1 **Long-Read Genome Assembly and Gene Model Annotations for the Rodent Malaria Parasite**

2 *Plasmodium yoelii* **17XNL**

3

4 Mitchell J. Godin[1]*, Aswathy Sebastian[2]*, Istvan Albert[1,2#], Scott E. Lindner[1#]

5 * Equal Contributions, # Co-Corresponding Authors

6 1. Department of Biochemistry and Molecular Biology, The Huck Center for Malaria Research, The

7 Center for Eukaryotic Gene Regulation, Pennsylvania State University, University Park, PA, 16802

8 2. Huck Institutes of the Life Sciences, Pennsylvania State University, University Park, PA, 16802

9
10 **Corresponding Authors:**

11 Istvan Albert, Ph.D.

12 206C Life Sciences Building

13 University Park, PA 16802

14 P: +1.814.865.2281

15 E: iua1@psu.edu

16 ORCID: 0000-0001-8366-984X

17

18 Scott E. Lindner, Ph.D.

19 W230B Millennium Science Complex

20 University Park, PA 16802

21 P: +1.814.867.4062

22 E: Scott.Lindner@psu.edu

23 ORCID: 0000-0003-1799-3726

24

**Abstract**

Malaria causes over 200 million infections and over 600 thousand fatalities each year, with most cases attributed to a human-infectious *Plasmodium* species, *Plasmodium falciparum*. Many rodent-infectious *Plasmodium* species, like *Plasmodium berghei*, *Plasmodium chabaudi*, and *Plasmodium yoelii,* have been used as genetically tractable model species that can expedite studies of this pathogen. In particular, *P. yoelii* is an especially good model for investigating the mosquito and liver stages of parasite development because key attributes closely resemble those of *P. falciparum*. Because of its importance to malaria research, in 2002 the 17XNL strain of *P. yoelii* was the first rodent malaria parasite to be sequenced. While sequencing and assembling this genome was a breakthrough effort, the final assembly consisted of >5000 contiguous sequences that impacted the creation of annotated gene models. While other important rodent malaria parasite genomes have been sequenced and annotated since then, including the related *P. yoelii* 17X strain, the 17XNL strain has not. As a result, genomic data for 17X has become the *de facto* reference genome for the 17XNL strain while leaving open questions surrounding possible differences between the 17XNL and 17X genomes. In this work, we present a high-quality genome assembly for *P. yoelii* 17XNL using HiFi PacBio long-read DNA sequencing. In addition, we use Nanopore long-read direct RNA-seq and Illumina short-read sequencing of mixed blood stages to create complete gene models that include not only coding sequences but also alternate transcript isoforms, and 5' and 3' UTR designations. A comparison of the 17X and this new 17XNL assembly revealed biologically meaningful differences between the strains due to the presence of coding sequence variants. Taken together, our work provides a new genomic and gene expression framework for studies with this commonly used rodent malaria model species.

**Introduction**

Malaria remains a major global health burden (WHO Malaria Report 2022, (1)), with most of the 600,000 fatalities resulting from infection by human-infectious *Plasmodium falciparum*. The use of rodent-infectious model species has been instrumental to better understand those species that cause human disease due to high levels of genetic and physiological conservation across species (2). Researchers have routinely used these rodent model species, such as *P. yoelii*, *P. berghei*, and *P. chabaudi*, to investigate the entire *Plasmodium* life cycle, as genetic manipulations have long been rapid and rigorous in these species (2). We and others study *P. yoelii*, which is an especially good model for the mosquito and liver stages of *P. falciparum* parasite development (2). This is partly because *P. yoelii* mosquito stage parasites develop at a similar pace as do those of *P. falciparum*, and their sporozoites are less promiscuous than *P. berghei* sporozoites (2). Because of this, many studies of genetically attenuated parasite (GAP) vaccine candidates based upon sporozoites have recently included the use of *P. yoelii* as a pre-clinical model system (3). In support of this, large-scale analyses of gene expression of *P. yoelii* now match those available for *P. berghei* in many ways (4-9). For these reasons, *P. yoelii* has been an important malaria parasite used as a proxy for *P. falciparum* in pre-clinical and discovery phase studies.

Intuitively, genetic studies of any species are best conducted with accurate genome assemblies and gene models. Therefore, several species of *Plasmodium* parasites were the subject of early whole-genome sequencing efforts in the late 1990s and early 2000s (10, 11). This work provided a genome assembly of the human-infectious *P. falciparum* parasite with 14 nuclear chromosomes and the two organellar genomes of its mitochondrion and apicoplast (11). In addition, gene models for *P. falciparum* were annotated with introns/exons, with further improvements establishing 5'/3' untranslated regions (UTRs) and transcript isoforms (12, 13). Similarly, the rodent-infectious *Plasmodium berghei* ANKA parasite was originally sequenced in 2005, resulting in a genome assembly with 7,497 contiguous

3

72     sequences (contigs) that were later reduced to 16 contigs with a hybrid Illumina and 454 sequencing

73     approach in 2014, and then further refined using PacBio sequencing in 2016 (14-16). Prior to this, the

74     non-lethal *P. yoelii* 17XNL strain was the first rodent malaria parasite sequenced in 2002, which used

75     ABI3700 sequencers and yielded a genome assembly of over 5,000 contigs (10). The *P. yoelii* 17X strain,

76     from which 17XNL was derived, was sequenced in 2014 alongside PbANKA using the same Illumina and

77     454 sequencing approach to similarly establish a 16 contig genome assembly (15). With the advent of

78     more accurate long-read sequencing technologies, there has been a renewed interest in sequencing the

79     *Plasmodium* genomes and transcriptomes, including those of another *P. yoelii* strain, PyN67, which has

80     been used to study genetic polymorphisms and drug responses (17). In addition, the genomes of other

81     apicomplexan parasites, such as *Cryptosporidium* and *Babesia* species, have now been established using

82     a combination of long-read Nanopore sequencing and short-read Illumina sequencing (18-20).

83         Although their genomes have been updated and are conveniently provided on PlasmoDB.org,

84     Py17X and PbANKA have gene models that largely reflect the coding sequences, but not their UTRs

85     despite the availability of RNA-seq data that could be used to approximate them (4-9, 15, 21-25). Finally,

86     while the 17XNL strain of *P. yoelii* remains a highly used laboratory strain worldwide, its reference

87     genome and gene models have not been revisited since 2002, and thus its genome assembly and gene

88     models remain highly fragmented and incomplete. As a result, most researchers use the genome

89     assembly and gene models of the related *P. yoelii* 17X strain as a proxy when working with the 17XNL

90     strain and must operate under the assumption that the genomes of the two strains are effectively the

91     same. However, this prompts a few important questions. How similar are the 17X and 17XNL strains? In

92     what ways are they truly suitable proxies for one another? Given the state of the 17XNL genome

93     assembly and the limited gene models available for both strains, these questions could not be accurately

94     addressed. However, these kinds of questions can now be more rigorously addressed with the inclusion

95     of long-read DNA sequencing. The long sequence reads produced by PacBio and Nanopore approaches

96   better facilitate the scaffolding of long, contiguous sequences in a *de novo* assembly, even for complex

97   genomes that have extreme AT-content and/or high degrees of repetitiveness, such as found with

98   *Plasmodium* species (26, 27). Additionally, by combining long-read and short-read sequencing, the

99   strengths of each can be used to polish the assembly to reduce systematic errors introduced by each of

100   the different methodologies.

101       Therefore, here we have created a high-quality reference genome and gene model annotation

102   for the *P. yoelii* 17XNL strain that we have used to address these outstanding questions. We utilized HiFi

103   PacBio DNA-seq to create a Py17XNL reference genome with 16 high confidence/high accuracy contigs.

104   Even without any polishing efforts, this approach outperformed a parallel effort using a hybrid

105   Nanopore long-read DNA-seq/Illumina short-read DNA-seq method by several key metrics, including its

106   assembly quality and the reduction of gaps. Furthermore, we created gene annotations for genes

107   transcribed in asexual and sexual blood stages using a combination of Nanopore direct RNA-seq and our

108   pre-existing Illumina RNA-seq datasets. These annotations include definitions of introns, exons, 5', and

109   3' UTRs, and transcript isoforms expressed in asexual and sexual blood stages. Using these data, we

110   compared the genomic variance between the Py17XNL and Py17X strains to gain insight into the

111   differences between the two strains and identified that most sequence variants reside in intergenic

112   regions, whilst variation in the coding sequence of a select few genes could result in meaningful changes

113   in Py17XNL parasite biology.

114

115 **Results**

116

117 <u>A Comparison of Genome Assembly Approaches: PacBio HiFi vs. Nanopore/Illumina Hybrid Sequencing</u>

118 *P. yoelii* 17XNL remains a commonly studied rodent malaria strain. Yet, its genome assembly

119 remains highly fragmented and consists of over 5000 contigs as generated in 2002 (10). Consequently,

120 most researchers use the reference genome of the related Py17X strain as a substitute for Py17XNL

121 without knowing how appropriate it is to use it as a genomic proxy (15). To resolve these questions, we

122 created a high-quality genome assembly of *P. yoelii* 17XNL Clone 1.1 obtained from BEI Resources, which

123 is the common origin of this strain of parasites for many laboratories. Because several sequencing

124 methodologies are now commonly used to assemble whole genomes, we used Nanopore, PacBio, and

125 Illumina sequencing with DNA- or RNA-based libraries to determine an optimal approach to create a

126 genome assembly with associated gene models for *P. yoelii* 17XNL. Nanopore ligation-based long-read

127 DNA sequencing is currently favored by many researchers as it can provide extremely long sequence

128 reads, resolve long stretches of repetitive regions, and assemble long structural variants in the genome

129 (26, 27). HiFi PacBio DNA sequencing provides very high accuracy due to the sequencing of circular

130 consensus sequences (ccs) of ~10kb DNA fragments, providing a middle ground between the sequencing

131 sizes provided by Illumina and Nanopore sequencing (27). We explored several data analysis protocols

132 for combining data from different platforms to optimize this genome assembly. As widely available

133 Nanopore sequencing chemistries (Q10) yield a systematic error, Nanopore data are often paired with

134 Illumina data to provide error correction. The sequencing error rates from Illumina are typically above

135 99.9% and can be used with polishing algorithms to identify errors in assemblies that were produced

136 with long, noisy reads (18, 19). A detailed outline of our experimental methods is included in

137 Supplemental Figure 1. Briefly, swiss webster outbred mice were infected with Py17XNL Clone 1.1

138 parasites that had been passaged only once following receipt from BEI Resources to create a genome

139 assembly reflective of the current stocks available in the depository. Upon reaching 1-3% parasitemia,

6

140    mice were euthanized, white cells were depleted by cellulose, and red blood cells (RBCs) were lysed by

141    saponin. The parasite pellets were used to produce high molecular weight genomic DNA using the NEB

142    Monarch HMW DNA Extraction Kit for Cells and Blood as previously described (28). DNA purity, quantity,

143    and fragment lengths were determined to all be high quality by NanoDrop, Qubit, and TapeStation

144    measurements, respectively (Supplementary Table 1, Supplementary Figure 2)). This approach yielded

145    DNA fragments of higher quality and higher molecular weight than the Qiagen QIAamp DNA Blood Mini

146    Kit that is routinely used in our laboratory and in others. Matched gDNA samples were sequenced on a

147    Nanopore MinION R9.4.1 flow cell using the ligation sequencing kit, as well as on an Illumina NextSeq

148    550 using the Illumina DNA PCR-Free kit. In parallel, Py17XNL HMW gDNA was sequenced on a PacBio

149    Sequel using a PacBio SMRT cell. The raw reads from the PacBio sequencing run were converted into

150    circular consensus sequences using the CCS algorithm.

151         To assess the quality of both Nanopore sequencing runs, we utilized Nanoplot, a quality control

152    plotting suite specifically for long-read sequencing data (Supplemental Figure 3, Supplemental Table 2))

153    (29). The Nanopore sequencing runs for both replicate one and two resulted in an overall average read

154    length of 16,706 bases, with an average Qscore of 11.3 (Supplemental Table 2). During this study, an

155    improved high-accuracy base calling algorithm from Nanopore was released, which we also tested to

156    see if it could improve our read quality. Upon re-basecalling the fast5 files, we saw a considerable

157    increase in Qscore, from 11.3 to 14.1, even without the quality score filter that is imposed with the

158    default, fast basecalling algorithm (Supplemental Figure 3, Supplemental Table 2). Despite this

159    improvement in quality, there were no significant differences in the mean read length or the throughput

160    (Supplemental Table 2). Using PacBio ccs reads, we saw an improvement in accuracy to an average

161    Qscore of 36.3 (Figure 1A). The biggest difference came with throughput, which increased to

162    1,660,222,360 bases from 707,945,539 bases whilst still maintaining an average read length of 5,712.7

163    bases (Figure 1B, Supplemental Table 2).

7

164    We generated genome assemblies from both long-read datasets with the bioinformatic

165    workflows described in Figure 2. To create the Nanopore/Illumina hybrid genome assembly, we

166    assembled the Py17XNL Nanopore data using Flye (30) and scaffolded the contigs using the Py17X

167    genome as a guide in conjunction with the RagTag scaffolding program (30, 31). Finally, we layered error

168    correction onto it in a multi-step approach, first using nextpolish, followed by multiple rounds of

169    consensus generation based on Illumina data alignment and variant calling (Figure 2) (32). Through this

170    process, we were able to reduce the number of contigs down to 16, but at the cost of covering less of

171    the genome (95%) and introducing 34 misassembles (Figure 1C) as defined by the assembly evaluator

172    program Quast (33). The PacBio-based genome assembly was generated with the HiCanu program (34),

173    which produced a *de novo* genome assembly with 132 contigs (Figure 2). The resulting contigs were

174    filtered to contain only the target species by aligning them against the Py17X genome using minimap2

175    (35). Contigs that had a primary alignment length of >2% of the 17X reference chromosome were

176    assigned the matching chromosome names. A consensus genome was then created by aligning these

177    contigs with the 17X reference genome and filling in the missing genomic regions, mainly chromosomal

178    ends. This resulted in a final assembly of 23.08 Mb with 16 contigs (Figure 1C, Table 1). We have

179    adopted the higher quality PacBio-based genome assembly for the *P. yoelii* 17XNL strain for the rest of

180    our analyses and for provision to the community on PlasmoDB.org, which we term Py17XNL_2 to

181    distinguish it from the original genome assembly (Py17XNL_1) (10, 36). However, as both the hybrid

182    Nanopore/Illumina and PacBio assemblies are potentially valuable to our research community, both

183    assemblies have been publicly deposited in NCBI.

184

185    Nanopore Direct RNA sequencing provides new information to pre-existing gene models

186    We also set out to create more comprehensive gene models to increase the utility of the new

187    Py17XNL_2 genome assembly for the P. yoelii 17XNL strain. In the currently available gene models for *P.*

8

188     *berghei* (ANKA) and *P. yoelii* (17X, 17XNL) on PlasmoDB, only the coding sequences of genes are

189     provided with no designation of untranslated regions (UTRs), and little information is provided about

190     alternatively spliced transcripts. We generated gene models that provide complete transcript

191     information, including start and stop codons, transcription start and stop sites, and UTRs.

192     Experimentally, we performed Nanopore direct RNA-seq in biological duplicate to generate long

193     sequence reads of asexual and sexual blood stage transcripts. Briefly, total RNA was extracted from

194     parasitized mouse blood to create an RNA-seq library that was sequenced with the Nanopore Direct

195     RNA-Sequencing Kit (Supplemental Figure 1, Supplemental Figure 4, Supplemental Table 1). These direct

196     RNA-sequencing reads were quality controlled using Nanoplot with the same parameters described for

197     Nanopore ligation DNA-sequencing for both "Fast" and "High Accuracy" basecalling approaches

198     (Supplemental Figure 5) (29). In total, 429,888,068 bases were sequenced after combining the

199     replicates, with an average Qscore of 12 after high accuracy basecalling, which again outperformed the

200     fast basecalling approach (Supplemental Figure 5, Supplemental Table 2)). The mean read length across

201     replicates was 858 bases, with the longest read being 8,789 bases (Supplemental Figure 5, Supplemental

202     Table 2).

203             Gene models were created with two alternative methodologies using Nanopore direct RNA-seq

204     long reads alone or in combination with our previously published Illumina short-read RNA-seq of mixed

205     asexual and sexual blood stage parasites (Figure 2) (37). These parallel approaches are both informative,

206     given the strengths and limitations of both sequencing techniques. Nanopore direct RNA sequencing

207     provides information that allows us to identify long/full-length sequencing reads that initiate at the 3'

208     end of mRNAs (38). However, when the full-length mRNA is not sequenced, less information is provided

209     for the 5' end (38). This limitation is remedied by the strong depth and breadth of sequencing coverage

210     provided via Illumina sequencing. For both approaches, Nanopore RNA-seq reads were aligned to our

211     Py17XNL_2 genome using minimap2 (35). For gene models created with both Nanopore and Illumina

212 RNA-seq data, in parallel, the Illumina short reads were aligned to the Py17XNL_2 genome using Hisat2

213 (31375807). To create the gene models and assign gene names/descriptions, we used Braker2, as well as

214 reciprocal blast searches using the blastp program of the BLAST suite (39, 40). The Nanopore-only

215 approach helped us identify 5,683 genes, 5,828 mRNAs, 66 tRNAs, and 40 rRNAs. Using the

216 Nanopore/Illumina hybrid approach, we found 6,077 genes, 7,047 mRNAs, 66 tRNAs, and 40 rRNAs

217 (Table 1). Gene models that were generated using both Nanopore and Illumina reads more closely

218 matched the anticipated UTR length that was defined in recent *Plasmodium falciparum* transcriptomics

219 data (13). A representative example of this more comprehensive gene model is illustrated in Figure 3A.

220 Nanopore direct RNA-seq initiates at the 3' end due to the use of poly(dT) sequencing primers.

221 As a result, significantly higher coverage was obtained for the 3' UTRs than for 5' UTRs (Supplemental

222 Table 3). The higher coverage allows us to further analyze the 3' UTR length distribution for *P. yoelii*

223 17XNL, which is of interest as *cis*-regulatory elements are often found in this portion of eukaryotic

224 mRNAs (41). The majority of reads have a UTR length between 100 and 200 bp, with a mean length of

225 364 bp. The largest UTR reported for the H2B.Z histone variant mRNA, with 1994 nt (Figure 3B,

226 Supplemental Table 3). Compared to the most up-to-date *P. falciparum* transcriptome, which used

227 DAFT-Seq to resolve UTRs, P. yoelii 17XNL's 3' UTRs appear slightly shorter on average (Figure 3C) (13).

228

229 Comparison between Py17XNL_2 and reference genomes demonstrates the completeness of the

230 assembly

231 Using the Py17XNL_2 genome assembly and associated gene models, we compared our results

232 to the original Py17XNL genome (Py17XNL_1) and the Py17X reference genome. As anticipated, there

233 was a substantial reduction in the number of gaps/misassembles and greater genome coverage when

234 comparing Py17XNL_2 vs. Py17XNL_1 (Table 1). Although our Nanopore/Illumina-based genome

235 assembly (Py17XNL Nanopore) contained the same number of contigs as the PacBio-based Py17XNL_2

10

236    assembly, the significantly fewer misassemblies generated in the PacBio-based assembly provided a

237    more accurate reference genome for future research uses. Additionally, our final PacBio-based assembly

238    (Py17XNL_2) closely resembles that of Py17X, which was also created using recently developed

239    sequencing technologies (15). However, when compared to Py17X, the Py17XNL_2 assembly has lower

240    coverage of repetitive sequences at the sub-telomeric regions, which precluded us from robustly

241    assembling these regions for some chromosomes, requiring consensus generation based on alignments

242    to the Py17X reference genome. Similarly, the 6,086 new gene annotations more accurately represent

243    the anticipated number of genes for Py17XNL and more closely match those annotated in other

244    *Plasmodium* species (Table 1). Moreover, these new gene models include both coding sequences, UTRs,

245    and transcript isoforms, which are lacking in the provided gene models currently available on PlasmoDB

246    for this specific species. In addition to these assessment metrics, we determined the completeness of

247    the reference genome based on marker genes. To quantify this, we used a Benchmarking Universal

248    Single-Copy Orthologs (BUSCO) analysis that detects whether a predefined set of single-copy marker

249    genes in the *Plasmodium* lineage are present in these data (Figure 4) (42). This BUSCO dataset contains

250    3642 BUSCO groups from 23 different species, including *P. falciparum* 3D7, *P. yoelii* 17XNL, *P. vivax, P.*

251    *berghei* ANKA, *P. chabaudi,* and others. From this search, 3556/3642 (97.6%) of complete and single

252    copy BUSCOs were found to be present, indicating that this genome assembly and gene annotation has

253    a high level of completeness (Figure 4).

254

255    Variation between the Py17XNL and Py17X reference genomes primarily resides in the intergenic

256    regions and the ends of chromosomes

257          As Py17X is commonly used as an interchangeable proxy genome for Py17XNL, we sought to

258    determine what similarities and differences exist between the strains and how the differences may

259    impact genetic studies. We performed chromosome-wide alignments between the Py17X and

11

260    Py17XNL_2 genomic builds using the minimap2 (35) program and assessed the genome-wide variants

261    with paftools. We observed extensive linear agreement between the two strains, with 99.9% of the

262    Py17X genome matching with the Py17XNL_2 genomic build (Figure 5A). We did not detect any large

263    structural variation between the strains, a finding also supported through our alternative Py17XNL_2

264    Nanopore/Illumina genome build. At the same time, we also identified a total of 1,955 potential single

265    nucleotide/short variants across the two strains, the majority (62%) of which were found in intergenic

266    regions (Figure 5B). We found that the apicoplast genome was identical between strains, whereas the

267    Py17XNL_2 mitochondrial genome has a 127 bp deletion in the middle of its sequence. Compared to

268    Py17X, the deletion is located in the intergenic region between *cox1* (PY17X_MIT00800) and a ribosomal

269    RNA fragment annotated as PY17X_MIT00700. Together, we conclude that while these two strains are

270    highly similar, there are sequence differences that may be functionally relevant.

271            To further determine the potential impacts of these genome variants, we characterized the

272    position of variants with respect to nearby genes and, when applicable, determined the specific DNA

273    and amino acid changes that would result from the change. Most variants were found to be located in

274    intergenic regions and were characterized as single base pair indels (Figure 6, Figure 7A,B). Of the 334

275    variants that fell within coding regions, we characterized the changes in nucleotide and amino acid

276    protein composition of the encoded proteins (representative examples are provided in Figure 6,

277    nucleotide and amino acid level changes are provided in Supplemental Table 4). Due to the over

278    representation of single bp indels, the majority of amino acid changes lead to frameshifts (Figure 7C).

279    Most of these frameshifts took place in genes that were unnamed with an unknown function, requiring

280    further investigation to determine the biological impacts of these differences.

281            To further interrogate those variants that occurred in well-characterized genes, we manually

282    curated the results and verified the variant calls via various quality measures. We checked if the variant

283    had sufficient PacBio ccs read support (80% of reads support the variant with a minimum of 5x coverage

284    at the region), and when possible, also determined if additional Nanopore/Illumina DNA and RNA

285    sequencing reads supported the variant (80% of reads support the variant with a minimum of 5x

286    coverage for Illumina sequencing and 2x coverage for Nanopore sequencing). Through manual curation,

287    a substantial number of variants had support from at least three sequencing methodologies. Due to the

288    strict thresholds of this variant calling process, some sequencing methods did not capture the variant

289    sufficiently enough to provide support, typically due to a lack of coverage at the position of the variant.

290    An example of this occurring is with CSP, which had a large deletion that was adequately supported by

291    PacBio and Nanopore DNA-seq data (Figure 6). Illumina DNA-seq reads, which should capture this

292    variant due to their high accuracy, instead have a complete loss of coverage, with only one read

293    correctly mapping to this repetitive region (Figure 6, Bottom Panel). As a result, we encourage the use of

294    long-read sequencing platforms to identify variants that may be missed when using Illumina sequencing.

295            Based upon these criteria, we identified if these genes were expressed in asexual/sexual blood

296    stages due to sequencing support from either Illumina or Nanopore direct RNA-seq and created

297    separate variant lists accordingly (Supplemental Table 4). Finally, we filtered out genes with no

298    annotated gene name and those that belong to a variable gene family (fam/pir gene families) (Figure 6).

299    After this filtering, we focused our analyses on the remaining 13 blood stage-expressed genes (Table 2).

300    Although the biological implications of the differences between Py17X and Py17XNL will need further

301    experimental validation, many variants could have interesting impacts. One such example is *ap2-sp*,

302    which has both synonymous and missense variants between the AT-hook and AP2 domain (43, 44). AP2-

303    SP is an ApiAP2 transcription factor with many target genes that are expressed specifically in the

304    sporozoite stage of the *Plasmodium* life cycle (43, 45, 46). It has also been shown that disruption of this

305    gene results in the loss of sporozoite formation entirely in the related *P. berghei* parasite and has

306    important activities in blood stages in *P. falciparum* (43, 45-47). Another affected gene is *pk4*, which

307    encodes an essential eIF2α kinase related enzyme and contains changes in its non-cytoplasmic domain

13

308     as determined by InterPro domain predictions across 17X and 17XNL strains (48-51). Further study of

309     these genes and several other candidates is warranted to understand the biological role these variants

310     may play across strains.

311         Among the non-blood stage expressed genes, *trap, lisp2,* and *csp* all had variants in their coding

312     sequences when comparing 17XNL to 17X. Of these, the most notable one is a large in-frame deletion

313     within the repeat region of CSP, leading to the loss of six of the repeating units of D/PQGPGA in Py17XNL

314     (Supplemental Figure 6). Similarly, the YM strain of *P. yoelii* is even shorter and lacks an additional

315     repeating unit compared to 17XNL (Supplemental Figure 6) (15). In *P. berghei*, it was found that 25% of

316     this repeat region could be eliminated before impacting parasite development, which is approximately

317     the length reduction observed in 17XNL and YM as compared to 17X (52). Therefore, this may reflect a

318     minimum repeat length for CSP functions.

**Discussion**

319

320    Here we have created a high-quality genome assembly with experimentally validated gene

321    models for the commonly used 17XNL strain of the *P. yoelii* malaria parasite species. We envision this

322    will be an important resource to the malaria research community, as it provides a much-needed update

323    to the Py17XNL_1 reference genome, which was among the first to be completed in the early days of the

324    genomics era 20 years ago (10). By directly comparing the strengths and genome assemblies created

325    from either PacBio HiFi sequencing or a combination of Nanopore DNA-seq and Illumina DNA-seq, we

326    identified that while the hybrid Nanopore/Illumina approach yielded a robust genome assembly, the

327    PacBio HiFi-based assembly consisted of fewer misassembles and covered a greater fraction of the

328    genome. Therefore, we have chosen the PacBio-based genome assembly as our new working reference

329    genome for *P. yoelii* 17XNL strain, which we have designated as Py17XNL_2 within this study. Our

330    findings align with many recent studies conducted to improve the reference information on *Plasmodium*

331    species, which also utilized an exclusively PacBio-based approach (12, 14, 17, 53). We also deemed it

332    important to conduct both approaches to leverage the strengths of Nanopore sequencing, which permit

333    greater detection of large-scale structural variants in the genome as compared to approaches with

334    shorter read lengths (26). The strengths of Nanopore sequencing have also been leveraged by other

335    sequencing efforts, most notably the recent telomere-to-telomere sequencing effort of the human

336    genome that used ultra-long read approaches (54). During this study, advances in Nanopore basecalling

337    software were made that enabled more accurate sequencing without the need for re-sequencing or new

338    hardware. We therefore directly compared the previous "fast" vs. new "high accuracy" basecalling

339    algorithms and observed a substantial increase in Qscores associated with the same DNA and RNA

340    sequencing data (Supplemental Table 2). However, even with the use of high accuracy basecalling,

341    PacBio data still enabled the most accurate Py17XNL genome assembly and covered the greatest

342    fraction of the genome.

15

343     To provide an even more useful genome reference, here we also provide new gene annotations

344     for the 17XNL strain of *P. yoelii* to facilitate more reliable forward and reverse genetic studies of this key

345     model malaria species. Regardless of the rodent malaria RNA-seq studies that have been performed,

346     gene models available on PlasmoDB for *P. berghei* and *P. yoelii* only consist of their putative coding

347     sequences. Here we have now added experimentally validated information on alternatively spliced

348     transcripts and untranslated regions (UTRs) of Py17XNL blood stage-expressed genes. To date, the only

349     other comparable efforts in our field have been applied to *P. falciparum* with a focus on either

350     identifying alternatively spliced transcripts or experimentally defining and annotating long noncoding

351     RNAs (lncRNAs) (55-57). Additionally, because Nanopore direct RNA-seq reads initiate at the 3' end of

352     mRNAs and progress toward the 5' end, it is also strong-suited in providing information about the 3'

353     UTRs of a population of mRNAs. From this, we created both a Nanopore-only and a Nanopore/Illumina

354     hybrid gene model annotation that can both be useful to researchers depending on the questions they

355     are pursuing. We are therefore providing both gene model files as resources to our community. These

356     gene models include well-defined 3'UTRs for Py17XNL that are in agreement with the length distribution

357     of those described for *P. falciparum* (13). Due to the strengths of this approach, we anticipate that

358     Nanopore direct RNA sequencing will become a useful tool for future work on *Plasmodium* parasites,

359     especially as sequencing chemistry and basecalling algorithms improve.

360     With this greatly improved Py17XNL reference genome, we also were able to critically analyze

361     genomic variation across the 17X and 17XNL strains of *P. yoelii*. As it is currently common practice to use

362     Py17X as a proxy genome for Py17XNL for genomic studies, we thought it was important to begin

363     addressing whether biologically relevant differences were present that would impact such efforts. By

364     aligning the two genomes, we saw that there was an excellent linear agreement between them, with

365     most variation taking place in intergenic regions. In total, there were 1,955 variants across the entire

366     sequence, with 334 of those being in the coding sequence of genes. Most of these variants were single

16

367    bp indels that most likely accounted for the overrepresentation of frameshift variants in the respective

368    amino acid sequence. Upon further analysis of these variants, some interesting questions arose

369    regarding the biological implications that these changes could have. Specific examples of genes with

370    impactful variants include the ApiAP2 transcription factor AP2-SP and PK4, which are essential for

371    *Plasmodium* development and warrant follow-up studies (43, 45, 46, 48, 50, 51). Aside from these blood

372    stage-expressed genes, it is also important to note the large-scale differences between the 17X, 17XNL,

373    and YM strains of *P. yoelii* in the central repeat of CSP, which are 150, 114, and 108 amino acids long,

374    respectively (Supplemental Figure 6). The in-frame deletions result in truncations of entire six amino

375    acid repeating units of D/PQGPGA, with 17XNL having six fewer units and YM having seven fewer than

376    17X. In *P. berghei*, it was found that 25% of this repeat region could be eliminated before impacting

377    parasite development, which reflects the approximate reduction in repeat length in 17XNL and YM

378    strains as compared to 17X (52). We anticipate that this may reflect a minimum repeat length that is

379    applicable to both highly related species. This is indirectly corroborated by the absence of any reports

380    that have documented significant differences in sporozoite development, functions, or transmissibility

381    between the 17X and 17XNL strains. Additionally, this particular variant was identified in both Nanopore

382    and PacBio long-read DNA-sequencing datasets, with Illumina short-read sequencing lacking coverage at

383    this site to accurately identify this deletion (Figure 6). This highlights the utility of long-read sequencing

384    technologies to resolve highly repetitive genomes.

385            This *P. yoelii* 17XNL_2 reference genome and its more comprehensive gene annotations provide

386    a resource that we believe will be helpful to the rodent malaria research community. We stress that

387    while most genes are identical between the 17X and 17XNL strains, there is appreciable genomic

388    variance in some important genes that should be considered when conducting genomic studies.

389    Therefore, we conclude that for many efforts, 17X is a suitable genomic proxy for 17XNL, but caution

390    against it for genes where variance exists, such as *csp, trap, lisp2, ap2-sp, pk4*, and others. Given the

17

391 improvements for both the Py17XNL_2 genome assembly and gene models presented here, we would

392 instead encourage their adoption as the working reference genome and gene annotation source for

393 studies of *P. yoelii* 17XNL.

394
395 **Materials and Methods**
396
397 Animal Experiments Statement
398
399     All animal care strictly followed the Association for Assessment and Accreditation of Laboratory

400 Animal Care (AAALAC) guidelines and was approved by the Pennsylvania State University Institutional

401 Animal Care and Use Committee (IACUC# PRAMS201342678). All procedures involving vertebrate

402 animals were conducted in strict accordance with the recommendations in the Guide for Care and Use

403 of Laboratory Animals of the National Institutes of Health with approved Office for Laboratory Animal

404 Welfare (OLAW) assurance.

405

406 Experimental Animals

407     Six-to-eight-week-old female swiss webster mice from Envigo were used for all experiments in

408 this work.

409

410 Parasite Preparation and Isolation

411     Mice infected with wild-type Py17XNL Clone 1.1 parasites from BEI Resources until a parasitemia

412 between 1-3% was reached. Approximately 1 mL of blood was collected from each euthanized mouse,

413 which was then added to 5 mL of heparinized (200 U) 1X PBS to prevent coagulation. The infected blood

414 was spun and the serum was aspirated to isolate the red blood cells (RBCs). Cells were resuspended in

415 10 mL 1X PBS and then passed through a cellulose column (Sigma #C6288) to remove mouse leukocytes.

416 The RBCs were then lysed in 0.1% w/v saponin in 1X PBS for 5 minutes at room temperature, and

417 parasite pellets were subsequently washed in 10 mL 1X PBS.

18

418

419 ## gDNA Preparation

420     All gDNA samples used for Nanopore, Illumina, and PacBio sequencing were prepared using the

421 NEB Monarch HMW DNA Extraction Kit for Cells and Blood (NEB #T3050) using the manufacturer's

422 protocol for fresh blood with slight modifications as we have previously described (28). Briefly, the

423 saponin lysed parasite pellet was resuspended in 150 μL of Nuclei Prep Buffer containing RNase A. After

424 resuspension of the pellet, 150 μL of Nuclei Lysis Buffer containing Proteinase K was added and mixed by

425 inversion. The sample was then placed in a thermal mixer at 56°C with an agitation speed of 1500 rpm

426 for 10 minutes. Next, 75 μL of precipitation enhancer was added and mixed by inversion. Two DNA

427 capture beads were added to the tube, along with 275 μL of isopropanol. The sample was then mixed 30

428 times with manual, slow, end-over-end inversions to ensure the gDNA stuck to the capture beads. The

429 supernatant was removed, and the beads were washed twice with 500 μL of gDNA wash buffer.

430 Subsequently, 100 μL of elution buffer II was added and the sample was incubated for five minutes at

431 56°C in a thermal mixer with agitation at 300 rpm. The beads were added to a bead retainer in an

432 Eppendorf tube, and the sample was spun down for 30 seconds at 12,000 $xg$. All samples were stored at

433 4°C to minimize shearing from freeze-thaw cycles. Fresh gDNA samples were made for replicate 1 and

434 replicate 2 for Nanopore sequencing and the sole sample for PacBio. The same gDNA samples used for

435 Nanopore replicates 1 and 2 were used for Illumina DNA sequencing replicates 1 and 2. Sample

436 concentration and purity were assessed via Qubit and Nanodrop, respectively (Thermo Fisher Scientific®

437 Nanodrop® 2000 and Qubit® instruments with the Qubit dsDNA BR Assay Kit (Cat #Q32853)). Fragment

438 length was assessed using an Agilent Technologies® TapeStation® 4200 system with Genomic DNA

439 ScreenTapes (Cat #5067-5366 and 5067-5365).

440

441 ## RNA Preparation

19

442          RNA samples were prepared from two biological replicates for Nanopore direct RNA sequencing.

443          RNA samples were produced using the Qiagen RNeasy kit using the manufacturer's protocol with slight

444          modifications to improve yield (Cat # 74104). Briefly, 350 μL of Buffer RLT was added to resuspend the

445          parasite pellet. The sample was passed through a 20-gauge needle five times and put back into the same

446          microfuge tube. Next, 350 μL of 70% ethanol was then added and mixed by pipetting using wide-bore

447          pipettes. The sample was then added to the spin column and was centrifuged for 15 seconds at 8,000

448          *xg*. The column was washed twice with 500 μL RPE buffer and was again centrifuged for 15-60 seconds

449          at 8,000 *xg*. Residual ethanol was removed by a final spin at these parameters. RNA was eluted from the

450          column into a fresh microfuge tube with 30 μL of DEPC-treated water. The sample was incubated for 15

451          minutes at room temperature to improve recovery yield. The sample was then collected by

452          centrifugation for 1 minute at 8,000 *xg*. A second elution with 30 μL DEPC-treated water was performed

453          as above to improve yield. To eliminate contaminating DNA, a Dnase I digestion was performed with

454          slight modifications to the Sigma #AMPD1 technical bulletin. Briefly, 8 μL of the prepared RNA was

455          mixed with 1 μL 10X Reaction buffer and 1 μL Dnase I, Amplification Grade, 1 unit/μL (Cat # AMPD1-KT).

456          The sample was gently mixed and incubated at room temperature for 30 minutes. The digestion was

457          terminated by the addition of 1 μL stop solution, followed by heat inactivation of the Dnase I. RNA was

458          precipitated with ethanol by adding 0.1 volume of 3M sodium acetate pH5.5@RT, four volumes of

459          reagent grade 200 proof ethanol, and 0.5 μL 20mg/ml glycogen. The solution was allowed to precipitate

460          overnight at -80°C. The solution was then spun down at 4°C at 12,000 *xg* for 10 minutes. The

461          supernatant was aspirated and 1 mL of 70% ethanol was added to wash the pellet. The pellet was spun

462          down as above and the supernatant was aspirated. The pellet was then allowed to air dry with the tube

463          inverted on a Kimwipe for 10 minutes. Sample concentration and purity were assessed via Qubit and

464          Nanodrop, respectively (Thermo Fisher Scientific® Nanodrop® 2000 and Qubit® instruments with the

465          Qubit dsDNA BR Assay Kit (Cat #Q32853)). RNA integrity was tested using an Agilent 2100 Bioanalyzer.

20

466

467 <u>Nanopore Ligation-based DNA Sequencing</u>

468 DNA sequencing of Nanopore replicates 1 and 2 was performed using the SQK-LSK110 Ligation

469 sequencing kit using the manufacturer's protocol. Genomic DNA (~1 µg as measured by Qubit) was

470 sequenced on an R9.4.1 (Cat # FLO-MIN106D) flow cell for 24 hours, washing between samples as per

471 manufacturer's recommendations (EXP-WSH003).

472

473 <u>Nanopore Direct RNA Sequencing</u>

474 RNA sequencing of Nanopore replicates 1 and 2 was performed using the SQK-RNA0002 Direct

475 RNA-sequencing kit from Oxford Nanopore Technologies using 500ng RNA. All sequencing was

476 performed on an R9.4.1 flow cell for 24 hours, washing between samples as per manufacturer's

477 recommendations (EXP-WSH004).

478

479 <u>Illumina DNA Sequencing</u>

480 Illumina DNA sequencing libraries were created using the Illumina DNA PCR-Free Kit with 100 ng

481 of total input (Cat # 20041794). Illumina libraries were sequenced on a NextSeq 550 Mid Output

482 150x150 paired-end sequencing run.

483

484 <u>PacBio Sequencing</u>

485 PacBio libraries were created using the PacBio SMRTbell Express Template Prep kit 2.0 (Cat #

486 TPK 2.0) using an input of 2 µg gDNA that was sheared with the Covaris g-TUBE to an average fragment

487 length of 10kb (Cat # 520079). The library was sequenced on a PacBio Sequel using a SMRT Cell 1M v3 LR

488 at a 10 pM library loading concentration with a 2-hour pre-extension time and a 20-hour movie time

489 (Cat # 101-531-000).

21

490

491     Data Analysis

492         High-quality HiFi reads were extracted from PacBio sequencing data requiring a minimum of 3

493     full passes in CCS command (v6.0.0) (31562484). The HiFi reads were *de novo* assembled using HiCanu,

494     specifying a genome size of 23 Mb. The resulting 132 contigs were aligned to the Py17X reference

495     genome using minimap2 (35), and all small contigs that had <2% alignment with a 17X chromosome

496     were filtered out. This reduced the contig number to 30 with a 22.2 Mb genome size. Chromosome

497     names were assigned based on alignment to Py17X. A consensus genome was generated based on the

498     alignment of these 30 contigs to the 17X reference genome using a custom program (Supp File 1). This

499     resulted in a final assembly of 16 contigs totaling 23.08 Mb. The apicoplast was circularized using

500     Circlator (v.1.5.5) (58), followed by manual correction of coordinates based on Py17X alignment.

501         The second assembly approach utilized both Nanopore and Illumina DNA-seq reads. All

502     Nanopore-based raw reads were first analyzed using Nanoplot (v1.33.0) (29) for quality control

503     purposes. Nanopore reads were assembled using Flye (v2.9) (30) software which resulted in 26 contigs.

504     The assembled contigs were scaffolded with Ragtag.py using the Py17X genome as a reference, which

505     resulted in 17 contigs (31). The resulting assembly was polished in a multi-step approach. First, the

506     assembly was polished using nextPolish (v1.3.1) (32). The homopolymer and long indel errors that were

507     still present were corrected in the second step. For these, 150 base pair reads were simulated from the

508     Py17X genome and mapped to the polished assembly in step 1 using bwa mem (59). Variants were

509     called with freebayes (v 1.2.0) (60), and a consensus was created with bcftools (v 1.15) (61), resulting in

510     a second round of polished assembly. To further correct errors, we mapped the Py17XNL Illumina

511     genomic DNA to the resulting assembly, called variants, and generated a consensus. This resulted in a

512     third-round polished assembly. Overlapping contigs were merged, and the apicoplast sequences were

513     circularized, resulting in a genome assembly consisting of 14 nuclear chromosomes and two organellar

514     chromosomes. The low complexity regions and tandem repeats in both assemblies were soft masked

515     using the tantan program (62). Assembly reports were done to compare Nanopore-based and PacBio-

516     based genomes using the Quast program. The variants between Py17X and Py17XNL genomes were

517     obtained from the minimap2 (v2.18) (35) whole genome alignment using paftools.js. The variants were

518     annotated, and variant effects were obtained using SnpEff (v.5.1d) (63). The assembly completeness was

519     assessed using BUSCO (42). For this assessment, the *Plasmodium* lineage database

520     (plasmodium_odb10), which contained 3642 sequences from 23 *Plasmodium* species, was searched to

521     check for the presence and completeness of the single-copy marker genes.

522          To create gene models and assign gene names, Braker2 (v2.1.6) (39) was first used to predict

523     genes, which was followed with the use of reciprocal blastp. Two sets of gene models were generated

524     for this assembly. For the first set, Nanopore dRNA-Seq reads were mapped to the assembled genome

525     with minimap2 (v2.18) (35). dRNA-Seq read alignments provided additional exon-intron evidence in

526     Braker2-based gene model predictions. Gene names were assigned by a reciprocal blast of the predicted

527     proteins against Py17X proteins. For the second set of gene-model predictions, both Nanopore dRNA-

528     Seq and Illumina RNA-Seq datasets were used. Illumina RNA-seq reads were mapped to the assembled

529     genome using Hisat2 (v.2.2.1) (64) and were merged with Nanopore dRNA-seq alignments, which were

530     then used for Braker2 gene-model prediction. Additionally, Prokka (v 1.14.6) was used to make gene

531     predictions in mitochondria and apicoplast. Finally, tRNAs were predicted using tRNASCAN-SE (v.2.0.9.)

532     (65), and rRNAs were identified by a blast search of the assembled genome using Py17X rRNAs.

533

534     Data Availability

535     Datasets associated with this study are available using the following identifiers: SRA BioProject:

536     PRJNA769959, Nanopore assembly accessions: CP086268-CP086283, PacBio assembly accessions:

23

537    CP115525-CP115540. All assembly files produced in this study are provided as Supplementary File 2, and

538    will be provided to VEuPathDB/PlasmoDB for integration and community use.

539

554

555    **Financial Disclosures**

556    We have no financial disclosures associated with this study.

557 **Figure Legends**

558

559 **Figure 1: PacBio HIFi high-quality long reads improve upon the pre-existing Py17XNL genome and**

560 **outperform a hybrid assembly approach with Nanopore and Illumina sequencing.** (A) QScore vs. read

561 length distribution for a PacBio sequencing run that was used to construct the final Py17XNL_2 genome

562 assembly is presented. Note: HiFi PacBio sequencing has a minimum QScore threshold of 20, and a

563 maximum QScore threshold of 93. (B) A histogram is plotted to illustrate the distribution of PacBio read

564 lengths. (C) A comparison of assembly statistics between Nanopore and PacBio sequencing runs is

565 provided. All statistics are based on contigs of size >=500 bp. (D) The cumulative length of contigs is

566 plotted from largest to smallest.

567

568 **Figure 2: Bioinformatics workflow used for genome assembly and annotation.** (Left) Genome

569 Assembly: High-accuracy ccs reads that were generated from PacBio subreads and trimmed Nanopore

570 reads were *de novo* assembled to create draft genomes. Contigs were selected, and chromosome names

571 were assigned based on the *P. yoelii* 17X reference genome alignment. Further processing of the

572 Nanopore + Illumina hybrid assembly involved implementing scaffolding and iterative polishing. (Right)

573 Gene-model prediction: A Nanopore dRNA-seq-based gene model and a hybrid gene model combining

574 both Nanopore dRNA-seq and Illumina RNA-seq data were generated using Braker2. The predicted

575 genes were annotated using reciprocal BLAST against *P. yoelii* 17X proteins. Illumina RNA-seq reads were

576 previously reported (37).

577

578 **Figure 3: Expanded *Plasmodium yoelii* 17XNL gene models leveraging RNA-seq data.** (A) An example

579 gene model depicting IMC1a and its respective sequence features is provided. (B) The 3'UTR length

580 distribution of all detected mRNAs is plotted as a histogram for chromosomal and mitochondrial genes.

581 Transcripts encoded by the apicoplast are not polyadenylated and were not detected by Nanopore

582 dRNA-seq. (C) The maximum, average, median, and mode of the 3' UTR lengths from all chromosomal

583 and mitochondrial transcripts are compared to those from a *Plasmodium falciparum* dataset (13).

584

585 **Figure 4: BUSCO analysis demonstrates genome assembly completeness.** Of the 3,642 BUSCO groups

586 that were searched, 3,556 single-copy BUSCOs were found to be present in the 17XNL_2 assembly

587 resulting in a completeness score of 97.6%. The BUSCO results for Py17XNL_1 (83.9%) and Py17X

588 (98.0%) reference genomes are also shown for comparison.

25

589

590    **Figure 5: Differences between the *P. yoelii* 17X and 17XNL_2 assemblies.** (A) The Py17XNL_2 reference

591    genome was mapped to Py17X to determine their degree of similarity. A dot plot depicting this

592    agreement is shown, with blue lines denoting unique alignments and orange lines depicting repeat

593    regions. (B) A circos plot is presented with the following tracks listed from outside to inside: 1) Py17X

594    reference genome, 2) Py17XNL_2 ccs read coverage in the natural log scale (minimum value of 0 and

595    maximum value of 8), 3) SNPs and indels between the two genomes are shown in light green, 4) SNPs

596    and indels in the coding sequence of genes are shown in orange. An expanded view that includes the

597    apicoplast and mitochondria is shown separately.

598

599    **Figure 6: Identification of blood stage-expressed variants between 17X and 17XNL_2.** (Left) Variants of

600    interest that are expressed in blood-stage parasites were chosen based on the presence of the variant

601    sequence within the coding sequence, the extent to which the variant calls are supported by sequencing

602    data, and if the gene has been named. Downselected genes are further described in Supplemental Table

603    4. To be considered, at least two sequencing methods needed to support the variant call, with at least

604    80% of the reads in agreement and a minimum of five reads at the position (three read minimum for

605    Nanopore). (Right) IGV snapshots with representative examples of different variants found in AP2-SP

606    (PY17XNL_1303202), RAD50 (PY17XNL_0104722), or CSP (PY17XNL_0404050) are presented top to

607    bottom.

608

609    **Figure 7: The location and potential impact on translation of variants between 17X and 17XNL_2**

610    **genome assemblies.** (A) The distribution of variant locations throughout the entire Py17XNL_2 genome

611    is shown. (B) The types of variants represented within the Py17XNL_2 genome with their respective

612    counts are plotted. (C) The distribution of variant types within coding sequences is depicted as a bar

613    graph.

614

615    **Table 1: Summary of finalized genome assembly and gene model creation statistics.** * Determined

616    using the Quast program by alignment to the 17X reference for this study. ** Determined using a BUSCO

617    analysis for this study. *** As reported in Carlton et al Nature 2002. N/D: Not determined.

618

619    **Table 2: Prioritized list of coding sequence variants between the Py17X and Py17XNL genome.**

620

621   **Supplementary Figure 1: Experimental workflow for all sequencing runs performed.** Four different

622   sample/sequencing types were generated. For each, mice were infected with Py17XNL strain parasites

623   until parasitemia reached 1-3%, at which point blood was collected, passed through a cellulose column,

624   and saponin lysed prior to DNA or RNA recovery. For Illumina, PacBio, and Nanopore DNA samples, the

625   NEB Monarch High Molecular Weight Blood Kit was used. For Nanopore RNA samples, a Qiagen RNeasy

626   Kit with subsequent DNaseI treatment was used. For quality control purposes, a Nanodrop and Qubit

627   were used to assess each biological sample. Additionally, TapeStation and Bionalyzer were used for DNA

628   and RNA samples, respectively. The library preparation methods and sequencing devices used for each

629   sample are also indicated. The Illumina RNA-seq data utilized in this study was previously published by

630   our laboratories and was retrieved from the GEO depository (Accession #GSE136674) (37).

631

632   **Supplementary Figure 2: Determination of gDNA fragment length by TapeStation.** (A) The Qiagen

633   Blood Amp Kit or the NEB Monarch High Molecular Weight Blood Kit were used to prepare gDNA and

634   samples were run in parallel on an Agilent TapeStation 4150. High molecular weight gDNA from lane C1

635   was used for Nanopore replicate one. (B) High molecular weight gDNA used for Nanopore replicate two

636   was run separately on the same Agilent TapeStation 4150 instrument. The hazard symbol in lane B1

637   indicates the sample was run outside of the manufacturer's recommended concentration. (C) High

638   molecular weight gDNA that was used for PacBio HiFi sequencing is shown in lane C2 on the right. All

639   other lanes were samples from unrelated experiments.

640

641   **Supplementary Figure 3: A comparison of fast basecalling and high accuracy basecalling for Nanopore**

642   **DNA sequencing.** (A and B) Nanopore ligation sequencing reads from replicate one were basecalled

643   using the fast basecalling algorithm (A) or the high accuracy basecalling algorithm (B). The Qscore vs.

644   read length distribution is depicted as a scatter plot (top), and the read length and their respective

645   counts are plotted as a histogram (bottom). (C and D) The same comparisons as described in A and B

646   were applied to replicate two.

647

648   **Supplementary Figure 4: Bioanalyzer results demonstrate that RNA samples are of high quality.** (A)

649   Total RNA isolated for replicate one of Nanopore direct RNA sequencing was run on a Bioanalyzer for

650   quality control purposes. The yellow hazard sign in lane B1 indicates that the markers ran outside of

651   their standard position, leading to an edited RIN. The sample was also run at a 1:5 dilution. (B) The RNA

652   sample used for replicate 2 of Nanopore direct RNA sequencing was run separately in the same way.

27

653

**Supplementary Figure 5: A comparison of fast basecalling and high accuracy basecalling for Nanopore direct RNA sequencing.** Nanopore direct RNA sequencing reads from replicate one (A and B) or two (C and D) were basecalled using the fast basecalling algorithm or the high accuracy basecalling algorithm. The Qscore vs. read length distribution is depicted as a scatter plot (top), and the read length and their respective counts are plotted as a histogram (bottom).

659

**Supplementary Figure 6: The CSP central repeat region length varies across sequenced *P. yoelii* strains.** The amino acid sequence for the central repeat region of circumsporozoite protein (CSP) is shown for *P. yoelii* 17X, *P. yoelii* 17XNL, and *P. yoelii* YM.

663

**Supplementary Table 1: Quality control measurements from Nanodrop, Qubit, TapeStation, and Bioanalyzer.**

666

**Supplementary Table 2: Nanopore sequencing statistics for replicates 1 and 2 with the fast basecaller (LA) and the high accuracy (HA) basecaller.**

669

**Supplementary Table 3: 5' and 3' untranslated region (UTR) information by gene name using Nanopore and Illumina ("hybrid") or Nanopore-only approaches.**

672

**Supplementary Table 4: Complete list of identified CDS variants with variant sequence and sequencing support information.**

675

**Supplementary File 1: Makefile that details the bioinformatics workflow used in this study.**

677

**Supplementary File 2: All assembly files generated in this study, including genome fasta, transcript fasta, cds fasta, protein fasta, and GFF3 files.**

**References:**

1.    Organization WH. World Malaria Report. 2022.

2.    De Niz M, Heussler VT. Rodent malaria models: insights into human disease and parasite biology. Curr Opin Microbiol. 2018;46:93-101.

3.    Voza T, Miller JL, Kappe SH, Sinnis P. Extrahepatic exoerythrocytic forms of rodent malaria parasites at the site of inoculation: clearance after immunization, susceptibility to primaquine, and contribution to blood-stage infection. Infect Immun. 2012;80(6):2158-64.

4.    Lindner SE, Swearingen KE, Shears MJ, Walker MP, Vrana EN, Hart KJ, et al. Transcriptomics and proteomics reveal two waves of translational repression during the maturation of malaria parasite sporozoites. Nat Commun. 2019;10(1):4964.

5.    Swearingen KE, Lindner SE. Plasmodium Parasites Viewed through Proteomics. Trends Parasitol. 2018;34(11):945-60.

6.    Munoz EE, Hart KJ, Walker MP, Kennedy MF, Shipley MM, Lindner SE. ALBA4 modulates its stage-specific interactions and specific mRNA fates during Plasmodium yoelii growth and transmission. Mol Microbiol. 2017;106(2):266-84.

7.    Li J, Cai B, Qi Y, Zhao W, Liu J, Xu R, et al. UTR introns, antisense RNA and differentially spliced transcripts between Plasmodium yoelii subspecies. Malar J. 2016;15:30.

8.    Lindner SE, Mikolajczak SA, Vaughan AM, Moon W, Joyce BR, Sullivan WJ, Jr., et al. Perturbations of Plasmodium Puf2 expression and RNA-seq of Puf2-deficient sporozoites reveal a critical role in maintaining RNA homeostasis and parasite transmissibility. Cell Microbiol. 2013;15(7):1266-83.

9.    Ogun SA, Tewari R, Otto TD, Howell SA, Knuepfer E, Cunningham DA, et al. Targeted disruption of py235ebp-1: invasion of erythrocytes by Plasmodium yoelii using an alternative Py235 erythrocyte binding protein. PLoS Pathog. 2011;7(2):e1001288.

10.    Carlton JM, Angiuoli SV, Suh BB, Kooij TW, Pertea M, Silva JC, et al. Genome sequence and comparative analysis of the model rodent malaria parasite Plasmodium yoelii yoelii. Nature. 2002;419(6906):512-9.

11.    Gardner MJ, Hall N, Fung E, White O, Berriman M, Hyman RW, et al. Genome sequence of the human malaria parasite Plasmodium falciparum. Nature. 2002;419(6906):498-511.

12.    Yang M, Shang X, Zhou Y, Wang C, Wei G, Tang J, et al. Full-Length Transcriptome Analysis of Plasmodium falciparum by Single-Molecule Long-Read Sequencing. Front Cell Infect Microbiol. 2021;11:631545.

13.    Chappell L, Ross P, Orchard L, Russell TJ, Otto TD, Berriman M, et al. Refining the transcriptome of the human malaria parasite Plasmodium falciparum using amplification-free RNA-seq. BMC Genomics. 2020;21(1):395.

14.    Fougere A, Jackson AP, Bechtsi DP, Braks JA, Annoura T, Fonager J, et al. Variant Exported Blood-Stage Proteins Encoded by Plasmodium Multigene Families Are Expressed in Liver Stages Where They Are Exported into the Parasitophorous Vacuole. PLoS Pathog. 2016;12(11):e1005917.

15.    Otto TD, Bohme U, Jackson AP, Hunt M, Franke-Fayard B, Hoeijmakers WA, et al. A comprehensive evaluation of rodent malaria parasite genomes and gene expression. BMC Biol. 2014;12:86.

722   16.    Hall N, Karras M, Raine JD, Carlton JM, Kooij TW, Berriman M, et al. A comprehensive
723   survey of the Plasmodium life cycle by genomic, transcriptomic, and proteomic analyses.
724   Science. 2005;307(5706):82-6.
725   17.    Zhang C, Oguz C, Huse S, Xia L, Wu J, Peng YC, et al. Genome sequence, transcriptome,
726   and annotation of rodent malaria parasite Plasmodium yoelii nigeriensis N67. BMC Genomics.
727   2021;22(1):303.
728   18.    Wang J, Chen K, Yang J, Zhang S, Li Y, Liu G, et al. Comparative genomic analysis of
729   Babesia duncani responsible for human babesiosis. BMC Biol. 2022;20(1):153.
730   19.    Menon VK, Okhuysen PC, Chappell CL, Mahmoud M, Mahmoud M, Meng Q, et al. Fully
731   resolved assembly of Cryptosporidium parvum. Gigascience. 2022;11.
732   20.    Baptista RP, Li Y, Sateriale A, Sanders MJ, Brooks KL, Tracey A, et al. Long-read assembly
733   and comparative evidence-based reanalysis of Cryptosporidium genome sequences reveal
734   expanded transporter repertoire and duplication of entire chromosome ends including
735   subtelomeric regions. Genome Res. 2022;32(1):203-13.
736   21.    Toro-Moreno M, Sylvester K, Srivastava T, Posfai D, Derbyshire ER. RNA-Seq Analysis
737   Illuminates the Early Stages of Plasmodium Liver Infection. mBio. 2020;11(1).
738   22.    Caldelari R, Dogga S, Schmid MW, Franke-Fayard B, Janse CJ, Soldati-Favre D, et al.
739   Transcriptome analysis of Plasmodium berghei during exo-erythrocytic development. Malar J.
740   2019;18(1):330.
741   23.    Kent RS, Modrzynska KK, Cameron R, Philip N, Billker O, Waters AP. Inducible
742   developmental reprogramming redefines commitment to sexual development in the malaria
743   parasite Plasmodium berghei. Nat Microbiol. 2018;3(11):1206-13.
744   24.    Yeoh LM, Goodman CD, Mollard V, McFadden GI, Ralph SA. Comparative
745   transcriptomics of female and male gametocytes in Plasmodium berghei and the evolution of
746   sex in alveolates. BMC Genomics. 2017;18(1):734.
747   25.    Mancio-Silva L, Slavic K, Grilo Ruivo MT, Grosso AR, Modrzynska KK, Vera IM, et al.
748   Nutrient sensing modulates malaria parasite virulence. Nature. 2017;547(7662):213-6.
749   26.    Wang Y, Zhao Y, Bollas A, Wang Y, Au KF. Nanopore sequencing technology,
750   bioinformatics and applications. Nat Biotechnol. 2021;39(11):1348-65.
751   27.    Amarasinghe SL, Su S, Dong X, Zappia L, Ritchie ME, Gouil Q. Opportunities and
752   challenges in long-read sequencing data analysis. Genome Biol. 2020;21(1):30.
753   28.    Godin MJL, S.E. NEB Monarch® HMW DNA Extraction Kit improves sample preparation
754   for Oxford Nanopore Technologies sequencing of malaria parasites2021. Available from:
755   https://www.neb.com/-/media/nebus/files/application-
756   notes/appnote_monarch_hmw_dna_improves_sample_prep_for-
757   ont_sequencing_of_malaria.pdf?rev=9c525dc516834ba684bf168dc13a164b
758   29.    De Coster W, D'Hert S, Schultz DT, Cruts M, Van Broeckhoven C. NanoPack: visualizing
759   and processing long-read sequencing data. Bioinformatics. 2018;34(15):2666-9.
760   30.    Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using
761   repeat graphs. Nat Biotechnol. 2019;37(5):540-6.
762   31.    Alonge M, Lebeigle L, Kirsche M, Aganezov S, Wang X, Lippman ZB, et al. Automated
763   assembly scaffolding elevates a new tomato system for high-throughput genome editing.
764   bioRxiv. 2021:2021.11.18.469135.

765    32.    Hu J, Fan J, Sun Z, Liu S. NextPolish: a fast and efficient genome polishing tool for long-
766    read assembly. Bioinformatics. 2020;36(7):2253-5.
767    33.    Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome
768    assemblies. Bioinformatics. 2013;29(8):1072-5.
769    34.    Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and
770    accurate long-read assembly via adaptive k-mer weighting and repeat separation. Genome Res.
771    2017;27(5):722-36.
772    35.    Li H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics.
773    2018;34(18):3094-100.
774    36.    Aurrecoechea C, Brestelli J, Brunk BP, Dommer J, Fischer S, Gajria B, et al. PlasmoDB: a
775    functional genomic database for malaria parasites. Nucleic Acids Res. 2009;37(Database
776    issue):D539-43.
777    37.    Hart KJ, Power BJ, Rios KT, Sebastian A, Lindner SE. The Plasmodium NOT1-G paralogue
778    is an essential regulator of sexual stage maturation and parasite transmission. PLoS Biol.
779    2021;19(10):e3001434.
780    38.    Parker MT, Knop K, Sherwood AV, Schurch NJ, Mackinnon K, Gould PD, et al. Nanopore
781    direct RNA sequencing maps the complexity of Arabidopsis mRNA processing and m(6)A
782    modification. Elife. 2020;9.
783    39.    Bruna T, Hoff KJ, Lomsadze A, Stanke M, Borodovsky M. BRAKER2: automatic eukaryotic
784    genome annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database. NAR
785    Genom Bioinform. 2021;3(1):lqaa108.
786    40.    Gish W, States DJ. Identification of protein coding regions by database similarity search.
787    Nat Genet. 1993;3(3):266-72.
788    41.    Gebauer F, Preiss T, Hentze MW. From cis-regulatory elements to complex RNPs and
789    back. Cold Spring Harb Perspect Biol. 2012;4(7):a012245.
790    42.    Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing
791    genome assembly and annotation completeness with single-copy orthologs. Bioinformatics.
792    2015;31(19):3210-2.
793    43.    Yuda M, Iwanaga S, Shigenobu S, Kato T, Kaneko I. Transcription factor AP2-Sp and its
794    target genes in malarial sporozoites. Mol Microbiol. 2010;75(4):854-63.
795    44.    Balaji S, Babu MM, Iyer LM, Aravind L. Discovery of the principal specific transcription
796    factors of Apicomplexa and their implication for the evolution of the AP2-integrase DNA binding
797    domains. Nucleic Acids Res. 2005;33(13):3994-4006.
798    45.    Martins RM, Macpherson CR, Claes A, Scheidig-Benatar C, Sakamoto H, Yam XY, et al. An
799    ApiAP2 member regulates expression of clonally variant genes of the human malaria parasite
800    Plasmodium falciparum. Sci Rep. 2017;7(1):14042.
801    46.    De Silva EK, Gehrke AR, Olszewski K, Leon I, Chahal JS, Bulyk ML, et al. Specific DNA-
802    binding by apicomplexan AP2 transcription factors. Proc Natl Acad Sci U S A.
803    2008;105(24):8393-8.
804    47.    Russell TJ, De Silva EK, Crowley VM, Shaw-Saliba K, Dube N, Josling G, et al. Inhibitors of
805    ApiAP2 protein DNA binding exhibit multistage activity against Plasmodium parasites. PLoS
806    Pathog. 2022;18(10):e1010887.

807  48.     Zhang M, Mishra S, Sakthivel R, Rojas M, Ranjan R, Sullivan WJ, Jr., et al. PK4, a
808  eukaryotic initiation factor 2alpha(eIF2alpha) kinase, is essential for the development of the
809  erythrocytic cycle of Plasmodium. Proc Natl Acad Sci U S A. 2012;109(10):3956-61.
810  49.     McDowall J, Hunter S. InterPro protein classification. Methods Mol Biol. 2011;694:37-47.
811  50.     Zhang M, Fennell C, Ranford-Cartwright L, Sakthivel R, Gueirard P, Meister S, et al. The
812  Plasmodium eukaryotic initiation factor-2alpha kinase IK2 controls the latency of sporozoites in
813  the mosquito salivary glands. J Exp Med. 2010;207(7):1465-74.
814  51.     Mohrle JJ, Zhao Y, Wernli B, Franklin RM, Kappes B. Molecular cloning, characterization
815  and localization of PfPK4, an eIF-2alpha kinase-related enzyme from the malarial parasite
816  Plasmodium falciparum. Biochem J. 1997;328 ( Pt 2)(Pt 2):677-87.
817  52.     Balaban AE, Kanatani S, Mitra J, Gregory J, Vartak N, Sinnis-Bourozikas A, et al. The
818  repeat region of the circumsporozoite protein is an elastic linear spring with a functional role in
819  <em>Plasmodium</em> sporozoite motility. bioRxiv. 2021:2021.05.12.443759.
820  53.     Bryant JM, Baumgarten S, Lorthiois A, Scheidig-Benatar C, Claes A, Scherf A. De Novo
821  Genome Assembly of a Plasmodium falciparum NF54 Clone Using Single-Molecule Real-Time
822  Sequencing. Genome Announc. 2018;6(5).
823  54.     Nurk S, Koren S, Rhie A, Rautiainen M, Bzikadze AV, Mikheenko A, et al. The complete
824  sequence of a human genome. Science. 2022;376(6588):44-53.
825  55.     Hoshizaki J, Adjalley SH, Thathy V, Judge K, Berriman M, Reid AJ, et al. A manually
826  curated annotation characterises genomic features of P. falciparum lncRNAs. BMC Genomics.
827  2022;23(1):780.
828  56.     Shaw PJ, Kaewprommal P, Wongsombat C, Ngampiw C, Taechalertpaisarn T,
829  Kamchonwongpaisan S, et al. Transcriptomic complexity of the human malaria parasite
830  Plasmodium falciparum revealed by long-read sequencing. PLoS One. 2022;17(11):e0276956.
831  57.     Lee VV, Judd LM, Jex AR, Holt KE, Tonkin CJ, Ralph SA. Direct Nanopore Sequencing of
832  mRNA Reveals Landscape of Transcript Isoforms in Apicomplexan Parasites. mSystems.
833  2021;6(2).
834  58.     Hunt M, Silva ND, Otto TD, Parkhill J, Keane JA, Harris SR. Circlator: automated
835  circularization of genome assemblies using long sequencing reads. Genome Biol. 2015;16:294.
836  59.     Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
837  Bioinformatics. 2009;25(14):1754-60.
838  60.     Said Mohammed K, Kibinge N, Prins P, Agoti CN, Cotten M, Nokes DJ, et al. Evaluating
839  the performance of tools used to call minority variants from whole genome short-read data.
840  Wellcome Open Res. 2018;3:21.
841  61.     Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve years of
842  SAMtools and BCFtools. Gigascience. 2021;10(2).
843  62.     Frith MC. A new repeat-masking method enables specific detection of homologous
844  sequences. Nucleic Acids Res. 2011;39(4):e23.
845  63.     Cingolani P. Variant Annotation and Functional Prediction: SnpEff. Methods Mol Biol.
846  2022;2493:289-314.
847  64.     Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. Graph-based genome alignment and
848  genotyping with HISAT2 and HISAT-genotype. Nat Biotechnol. 2019;37(8):907-15.
849  65.     Lowe TM, Eddy SR. tRNAscan-SE: a program for improved detection of transfer RNA
850  genes in genomic sequence. Nucleic Acids Res. 1997;25(5):955-64.

**Godin and Sebastian *et al*. Figure 1**

A.

**PacBio QScore vs. Read Length**

B.

**PacBio Read Length Distribution**

C.

Aligned to "17X" | 23 083 521 bp | 16 fragment s | 21.74 % G+C

Worst Median Best

| Genome statistics | Nanopore | PacBio |
|---|---|---|
| Genome fraction (%) | 95.314 | 99.999 |
| Duplication ratio | 1.021 | 1 |
| Largest alignment | 2 834 378 | 3 033 452 |
| Total aligned length | 22 442 174 | 23 084 626 |
| NGA50 | 1 733 329 | 2 046 261 |
| LGA50 | 6 | 5 |
| Misassemblies | | |
| # misassemblies | 34 | 0 |
| Misassembled contigs length | 15 930 248 | 0 |
| Mismatches | | |
| # mismatches per 100 kbp | 28.32 | 1.48 |
| # indels per 100 kbp | 24.83 | 6.99 |
| # N's per 100 kbp | 3.55 | 0 |
| Statistics without reference | | |
| # contigs | 16 | 16 |
| Largest contig | 2 895 268 | 3 033 452 |
| Total length | 22 515 783 | 23 084 626 |
| Total length (>= 1000 bp) | 22 515 783 | 23 084 626 |
| Total length (>= 10000 bp) | 22 509 826 | 23 078 671 |
| Total length (>= 50000 bp) | 22 475 516 | 23 044 347 |

D.

17XNL Nanopore

17XNL PacBio

Reference

Contigs are ordered from largest (contig #1) to smallest.

**Figure 1: PacBio HIFi high-quality long reads improve upon the pre-existing Py17XNL genome and outperform a hybrid assembly approach with Nanopore and Illumina sequencing.** (A) QScore vs. read length distribution for a PacBio sequencing run that was used to construct the final Py17XNL_2 genome assembly is presented. Note: HiFi PacBio sequencing has a minimum QScore threshold of 20, and a maximum QScore threshold of 93. (B) A histogram is plotted to illustrate the distribution of PacBio read lengths. (C) A comparison of assembly statistics between Nanopore and PacBio sequencing runs is provided. All statistics are based on contigs of size >=500 bp. (D) The cumulative length of contigs is plotted from largest to smallest.
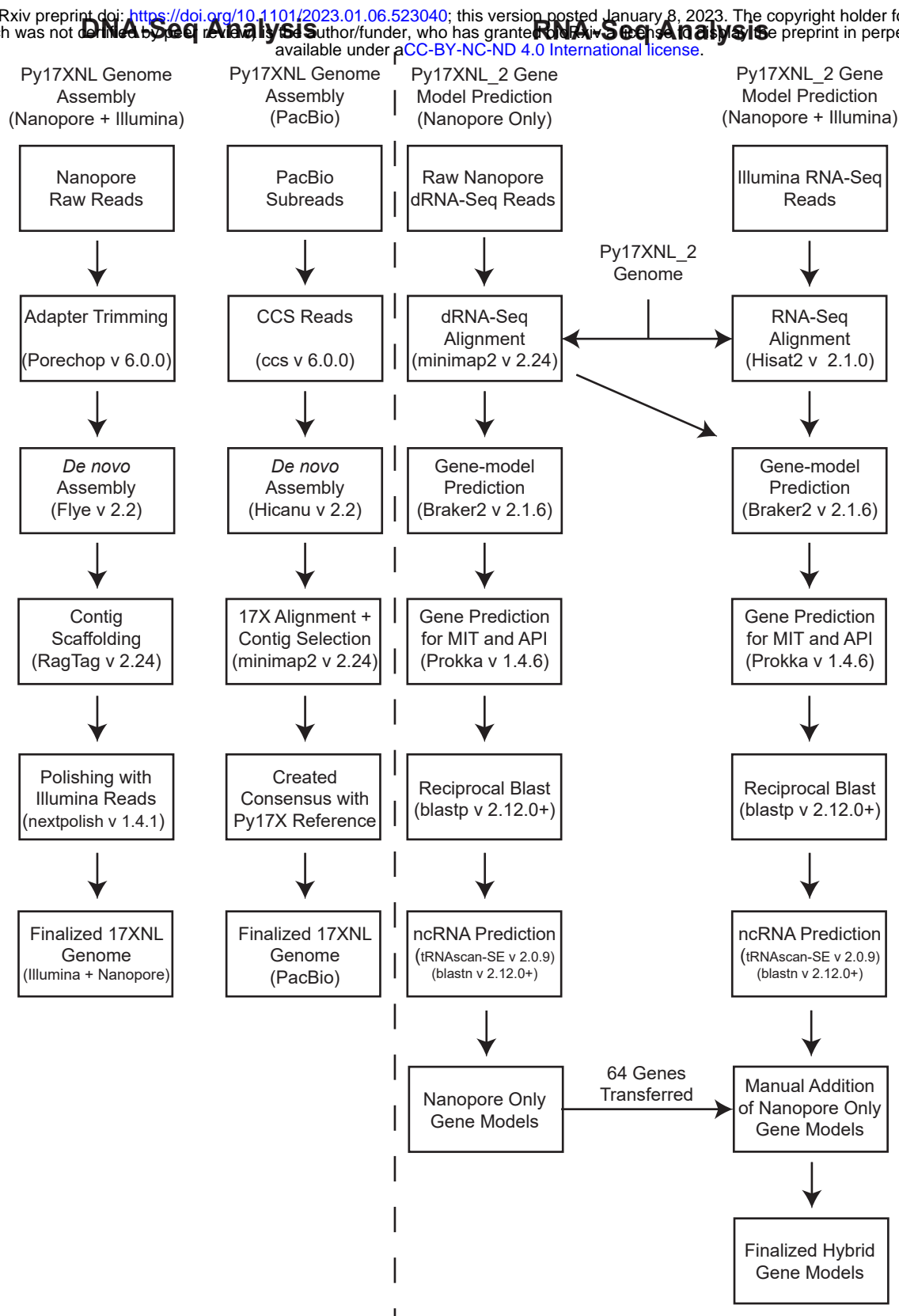
**Godin and Sebastian *et al*. Figure 2**

**Figure 2: Bioinformatics workflow used for genome assembly and annotation.**
(Left) Genome Assembly: High-accuracy ccs reads that were generated from PacBio subreads and trimmed Nanopore reads were de novo assembled to create draft genomes. Contigs were selected, and chromosome names were assigned based on the *P. yoelii* 17X reference genome alignment. Further processing of the Nanopore + Illumina hybrid assembly involved implementing scaffolding and iterative polishing. (Right) Gene-model prediction: A Nanopore dRNA-seq-based gene model and a hybrid gene model combining both Nanopore dRNA-seq and Illumina RNA-seq data were generated using Braker2. The predicted genes were annotated using reciprocal BLAST against *P. yoelii* 17X proteins. Illumina RNA-seq reads were previously reported (37).

# Godin and Sebastian *et al*. Figure 3

A.



B.



**3'UTR Distribution**

C.

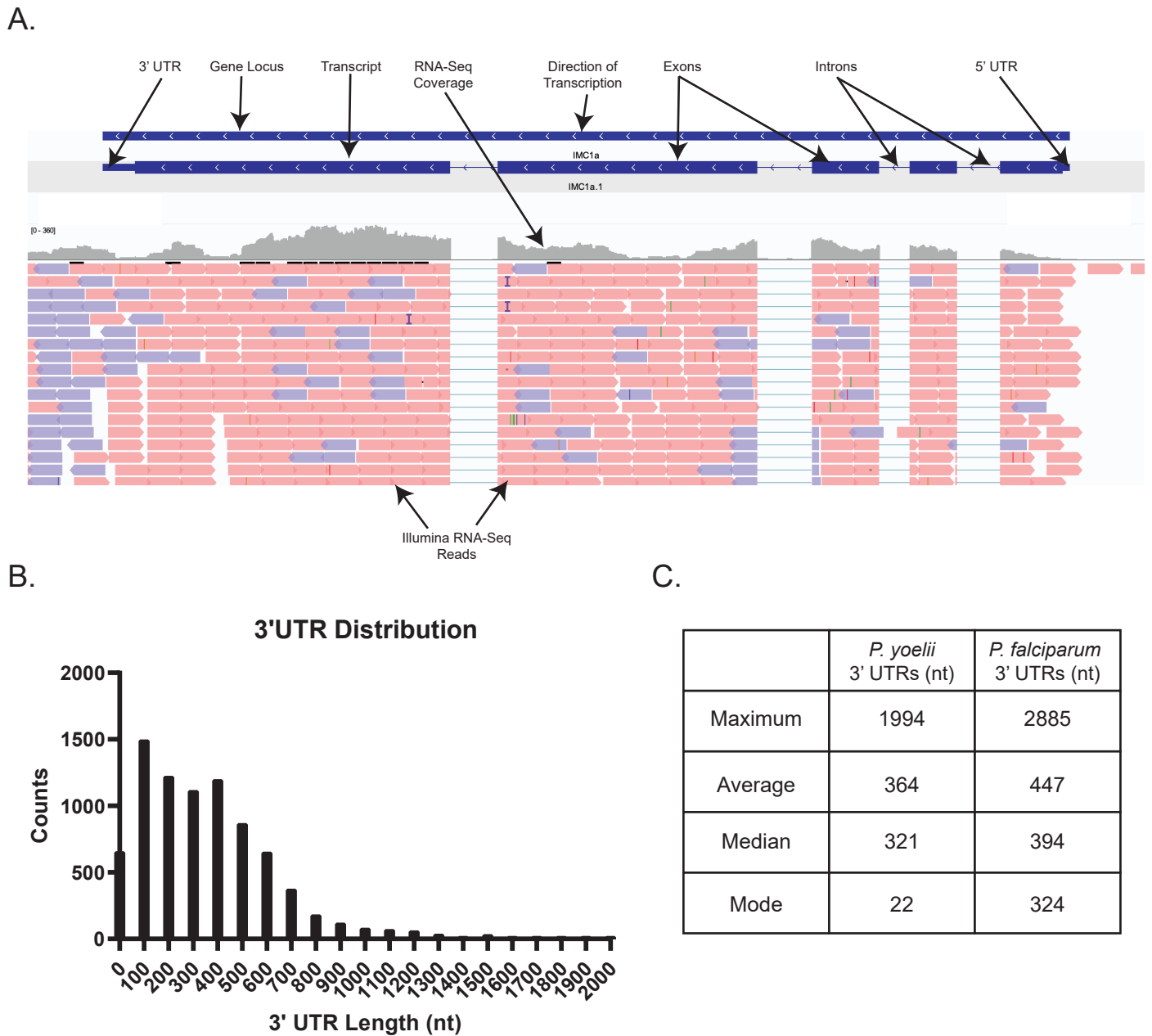|  | *P. yoelii*<br>3' UTRs (nt) | *P. falciparum*<br>3' UTRs (nt) |
|---|---|---|
| Maximum | 1994 | 2885 |
| Average | 364 | 447 |
| Median | 321 | 394 |
| Mode | 22 | 324 |

**Figure 3: Expanded *Plasmodium yoelii* 17XNL gene models leveraging RNA-seq data.**
(A) An example gene model depicting IMC1a and its respective sequence features is provided.
(B) The 3'UTR length distribution of all detected mRNAs is plotted as a histogram for chromosomal and mitochondrial genes. Transcripts encoded by the apicoplast are not polyadenylated and were not detected by Nanopore dRNA-seq. (C) The maximum, average, median, and mode of the 3' UTR lengths from all chromosomal and mitochondrial transcripts are compared to those from a *Plasmodium falciparum* dataset (13).
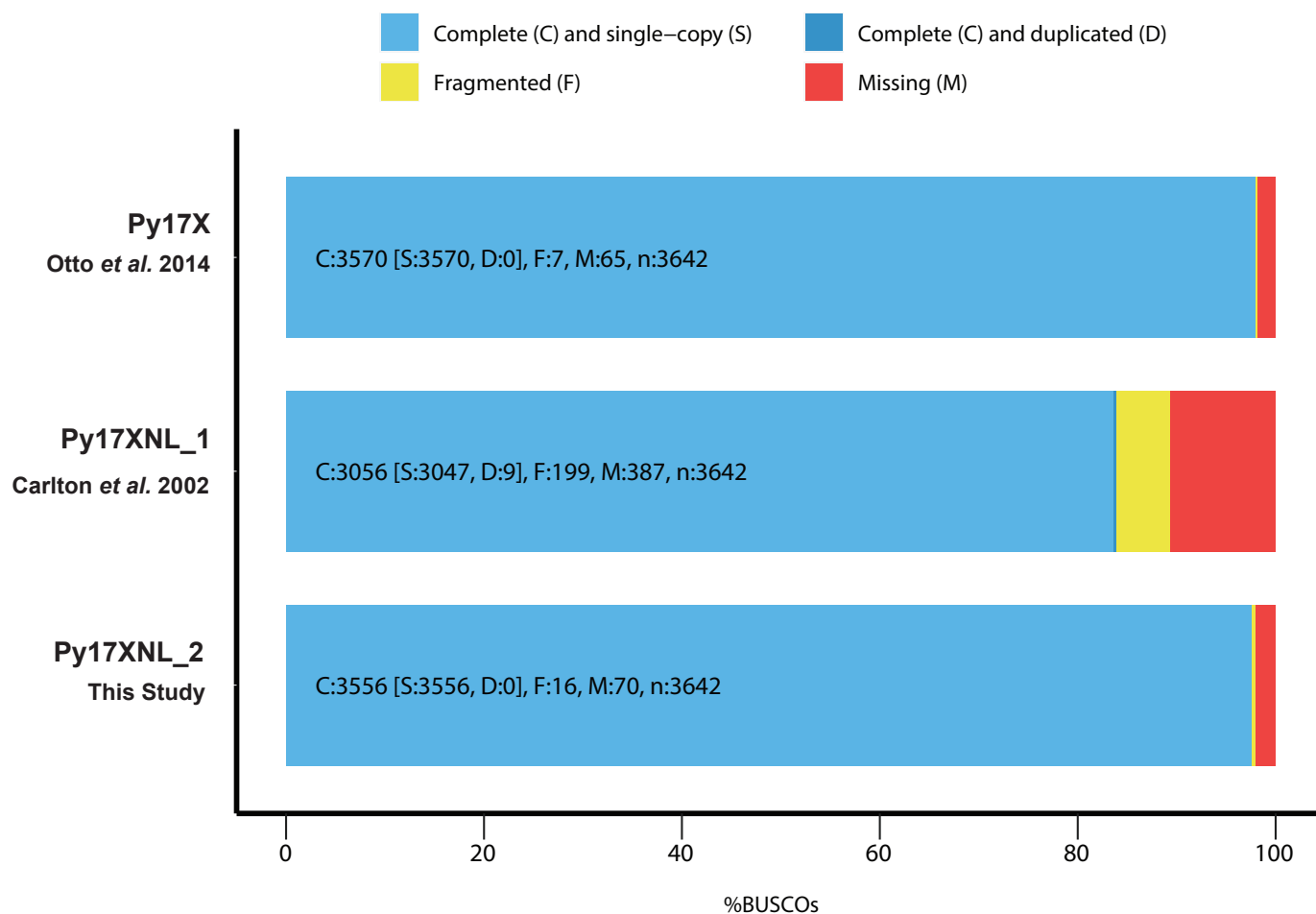
**Figure 4: BUSCO analysis demonstrates genome assembly completeness.**
Of the 3,642 BUSCO groups that were searched, 3,556 single-copy BUSCOs were found to be present in the 17XNL_2 assembly resulting in a completeness score of 97.6%. The BUSCO results for Py17XNL_1 (83.9%) and Py17X (98.0%) reference genomes are also shown for comparison.
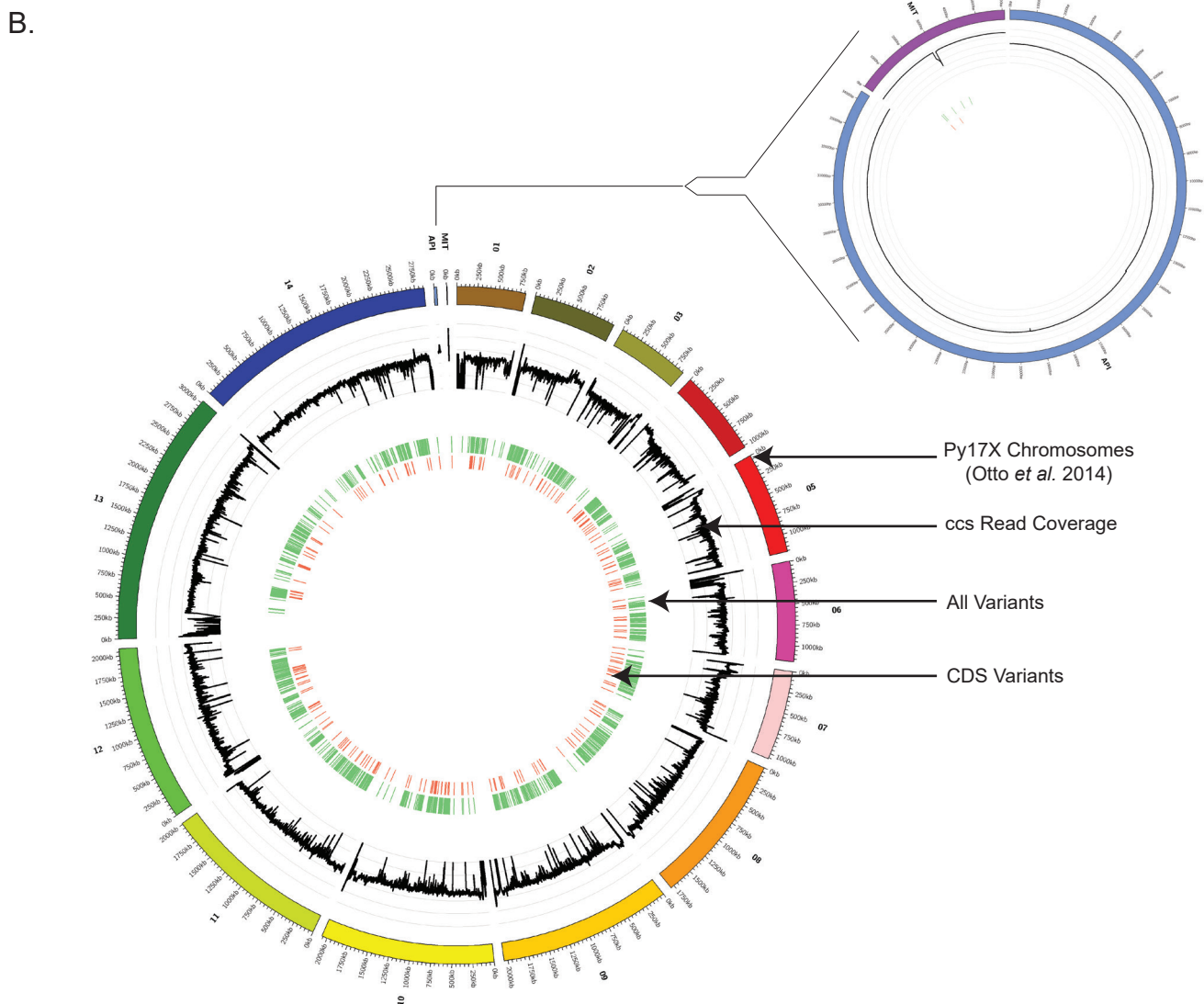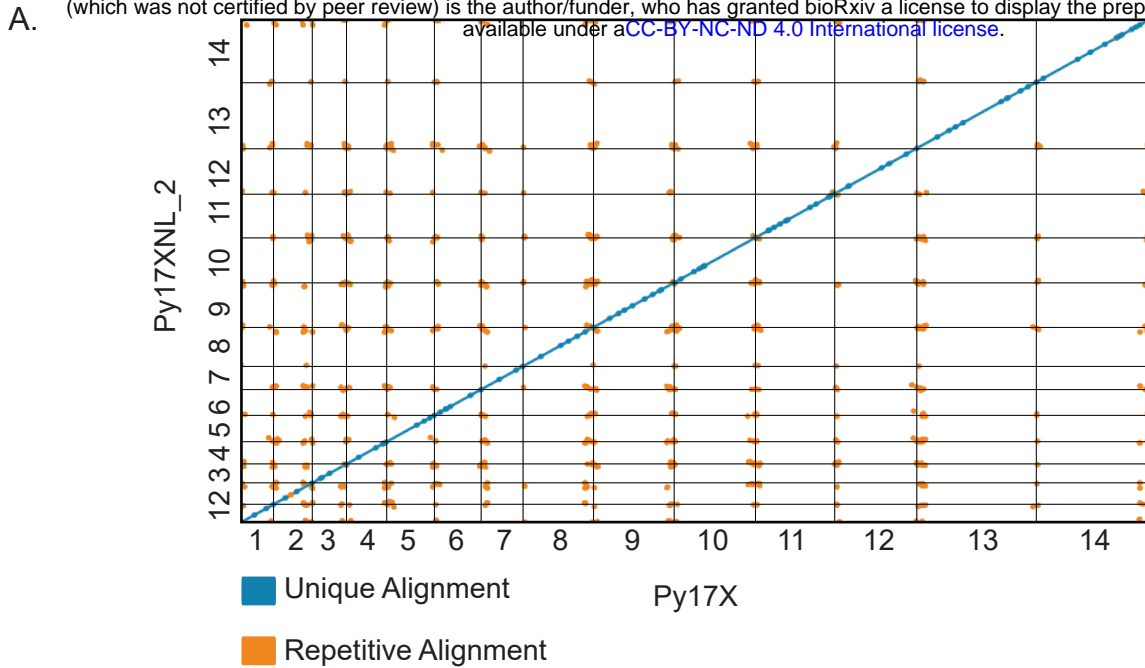
# Godin and Sebastian *et al*. Figure 5

**Figure 5: Differences between the *P. yoelii* 17X and 17XNL_2 assemblies.** (A) The Py17XNL_2 reference genome was mapped to Py17X to determine their degree of similarity. A dot plot depicting this agreement is shown, with blue lines denoting unique alignments and orange lines depicting repeat regions. (B) A circos plot is presented with the following tracks listed from outside to inside: 1) Py17X reference genome, 2) Py17XNL_2 ccs read coverage in the natural log scale (minimum value of 0 and maximum value of 8), 3) SNPs and indels between the two genomes are shown in light green, 4) SNPs and indels in the coding sequence of genes are shown in orange. An expanded view that includes the apicoplast and mitochondria is shown separately.

**Figure 6: Identification of blood stage-expressed variants between 17X and 17XNL_2.** (Left) Variants of interest that are expressed in blood-stage parasites were chosen based on the presence of the variant sequence within the coding sequence, the extent to which the variant calls are supported by sequencing data, and if the gene has been named. Downselected genes are further described in Supplemental Table 4. To be considered, at least two sequencing methods needed to support the variant call, with at least 80% of the reads in agreement and a minimum of five reads at the position (three read minimum for Nanopore). (Right) IGV snapshots with representative examples of different variants found in AP2-SP (PY17XNL_1303202), RAD50 (PY17XNL_0104722), or CSP (PY17XNL_0404050) are presented top to bottom.

# Godin and Sebastian *et al*. Figure 7

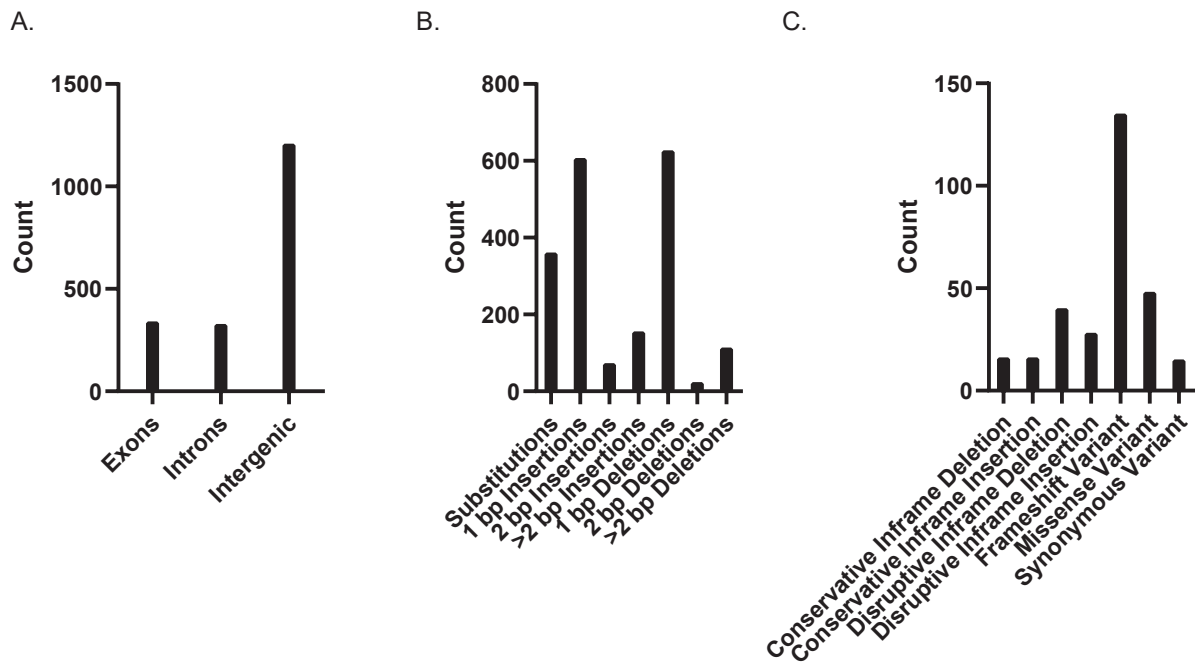**Figure 7: The location and potential impact on translation of variants between 17X and 17XNL_2 genome assemblies.** (A) The distribution of variant locations throughout the entire Py17XNL_2 genome is shown. (B) The types of variants represented within the Py17XNL_2 genome with their respective counts are plotted. (C) The distribution of variant types within coding sequences is depicted as a bar graph.

**Table 1.**

| | 17X (Otto *et al.* 2014) | 17XNL_1 (Carlton *et al.* 2002) | 17XNL_2 (This Study) |
|---|---|---|---|
| **Quast summary** | | | |
| Number of contigs | 16 | 5617 (5687***) | 16 |
| Assembled genome size | 23.08Mb | 23.1 Mb | 23.08 Mb |
| Genome Fraction (%) | - | 88.68% * | 99.90% |
| Largest Alignment | - | 51.5 Kb | 3,033,452 |
| NGA50 | - | 7,668 * | 2,046,261 |
| LGA50 | - | 851 * | 5 |
| Mismatches per 100 kb | - | 72.72 * | 1.48 |
| Indels per 100 kb | - | 44.07 * | 6.99 |
| N's per 100 kb | - | 50 * | 0 |
| | | | |
| **BUSCO summary** | | | |
| BUSCO genome completeness | 98% ** | 83.9% ** | 97.60% |
| BUSCO protein completeness | 99% ** | 64.7% ** | 90.00% |
| | | | |
| **Annotation summary** | | | |
| Number of genes | 6,263 | 7774 (5878***) | 6,086 |
| Number of mRNAs | 6,041 | 7,724 | 7,052 |
| Number of tRNAs | 79 | 50 (39***) | 66 |
| Number of rRNAs | 40 | 0 (7***) | 40 |
| | | | |
| **Variant summary** | | | |
| Variants with respect to 17X | - | N/D | 1,955 |
| Number of substitutions | - | N/D | 360 |
| Number of insertions | - | N/D | 833 |
| Number of deletions | - | N/D | 762 |
| Number of variants in CDS | - | N/D | 334 |

\* Determined using the Quast program for this study.

\** Determined using a BUSCO analysis for this study.

\*** As reported in Carlton et al Nature 2002

N/D Not Determined

**Table 2**

| Py17X Gene ID | Py17XNL Gene ID | Gene Name | Chromosome Position | Py17X DNA | Py17XNL DNA | Mutation Type | AA Change | Protein Alignment |
|---|---|---|---|---|---|---|---|---|
| PY17X_0811400 | Py17XNL_0801678 | proteasome subunit alpha type-3.1 | 529296 | G | A | Missense Variant | Val24Ile | Py17X MAGLSAGYDLSVSTFSPDGRLYQVEYIYKAINNNNTSISLECKDGVISCSINTSLEKNKMIKKNSYNRIYYV / 17XNL MAGLSAGYDLSVSTFSPDGRLYQIEYIYKAINNNNTSISLECKDGVISCSINTSLEKNKMIKKNSYNRIYYV / Cons ***********************:************************************************* |
| PY17X_0316900 | Py17XNL_0303734 | Plasmodium exported protein | 625346 | G | A | Missense Variant | Gly1324Ser | Py17X PEQKENGDIGEASNNAAELKENMNDLLKDTIEISKESIKEHDAQSIMFTRKFIKHVSGYDIQKAKDHPTDED / 17XNL PEQKENGDIGEASNNAAELKENMNDLLKDTIEISKESIKEHDAQSIMFTRKFIKHVSSYDIQKAKDHPTDED / Cons ********************************************************.*************** |
| PY17X_1022000 | Py17XNL_1002268 | PP7 | 964836 | C | A | Missense Variant | Leu342Phe | Py17X FAFKLSNYDSVIINRGNHECSYMNEVYGFHNEVLSKYDESVFDIFQEIFELLSLSVNIQNQIFVVHGGLSRY / 17XNL FAFKLSNYDSVIINRGNHECSYMNEVYGFHNEVLSKYDESVFDIFQEIFELLSFSVNIQNQIFVVHGGLSRY / Cons ****************************************************:****************** |
| PY17X_1419800 | Py17XNL_1401046 | ACDC domain-containing protein | 839701 | G | T | Missense Variant | Arg1562Ile | Py17X VNENFTAELNGVQMYNGNEKKKKKKNYSLSINKNNGNIKDNENTNEILLRYENEVYAPNNDVEKNLIEDNNI / 17XNL VNENFTAELNGVQMYNGNEKKKKKKNYSLSINKNNGNIKDNENTNEILLIYENEVYAPNNDVEKNLIEDNNI / Cons ************************************************:****************** |
| PY17X_0106400 | Py17XNL_0104725 | RNA-binding protein | 339741 | AGATAGG | A | Disruptive Inframe Deletion | Asp578_Arg579del | Py17X RDRDRDRDRDRDRDRDRDRDRR / 17XNL RDRDR--DRDRDRDRDRDRDRR / Cons *****  ************** |
| PY17X_1206500 | Py17XNL_1204935 | UTP25 | 344114 | GGAAAATGGGAAT | G | Disruptive Inframe Deletion | Gly163_Asn166del | Py17X ENGNENGNENGNENGNENGNENDKNGNDKNGNDKNGNDKNEASSFQSKDEIYMNILINNIKSQNEDFLNVKE / 17XNL ENGNENGNENGNENGNEN----DKNGNDKNGNDKNGNDKNEASSFQSKDEIYMNILINNIKSQNEDFLNVKE / Cons *****************    ********************** |
| PY17X_1110800 | Py17XNL_1105517 | KH domain-containing protein | 539935 | G | GTTA | Disruptive Inframe Insertion | Asn1809dup | Py17X NNNVGRDNIIRKENKGIMMHDDKDKFSKGGNNRYFGDKTNNFNNKN-NNNNNNNNNNNNNNNNNAKNNYLSKDSMI / 17XNL NNNVGRDNIIRKENKGIMMHDDKDKFSKGGNNRYFGDKTNNFNNKNNNNNNNNNNNNNNNNNNNAKNNYLSKDSMI / Cons **********************************************  ************************* |
| PY17X_1334500 | Py17XNL_1303202 | AP2-SP | 1594433 | AAC, T, AC | GCT, C, TA | Synonymous and Missense Variant | Val133Tyr, Val136Ser | Py17X QINYNISNDIMNTVPSTNCDVTHDSVSSVPNNAFENVENVKNVENVENVKNVENVENVENYENVENVENVENYEN / 17XNL QINYNISNDIMNTVPSTNCDVTHDSVSSVPNNAFENVENVKNVENVENVKNVENVENVENYENSENVENYEN / Cons *********************************************************** ** ******** |
| PY17X_1128400 | Py17XNL_1105678 | PK4 | 1187719 | T | TGAA | Conservative Inframe Insertion | Glu266dup | Py17X FYNSYNYCNNNNSKRDEKIEKNIVEKNIENKYNIKEYDKTNKSILFPIE-EFKIIQIENNIERNYIVPKES / 17XNL FYNSYNYCNNNNSKRDEKIEKNIVEKNIENKYNIKEYDKTNKSILFPIEEEFKIIQIENNIERNYIVPKES / Cons ************************************************ ****************** |
| PY17X_1451200 | Py17XNL_1401341 | BDP5 | 2018932 | A | AAATATAAAC | Disruptive Inframe Insertion | Asn320_Asn321insAsnIleAsn | Py17X NKIRSKNEINNSPNTDKVEKNIN---NINNINNINNNTNNNNVHEYVPNNLNDEFIEEKKLDKNKFNEYKNN / 17XNL NKIRSKNEINNSPNTDKVEKNINNINNINNINNINNNTNNNNVHEYVPNNLNDEFIEEKKLDKNKFNEYKNN / Cons *********************   ************************************ |
| PY17X_0942100 | Py17XNL_0900429 | PAIP1 | 1662686 | A | AAAT | Disruptive Inframe Insertion | Asn1941dup | Py17X NVNKNKEIGKDEIQINSQINNLDDNAKGKKSNIFNQAKSSYKYPAEEGENNSNTNTSTEN-NNNNNNNNNKT / 17XNL NVNKNKEIGKDEIQINSQINNLDDNAKGKKSNIFNQAKSSYKYPAEEGENNSNTNTSTENNNNNNNNNNNKT / Cons ***********************************************************  ********* |
| PY17X_0106100 | Py17XNL_0104722 | RAD50 | 324089 | CGTT | C | Disruptive Inframe Deletion | Gln834del | Py17X ENITNCVNKNEDILSDNLIKLESKKRVTAHFEELENGMKKKQRQEQDKFETVQKMKIEKIEKISKIEKINKI / 17XNL ENITNCVNKNEDILSDNLIKLESKKRVTAHFEELENGMKKK-RQEQDKFETVQKMKIEKIEKISKIEKINKI / Cons *****************************************  ***************************** |
| PY17X_1001900 | Py17XNL_1204888 | erythrocyte membrane antigen 1 | 174565 | TAAATGA | T | Conservative Inframe Deletion | Asn287_Glu288del | Py17X SYLNNGENAEDQELDDEVASCFADGENVNDKELDEVISYLANGENVNVNVNENVNENVNENVNENENENENE / 17XNL SYLNNGENAEDQELDDEVASCFADGENVNDKELDEVISYLANGENVNVNVNENVNENVNENVNE--NENENE / Cons ***************************************************************  ****** |