1  **Long transposon-rich centromeres in an oomycete reveal divergence of**

2  **centromere features in Stramenopila-Alveolata-Rhizaria lineages**

3

4

5

6  Yufeng "Francis" Fang[a], Marco A. Coelho[a], Haidong Shu[b], Klaas Schotanus[a], Bhagya C.

7  Thimmappa[c,1], Vikas Yadav[a], Han Chen[b], Ewa P. Malc[d], Jeremy Wang[d], Piotr A. Mieczkowski[d],

8  Brent Kronmiller[e], Brett M. Tyler[e], Kaustuv Sanyal[c], Suomeng Dong[b], Minou Nowrousian[f], and

9  Joseph Heitman[a,2]

10

11

12

13  [a]Department of Molecular Genetics and Microbiology, Duke University Medical Center, Durham,

14  North Carolina 27710; [b]College of Plant Protection, Nanjing Agricultural University, Nanjing,

15  China; [c]Molecular Biology and Genetics Unit, Jawaharlal Nehru Centre for Advanced Scientific

16  Research, Bangalore, India; [d]Department of Genetics, University of North Carolina, Chapel Hill,

17  North Carolina 27599; [e]Center for Genome Research and Biocomputing and Department of

18  Botany and Plant Pathology, Oregon State University, Corvallis, Oregon 97331; [f]Lehrstuhl fuer

19  Molekulare und Zellulaere Botanik, Ruhr-Universitaet Bochum, Bochum, Germany

20

21

22

23  [1]Present address: Department of Biochemistry, Robert-Cedergren Centre for Bioinformatics and

24  Genomics, University of Montreal, 2900 Edouard-Montpetit, Montreal, H3T1J4, QC, Canada.

25  [2]To whom correspondence should be addressed. Email: heitm001@duke.edu.

26    **Abstract**

27    Centromeres are chromosomal regions that serve as platforms for kinetochore assembly

28    and spindle attachments, ensuring accurate chromosome segregation during cell division.

29    Despite functional conservation, centromere DNA sequences are diverse and often repetitive,

30    making them challenging to assemble and identify. Here, we describe centromeres in an

31    oomycete *Phytophthora sojae* by combining long-read sequencing-based genome assembly and

32    chromatin immunoprecipitation for the centromeric histone CENP-A followed by high-throughput

33    sequencing (ChIP-seq). *P. sojae* centromeres cluster at a single focus at different life stages and

34    during nuclear division. We report an improved genome assembly of the *P. sojae* reference strain,

35    which enabled identification of 15 enriched CENP-A binding regions as putative centromeres.

36    By focusing on a subset of these regions, we demonstrate that centromeres in *P. sojae* are

37    regional, spanning 211 to 356 kb. Most of these regions are transposon-rich, poorly transcribed,

38    and lack the histone modification H3K4me2 but are embedded within regions with the

39    heterochromatin marks H3K9me3 and H3K27me3. Strikingly, we discovered a Copia-like

40    transposon (CoLT) that is highly enriched in the CENP-A chromatin. Similar clustered elements

41    are also found in oomycete relatives of *P. sojae*, and may be applied as a criterion for prediction

42    of oomycete centromeres. This work reveals a divergence of centromere features in oomycetes

43    as compared to other organisms in the Stramenopila-Alveolata-Rhizaria (SAR) supergroup

44    including diatoms and *Plasmodium falciparum* that have relatively short and simple regional

45    centromeres. Identification of *P. sojae* centromeres in turn also augments the genome assembly.

46

47    Key words: Stramenopila; genome assembly; CENP-A; Copia-like transposon

48

49 **Significance Statement**

50    Oomycetes are fungal-like microorganisms that belong to the stramenopiles within the

51   Stramenopila-Alveolata-Rhizaria (SAR) supergroup. The *Phytophthora* oomycetes are infamous

52   as plant killers, threatening crop production worldwide. Because of the highly repetitive nature of

53   their genomes, assembly of oomycete genomes presents challenges that impede identification of

54   centromeres, which are chromosomal sites mediating faithful chromosome segregation. We

55   report long-read sequencing-based genome assembly of the *Phytophthora sojae* reference strain,

56   which facilitated the discovery of centromeres. *P. sojae* harbors large regional centromeres

57   enriched for a Copia-like transposon that is also found in discrete clusters in other oomycetes.

58   This study provides insight into the oomycete genome organization, and broadens our knowledge

59   of the centromere structure, function and evolution in eukaryotes.

3

## Introduction

Accurate segregation of chromosomes during mitosis and meiosis is critical for the development and reproduction of all eukaryotic organisms. Centromeres are specialized regions of chromosomes that mediate kinetochore formation, spindle attachment, and sister chromatid segregation during cell division (1, 2). The DNA coincident with functional centromeres typically consists of unusual sequence composition (e.g. AT-rich) and structure (e.g. repeats, transposable elements), low gene density, and transcription of non-coding RNA (ncRNA) as well as heterochromatic nature (3). However, an active centromere is defined not by DNA sequences but by the deposition of a centromere-associated protein called centromere protein A (CENP-A, also known as CenH3) (1, 4). CENP-A is a histone H3 variant, which replaces the canonical H3 in the nucleosomes at centromeres and provides the foundation for kinetochore assembly (1, 5, 6).

Despite the fact that centromere function is broadly conserved, centromeric sequences vary greatly in size and composition, ranging from "point" centromeres of 125 bp in length to "regional" centromeres consisting of up to megabases of repeated sequences to holocentromeres that extend along the entire length of the chromosome (1, 3). To date, point centromeres have been only reported in the budding yeast *Saccharomyces cerevisiae* and its close relatives, holocentromeres have been identified in some insects, plants and nematodes, represented by *Caenorhabditis elegans*, while regional centromeres are the most common type and found in nearly all eukaryotic phyla (1, 3). Most animals and plants have large regional centromeres composed of satellite sequences that are organized into a variety of different higher order repeats (4, 7, 8). Some plant centromeres also possess a different type of repeat called centromere-specific retroelements (CR) (9). In comparison, all fungal centromeres identified to date do not contain satellite repeats and have diverse organizations. The size of fungal regional centromeres ranges from several kilobases, such as in *Candida albicans*, to hundreds of kilobases in *Neurospora crassa* (10, 11). The centromeric sequences of fungal regional

4

85   centromeres can be composed of active or inactive clusters of transposable elements and thus

86   very repetitive, such as in *Cryptococcus* spp. and *N. crassa* (12, 13), or can be nonrepetitive and

87   very short, such as in the wheat pathogen *Zymoseptoria tritici* (14) and *C. albicans* (15).

88   Information on centromeres is limited in other eukaryotic lineages. The malaria pathogen

89   *Plasmodium falciparum* and the diatom *Phaeodactylum tricornutum* CENP-A binding regions are

90   characterized by short simple AT-rich sequences (16, 17), while the parasite *Toxoplasma gondii*

91   has a simple centromere without nucleotide bias (18).

92        Due to their highly repetitive nature, assembly of large regional centromeres presents a

93   significant challenge. Emerging long-read sequencing technologies, such as Pacific Bioscience

94   (PacBio) and Oxford Nanopore Technology (ONT), have led to substantial advances in resolution

95   of chromosomal structures including highly repetitive sequences such as centromeres. Using

96   these technologies, centromeres that were difficult to resolve using short-read sequencing, were

97   defined in various organisms, from fungi (12, 19, 20) to insects (21), plants (22) and humans (23).

98        Oomycetes are fungal-like organisms but belong to the stramenopila kingdom within the

99   Stramenopila-Alveolata-Rhizaria (SAR) supergroup (24, 25). The SAR supergroup contains a

100  high diversity of lineages that include many important photosynthetic lineages (e.g. diatoms and

101  kelp), and important parasites of animals (e.g., *Plasmodium*, the causative agent of malaria) and

102  plants (e.g., oomycetes, or water molds) (26). *Phytophthora* is a large oomycete genus (>160

103  species found to date) and contains some of the most devastating plant pathogens that destroy

104  a wide range of plants important in agriculture, forestry, ornamental and recreational plantings,

105  and natural ecosystems (27). One notorious example is *Phytophthora infestans*, which caused

106  the great Irish potato famine of the mid-1840s (28). Today, *Phytophthora* species remain

107  significant threats to major food crops, causing multi-billion US dollars losses annually throughout

108  the world (27, 29). *Phytophthora sojae* is a widespread soil-borne pathogen of soybean. Because

109  of its economic impact, and tractable genetic manipulation (30-32), *P. sojae* has become a model

110  species to study oomycete genetics, biology, and interactions with plants.

111      To date, the genomes of more than 20 *Phytophthora* species have been sequenced (33).

112      Their genomes are generally large and display complex features: they are diploid, highly

113      heterozygous for heterothallic species, and very repetitive, which makes genome assembly

114      challenging. The most contiguous oomycete genome assembly published to date is of the *P. sojae*

115      reference genome, which was generated based on Sanger random shotgun sequencing and

116      subsequent improvements involving gap closure and BAC sequencing (25, 34). *P. sojae* genome

117      assembly v3.0 (www.jgi.doe.gov) spans ~82 Mb and contains 82 scaffolds; however, there are

118      ~3 Mb of unresolved gaps (N's) persisting in the assembly. Recently, significant progress has

119      been made in genome assemblies of oomycetes based on long-read sequencing (35, 36);

120      however, the identity or the nature of the DNA sequences that form essential chromosomal

121      elements such as centromeres, remain unknown. In this study, using the evolutionarily conserved

122      kinetochore protein CENP-A as a tool, we investigated cellular dynamics of the kinetochore

123      complex in *P. sojae*, and uncovered the nature of the oomycete centromeres with the aid of long-

124      read genome sequencing and ChIP-seq technologies. Our findings suggest that the centromeres

125      of *P. sojae* are divergent from those reported in other SAR lineages, and their features may be

126      used to predict centromeres in other oomycetes.

127

128      **Results**

129      **GFP-tagging of CENP-A in *P. sojae* reveals clustered centromeres in different life stages**

130      **and throughout hyphal growth**

131      Kinetochore protein homologs have been predicted in diverse eukaryotic lineages

132      including oomycete species (37). To identify kinetochore proteins in *P. sojae*, we conducted

133      BLAST searches against the existing *P. sojae* genome database using the predicted oomycete

134      orthologs as query. Gene models of *P. sojae* kinetochore proteins were examined and corrected

6

135    based on RNA-seq data when necessary. Protein sequences were verified based on the presence

136    of corresponding motifs (Fig. S1 and Dataset S1).

137         To examine centromere/kinetochore organization and localization in *P. sojae*, we selected

138    CENP-A, the hallmark of centromere identity in most organisms. The RNA-seq data did not

139    support the gene models of *CENP-A* that was instead verified by 3'-RACE and RT-PCR, followed

140    by Sanger sequencing (Fig. S2 *A* and *B*). *P. sojae* CENP-A (PsCENP-A) has a conserved C-

141    terminus including the "CENP-A targeting domain" (CATD) (Fig. S2*C*). *GFP* was fused to *CENP-*

142    *A* at the N-terminus and transiently expressed in *P. sojae* transformants with a constitutive

143    promoter derived from the *Bremia lactucae HAM34* gene (Fig. S2*D*). Overexpressed GFP-CENP-

144    A exhibited nuclear localization with a single fluorescent focus in the nucleus (Fig. S2*D*),

145    suggesting that *P. sojae* has a clustered centromere organization.

146         We also generated GFP labeled CENP-A expressed from the endogenous locus utilizing

147    CRISPR/Cas9-mediated gene replacement (Figs. 1A and S3). Homokaryotic GFP-CENP-A

148    strains exhibited single GFP foci within nuclei from different *P. sojae* life stages (Figs. 1B)*,*

149    confirming that the clustered centromere organization is a feature in *P. sojae.* In addition, we

150    tracked the centromere dynamics during hyphal growth. Intriguingly, the clustered centromere

151    pattern was maintained throughout *P. sojae* nuclear division (Fig. 1C and Movie S1).

152

**Identification of centromeres in a long-read Nanopore-based assembly**

154         To identify *P. sojae* centromeres, we performed native chromatin immunoprecipitation (N-

155    ChIP) using an anti-GFP antibody against the GFP-CENP-A fusion, followed by high-throughput

156    Illumina DNA sequencing. ChIP-seq reads were mapped to the latest Sanger genome assembly

157    (*P. sojae* V3 from JGI), which identified 12 scaffolds that showed relatively concentrated

158    enrichment of CENP-A reads (Fig. S4*A*). CENP-A peaks appeared scattered in Scaffold 1 and

7

159    Scaffold 11, while more clustered in the other 10 scaffolds. However, further examination of each

160    CENP-A binding region revealed that most of the regions were interrupted by many sequence

161    gaps, which hampered analysis of the sequence features of the candidate centromeres. Thus, we

162    processed to re-sequence and re-assemble the reference *P. sojae* genome.

163    To improve the genome assembly of *P. sojae* reference strain P6497, we applied

164    Nanopore long-read sequencing and generated a *de novo* genome assembly with SMARTdenovo

165    together with polishing from PacBio and Sanger reads (Fig. S5*A* and *Appendix SI Text*). The

166    resulting assembly of the nuclear genome (Psojae2019.1) has a size of 86 Mb contained in 70

167    contigs, with a contig N50 of 2 Mb (Fig. S5*C*). Comparison of Psojae2019.1 to the Sanger

168    assembly indicated that Psojae2019.1 has more repetitive sequences and most regions were

169    colinear (Fig. S5 *B* and *C*, also see *Appendix SI Text* for details). We also checked telomere

170    repeats using a motif proposed for oomycetes (38), and found 13 contigs (versus 7 in Sanger)

171    that harbor telomeric sequences at single ends (*Appendix SI Text* and Dataset S2).

172    ChIP-seq reads derived from PsCENP-A were mapped to the new genome assembly

173    Psojae2019.1 (Table S2), which initially revealed 16 regions exhibiting CENP-A enrichment. On

174    closer analysis, we found that the unassembled centromere in contig 20 was an artifact caused

175    by inaccurate genome assembly, as this region was duplicated with a centromere-containing

176    region in contig34 (Fig. S6*F*). Of the 15 remaining CENP-A binding regions, 11 regions were

177    assembled within contigs, whereas four regions were disrupted at the edge of contigs (Fig. 2).

178    Long-read coverage analysis verified the integrity of 10 centromeres (Fig. S7), while the CENP-

179    A peaks in Contig 37 and three broken ones (in Contigs 9, 10, 57) lacked sufficient long-read

180    coverage. We focused on the 10 verified CENP-A regions for the further studies (Table 1).

181    RNAseq analysis indicated that all of the 10 CENP-A regions exhibited low transcription, except

182    the region in Contig 11. Contig 11 contained two adjacent ChIP-seq peaks, one was 19 kb and

183    the other was 114 kb, which were interrupted by a 21 kb transcriptionally active region (Fig. 4*C*).

184    Here, we define it as one centromere (*CEN4*). Among the 10 CENP-A regions, five have a length

185    of ~190 kb, and three are ~160 kb, while *CEN3* and *CEN6* are significant larger (>270 kb) (Table

186    1). All of these centromeres have a GC content comparable to the whole genome (52.16 - 58.13%

187    vs. 54.6%) (Table 1). Taken together, our CENP-A ChIP-seq analysis utilizing the newly

188    assembled genome indicates that *P. sojae* CENP-A prefers to bind large poorly transcribed

189    genomic regions with no specific DNA sequence bias.

190          To examine the correlation between the centromere regions identified in the new genome

191    assembly and in the Sanger assembly, we conducted synteny analysis using the genomic regions

192    flanking the centromeres. The locations of CENP-A found in the Psojae2019.1 assembly were

193    highly correlated with those in the Sanger assembly, except *CEN10* (Table 1, Figs. 3 and S6).

194    Contig 51 was colinear with the Sanger scaffold 23; however, no enriched CENP-A signal was

195    detected for this scaffold, probably because the region corresponding to *CEN10* is interrupted by

196    gaps. Notably, the two CENP-A binding regions in Sanger Scaffold 1 were found to correspond

197    to *CEN8* and *CEN9*, and the smaller one was expanded from 20 kb to 188 kb corresponding to

198    *CEN8* (Table 1, Figs. S4*B* and S6*H*). In addition, four contigs of the Psojae2019.1 assembly

199    (contigs 4, 38, 23, and 58) are collinear with Sanger Scaffold 1, and telomere repeats are found

200    at the ends of Contigs 4 and Contig 58, further suggesting that Scaffold 1 of the Sanger genome

201    is assembled incorrectly and should be split into two scaffolds (Fig. S6*H*). Overall, comparison of

202    centromeres identified in the Sanger and Psojae2019.1 assemblies further confirms their

203    authenticity and reflects some misassemblies that are present in the Sanger genome assembly.

204

205    ***P. sojae* CENP-A regions are embedded within heterochromatin**

206          To define the epigenetic state of *P. sojae* centromeric regions, we performed ChIP-seq

207    with antibodies against two heterochromatin marks (H3K9me3, trimethylation of lysine 9 of

208    histone H3, and H3K27me3, trimethylation of lysine 27 of histone H3) and one euchromatin mark

9

209    (H3K4me2, dimethylation of lysine 4 of histone H3). The distribution of H3K9me3 and H3K27me3

210    is generally coincident throughout the genome, and both were colocalized with the CENP-A

211    binding regions (Figs. 4 and S7). Intriguingly, the heterochromatic region extended 8kb to 64 kb

212    beyond each CENP-A binding region (Fig. 4A and Table 1), similar to pericentromeric

213    heterochromatin regions described in other species (13, 21, 39). In contrast, the euchromatic

214    mark H3K4me2 was excluded from the CENP-A region and its flanking pericentric regions, and

215    generally overlapped with the mRNA transcriptional profile (Figs. 4 *B-C* and S7). Thus, distribution

216    of histone modifications suggests that the CENP-A and heterochromatin regions are not spatially

217    distinct, and we define the latter as pericentric regions.

218

**A Copia-like transposon (CoLT) is highly enriched in the *P. sojae* centromeres**

220          The Psojae2019.1 genome assembly contains 31% repetitive sequences, the majority of

221    which are transposable elements (TEs) (Fig. S5*D*). Our analysis showed that centromeres are

222    also composed of many repetitive elements, mostly LTR-retrotransposons (Figs. 3 and S6). To

223    identify whether the centromeres in *P. sojae* possess any common sequences or repeat elements,

224    all identified CENP-A regions were subject to multiple sequence alignment. This analysis found

225    an ~5 kb sequence that is highly similar (>98%) and shared among 10 centromeres (Fig. S8 and

226    Dataset S3). BLAST analyses with the consensus 5 kb sequence against the genome revealed

227    that although this element is not exclusive to centromeres, it is significantly enriched in

228    centromeres: approximately 90% of all genomic copies of this element localized to centromeres

229    (Fig. 5A). Moreover, this element is present as clusters in centromeric regions, and only sparsely

230    found in other regions of the genome, further strengthening its association with centromeres.

231    Further examination of the sequence indicates that it resembles a Copia transposon-like

232    transposon, and we named it CoLT for Copia Like Transposon (Fig. 5*B*, Dataset S3).

233

234 **CoLT clusters are conserved in two *P. sojae* oomycete relatives and may be a hallmark of**

235 **oomycete centromeres**

236       To examine if clustered CoLT elements found in *P. sojae* centromeres are also present in

237 other oomycete genomes, we conducted BLAST searches using the 5 kb consensus sequence

238 derived from *P. sojae* centromeres against the genome assemblies of two *P. sojae* relatives,

239 *Bremia lactucae* (downy mildew, lettuce pathogen) and *Phytophthora citricola* (citrus pathogen)*,*

240 which have relatively contiguous genome assemblies. Interestingly, similar CoLTs clusters were

241 observed in these genomes, and usually appeared once per contig (Figs. 5*C* and S9*A*). To assess

242 if these clustered CoLTs were syntenic with the *P. sojae* centromere-containing contigs, we

243 examined the CoLT clusters that were present within Mb-long scaffolds/contigs. Synteny analysis

244 demonstrated that five regions in the *B. lactucae* genome that had CoLT clusters were collinear

245 with *P. sojae* centromeres (Figs. *5C* and *D).* Unexpectedly, Scaffold 2 (original name,

246 SHOA01000004.1, see Dataset S4 for details) contained two CoLT clusters that were syntenic

247 with *P. sojae CEN3* and *CEN5* (Fig. 5D), indicating that scaffold 2 may be incorrectly assembled

248 (Fig. 5D). It should be noted that the *B. lactucae* genome assembly still has a large percentage

249 of unresolved gaps likely due to its highly heterozygous nature (36). In comparison, all three

250 selected regions that had clustered CoLT clusters within *P. citricola* contigs (PcContigs) were

251 syntenic with *P. sojae* centromeres (*CEN3*/PcContig2, *CEN9*/PcContig1, *CEN5*/PcContig26)

252 (Figs. S9 *B-D*). However, a large number of the CoLT clusters localized at contig ends, or were

253 distributed across the length of short contigs (Fig. S9*A*). This suggests that many of the

254 centromeric regions in *P. citricola* were not fully assembled. Taken together, we propose that the

255 clustered CoLT elements may be used as criteria to predict centromere regions in other

256 *Phytophthora* species and possibly in other oomycetes.

11

**Discussion**

257

258      In this study, we identified centromeres in the oomycete plant pathogen *P. sojae* by

259    combining long-read sequencing and ChIP-seq with the GFP tagged kinetochore protein CENP-

260    A. Cellular dynamics analysis revealed that *P. sojae* centromeres were clustered within nuclei in

261    different life stages and during vegetative growth. 10 fully assembled and five incompletely

262    assembled CENP-A binding regions were identified. The common features shared by these

263    regions include: a) a low level of transcription; b) a GC content similar to that of the whole genome;

264    c) repetitive sequences; d) enrichment for a specific Copia-like transposon; e) overlapping and

265    surrounding heterochromatin; and f) lack of H3K4me2.

266      While CENP-A is conserved among different organisms, centromere sequences evolve

267    rapidly (1, 40). Although the filamentous fungal-like oomycetes are classified in the stramenopiles

268    of the SAR supergroup, it is intriguing to observe that the centromeres that we identified in *P.*

269    *sojae* are much larger and more complex, comparing to those reported in its stramenopile relative,

270    the diatom *P. tricornutum*, and those found in the parasites (*P. falciparum* and *T. gondii*) of the

271    alveolates (Fig. 6). In the latter three cases, all centromeres are composed of non-repetitive

272    sequences. Surprisingly, *P. sojae* centromeres show structural similarity to several, only distantly

273    related, fungal species, such as *N. crassa* (13) and *Cryptococcus neoformans* (12). These

274    features include an enrichment of transposons (or their remnants), and overlap with the

275    constitutive heterochromatin mark H3K9me2/3. Remarkably, the euchromatin mark H3K4me2

276    has been shown to be associated with centromeres in humans, mouse, *Drosophila*, *S. pombe*,

277    and rice (39, 41-43), but is excluded from other fungal regional centromeres reported to date and

278    in *P. sojae*. In humans and *D. melanogaster*, the CENP-A and pericentromeric heterochromatin

279    domains are spatially distinct, and the CENP-A domain is flanked by but does not overlap with

280    heterochromatin (39, 43, 44). In contrast, the entire centromere of *P. sojae* is embedded in

281    heterochromatin. It is unknown if the distribution of heterochromatin regions affects centromere

282    distribution in *P. sojae*, but heterochromatin has been shown to be important for centromere

283    function and kinetochore assembly in *N. crassa* and *S. pombe* (13, 45, 46). In addition, it is of

284    interest that *P. sojae* H3K9me3 and H3K27me3 fully overlap with the centromeric regions, which

285    have not been observed in centromeres of other species thus far, but was shown in human and

286    mouse pericentromeres (8, 47). On the other hand, these two epigenetic marks generally coexist

287    throughout the entire genome, suggesting it might be just a general profile of H3K27me3 and

288    H3K9me3 in *P. sojae*.

289        Transposable elements (and their relics) have been known as residents of the

290    centromeres and pericentromeres of many animals, plants, and fungi (48). While animal

291    centromeres are associated with both satellite DNA and retroelements, satellite DNA is usually

292    regarded as the main sequence components (49). Centromeres of many plants, such as maize

293    and rice, are built on centromere-specific retrotransposons (CR), and a certain CR is usually

294    unique to a particular chromosome (7). Centromeres of *N. crassa* (13) and *C. neoformans* (12)

295    are composed of retrotransposons, and the retroelements in *C. neoformans* are centromere-

296    specific (12). In comparison, although *P. sojae* regional centromeres include various transposons,

297    many of these elements are not limited to this region and can also be found in other genomic

298    areas. Our study shows that a specific Copia-like transposon (CoLT) is highly enriched in the *P.*

299    *sojae* centromeric regions and confines the CENP-A binding regions (Figs. 4 *B-C* and S7). A

300    similar distribution pattern of centromere-associated retrotransposons was recently found in

301    *Drosophila melanogaster* (21). In *D. melanogaster*, a non-LTR retroelement named *G2*/*Jocky-3*

302    was found to be enriched in CENP-A chromatin, and this element is also associated with

303    centromeres in its sister species *D. simulans* (21). Strikingly, the CoLT elements were found to

304    be clustered in the genomes of *P. sojae* oomycete relatives, and some of those regions were

305    syntenic with *P. sojae* centromeric regions. As most of the oomycete genome assemblies were

306    not based on long-read sequencing technology, and thus are very fragmented, it remains to be

13

307     seen if the CoLT elements have evolved to be widely utilized by oomycetes as a platform for

308     CENP-A loading.

309         Due to large genome scales and potentially similar chromosome sizes, the karyotypes of

310     *Phytophthora* species cannot be well resolved by pulse field gel electrophoresis (31, 50). The

311     chromosome number of *P. sojae* is not yet accurately known, but has been estimated to be

312     between 10 and 15 based on an earlier cytological study (51). By comparing the location of

313     centromeres in the Sanger and Psojae2019.1 assembly, we can validate and predict the

314     configuration of 11 centromeres, namely *CEN1-CEN10*, and *CEN_C9* + *CEN_C48* (Table 1 and

315     Table S4). Three centromeres, namely *CEN_C37*, *CEN_C10* and *CEN_C57*, are not fully

316     assembled. Thus, our results offer a new estimate 12-14 chromosomes in *P. sojae*.

317         N-ChIP was implemented for this study, because several attempts to perform ChIP

318     analysis based on traditional formaldehyde-cross-linking strategies were unsuccessful. Cross-

319     linking with 1% formaldehyde caused degradation of DNA and failure of ChIP. *P. sojae*

320     transformants expressing GFP tagged CENP-A and CENP-C were both used for N-ChIP-seq.

321     However, only the GFP-CENP-A transformant produced significant enrichment, indicating that the

322     binding of CENP-C to chromosomes may be too weak to recover target DNA under native

323     conditions without cross-linking.

324         Our analysis showed that having an improved reference genome assembly based on long-

325     read sequencing technologies was crucial to the identification and characterization of

326     centromeres in *P. sojae*. Our attempt to characterize centromere sequences using the classical

327     Sanger assembly was not successful because most of the non-coding repetitive regions were not

328     assembled. While the N50 of the new genome assembly Psojae2019.1 is lower than that of the

329     Sanger assembly, the contigs do not contain gaps and many of the gaps present in the Sanger

330     assembly have been closed (Fig. S6). We tried to scaffold the assembly with different scaffolding

331     programs such as npScarf (52), SSPACE (53), LINKS (54) and the optical BioNano mapping

332    (*Appendix SI Text* and Fig. S10). Although these scaffolders improved the contiguity (up to 35

333    scaffolds using SSPACE), they also generated multiple conflicts with the Sanger assembly, and

334    most of the joins could not be supported by evidence such as long read coverage (Fig. S11 and

335    Table S3). Thus, we opted to retain the contig-level assembly in our study. However, identification

336    of centromeres helped to resolve several structural problems present in the "classical" *P. sojae*

337    Sanger assembly, and revealed potential structural problems in other oomycete genome

338    assemblies. On basis of the presence of centromeres and predicted telomeres together with

339    synteny analyses, we found that three Sanger scaffolds/contigs may represent full length

340    chromosomes, namely Scaffold 2/Contigs [26+1+35+6] (Fig. S6*A*), Scaffold 5/Contigs

341    [17+36+7+49+45] (Fig. S6*G*); and partial Scaffold 1/Contigs [58+38+4] (Fig. S6*H*). Notably,

342    telomeres appear on the both ends of Sanger Scaffold 5 and its syntenic contigs in Psojae2019

343    (Fig. S6G). There are five *P. sojae* centromeres that are not fully assembled. With the

344    development of sequencing and assembly technologies, a finalized chromosome-level genome

345    assembly could help to assemble those broken centromeres, and refine the centromere

346    sequences that we identified.

347      Centromeres and their associated kinetochore network serve critical functions in genome

348    stability and replication. Failures in kinetochore assembly and attachment increase the probability

349    of chromosome mis-segregation leading to aneuploidy (55). While these drastic genome changes

350    can be detrimental to the organism, formation of aneuploidy and polyploidy is an important

351    strategy orchestrated by pathogens to adapt to the environment during periods of stress (56).

352    Polyploidy and aneuploidy are prevalent in *Phytophthora* natural isolates and progeny from sexual

353    reproduction (35, 57-60). Interestingly, plant hosts can induce aneuploidy of the sudden oak death

354    pathogen *P. ramorum*, which enhances its phenotypic diversity and increases its adaption to the

355    environment (59). Recently, a phenomenon termed dynamic extreme aneuploidy (DEA) was

356    described in a vegetable oomycete pathogen, *P. capsici*, in which high variability among progeny

357    produced by asexual spores was caused by ploidy variation (61). However, the mechanisms

15

358    resulting in oomycete aneuploidy and/or polyploidy is understudied. As centromeres are the

359    functional and structural foundation for kinetochore assembly and proper chromosome

360    segregation, identification of centromeres and kinetochore proteins in *P. sojae* may help to

361    illuminate the mechanisms underlying oomycete genetic, genomic, and phenotypic diversification.

**Materials and methods**

### *P. sojae* culture and transformation

All the strains used in this study are listed in Table S5. The reference *P. sojae* isolate P6497 (race 2) used in this study was routinely grown and maintained in cleared V8 media at 25 °C in the dark. Transient gene expression assays based on an optimized polyethylene glycol (PEG) mediated protoplast transformation protocol (30) was applied to examine the nuclear localization of CENP-A. Stable and homokaryotic transformants were chosen for ChIP-seq, which were generated by passaging on V8 supplemented with 50 µg/mL G418 (Geneticin, AG Scientific, San Diego, California, USA) for at least 5 times followed by zoospore isolation. Co-transformation was employed to generate strains expressing both H2B-mCherry and GFP-CENP-A. Transformation was performed as previously described (30). Sporangia and zoospores were induced by water flooding according to a method described previously (62).

### Construction of plasmids

All the primers used in this study are listed in Table S6. All GFP fusion constructs were generated based on the plasmid backbone pYF3-GFP (63), in which StuI was used for the N-terminal fusions, and HpaI was used for the C-terminal fusions.

3'-RACE was conducted to validate the gene model of CENP-A, according to the manufacturer instruction (Invitrogen, Cat. no. 18373-019). All PCR-amplifications were performed using Phusion High-Fidelity DNA Polymerase (NEB, M0530S).

### CRISPR-mediated gene replacement

A sgRNA guide sequence whose PAM sequence overlapped with the start codon of CENP-A was selected as the CRISPR/Cas9 targets. An oligo annealing strategy was used for assembly of the sgRNA expression cassettes according to previously described methods (30).

17

387     HDR templates for CENP-A was assembled using NEBuilder® HiFi DNA Assembly. 5'-junction,

388     3'-junction and spanning diagnostic PCR were performed to genotype mutants, utilizing the

389     primers listed in Table S6.

390

391     **Microscopy imaging of *P. sojae* transformants**

392          A Zeiss 780 inverted confocal microscope was adopted to examine the subcellular

393     localization of GFP tagged CENP-A driven by strong promoters. Images were captured using a

394     63 X oil objective with excitation/emission settings (in nm) 488/504-550 for GFP, and 561/605-

395     650 for mCherry. DeltaVision elite deconvolution microscope (Olympus IX-71 base) equipped with

396     Coolsnap HQ2 high resolution CCD camera was employed to examine the subcellular localization

397     of GFP tagged CENP-A produced from the native loci. Images were captured using a 100 X oil

398     objective (100x/1.40 oil UPLSAPO100X0 1-U2B836 WD 120 micron DIC ∞/0.17/FN26.5, UIS2)

399     with an excitation filter, 475/28 and an emission filter, 525/50 for GFP. Time-lapse experiments

400     were performed with 40 X oil objective (40x/0.65-1.35 oil UAPO40XOI3/340 1-UB768R WD 100

401     micron DIC ∞/0.17/FN22, UIS2, BFP1), with the same filters. Confocal images were edited using

402     microscope's built-in Zen 2012 software (Blue and/or Black edition according to different

403     purposes). DeltaVision images were edited using Fiji-ImageJ and Photoshop.

404

405     **High molecular weight genomic DNA extraction and ONT sequencing**

406          High molecular weight (HMW) genomic DNA (gDNA) from *P. sojae* was isolated by the

407     CTAB DNA extraction method. 1 g 3-day old fresh *P. sojae* liquid cultures were collected by

408     filtration and washed twice with sterile water. The resulting damp mycelial pads were frozen

409     immediately in liquid nitrogen in a pre-cooled mortar, then ground by a pestle. Mycelial powder

410     was transferred to a 50 ml Falcon tube and mixed gently with 10 ml room temperature *P. sojae*

411     CTAB extraction buffer (200 mM Tris·HCl pH=8.5, 250 mM NaCl, 25 mM EDTA pH=8.0, 2% SDS,

412   1% CTAB). The suspension was incubated in 65℃ for 15 minutes with mixing every 5 minutes.

413   An equal volume of phenol/chloroform/isoamyl alcohol (25:24:1, saturated with 10 mM Tris,

414   pH=8.0 and 1 mM EDTA) was added to the suspension and mixed gently by inverting the tube,

415   then centrifuge at4℃, 5000 g for 15min. The supernatant was transferred to a new 50 ml tube and

416   treated with RNase A (final concentration, 100 µg/ml) at 37℃ for about 1 hour, followed by

417   proteinase K treatment (final concentration 200 µg/ml) at 50℃ for 2 hours. An equal volume of

418   chloroform was added to the solution and mixed gently by inverting the tube, then centrifuge, 4℃,

419   5000 g for 15min. The supernatant was transferred to a new 50 ml Falcon tube and DNA

420   precipitated by addition of an equal volume of isopropanol. The tube was mixed gently and

421   incubated on ice for 6 hours. The resulting white clump of DNA was spooled by a pipette tip and

422   washed once with 70% ethanol. The gDNA was air-dried for 15 minutes at room temperature and

423   dissolved in 100 µl sterile water. The quantity of DNA was examined by Qubit and the quality was

424   checked by pulse field gel electrophoresis (PFGE).

425      1D Genomic DNA by Ligation kits (SQK-LSK108, for MinION); SQK-LSK109, for GridION)

426   were used to prepare the Oxford Nanopore library. Oxford Nanopore sequencing runs was

427   performed on SpotON R9.4 flow cells with MinKNOW V1.11.5 using MinION or SpotON R9.4.1

428   flow cells with MinKNOW V3.1.20 using GridION. All of the GridION sequence were basecalled

429   (on GridION, in real time) using Guppy v2.0.5.

430

431   **Native ChIP-seq**

432      Native ChIP was performed according to the ChIP protocol accompanying Gent, Wang

433   and Dawe (64) with modifications. Briefly, 1-3 mg mycelia were collected from 1-1.5 L of ~3-day

434   culture by filtration system, and ground into fine powder in liquid nitrogen with pre-chilled mortars

435   and pestles. Nuclei were isolated and digested by micrococcal nuclease (M0247S, NEB) at 37 for

436   6 min. An antibody against GFP (Abcam, ab290) was used to immunoprecipitate single

19

437    nucleosomes containing the GFP-CENP-A fusion (driven by the strong promoter derived from

438    *HAM34* gene). Antibodies H3K9me3 (Abcam, ab8898), H3K27me3 (Active Motif, 39157), and

439    H3K4me2 (Millipore, 07-030) were used to immunoprecipitate nucleosomes with relevant

440    modifications. ChIP-seq of GFP-CENP-A and H3K27me3 were performed by Genewiz using

441    Illumina NextSeq500 that generated 150 nucleotide paired-end reads; ChIP-seq of H3K9me3 and

442    H3K4me2 were conducted by BGI using Illumina Hiseq 4000 that produced 50 nucleotide single-

443    end reads. Numbers of reads for each sample are listed in Table S2.

444

445    **Analysis of ChIP-seq and RNA-seq**

446        To map ChIP-seq reads to the genomes, the quality of raw ChIP-seq reads were first

447    assessed by FastQC (v0.11.6). For ChIP-seq of CENP-A, H3K27me3, the resulting reads were

448    trimmed with fastx-clipper and mapped with Bowtie2 with default parameters (65) and aligned to

449    the genome assemblies. For H3K9me3 and H3K4me2, the ChIP-seq reads were polished by BGI

450    prior to be released and thus mapped to the genomes directly using the same Bowtie2 setup. The

451    aligned file (.bam) was sorted and indexed by samtools (version 1.9). Subsequently the ChIP-ed

452    and input samples were analyzed with DeepTools(v3.2.0) "bamCompare" to calculate normalized

453    ChIP signal ($log2[ChIP_{RPKM}/Input_{RPKM}]$) and bigwig files were generated . Then .bw files were

454    visualized using the Integrative Genome Viewer (IGV). (https://software.

455    broadinstitute.org/software/igv/). To get profile mRNA, the existing RNA-Seq reads (FungiDB,

456    https://fungidb.org/fungidb/) were aligned to the genomes using HISAT2 (version 2.1.0), and the

457    resulting files (.bam) were sorted and indexed by samtools (version 1.9). The .bam file was

458    converted to .tdf for visualization using IGV.

459

460

461

**Genome assembly, analysis of genomic features and synteny comparison**

462  
463  Details of the *de novo* genome assembly is described in SI Text. To predict gene models,

464  first, the assembly Psojae2019.1 was subjected to repeat masking utilizing RepeatMasker (66)

465  based on a library of *de novo*-identified repeat consensus sequences that was generated by

466  RepeatModeler (www.repeatmasker.org/RepeatModeler.html). Next, the repeat-masked

467  assembly was used to predict gene models *ab initio* based on MAKER (v2.31.18) (67) with

468  predicted proteins from available *P. sojae* and *P. infestans* genome annotations as input (25, 68).

469  GC content was calculated in non-overlapping 5-kb windows using a modified Perl script

470  (gcSkew.pl, https://github.com/Geo-omics/scripts/blob/master/AssemblyTools/gcSkew.pl) and

471  plotted as the deviation from the genome average for each contig. Genes encoding ribosomal

472  RNA (18S, 5.8S, 25S, and 5S) and tRNA were inferred and annotated based on RNAmmer (v1.2)

473  (69) and tRNAscan-SE (v2.0) (70), respectively. To find telomeres, a custom-made Perl script

474  was used to search for the sequence "TTTAGGG" that was proposed for oomycetes telomeric

475  sequences (38). Pairwise synteny comparison between the two *P. sojae* genome assemblies (i.e.

476  P. sojae V3 and Psojae2019.1) or between different oomycete species was conducted using

477  BLASTn. BLASTn hits and other genomic features were plotted using Circos (v0.69-6) (71).

478  Whole-genome alignment was computed with MashMap (https://github.com/marbl/MashMap)

479  employing default settings, and was visualized as a dot plot (72).

480

**Bionano mapping**

482  *P. sojae* protoplasts were generated from 2.5-day old mycelial and were embedded into

483  agarose. Bionano Prep Cell Culture DNA Isolation Protocol was employed for extracting the high

484  molecular weight DNA. DNA labelling with DLE-1 was performed according to the standard

485  protocols provided by Bionano Genomics (Document number 30206, version F). Labelled DNA

486  samples were loaded into two flow cells and run on a Saphyr system (Bionano Genomics). The

487  *de novo* assembly was performed using Bionano Solve 3.3. Standard parameters for Saphyr data

488    were used without "extend and split" and without haplotype refinement in order to create a single

489    map for each allele ("optArguments_nonhaplotype_noES_DLE1_saphyr.xml"). In the process of

490    *de novo* assembly, data generated from two flow cells were merged. An assembly graph was

491    generated during a pairwise comparison of all of the molecules with a p value threshold of 1e-11,

492    and was refined based on molecules aligned to the assembled maps with a p value threshold of

493    1e-12. After five rounds of extension and refinement, a final refinement was conducted with a p

494    value threshold of 1e-16. Then, the *de novo* assembled map was used to scaffold the sequence

495    assembly. When using the hybrid scaffold module of Bionano Solve 3.3 pipeline, the option of

496    "resolve conflicts" for sequence contigs and Bionano maps was selected. The standard hybrid

497    scaffold settings with a modified parameter (-E 0) was applied to remove discrepancies between

498    sequence assembly and Bionano *de novo* assembly. Sequence contigs were *in silico* digested,

499    based on the recognition sequence (CTTAAG) of DLE-1. Conflicts detection was accomplished

500    by aligning contig maps to Bionano maps with p value threshold of 1e-10. When divergence was

501    identified, the conflicts were resolved by cutting either the contig or the map, depending on the

502    quality of the genome map at the divergent position.

503

504    **Analysis of transposable elements and identification of CoLT**

505    To identify transposable elements in *P. sojae*, the new genome assembly was subjected to

506    RepeatMasker (Repbase v23.09) analysis and hits were mapped to this genome assembly. The

507    Copia-like transposon (CoLT) element was identified in a stepwise way by multiple sequence

508    alignments followed by extraction of a consensus sequence and BLASTn analyses. Specifically,

509    an approximately 5 kb consensus sequence was identified in the alignment of centromere

510    sequences (including incompletely assembled ones) utilizing the multiple alignment program

511    MAFFT, a plug-in in the Geneious R9 software (http://www.geneious.com), with default

512    parameters. Then the consensus sequence was used as a query to perform a BLASTn search

513    against the Psojae2019.1 genome assembly. The resulting sequence hits were used to map

22

514    against the genome, and hits longer than 500 bp were used for representing in the figures. The

515    longest sequence hit with highest identity was retrieved, and was used as a query to execute a

516    second round of BLASTn search against the NCBI database to further characterize the

517    sequence. The results of BLASTn analysis indicated that that the sequence was highly similar to

518    a Copia-like transposable element. To define the domains of the CoLT, this sequence was

519    further analyzed by repeat identification (utilizing a bioinformatics software Unipro UGENE(73)),

520    and by searches utilizing the Repbase database (https://www.girinst.org/) and NCBI CD-search

521    (https://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi). Sequences of the 5 kb consensus and

522    the top hit in the Psojae2019.1 genome assembly are shown in Dataset S3.

523

524    **Prediction of centromeric regions in *P. sojae* closely related species**

525         To predict centromeres of the two oomycete species, namely *Phytophthora citricola*

526    P0716, (Genbank: GCA_007655245.1, with permission of the author) and *Bremia lactucae* SF5,

527    (GenBank: GCA_004359215.1) (36), BLASTn searches were conducted utilizing the *P. sojae*

528    Copia-like transposon (CoLT) as a query. Significant hits (>90% identity and > 500 bp) were

529    retrieved, and were plotted to all scaffolds of the *B. lactucae* assembly and to contigs > 10 kb of

530    the *P. citricola* assembly. For CoLT clusters that were localized within scaffolds or contigs, their

531    collinearities with the Psojae2019.1 assembly were further examined with BLASTn, and

532    visualized by Circos.

533

534    **Data availability**

535         All raw data of ChIP-seq and Nanopore sequencing and related processed files are

536    available in the NCBI under the BioProject PRJNA563922.

537

538 **Acknowledgments**

548 **References**

549 1.    L. E. Kursel, H. S. Malik, Centromeres. *Curr Biol* **26**, R487-R490 (2016).
550 2.    K. M. Stimpson, B. A. Sullivan, Epigenomics of centromere assembly and function. *Curr Opin Cell*
551       *Biol* **22**, 772-780 (2010).
552 3.    A. Buscaino, R. Allshire, A. Pidoux, Building centromeres: home sweet home or a nomadic
553       existence? *Curr Opin Genet Dev* **20**, 118-126 (2010).
554 4.    N. Wang, R. K. Dawe, Centromere size and its relationship to haploid formation in plants. *Mol*
555       *Plant* **11**, 398-406 (2018).
556 5.    B. E. Black *et al.*, Structural determinants for generating centromeric chromatin. *Nature* **430**,
557       578-582 (2004).
558 6.    A. Guse, C. W. Carroll, B. Moree, C. J. Fuller, A. F. Straight, In vitro centromere and kinetochore
559       assembly on defined chromatin templates. *Nature* **477**, 354-358 (2011).
560 7.    L. Comai, S. Maheshwari, M. P. A. Marimuthu, Plant centromeres. *Curr Opin Plant Biol* **36**, 158-
561       167 (2017).
562 8.    S. M. McNulty, B. A. Sullivan, Alpha satellite DNA biology: finding function in the recesses of the
563       genome. *Chromosome Res* **26**, 115-138 (2018).
564 9.    W. Jin *et al.*, Maize centromeres: organization and functional adaptation in the genetic
565       background of oat. *Plant Cell* **16**, 571-581 (2004).
566 10.   K. M. Smith, J. M. Galazka, P. A. Phatale, L. R. Connolly, M. Freitag, Centromeres of filamentous
567       fungi. *Chromosome Res* **20**, 635-656 (2012).
568 11.   V. Yadav, L. Sreekumar, K. Guin, K. Sanyal, Five pillars of centromeric chromatin in fungal
569       pathogens. *PLoS Pathog* **14**, e1007150 (2018).
570 12.   V. Yadav *et al.*, RNAi is a critical determinant of centromere evolution in closely related fungi.
571       *Proc Natl Acad Sci U S A* **115**, 3108-3113 (2018).
572 13.   K. M. Smith, P. A. Phatale, C. M. Sullivan, K. R. Pomraning, M. Freitag, Heterochromatin is
573       required for normal distribution of *Neurospora crassa* CenH3. *Mol Cell Biol* **31**, 2528-2542
574       (2011).
575 14.   K. Schotanus *et al.*, Histone modifications rather than the novel regional centromeres of
576       *Zymoseptoria tritici* distinguish core and accessory chromosomes. *Epigenetics & Chromatin* **8**, 41
577       (2015).
578 15.   K. Sanyal, M. Baum, J. Carbon, Centromeric DNA sequences in the pathogenic yeast *Candida*
579       *albicans* are all different and unique. *Proc Natl Acad Sci U S A* **101**, 11374-11379 (2004).
580 16.   W. A. Hoeijmakers *et al.*, *Plasmodium falciparum* centromeres display a unique epigenetic
581       makeup and cluster prior to and during schizogony. *Cell Microbiol* **14**, 1391-1401 (2012).
582 17.   R. E. Diner *et al.*, Diatom centromeres suggest a mechanism for nuclear DNA acquisition. *Proc*
583       *Natl Acad Sci U S A* **114**, E6015-E6024 (2017).
584 18.   C. F. Brooks *et al.*, *Toxoplasma gondii* sequesters centromeres to a specific nuclear region
585       throughout the cell cycle. *Proc Natl Acad Sci U S A* **108**, 3767-3772 (2011).
586 19.   V. Yadav *et al.*, Cellular dynamics and genomic identity of centromeres in cereal blast fungus.
587       *mBio* **10**, e01581-01519 (2019).
588 20.   M. I. Navarro-Mendoza *et al.*, Early diverging fungus *Mucor circinelloides* lacks centromeric
589       histone CENP-A and displays a mosaic of point and regional centromeres. *bioRxiv*, 706580
590       (2019).
591 21.   P. B. Becker *et al.*, Islands of retroelements are major components of *Drosophila* centromeres.
592       *PLOS Biology* **17** (2019).

25

593    22.    T. Lan *et al.*, Long-read sequencing uncovers the adaptive topography of a carnivorous plant
594           genome. *Proc Natl Acad Sci U S A* **114**, E4435-E4441 (2017).
595    23.    M. Jain *et al.*, Linear assembly of a human centromere on the Y chromosome. *Nat Biotechnol* **36**,
596           321-323 (2018).
597    24.    P. J. Keeling, F. Burki, Progress towards the Tree of Eukaryotes. *Curr Biol* **29**, R808-R817 (2019).
598    25.    B. M. Tyler *et al.*, *Phytophthora* genome sequences uncover evolutionary origins and
599           mechanisms of pathogenesis. *Science* **313**, 1261-1266 (2006).
600    26.    J. D. Grattepanche *et al.*, Microbial diversity in the eukaryotic SAR clade: Illuminating the
601           darkness between morphology and molecular data. *Bioessays* **40**, e1700198 (2018).
602    27.    D. C. Erwin, O. K. Ribeiro, *Phytophthora diseases worldwide* (American Phytopathological Society
603           (APS Press), 1996).
604    28.    R. H. Jiang, B. M. Tyler, Mechanisms and evolution of virulence in oomycetes. *Annu Rev
605           Phytopathol* **50**, 295-318 (2012).
606    29.    S. Savary *et al.*, The global burden of pathogens and pests on major food crops. *Nat Ecol Evol* **3**,
607           430-439 (2019).
608    30.    Y. Fang, L. Cui, B. Gu, F. Arredondo, B. M. Tyler, Efficient genome editing in the oomycete
609           *Phytophthora sojae* using CRISPR/Cas9. *Curr Protoc Microbiol* **44**, 21A-21 (2017).
610    31.    H. S. Judelson, M. D. Coffey, F. R. Arredondo, B. M. Tyler, Transformation of the oomycete
611           pathogen *Phytophthora-megasperma* f.sp. *glycinea* occurs by DNA integration into single or
612           multiple chromosomes. *Current Genetics* **23**, 211-218 (1993).
613    32.    Y. Fang, B. M. Tyler, Efficient disruption and replacement of an effector gene in the oomycete
614           *Phytophthora sojae* using CRISPR/Cas9. *Molecular Plant Pathology* **17**, 127-139 (2016).
615    33.    J. McGowan, D. A. Fitzpatrick, Genomic, network, and phylogenetic analysis of the oomycete
616           effector arsenal. *mSphere* **2**, e00408-00417 (2017).
617    34.    B. M. Tyler, M. Gijzen, "The *Phytophthora sojae* genome sequence: foundation for a revolution"
618           in Genomics of plant-associated fungi and oomycetes: dicot pathogens. (Springer, 2014), pp.
619           133-157.
620    35.    C. M. Malar *et al.*, Haplotype-phased genome assembly of virulent *Phytophthora ramorum*
621           isolate ND886 facilitated by long-read sequencing reveals effector polymorphisms and copy
622           number variation. *Mol Plant Microbe Interact* **32**, 1047-1060 (2019).
623    36.    K. Fletcher *et al.*, Genomic signatures of heterokaryosis in the oomycete pathogen *Bremia
624           lactucae*. *Nat Commun* **10**, 2645 (2019).
625    37.    J. J. van Hooff, E. Tromer, L. M. van Wijk, B. Snel, G. J. Kops, Evolutionary dynamics of the
626           kinetochore network in eukaryotes as revealed by comparative genomics. *EMBO Rep* **18**, 1559-
627           1571 (2017).
628    38.    J. Fulneckova *et al.*, A broad phylogenetic survey unveils the diversity and evolution of telomeres
629           in eukaryotes. *Genome Biol Evol* **5**, 468-483 (2013).
630    39.    B. A. Sullivan, G. H. Karpen, Centromeric chromatin exhibits a histone modification pattern that
631           is distinct from both euchromatin and heterochromatin. *Nature Structural & Molecular Biology*
632           **11**, 1076-1083 (2004).
633    40.    S. Henikoff, K. Ahmad, H. S. Malik, The centromere paradox: stable inheritance with rapidly
634           evolving DNA. *Science* **293**, 1098-1102 (2001).
635    41.    T. A. Volpe *et al.*, Regulation of heterochromatic silencing and histone H3 lysine-9 methylation
636           by RNAi. *Science* **297**, 1833-1837 (2002).
637    42.    X. Y. Li *et al.*, High-resolution mapping of epigenetic modifications of the rice genome uncovers
638           interplay between DNA methylation, histone methylation, and gene expression. *Plant Cell* **20**,
639           259-276 (2008).

640 43. R. C. Allshire, G. H. Karpen, Epigenetic regulation of centromeric chromatin: old dogs, new
641 tricks? *Nat Rev Genet* **9**, 923-937 (2008).
642 44. M. D. Blower, B. A. Sullivan, G. H. Karpen, Conserved organization of centromeric chromatin in
643 flies and humans. *Dev Cell* **2**, 319-330 (2002).
644 45. K. C. Scott, S. L. Merrett, H. F. Willard, A heterochromatin barrier partitions the fission yeast
645 centromere into discrete chromatin domains. *Curr Biol* **16**, 119-129 (2006).
646 46. R. C. Allshire, H. D. Madhani, Ten principles of heterochromatin formation and function. *Nat Rev*
647 *Mol Cell Biol* **19**, 229-244 (2018).
648 47. I. K. Greaves, D. Rangasamy, P. Ridgway, D. J. Tremethick, H2A.Z contributes to the unique 3D
649 structure of the centromere. *Proc Natl Acad Sci U S A* **104**, 525-530 (2007).
650 48. S. Friedman, M. Freitag, "Evolving centromeres and kinetochores" in Advances in genetics.
651 (Elsevier, 2017), vol. 98, pp. 1-41.
652 49. J. D. Brown, R. J. O'Neill, The evolution of centromeric DNA sequences. *eLS*  (2001).
653 50. P. W. Tooley, M. M. Carras, Separation of chromosomes of *Phytophthora* species using CHEF gel
654 electrophoresis. *Experimental Mycology* **16**, 188-196 (1992).
655 51. E. Sansome, C. Brasier, Polyploidy associated with varietal differentiation in the megasperma
656 complex of *Phytophthora*. *Transactions of the British Mycological Society* **63**, 461-IN411 (1974).
657 52. M. D. Cao *et al.*, Scaffolding and completing genome assemblies in real-time with nanopore
658 sequencing. *Nat Commun* **8**, 14515 (2017).
659 53. M. Boetzer, W. Pirovano, SSPACE-LongRead: scaffolding bacterial draft genomes using long read
660 sequence information. *BMC Bioinformatics* **15**, 211 (2014).
661 54. R. L. Warren *et al.*, LINKS: Scalable, alignment-free scaffolding of draft genomes with long reads.
662 *Gigascience* **4**, 35 (2015).
663 55. D. A. Compton, Mechanisms of aneuploidy. *Curr Opin Cell Biol* **23**, 109-113 (2011).
664 56. R. T. Todd, A. Forche, A. Selmecki, Ploidy Variation in Fungi: Polyploidy, Aneuploidy, and
665 Genome Evolution. *Microbiol Spectr* **5**, 599-618 (2017).
666 57. M. Elliott *et al.*, Characterization of phenotypic variation and genome aberrations observed
667 among *Phytophthora ramorum* isolates from diverse hosts. *BMC Genomics* **19**, 320 (2018).
668 58. M. P. Dobrowolski, I. C. Tommerup, H. D. Blakeman, P. A. O'Brien, Non-Mendelian inheritance
669 revealed in a genetic analysis of sexual progeny of *Phytophthora cinnamomi* with microsatellite
670 markers. *Fungal Genet Biol* **35**, 197-212 (2002).
671 59. T. Kasuga *et al.*, Host-induced aneuploidy and phenotypic diversification in the sudden oak
672 death pathogen *Phytophthora ramorum*. *BMC Genomics* **17**, 385 (2016).
673 60. T. van der Lee, A. Testa, A. Robold, J. van 't Klooster, F. Govers, High-density genetic linkage
674 maps of *Phytophthora infestans* reveal trisomic progeny and chromosomal rearrangements.
675 *Genetics* **167**, 1643-1661 (2004).
676 61. K. Lamour, S. Shresthra, Y. Zhou, X. Liu, J. Hu, Dynamic extreme aneuploidy (DEA) in the
677 vegetable pathogen *Phytophthora capsici* sheds light on instant evolution and intractability.
678 *bioRxiv*, 297788 (2018).
679 62. L. Lin *et al.*, The MADS-box transcription factor PsMAD1 is involved in zoosporogenesis and
680 pathogenesis of *Phytophthora sojae*. *Frontiers in Microbiology* **9** (2018).
681 63. Y. Fang, H. S. Jang, G. W. Watson, D. P. Wellappili, B. M. Tyler, Distinctive nuclear localization
682 signals in the oomycete *Phytophthora sojae*. *Front Microbiol* **8**, 10 (2017).
683 64. J. I. Gent, N. Wang, R. K. Dawe, Stable centromere positioning in diverse sequence contexts of
684 complex and satellite centromeres of maize and wild relatives. *Genome Biol* **18**, 121 (2017).
685 65. A. R. Quinlan, I. M. Hall, BEDTools: a flexible suite of utilities for comparing genomic features.
686 *Bioinformatics* **26**, 841-842 (2010).
687 66. A. Smit, R. Hubley, P. Green (2015) RepeatMasker Open-4.0. 2013–2015.

27

688   67.   B. L. Cantarel *et al.*, MAKER: An easy-to-use annotation pipeline designed for emerging model
689         organism genomes. *Genome Res.* **18**, 188-196 (2008).
690   68.   B. J. Haas *et al.*, Genome sequence and analysis of the Irish potato famine pathogen
691         *Phytophthora infestans*. *Nature* **461**, 393-398 (2009).
692   69.   K. Lagesen *et al.*, RNAmmer: consistent and rapid annotation of ribosomal RNA genes. *Nucleic*
693         *Acids Res* **35**, 3100-3108 (2007).
694   70.   T. M. Lowe, P. P. Chan, tRNAscan-SE On-line: integrating search and context for analysis of
695         transfer RNA genes. *Nucleic Acids Res* **44**, W54-57 (2016).
696   71.   M. Krzywinski *et al.*, Circos: an information aesthetic for comparative genomics. *Genome Res* **19**,
697         1639-1645 (2009).
698   72.   C. Jain, S. Koren, A. Dilthey, A. M. Phillippy, S. Aluru, A fast adaptive algorithm for computing
699         whole-genome homology maps. *Bioinformatics* **34**, i748-i756 (2018).
700   73.   K. Okonechnikov, O. Golosova, M. Fursov, U. team, Unipro UGENE: a unified bioinformatics
701         toolkit. *Bioinformatics* **28**, 1166-1167 (2012).
702   74.   S. Kumar, G. Stecher, M. Suleski, S. B. Hedges, TimeTree: A resource for timelines, timetrees,
703         and divergence times. *Mol Biol Evol* **34**, 1812-1819 (2017).

704

705

**Figures and tables**

**Figure 1.** Subcellular localization of CENP-A in *P. sojae* at different life stages and during vegetative growth. (A) A schematic showing the generation of *GFP*-fused *CENP-A* utilizing CRISPR/Cas9 mediated gene replacement. (B) Subcellular localization of GFP-tagged CENP-A (expressed from the endogenous locus) in *P. sojae* hyphae, sporangia, and encysted zoospores. (C) Time-lapse images illustrating localization of GFP tagged CENP-A during hyphal growth. Dashed squares denote occurrence of nuclear division. Representative images are shown. Scale bars in all images, 5 µm.

**Figure 2.** Contigs in the Psojae2019.1 assembly demonstrating CENP-A enrichment based on ChIP-seq. (A) 10 fully assembled CENP-A binding sites presented contigs. (B) 5 incompletely assembled CENP-A binding regions. All contigs are drawn to scale and the ruler indicates the length of the contigs. All CENP-A profiles shown were normalized to input DNA. mRNA profiles are shown as log-scales. Solid stars indicate the CENP-A enriched regions within contigs; hollow stars denote broken centromeres at the edge.

**Figure 3.** A representative Circos visualization comparing centromere-containing genomic regions between the Sanger V3 Scaffold 2 and the Psojae2019.1 assembly. The outer tracks illustrate assembled contigs (in Psojae2019.1) or scaffold (in P. sojae V3) and are color coded as listed in the key on the bottom. Yellow regions on the outer tracks indicate the locations of centromeres (CENP-A binding regions). Blue and orange lines link regions with collinearity across >2 kb, with orange lines corresponding to inversion. Grey box-shaded centromere-containing regions are magnified for detailed visualization.

730   **Figure 4.** *P. sojae* centromeres display heterochromatin marks and are enriched for a Copia-like

731   transposon (CoLT). (A) Schematics showing the *P. sojae* core centromeres (CENP-A binding

732   regions) and the pericentric regions of various lengths. Dark and light grey bars indicate core

733   centromeric and pericentric regions. Digits at the center indicate the size of core centromeres;

734   digits on the left denote the full length of the centromeres (a combination of core centromere and

735   pericentromeric region). The right pericentric region of *CEN5* and the left pericentric region of

736   *CEN10* are not fully assembled, and are indicated by dashed bars. Their full lengths labeled with

737   question marks. (B-C) Two centromeres (*CEN1* and *CEN4*) are shown as representatives to

738   compare CENP-A localization to the distributions of modified histones and CoLT elements. A 400

739   kb region harboring the centromeric region is shown for *CEN1* and *CEN4*. Cyan block, a

740   transcriptionally active region that interrupts *CEN4*. Profiles of CENP-A, H3K9me3, H3K29me3

741   and H3K4me2 shown were normalized to input. mRNA profiles are shown as log-scales.

742

743   **Figure 5.** Genomic distribution of CoLT in the *P. sojae* and *B. lactucae* genomes. (A) Location of

744   CoLT elements across all Psojae2019.1 contigs. (B) Diagram of a representative CoLT structure.

745   CoLT shown here represents the best hit of BLAST using the consensus sequence obtained from

746   alignment of the centromeres (Dataset S3). CoLT mainly contains two parts annotated (by

747   Repbase) as Copia-24_PIT-LTR, and Copia-24_PIT-I comprising gag, PR (protease), IN

748   (integrase), RT (reverse transcriptase) and RH (RNase H) domains. Other CoLT elements show

749   similar structure, except different lengths of LTR and other domains. (C) Location of CoLT

750   elements across all *B. lactucae* scaffolds >100 kb. For ease of analysis, scaffolds in *B. lactucae*

751   assembly were sorted and re-named based on sizes (large to small). See Dataset S4 for the

752   original scaffold names. (D) A representative Circos plot comparing *B. lactucae* Scaffold 2 that

753   has clustered CoLTs with the corresponding Psojae2019 contigs. Regions underscored by green

754   lines indicate both sides of the CoLT clusters were syntenic with the regions surrounding *P. sojae*

30

755 centromeres. Regions underlined by blue indicate only one side of the CoLT clusters were found

756 to be syntenic with *P. sojae* centromere flanking sequences.

757

758 **Figure 6.** Diversity of centromere features within the SAR supergroup. Simplified schematics (not

759 to scale) showing the structure, epigenetic modifications, size and composition of centromeres

760 across SAR lineages. *Epigenetic state was not examined in the diatom centromeric regions;

761 however, several AT-rich DNA sequence can be employed for episome maintenance, suggesting

762 diatom centromere might not be epigenetically dependent. Histone modification H3K27me2 was

763 only tested in *P. sojae*. The phylogeny was constructed using TimeTree (74). *Homo sapiens*,

764 *Arabidopsis thaliana* and *Neurospora crassa* were used as representatives of animals, plants and

765 for the phylogeny analysis, and are used as outgroups illustrating the evolutionary status of the

766 SAR supergroup.

767

768 **Table 1.** Centromeres identified in the Psojae2019.1 assembly and their counterparts in the
769 Sanger assembly

| | | Psojae2019.1 | | | Sanger V3 | |
|---|---|---|---|---|---|---|
| Name | Contig | Position of core centromere, kb (size, kb) | GC% of *CEN* | Position of pericentric region, kb (size, kb) | Scaffold | Position of *CEN* (kb) |
| *CEN1* | 1 | 415-579 (164) | 56.91 | 361-415 (54) 579-650 (71) | 2 | 7086-7267 |
| *CEN2* | 2 | 4102-4296 (194) | 55.51 | 4094-4102 (8) 4296-4305 (9) | 8 | 2374-2420 |
| *CEN3* | 3 | 3138-3412 (274) | 54.81 | 3223-3138 (15) 3412-3560 (48) | 9 | 3075-3286 |
| *CEN4** | 11 | 991-1175 (184) | 58.13 | 944-991 (47) 1175-1205 (30) | 4 | 995-1246 |
| *CEN5* | 18 | 1556-1709 (153) | 57.93 | 1492-1556 (64) 1709-? (>22)† | 3 | 4078-4138 |
| *CEN6* | 34 | 696-880 (184) | 57.65 | 643-696 (53) 882-904 (22) | 6 | 1521-1726 |
| CEN7 | 36 | 433-708 (275) | 52.16 | 376-433 (57) 708-732 (24) | 5 | 2049-2310 |
| *CEN8* | 38 | 302-490 (188) | 57.01 | 288-302 (14) 490-515 (25) | 1 | 9667-9688 |
| *CEN9* | 41 | 153-342 (189) | 57.40 | 101-153 (52) 342-358 (16) | 1 | 2921-3079 |
| *CEN10* | 51 | 36-152 (116) | 57.93 | ?-36 (>36)† 152-216 (64) | - | - |

770

771 *Contig11 contains a minor (coordinate, 990,422-1,009,603, 19 kb) and a major (coordinate, 1,060,108-1,174,601, 114 kb) peak that
772 are separated by a 21 kb transcriptionally active region.
773 †One side of pericentric heterochromatin region is not fully assembled.

774 **Supplemental Information**

775 **SI Text:** Nanopore sequencing and *de novo* assembly of the reference *P. sojae* genome
776 **Figs. S1.** Summary of the presence and absence of putative core kinetochore proteins identified
777 in *P. sojae*
778 **Fig. S2.** Identification and expression of CENP-A in *P. sojae*.
779 **Fig. S3.** Generation of *P. sojae* strains expressing *GFP* tagged *CENP-A* utilizing CRIPSR/Cas9
780 mediated genome editing.
781 **Fig. S4.** Scaffolds in the Sanger assembly that are suggested to harbor putative centromeres.
782 **Fig. S5.** Pipeline used for *de novo* assembly and metrics of the *P. sojae* genome assembly
783 Psojae2019.1.
784 **Fig. S6.** Comparison of centromere-containing genomic regions between the Sanger (P. sojae
785 V3) and the Psojae2019.1 assemblies.
786 **Fig. S7.** Summary of features of each intact centromere and read coverage analysis of
787 centromere.
788 **Fig. S8.** MAFFT-based alignment of CENP-A binding regions reveals a 5 kb sequence that are
789 highly similar among *P. sojae* centromeres.
790 **Fig. S9.** Genomic distribution of CoLT in the *P. citricola* genome.
791 **Fig. S10.** Representative contigs that are anchored by BioNano mapping and contigs that are
792 suggested to be joined.
793 **Fig. S11.** Dot plot comparison of scaffolded assemblies against the original Psojae2019.1
794 assembly and the Sanger assembly.
795 **Table S1.** Metrics of ONT sequencing
796 **Table S2**. Statistics of ChIP-seq samples
797 **Table S3.** Metrics of scaffolded assemblies and their comparison to the Sanger V3 and
798 Psojae2019.1 assembly.
799 **Table S4.** Five incompletely assembled centromeres in the Psojae2019 assembly and their
800 corresponding CENP-A regions mapped in the Sanger assembly
801 **Table S5.** *P. sojae* strains used in the study
802 **Table S6.** Primers used in this study.
803 **Movie S1 (separate file).** Time-lapse experiment showing cellular dynamics of CENP-A during
804 *P. sojae* vegetative growth.
805 **Dataset S1 (separate file)**. Sequences of kinetochore orthologs identified in *P. sojae.*
806 **Dataset S2 (separate file)**. 13 telomeres predicted in the Psojae2019.1 assembly.
807 **Dataset S3 (separate file)**. DNA sequence of the CoLT consensus sequence and the best hit
808 **Dataset S4 (separate file)**. Original names of the sorted *B. lactucae* scaffolds.
809 **Dataset S5 (separate file)**. Bionano mapping report.

Figure 1

Figure 2

A

**Contig 1**
(*CEN1*)

1 Mb   2 Mb   3 Mb   4 Mb   5 Mb   6 Mb

6.4 Mb

CENP-A    [1 - 6.0]

mRNA    [1 - 10955]

**Contig2**
(*CEN2*)    4.9 Mb

**Contig 3**
(*CEN3*)    3.6 Mb

**Contig 11**
(*CEN4*)    2.2 Mb

**Contig18**
(*CEN5*)    1.7 Mb

**Contig 34**
(*CEN6*)    0.9 Mb

**Contig 36**
(*CEN7*)    0.9 Mb

**Contig 38**
(*CEN8*)    0.5 Mb    0.8 Mb

**Contig 41**
(*CEN9*)    0.6 Mb

**Contig 51**
(*CEN10*)    0.4 Mb

B

**Contig 9**
(*CEN_C9*)    1 Mb   2 Mb    2.8 Mb

**Contig 10**
(*CEN_C10*)    2.3 Mb

**Contig 37**
(*CEN_C37*)    0.5 Mb    0.8 Mb

**Contig 48**
(*CEN_C48*)    0.5 Mb

**Contig 57**
(*CEN_C57*)    0.3 Mb

Figure 3

**Figure 3**

key:

(A) contigs: | Telomeric repeats | Centromeres (yellow) ··· assembly gaps

(D) GC content — red above / blue below genome average (5 kb non-overlapping window)

(B) | Copia-like TE (CoLT) enriched in the centromeres

(E) tRNA genes

(C) Transposable elements (from the outside inward):
LTR retrotransposons    DNA transposons    Other transposons

(F) BLASTn links (> 2 kb)

P. sojae V3    Psojae2019.1

Figure 4

**Figure 5**

# Figure 6



| | Species | Centromere | | | Ref. |
|---|---|---|---|---|---|
| | | Structure* | Length of CENP-A domain | Composition | |
| **SAR** | | | | | |
| **Stramenopiles** | *Phytophthora sojae* | | 116-275 kb | transposons no nucleotide bias | This study |
| | *Phaeodactylum tricornutum* | | 2-6 kb | simple AT-rich sequence | 17 |
| **Alveolates** | *Plasmodium falciparum* | | 4-5 kb | simple AT-rich sequence | 16 |
| | *Toxoplasma gondii* | | 10 -20 kb | simple sequence no nucleotide bias | 18 |
| | Fungi | | | | |
| | Animals | | | | |
| | Plants | | | | |

Time (MYA)

1768 1600 1200 800 400 0

**Keys:**

Centromere structure
- ▮ Core centromere
- ▮ Pericentric region
- ||||| Transposable elements

Epigenetic state
- 🟠 CENP-A
- ⚫ H3K9me2/3
- 🟣 H3K27me2/3