

GigaScience

Access to RNA seq data from 1,173 green plant species: the 1000 Plant Transcriptomes Initiative (1KP) --Manuscript Draft--

Manuscript Number:	GIGA-D-19-00241R2	
Full Title:	Access to RNA seq data from 1,173 green plant species: the 1000 Plant Transcriptomes Initiative (1KP)	
Article Type:	Data Note	
Funding Information:	Alberta Innovates - Technology Futures (RES0010334)	Prof Gane Ka-Shu Wong
Abstract:	<p>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.</p>	
Corresponding Author:	Gane Ka-Shu Wong 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	
	Naim Matasci	
	Saravananaraj Ayyampalayam	
	Shuangxiu Wu	
	Jing Sun	
	Jun Yu	
	Fabio Rocha Jimenez Vieira	
	Chris Bowler	
	Richard G. Dorrell	
	Matthew A. Gitzendanner	
	Ling Li	
	Wensi Du	
	Kristian Ullrich	
	Norman J. Wickett	
	Todd J. Barkmann	
	Michael S. Barker	
	James H. Leebens-Mack	
	Gane Ka-Shu Wong	

Order of Authors Secondary Information:	
Response to Reviewers:	Final revisions made as requested.
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>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.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>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.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	No
<p>If not, please give reasons for any omissions below.</p> <p>as follow-up to "Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough</p>	The data is derived from plant samples for which no attempt was made to identify an age or sex for the source.

<p>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.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p> <p>"</p>	
<p>Availability of data and materials</p> <p>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.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>No</p>
<p>If not, please give reasons for any omissions below.</p> <p>as follow-up to "Availability of data and materials</p> <p>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.</p> <p>Have you have met the above</p>	<p>Additional data (contamination analysis, etc) will be submitted to GigaDB after this online process, as per the journal instructions.</p>

requirement as detailed in our [Minimum Standards Reporting Checklist?](#)

"

[Click here to view linked References](#)

1 Access to RNA seq data from 1,173 green plant species: the 1000 Plant 2 Transcriptomes Initiative (1KP).

3
4 Authors:

- 5 Eric J. Carpenter ejc@ualberta.ca, ORCID: 0000-0001-6267-7082 [1]
6 Naim Matasci <nmatasci@usc.edu>, ORCID: 0000-0003-4416-048X [2]
7 Saravanaraj Ayyampalayam <raj@plantbio.uga.edu> [3]
8 Shuangxiu Wu <wushx@big.ac.cn> [4]
9 Jing Sun <jsun@genetics.ac.cn> [4]
10 Jun Yu <junyu@big.ac.cn> [4]
11 Fabio Rocha Jimenez Vieira rocha@biologie.ens.fr, ORCID: 0000-0001-9872-0337 [5]
12 Chris Bowler <cbowler@biologie.ens.fr> [5]
13 Richard G. Dorrell <dorrell@biologie.ens.fr> [5]
14 Matthew A. Gitzendanner magitz@ufl.edu, ORCID: 0000-0002-7078-4336 [6]
15 Ling Li <liling3@cngb.org> [7]
16 Wensi Du <duwensi@cngb.org> [7]
17 Kristian Ullrich <ullrich@evolbio.mpg.de> [8]
18 Norm J. Wickett <norman.wickett@gmail.com> [9]
19 Todd J. Barkmann <todd.barkman@wmich.edu> [10]
20 Michael S. Barker <msbarker@email.arizona.edu> [11]
21 James H. Leebens-Mack <jleebensmack@uga.edu>, ORCID: 0000-0003-4811-2231 [12]
22 Gane Ka-Shu Wong <gane@ualberta.ca>* contact author, ORCID: 0000-0001-6108-5560 [1,7,13]
23

- 24 1. Department of Biological Sciences, University of Alberta, Edmonton, Alberta, T6G 2E9, Canada.
25 2. CyVerse, University of Arizona, Arizona, U.S.A.; Current address: Lawrence J. Ellison Institute for
26 Transformative Medicine, University of Southern California, Los Angeles, CA 90033, U.S.A.
27 3. Georgia Advanced Computing Resource Center, University of Georgia, Athens GA 30602, USA. 4.
28 CAS Key Laboratory of Genome Sciences and Information, Beijing, Institute of Genomics, Chinese
29 Academy of Sciences, Beijing 100101, China.
30 5. Institut de Biologie de l'ENS (IBENS), Département de biologie, École normale supérieure, CNRS,
31 INSERM, Université PSL, 75005 Paris, France
32 6. Department of Biology, University of Florida, Gainesville, Florida 32611, USA.
33 7. BGI-Shenzhen, Beishan Industrial Zone, Yantian District, Shenzhen 518083, China.
34 8. Department of Evolutionary Genetics, Max Planck Institute for Evolutionary Biology, Plön,
35 Germany.
36 9. Chicago Botanic Garden, Glencoe, IL 60022, and Program in Biological Sciences, Northwestern
37 University, Evanston, IL 60208 USA.
38 10. Department of Biological Sciences, Western Michigan University, Kalamazoo MI 49008-5410
39 USA.
40 11. Department of Ecology & Evolutionary Biology, University of Arizona, Tucson, AZ 85721 USA.
41 12. Department of Plant Biology, University of Georgia, Athens, GA 30602, USA.
42 13. Department of Medicine, University of Alberta, Edmonton, Alberta, T6G 2E1, Canada.
43

44 Abstract

45 The 1000 Plant Transcriptomes Initiative (1KP) explored the genetic diversity of green plants
46 (Viridiplantae) by sequencing RNA from 1,342 samples representing 1,173 species. All of the analyses

47 done for the 1KP capstone, and previous studies on subsets of these data, are based on a series of de
48 novo transcriptome assemblies and related outputs that will be described in this publication. We also
49 describe assessments of the data quality and an analysis to remove cross-contamination between the
50 samples. These data will be useful to researchers with interests in specific gene families, either across
51 the green plant tree of life or in more focused lineages.

52

53

54 **Keywords**

55

56 RNA, plants, assemblies, genes, contamination, completeness

57

58

59

60 **Data Description**

61

62 1KP has sequenced and analysed RNA from 1,342 RNA samples of 1,173 green plant species
63 representing all major taxa within the Viridiplantae, including streptophyte and chlorophyte green
64 algae, bryophytes, ferns, angiosperms, and gymnosperms. Importantly, our selection criteria eschewed
65 the model organisms and crop species where other plant sequencing efforts have historically been
66 concentrated. While many of the samples were selected for the phylogenomic analyses, others were
67 motivated by different subprojects.

68

69 Major papers describing the project have been published elsewhere [1, 2]. The most recent papers [1, 3]
70 are focused on large-scale phylogenetic analyses made possible by the breadth of this data set. While all
71 of the 1,342 samples were used in one analysis or another, not all of them were judged of adequate
72 quality for every analysis. As each paper uses different analyses, appropriate criteria for sample quality
73 are different, and thus each uses a different subset of the sample data. This Data Note describes the
74 whole data set and provides additional details on the sample and sequence processing as well as quality
75 assessments of these data. This supplements and replaces our earlier work [4] outlining plans for the
76 1KP efforts.

77

78 **Methods**

79

80 **Sampling strategy**

81

82 Because of the diversity and the number of species analyzed, no one source could be used. Samples
83 were provided by a global network of collaborators who obtained materials from a variety of sources,
84 including field collection of wild plants, greenhouses, botanical gardens, laboratory specimens, and
85 algal culture collections. To ensure an abundance of expressed genes, we preferred live growing cells,
86 e.g. young leaves, flowers, or shoots, although many samples were also from roots, or other tissues.
87 Because of the sample diversity, we did not attempt to define specific standards on growth conditions,
88 time of collection, or age of tissue. For more details, see the supplemental methods in the major
89 analysis paper [1].

90

91 **RNA extraction**

92

93 Given the biochemical diversity of these samples, no one RNA extraction protocol was appropriate for
 94 all samples. Most samples were extracted using commonly known protocols or using commercial kits.
 95 For complete details of the many specific protocols used, please see Appendix S1 of Johnson et al. [5]
 96 and Jordon-Thaden et al. [6]. The individual protocols are also available via a protocols.io collection
 97 [7]. Depending on the sample, RNA extractions might have been done by the sample provider, a
 98 collaborator near the provider, or the sequencing lab (BGI-Shenzhen).

99

100

101 Sequencing at BGI

102

103 Samples of extracted RNA or frozen tissues were sent to the sequencing lab, BGI-Shenzhen. Prior to
 104 library construction, RNA samples were screened by Agilent Bioanalyzer RIN scores [8] and basic
 105 photometry; obvious low-quality outliers (e.g., RIN scores less than 6 and/or loss of distinct
 106 electropherogram peaks) were excluded. Libraries for Illumina sequencing were constructed using
 107 Illumina’s standard procedures. Some samples for which only a small amount of RNA was available
 108 were processed using TruSeq kits.

109

110 Initially, sequencing was done on the Illumina GAII platform, but later samples were run on the HiSeq
 111 platform. Associated with this change was a shift from ~72 bp read lengths to 90 bp read lengths (both
 112 cases paired-end). Libraries were indexed and multiplexed in the sequencer lanes to a target
 113 sequencing depth of 2 Gbp per sample. Average depth achieved was 1.99 Gbp of sequence with Phred
 114 quality 30 (1 error per thousand bases) or better, and varies across samples with half of samples in the
 115 1.9–2.5 Gbp range as summarized by Table 1.

116

117 Table 1 Distribution in amount of sequence data per sample library. Summary percentiles
 118 characterising the sizes of the datasets in giga-basepairs of sequence.

119

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

120 .

121

122 The data was cleaned by eliminating reads containing adapter-primer sequences or high numbers of
 123 low quality bases (i.e. more than half of Phred quality below 5 [32 % error rate] or more than 10%
 124 uncalled). Sequencing and transcriptome assembly protocols are available in protocols.io [9].

125

126

127 *De novo* assembly

128

129 Once the data was transferred from BGI, the FastQ files were given a uniform name based on a quasi-
 130 random four-letter identification code. A list of all the samples and their ID code is included in the
 131 associated data. These identifiers also distinguish otherwise identical repeated samples, and provide a
 132 stable reference when a sample's species identification is changed.

133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179

Quality filtered reads were assembled using the SOAPdenovo-Trans transcript assembler (version 2012-04-05) [10]. No additional pre-processing of the data was performed. This largely used the program defaults, with the slight modification of increasing the k -mer length to 25 bp and reducing the number of processor threads to one. This reduced thread count allowed us to more efficiently use our computer resources. Both the internal FillGap module and the external GapCloser post-processor (supplied with SOAPdenovo-Trans) were run. An example of the commands used for one of the assemblies (dataset AEPI, *Lineum leonii*):

```
SOAPdenovo-Trans-31kmer all -s config -p 1 -K 25 -e 2 -F -L 100 -t 5 -o AEPI
GapCloser -a AEPI.scafSeq -b config -o AEPI.GapCloser.fa -l 100 -p 25 -t 1
```

These commands refer to a configuration file named config, which specified the expected insert size, maximum read length, and read-sequence filenames. The contents of this file were:

```
max_rd_len=120
[LIB]
avg_ins=200
rank=1
q1=AEPI-read_1.fq
q2=AEPI-read_2.fq
```

When multiple samples from the same species were co-assembled, the last five lines were repeated for each data source with the appropriate filenames. Such assemblies were also assigned unique four-letter identifiers. After assembly the output contigs/scaffolds were renamed giving each a unique name containing the assembly's four-letter identifier, a number within the assembly, and as a descriptive name the species, with additional description of the tissue or other identifier when multiples samples of the same species where sequenced.

Identification of coding regions and protein translation

To identify likely proteins within the assembled transcripts, sequences were passed through TransPipe [11], which identified reading frames and protein translations by comparison to protein sequences from 22 sequenced and annotated plant genomes in Phytozome (RRID:SCR_006507)[12]. Using BLASTX (RRID:SCR_001653)[13], best hit proteins were paired with each assembled scaffold at a threshold of $1E-10$ expectation-value and a minimum length of 100 amino acid residues. Scaffolds that did not have a best hit protein at this level were removed. These removed scaffolds are predominantly from the numerous short and likely fragmentary sequences; however some complete genes will have been lost. To determine reading frames and estimate amino acid sequences, each gene is aligned against its best hit protein by Genewise 2.2.0 (RRID:SCR_015054)[14]. Using the highest scoring Genewise DNA-protein alignments, stop codons and those codons containing ambiguous nucleotides were removed to produce an amino acid sequence for each gene. Outputs in the associated data are paired DNA and protein sequences.

BLAST searches

180 Thanks to the support of China National GeneBank (CNGB), a BLAST search service
181 (<http://db.cngb.org/onekp/>) allows public searches against the assemblies and protein translations.
182 CNGB developed the service using NCBI BLAST+ (version 2.6.0) [15]. It integrates all public
183 datasets from CNGB applications, BGI projects and external data sources, and provides a
184 comprehensive and convenient sequence searching. A specialized interface for BLAST searching the
185 1KP dataset allows limiting the search to specific families, orders, or 25 higher-level clades. For
186 assemblies, there are 21,398,790 nucleotide sequences, 6,188,419,272 bases in total. And for the
187 Transpipe protein translations, there are 103 million protein sequences comprising over 47 billion
188 amino acids in total.

189

190

191 **Validation**

192

193 Purity and contamination

194

195 High throughput sequencing methods are always at risk of contamination [16]. In 1KP, the diversity of
196 sources for the samples, and especially the fact that axenic cultures are not a viable option in most
197 instances, ensure that there will always be some contamination of the plant tissue by other
198 environmental nucleic acids. These can reasonably be expected to include bacterial, fungal, and insect
199 species that live in and on the plant tissues, and more rarely, from contact with larger species such as
200 frogs, mice, birds and humans.

201

202 For most analyses, these minor contaminants are not expected to matter, as only the most abundant of
203 such contaminants will be present in sufficient quantities to assemble. In many cases, they are also
204 sufficiently diverged from the intended species that they can be easily recognised as non-plant genes.
205 Unfortunately, this is not always the case. Some analyses are further protected by looking at the whole
206 of the available transcriptome, whereby the many genes from the target species will overpower a few
207 contaminants. Single gene family analyses do not have this advantage and must rely on other methods
208 to reject non-plant genes.

209

210 Another possibility is significant contamination during sample processing when plant RNA is
211 transferred between adjacent samples, or when whole samples are accidentally mislabelled.

212

213 Given the potential contamination problems, we tried to identify them in the sequence data by
214 comparing the assembled sequences by BLASTn to a reference set of nuclear 18S rRNA sequences
215 from the SILVA SSU rRNA database (<http://www.arb-silva.de>) [17]. The BLASTn alignment to an
216 assembly with the lowest expectation-value is taken to indicate the assembly has a similar taxonomic
217 origin as the reference sequence. However, alignments of less than 300 bp or expectation-values above
218 $1E-9$ often align to several distantly related species and were ignored.

219

220 For most samples, we found an 18S sequence most-similar to a SILVA sequence from the same
221 taxonomic family as the expected sample species. This is not true for all our samples, and may indicate
222 a failure to assemble the 18S sequence, limitations in the taxonomic identification from the BLASTn
223 results, or mislabelling of sample. In a few cases, additional (and possibly contaminant) 18S sequences
224 were found. Because the 18S rRNA sequence is highly expressed, we expect that this method is likely

225 to be sensitive to low levels of contamination. In a few cases, the taxonomic irregularities were judged
226 sufficiently severe that samples were excluded from various analyses.

227
228 The accompanying data includes two accessory files containing details of this SILVA based SSU
229 validation for each sample [18]. The first lists whether the sample is overall judged to be validated as
230 containing the expected taxon, and whether it had alignments to any other plant sequences (described
231 as “worrisome contamination”). The second file, more detailed, lists each scaffold identified as being
232 18S-like sequence, and which reference sequence it matched against.

233
234 It must be emphasized, however, that these files (and indeed this entire section) describe how we
235 removed contaminations from the final analyses. The published data, 1,342 RNA samples from 1,173
236 green plant species, does not include the worst contaminations.

237
238 Pairwise Cross-contamination of Assemblies

239
240 Cross contamination between datasets was also identified by a genome-scale sequence search pipeline,
241 adapted from previous studies [19-21]. Briefly, each pair of assemblies (nucleotide) was compared and
242 a threshold identity level established, above which sequences are likely to be contamination between
243 the pair. While best for identifying technical contamination between libraries (e.g. due to mixing of
244 RNA samples), this technique could also detect other biological contamination events (e.g.
245 contamination of pairs of libraries with common commensal organisms). An additional search step,
246 using the entire 1KP sequence library, identified the probable evolutionary origin of each sequences.

247
248 The pair-wise comparison used LAST v. 963 (RRID:SCR_006119)[22] with the --cR01 option, and the
249 respective matches were grouped and ordered by similarity. To avoid artifactually excluding sequences
250 between closely related species, which may have very high degrees of similarity [16], pairs of libraries
251 from the same family, along with pairs of libraries separated by two or fewer branches in the consensus
252 1KP multigene phylogeny, were excluded from the searches [2].

253
254 The expected distribution of the matched sequence identities has a maximum at the pairwise identity
255 reflecting the evolutionary distance between the two species [20, 21]. In contrast, a cross-contaminated
256 pair should contain many sequences of near 100% similarity, and the similarity value which has the
257 first minimum number of sequences below this level (i.e. the first inflexion point in a curve plotting the
258 total number of sequences of each percentage similarity value) can be used as a threshold for
259 discriminating contaminating sequences [20, 21]. The code is available at [https://github.com/Plant-
260 and-diatom-genomics-IBENS-Paris/Decontamination-pipeline](https://github.com/Plant-and-diatom-genomics-IBENS-Paris/Decontamination-pipeline).

261
262 The output of this analysis is pairs of apparent orthologs whose sequence similarities are higher than
263 the cut-off in one or both libraries, i.e. potential contamination. To discriminate donors and recipients
264 in each contaminant pair, each of these potential contaminants was searched against all the non-
265 contaminant assemblies by BLASTn, using the option -max_target_seqs 3 [23]. Queries with at least
266 one of the three best alignments against a sequence from the same family, or from a taxon separated by
267 fewer than two branches within the 1kp tree [2], were excluded from the list of potential contaminants;
268 whereas sequences that yielded best hits exclusively against more distantly related taxa, were verified
269 as potential contaminants. Clean and contaminant FASTA sequence files for each library are available
270 in the accompanying data.

271

272 An overview of the results is presented in Fig. 1. In total, we identified 79,175 nucleotide sequences
273 (0.3 %) of a total 23,436,405 searched as being clearly of contaminant origin (Fig. 1A). A further
274 1,477,637 (6.3%) of the sequences might either occur as contaminants in other libraries, or could not
275 clearly be identified as being of vertical origin via the search pipeline used. The results obtained were
276 concordant with the other contamination analyses. For example, libraries known to have aberrant 18S
277 sequences contained a much larger average proportion of contaminant sequences (5.890/217,270
278 sequences, 2.7 %), but contained very few sequences that were identified as contaminants in other
279 libraries (252 sequences, 0.1%, Fig. 1A). A similar, but smaller enrichment in contaminants was
280 identified in libraries identified through 18S sequences as containing unconfirmed contamination
281 (16,871/ 912139 sequences; 1.8%), suggesting that at least some of these libraries are genuinely
282 biologically contaminated (Fig. 1A).

283

284 Specific libraries contained a much larger proportion of contaminant sequences, with 57.8% of the
285 *Deutzia scabia* (OTAN) found to be contaminant (Fig. 1B). These specific contaminations are from
286 *Gunnera manicata* (XMQO) (Fig. 1C), in line with the 18S based finding. Other cross-contamination
287 events found by this method include *Pseudolarix amabilis* found in *Monoclea gottschei* and *Galium*
288 *boreale* in *Impatiens balsamifera*. We also, however, identified examples of widespread contamination
289 in libraries that had previously not been detected, for example over 35% of the sequences detected in
290 two libraries of the green alga *Olltmansiellopsis viridis* (Fig. 1B). These may relate to contaminants
291 that do not produce 18S sequences, as evidenced by the recent detection of Rhodobacteralean
292 commensal sequences in 1kp libraries from *Mantoniella squamata* (QXSZ), *Bathycoccus prasinus*
293 (MCPK) and *Nannochloropsis oculata* (JCFK) [24]. Additional results are provided in the associated
294 data release.[18]

295

296

297

298

299 Assembly qualities

300

301 We assessed the quality of each assembled scaffold/contig using the read mapping mode of Transrate
302 [25], which detects several classes of common assembly errors and assigns a quality score to each
303 scaffold. Users of the data may choose to omit those portions of the assembly judged as low-quality
304 when doing their own analyses. While the assemblies for each sample vary in assessed quality (Table
305 2), there are thousands of good scaffolds in even the worst of them.

306

307 Table 2. Assembly quality assessment by Transrate. Characteristic percentiles summarising the per
308 sample distributions of high-quality scaffolds for both total counts and fractions of the sample.

309

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

310 .
311
312 Completeness of gene set
313

314 Two different approaches were used to estimate transcriptome completeness. Firstly, BUSCO v1 [26]
315 was applied with default settings, using the eukaryote and embryophyte conserved gene data sets
316 (eukaryota_odb9, embryophyta_odb9) as the query databases. Secondly, conditional reciprocal best
317 BLAST (CRBB) hits were calculated using CRB-BLAST [27] with default parameters. The predicted
318 coding sequences were used as queries against the set of 248 core eukaryotic genes (CEGs) distributed
319 with the CEGMA software (Core Eukaryotic Genes Mapping Approach); these 248 genes are highly
320 conserved in eukaryotic genomes [28] and hence should be present in most transcriptomes.

321
322 As with all RNA-Seq data, some genes are more highly expressed than others. While the CEGMA and
323 BUSCO gene sets are intended to demonstrate the completeness of the transcriptomes, they are
324 sensitive to the expression of these genes. Not all these genes will be expressed in the sample's tissues
325 at sufficiently high levels to be assembled. A plot of the number of assembled scaffolds vs. the fraction
326 of the three gene sets found in the assembled scaffolds shows an increase in the gene fractions found as
327 the number of assembled scaffolds increases (Fig. 2). However, these quickly saturate at 80+% for the
328 CEGMA and BUSCO-eukaryote sets, with a continuing rise over a larger range for the BUSCO-
329 embryophyte set.

330
331 This shows that the three gene sets have somewhat different expression patterns, with the CEGMA and
332 BUSCO-eukaryotic sets comprising genes that are more readily detected in our RNA samples. Some
333 of the weaker sensitivity to the BUSCO-embryophyte set is attributable to our sampling species outside
334 of this phylum, which may not have the homologous genes; however, the difference is present when
335 only the embryophyte samples are considered (not shown).

336
337 Percentage CEG abundance was calculated as number of CEGs with a CRBB hit divided by 248, the
338 number of CEGs used. The percentage BUSCO abundance was calculated as 100% minus the missing
339 percentage. Samples with low abundance by these measures should be treated with caution because the
340 observed transcriptome incompleteness may indicate problems in library preparation or other types of
341 poor sample quality. For these reasons the taxonomic analyses in Ref. 1 excluded samples with less
342 than 57.5% BUSCO abundance. The table below shows the percentages of complete genes found for
343 each of the three references at several percentile of the whole dataset.

344
345 Table 3. Completeness of gene sets. Characteristic percentiles summarizing the distributions of the
346 CEGMA 248 and BUSCO genome completeness scores. *BUSCO numbers are the sum of the
347 complete and fragment assembly counts reported, with numbers based on the complete sequence
348 numbers alone given in parentheses.

349
350

Percentile	CEGMA 248	BUSCO – Embryophyta*	BUSCO – Eukaryota*
5th	79.03	11.2 (8.5)	66.0 (37.3)
25th	89.92	44.1 (29.8)	84.9 (64.4)
50th	92.34	62.5 (48.2)	90.4 (75.9)
75th	93.55	75.2 (59.6)	93.7 (84.1)

95th	94.76	82.6 (73.2)	96.1 (91.0)
------	-------	-------------	-------------

351
352
353
354
355
356
357
358
359

Re-use potential

Since many of the samples are from poorly sequenced clades, the Thousand Plant sequence data is the first-large scale sequence data available for many species. We expect these sequences to be of broad interest to the plant sciences community, whether researchers merely use our sequences, supplement them with their own sequences, or develop PCR primer and probe sets to collect entirely new sequence data.

360 Availability of supporting source code and requirements

- 361 • Project name: 1KP Decontamination-pipeline
- 362 • Project home page: <https://github.com/Plant-and-diatom-genomics-IBENS-Paris/Decontamination-pipeline>
- 363 • Operating system: linux
- 364 • Programming language: bash
- 365 • Other requirements: LAST, join C++ libraries
- 366 • License: GNU GPL v3

368
369

370 Availability of Supporting Data

371

372 Sequencing data is available from EBI BioProject’s: PRJEB4921, PRJEB8056, PRJEB21674,
373 PRJNA163187 and STUDY: SRP012845. Data for the 1KP project is available in Cyverse Data
374 Commons[29]. All the other supporting data presented here is associated with a GigaDB
375 submission[16]. These include:

376

- 377 1. Tables with list of samples/assemblies (Sample-List-with-Taxonomy.tsv) and corresponding
378 ENA/NCBI references (NCBI-ENA-Sequence-Identifiers.csv) and GigaDB links (to be added).
379
- 380 2. The major part of the provided data has for a directory for each assembly. This is named based on
381 the four-letter code and a species name. Within the directory are a FASTA file containing the
382 SOAPdenovo-Trans assembly, translations of the scaffolds to amino acids, the subset of the nucleotide
383 sequence corresponding to the translation, and tab-separated (text) files with tables of Transrate outputs
384 assessing the assemblies and lists of the reference sequence each translation is based on. These are
385 available for each of the assemblies listed in the supplemental table. (onekp-data directory)
386
387 e.g. in directory AALA-Meliosma_cunifolia are AALA-SOAPdenovo-Trans-assembly.fa.bz2, AALA-
388 translated-protein.fa.gz, AALA-translated-nucleotides.fa.gz, AALA- Transrate-assembly-stats.tsv.gz,
389 and AALA-translated-reference-names.tsv.gz

390

- 391 3. Two accessory tables containing details of the SILVA based SSU validation for each sample. The
392 first (18S-analysis-Sample-Summary.xlsx) lists whether the sample is overall judged to be validated as

393 containing the expected sequence, and whether it had alignments to any other plant sequences
394 (described as worrisome contamination). The second file (18S-analysis-Scaffold-Results.xlsx), has
395 more details listing each scaffold identified as being an 18S sequence, and which reference sequence it
396 matched against.

397
398 4. The cross-contamination details. A summary file (Cross-contamination-Details.xlsx) includes a
399 table (sheet Contamination Frequencies) with the number of contaminants, number of non-contaminant
400 sequences, and the number of sequences inferred to be contaminants in other taxa for each sequence
401 library.. Also included (sheet Contaminant Pairs) is a list of each pair of contaminant sequences
402 identified, with the first column showing the contaminant sequence, and the second column the
403 sequence corresponding to the orthologous contaminating partner against which the sequence was
404 identified. Also included is a list of taxonomically close sample pairs which were not compared (sheet
405 Excluded Taxa). Clean and contaminant FASTA sequence files for each library are available in the
406 accompanying data (1kp_decontamination_libraries.gz.zip).

407
408

409 **Declarations**

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

412

413 **References**

414

415 1. One Thousand Plant Transcriptomes Initiative. One Thousand Plant Transcriptomes and
416 Phylogenomics of Green Plants. *Nature*. in press, 2019.

417

418 2. Wickett NJ, Mirarab S, Nguyen N, Warnow T, Carpenter E, Matasci N, Ayyampalayam S, Barker
419 MS, Burleigh JG, Gitzendanner MA, Ruhfel BR, Wafula E, Der JP, Graham SW, Mathews S,
420 Melkonian M, Soltis DE, Soltis PS, Miles NW, Rothfels CJ, Lisa Pokorny, Shaw AJ, DeGironimo L,
421 Stevenson DW, Surek B, Villarreal J-C, Roure B, Philippe H, dePamphilis CW, Chen T, Deyholos
422 MK, Baucom RS, Kutchan TM, Augustin MM, Wang J, Zhang Y, Tian Z, Yan Z, Wu X, Sun X, Wong
423 GK-S, Leebens-Mack J. Phylotranscriptomic analysis of the origin and early diversification of land
424 plants. *Proc. Natl. Acad. Sci. USA* 2014;111:E4859–E4868 doi:10.1073/pnas.1323926111

425

426 3. Li Z, Barker MS. Inferring putative ancient whole genome duplications in the 1000 Plants (1KP)
427 initiative: access to gene family phylogenies and age distributions. bioRxiv 735076; doi:
428 <https://doi.org/10.1101/735076>

429

430 4. Matasci N, Hung L-H, Yan Z, Carpenter EJ, Wickett NJ, Mirarab S, Nguyen N, Warnow T,
431 Ayyampalayam S, Barker M, Burleigh JG, Gitzendanner MA, Wafula E, Der JP, dePamphilis CW,
432 Roure B, Philippe H, Ruhfel BR, Miles NW, Graham SW, Mathews S, Surek B, Melkonian M, Soltis
433 DE, Soltis PS, Rothfels C, Pokorny L, Shaw JA, DeGironimo L, Stevenson DW, Villarreal JC, Cheni
434 T, Kutchan TM, Rolf M, Baucom RS, Deyholos MK, Samudrala R, Tian Z, Wu X, Sun X, Zhang Y,
435 Wang J, Leebens-Mack J, Wong GK-S. Data access for the 1,000 Plants (1KP) project. *GigaScience*
436 2014;3 doi:10.1186/2047-217X-3-17

437

- 438 5. Johnson MTJ, Carpenter EJ, Tian Z, Bruskiwich R, Burris JN, Carrigan CT, Chase MW, Clarke
439 ND, Covshoff S, dePamphilis CW, Edger PP, Goh F, Graham S, Greiner S, Hibberd JM, Jordon-
440 Thaden I, Kutchan TM, Leebens-Mack J, Melkonian M, Miles N, Myburg H, Patterson J, Pires JC,
441 Ralph P, Rolf M, Sage RF, Soltis D, Soltis P, Stevenson S, Stewart CN Jr, Surek B, Thomsen CJM,
442 Villarreal JC, Wu X, Zhang Y, Deyholos MK, Wong GK-S. Evaluating Methods for Isolating Total
443 RNA and Predicting the Success of Sequencing Phylogenetically Diverse Plant Transcriptomes. *PLOS*
444 *One* 2012; doi:10.1371/journal.pone.0050226.
445
- 446 6. Jordon-Thaden IE, Chanderbali AS, Gitzendanner MA, Soltis DE. Modified CTAB and TRIzol
447 Protocols Improve RNA Extraction from Chemically Complex Embryophyta. *Appl in Plant Sci*
448 2015;3:1400105 doi:10.3732/apps.1400105.
449
- 450 7. Marc T. J. Johnson, et al. (2019). RNA Isolation from Plant Tissue. *protocols.io*
451 [dx.doi.org/10.17504/protocols.io.439gyr6](https://doi.org/10.17504/protocols.io.439gyr6)
452
- 453 8. Mueller O, Lightfoot S, Schroeder A. Agilent Technologies Application Note: RNA Integrity
454 Number (RIN) – Standardization of RNA Quality Control. 2016.
455 <https://www.agilent.com/cs/library/applications/5989-1165EN.pdf>
456
- 457 9. Eric J. Carpenter et al. (2019). Sequencing Protocols for the One Thousand Plant Transcriptomes
458 Initiative. *protocols.io* <http://dx.doi.org/10.17504/protocols.io.38jgrun>
459
- 460 10. Xie Y, Wu G, Tang J, Luo R, Patterson J, Liu S, Huang W, He G, Gu S, Zhou SLX, Lam T-W, Li
461 Y, Xu X, Wong GK-S, Wang J. SOAPdenovo-Trans: de novo transcriptome assembly with short
462 RNA-Seq reads. *Bioinformatics* 2014;30:1660–1666 doi:10.1093/bioinformatics/btu077.
463
- 464 11. Barker MS, Dlugosch KM, Dinh L, Challa RS, Kane NC, King MG, Rieseberg LH. EvoPipes.net:
465 Bioinformatic tools for ecological and evolutionary genomics. *Evol. Bioinfo.* 2010;6:143–149
466 doi:10.4137/EBO.S5861.
467
- 468 12. Goodstein DM, Shu S, Howson R, Neupane R, Hayes RD, Fazo J, Mitros T, Dirks W, Hellsten U,
469 Putnam N, Rokhsar DS. Phytozome: a comparative platform for green plant genomics. *Nucl. Acids*
470 *Res.* 2012;40:D1178–D1186 doi:10.1093/nar/gkr944.
471
- 472 13. Wheeler DL, Barrett T, Benson DA, Bryant SH, Canese K, Chetvernin V, Church DM, Dicuccio
473 M, Edgar R, Federhen S, Feolo M, Geer LY, Helmberg W, Kapustin Y, Khovayko O, Landsman D,
474 Lipman DJ, Madden TL, Maglott DR, Miller V, Ostell J, Pruitt KD, Schuler GD, Shumway M,
475 Sequeira E, Sherry ST, Sirotkin K, Souvorov A, Starchenko G, Tatusov RL, Tatusova TA, Wagner L,
476 Yaschenko E. Database resources of the National Center for Biotechnology Information. *Nucl. Acids*
477 *Res.* 2008;36:D13–D21 doi:10.1093/nar/gkm1000.
478
- 479 14. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. *Genome Res.* 2004;14:988–995
480 doi:10.1101/gr.1865504.
481
- 482 15. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. BLAST+:
483 architecture and applications. *BMC Bioinformatics* 2009;10:421. doi:10.1186/1471-2105-10-421.

484
485 16. Lusk RW. Diverse and Widespread Contamination Evident in the Unmpped Depths of High
486 Throughput Sequencing Data. PLoS ONE 2014;9(10) e110808 doi:10.1371/journal.pone.0110808.
487
488 17. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO. The SILVA
489 ribosomal RNA gene database project: improved data processing and web-based tools. Nucl. Acids
490 Res. 2013;41:D590–D596 doi:10.1093/nar/gks1219.
491
492 18. Carpenter EJ; Matasci N; Ayyampalayam S; Wu S; Sun J; Yu J; Jimenez Vieira FR; Bowler C;
493 Dorrell RG; Gitzendanner MA; Li L; Du W; Ullrich K; Wickett NJ; Barkmann TJ; Barker MS;
494 Leebens-Mack JH; Wong GK (2019): Data and results from RNA-sequencing of 1,173 species for the
495 1000 Plants (1KP) initiative GigaScience Database. <http://dx.doi.org/10.5524/100627>
496 19. Dorrell RG, Gile G, McCallum G, Méheust R, Bapteste EP, Klinger CM, Brillet-Guéguen L,
497 Freeman KD, Richter DJ, Bowler C. Chimeric origins of ochrophytes and haptophytes revealed
498 through an ancient plastid proteome. Elife 2007; 6, 23717 doi:10.7554/eLife.23717.
499
500 20. Dorrell RG, Azuma T, Nomura M, de Kerdrel GA, Paoli L, Yang S, Bowler C, Ishii K,
501 Miyashita H, Gile GH, Kamikawa R. Principles of plastid reductive evolution illuminated by
502 nonphotosynthetic chrysophytes. Proc. Natl. Acad. Sci. 2019;116:6914–6923
503 doi:10.1073/pnas.1819976116
504
505 21. Marron AO, Ratcliffe S, Wheeler GL, Goldstein RE, King N, Not F, de Vargas C, Richter DJ. The
506 Evolution of Silicon Transport in Eukaryotes. Mol Biol Evol 2016;33(12):3226–3248
507 doi:10.1093/molbev/msw209.
508
509 22. Kiełbasa SM, Wan R, Sato K, Horton P, Frith MC. Adaptive seeds tame genomic sequence
510 comparison. Genom Res 2011;21(3):487–493 doi:10.1101/gr.113985.110.
511
512 23. Moreno-Hagelsieb G, Latimer K. Choosing BLAST options for better detection of orthologs as
513 reciprocal best hits. Bioinformatics 2008;24(3):319–324 doi:10.1093/bioinformatics/btm585.
514
515 24. Sato S, Nanjappa D, Dorrell RG, Jimenez Vieira FR, Kazamia E, Tirichine L, Veluchamy A, Jaillon
516 O, Wincker P, Fussy Z, Kuo A, Obornik M, Munoz-Gomez SA, Mann DG, Bowler C, Zingone A.
517 Genome-enabled phylogenetic and functional reconstruction of an araphid pennate diatom CCMP470,
518 previously assigned as a radial centric diatom, and its bacterial commensal. 2019. The molecular life of
519 diatoms EMBO workshop. Poster Abstract 2. <http://meetings.embo.org/event/19-diatoms>
520
521 25. Smith-Unna R, Bournnell C, Patro R, Hibberd J, Kelly S. TransRate: reference free quality
522 assessment of de novo transcriptome assemblies. Genome Res. 2016;26:1134–1144;
523 doi:10.1101/gr.196469.115.
524
525 26. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing
526 genome assembly and annotation completeness with single-copy orthologs. Bioinformatics
527 2015;31:3210–3212 doi:10.1093/bioinformatics/btv351.
528

- 529 27. Aubry S, Kelly S, Kumpers BMC, Smith-Unna RD, Hibberd JM. Deep Evolutionary Comparison
530 of Gene Expression Identifies Parallel Recruitment of Trans-Factors in Two Independent Origins of C4
531 Photosynthesis. PLOS Genetics 2014 doi:10.1371/journal.pgen.1004365.
532
- 533 28. Parra G, Bradnam K, Ning Z, Keane T, Korf I. Assessing the gene space in draft genomes. Nucl.
534 Acids Res. 2009;37:289–297 doi:10.1093/nar/gkn916.
535
- 536 29. One Thousand Plant Transcriptomes Initiative. Data Resources for One Thousand Plant
537 Transcriptomes Elucidate Green Plant Phylogenomics. CyVerse Data Commons. 2019.
538 <https://doi.org/10.25739/8m7t-4e85>
539

540
541

542 Figure Captions:

543

544 Fig. 1. Panel A provides an overview of the total sequence percentage verified to be of contaminant
545 origin (red), or inferred to be possible contaminants in other sequence libraries (grey) in all 1KP
546 libraries, and libraries inferred to be contaminated through the 18S phylogenetic placement. Panel B
547 lists 21 libraries in which more than 6% of the total sequences are potential contaminants. Panel C
548 shows a heatmap of inferred contaminant interactions between pairs of species; contaminated species
549 are shown on the vertical axis, and contaminating species on the horizontal axis.

550

551

552 Fig. 2. Fraction of the gene sets found (complete + fragments) versus the number of scaffolds
553 (300+ bp) in the assemblies. For each sample, the fraction of the eukaryota and embryophyta sets
554 found in the assemblies are calculated with BUSCO and the fraction of the CEGMA 248 set with the
555 CRBB tool. All three sets are more completely recovered at higher scaffold counts, but the BUSCO
556 embryophyta set is less complete in our samples.



