

1 ***Plectus sambesii* assembly methods: Supplementary Note 1**

2
3 We acquired ~250M 150bp paired-end genome sequencing read pairs with 550bp
4 insert size through Edinburgh Genomics. Quality control of the genome sequencing
5 data was performed using Fastqc v0.10.1 (Andrews, 2010), and reads were quality
6 trimmed using skewer v0.2.2 (Jiang et al., 2014) with parameters `-q 30 -l 51`. K-mer
7 plots were generated with kmc, which revealed extensive heterozygosity. A
8 preliminary single-end assembly was generated with Velvet (Zerbino and Birney,
9 2008) to screen for contaminants but no significant hits were found using Blobtools
10 (Kumar et al., 2013; Laetsch, 2016). The dataset was digitally normalised to 40x
11 coverage and assembly was carried out with SPAdes v3.5.1 (Bankevich et al., 2012)
12 with parameters `--k 21,33,55,77,99 --cov-cutoff auto --careful`. Error correction of the
13 sequencing reads was done with BayesHammer that runs as part of the SPAdes
14 pipeline. Examination of the resulting distribution of contig coverages revealed a
15 bimodal distribution of coverage indicating again heterozygosity. Haplocontigs were
16 collapsed and postprocessed with Redundans (Pryszcz and Gabaldón, 2016), that
17 runs SSPACE3 (Boetzer et al., 2011) and SOAPdenovo Gapcloser internally. The
18 identity percentage for the collapse of redundant contigs was chosen to be of 90%,
19 ensuring the presence of 95.97% of CEGMA KOGs (Parra et al., 2007). Lower
20 identity percentage values did not result in an appreciable reduction of the span
21 suggesting that most heterozygous regions are collapsed at that threshold. Blobplot
22 analysis after postprocessing revealed ~20Mb of bacterial contaminant sequence at
23 60% GC and lower coverage than the nematode contigs, which we therefore
24 excluded from the final assembly.

25
26 To annotate genes, we first masked repeats in the assembly. We used
27 RepeatModeler v1.0.8 (Smit and Hubley, 2008) to generate a species-specific repeat
28 library that was concatenated with a nematode repeat library extracted from
29 RepBase RepeatMasker edition update 27-01-2017. The combined library was used
30 to mask the genome with RepeatMasker v4.0.7 (Smit et al., 2013). To annotate the
31 genome, we acquired ~100M RNA-seq reads from rRNA-depleted paired-end RNA-
32 seq libraries generated in the LMS Genomics Facility. Quality control of the RNA-seq
33 reads was performed with Fastqc and Skewer with parameters `-q 30 -l 50`. Spliced
34 alignment of the reads to the masked assembly was done with STAR v2.5.2 (Dobin
35 et al., 2013) with default parameters except for `--twopassMode Basic`. RNA-seq
36 alignments were used for automated HMM training and annotation using BRAKER
37 v1.9 (Hoff et al., 2016) with parameters `--softmasking 1 --UTR off --gff3`, running
38 internally GeneMark-ET v4.21 (Lomsadze et al., 2014) and augustus v2.3.2 (Stanke
39 and Morgenstern, 2005).

41 **References**

- 42
43
44 Andrews, S. (2010). FastQC A Quality Control tool for High Throughput Sequence
45 Data. [Http://www.bioinformatics.babraham.ac.uk/projects/fastqc/](http://www.bioinformatics.babraham.ac.uk/projects/fastqc/).
46 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. a., Dvorkin, M., Kulikov, A.S.,
47 Lesin, V.M., Nikolenko, S.I., Pham, S., Pribelski, A.D., et al. (2012). SPAdes: A New
48 Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. *J.*
49 *Comput. Biol.* 19, 455–477.
50 Boetzer, M., Henkel, C. V, Jansen, H.J., Butler, D., and Pirovano, W. (2011).

51 Scaffolding pre-assembled contigs using SSPACE. *Bioinformatics* 27, 578–579.
52 Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P.,
53 Chaisson, M., and Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner.
54 *Bioinformatics* 29, 15–21.
55 Hoff, K.J., Lange, S., Lomsadze, A., Borodovsky, M., and Stanke, M. (2016).
56 BRAKER1: Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-ET
57 and AUGUSTUS. *Bioinformatics* 32, 767–769.
58 Jiang, H., Lei, R., Ding, S.-W., and Zhu, S. (2014). Skewer: a fast and accurate
59 adapter trimmer for next-generation sequencing paired-end reads. *BMC*
60 *Bioinformatics* 15, 182.
61 Kumar, S., Jones, M., Koutsovoulos, G., Clarke, M., and Blaxter, M. (2013).
62 Blobology: exploring raw genome data for contaminants, symbionts and parasites
63 using taxon-annotated GC-coverage plots. *Front. Genet.* 4, 237.
64 Laetsch, D. (2016). Blobtools. GitHub Repos.
65 Lomsadze, A., Burns, P.D., and Borodovsky, M. (2014). Integration of mapped RNA-
66 Seq reads into automatic training of eukaryotic gene finding algorithm. *Nucleic Acids*
67 *Res.* 42, e119.
68 Parra, G., Bradnam, K., and Korf, I. (2007). CEGMA: A pipeline to accurately
69 annotate core genes in eukaryotic genomes. *Bioinformatics* 23, 1061–1067.
70 Prysycz, L.P., and Gabaldón, T. (2016). Redundans: an assembly pipeline for highly
71 heterozygous genomes. *Nucleic Acids Res.* 44, e113–e113.
72 Smit, A., and Hubley, R. (2008). RepeatModeler Open-1.0.
73 Smit, A., Hubley, R., and Green, P. (2013). RepeatMasker Open-4.0.
74 Stanke, M., and Morgenstern, B. (2005). AUGUSTUS: A web server for gene
75 prediction in eukaryotes that allows user-defined constraints. *Nucleic Acids Res.* 33,
76 465–467.
77 Zerbino, D.R., and Birney, E. (2008). Velvet: Algorithms for de novo short read
78 assembly using de Bruijn graphs. *Genome Res.* 18, 821–829.
79

1 **Supplementary Note 2: Phylogeny References**

2 **Protozoa**

3 Sierra, R., Matz, M. V., Aglyamova, G., Pillet, L., Decelle, J., Not, F., et al. (2013).
4 Deep relationships of Rhizaria revealed by phylogenomics: A farewell to Haeckel's
5 Radiolaria, *Molecular Phylogenetics and Evolution*, 67(1), 53–59.

6 <http://doi.org/10.1016/j.ympev.2012.12.011>

7

8 Cavalier-Smith, T., Fiore-Donno, A. M., Chao, E., Kudryavtsev, A., Berney, C.,
9 Snell, E. A., & Lewis, R. (2015). Multigene phylogeny resolves deep branching of the
10 Amoebozoa. *Molecular Phylogenetics and Evolution*, 83(C), 293–304.

11 <http://doi.org/10.1016/j.ympev.2014.08.011>

12

13 Derelle, R., Lopez-Garcia, P. N., Timpano, H. L. N., & Moreira, D. (2016). A
14 Phylogenomic Framework to Study the Diversity and Evolution of Stramenopiles
15 (=Heterokonts). *Molecular Biology and Evolution*, 33(11), 2890–2898.

16 <http://doi.org/10.1093/molbev/msw168>

17

18 Jinkerson, R. E., Radakovits, R., & Posewitz, M. C. (2014). Genomic insights from
19 the oleaginous model alga *Nannochloropsis gaditana*. *Bioengineered*, 4(1), 37–43.

20 <http://doi.org/10.4155/bfs.11.7>

21

22 Fritz-Laylin, L. K., Prochnik, S. E., Ginger, M. L., Dacks, J. B., Carpenter, M. L.,
23 Field, M. C., et al. (2010). The Genome of *Naegleria gruberi* Illuminates Early
24 Eukaryotic Versatility. *Cell*, 140(5), 631–642.

25 <http://doi.org/10.1016/j.cell.2010.01.032>

26

27 Read, B. A., Kegel, J., Klute, M. J., Kuo, A., Lefebvre, S. C., Maumus, F., et al.
28 (2014). Pan genome of the phytoplankton *Emiliana underpins* its global distribution.

29 *Nature*, 499(7457), 209–213. <http://doi.org/10.1038/nature12221>

30

31 **Microsporidia**

32

33 Vossbrinck, C.R., Debrunner-Vossbrinck, B.A., & Weiss, L.M., 2014 Phylogeny of
34 the Microsporidia (2014). In *Microsporidia, Pathogens of opportunity*, Edited Weiss,
35 L and Becnel, J, Wiley.

36

37 **Fungi**

38

39 Sugiyama, I., Hosaka, K., Suh S.O (2006). Early diverging Ascomycota: phylogenetic
40 divergence and related evolutionary enigmas. *Mycologia*, 98(6), 996-1005.

41

42 Hibbert D.S., (2006). A phylogenetic overview of the Agaricomycotina, *Mycologia*.
43 98(6):917-25.

44

45 James, T. Y., Kauff, F., Schoch, C. L., Matheny, P. B., Hofstetter, V. R., Cox, C. J., et
46 al. (2006). Reconstructing the early evolution of Fungi using a six-gene phylogeny.

47 *Nature*, 443(7113), 818–822. <http://doi.org/10.1038/nature05110>

48

49 Masneuf-Pomarede, I., Bely, M., Marullo, P., & Albertin, W. (2016). The Genetics of
50 Non-conventional Wine Yeasts: Current Knowledge and Future Challenges. *Frontiers*
51 *in Microbiology*, 6(1569), 702. <http://doi.org/10.1016/j.ijfoodmicro.2011.08.026>

Supplementary Note 3: sample numbers for genome-wide comparisons

CGs covered by >10 reads used for significance testing in Figure 1c

Plectus sambesii

gene 1230239

rRNA 1193

LTR 7069

DNA 1599

unannotated 727548

Romanomermis culicivorax

gene 656710

rRNA 561

LTR 4897

DNA 1830

unannotated 124399

Trichuris muris

gene 2115909

rRNA 3183

LTR 79067

DNA 51875

unannotated 1247125

Trichinella spiralis

gene 738591

rRNA 291

LTR 21480

DNA 3540

unannotated 193453

Caenorhabditis briggsae

gene 82526

rRNA 708

LTR 590

DNA 3725

unannotated 59570

Nippostrongylus brasiliensis

Gene 20765

rRNA 926

DNA 1199

LTR 2879

unannotated 12940

Analysed genomes: number of genome-wide Cs and CGs

P. sambesii

total_C	29615596
total_CG	7728644
total_CGA	2355653
total_CGC	2169535
total_CGG	1402387
total_CGT	1793453

R. culicivorax

total_C	48473942
total_CG	7829396
total_CGA	2785794
total_CGC	1457635
total_CGG	1529381
total_CGT	2050824

T. spiralis

total_C	9943967
total_CG	1710706
total_CGA	539942
total_CGC	348898
total_CGG	309852
total_CGT	511933

T. muris

total_C	18839135
total_CG	4505970
total_CGA	1184822
total_CGC	1125042
total_CGG	921658
total_CGT	1274273

N. brasiliensis

total_C	58220992
total_CG	12740122
total_CGA	4131434
total_CGC	2803568
total_CGG	2582339
total_CGT	3221088

C. briggsae

total_C	19684666
total_CG	3455622
total_CGA	1266130
total_CGC	588973
total_CGG	765446
total_CGT	834762

B. mori

total_C	79039767
total_CG	16946963
total_CGA	4852595
total_CGC	3693029
total_CGG	3551168
total_CGT	4848034

M. musculus

total_C	38179275
total_CG	1428189
total_CGA	354635
total_CGC	299758
total_CGG	374931
total_CGT	398865

Analysed genomes: CGs within ± 1000 bp of TSS with sufficient coverage to estimate methylation (see Supplemental Figure 3)

P. sambesii 2002411
*R. culicivora*x 6231622
T. spiralis 1055855
T. muris 1632567

CGs in 1st exon vs introns analysed (Supplemental Figure 3)

P. sambesii 480833e,4352323i
*R. culicivora*x 2506380e,5629056i
T. spiralis 502608e,1992647i
T. muris 287818e,2699672i

total number of genes analysed with DNA methylation coverage for repeat homology (Supplemental Figure 2)

P. sambesii 38467
*R. culicivora*x 43622
T. spiralis 12913
T. muris 10932