1  **A strategy for complete telomere-to-telomere assembly of ciliate**

2  **macronuclear genome using ultra-high coverage Nanopore data**

3

4  Guangying Wang[1], Xiaocui Chai[1,2], Jing Zhang[1], Wentao Yang[1,2],

5  Chuanqi Jiang[1], Kai Chen[1], Wei Miao[1,3,4]*, Jie Xiong[1]*

6

7  [1]Key Laboratory of Aquatic Biodiversity and Conservation, Institute of

8  Hydrobiology, Chinese Academy of Sciences, Wuhan 430072, China

9  [2]University of Chinese Academy of Sciences, Beijing 100049, China

10  [3]CAS Center for Excellence in Animal Evolution and Genetics, Kunming

11  650223, China

12  [4]State Key Laboratory of Freshwater Ecology and Biotechnology of China,

13  Wuhan 430072, China

14

15  *Corresponding author

16  Email:

17  miaowei@ihb.ac.cn

18  xiongjie@ihb.ac.cn

19

20

21

## ABSTRACT

Ciliates contain two kinds of nuclei: the germline micronucleus (MIC) and the somatic macronucleus (MAC) in a single cell. The MAC usually have fragmented chromosomes. These fragmented chromosomes, capped with telomeres at both ends, could be gene size to several megabases in length among different ciliate species. So far, no telomere-to-telomere assembly of entire MAC genome in ciliate species is finished. Development of the third generation sequencing technologies allows to generate sequencing reads up to megabases in length that could possibly span an entire MAC chromosome. Taking advantage of ultra-long Nanopore reads, we established a simple strategy for the complete assembly of ciliate MAC genomes. Using this strategy, we assembled the complete MAC genomes of two ciliate species *Tetrahymena thermophila* and *Tetrahymena shanghaiensis*, composed of 181 and 214 chromosomes telomere-to-telomere respectively. The established strategy as well as the high-quality genome data will provide a useful approach for ciliate genome assembly, and a valuable community resource for further biological, evolutionary and population genomic studies.

**Introduction**

Ciliate separates its germline and somatic genetic information by maintaining two kinds of functionally distinct nuclei: the diploid micronucleus (MIC), and the polyploid macronucleus (MAC) (Gorovsky, 1973; Lynn, 2008). The MIC, like other eukaryotes, usually contains long chromosomes with centromeres and capped by telomeres. In general, the MAC genome comes from the MIC genome through a so-called MAC differentiation process in the sexual stage (conjugation) of ciliate (Orias, 2000). In MAC differentiation, the MIC-like chromosomes are fragmented into small pieces at the chromosome breakage sites (CBSs), and the internal eliminated sequences (IESs), which contain transposable elements, are removed (Orias, 2000). This process finally results in the MAC containing fragmented chromosomes with length range from gene size to several megabases, and capped by telomere sequences at both ends but without centromeres.

Development of the third generation sequencing technologies, e.g. the Nanopore sequencing, allows to generate sequencing reads up to megabases in length (Jain et al., 2018), and thus could sometimes sequence an entire MAC chromosome of ciliate by a single read. The generation of such long sequencing reads gives the opportunity to assemble more complete MAC genomes of ciliates.

3

64    Here, we reported a simple strategy which was used to assemble the

65    complete genome of *T. thermophila* and *T. shanghaiensis* using high

66    coverage Nanopore sequencing data.

67

68

## 69    Materials and Methods

### 70    Cell culture and DNA extraction

71    *T. thermophila* SB210 and *T. shanghaiensis* (ATCC accession:

72    205039) cells were grown in SPP medium (Cassidy-Hanley, 2012) and

73    harvested at a density of 250,000 cells/ml. The total DNA was extracted

74    using the Blood & Cell Culture DNA Midi Kit (Q13343, Qiagen, CA, USA)

75    following the manufacturer's protocol. The DNA was then purified using

76    the Agencourt AMPure XP beads (A63881, BECKMAN), and the DNA

77    quality and quantity was tested using both NanoDrop One UV-Vis

78    spectrophotometer (Thermo Fisher Scientific, USA) and Qubit 3.0

79    Fluorometer (Invitrogen, USA).

80

### 81    Nanopore sequencing

82    Approximately 10 µg of DNA was size-selected (10-50 Kb) using Blue

83    Pippin (Sage Science, Beverly, MA), and sequencing library was

84    constructed using the Ligation sequencing 1D kit (SQK-LSK108, ONT, UK)

85   according to the manufacturer's instructions. Each library was sequenced

86   on R9.4 FlowCells using the PromethION sequencer (ONT, UK) for 48

87   hours. Base calling was subsequently performed on fast5 files using the

88   ONT Albacore software (v0.8.4), and the "passed filter" reads (high

89   quality data) were used for downstream analysis.

90

91   **Genome assembling and polishing**

92   Genome assembling was performed using 60X Nanopore datasets.

93   Assemblers, including CANU (Koren et al., 2017), NECAT

94   (https://github.com/xiaochuanle/NECAT), SHASTA

95   (https://github.com/chanzuckerberg/shasta), Flye (Kolmogorov, Yuan, Lin,

96   & Pevzner, 2019), and wtdbg2 (Ruan & Li, 2019), were used. The

97   parameters for the assemblers are listed as follows: 1) CANU, -fast

98   genomeSize=100m; 2) NECAT, GENOME_SIZE=100000000

99   MIN_READ_LENGTH=3000; 3) SHASTA, default settings; 4) Flye, -g

100   100m; 5) wtdbg2, default settings. The performance of CANU and NECAT

101   far better than three other assemblers in assembling the MAC

102   chromosomes capped with telomere sequences in both ends. Comparing

103   to CANU, the time cost of NECAT was far less than CANU, and thus

104   NECAT was recommended. Quickmerge

105   (https://github.com/mahulchak/quickmerge) was used to merge the

106  un-closed scaffolds to the 60X genome assemblies (command line:

107  merge_wrapper.py un-closed_scaffolds 60X_assembly). After each round

108  of merging, the closed scaffolds (MAC chromosomes) were extracted,

109  and the left un-closed scaffolds were used to perform the next round of

110  merging. After that, an addition round of merging between the un-closed

111  scaffolds and error corrected telomere-sequences-containing reads was

112  performed using miniasm (-1 -2 -c 1) (Li, 2016). Final genome polishing

113  was performed based on the Illumina sequencing data using Pilon

114  (https://github.com/broadinstitute/pilon).

115

116

## Results and Discussion

118  *T. thermophila* is a very useful unicellular model organism for

119  molecular and cellular biology (Ruehle, Orias, & Pearson, 2016). In 2006,

120  the MAC genome of *T. thermophila* has been sequenced using the

121  Sanger method (Coyne et al., 2008; Eisen et al., 2006), which greatly

122  accelerated the studies using *Tetrahymena* system. The current MAC

123  genome assembly (103.0 Mb,

124  http://ciliate.org/index.php/home/downloads) of *T. thermophila* has 1158

125  scaffolds, among which 128 (~58.9 Mb) were capped by telomeres with

126  C4A2 repeats at 5'-end and G4T2 repeats at 3'-end (hereafter defined as

127 closed scaffolds) and could be regarded as complete MAC chromosomes.

128 However, about a half of genome sequences, composed of 1030

129 scaffolds, are still not assembled as complete MAC chromosomes

130 (hereafter defined as un-closed scaffolds).

131 About 1000X Nanopore sequencing data (total DNA of both MAC and

132 MIC, reads N50: 25.8 Kb) were obtained to finish the MAC genome

133 assembly. Comparison of different third-generation sequencing data

134 assemblers , including CANU, NECAT, SHASTA, Flye and wtdbg2, were

135 performed. In practice, CANU and NECAT showed better performance on

136 assembling closed scaffolds compared to other assemblers. We divided

137 the ~1000X Nanopore data into different parts, each with ~60X data, and

138 individually assembled them (Figure 1). We have two reasons to do this

139 division: 1) The MIC reads (contaminations) could be limited below 3X

140 (the copy number ratio between MAC and MIC is 45:2), which will usually

141 be  filtered by  genome assemblers (Jain et al., 2018); 2) At 60X

142 coverage, CANU and NECAT already have good assembling

143 performance and the time cost of assembling could be greatly reduced.

144 We started from the 1158 scaffolds in current genome assembly of *T.*

145 *thermophila* (Figure 1), and divided these scaffolds into two parts: 1) 128

146 closed scaffolds which assembled as complete MAC chromosomes; 2)

147 1030 un-closed scaffolds which have not been assembled as MAC

148    chromosomes. For the 128 closed scaffolds, three of them still have gaps

149    (one per each). These gaps were easily closed by aligning the three

150    scaffolds to the 60X Nanopore data assemblies. The left 1030 un-closed

151    scaffolds were iteratively merged with each assembled genome using

152    60X Nanopore data (Figure 1). After six rounds of merging using

153    quickmerge, 34 closed scaffolds were newly obtained. After that, we

154    extracted the 256,181 raw telomere-sequence-containing reads (TSCR,

155    reads N50, 28.5 Kb) from Nanopore data (Figure 1), and sequencing

156    errors were corrected using NECAT. These error corrected TSCR were

157    aligned to the left scaffolds using minimap2, and followed by a new round

158    of assembly using miniasm (Figure 1), and additional 12 scaffolds with

159    telomere sequences capped at both ends were obtained, and only six

160    scaffolds (3.3 Mb) could not be resolved. To close these six scaffolds, we

161    manually checked the overlaps between TSCR and these scaffolds

162    (Figure 1), and all of them were closed by trimming their terminal

163    sequences and re-merging with TSCR.

164       In summary, the complete MAC genome (102.9 Mb) with a total of 181

165    MAC chromosomes (including rDNA mini-chromosome) were obtained.

166    These MAC chromosomes were re-named from 1 to 181 by their order

167    along the five MIC chromosomes. Figure 2 showed the full panel of the

168    181 MAC chromosomes. The longest MAC chromosome is 3.26 Mb in

169    length, and the shortest one (excluding rDNA mini-chromosome) is 38 Kb

170    in length. The real N50 of the MAC genome is ~891 Kb. A total of 22

171    classes of repetitive sequences, which masked 5.2% MAC genome, were

172    identified by RepeatModeler. The repetitive sequences in the MAC are

173    not randomly distributed, most of them are enriched in the MAC

174    chromosomes and derived from the pericentromeric and subtelomeric

175    regions of MIC chromosomes (Hamilton et al., 2016; Xiong et al., 2019).

176    In particular, we also found some new genes which missed in the current

177    genome assembly, for example, the alpha 2 subunit of the proteasome.

178      To test the applicability of this strategy, we generated ~900X

179    Nanopore sequencing data (reads N50: 30.8 Kb) of *T. shanghaiensis*.

180    Instead using pre-existed assembly, we started from a 60X *de novo*

181    assembly by NECAT, and then followed the strategy showing in Figure 1.

182    After eight rounds of merging using quickmerge and a round of assembly

183    using miniasm, and followed by additional manual checking, we finally got

184    the complete genome of *T. shanghaiensis* with 214 MAC chromosomes

185    (92.0 Mb) which capped with telomere sequences at both ends. Genome

186    assembly statistics of the two *Tetrahymena* species are shown in Table 1.

187    We anticipate that the established strategy can probably be used directly

188    or with a slight adaptation to assemble complete MAC genomes of other

189    ciliate species.

190

## **Acknowledgments**

199

## **Author contributions**

W.M. and J.X. designed the project. J.X., G.W. and W.Y. assembled and annotated the genome. X.C., J.Z., C.J. and K.C. prepared DNA samples for sequencing. J.X. and G.W. wrote the manuscript. All authors read, revised and approved the final manuscript.

205

## **Data accessibility**

The complete genome sequences of *T. thermophila* and *T. shanghaiensis* can be accessed from http://ciliate.ihb.ac.cn/tcgd/download.html.

209

210

211

## References

213 Cassidy-Hanley, D. M. (2012). *Tetrahymena* in the laboratory: strain

214 resources, methods for culture, maintenance, and storage. *Methods in*

215 *Cell Biology, 109*, 237-276. doi:10.1016/B978-0-12-385967-9.00008-6

216 Coyne, R. S., Thiagarajan, M., Jones, K. M., Wortman, J. R., Tallon, L. J.,

217 Haas, B. J., . . . Methe, B. A. (2008). Refined annotation and assembly

218 of the *Tetrahymena thermophila* genome sequence through EST

219 analysis, comparative genomic hybridization, and targeted gap closure.

220 *BMC Genomics, 9*, 562. doi:10.1186/1471-2164-9-562

221 Eisen, J. A., Coyne, R. S., Wu, M., Wu, D. Y., Thiagarajan, M., Wortman,

222 J. R., . . . Orias, E. (2006). Macronuclear genome sequence of the

223 ciliate *Tetrahymena thermophila*, a model eukaryote. *PloS Biology,*

224 *4*(9), 1620-1642. doi:10.1371/journal.pbio.0040286

225 Gorovsky, M. A. (1973). Macro - and micronuclei of *Tetrahymena*

226 *pyriformis*: A model system for studying the structure and function of

227 eukaryotic nuclei. *The Journal of Protozoology, 20*(1), 19-25.

228 doi:10.1111/j.1550-7408.1973.tb05995.x

229 Hamilton, E. P., Kapusta, A., Huvos, P. E., Bidwell, S. L., Zafar, N., Tang,

230 H. B., . . . Coyne, R. S. (2016). Structure of the germline genome of

231 *Tetrahymena thermophila* and relationship to the massively rearranged

232    somatic genome. *Elife, 5*. doi:10.7554/eLife.19090.001

233    Jain, M., Koren, S., Miga, K. H., Quick, J., Rand, A. C., Sasani, T. A., . . .

234    Fiddes, I. T. (2018). Nanopore sequencing and assembly of a human

235    genome with ultra-long reads. *Nature Biotechnology, 36*, 338-345.

236    doi:10.1038/nbt.4060

237    Kolmogorov, M., Yuan, J., Lin, Y., & Pevzner, P. A. (2019). Assembly of

238    long, error-prone reads using repeat graphs. *Nature Biotechnology, 37*,

239    540-546. doi:10.1038/s41587-019-0072-8

240    Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., &

241    Phillippy, A. M. (2017). Canu: scalable and accurate long-read

242    assembly via adaptive k-mer weighting and repeat separation.

243    *Genome research, 27*(5), 722-736. doi:10.1101/gr.215087.116

244    Li, H. (2016). Minimap and miniasm: fast mapping and de novo assembly

245    for noisy long sequences. *Bioinformatics, 32*(14), 2103-2110.

246    doi:10.1093/bioinformatics/btw152

247    Lynn, D. (2008). *The ciliated protozoa: characterization, classification,*

248    *and guide to the literature* (3rd ed.). Berlin: Springer.

249    Orias, E. (2000). Toward sequencing the Tetrahymena genome:

250    exploiting the gift of nuclear dimorphism. *Journal of Eukaryotic*

251    *Microbiology, 47*(4), 328-333. doi:10.1111/j.1550-7408.2000.tb00057.x

252    Ruan, J., & Li, H. (2019). Fast and accurate long-read assembly with

253     wtdbg2. *Nature Methods*. doi:10.1038/s41592-019-0669-3

254   Ruehle, M. D., Orias, E., & Pearson, C. G. (2016). *Tetrahymena* as a

255     Unicellular Model Eukaryote: Genetic and Genomic Tools. *Genetics,*

256     *203*(2), 649-665. doi:10.1534/genetics.114.169748

257   Xiong, J., Lu, X., Zhou, Z., Chang, Y., Yuan, D., Tian, M., . . . Orias, E.

258     (2012). Transcriptome analysis of the model protozoan, *Tetrahymena*

259     *thermophila*, using deep RNA sequencing. *PloS one, 7*(2), e30630.

260     doi:10.1371/journal.pone.0030630

261   Xiong, J., Yang, W., Chen, K., Jiang, C., Ma, Y., Chai, X., . . . Miao, W.

262     (2019). Hidden genomic evolution in a morphospecies-The landscape

263     of rapidly evolving genes in *Tetrahymena*. *PloS Biology, 17*(6),

264     e3000294. doi:10.1371/journal.pbio.3000294

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284
285

## **Figure legends**

287

**Figure 1. Diagram showing the strategy to assemble complete MAC genome of ciliate.** M1 to Mn, the un-closed scaffolds in each round (1 to n) which do not have telomere sequences in both ends. M_miniasm means the un-closed scaffolds after merging using miniasm. C1 to Cn, the closed scaffolds (MAC chromosomes) in each round (1 to n) which have telomere sequences in both ends. C_miniasm means the closed scaffolds (MAC chromosomes) after merging using miniasm. C_manual means the closed scaffolds after the manual checking of overlaps between TSCR and un-closed scaffolds (trimmed).

297

**Figure 2. A full panel of 181 MAC chromosomes of *T. thermophila*.** For each MAC chromosome, the pink boxes represents the predicted genes; the red boxes represent all the genes that have been named in TGD wiki (http://ciliate.org/); the blue histogram represents the gene expression profile across the chromosome in vegetative growth (Xiong et al., 2012).

304

305

306

307

308

309 **Table 1.** Genome assembly statistics of *T. thermophila* and *T.*

310 *shanghaiensis*

| | Assembly in this study | | Current assembly | |
|---|---|---|---|---|
| | *T. thermophila* | *T. shanghaiensis* | *T. thermophila*[†] | *T. shanghaiensis*[‡] |
| Assembly size (Mb) | 102.9 | 92.0 | 103.0 | 95.6 |
| Number of scaffolds | 181 | 214 | 1,158 | 2,660 |
| Closed Scaffolds[§] | 181 | 214 | 128 | 31 |
| Scaffold N50 (Kb) | 891.3 | 620.0 | 520.9 | 153.6 |
| Longest scaffold size (Mb) | 3.26 | 1.98 | 2.22 | 0.79 |
| Mean scaffold size (Kb) | 568.5 | 430.0 | 89.0 | 36.0 |
| "N" gaps (Kb)[¶] | 0 | 0 | 63.7 | 90.0 |

311

312 [†]Genome data from website: http://ciliate.org/index.php/home/downloads
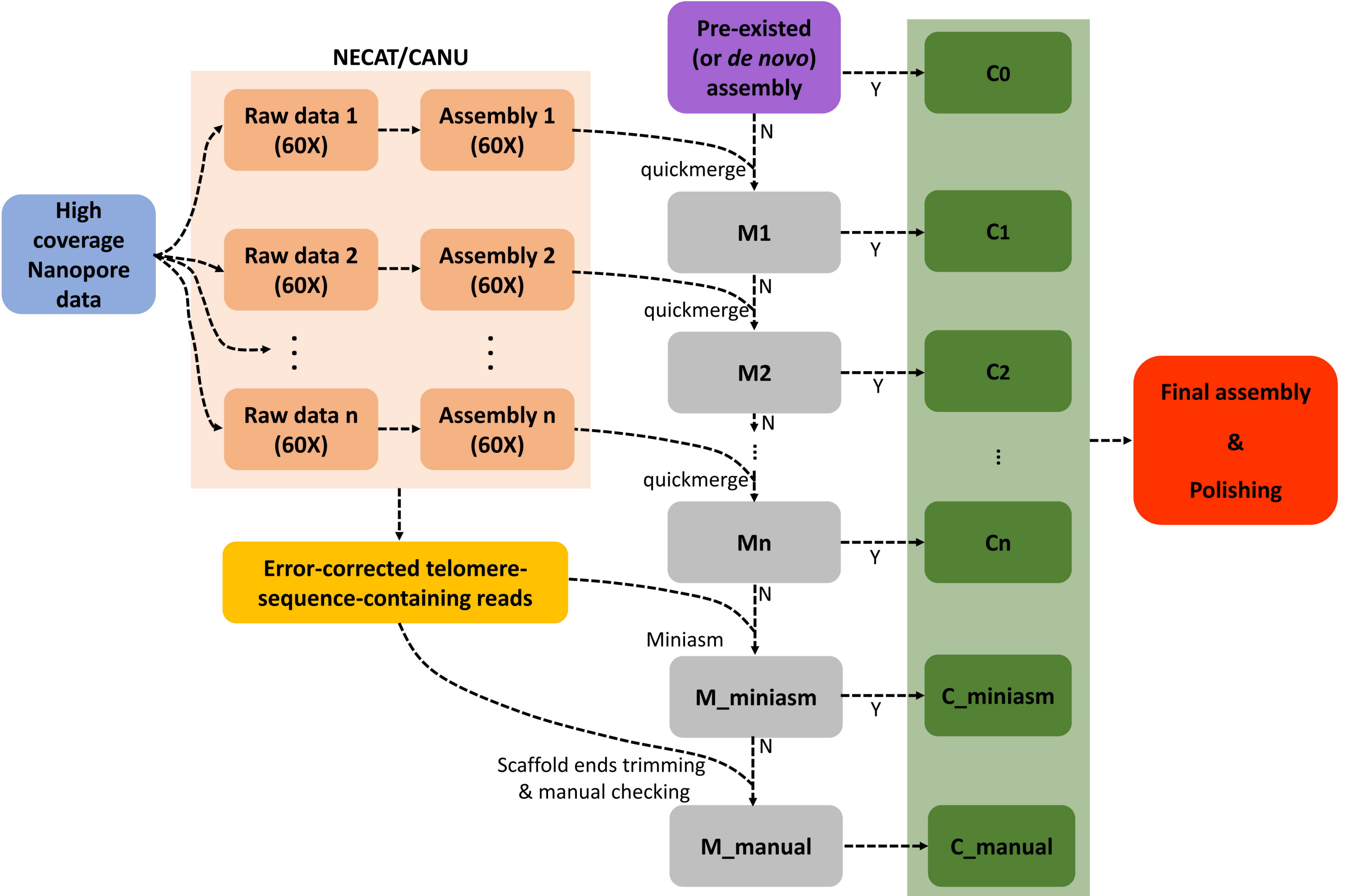
313 [‡]Genome data from (Xiong et al., 2019)

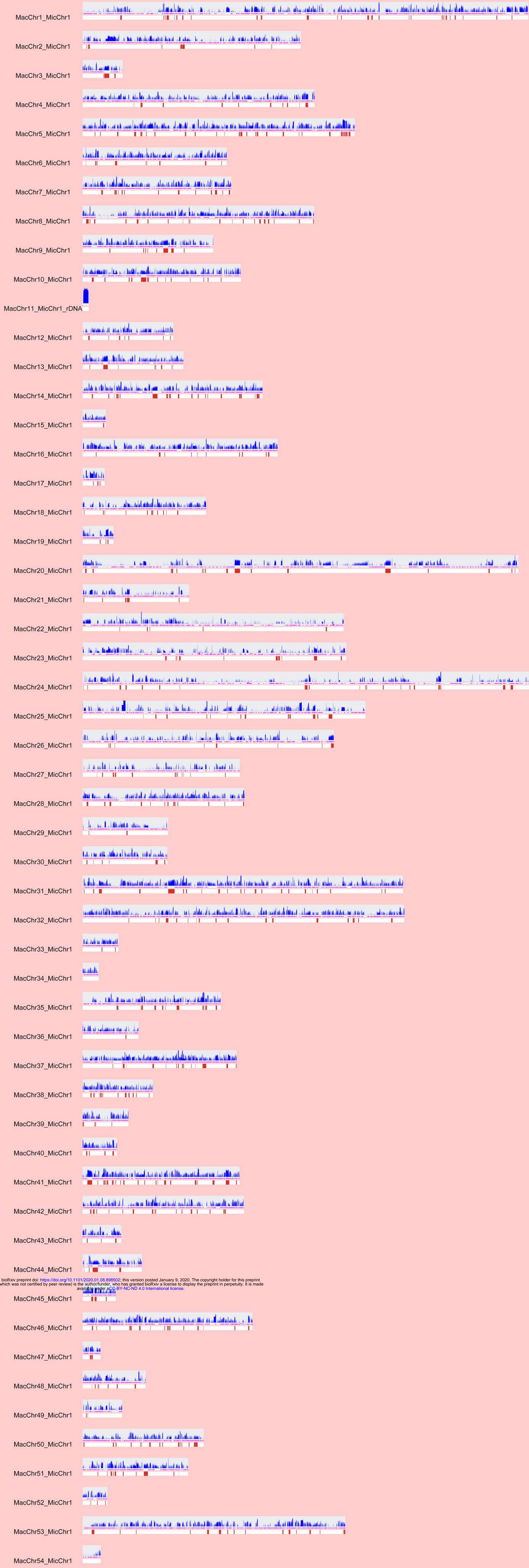314 [§]Scaffolds capped with telomeres at both ends

315 [¶]Sum of all "N" nucleotides present in the genome assembly
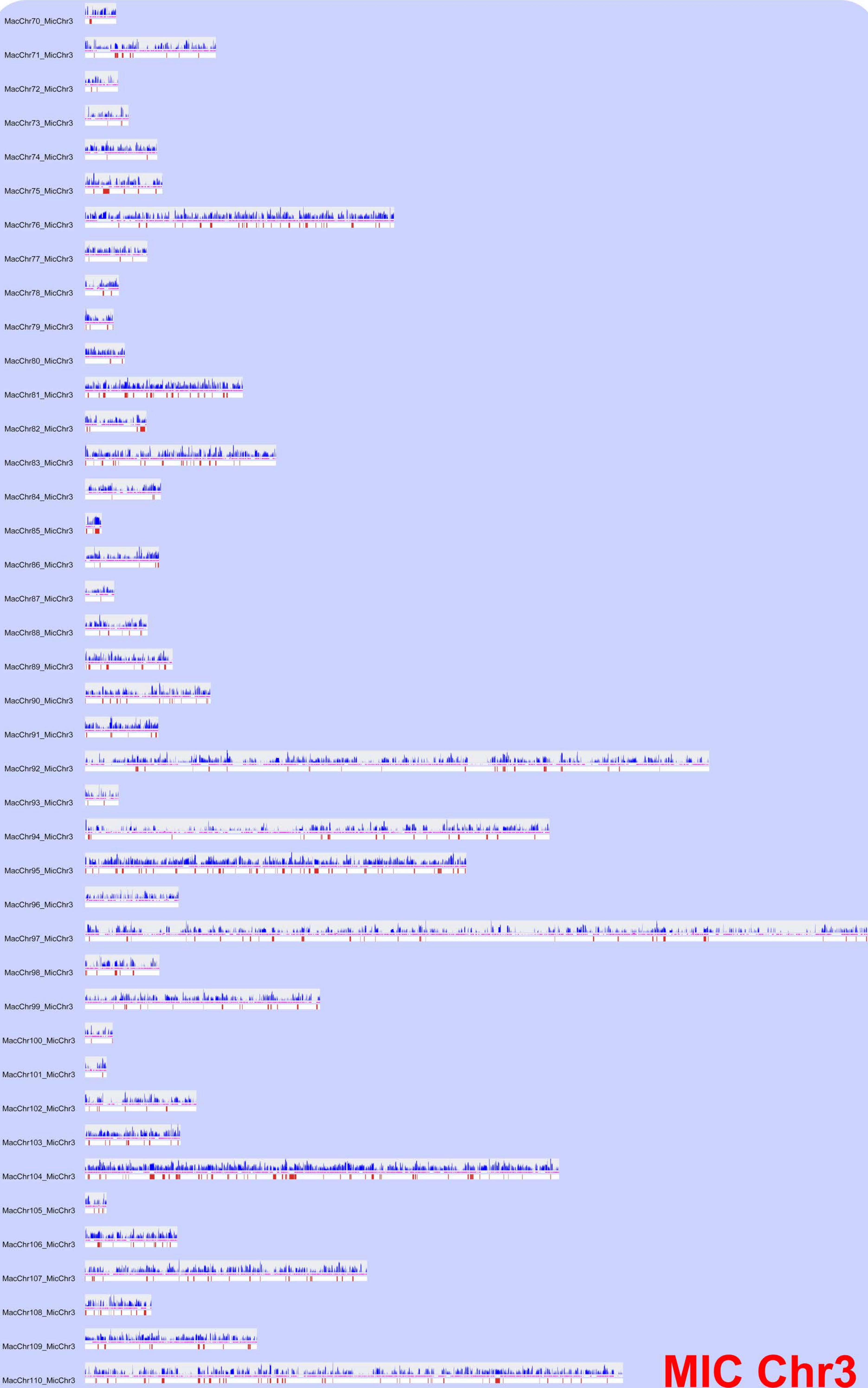
316

The complete MAC genome of *Tetrahymena thermophila*
181 chromosomes