# GigaScience
## Data for "A Phylogenomic View of Evolutionary Complexity in Green Plants"
### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-19-00241 |
| Full Title: | Data for "A Phylogenomic View of Evolutionary Complexity in Green Plants" |
| Article Type: | Data Note |
| Funding Information: | Alberta Innovates - Technology Futures (RES0010334) — Prof Gane Ka-Shu Wong |

| | |
|---|---|
| Abstract: | The 1000 Plants (1KP) initiative explored the genetic diversity of green plants (Viridiplantae) by sequencing RNA from 1,342 samples representing 1,173 species. All of the analyses done for the 1KP capstone, and previous studies on subsets of these data, are based on a series of de novo transcriptome assemblies and related outputs that will be described in this publication. We also describe assessments of the data quality and an analysis to remove cross-contamination between the samples. These data will be useful to researchers with interests in specific gene families, either across the green plant tree of life or in more focused lineages. |

| | |
|---|---|
| Corresponding Author: | Gane Ka-Shu Wong<br><br>CANADA |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | |
| Corresponding Author's Secondary Institution: | |
| First Author: | Eric J. Carpenter |
| First Author Secondary Information: | |
| Order of Authors: | Eric J. Carpenter |
| | Gane Ka-Shu Wong |
| Order of Authors Secondary Information: | |

| | |
|---|---|
| Additional Information: | |

| Question | Response |
|---|---|
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |

| | |
|---|---|
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | No |
| If not, please give reasons for any omissions below.<br><br>as follow-up to "**Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?<br><br>" | The data is derived from plant samples for which no attempt was made to identify an age or sex for the source. |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using | No |

| | |
|---|---|
| a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br><br>Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)? | |
| If not, please give reasons for any omissions below.<br><br><br>   as follow-up to **"Availability of data and materials**<br><br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in [publicly available repositories](#) (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br><br>Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?<br><br>**"** | Additional data (contamination analysis, etc) will be submitted to GigaDB after this online process, as per the journal instructions. |

1    Title: The Data for One Thousand Plant Transcriptomes Initiative: A Phylogenomic View of

2    Evolutionary Complexity in Green Plants

3

4    Authors:

5

6    Eric J. Carpenter <ejc@ualberta.ca> [1]

7    Naim Matasci <nmatasci@usc.edu> [2]

8    Saravanaraj Ayyampalayam <raj@plantbio.uga.edu> [3]

9    Shuangxiu Wu <wushx@big.ac.cn> [4]

10   Jing Sun <jsun@genetics.ac.cn> [4]

11   Jun Yu <junyu@big.ac.cn> [4]

12   Fabio Rocha Jimenez Vieira <rocha@biologie.ens.fr> [5]

13   Chris Bowler <cbowler@biologie.ens.fr> [5]

14   Richard G. Dorrell <dorrell@biologie.ens.fr> [5]

15   Matthew A. Gitzendanner <magitz@ufl.edu> [6]

16   Ling Li <liling3@cngb.org> [7]

17   Wensi Du <duwensi@cngb.org> [7]

18   Kristian Ullrich <ullrich@evolbio.mpg.de> [8]

19   Norm J. Wickett <norman.wickett@gmail.com> [9]

20   Todd J. Barkmann <todd.barkman@wmich.edu> [10]

21   Michael S. Barker <msbarker@email.arizona.edu> [11]

22   James H. Leebens-Mack <jleebensmack@uga.edu> [3]

23   Gane Ka-Shu Wong <gane@ualberta.ca>* contact author [1,7,12]

24

25   1. Department of Biological Sciences, University of Alberta, Edmonton, Alberta, T6G 2E9, Canada.

26    2. CyVerse, University of Arizona, Arizona, U.S.A.; Current address: Lawrence J. Ellison Institute for

27    Transformative Medicine, University of Southern California, Los Angeles, CA 90033, U.S.A.

28    3. Department of Plant Biology, University of Georgia, Athens, GA 30602, USA.

29    4. CAS Key Laboratory of Genome Sciences and Information, Beijing, Institute of Genomics, Chinese

30    Academy of Sciences, Beijing 100101, People's Republic of China.

31    5. Institut de Biologie de l'ENS (IBENS), Département de biologie, École normale supérieure, CNRS,

32    INSERM, Université PSL, 75005 Paris, France

33    6. Department of Biology, University of Florida, Gainesville, Florida 32611, USA.

34    7. BGI-Shenzhen, Beishan Industrial Zone, Yantian District, Shenzhen 518083, People's Republic of

35    China.

36    8. Department of Evolutionary Genetics, Max Planck Institute for Evolutionary Biology, Plön,

37    Germany.

38    9. Chicago Botanic Garden, Glencoe, IL 60022, and Program in Biological Sciences, Northwestern

39    University, Evanston, IL 60208 USA.

40    10. Department of Biological Sciences, Western Michigan University, Kalamazoo MI 49008-5410

41    USA.

42    11. Department of Ecology & Evolutionary Biology, University of Arizona, Tucson, AZ 85721 USA.

43    12.  Department of Medicine, University of Alberta, Edmonton, Alberta, T6G 2E1, Canada.

44

45    **Abstract**

46

47    The 1000 Plants (1KP) initiative explored the genetic diversity of green plants (Viridiplantae) by

48    sequencing RNA from 1,342 samples representing 1,173 species. All of the analyses done for the 1KP

49    capstone, and previous studies on subsets of these data, are based on a series of de novo transcriptome

50    assemblies and related outputs that will be described in this publication. We also describe assessments

51    of the data quality and an analysis to remove cross-contamination between the samples. These data will

52 be useful to researchers with interests in specific gene families, either across the green plant tree of life

53 or in more focused lineages.

54

55

56 **Keywords**

57

58 RNA, plants, assemblies, genes, contamination, completeness

59

60

61

62 **Data Description**

63

64 1KP has sequenced RNA from 1,342 RNA samples of 1,173 green plant species representing all major

65 taxa within the Viridiplantae, including streptophyte and chlorophyte green algae, bryophytes, ferns,

66 angiosperms, and gymnosperms.  Importantly, our selection criteria eschewed the model organisms and

67 crop species where other plant sequencing efforts have historically been concentrated.

68

69 Major papers describing the project have been published elsewhere [1,2].  This Data Note describes the

70 sequence data set and provides additional details on the sample and sequence processing as well as

71 quality assessments of these data.

72

73 **Methods**

74

75 Sampling strategy

76

77    Because of the diversity and the number of species analyzed, no one source could be used.  Samples

78    were provided by a global network of collaborators who obtained materials from a variety of sources,

79    including field collection of wild plants, greenhouses, botanical gardens, laboratory specimens, and

80    algal culture collections.  To ensure an abundance of expressed genes, we preferred live growing cells,

81    e.g. young leaves, flowers, or shoots, although many samples were also from roots, or other tissues.

82    Because of the sample diversity, we did not attempt to define specific standards on growth conditions,

83    time of collection, or age of tissue.  For more details, see the supplemental methods in the capstone

84    paper [1].

85

86    RNA extraction

87

88    Given the biochemical diversity of these samples, no one RNA extraction protocol was appropriate for

89    all samples.  Most samples were extracted using commonly known protocols or using commercial kits.

90    For complete details of the many specific protocols used, please see Appendix S1 of Johnson et al. [3]

91    and Jordon-Thaden et al. [4].  Depending on the sample, RNA extractions might have been done by the

92    sample provider, a collaborator near the provider, or the sequencing lab (BGI-Shenzhen).

93

94

95    Sequencing at BGI

96

97    Samples of extracted RNA or frozen tissues were sent to the sequencing lab, BGI-Shenzhen.  Prior to

98    library construction, RNA samples were screened by Agilent Bioanalyzer RIN scores [5] and basic

99    photometry; obvious low-quality outliers (e.g., RIN scores less than 6 and/or loss of distinct

100    electropherogram peaks) were excluded.  Libraries for Illumina sequencing were constructed using

101    Illumina's standard procedure.  Some samples for which only a small amount of RNA was available

102    were processed using TruSeq kits.

103

104 Initially, sequencing was done on the Illumina GAII platform, but later samples were run on the HiSeq

105 platform. Associated with this change was a shift from ~72 bp read lengths to 90 bp read lengths (both

106 cases paired-end). Libraries were indexed and multiplexed in the sequencer lanes to a target

107 sequencing depth of 2 Gbp per sample. Average depth achieved was 1.99 Gbp of sequence of better

108 than Phred quality 30 (1 error per thousand bases).

109

| Percentile | Dataset Size (all base qualities) |
|---|---|
| 5th | 1.3 Gbp |
| 25th | 1.9 Gbp |
| 50th | 2.2 Gbp |
| 75th | 2.5 Gbp |
| 95th | 3.0 Gbp |

110

111

112 The data was cleaned by eliminating reads with excessive adapter-primer sequences or high numbers of

113 low quality bases (i.e. more than half of Phred quality 5 or lower [32 % error rate] or more than 10%

114 uncalled).

115

116

117 *De novo* assembly

118

119 Quality filtered reads were assembled using the SOAPdenovo-Trans transcript assembler (version

120 2012-04-05) [6]. No additional pre-processing of the data was performed. This largely used the

121 program defaults, with the slight modification of increasing the *k*-mer length to 25 bp and reducing the

122 number of processor threads to one. This reduced thread count allowed us to more efficiently use our

123 computer resources. Both the internal FillGap module and the external GapCloser post-processor

124 (supplied with SOAPdenovo-Trans) were run. An example of the commands used for one of the

125 assemblies (dataset AEPI):

126

```
127  SOAPdenovo-Trans-31kmer all -s config -p 1 -K 25 -e 2 -F -L 100 -t 5 -o AEPI

128  GapCloser -a AEPI.scafSeq -b config -o AEPI.GapCloser.fa -l 100 -p 25 -t 1
```

129

130 These commands refer to a configuration file named config, which specified the expected insert size,

131 maximum read length, and read-sequence filenames. The contents of this file were:

132

```
133  max_rd_len=120

134  [LIB]

135  avg_ins=200

136  rank=1

137  q1=AEPI-read_1.fq

138  q2=AEPI-read_2.fq
```

139

140 When multiple samples from the same species were co-assembled, the last five lines were repeated for

141 each data source with the appropriate filenames. See the supplemental files in the accompanying

142 analysis paper [1] and protocols in protocols.io for more details [cite].

143

144 Protein translation

145

146 To identify likely proteins within the assembled transcripts, sequences were passed through TransPipe

147 [7], which identified reading frames and protein translations by comparison to protein sequences from

148 22 sequenced and annotated plant genomes in Phytozome [8]. Using BLASTX [9], best hit proteins

149     were paired with each assembled scaffold at a threshold of 1E-10 expectation-value and a minimum

150     length of 100 amino acid residues.  Scaffolds that did not have a best hit protein at this level were

151     removed.  To determine reading frames and estimate amino acid sequences, each gene is aligned

152     against its best hit protein by Genewise 2.2.0 [10].  Using the highest scoring Genewise DNA-protein

153     alignments, stop codons and  those containing ambiguous nucleotides were removed to produce an

154     amino acid sequence for each gene.  Outputs include paired DNA and protein sequences.

155

156

157     BLAST searches

158

159     Thanks to the support of China National GeneBank (CNGB), a BLAST search service

160     (http://db.cngb.org/onekp/) allows public searches against the assemblies and protein translations.

161     CNGB developed the service using NCBI BLAST+ (version 2.6.0) [11].  It integrates all public

162     datasets from CNGB applications, BGI projects and external data sources, and provides a

163     comprehensive and convenient sequence searching.  A specialized interface for BLAST searching the

164     1KP dataset allows limiting the search to specific families, orders, or 25 higher-level clades.  For

165     assemblies, there are 21,398,790 nucleotide sequences, 6,188,419,272 bases in total. And for the

166     Transpipe protein translations, there are 103 million protein sequences comprising over 47 billion

167     amino acids in total.

168

169

170     **Validation**

171

172     Purity and contamination

173

174    High throughput sequencing methods are always at risk of contamination, as even a 1 ppm contaminant

175    produces multiple reads.  In practice, data has been found to often include sequences best attributed to

176    additional contaminating sources [12].  For 1KP, the diversity of sources for the samples, and

177    especially the fact that axenic cultures are not a viable option in most instances, ensures that there will

178    always be some contamination of the plant tissue by other environmental nucleic acids.  These can

179    reasonably be expected to include bacterial, fungal, and insect species that live in and on the plant

180    tissues, and more rarely, from contact with larger species such as frogs, mice, birds and humans.

181

182    For most analyses, these minor contaminants are not expected to matter, as only the most abundant of

183    such contaminants will be present in sufficient quantities to assemble. In many cases, they are also

184    sufficiently diverged from the intended species that they can be easily recognised as non-plant genes.

185    Unfortunately, this is not always the case.  Some analyses are further protected by looking at the whole

186    of the available transcriptome, whereby the many genes from the target species will overpower a few

187    contaminants.  Single gene family analyses do not have this advantage and must rely on other methods

188    to reject non-plant genes.

189

190    Another possibility is significant contamination during sample processing when plant RNA is

191    transferred between adjacent samples, or when whole samples are accidentally mis-labeled.

192

193    We tried to guard against these problems by several analyses, one of which compared the assembled

194    sequences by BLASTn to a reference set of nuclear 18S rRNA sequences from the SILVA SSU rRNA

195    database (http://www.arb-silva.de) [13].  The BLASTn alignment to an assembly with the lowest

196    expectation-value is taken to indicate the assembly has a similar taxonomic origin as the reference

197    sequence.  However, alignments of less than 300 bp or expectation-values above 1E-9 often align to

198    several distantly related species and were ignored.

199

200  For most samples we found an 18S sequence most-similar to a SILVA sequence from the same

201  taxonomic family as the expected sample species.  This is not true for all our samples, and may indicate

202  a failure to assemble the 18S sequence, limitations in the taxonomic identification from the BLASTn

203  results, or mis-labelling of sample.  In a few cases, additional (and possibly contaminant) 18S

204  sequences were found.  Because the 18S rRNA sequence is highly expressed, we expect that this

205  method is likely to be sensitive to low levels of contamination.  In a few cases, the taxonomic

206  irregularities were judged sufficiently severe that samples were excluded from various analyses.

207

208  The accompanying data includes two accessory files containing details of this SILVA based SSU

209  validation for each sample.  The first lists whether the sample is overall judged to be validated as

210  containing the expected taxon, and whether it had alignments to any other plant sequences (described

211  as "worrisome contamination").  The second file, more detailed, lists each scaffold identified as being

212  18S-like sequence, and which reference sequence it matched against.

213

214

215  Pairwise Cross-contamination of Assemblies

216

217  Cross contamination between the datasets was identified by using a genome-scale sequence search

218  pipeline, adapted from previous studies [14-16]. Briefly, each pair of assemblies (nucleotide) was

219  compared and a threshold identity level established, above which sequences are likely to be

220  contamination between the pair.  While best for identifying technical contamination between libraries

221  (e.g. due to mixing of RNA samples), this technique could also detect other biological contamination

222  events (e.g. contamination of pairs of libraries with common commensal organisms).  An additional

223  search step, using the entire 1KP sequence library, identified the probable evolutionary origin of each

224  sequences.

225

226    The pair-wise comparison used LAST v. 963 [17] using the --cR01 option, and the respective matches

227    were grouped and ordered by similarity. To avoid artifactually excluding sequences between closely

228    related species, which may have very high degrees of similarity [13], pairs of libraries from the same

229    family, along with pairs of libraries separated by two or fewer branches in the consensus 1kp multigene

230    phylogeny, were excluded from the searches [2].

231

232    The expected distribution of the matched sequence identities has a maximum at the pairwise identity

233    reflecting the evolutionary distance between the two species [15, 16].  In contrast, a cross-contaminated

234    pair should contain  many sequences of near 100% similarity, and the similarity value which has the

235    first minimum number of sequences below this level (i.e. the first inflexion point in a curve plotting the

236    total number of sequences of each percentage similarity value) can be used as a  threshold for

237    discriminating contaminating sequences [15, 16].  The code is available at https://github.com/Plant-

238    and-diatom-genomics-IBENS-Paris/Decontamination-pipeline.

239

240    The output of this analysis is pairs of apparent orthologs whose sequence similarities are higher than

241    the cut-off in one or both libraries, i.e. potential contamination.  To discriminate donors and recipients

242    in each contaminant pair, each of these potential contaminants was searched against all the non-

243    contaminant assemblies by BLASTn, using the option -max_target_seqs 3 [18].  Queries with at least

244    one of the three best alignments against a sequence from the same family, or from a taxon separated by

245    fewer than two branches within the 1kp tree [2], were excluded from the list of potential contaminants;

246    whereas sequences that yielded best hits exclusively against more distantly related taxa, were verified

247    as potential contaminants. Clean and contaminant FASTA sequence files for each library are available

248    in the accompanying data.

249

250    An overview of the results is presented in Fig. 1.  In total, we identified 79,175 nucleotide sequences

251    (0.3 %) of a total 23,436,405 searched as being clearly of contaminant origin (Fig. 1A). A further

252  1,477,637 (6.3%) of the sequences might either occur as contaminants in other libraries, or could not

253  clearly be identified as being of vertical origin via the search pipeline used. The results obtained were

254  concordant with the other contamination analyses. For example, libraries known to have aberrant 18S

255  sequences contained a much larger average proportion of contaminant sequences (5.890/217,270

256  sequences, 2.7 %), but contained very few sequences that were identified as contaminants in other

257  libraries (252 sequences, 0.1%, Fig. 1A). A similar, but smaller enrichment in contaminants was

258  identified in libraries identified through 18S sequences as containing unconfirmed contamination

259  (16,871/ 912139 sequences; 1.8%), suggesting that at least some of these libraries are genuinely

260  biologically contaminated (Fig. 1A).

261

262  Specific libraries contained a much larger proportion of contaminant sequences, with 57.8% of the

263  *Deutzia scabia* (OTAN) found to be contaminant (Fig. 1B). These specific contaminations are from

264  *Gunnera manicata* (XMQO) (Fig. 1C), in line with the 18S based finding. Other cross-contamination

265  events found by this method include *Pseudolarix amabilis* found in *Monoclea gottschei* and *Galium*

266  *boreale* in *Impatien balsamifera*. We also, however, identified examples of widespread contamination

267  in libraries that had previously not been detected, for example over 35% of the sequences detected in

268  two libraries of the green alga *Olltmansiellopsis viridis* (Fig. 1B). These may relate to contaminants

269  that do not produce 18S sequences, as evidenced by the recent detection of Rhodobacteralean

270  commensal sequences in 1kp libraries from *Mantoniella squamata* (QXSZ), *Bathycoccus prasinos*

271  (MCPK) and *Nannochloropsis oculata* (JCFK) [19].   Additional results are provided in the associated

272  data release.

273

274

275

276

277  Assembly qualities

278

279   We assessed the quality of each assembled scaffold using Transrate [20], which detects several classes

280   of common assembly errors and assigns a quality score to each scaffold.  Users of the data may choose

281   to omit those portions of the assembly judged as low-quality when doing their own analyses.

282

| Percentile | Good Contigs (all sizes) | Good Contigs - Percentage |
|---|---|---|
| 5 | 19,355 | 32.47% |
| 25 | 30,755 | 44.83% |
| 50 | 37,983 | 53.65% |
| 75 | 47,608 | 62.93% |
| 95 | 71,368 | 74.87% |

283

284

285   Completeness of gene set

286

287   Two different approaches were used to estimate transcriptome completeness. Firstly, BUSCO v1 [21]

288   was applied with default settings, using the eukaryote and embryophyte conserved gene data sets

289   (eukaryota_odb9, embryophyta_odb9) as the query databases.  Secondly, conditional reciprocal best

290   BLAST (CRBB) hits were calculated using CRB-BLAST [22] with default parameters. The predicted

291   coding sequences were used as queries against the set of 248 core eukaryotic genes (CEGs) distributed

292   with the CEGMA software (Core Eukaryotic Genes Mapping Approach); these 248 genes are highly

293   conserved in eukaryotic genomes [23] and hence should be present in most transcriptomes.

294

295   As with all RNA-seq data, some genes are more highly expressed than others.  While the CEGMA and

296   BUSCO gene sets are intended to demonstrate the completeness of the transcriptomes, they are

297   sensitive to the expression of these genes.  Not all these genes will be expressed in the sample's tissues

298 at sufficiently high levels to be assembled.  A plot of the number of assembled scaffolds vs. the fraction

299 of the three gene sets found in the assembled scaffolds shows an increase in the gene fractions found as

300 the number of assembled scaffolds increases (Fig. 2).  However, these quickly saturate at 80+% for the

301 CEGMA and BUSCO-eukaryote sets, with a continuing rise over a larger range for the BUSCO-

302 embryophyte set.

303

304 This shows that the three gene sets have somewhat different expression patterns, with the CEGMA and

305 BUSCO-eukaryotic sets comprising genes that are more readily detected in our RNA samples.  Some

306 of the weaker sensitivity to the BUSCO-embryophyte set is attributable to our sampling species outside

307 of this phylum, which may not have the homologous genes; however, the observed effect is larger than

308 this and is also present when only the embryophyte samples are considered (not shown).

309

310 Percentage CEG abundance was calculated as number of CEGs with a CRBB hit divided by 248, the

311 number of CEGs used.  The percentage BUSCO abundance was calculated as 100% minus the missing

312 percentage. Samples with low abundance by these measures should be treated with caution because the

313 observed transcriptome incompleteness may indicate problems in library preparation or other types of

314 poor sample quality. For these reasons the taxonomic analyses in Ref. 1 excluded samples with less

315 than 57.5% BUSCO abundance.  The table below shows the percentages of complete genes found for

316 each of the three references at several percentile of the whole dataset.

317

| Percentile | CEGMA 248 | BUSCO – Embryophyta* | BUSCO – Eukaryota* |
|---|---|---|---|
| 5 | 79.03 | 11.2 (8.5) | 66.0 (37.3) |
| 25 | 89.92 | 44.1 (29.8) | 84.9 (64.4) |
| 50 | 92.34 | 62.5 (48.2) | 90.4 (75.9) |
| 75 | 93.55 | 75.2 (59.6) | 93.7 (84.1) |
| 95 | 94.76 | 82.6 (73.2) | 96.1 (91.0) |

318   *Complete+fragment assemblies reported with complete sequences in parentheses.

319

320   Re-use potential

321

322   Since many of the samples are from poorly sequenced clades, the Thousand Plant sequence data is the

323   first-large scale sequence data available for many species.  We expect these sequences to be of broad

324   interest to the plant sciences community, whether researchers merely use our sequences, supplement

325   them with their own sequences, or develop PCR primer and probe sets to collect entirely new sequence

326   data.

327

328

329

330   **Availability of Supporting Data**

331

332   Data to be in an associated *Gigascience*/GigaDB submission: [A copy of this is currently available at:

333   https://drive.google.com/drive/folders/175nB8kf1UQushuEzv7UaJLPNNwdOrxh5?usp=sharing ]

334

335   1. Tables with list of samples/assemblies (Sample-List-with-Taxonomy.tsv) and corresponding

336   ENA/NCBI references (NCBI-ENA-Sequence-Identifiers.csv) and GigaDB links (to be added).

337

338   2. The major part of the provided data includes a FASTA files containing the SOAPdenovo-Trans

339   assembly, the translation of the scaffolds to amino acids, the subset of the nucleotide sequence

340   corresponding to the translation, and tab-separated (text) files with tables of Transrate outputs.  These

341   are available for each of the assemblies listed in the supplemental table.  (onekp-data directory)

342

343   e.g. AALA-SOAPdenovo-Trans-assembly.fa.bz2, AALA-SOAPdenovo-Trans-translated.tar.bz2, AALA-

3. Two accessory tables containing details of the SILVA based SSU validation for each sample. The first (18S-analysis-Sample-Summary.xlsx) lists whether the sample is overall judged to be validated as containing the expected sequence, and whether it had alignments to any other plant sequences (described as worrisome contamination). The second file (18S-analysis-Scaffold-Results.xlsx), has more details listing each scaffold identified as being an 18S sequence, and which reference sequence it matched against.

4. The cross-contamination details. A summary file (Cross-contamination-Details.xlsx) includes a table (sheet Contamination Frequencies) with the number of contaminants, number of non-contaminant sequences, and the number of sequences inferred to be contaminants in other taxa for each sequence library.. Also included (sheet Contaminant Pairs) is a list of each pair of contaminant sequences identified, with the first column showing the contaminant sequence, and the second column the sequence corresponding to the orthologous contaminating partner against which the sequence was identified. Also included is a list of taxonomically close sample pairs which were not compared (sheet Excluded Taxa). Clean and contaminant FASTA sequence files for each library are available in the accompanying data (1kp_decontamination_libraries.gz.zip).

## Declarations

The authors declare that they have no conflicting interests, and that they believe that all the plant tissues were collected in accordance with applicable regulations and laws.

## References

370   1. One Thousand Plant Transcriptomes Initiative.  A Phylogenomic View of Evolutionary Complexity

371   in Green Plants.  In review, 2019.

372

373   2. Wickett NJ, Mirarab S, Nguyen N, Warnow T, Carpenter E, Matasci N, Ayyampalayam S, Barker

374   MS, Burleigh JG, Gitzendanner MA, Ruhfel BR, Wafula E, Der JP, Graham SW, Mathews S,

375   Melkonian M, Soltis DE, Soltis PS, Miles NW, Rothfels CJ, Lisa Pokorny, Shaw AJ, DeGironimo L,

376   Stevenson DW, Surek B, Villarreal J-C,  Roure B, Philippe H, dePamphilis CW, Chen T, Deyholos

377   MK, Baucom RS, Kutchan TM, Augustin MM, Wang J, Zhang Y, Tian Z, Yan Z, Wu X, Sun X, Wong

378   GK-S,  Leebens-Mack J. Phylotranscriptomic analysis of the origin and early diversification of land

379   plants. Proc. Natl. Acad. Sci. USA 2014;111:E4859–E4868 doi:10.1073/pnas.1323926111

380

381   3. Johnson MTJ, Carpenter EJ, Tian Z, Bruskiewich R, Burris JN, Carrigan CT, Chase MW, Clarke

382   ND, Covshoff S, dePamphilis CW, Edger PP, Goh F, Graham S, Greiner S, Hibberd JM, Jordon-

383   Thaden I, Kutchan TM, Leebens-Mack J, Melkonian M, Miles N, Myburg H, Patterson J, Pires JC,

384   Ralph P, Rolf M, Sage RF, Soltis D, Soltis P, Stevenson S, Stewart CN Jr, Surek B, Thomsen CJM,

385   Villarreal JC, Wu X, Zhang Y, Deyholos MK, Wong GK-S.  Evaluating Methods for Isolating Total

386   RNA and Predicting the Success of Sequencing Phylogenetically Diverse Plant Transcriptomes.  PLOS

387   One 2012; doi:10.1371/journal.pone.0050226.

388

389   4. Jordon-Thaden IE, Chanderbali AS, Gitzendanner MA, Soltis DE.  Modified CTAB and TRIzol

390   Protocols Improve RNA Extraction from Chemically Complex Embryophyta.  Appl in Plant Sci

391   2015;3:1400105 doi:10.3732/apps.1400105.

392

393   5. Mueller O, Lightfoot S, Schroeder A. Agilent Technologies Application Note: RNA Integrity

394   Number (RIN) – Standardization of RNA Quality Control. 2016.

395   https://www.agilent.com/cs/library/applications/5989-1165EN.pdf

396

397 6. Xie Y, Wu G, Tang J, Luo R, Patterson J, Liu S, Huang W, He G, Gu S, Zhou SLX, Lam T-W, Li Y,

398 Xu X, Wong GK-S, Wang J. SOAPdenovo-Trans: de novo transcriptome assembly with short RNA-

399 Seq reads. Bioinformatics 2014;30:1660–1666 doi:10.1093/bioinformatics/btu077.

400

401 7. Barker MS, Dlugosch KM, Dinh L, Challa RS, Kane NC, King MG, Rieseberg LH. EvoPipes.net:

402 Bioinformatic tools for ecological and evolutionary genomics. Evol. Bioinfo. 2010;6:143–149

403 doi:10.4137/EBO.S5861.

404

405 8. Goodstein DM, Shu S, Howson R, Neupane R, Hayes RD, Fazo J, Mitros T, Dirks W, Hellsten U,

406 Putnam N, Rokhsar DS. Phytozome: a comparative platform for green plant genomics. Nucl. Acids

407 Res. 2012;40:D1178–D1186 doi:10.1093/nar/gkr944.

408

409 9. Wheeler DL, Barrett T, Benson DA, Bryant SH, Canese K, Chetvernin V, Church DM, Dicuccio M,

410 Edgar R, Federhen S, Feolo M, Geer LY, Helmberg W, Kapustin Y, Khovayko O, Landsman D,

411 Lipman DJ, Madden TL, Maglott DR, Miller V, Ostell J, Pruitt KD, Schuler GD, Shumway M,

412 Sequeira E, Sherry ST, Sirotkin K, Souvorov A, Starchenko G, Tatusov RL, Tatusova TA, Wagner L,

413 Yaschenko E. Database resources of the National Center for Biotechnology Information. Nucl. Acids

414 Res. 2008;36:D13–D21 doi:10.1093/nar/gkm1000.

415

416 10. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. Genome Res. 2004;14:988–995

417 doi:10.1101/gr.1865504.

418

419 11. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. BLAST+:

420 architecture and applications. BMC Bioinformatics 2009;10:421. doi:10.1186/1471-2105-10-421.

421

422 12. Lusk RW. Divese and Widespread Contamination Evident in the Unmpped Depths of High

423 Throughput Sequencing Data. PLoS ONE 2014;9(10) e110808 doi:10.1371/journal.pone.0110808.

424

425 13. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO. The SILVA

426 ribosomal RNA gene database project: improved data processing and web-based tools. Nucl. Acids

427 Res. 2013;41:D590–D596 doi:10.1093/nar/gks1219.

428

429 14. Dorrell RG, Gile G, McCallum G, Méheust R, Bapteste EP, Klinger CM, Brillet-Guéguen L,

430 Freeman KD, Richter DJ, Bowler C. Chimeric origins of ochrophytes and haptophytes revealed

431 through an ancient plastid proteome. Elife 2007; 6, 23717 doi:10.7554/eLife.23717.

432

433 15. Dorrell RG, et al. (2019) Contrasting evolutionary fates accompany the loss of photosynthesis in

434 different heterotrophic chrysophytes. Proc Natl Acad Sci USA, in press.

435

436 16. Marron AO, Ratcliffe S, Wheeler GL, Goldstein RE, King N, Not F, de Vargas C, Richter DJ. The

437 Evolution of Silicon Transport in Eukaryotes. Mol Biol Evol 2016;33(12):3226-3248

438 doi:10.1093/molbev/msw209.

439

440 17. Kiełbasa SM, Wan R, Sato K, Horton P, Frith MC. Adaptive seeds tame genomic sequence

441 comparison. Genom Res 2011;21(3):487-493 doi:10.1101/gr.113985.110.

442

443 18. Moreno-Hagelsieb G, Latimer K. Choosing BLAST options for better detection of orthologs as

444 reciprocal best hits. Bioinformatics 2008;24(3):319-324 doi:10.1093/bioinformatics/btm585.

445
446 19. Sato S, Nanjappa D, Dorrell RG, Jimenez Vieira FR, Kazamia E, Tirichine L, Veluchamy A, Jaillon

447 O, Wincker P, Fussy Z, Kuo A, Obornik M, Munoz-Gomez SA, Mann DG, Bowler C, Zingone A.

448 Genome-enabled phylogenetic and functional reconstruction of an araphid pennate diatom CCMP470,

449 previously assigned as a radial centric diatom, and its bacterial commensal. Manuscript submitted.

450

451 20. Smith-Unna R, Boursnell C, Patro R, Hibberd J, Kelly S.  TransRate: reference free quality

452 assessment of de novo transcriptome assemblies.  Genome Res. 2016;26:1134–1144;

453 doi:10.1101/gr.196469.115.

454

455 21. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM.  BUSCO: assessing

456 genome assembly and annotation completeness with single-copy orthologs.  Bioinformatics

457 2015;31:3210–3212 doi:10.1093/bioinformatics/btv351.

458

459 22. Aubry S, Kelly S, Kümpers BMC, Smith-Unna RD, Hibberd JM.  Deep Evolutionary Comparison

460 of Gene Expression Identifies Parallel Recruitment of Trans-Factors in Two Independent Origins of C4

461 Photosynthesis.  PLOS Genetics 2014 doi:10.1371/journal.pgen.1004365.

462

463 23. Parra G, Bradnam K, Ning Z, Keane T, Korf I.  Assessing the gene space in draft genomes.  Nucl.

464 Acids Res.  2009;37:289–297 doi:10.1093/nar/gkn916.

465

466 Figure Captions:

467

468 Fig. 1. Panel A provides an overview of the total sequence percentage verified to be of contaminant

469 origin (red), or inferred to be possible contaminants in other sequence libraries (grey) in all 1kp

470 libraries, and libraries inferred to be contaminated through other techniques (e.g. 18S phylogenetic

471 placement).  Panel B lists 21 libraries in which > 6% of the total sequences are potential contaminants.

472 Panel C shows a heatmap of inferred contaminant interactions between pairs of species; contaminated

473 species are shown on the vertical axis, and contaminating species on the horizontal axis.

474

475

476   Fig. 2.  Fraction of the gene sets found (complete + fragments) versus the number of scaffolds

477   (300+ bp) in the assemblies.  For each sample, the fraction of the eukaryota and embryophyta sets

478   found in the assemblies are calculated with BUSCO and the fraction of the CEGMA 248 set with the

479   CRBB tool.  All three sets are more completely recovered at higher scaffold counts, but the BUSCO

480   embryophyta set is less complete in our samples.

Figure 1

Figure 2

Figure 2. Effect of Scaffold Numbers on Gene Set Completeness