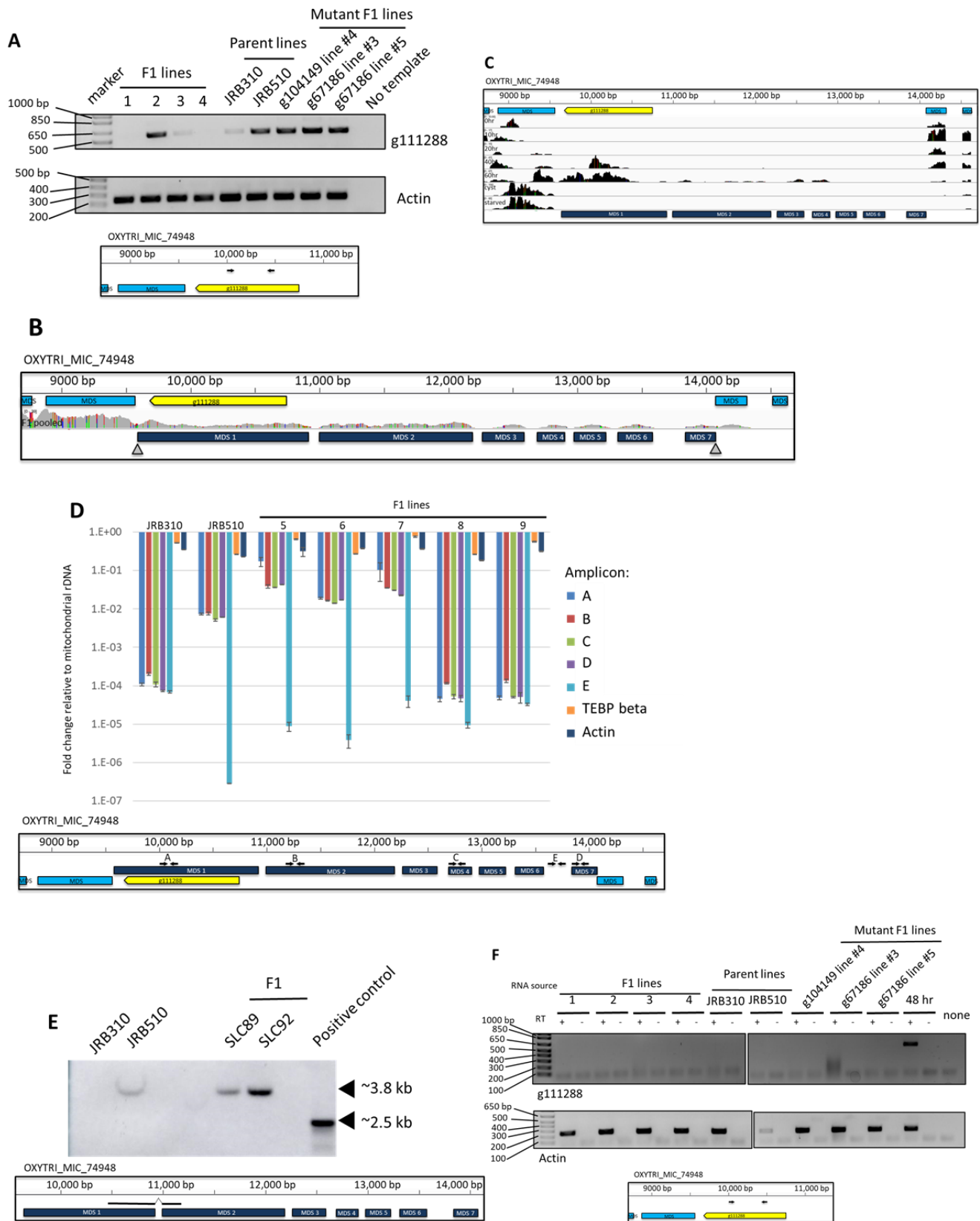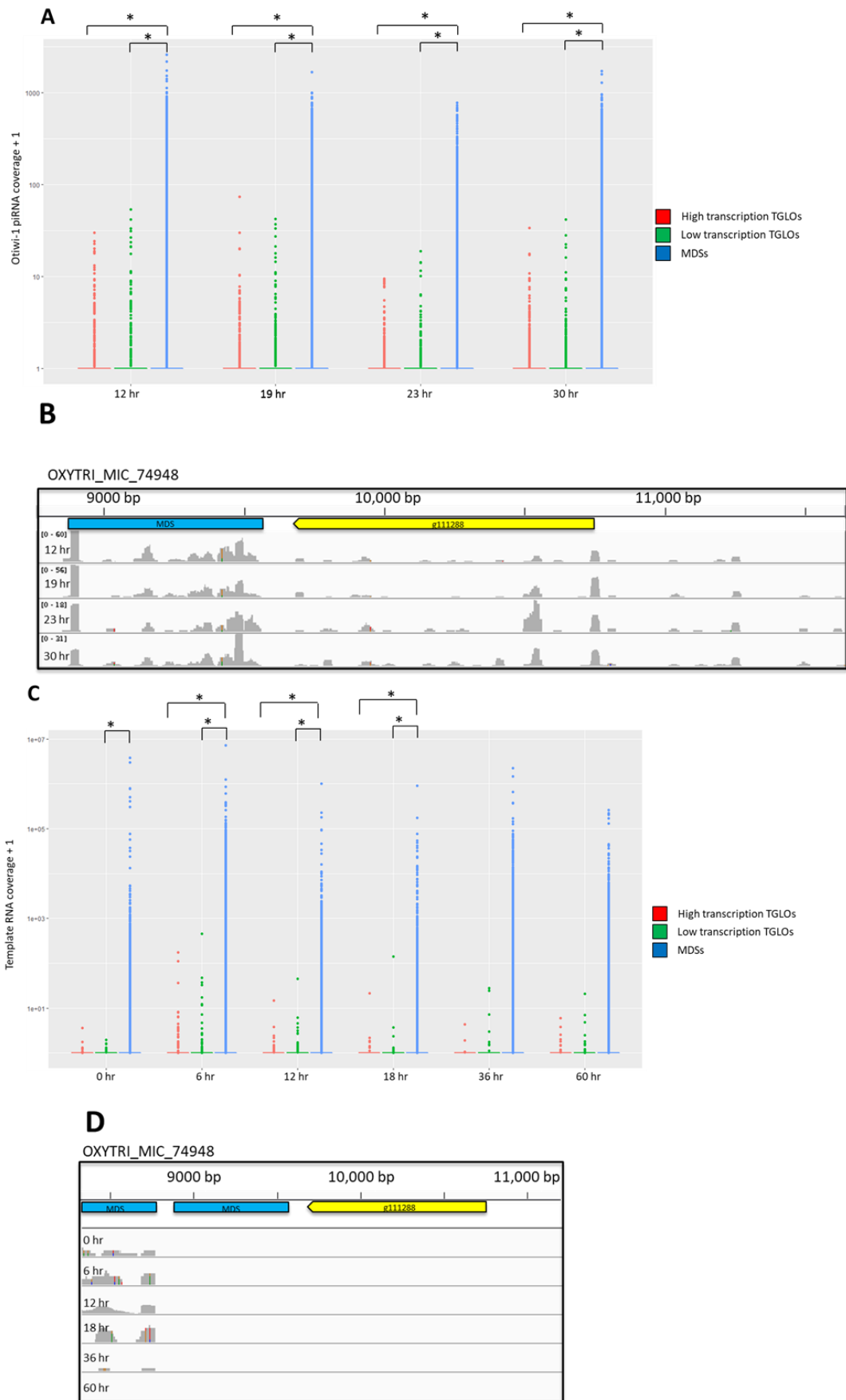Developmental time (hr)   0  33  51  59  74 101

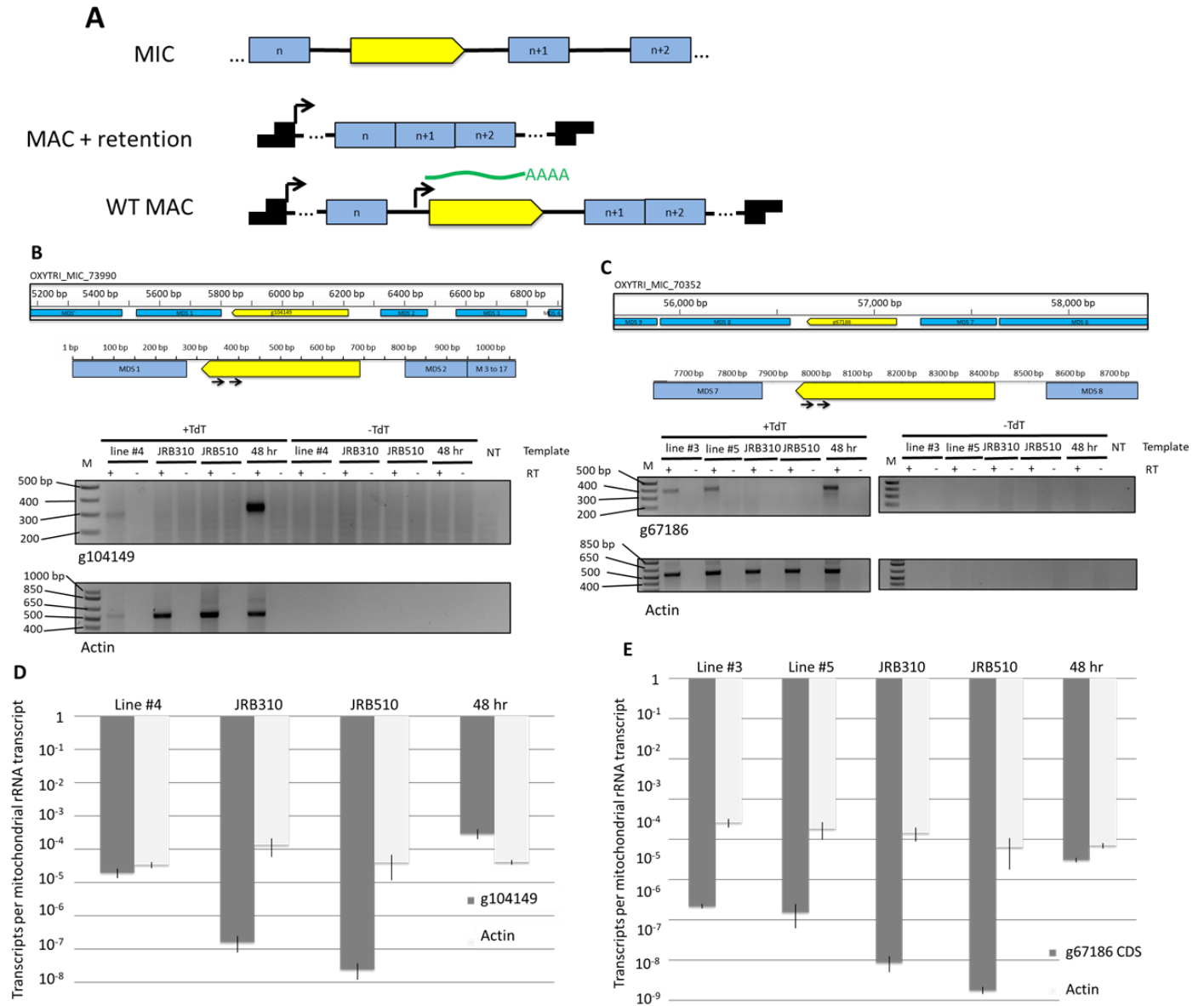g43073

g104149

g93797



Figure 2

Figure 3

Figure 4

Figure 5

Figure 6