1    **BETHUNE et al. Long fragment capture and sequencing**

2    **Long-fragment targeted capture for long read sequencing of**

3    **plastomes**

4    Bethune Kevin[1*]; Mariac Cedric[1*]; Couderc Marie[1]; Scarcelli Nora[1]; Santoni Sylvain[2]; Ardisson Morgane[2];

5    Martin Jean-Francois[3]; Montufar Rommel[4]; Valentin Klein[1]; Francois Sabot[1]; Vigouroux Yves[1]; Couvreur

6    L.P. Thomas[1]

7    [1]IRD, DIADE, Univ Montpellier, Montpellier, France

8    [2]UMR AGAP, Equipe Diversitée et Adaptation de la Vigne et des Espèces Méditerranéennes, INRA, 2

9    Place Viala, 34060 Montpellier, France

10   [3]CBGP, Montpellier SupAgro, INRA, CIRAD, IRD, Univ Montpellier, Montpellier, France

11   [4] Facultad de Ciencias Exactas y Naturales, Pontificia Universidad Católica del Ecuador, Quito, Ecuador

12   * both authors contributed equally to the work

15   **ABSTRACT**

16   Third generation sequencing methods generate significantly longer reads than those produced using

17   alternative sequencing methods. This provides increased possibilities to better study biodiversity,

18   phylogeography and population genetics. We developed a protocol for in-solution enrichment

19   hybridization capture of long DNA fragments applicable to complete chloroplast genomes. The protocol

1

20    uses cost effective in-house probes developed via long-range PCR and was used in six non-model

21    monocot species (Poaceae: African rice, pearl millet, fonio; and three palm species). DNA was extracted

22    from fresh and silicagel dried leaves. Our protocol successfully captured long read chloroplast fragments

23    (up to 4 264 bp median) with an enrichment rate ranging from 15% to 98%. DNA extracted from silicagel

24    dried leaves led to low quality plastome assemblies when compared to freshly extracted DNA. Our

25    protocol could also be generalized to capture long sequences from specific nuclear fragments.

26

27    Keywords: MinION, DNA probes, long-range PCR, whole chloroplast sequencing, De Novo assembly

28    **INTRODUCTION**

29    High throughput sequencing is revolutionizing research in plant evolutionary biology. The development

30    of second generation sequencing (SGS) led to a massive amount of sequence data to be generated in a

31    cost effective way (Straub et al. 2012). Besides the many advantages of SGS one shortcoming is that they

32    generate short reads (between 100-400 base pairs (bp)). This is problematic for *de novo* assemblies of

33    plant genomes that prove difficult in resolving repetitive sequences due to transposable elements,

34    polyploidy and large genome sizes.

35        In contrast to SGS, third generation sequencing (TGS) directly targets single DNA molecules

36    without prior PCR, enabling "real time sequencing" (Bleidorn 2016). The main improvement of TGS is the

37    significant increase in read length from tens to tens of thousands of bases per single read (termed 'long

38    reads'). This provides important advantages to improve *de novo* assemblies (Jiao and Schneeberger

39    2017), gap filling (Eckert et al. 2016) or phasing (Laver et al. 2016). Technologies such as Pacific

40    Biosciences (PacBio) and Oxford Nanopore Technologies (ONT) are able to generate mean read lengths

41    ranging from 5kbp to 200kbp in standard analyses (and peak up to 2 mbp) depending on the quality of

2

the DNA (Lee et al. 2016). One drawback is that most TGS technologies have high error rates when compared to SGS (~ 10% for ONT MinION versus 0.1% for Illumina, Goodwin et al. 2016). However, new base calling algorithms, associated with *a posteriori* corrections, allow for a significant decrease sequence errors. With sufficient coverage and proper algorithms, TGS can lead to assemblies with consensus nucleotide accuracy of 99.90% (Lee et al. 2016).

The application, however, of TGS using MinION to complex genomes such as plants is problematic mainly because of the generally low output of data currently available (10-20 Gb versus 1,500 Gb for a HiSeq4000, Illumina). Thus, efficiently sequencing specific regions will depend on genome reduction approaches, such as targeted sequencing (Cronn et al. 2012; Jones and Good 2016). Genome reduction *via* sequence capture refers to DNA fragments (nuclear, ribosomal or plastid) that are directly captured from a total genomic library using probes binding to the complementary DNA sequences. This approach has the advantage of being cost effective, optimizes read depth on the targeted region and allows to analyze more samples per run. However, sequence capture is only routinely undertaken on short DNA fragments (Mamanova et al. 2010; Cronn et al. 2012), limiting its usefulness for long read based TGS.

Sequencing of complete chloroplasts or plastomes have been shown to be a marker of choice for the study plant evolution (Mariac et al. 2014; Twyford & Ness 2017). *De novo* assembly of plastomes based on short reads can be problematic (Mariac et al. 2014) leading to low quality reference plastomes. This is especially true for non-model taxa where no high quality reference genomes are available. Given the low output of data from MinION, this technology cannot be easily used to sequence plastomes directly from genomic DNA (e.g. genome skimming). The main challenge in order to efficiently apply TGS to the study of plant evolution will be based on our ability to capture long DNA fragments. To date long read targeted capture has mainly been undertaken on simple organism such as bacteria or virus (e.g.

3

PeerJ Preprints | https://doi.org/10.7287/peerj.preprints.27411v1 | CC BY 4.0 Open Access | rec: 5 Dec 2018, publ: 5 Dec 2018

65    Eckert et al. 2016) and rarely in complex organisms such as plants. Protocols for DNA enrichment for

66    segments in excess of 20kbp in length have also been developed (Dapprich et al. 2016). In plants, few

67    studies have undertaken long read targeted capture (Giolai et al. 2016, 2017). These protocols prove that

68    capturing long fragments is possible but has yet to be routinely developed for plants.

69         Here, we present a protocol to capture long reads for plastome sequencing and reassembling

70    using ONT MinION technology. We first developed our protocol for the model plant species *Oryza sativa*

71    (Asian rice). We then applied the protocol to sequence plastomes in several wild species and non-model

72    but economically important crops. Finally, we tested the ability to capture and assemble plastomes from

73    DNA extracted silicagel dried leaves.

74

75    **MATERIAL AND METHODS**

76    *Sampling strategy and DNA extraction*

77    For this study, we focused on seven economically important plant species from Asia, Africa and South

78    America. First, we developed and validated our long read capture protocol using the model plant species

79    *Oryza sativa* (Asian rice). We then applied our protocol to several other plant species from the same

80    genus (*Oryza*), family (*Poaceae*) and finally super-order (*Lilianae* or Monocotyledons): African rice (*Oryza*

81    *glaberrima* Steud.), a close relative to Asian rice, Pearl Millet (*Cenchrus americanus* (L.) Morrone

82    (*Pennisetum glaucum*)), Fonio (*Digitaria exilis* Stapf.), and three species of palms: *Podococcus acaulis*

83    Hua, *Raphia textilis* Welw. and *Phytelephas aequatoralis* Spruce (Table 1, Table S1). Export of *Podococcus*

84    *acaulis* and *Phytelephas aequatoralis* silicagel dried leaves were authorized by the Centre national de la

85    recherche scientifique (CENAREST, Gabon) and the Ministerio del Ambiente (Ecuador), respectively.

4

86    DNA was extracted from fresh leaves for *O. sativa*, *O. glaberrima*, *C. americanus* and *D. exilis*; while

87    silicagel dried leaves were used for DNA extraction for *Podococcus aucaulis*, *Raphia textilis* and

88    *Phytelephas aequatoralis*. In both cases DNA extraction was performed using a MATAB lysis buffer and

89    chloroform isoamyl alcohol (24:1) purification method following Mariac et al. (2006).

90    *General probe design*

91    Long fragment chloroplast sequences were captured from the total genomic DNA extracts using two

92    different sets of biotinylated probes: one based on *O. sativa* and used on related Poaceae species (*O.*

93    *glaberrima, C. americanus*, *D. exilis*) and one based on *P. barteri* Mann & H.Wendl. and used for *P.*

94    *aucaulis*, *R. textilis* and *P. aequatoralis*. Probe production (Figure 1, Supplementary file for detailed steps)

95    was undertaken following the protocol described elsewhere (Cronn et al. 2012; Mariac et al. 2014) and

96    lead to an average probe size of 300 bp: first, an initial full length chloroplast was amplified by long range

97    PCR (LRPCR), using 11 primer pairs taken from Scarcelli et al. (2011) for *O. sativa* (Table S2), and another

98    set of 11 primer pairs taken from Faye et al (2016) for *P. barteri* (Table S2). LR-PCR were carried out using

99    the LongAmp Taq PCR kit (New England BioLabs® Inc., #E5200S) following the manufacturer's instruction

100   in a final volume of 50 µL and using 300 ng of DNA. For each probe set LR-PCR amplicons were

101   equimolary pooled and sheared to reach a mean size fragment of 300 bp, then ligated to adapters so

102   that they can be PCR amplified with biotinyladed primers.

103   *Library preparation, in-solution hybridization, multiplexing and sequencing*

104   Illumina type libraries were constructed following the Rohland & Reich (2012) protocol using 6-bp

105   barcodes and Illumina indexes with some extra steps added to allow for amplification and in-solution

106   hybridization (Figure 1, Table S3). Briefly, each high molecular weight DNAs were sheared using G-tubes

107   (Covaris®) to a mean target size of 10 Kb. DNA fragments below 2,000 bp were removed by a sizing step

108   performed with 0.4X ampure beads. DNA was then end-repaired, ligated with adapters (allowing PCR

5

109    amplifications) and then nick filled-in before performing a pre-hybridation PCR. Optimal cycle number

110    (ranging from 5 to 12) was defined by real-time amplification (KAPA Biosystems, KK2700). After clean-up

111    and quantification using NanoQuant and QIAxcel, library preparations were mixed with biotin-labelled

112    probes for hybridization of the targeted regions. DNA-probe hybridization complexes were then

113    immobilized with 100 µg of streptavidine coated magnetic beads. This step was performed using the

114    Dynabeads™ M-280 kilobaseBINDER™ Kit (Invitrogen, ThermoFisher Scientific, #60101), which is

115    designed for immobilizing double stranded DNA molecules longer than 2 kbp.

116    A magnetic field was applied to the resulting solution and the supernatant containing unbounded DNA

117    was discarded. Enriched DNA fragments were then dehybridized from the beads and amplified in a 12 to

118    15 cycles real-time PCR in order to obtain requested quantity for the Nanopore library preparation. The

119    final libraries were then constructed following the Nanopore library preparation detailed in the 1D

120    Amplicon by ligation (SQK-LSK108) protocol for single samples and also in the 1D Native barcoding

121    genomic DNA (with EXP-NBD103 and SQK-LSK108) protocol. Briefly 1µg of enriched DNA was end-

122    repaired, extended with a dA-tailing, ligated with Nanopore barcodes and then with Nanopore tether-

123    adapter required previous to loading and sequencing on the MinION flowcell. To benefit from

124    multiplexing and limit costs and workload, up to four individuals were equimolarily pooled using Oxford

125    Nanopore barcodes. Prior to each run, flowcells (FLO-MIN106 R9.4 Version) were quality-tested using

126    the MinKNOW software version-1.2.8 to ensure the presence of at least 50% (256) of active channels.

127    Flow cells were loaded with around 275±100 fmol of capture-amplified DNA libraries.

128    *Non-enriched MiSeq data*

129    To estimate enrichment rate, we used single sample non-enriched library datasets originating from

130    various Illumina MiSeq sequencing runs for *O. sativa*, *O. glaberrima*, *C. americanus* and *P. aequatoralis*.

131    For *D. exilis*, *P. aucaulis*, and the *R. textilis*, we merged 10, 2 and 16 samples, respectively, of non-

6

132    enriched libraries to provide adequate read counts. Forward sequencing read outputs from each MiSeq

133    runs, namely R1 files, were first demultiplexed using demultadapt script

134    (https://github.com/Maillol/demultadapt) to sort reads according to a given barcodes list. Adapters at

135    the beginning of each read from the R2 and demultiplexed R1 files were removed using cutadapt-1.2.1

136    software (Martin 2011) with the default parameters. Reads were then filtered on their length (size >

137    35bp) and mean quality values (Q > 30) before being paired using compare_fastq_paired_v5.pl,

138    (https://github.com/SouthGreenPlatform/arcad-

139    hts/blob/master/scripts/arcad_hts_3_synchronized_paired_fastq.pl                        and

140    /arcad_hts_2_Filter_Fastq_On_Mean_Quality.pl). A last trimming step using the fastx-trimmer command

141    from the FASTX-Toolkit (http://hannonlab.cshl.edu/fastx_toolkit/) was undertaken onto the R2 paired

142    files to remove the last six bases of each read to ensure removal of any possible barcode present on

143    short reads.

144    *Bioinformatics*

145    All command lines are available in the Supplementary file.

146    Using the MinION Fast5 output format, base-calling and demultiplexing were undertaken using the

147    Albacore program v2.5.11 (https://github.com/Albacore/albacore). This generated a fastq file from

148    which reads were filtered out. The average quality score was lower than 7. For each barcode a quality

149    control using the MinionQC R script (https://github.com/roblanf/minion_qc) was done to check for read

150    mean length and quality scores. Reads were then trimmed using Porechop

151    (https://github.com/rrwick/Porechop) in order to remove the sequencing adapters and barcodes. The

152    only non-default setting is that splitting reads containing middle adapters was disabled, in order to avoid

153    issues during the polishing step using Nanopolish (see below).

154    For each library the percentage of chloroplastic reads was estimated by mapping reads to a reference

155    chloroplast genome using the Burrows-Wheeler alignment tool (bwa mem, https://github.com/lh3/bwa)

156    with "-B 1" option for non-enriched short reads data and "-x ont2d" option for long reads data (Li and

157    Durbin 2009). We then calculated the X-fold enrichment to evaluate capture efficiency (the ratio of

158    chloroplastic reads obtained with capture relative to chloroplastic reads obtained without capture).

159    Coverage and depth values were calculated using Bedtools (Quinlan and Hall 2010) genomecov

160    (https://github.com/arq5x/bedtools2). Mismatch percentage values between mapped reads and

161    references were recovered using Tablet v. 1.17.08.17 (Milne et al. 2010).

162    *De novo assembly of chloroplast genomes*

163    *De novo* assembly of plastomes based on the long MinION reads, we used the Flye assembler version 2.3

164    (Kolmogorov et al. 2018). For *O. sativa*, all available reads (17,129) were assembled. For the other

165    species, the number of reads was too high, in excess of 3000x of the reference coverage for some

166    datasets, which caused memory usage issues. To alleviate this, the reads were randomly split into sets of

167    approximately equal size. Each set was then assembled individually using the raw nanopore reads mode.

168    The "min_overlap" parameter (i.e. the minimum overlap between reads), in Flye was adjusted on a

169    species by species basis ranging from 3 000 pb (the default value for our genome size) to 1 000 bp,

170    depending on the medium read length for each species. This was done in order to ensure that a

171    sufficient amount of overlaps were detected for the assembly. The draft assemblies were then polished

172    using Nanopolish version 0.9.1 (https://github.com/jts/nanopolish), using minimap2 on the "map-ont"

173    preset for the overlapping step. Finally, the assemblies were mapped on the reference sequence of each

174    species using the dnadiff tool of MUMmer version 4.0beta2 (Kurtz et al. 2004), which directly provides

175    alignment coordinates and global statistics such as the mean identity percentage of alignments.

8

176    Besides read length, the uniformity of coverage of the reference by the reads could also have an impact

177    on the correct assembly of plastomes. This is especially problematic for low molecular weight DNA

178    extractions (in our case from silicagel dried leaves) which resulted in shorter read lengths on average

179    when compared to low molecular weight DNA extractions (in our case from fresh tissue). To test for the

180    impact of the uniformity of reference coverage on the assembly, simulated reads for *P. aequatoralis*

181    (DNA extracted from silicagel dried leaves) were generated using NanoSim v2.1.0 (Yang et al. 2017). A

182    model was first trained on the raw real reads and then 40'000 simulated reads were generated, ensuring

183    they have approximately the same length distribution and error model as the real reads (see results).

184    The simulated reads were then assembled using the same workflow as above.

185

186    **RESULTS**

187    *Plastome enrichment protocol validation on* Oryza sativa

188    After read filtering (Q>7) the median length of the 12 227 mapped plastome reads was of 4 264 bp

189    (Table 1, Figure 2A). We recovered the whole plastome with an average coverage depth of 364X for the

190    enriched MinION library, with a standard deviation increasing from 0.25 to 0.37 between enriched and

191    non-enriched libraries (Figure 2A, B, Table S1). The average mismatch was 11.80% (Table S1). Finally,

192    70.8% of the reads mapped to the reference chloroplast (Figure 2D, Table S1) representing a ~5-fold

193    increase in chloroplast reads when compared to the non-enriched MiSeq sequenced library (13.32%

194    mapped, Table S1). The longest plastome read recovered was 25 828 bp long (Table 1).

195    *Plastome enrichment protocol applied to non-model species*

196    DNA extraction qualities were variable deepening on the source of the leaf material used. Freshly

197    extracted DNA always produced single bands (not degraded) with fragments higher than 20 kb (Table

198    S1). Silicagel extracted DNA in contrast was of lower quality generally degraded (smear present) with

199    fragments under 20 kb long (Table S1). For the six non-model species, sequencing of the non-enriched

200    libraries resulted in 0.63% to 7.94% of chloroplast reads (Table S1). In contrast, enriched libraries

201    resulted in 15.7% to 98.2% chloroplast reads, corresponding to a 12 to 161-fold increase in chloroplast

202    DNA sequences (Figure 3, Table 1, Table S1). The mean average of fragments sequenced from freshly

203    extracted DNA was 4 279 bp versus 2 525 bp for DNA extracted from silicagel dry leaves (Figure 4A).

204    Sequences mapped to the reference plastomes ranged mainly from 2 kb and 8 kb, depending on the

205    species (Figure 4A). Average coverage depth was 1,988X for enriched libraries (Figure 4B, Table S1). The

206    longest read mapped to the plastome ranged from 10 405 bp to 25 167 bp for *R. textilis* and *C.*

207    *americanus*, respectively (Table 1).

208    *De novo assembly of chloroplast genome*

209    When DNA was extracted from fresh leaves, the chloroplast was assembled in two contigs covering most

210    of the reference (Table 1, Figure S1 for a visual exemple in *C. americanus*). Assembled contig lengths

211    varied from 81 053 to 12 5727 bp long. However, the assembler never managed to achieve the full

212    assembly and circularization of a chloroplast into a single contig. For DNA extracted from silicagel dried

213    leaves, where reads were shorter and the coverage more heterogeneous, assembly was suboptimal

214    (Table 1, Figure S2 for a visual exemple in *Phytelephas aequatorialis*) with more final contigs (10-17),

215    uncovered regions and sometimes misassemblies. The longest assembled contigs were also much short

216    than for fresh DNA (Table 2). In addition, the Inverted Repeats (IRs) were also often not differentiated.

217    Using a simulated dataset of reads uniformly distributed across the chloroplast (Figure S3) and based on

218    the same quality as *P. aequatorialis* significantly improved assembly (Table 2). The assembler resulted in

219    four contigs (versus 13) covering almost 99.72% (versus 87.60%) of the reference and the longest of

10

220    107 633 bp (versus 21 797 bp). However, the existence of two distinct repeated regions was still not

221    resolved.

222

223    **CONCLUSIONS**

224        Here, we show that targeted capture hybridization of long plastome DNA fragments with good

225    coverage (362 x to 3318 x) is possible in plants (Table 1, Table S1). In addition, we show a significant

226    enrichment of our target region (the plastome) when compared to non-enriched data (Figure 2D, 3,

227    Table S1). The different steps of our protocol (Figure 1, Supplementary file) are not fundamentally

228    different from previous chloroplast short read capture protocols (e.g. Mariac et al. 2014) based on in-

229    house probe preparation, shearing, adapter ligation, hybridization and finally capture (Figure 1,

230    Supplementary file). Thus, our approach requires minimal adaptation from previous cost and time

231    effective protocols and should therefore be of broad interest. The main technical change focused on the

232    beads used to capture long DNA fragments. For that we used the kilobaseBINDER$^{TM}$ Kit of Invitrogen

233    which is said to capture DNA fragments longer than 2 kbp. The sizing step we performed at 0.4X using

234    ampure removes fragments smaller than 2 kb and corresponds to the maximum allowed size with the

235    ampure beads. However, other approaches are possible to achieve sizing with higher molecular weight

236    and could be tested (e.g. gel extraction, Automated Size Selection System).

237        When capturing plastomes across a range of difference species, we find a difference in enrichment

238    percentage ranging from 15.7% to 98.2% of useful reads (Table 1, Figure 4). Differences in genome

239    versus plastome ratios between species can explain the variation of on-target mapped reads percentage

240    compared to non-enriched libraries. In general, species with smaller genomes show higher mapped read

241    percentages. Alternatively, the material used for DNA extractions, the cellular type and the degradation

11

242    state, can also explain such variations. The low enrichment observed for *P. acaulis* (Table 1, Figure 3)

243    could potentially be linked to a large genome size, although we do not have an estimate of its genome.

244        A common coverage gap is observed among the chloroplasts of the three palm species due to a

245    region that wasn't covered by the probes (see Faye et al. 2016). Coverage depressiveness of other

246    regions can be explained with biases that occur during DNA shearing, PCR amplification and hybridization

247    capture, considering a CG content effect. Probe bulk normalization from long-range PCR also has to be

248    taken into account. However, global decrease of standard deviation of enriched libraries proved a slight

249    increase of the whole target coverage homogeneity compared with non-enriched libraries. This means

250    that the hybridization capture performed in our protocol didn't introduce more on-target coverage

251    heterogeneity. Nevertheless, applying an alternative capture method such as region-specific extraction

252    (Dapprich et al. 2016) could help maintaining an overall good coverage by accessing high complex,

253    variable, repeat-masked or unknown regions that forbid adequate probe binding.

254        Probes were designed to hybridize across the whole targeted region (Figure 1), as is generally

255    done using short read approaches (Stull et al. 2013; Mariac et al. 2014). However, a recent study showed

256    that probes targeting small regions are also effective to capture long reads surrounding the targeted

257    region. Indeed, Gasc & Peyret (2017) were able to reconstruct a 21.6 kbp fragment using probes

258    designed for a small 471 bp microbial gene target. This shows that long read capture will also be very

259    useful for targeted sequence capture of nuclear regions.

260        We demonstrated the capacity of heterologous plastome probes to capture target DNA in other

261    species or genera in Arecacece and Poaceae. For example, probes designed on *P. barteri* hybridized well

262    to other palm genera in different sub families. This underlines the good portability of probes for

263    capturing plastomes across a broad evolutionary spectrum (Stull et al. 2013), even for long fragment

264    capture.

12

265 **Limits and challenges**

266 Although we were able to successfully capture long plastid fragments using our enrichment protocol,

267 assembling plastomes from this data remains challenging. Indeed, the best assembly resulted in 2

268 mapped contigs, and the worst one in 17 (Table 2). Assembly of plastomes is well known to be

269 problematic (Twyford & Ness 2017) mainly because of the presence of near identical inverted repeats

270 (IR). Indeed, the similarity of the two IRs is too high for assemblers to decipher between IRs when

271 resolving the assembly graph for the entry and exit point of those sequences. Thus, when the size of the

272 sequenced reads are shorter than the IRs themselves it becomes hard to the correctly assemble of the

273 plastome into a single contig. This is visible for exemple in *C. americanus* (Figure S1) where the resulting

274 two contigs do not across one of the IR regions leading to a failure in reaching a single contig. Of course,

275 this problem is enhanced when dealing with overall shorter reads sequenced from low molecular weight

276 DNA (see Figure S2 for an exemple). In our case, fragment length recovered of DNA extracted from

277 silicagel dried leaves was shorter than those extracted from fresh leaves (Table 1). Moreover, we

278 observed a decrease of the average library fragment size during preparation steps and mainly after PCR

279 because of preferential amplification of shorter fragments, as observed by Giolai et al. (2016) and Eckert

280 et al. (2016).

281 Optimizing read length in such a way that single reads are longer than the entire IR region should

282 significantly help in the assembly process. In this sense, DNA shearing could be removed in order to

283 increase the average size of the reads. Technical limitations would however be 1) the ability of

284 streptavidin beads to immobilize fragments of tens of thousands of base pairs and 2) the long range PCR

285 amplification step of the enriched fragments which is necessary to produce an input of several hundred

286 ng for the construction of nanopore libraries. The latter is probably the most limiting because it is

287 difficult to produce amplicons of several tens of Kb, and even if we achieve this, representation bias are

13

288    to be considered. Finally, we show, via simulations, that the uniformity of read coverage across the

289    reference are important for assembly (Table 2). Indeed, uniformly distributed reads, even of lower

290    quality, lead to better assemblages than poor coverage of the reference (Table 2). Therefore, uniform

291    coverage of the reference by the captured reads plays a big role in the correct and improved assembly

292    even for suboptimal DNA extractions.

293

## 294    ACKNOWLEDGEMENTS

301

## 302    AUTHOR CONTRIBUTIONS

303    CM, YV, TLPC conceived the idea; RM, TLPC, CM, YV provided material; CM, YV, JFM, SS, AM designed the

304    protocol; KB, CM, MC undertook the experiments; CM, KB, FS, VK analyzed the data. KB, TLPC led the

305    writing; all authors read and commented on the final version.

306

## 307    LITERATURE CITED

14

308    BLEIDORN C. 2016. Third generation sequencing: technology and its potential impact on evolutionary

309         biodiversity research. Syst. Biodivers. 14:1–8.

310    CRONN R., KNAUS B.J., LISTON A., MAUGHAN P.J., PARKS M., SYRING J.V., UDALL J. 2012. Targeted enrichment

311         strategies for next-generation plant biology. Am. J. Bot. 99:291–311.

312    DAPPRICH J., FERRIOLA D., MACKIEWICZ K., CLARK P.M., RAPPAPORT E., D'ARCY M., SASSON A., GAI X., SCHUG J.,

313         KAESTNER K.H., others. 2016. The next generation of target capture technologies-large DNA

314         fragment enrichment and sequencing determines regional genomic variation of high complexity.

315         BMC Genomics. 17:486.

316    ECKERT S.E., CHAN J.Z.-M., HOUNIET D., THE PATHSEEK CONSORTIUM, BREUER J., SPEIGHT G. 2016. Enrichment by

317         hybridisation of long DNA fragments for Nanopore sequencing. Microb. Genomics. 2.

318    FAYE A., DEBLAUWE V., MARIAC C., RICHARD D., SONKÉ B., VIGOUROUX Y., COUVREUR T.L.P. 2016. Phylogeography

319         of the genus Podococcus (Palmae/Arecaceae) in Central African rain forests: Climate stability

320         predicts unique genetic diversity. Mol. Phylogenet. Evol. 105:126–138.

321    GASC C., PEYRET P. 2017. Revealing large metagenomic regions through long DNA fragment hybridization

322         capture. Microbiome. 5:33.

323    GIOLAI M., PAAJANEN P., VERWEIJ W., PERCIVAL-ALWYN L., BAKER D., WITEK K., JUPE F., BRYAN G., HEIN I., JONES

324         J.D.G., CLARK M.D. 2016. Targeted capture and sequencing of gene-sized DNA molecules.

325         BioTechniques. 61:315–322.

326    GIOLAI M., PAAJANEN P., VERWEIJ W., WITEK K., JONES J.D.G., CLARK M.D. 2017. Comparative analysis of

327         targeted long read sequencing approaches for characterization of a plant's immune receptor

328         repertoire. BMC Genomics. 18:564.

15

329    GOODWIN, S., MCPHERSON, J.D., & MCCOMBIE, W.R. 2016. Coming of age: ten years of next-generation

330        sequencing technologies. Nature Reviews Genetics, 17, 333–351.

331    JIAO W.-B., SCHNEEBERGER K. 2017. The impact of third generation genomic technologies on plant genome

332        assembly. Curr. Opin. Plant Biol. 36:64–70.

333    KARAMITROS T., MAGIORKINIS G. 2015. A novel method for the multiplexed target enrichment of MinION

334        next generation sequencing libraries using PCR-generated baits. Nucleic Acids Res. 43:e152–

335        e152.

336    KOLMOGOROV M., YUAN J., LIN Y., PEVZNER P. 2018. Assembly of Long Error-Prone Reads Using Repeat

337        Graphs. bioRxiv.:247148.

338    KOREN S., WALENZ B.P., BERLIN K., MILLER J.R., BERGMAN N.H., PHILLIPPY A.M. 2017. Canu: scalable and

339        accurate long-read assembly via adaptive k-mer weighting and repeat separation. Genome Res.

340        27:722–736.

341    KURTZ S., PHILLIPPY A., DELCHER A.L., SMOOT M., SHUMWAY M., ANTONESCU C., SALZBERG S.L. 2004. Versatile and

342        open software for comparing large genomes. Genome Biol. 5:R12.

343    LEE H., GURTOWSKI J., YOO S., NATTESTAD M., MARCUS S., GOODWIN S., MCCOMBIE W.R., SCHATZ M. 2016. Third-

344        generation sequencing and the future of genomics. bioRxiv.:048603.

345    LI H., DURBIN R. 2009. Fast and accurate short read alignment with Burrows–Wheeler transform.

346        Bioinformatics. 25:1754–1760.

347    LI H., HANDSAKER B., WYSOKER A., FENNELL T., RUAN J., HOMER N., MARTH G., ABECASIS G., DURBIN R. 2009. The

348        Sequence Alignment/Map format and SAMtools. Bioinformatics. 25:2078–2079.

16

349     MAMANOVA L., COFFEY A.J., SCOTT C.E., KOZAREWA I., TURNER E.H., KUMAR A., HOWARD E., SHENDURE J., TURNER

350           D.J. 2010. Target-enrichment strategies for next-generation sequencing. Nat Meth. 7:111–118.

351     MARIAC C., SCARCELLI N., POUZADOU J., BARNAUD A., BILLOT C., FAYE A., KOUGBEADJO A., MAILLOL V., MARTIN G.,

352           SABOT F., SANTONI S., VIGOUROUX Y., COUVREUR T.L.P. 2014. Cost effective enrichment hybridization

353           capture of chloroplast genomes at deep multiplexing levels for population genetics and

354           phylogeography studies. Mol. Ecol. Resour. 14:1103–1113.

355     MILNE I., BAYER M., CARDLE L., SHAW P., STEPHEN G., WRIGHT F., MARSHALL D. 2010. Tablet—next generation

356           sequence assembly visualization. Bioinformatics. 26:401–402.

357     QUINLAN A.R., HALL I.M. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.

358           Bioinformatics. 26:841–842.

359     ROHLAND N., REICH D. 2012. Cost-effective, high-throughput DNA sequencing libraries for multiplexed

360           target capture. Genome Res. 22:939–946.

361     SCARCELLI N., BARNAUD A., EISERHARDT W., TREIER U.A., SEVENO M., d'ANFRAY A., VIGOUROUX Y., PINTAUD J.-C.

362           2011. A Set of 100 Chloroplast DNA Primer Pairs to Study Population Genetics and Phylogeny in

363           Monocotyledons. PLoS ONE. 6:e19954.

364     STRAUB S.C.K., PARKS M., WEITEMIER K., FISHBEIN M., CRONN R.C., LISTON A. 2012. Navigating the tip of the

365           genomic iceberg: Next-generation sequencing for plant systematics. Am. J. Bot. 99:349–364.

366     STULL G.W., MOORE M.J., MANDALA V.S., DOUGLAS N.A., KATES H.-R., QI X., BROCKINGTON S.F., SOLTIS P.S., SOLTIS

367           D.E., GITZENDANNER M.A. 2013. A Targeted Enrichment Strategy for Massively Parallel Sequencing

368           of Angiosperm Plastid Genomes. Appl. Plant Sci. 1:1200497.

17

369     YANG, C., CHU, J., WARREN, R. L., & BIROL, I. 2017. NanoSim: Nanopore sequence read simulator based on

370         statistical characterization. GigaScience.

371

372

## 373    **Tables**

374    Table 1: MinION plastome enriched library output data. The percentage of "plastome mapped reads"

375    was calculated using BWA to indicated reference plastomes.

| Species | DNA | probe origin | total number of reads | Median read length | longest read | % of plastid reads | longest plastome read | Median plastome read length |
|---|---|---|---|---|---|---|---|---|
| *Oryza sativa* | fresh | *O. sativa* | 17129 | 4627 | 26128 | 70,8 | 25828 | 4264 |
| *Oryza glaberrima* | fresh | *O. sativa* | 81361 | 3695 | 24804 | 98,2 | 24504 | 3398 |
| *Cenchrus americanus* | fresh | *C. americanus* | 105760 | 4914 | 25468 | 97,0 | 25167 | 4623 |
| *Digitaria exilis* | fresh | *C. americanus* | 141250 | 3783 | 19378 | 94,4 | 19078 | 3489 |
| *Podococcus acaulis* | silicagel | *P. barteri* | 202924 | 2486 | 13103 | 15,7 | 12805 | 2129 |
| *Raphia textilis* | silicagel | *P. barteri* | 83833 | 2322 | 10705 | 87,5 | 10405 | 1997 |
| *Phytelephas aequatorialis* | silicagel | *P. barteri* | 202925 | 2437 | 15132 | 79,0 | 14832 | 2158 |

376

377

378    Table 2: De novo assembly results from real (6) and simulated (1) data in number of contigs, and

379    coverage and identity percentages to the respective reference plastome genomes (see Table 1). Min

380    overlap is the minimum overlap between reads as defined in Flye. The simulated data was based on the

381    out pu results of *P. aequatorialis*.

| Species | Min overlap | Plastid contigs | Coverage % | Identity % | Longest contig |
|---|---|---|---|---|---|
| *Oryza glaberrima* | 3000 | 2 | 92.32 | 99.14 | 109087 |
| *Cenchrus americanus* | 3000 | 2 | 99.91 | 98.52 | 81053 |
| *Digitaria exilis* | 3000 | 2 | 99.97 | 99.18 | 125727 |
| *Podococcus acaulis* | 1000 | 17 | 81.24 | 98.86 | 22803 |
| *Raphia textilis* | 1000 | 10 | 83.87 | 98.84 | 21797 |

18

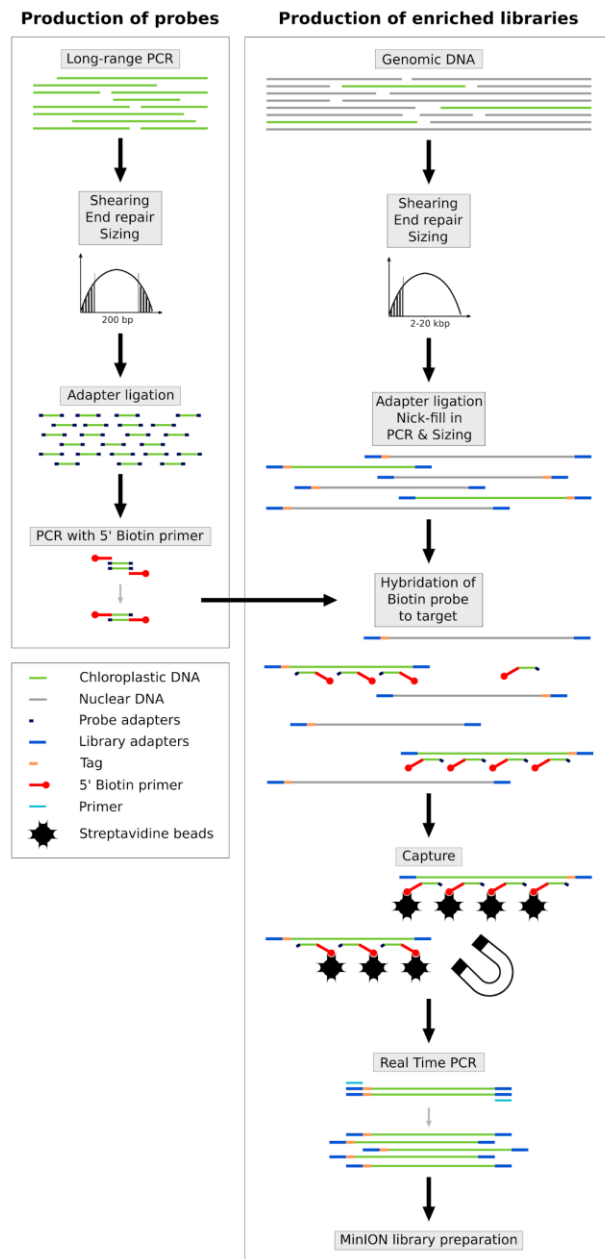| | | | | | |
|---|---|---|---|---|---|
| *Phytelephas aequatorialis* | 1000 | 13 | 87.60 | 98.31 | 20700 |
| *Simulated assembly* | 1000 | 4 | 99.72 | 99.05 | 107633 |

382

383    Figures (next page)



384

385    Figure 1: Schematic representation of the protocol used for long sequence capture of plastomes
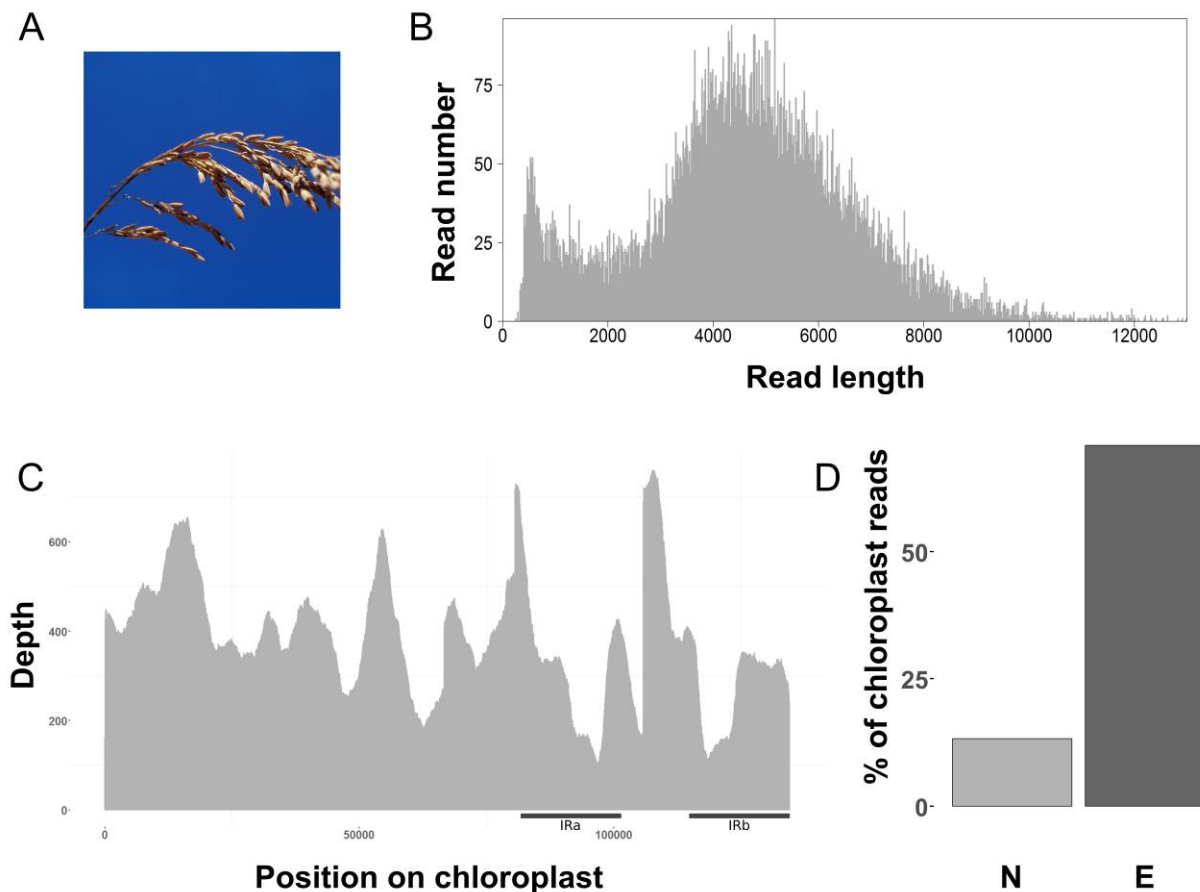
386    (modified from Mariac et al. 2014).

387

19

388

389 Figure 2: Long fragment capture results for *Oryza sativa*. A: Panicule of *Oryza sativa* (Jean-Pierre

390 Montoroi IRD ©). B: number of reads per read length before mapping. C: Plastome coverage after

391 mapping. Black bars indicate approximate position of both inverted repeats (IR). D: Percentage of useful

392 reads mapped to *Oryza sativa* reference plastome (KT289404.1) between non enriched (N, light grey)

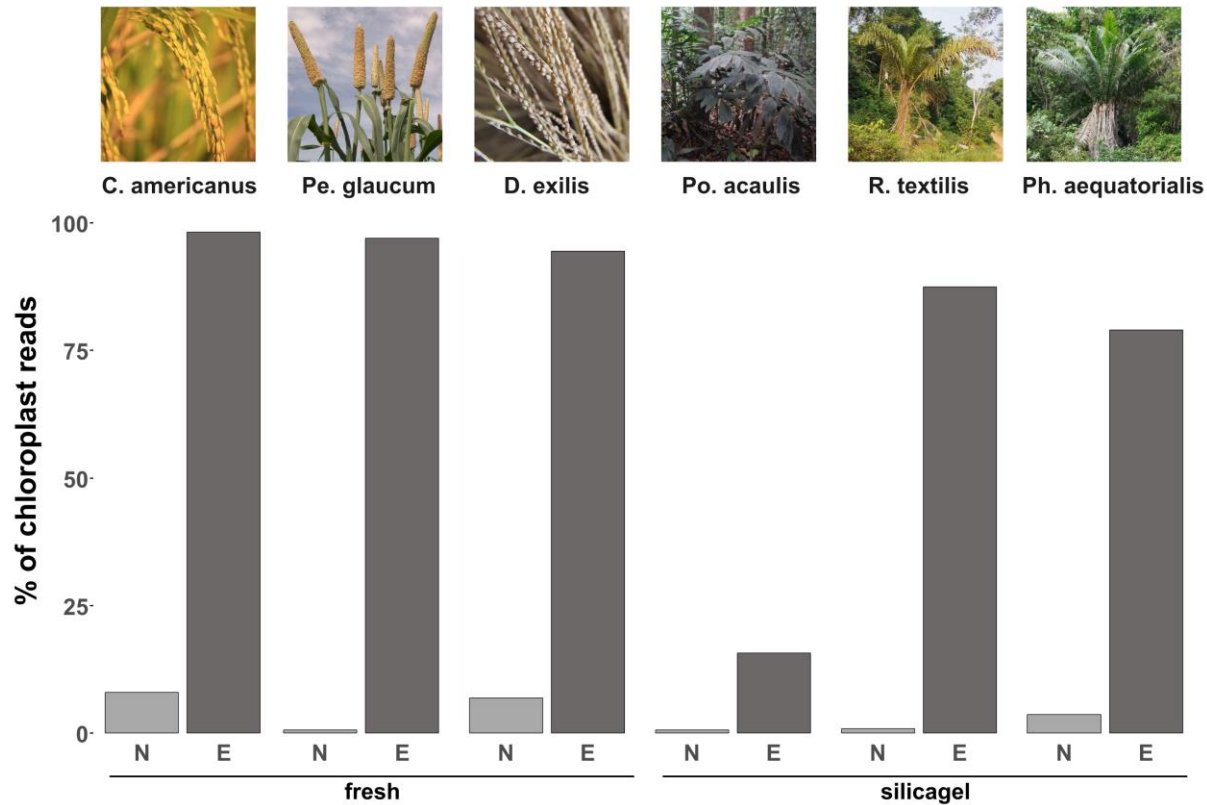393 and enriched (E, dark grey) libraries.

394

Figure 3: Percentage of useful reads mapped to their respective reference plastome (see Table 1) between Illumina non enriched (N, light grey) and MinION enriched (E, dark grey) protocols for the 6 non model species in our study. Photos: *O. glaberrima*: https://pxhere.com/fr/photo/706162; CC0 public domain; *Cenchrus americanus*: C. Mariac IRD ©; *D. exilis*: A. Barnaud IRD © ; *Po. acaulis, R. textilis, Ph. aequatorials*: TLP Couvreur IRD ©.
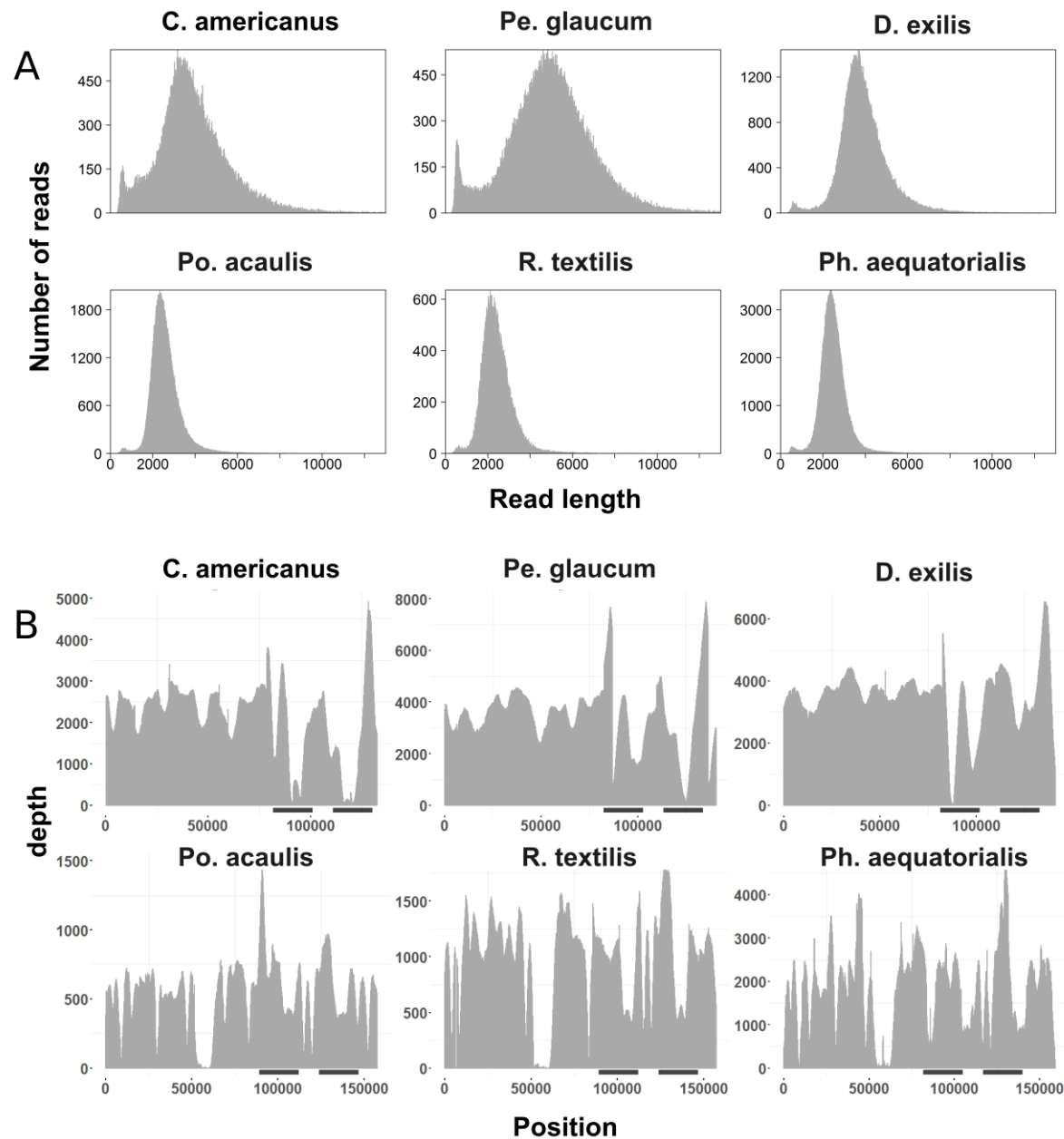
Figure 4: Long fragment capture results for six non model plant species. A: Number of reads per read
length before mapping. B: Plastome coverage results from the enriched long read capture protocol. Black
bars indicate approximate position of both inverted repeats (IR).