1

**Running head: An Alternative Method to Enrich *Tetrahymena* Micronuclear DNA**

**Title: An Alternative Method to Enrich *Tetrahymena* Micronuclear DNA**

**Authors' names:** Abigail A. Howell[a,b], Brandon Neldner[c], Jonathon T. Hileman[a], Timothy J. Licknack[a,b], Wesley E. Swenson[a], Gillian H. Gile[a], Rebecca A. Zufall[d], Reed A. Cartwright[a,b]

**Authors' addresses:**
a Arizona State University, School of Life Sciences, 427 E Tyler Mall Tempe, AZ, USA 85287-4501
b Arizona State University, The Biodesign Institute, Tempe, AZ, USA
c Thermo Fisher Scientific, Tempe, AZ, USA
d University of Houston, Department of Biology and Biochemistry, Houston, TX, USA

**Correspondence**
A. Howell, School of Life Sciences, Arizona State University, 427 E Tyler Mall Tempe, AZ, USA 85287-4501, USA Telephone number: 480-292-2575; e-mail: aahowel3@asu.edu

**Abstract**

The ciliate *Tetrahymena thermophila* is a unicellular eukaryotic model organism with two distinct nuclei: a diploid germline micronucleus and a ~45x somatic macronucleus. Centrifugation methods used to isolate micronuclei are time and resource intensive. An alternative approach to study the genome of micronuclei is genomic exclusion; however it results in a significant portion of the original micronucleus (~30%) being lost during chromosome fragmentation and removal of Internally Eliminated Sequences (IES). Whole cell sequencing is also not a viable option for studying micronuclear specific evolution, as ~95% of extracted DNA is derived from the macronucleus. Through detergent-free nuclei isolation and flow cytometry, we have developed a successful enrichment method for the micronuclei of *Tetrahymena* for genomic analysis. We have validated MIC enrichment during flow sorting through (a) Fisher's exact tests of uniquely mapped reads to the micro and macronuclear reference genomes, (b) mean coverage depth of IES and Macronuclear-destined regions after alignment to the micronuclear reference genome, and (c) IES retention scores.

**Keywords:** Flow Cytometry, Tetrahymena, Mutation, Internally Eliminated Sequences (IES), Micronuclei

## 1. Introduction

*Tetrahymena thermophila* is a microbial eukaryote with an extensive genetic toolkit, contributing to the discovery of the microtubule motor, dynein (Gibbons and Rowe 1965), catalytic RNA (Kruger et al. 1982), telomere structure and telomerase (Greider and Blackburn 1985), histone acetyl transferase (Brownell et al. 1996), and programmed excision of transposon-related DNA from the somatic genome (Taverna et al. 2002), among numerous other areas of research. As with all ciliates, *Tetrahymena* exhibit nuclear dimorphism, containing a diploid germline micronucleus (MIC) and ~45x somatic macronucleus (MAC). The life cycle of *T. thermophila* is made up of two stages: asexual reproduction by binary fission and sexual conjugation between individuals of different mating types. During normal vegetative growth, transcription occurs in the MAC, and the MIC remains transcriptionally silent. When dividing asexually, the MIC divides mitotically as most eukaryotic cells would, while the MAC divides through a process called amitosis, where chromosomes are randomly segregated, potentially leading to unequal chromosome numbers during MAC division (Allen and Nanney 1958; Doerder *et al.* 1975; Orias and Flacks 1975; Nanney and Preparata 1979). Through amitosis, a heterozygous MAC can become homozygous after multiple rounds of division through the stochastic fixation of one allele, a process known as phenotypic assortment (Orias and Flacks 1975; Nanney and Preparata 1979; Merriam and Bruns 1988). When stressed, *Tetrahymena* undergo meiosis, mate, and form new zygotic nuclei from which new MIC and MAC develop. The MAC is not a direct copy of the MIC, and during the development of the MAC chromosomes are fragmented and thousands of Internally Eliminated Sequences (IESs) are removed in the MAC but retained in the transcriptionally silent MIC (Yao *et al.* 1987). This dual life cycle of sexual and asexual

83 reproduction is thought to be an adaptive response to stress, increasing genetic diversity in an
84 individual to increase the odds of survival (Fjerdingstad et al. 2007).

85 The nuclear dimorphism and dual lifecycle of *T. thermophila* provide multiple experimental
86 advantages. The separate micronuclear and macronuclear genomes allow for lethal mutations and
87 segmental deletions to be maintained in the MIC until mating, which makes physical mapping of
88 mutations, DNA polymorphisms, or MIC-limited DNA elements possible without their
89 expression in the MAC. Phenotypic assortment, another unique attribute of *T. thermophila*,
90 allows for recessive mutations to be fully expressed after a series of asexual generations and
91 facilitates knockdown/knockout experiments of otherwise essential genes in the MAC (Hai *et al.*
92 2000). When combined with a drug resistance gene, MAC knockouts can be driven to near
93 homozygosity through stepwise increases in drug concentrations, a model which is difficult to
94 obtain in other organisms without an additional Cre-lox system to selectively knock-down
95 essential genes in certain tissues. Additionally, *Tetrahymena's* fast growth rate (~2- to 3-hr
96 doubling time), sequenced MIC and MAC genomes, and established mapped genetic markers
97 make them an ideal model system for forward and reverse genetics studies for gene discovery
98 (Ruehle *et al*. 2016). Despite their ubiquity as a genetic model system, there is no simple,
99 efficient method for separating the micro and macronuclei of *Tetrahymena.*

100 For some studies, it is desirable to isolate the micronuclei of *Tetrahymena*. This includes
101 research on scan RNAs (scnRNAs) produced in the MIC that regulate DNA elimination during
102 conjugation (Schoeberl et al., 2012), the enzymatic and chemical mapping of nucleosome
103 distribution in the micronuclei (Chen et al., 2016), or the role of heterochromatic histone
104 posttranslational modifications (PTMs) during meiosis and mitosis (Papazyan et al., 2014).
105 However, the centrifugation protocol used in these studies to isolate the micronuclei requires at
106 minimum 1-2 liters of culture for a sufficient yield (Sweet and Allis 2006; Duan et al., 2021).

107 Isolating the micronuclei of *Tetrahymena* would also be beneficial for mutation accumulation
108 studies that seek to characterize complex mutations events (insertions, deletions, and copy
109 number variants) which could be responsible for the fitness decline observed in previous
110 experiments (Long et al., 2013; Long et al., 2016). Growing large volumes of culture for
111 centrifugation and micronuclei isolation is undesirable in these studies, as increasing the number
112 of generations and population size could bias the mutation results. To generate sufficient
113 micronuclear material for sequencing, these studies utilized a process known as genomic
114 exclusion (GE). In GE, a *Tetrahymena* cell is mated with a star strain of *Tetrahymena* which
115 lacks a MIC, causing the resulting daughter cells to be generated from a single meiotic product
116 of the original cell. The MAC in these daughter cells are descended from the original
117 micronucleus, making all data reflective of only the MIC (Allen 1963). However, a significant
118 portion of the original MIC, and potential mutations, are lost during somatic genome
119 rearrangement. Whole cell sequencing is also not a viable option for studying micronuclear
120 specific evolution, as approximately 95% of DNA in the cell consists of macronuclear DNA.

121 In this study we present an alternative method to enrich samples for *Tetrahymena* micronuclei
122 using propidium iodine staining and flow sorting (Figure 1). Flow cytometry has been previously

123 demonstrated to sort subpopulations of nuclei to high purity in *Paramecium*, a close relative of
124 *Tetrahymena* (Guérin et al. 2017). Further, propidium iodine has shown to be successful in
125 staining both nuclei in previous experiments with *Tetrahymena* (Po-Hsuen et al. 2015). This
126 protocol minimizes the time and resources associated with growing large amounts of culture
127 required in previous nuclei isolation protocols (Sweet and Allis 2006; Duan et al., 2021). We
128 validate the flow sorted samples using high throughput DNA sequencing on the Illumina MiSeq
129 platform. Our findings demonstrate that flow cytometry is a viable method for enriching samples
130 for micronuclei in *Tetrahymena* that utilizes less than 25mL of culture compared to the liters of
131 culture required of other methods.

132 **2. Materials and Methods**

133 **2.1 Strains and Media**

134 The *T. thermophila* strain used in this study was SB210-E from the *Tetrahymena* Stock Center
135 (Cornell University). Cultures were grown in 25mL of Neff Medium (Cassidy-Hanley 2012)
136 with penicillin and streptomycin (250 µg/mL each) and amphotericin B (0.25 µg/mL) at 30°C
137 shaken at 140rpm for 5 days or until cell concentrations reached ~$10^6$-$10^7$cells/mL.

138 **2.2 Cell Preparation**

139 To prepare cells for flow sorting, one 25mL culture was centrifuged at 10,000rpm for 1 minute to
140 concentrate. The bottom 7.5mL of the concentrated culture and pelleted cells were removed and
141 added to a fresh tube with 1.125mL Galbraith's solution (Galbraith et al., 1983). The solution
142 was vortexed briefly for 30 seconds and then passed through a 40µm filter to remove debris and
143 washed with an additional 500µl of Galbraith's solution. 1 µl of Propidium Iodide (PI)
144 (1mg/mL) was then added per mL of solution to stain the nuclei. Here, a 10µl sample was
145 observed under a fluorescent microscope at 535nm (40x) to confirm staining with PI. Cells were
146 then lysed with a tight pestle Dounce homogenizer for 15 turns. A 10µl sample of the
147 homogenate was again observed under a fluorescent microscope at 535nm (40x) to ensure the
148 micronuclei had been removed from their "cup" next to the macronuclei and that nuclei remained
149 intact.

150 **2.3 Flow Cytometry and Cell Sorting**

151 Stained samples were sorted on a BD Biosciences FACS (Fluorescence-Activated Cell Sorting)
152 Aria IIu (San Jose, CA) utilizing a 100µm nozzle and Beckman Coulter IsoFlow Sheath Fluid
153 (Brea, CA) with a sheath pressure of 20 psi. Forward Scatter (FSC) and Side Scatter (SSC) were
154 measured using standard filters off the 488 nm laser for approx. 1 hour. Propidium iodide (PI)
155 was excited and measured off the 561nm laser. Nuclei were sorted into MAC-enriched and MIC-
156 enriched samples based on FSC, SSC, and intensity of PI signals (Figure 2) into 1.5mL
157 microcentrifuge tubes containing phosphate buffered saline (PBS). MIC and MAC are primarily
158 distinguished by the PI signal. Data were collected using BD FACS Diva 8.0.1 software (San
159 Jose, CA). Data were analyzed with FlowJo v10.6 (BD Biosciences, San Jose CA).

160 **2.4 Genomic DNA extraction and sequencing**

161  After sorting, genomic DNA was extracted from both samples using phenol-chloroform,
162  following a protocol provided by Pacific Biosciences (http://www.pacb.com). Samples were
163  concentrated to 30µl at 11.2 and 10.5 ng/µl for the MIC and MAC, respectively. Paired-end
164  sequencing was performed on the Illumina MiSeq Nano V2 platform (250 cycles) at the DNASU
165  core facility at the Biodesign Institute at Arizona State University. Samples were multiplexed
166  with the final number of reads per sample being 1,048,024 reads for the MAC FACS sample and
167  904,282 reads for the MIC FACS sample. Genomic DNA for whole cell control samples was
168  extracted by phenol-chloroform, and one paired-end library (average insert size 280bp) was
169  generated on the Apollo 384 liquid handler using KAPA Biosystem's LTP library preparation kit
170  (KK8232) following the manufacturer's instructions. Sequencing was performed on the Illumina
171  NEXseq platform (150 cycles) at the DNASU core facility at the Biodesign Institute at Arizona
172  State University. The final number of reads for the whole cell control sample was 132,788,552.
173  Sequencing reads are available from the NCBI's SRA database under a BioProject with
174  accession number PRJNA735576.

175  **2.5 Bioinformatic Analyses**

176  **Overview**

177  The goal of the following bioinformatic analyses is to quantify the enrichment of micronuclei.
178  Our first analysis compared the proportion of uniquely mapped reads to the MIC vs. MAC
179  reference genomes for each FACS sample to a whole cell sample. Unique regions in the MIC
180  reference are IESs, while unique regions in the MAC reference are IES excision boundaries
181  (Figure 3). The MIC-enriched FACS sample is expected to have a greater proportion of uniquely
182  MIC-mapped reads than the whole cell data. To estimate contamination from the opposite
183  genome in the FACS samples, we compared the proportion of uniquely mapped reads to the MIC
184  vs. MAC reference genomes in simulated MIC and MAC reads relative to simulated whole cell
185  reads. The simulated whole cell data also addresses any potential biases in the original whole cell
186  sample, as we would expect the proportion of uniquely mapped reads to the MIC vs. MAC
187  reference genomes for both the simulated and actual whole cell reads to be similar.

188  In our second analysis we compared unique and shared regions between the MIC and MAC
189  through coverage levels of IESs, which are unique to the MIC, and Macronuclear-Destined
190  Sequences (MDSs), which are found in both MIC and MAC, per FACS sample.

191  For our final analysis we modified an established method of validating flow sorting, IES
192  retention scores (IRSs), for use in *Tetrahymena*. The retention score of an IES is given by the
193  equation: IRS=IES+/(IES+ + IES−) (Swart et al., 2014), where IES+ represents the number of
194  reads that contain the IES sequence and IES- represents the number of reads that contain the
195  MAC excision boundary of the corresponding IES. Originally developed for *Paramecium*, the
196  IRS- score can be easily calculated by counting the number of reads that contain a conserved TA
197  dinucleotide found at the ends of Parameicum IESs (Arnaiz et al., 2012; Gratias and Bétermier,
198  2003). However, as *Tetrahymena* exhibits sequence diversity at IES excision sites, this method
199  can not be exactly replicated. To modify the IES retention scores for *Tetrahymena*, we calculated
200  the IES- as the number of reads that contain the MAC excision boundary of the corresponding

201 IES, which we determined as the regions in the MAC reference genome immediately adjacent to
202 a known IES region in a pairwise alignment MIC/MAC chain file generated by the software
203 transanno (github.com/informationsea/transanno) and minimap2 (Li 2018). Only viable IESs
204 were used for this analysis, which included those that are within 10bps of the adjacent MAC
205 scaffolds in a MIC/MAC chain file. IESs that overlap with MAC scaffolds were discarded.

## Fisher's exact tests of uniquely mapped reads

207 To validate enrichment of micronuclei by flow sorting, reads from the MIC-enriched FACS and
208 MAC-enriched FACS samples were trimmed with trimmomatic v0.38 (Bolger et al. 2014) and
209 aligned to a combined MIC
210 (http://datacommons.cyverse.org/browse/iplant/home/rcoyne/public/tetrahymena/MIC), MAC
211 (*Tetrahymena* Genome Database http://ciliate.org/), and mitochondrial (NCBI Reference
212 Sequence: NC_003029.1) reference using BWA mem v0.7.12 (Li and Durbin 2010). The
213 sequenced MIC genome (Hamilton et al. 2016) consists of 5 chromosomes totaling 157Mb, and
214 was generated using Illumina whole genome shotgun sequencing with PCR-free fragment
215 libraries. The micronuclei in Hamilton et al. (2016) were isolated using differential
216 sedimentation as described in Gorovsky et al. (1975). The MAC genome sequence (Sheng et al.
217 2020) is made up of 181 chromosomes capped with two telomeres with all gaps entirely closed
218 (103.3Mb), and was sequenced using $300\times$ long Single Molecule, Real-Time reads. The
219 macronuclei in Sheng et al. (2020) were isolated using a modified differential centrifugation
220 protocol described in Chen et al. (2016). The 47,577 bp mitochondrial genome of *T. thermophila*
221 has also been sequenced (NCBI Reference Sequence: NC_003029.1) (Brunk et al. 2003). All
222 unmapped reads were removed from the BAM files of each dataset using SAMtools v1.10 (Li et
223 al. 2009). Additionally, all secondary and chimeric alignments, which could indicate mapping to
224 a region shared between the two genomes, were also removed. The number of reads from each
225 BAM file that aligned uniquely to a region in the MIC or MAC reference genome were then
226 calculated using SAMtools. As a control, whole cell reads, equivalent in read number to the
227 MIC-enriched and MAC-enriched FACS samples, were randomly sampled from a previous
228 sequencing run (BioProject PRJNA735576) using seqtk (Shen et al., 2016). The proportion of
229 uniquely mapped reads from the whole cell data to each reference genome was also calculated
230 using SAMtools. Two Fisher's exact tests were performed using the whole cell reads compared
231 with the MIC-enriched FACS sample reads and the whole cells reads compared with the MAC-
232 enriched FACS sample reads.

## Simulated MIC-enriched FACS, and MAC-enriched FACS, and whole cell reads

234 To estimate contamination from the opposite genome in the FACS samples, we compared the
235 proportion of uniquely mapped reads to the MIC vs. MAC reference genomes in simulated MIC
236 and MAC reads relative to simulated whole cell reads. Simulated MIC-enriched FACS and
237 MAC-enriched FACS reads were sampled from their respective reference genomes using ART
238 (v1.5.0) (Huang et al., 2012) and aligned to the combined MIC, MAC, and mitochondrial
239 reference using BWA mem v0.7.12 (Li and Durbin 2010). Simulated 250-bp paired-end whole
240 cells reads were created also using ART (v1.5.0). DNA fragment size and standard deviation

241   were estimated from the whole cell sample reads using deepTools bamPEFragmentSize
242   (Ramírez et al., 2014). While the usual whole cell ploidy of *T. thermophila* is 2:45 MIC:MAC,
243   considering we used an asynchronized population of cells in our experiments the MAC is
244   estimated to be at an average of 64x (between 45x in the G1 phase and 90x in the G2 phase)
245   (Woodard et al., 1972), while the MIC is essentially 4x with no apparent G1 phase (Cole and
246   Sugai, 2012). Therefore, we simulated the whole cell ploidy by sampling the MIC reference at 4x
247   coverage and the MAC reference at 64x coverage. Simulated reads were aligned to the combined
248   MIC, MAC, and mitochondrial reference using BWA mem v0.7.12 (Li and Durbin 2010). The
249   simulated whole cell data also addresses any potential biases in the original whole cell sample, as
250   we would expect the proportion of uniquely mapped reads to the MIC vs. MAC reference
251   genomes for both the simulated and actual whole cell reads to be similar.

252   **Mean read depth of IESs and MDSs**

253   To determine coverage levels of IES, which are unique to the MIC, and Macronuclear-Destined
254   Sequences (MDSs), which are found in both MIC and MAC, per FACS sample, we used the
255   locations of IESs in the MIC supercontigs from Hamilton et al. (2016), Supplementary file 1C
256   and Supplementary file 3A. The depth of coverage for the IESs and MDSs was calculated using
257   Samtools depth (Li et al. 2009) for each sample.

258   **IES retention scores**

259   We further verified enrichment of micronuclei by calculating a retention score for each viable
260   IES. Viable IESs include those that are within 10bps of the adjacent MAC scaffolds in a
261   MIC/MAC chain file created using the software transanno
262   (github.com/informationsea/transanno) and minimap2 (Li 2018). IESs that overlap with MAC
263   scaffolds were discarded. MIC-enriched and MAC-enriched FACS data were aligned to a
264   reference sequence consisting of the MAC reference and the sequence of each individual IES
265   using BWA mem v0.7.12 (Li and Durbin 2010). IES retention scores (IRSs) were determined for
266   each IES by counting the number of reads that contain the IES sequence (IES+) and the number
267   of reads that contain the MAC excision boundary of the corresponding IES (IES-). The retention
268   score of an IES is given by the equation: IRS=IES+/(IES+ + IES−) (Swart et al., 2014).

269   **3. Results and Discussion**

270   **3.1 Flow-cytometric assay**

271   The MIC and MAC of *T. thermophila* are distinct in both the physical size of their nuclei (~3 µm
272   and ~10-15µm in diameter, respectively) (Figure 4) as well as their genome size (a diploid 157
273   Mb genome and a 45x ploid 103 Mb genome, respectively). Therefore, we FACS-sorted the
274   nuclei based on forward scatter (FSC) and side scatter (SSC) values as well as the intensity of the
275   PI signal. The MIC fraction of the sample represented 4.87% of the total events of the flow
276   sorting and 76.55% of the nuclear sample (circled points in Fig. 2). The MAC fraction of the
277   sample represented 1.46% of the total events of the flow sorting and 23.44% of the nuclear

278   sample (Figure 2). We then validated the FACS enrichment of MIC through comparisons of
279   uniquely mapped reads, mean coverage depth of IESs and MDSs, and IES retention scores.

**3.2 Percentage of uniquely mapped reads**

281   The proportion of uniquely mapped reads (i.e., reads that aligned solely to either the MIC or
282   MAC reference in the combined MIC, MAC, and mitochondrial reference) for the whole cell
283   sample was 0.36 MIC and 0.64 MAC (Table 1). The 36/64 proportion of uniquely mapped reads
284   to the MIC and MAC references represents the baseline proportion that the MIC and MAC
285   FACS data were compared to in order to validate enrichment.

286   Using the whole cell proportion of uniquely mapped reads to the MIC and MAC references as a
287   baseline, we observe that there is clear enrichment for MIC sequences in the MIC-enriched
288   FACS sample based on the sample's 83/17 proportion of uniquely mapped reads (Fisher's exact
289   test, $p < 0.001$). There is also enrichment for MAC sequences in the MAC-enriched FACS
290   sample based on the sample's proportion of 14/86 uniquely mapped reads to the MIC and MAC
291   references (Fisher's exact test, $p < 0.001$).

**3.3 Simulated whole cell, MIC-enriched FACS, and MAC-enriched FACS reads**

293   As there are far more unique regions in the MIC in the form of IESs (approx. 54 Mb) compared
294   to the unique regions of the MAC (IES excision junctions), we found it surprising that there was
295   a higher proportion of uniquely mapped reads to the MAC reference from the whole cell sample
296   in our baseline proportion (36/64 MIC/MAC). To understand the 36/64 baseline proportion of
297   uniquely mapped reads to the MIC and MAC references from the whole cell samples, we
298   simulated whole cell reads using ART (v1.5.0) (Huang et al., 2012) to reflect the 4:64
299   MIC:MAC ploidy of a *T. thermophila* cell. After alignment to the combined MIC, MAC,
300   mitochondrial, and rDNA chromosome reference using BWA mem v0.7.12 (Li and Durbin
301   2010), the proportion of uniquely mapped reads to the MIC and MAC references in the simulated
302   whole cell reads was 34/66 (Table 1) compared to the 36/64 (Table 1) of the original whole cell
303   samples, confirming the original proportion as an accurate baseline measurement (Fisher's exact
304   test, p=1). To investigate further why the number of uniquely mapped reads to the MAC
305   reference was higher than the uniquely mapped reads to the MIC reference in the whole cell
306   samples, despite all of the MAC's genomic content stemming from the MIC with the exception
307   of IES excision sites, we again simulated whole cell reads but at a 1:1 MIC:MAC ploidy. Table 1
308   illustrates that the majority of uniquely mapped reads after alignment to the combined reference
309   genome do originate from the MIC (90%), likely due to the thousands of IESs present in the MIC
310   which are absent in the MAC. The high number of uniquely mapped reads to the MAC reference
311   in normal whole cell samples is attributed to the high ploidy of the MAC. The number of reads
312   mapping uniquely to the MIC and MAC references from the 1:1 MIC:MAC ploidy simulated
313   whole cell reads (146028:17558) when multiplied by the 4:64 MIC:MAC ploidy, equal a
314   proportion of 34/66, identical to the 34/66 proportion of the 4:64 MIC:MAC ploidy simulated
315   whole cell reads.

316   Finally, we also simulated reads generated only from the MIC reference and reads generated
317   only from the MAC reference to estimate contamination from the opposite genome in the FACS
318   samples. The difference between the proportion of uniquely mapped reads to the MIC and MAC
319   references in the simulated MIC-enriched FACS samples (99.93/0.03, Table 1) and the
320   proportion of the actual MIC-enriched FACS sample (83/17, Table 1) indicates there is
321   contamination from the MAC. Possibly, this is due to the degradation of the MAC observed
322   during sorting or damage to the MAC during homogenization. Modifying the atmospheric
323   conditions in the flow sorter to adjust the pH can also be explored as a means to limit cross
324   contamination of nuclei due to the degrading MAC, as pH has been previously implicated in
325   reducing cell viability (Cossarizza et al. 2017). We observed a similar level of MIC
326   contamination in the MAC-enriched FACS sample based on the proportion of uniquely mapped
327   reads to the MIC and MAC references in the simulated MAC-enriched FACS samples
328   (0.05/99.5, Table 1) compared to the proportion of the actual MAC-enriched FACS sample
329   (14/86, Table 1). Compared to previously published studies utilizing differential centrifugation
330   (contamination of MAC ranging from 1-3%, Xiong et. al 2015; Xiong et. al 2016), our MIC-
331   MAC fraction cross contamination is significantly higher. This could be due to the length of time
332   the cultures were grown, as the MIC is more tightly attached to the MAC of stationary-phased
333   cells which can reduce the purity of the fractions (Gorovsky et al., 1975). For MIC purification it
334   is ideal that cell density does not exceed $2.5\times10^5$ cells mL$^{-1}$ (mid-log phase) (Chen et al. 2016),
335   therefore we might expect better results from flow cytometry had cultures been prevented from
336   growing to such high density. However, for many purposes, e.g. mutation detection, absolute
337   purity of the FACS samples is not required, as expectations of read origin and the likelihood of
338   *de novo* mutations can be adjusted based on the estimated level of contamination and further
339   supported by sequencing of both fractions. Beyond measuring the proportion of uniquely
340   mapped reads, MIC purity could also be measured by a second round of flow cytometry,
341   comparing the percentage of MIC in the total nuclear sample after the first and second round of
342   sorting as done in Guérin et al. (2017).

**3.4 Mean coverage depth of IESs and MDSs**

344   For the second validation test we calculated sequencing coverage depth in MDSs vs. IESs for the
345   whole cell sample, MIC-enriched FACS sample, and MAC-enriched FACS sample (Table 2).
346   This analysis compares unique and shared regions between the MIC and MAC as opposed to
347   only unique regions in the first validation test. For the whole cell sample we would expect an
348   IES:MDS ratio of 4:68 or 0.06:1, as IESs occur only in the diploid MIC (4x in G1 phase) while
349   Mac-destined regions occur in the MAC at 64x coverage (between 45x in the G1 phase and 90x
350   in the G2 phase) and in the MIC (4 + 64 = 68). Our actual ratio in the whole cell sample may
351   vary based on the number of IESs that are accidentally maintained in the MAC during genome
352   rearrangement or missing from the data due to insufficient sequencing depth. Using the whole
353   cell ratios as a baseline, we observed that in the MIC FACS data there was an IES:MDS read
354   depth ratio of 0.45:1 which shows enrichment in the MIC. While we would expect this ratio to be
355   1:1 as all regions of the MIC should be diploid, there could be MAC contamination or IESs
356   could be located in poorly assembled regions in the reference. For the MAC-enriched FACS
357   sample there was minimal IES coverage with a IES:MDS coverage ratio of 0.084:1, as all IESs

358    should be eliminated from the MAC. MIC contamination and IES retention are possible sources
359    of sequencing coverage in IES regions.

360    In our FACS samples there was also evidence of bacterial contamination. 71% of the total
361    number of reads in the MIC-enriched FACS sample and 34% of the total number of reads in the
362    MAC-enriched FACS sample did not map to the combined MIC, MAC, and mitochondrial
363    reference, compared with only 1% of reads in the whole cell data when mapped to the combined
364    MIC, MAC, and mitochondrial reference, which was from a separate experiment with deeper
365    sequencing. The high percentage of unmapped reads in the MIC-enriched and MAC-enriched
366    FACS samples prompted us to BLAST the unmapped reads from each sample to explore
367    potential sources of contamination. From the blast results we found that the major source of
368    contamination was bacterial with a mix of *Azoarcus*, *Acidovorax*, *Alicycliphilus*,
369    *Diaphorobacter*, *Alicycliphilus*, and *Pseudomonas*. Contamination could have occurred in
370    cultures, during flow sorting, or during DNA extraction. While contamination should not affect
371    the IES/MDS measures presented here, high contamination levels would be a considerable
372    drawback for costlier long-read sequencing. Bacterial contamination can be limited by filtering
373    with a <5µm filter before flow sorting as well incorporating the antibiotics neomycin,
374    kanamycin, or tetracycline (100 µg/mL each) into the cell culture as suggested by Cassidy-
375    Hanley (2012).
376
377    Also important for utilizing flow sorting for long read sequencing is the quality of DNA. While
378    the quality of the DNA collected from the MIC and MAC FACS samples (shearing, degradation)
379    was not collected for this study it can be obtained in future experiments using a Bioanalyzer.
380    Suggested methods of maintaining DNA quality during flow sorting include increasing the
381    surface area of single cell suspensions in order to maximize contact between cells and digestive
382    enzymes (Reichard and Asosingh 2018) and sanitizing all equipment to prevent degradation by
383    contaminating nucleases (Ormerod and Imrie 1990).

384    **3.5 IES retention scores**

385    The final metric used to validate MIC-enrichment from flow sorting was the IES retention scores
386    (IRS) of viable IESs. Viable IESs included those that are within 10bps of the adjacent MAC
387    scaffolds in a MIC/MAC chain file. IESs that overlap with MAC scaffolds were discarded. The
388    fraction of viable IESs to the potential total number of IESs (unverified) from Hamilton et al.
389    (2016) are included in Table 3. The location of the IESs in the fully assembled MIC
390    chromosomes were extrapolated from the position of the IESs in the MIC scaffolds
391    (Supplementary file 3A) and the position of the MIC scaffolds within the MIC chromosomes
392    (Supplementary file 1C). As some MIC scaffolds were incorporated multiple times into the MIC
393    chromosome assemblies there are repeats of identified IESs, bringing the original total of 7551
394    IESs from Hamilton et al. (2016) to 8171. For the MIC-enriched FACS sample if sorting were
395    perfect (and IES excision were perfect) the expected IRS is 1, as there should be no short reads
396    that span both the left and right excision boundary while there will be reads that map to the IES
397    itself. The MAC-enriched FACS sample IRS should be 0, as there will be reads that span both
398    the left and right excision boundary after the IES is removed and no reads that map uniquely to
399    any IES. Table 3 shows the average IRS for the MIC-enriched FACS data per chromosome as

400 closer to 1 than the MAC IRSs. The IRSs of the MIC-enriched FACS sample are also
401 demonstrated to skew towards 1 in Figure 5 (blue) while the MAC-enriched FACS sample IRSs
402 skew towards 0 (red). This indicates that there is enrichment for MIC DNA in the MIC-enriched
403 FACS sample and limited MIC contamination in the MAC-enriched FACS sample. The not
404 insignificant fraction of 0 scores for the MIC-enriched FACS sample (1068) could be from IESs
405 identified in the Hamilton et al. (2016) study that do not actually exist or were not present in our
406 sample. While the IESs used in this study are classified as "high-confidence" after verification
407 from at least two different identification methods (MAC read alignment to MIC reference, MIC
408 read alignment to MAC reference, and MIC-MAC cross-assembly alignment), there is still a
409 degree of uncertainty regarding the location of the identified IESs. This is due to the repetitive
410 sequences within the IESs and minor contamination of the MIC sequencing libraries with MAC
411 DNA during the assembly in Hamilton et al. (2016), which could create a mixture of inconsistent
412 short reads at IES/MDS junctions creating false breaks in the assembly that are not the result of
413 IESs. Additionally, research suggests (Feng et al. 2017; Jaspan et al. 2019) that the excision
414 boundaries of IESs exhibit variability during conjugation, which could account for the IRS
415 scores in the MIC-enriched FACS sample that are closer to 0 than 1 if the predicted IES/MDS
416 junction is incorrect due to excision variability. This validation technique could be improved
417 through the use of de novo IES detection in our own samples as opposed to using previously
418 published IES locations.

419 **Conclusions and Future Applications**

420 The use of flow sorting that we describe here to enrich for micronuclei in *T. thermophila* is
421 simple and requires only 25mL of culture compared to the liters required for centrifugation
422 methods. Our validation approaches, (a) Fisher's exact tests of uniquely mapped reads to the
423 micro and macronuclear reference genomes (b) mean coverage depth of IES and Macronuclear-
424 destined regions after alignment to the micronuclear reference genome and (c) IES retention
425 scores, all support MIC-enrichment from flow sorting (approximately 9x more MIC DNA
426 compared to whole cell sequencing based on validation approach (a)), but also suggest cross
427 contamination of nuclei between the MIC and MAC FACS samples. To improve MIC
428 purification in future iterations of flow sorting, cell cultures can be grown to no greater than mid-
429 log phase to prevent MICs from tightly attaching to the MAC.
430
431 While the quality of the DNA collected from the MIC and MAC FACS samples (shearing,
432 degradation) was not collected for this study it can be obtained in future experiments using a
433 Bioanalyzer. Suggested methods of maintaining DNA quality during flow sorting include
434 increasing the surface area of single cell suspensions in order to maximize contact between cells
435 and digestive enzymes (Reichard and Asosingh 2018) and sanitizing all equipment to prevent
436 degradation by contaminating nucleases (Ormerod and Imrie 1990).
437
438 This method will allow for rapid sequencing of MIC-enriched DNA for comparison with existing
439 MAC sequences to help elucidate the evolution and molecular mechanisms of genome
440 rearrangement in ciliates. Additionally, flow sorting allows for more detailed mutation detection
441 after sequencing than genomic exclusion. Using sequencing data from flow-sorted micronuclei,
442 we will be able to characterize complex mutations events (insertions, deletions, and copy number

443 variants) in *Tetrahymena* after mutation accumulation, which have previously been difficult to
444 detect bioinformatically after whole cell sequencing. Further, this method can be used to
445 sequence the MIC in species that are unable to be grown to sufficient volumes for centrifugation,
446 for example species of Karyorelicteans, which are notoriously difficult to culture.
447

448 **Data Availability**
449 All scripts used in the bioinformatics analysis presented here are available at
450 https://github.com/aahowel3/An-Alternative-Method-to-Enrich-Tetrahymena-Micronuclear-
451 DNA. Sequencing reads are available from the NCBI's SRA database under a BioProject with
452 accession number PRJNA735576.

453 **Acknowledgements**

458 **Literature Cited**

459 Allen, S. L. (1963). Genomic Exclusion in Tetrahymena: Genetic Basis*. *The Journal of
460     Protozoology*, *10*(4), 413–420. https://doi.org/10.1111/j.1550-7408.1963.tb01699.x

461 Allen, S. L., & Nanney, D. L. (1958). An Analysis of Nuclear Differentiation in the Selfers of
462     Tetrahymena. *The American Naturalist*, *92*(864), 139–160. https://doi.org/10.1086/282022

463 Bolger, A. M., Lohse, M., & Usadel, B. (2014). *Genome analysis Trimmomatic: a flexible
464     trimmer for Illumina sequence data*. *30*(15), 2114–2120.
465     https://doi.org/10.1093/bioinformatics/btu170

466 Brownell, J. E., Zhou, J., Ranalli, T., Kobayashi, R., Edmondson, D. G., Roth, S. Y., & Allis, C.
467     D. (1996). Tetrahymena histone acetyltransferase A: A homolog to yeast Gcn5p linking
468     histone acetylation to gene activation. *Cell*, *84*(6), 843–851. https://doi.org/10.1016/S0092-
469     8674(00)81063-6

470 Cassidy-Hanley, D. M. (2012). Tetrahymena in the Laboratory: Strain Resources, Methods for
471     Culture, Maintenance, and Storage. *Methods in Cell Biology*, *109*, 237.
472     https://doi.org/10.1016/B978-0-12-385967-9.00008-6

473 Chen, X., Gao, S., Liu, Y., Wang, Y., Wang, Y., & Song, W. (2016). Enzymatic and chemical
474     mapping of nucleosome distribution in purified micro- and macronuclei of the ciliated
475     model organism, Tetrahymena thermophila. *Science China Life Sciences*, *59*(9), 909–919.
476     https://doi.org/10.1007/S11427-016-5102-X

Doerder, F. P., Lief, J. H., & Doerder, L. E. (1975). A corrected table for macronuclear assortment in Tetrahymena pyriformis, syngen 1. *Genetics*, *80*(2), 263–265. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1213326/

Feng, L., Wang, G., Hamilton, E. P., Xiong, J., Yan, G., Chen, K., Chen, X., Dui, W., Plemens, A., Khadr, L., Dhanekula, A., Juma, M., Dang, H. Q., Kapler, G. M., Orias, E., Miao, W., & Liu, Y. (2017). A germline-limited piggyBac transposase gene is required for precise excision in Tetrahymena genome rearrangement. *Nucleic Acids Research*, *45*(16), 9481–9502. https://doi.org/10.1093/nar/gkx652

Fjerdingstad, E. J., Schtickzelle, N., Manhes, P., Gutierrez, A., & Clobert, J. (2007). Evolution of dispersal and life history strategies - Tetrahymena ciliates. *BMC Evolutionary Biology*, *7*(1), 133. https://doi.org/10.1186/1471-2148-7-133

Galbraith, D. W., Harkins, K. R., Maddox, J. M., Ayres, N. M., Sharma, D. P., & Firoozabady, E. (1983). Rapid flow cytometric analysis of the cell cycle in intact plant tissues. *Science*, *220*(4601), 1049–1051. https://doi.org/10.1126/science.220.4601.1049

Gibbons, I. R., & Rowe, A. J. (1965). Dynein: A protein with adenosine triphosphatase activity from cilia. *Science*, *149*(3682), 424–426. https://doi.org/10.1126/science.149.3682.424

Gorovsky, M. A., Yao, M. C., Keevert, J. B., & Pleger, G. L. (1975). Isolation of micro- and macronuclei of Tetrahymena pyriformis. *Methods in Cell Biology*, *9*(0), 311–327. https://doi.org/10.1016/S0091-679X(08)60080-1

Greider, C. W., & Blackburn, E. H. (1985). Identification of a Specific Telomere Terminal Transferase Activity in Tetrahymena Extracts. In *Cell* (Vol. 43).

Guérin, F., Arnaiz, O., Boggetto, N., Denby Wilkes, C., Meyer, E., Sperling, L., & Duharcourt, S. (2017). Flow cytometry sorting of nuclei enables the first global characterization of Paramecium germline DNA and transposable elements. *BMC Genomics*, *18*(1), 327. https://doi.org/10.1186/s12864-017-3713-7

Hai, B., Gaertig, J., & Gorovsky, M. A. (2000). Knockout heterokaryons enable facile mutagenic analysis of essential genes in Tetrahymena. *Methods in Cell Biology*, *62*(62), 513–531. https://doi.org/10.1016/s0091-679x(08)61554-x

Huang, W., Li, L., Myers, J. R., & Marth, G. T. (2012). ART: A next-generation sequencing read simulator. *Bioinformatics*, *28*(4), 593–594. https://doi.org/10.1093/bioinformatics/btr708

Jaspan, V. N., Taye, M. E., Carle, C. M., Chung, J. J., & Chalker, D. L. (2019). Boundaries of eliminated heterochromatin of Tetrahymena are positioned by the DNA-binding protein Ltl1. *Nucleic Acids Research*, *47*(14), 7348–7362. https://doi.org/10.1093/nar/gkz504

510 Kruger, K., Grabowski, P. J., Zaug, A. J., Sands, J., Gottschling, D. E., & Cech, T. R. (1982).
511     Self-splicing RNA: Autoexcision and autocyclization of the ribosomal RNA intervening
512     sequence of tetrahymena. *Cell*, *31*(1), 147–157. https://doi.org/10.1016/0092-
513     8674(82)90414-7

514 Lee, P.-H., Meng, X., & Kapler, G. M. (2015). Developmental Regulation of the Tetrahymena
515     thermophila Origin Recognition Complex. *PLoS Genetics*, *11*(1), e1004875.
516     https://doi.org/10.1371/journal.pgen.1004875

517 Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, *34*(18),
518     3094–3100. https://doi.org/10.1093/BIOINFORMATICS/BTY191

519 Li, H., & Durbin, R. (2010). *Fast and accurate long-read alignment with Burrows-Wheeler*
520     *transform*. *26*(5), 589–595. https://doi.org/10.1093/bioinformatics/btp698

521 Long, H. A., Paixão, T., Azevedo, R. B. R., & Zufall, R. A. (2013). Accumulation of
522     spontaneous mutations in the ciliate Tetrahymena thermophila. *Genetics*, *195*(2), 527–540.
523     https://doi.org/10.1534/genetics.113.153536

524 Long, H., Winter, D. J., Chang, A. Y. C., Sung, W., Wu, S. H., Balboa, M., Azevedo, R. B. R.,
525     Cartwright, R. A., Lynch, M., & Zufall, R. A. (2016). Low base-substitution mutation rate
526     in the germline genome of the ciliate tetrahymena thermophila. *Genome Biology and*
527     *Evolution*, *8*(12), 3629–3639. https://doi.org/10.1093/gbe/evw223

528 Long, H., Winter, D. J., Y-C Chang, A., Sung, W., Wu, S. H., Balboa, M., R Azevedo, R. B.,
529     Cartwright, R. A., Lynch, M., & Zufall, R. A. (n.d.). *Low Base-Substitution Mutation Rate*
530     *in the Germline Genome of the Ciliate Tetrahymena thermophila*.
531     https://doi.org/10.1093/gbe/evw223

532 Merriam, E. V., & Bruns, P. J. (1988). Phenotypic assortment in Tetrahymena thermophila:
533     assortment kinetics of antibiotic-resistance markers, tsA, death, and the highly amplified
534     rDNA locus. *Genetics*, *120*(2), 389–395. /pmc/articles/PMC1203518/?report=abstract

535 Nanney, D. L., & Preparata, R. M. (1979). Genetic Evidence Concerning the Structure of the
536     Tetrahymena thermophila Macronucleus*†. *The Journal of Protozoology*, *26*(1), 2–9.
537     https://doi.org/10.1111/j.1550-7408.1979.tb02722.x

538 Orias, E., & Flacks, M. (1975). Macronuclear genetics of Tetrahymena. I. Random distribution
539     of macronuclear gene copies in T. pyriformis, syngen 1. *Genetics*, *79*(2), 187–206.
540     /pmc/articles/PMC1213266/?report=abstract

541 Ormerod, M. G., & Imrie, P. R. (1990). Flow cytometry. *Methods in Molecular Biology*, *5*, 543–
542     558.

543 Papazyan, R., Voronina, E., Chapman, J. R., Luperchio, T. R., Gilbert, T. M., Meier, E.,
544     Mackintosh, S. G., Shabanowitz, J., Tackett, A. J., Reddy, K. L., Coyne, R. S., Hunt, D. F.,
545     Liu, Y., & Taverna, S. D. (2014). Methylation of histone H3K23 blocks DNA damage in
546     pericentric heterochromatin during meiosis. *ELife*, *2014*(3).
547     https://doi.org/10.7554/eLife.02996

548 Ramírez, F., Dündar, F., Diehl, S., Grüning, B. A., & Manke, T. (2014). DeepTools: A flexible
549     platform for exploring deep-sequencing data. *Nucleic Acids Research*, *42*(W1), W187.
550     https://doi.org/10.1093/nar/gku365

551 Reichard, A., & Asosingh, K. (2019). Best Practices for Preparing a Single Cell Suspension from
552     Solid Tissues for Flow Cytometry. *Cytometry Part A*, *95*(2), 219–226.
553     https://doi.org/10.1002/CYTO.A.23690

554 Ruehle, M. D., Orias, E., & Pearson, C. G. (2016). Tetrahymena as a unicellular model
555     eukaryote: Genetic and genomic tools. In *Genetics* (Vol. 203, Issue 2, pp. 649–665).
556     Genetics. https://doi.org/10.1534/genetics.114.169748

557 Schoeberl, U. E., Kurth, H. M., Noto, T., & Mochizuki, K. (2012). Biased transcription and
558     selective degradation of small RNAs shape the pattern of DNA elimination in Tetrahymena.
559     *Genes and Development*, *26*(15), 1729–1742. https://doi.org/10.1101/gad.196493.112

560 Shen, W., Le, S., Li, Y., & Hu, F. (2016). SeqKit: A Cross-Platform and Ultrafast Toolkit for
561     FASTA/Q File Manipulation. *PLOS ONE*, *11*(10), e0163962.
562     https://doi.org/10.1371/journal.pone.0163962

563 Sheng, Y., Duan, L., Cheng, T., Qiao, Y., Stover, N. A., & Gao, S. (2020). The completed
564     macronuclear genome of a model ciliate Tetrahymena thermophila and its application in
565     genome scrambling and copy number analyses. *Science China. Life Sciences*, *63*(10), 1534.
566     https://doi.org/10.1007/S11427-020-1689-4

567 Swart, E. C., Wilkes, C. D., Sandoval, P. Y., Arambasic, M., Sperling, L., & Nowacki, M.
568     (2014). Genome-wide analysis of genetic and epigenetic control of programmed DNA
569     deletion. *Nucleic Acids Research*, *42*(14), 8970–8983. https://doi.org/10.1093/nar/gku619

570 Sweet, M. T., & Allis, C. D. (2006). Isolation and Purification of Tetrahymena Nuclei. *Cold
571     Spring Harbor Protocols*, *2006*(23), pdb.prot4500-pdb.prot4500.
572     https://doi.org/10.1101/pdb.prot4500

573 Taverna, S. D., Coyne, R. S., & Allis, C. D. (2002). Methylation of histone H3 at lysine 9 targets
574     programmed DNA elimination in Tetrahymena. *Cell*, *110*(6), 701–711.
575     https://doi.org/10.1016/S0092-8674(02)00941-8

576 Xiong, J., Gao, S., Dui, W., Yang, W., Chen, X., Taverna, S. D., Pearlman, R. E., Ashlock, W.,
577     Miao, W., & Liu, Y. (2015). Distinct nucleosome distribution patterns in two structurally

578      and functionally differentiated nuclei of a unicellular eukaryote. *BioRxiv*, 018754.
579      https://doi.org/10.1101/018754

580 Xiong, J., Gao, S., Dui, W., Yang, W., Chen, X., Taverna, S. D., Pearlman, R. E., Ashlock, W.,
581      Miao, W., & Liu, Y. (2016). Dissecting relative contributions of cis- and trans-determinants
582      to nucleosome distribution by comparing Tetrahymena macronuclear and micronuclear
583      chromatin. *Nucleic Acids Research*, *44*(21), 10091–10105.
584      https://doi.org/10.1093/NAR/GKW684

585 Yao, M. C., Zheng, K., & Yao, C. H. (1987). A conserved nucleotide sequence at the sites of
586      developmentally regulated chromosomal breakage in tetrahymena. *Cell*, *48*(5), 779–788.
587      https://doi.org/10.1016/0092-8674(87)90075-4

588

589 **Figure Legends**

590 **Figure 1.** Flowchart summarizing the isolation of *Tetrahymena thermophila* MIC and MAC (see
591 text for details).

592 **Figure 2.** MIC and MAC are primarily distinguished by the PI signal. This flow cytometry dot
593 plot demonstrates the profiles for MIC and MAC nuclei of *Tetrahymena thermophila* after
594 staining with propidium iodide. The y-axis represents an autofluorescent signal generated in the
595 blank channel and the x-axis represents the fluorescent signal of PI stained samples. Signals in
596 the lower left hand corner of the dot plot outside the gated MIC and MAC signals (circled)
597 represent debris in the samples.

598 **Figure 3**. **Figure 3**. MIC enrichment can be quantified through uniquely mapped reads. A.
599 Schematic depiction of a *Tetrahymena* cell, with micronucleus indicated by red fill and
600 macronucleus indicated by green fill. B. Alignment of reads to the MIC (top) and MAC (bottom)
601 reference genomes. Reads derived from internally eliminated sequences (IES) map uniquely to
602 the MIC genome (pink regions) while reads that span the MAC excision boundary of an IES map
603 uniquely to the MAC genome.

604 **Figure 4.** The MIC and MAC of *Tetrahymena thermophila* are distinct in both physical size

605 of their nuclei as well as their genome size. This figure shows whole cell *Tetrahymena*
606 *thermophila* stained MIC and MAC (left) and MIC and MAC after homogenization (right).

607 **Figure 5.** Histogram of IES Retention scores for the MIC-enriched FACS Sample (blue) and
608 MAC-enriched FACS Sample (red). IES Retention scores for the Micronuclear FACS Sample
609 skew towards 1, indicating there are a high number of IESs in the sample. IES Retention scores
610 for the Macronuclear FACS Sample skew towards 0, indicating there are a low number of IESs
611 in the sample.