1    **Targeted sequence capture array for phylogenetics and population genomics in the**

2    **Salicaceae[1]**

3    Brian J. Sanderson[2,3,6], Stephen P. DiFazio[3], Quentin C. Cronk[4], Tao Ma[5], Matthew S. Olson[2]

4    [2]Department of Biological Sciences, Texas Tech University, Lubbock, TX 79409-3131 USA

5    [3]Department of Biology, West Virginia University, Morgantown, WV, 26506 USA

6    [4]Department of Botany, University of British Columbia, Vancouver, BC, V6T 1Z4 Canada

7    [5]Key Laboratory of Bio-Resource and Eco-Environment of Ministry of Education, College of

8    Life Sciences, Sichuan University, Chengdu 610065, People's Republic of China

9

10    Email addresses: BJS: brian@biologicallyrelevant.com

11                    SPD: spdifazio@mail.wvu.edu

12                    QC: quentin.cronk@ubc.ca

13                    TM: matao.yz@gmail.com

14                    MSO: matt.olson@ttu.edu

15    [6]Author for correspondence: brian@biologicallyrelevant.com

16

17    Number of words: 3831

18    [1]Manuscript received: _____; revision accepted _____.

19   **Abstract**

20   •   **Premise of the study**: The family Salicaceae has proved taxonomically challenging,

21       especially in the genus *Salix*, which is speciose and features frequent hybridization and

22       polyploidy. Past efforts to reconstruct the phylogeny with molecular barcodes have failed to

23       resolve the species relationships of many sections of the genus.

24   •   **Methods**: We used the wealth of sequence data in the family to design sequence capture

25       probes to target regions of 300-1200 base pairs of exonic regions of 972 genes.

26   •   **Results**: We recovered sequence data for nearly all of the targeted genes in three species of

27       *Populus* and three species of *Salix*. We present a species tree, discuss concordance among

28       gene trees, as well as some population genomic summary statistics for these loci.

29   •   **Conclusions**: Our sequence capture array has extremely high capture efficiency within the

30       genera *Populus* and *Salix*, resulting in abundant phylogenetic information. Additionally,

31       these loci show promise for population genomic studies.

32   **Key words**: Phylogenetics; *Populus*; Salicaceae; *Salix*; targeted sequence capture

33 **Introduction**

34 Although the cost of whole-genome sequencing has continued to dramatically decrease over the

35 past decade, the cost and complexity of whole-genome analyses still limit their utility and

36 accessibility for answering evolutionary questions in novel taxa (Richards, 2018). However, a

37 polished genome assembly is not necessary to address many questions. In this context, several

38 methods have been developed to reduce the cost and effort required to obtain genomic

39 information in novel species (McKain et al., 2018). The recent development of targeted sequence

40 capture presents an affordable method for consistently isolating specific, long, phylogenetically

41 informative regions in the taxa of interest (Gnirke et al., 2009; Mamanova et al., 2010; Hale et

42 al., 2020). Targeted sequence capture uses biotinylated RNA baits to target prepared sequencing

43 library fragments. The baited library fragments can then be pulled out of solution with

44 streptavidin-coated magnetic beads to selectively enrich the fragments that contain loci of

45 interest, while discarding the majority of library fragments that do not. This method offers many

46 advantages over other methods of genome sequence partitioning, such as genome skimming and

47 RAD-seq. It does not necessarily depend on a highly polished, annotated reference genome.

48 Additionally, the same loci can be consistently sequenced at a high depth across individuals

49 without requiring comprehensive, concurrent sequencing of all individuals (Mamanova et al.,

50 2010; Grover et al., 2012; Jones and Good, 2016).

51 In this paper, we report on the design and implementation of a targeted sequence capture

52 array to collect data for phylogenetic analysis within the Salicaceae, the plant family that

53 includes poplars and willows. Understanding species relationships within this family, and in

54 particular the genus *Salix*, has presented challenges to taxonomists as early as Linnaeus, who

55 noted that "species of this genus are extremely difficult to clarify" (Linnaeus, 1753; Skvortsov,

56    1999). *Salix* species present challenges to classification due to their wide geographic ranges,

57    hysteranthous phenology, extensive interspecific hybridization, polyploidy, and the lack of well-

58    defined and variable flower characters for morphological circumscription of taxa (Raup, 1959;

59    Skvortsov, 1999; Percy et al., 2014; Wang et al., 2020). Species of *Salix* exhibit holarctic

60    distributions, and there are several classifications which differ among continents and are

61    challenging to synthesize due to non-overlapping taxonomic treatment of species (Dickmann and

62    Kuzovkina, 2014). Past efforts to reconstruct the phylogeny of *Salix* using nuclear AFLP and

63    plastid barcode sequences have resulted in a lack of clearly resolved species relationships,

64    especially in the subgenus *Vetrix* (Trybush et al., 2008; Percy et al., 2014). A more recent study

65    using a supermatrix approach with RAD-seq data showed resolution within a subset of species of

66    the subgenera *Vetrix* and *Chamaetia*, highlighting the potential of large-scale molecular data to

67    resolve this phylogenetically challenging group (Wagner et al., 2018).

68           The utility of RAD-seq for collecting data for phylogeny, however, is limited by several

69    issues. First, RAD-seq does not consistently screen homologous regions across species and

70    across different experiments, which limits its utility for adding species to a phylogeny at a later

71    time. Second, because RAD-seq assesses diversity in very short segments of the genome, the

72    concatenation of this sort of data and supermatrix approaches are required for its use in

73    phylogenetic analyses requires (de Queiroz and Gatesy, 2007), which does not allow separate

74    exploration of gene and species phylogenies using super tree methods (Sanderson et al., 1998).

75    Additionally, concatenation approaches are likely to exacerbate problems associated with

76    maximum-likelihood methods for species with rapid diversification (Edwards et al., 2007;

77    Edwards, 2009). Targeted sequence capture does not have these limitations, and thus may be a

78    more appropriate genotyping platform for phylogenetics.

79    Species of *Populus* and *Salix* have been of great interest for the development of forestry

80    and biofuel products, resulting in polished reference genomes for *P. trichocarpa*, *P. tremula*, *P.*

81    *euphratica*, *S. purpurea* and *S. suchowensis*., as well as shallow resequencing data for many

82    additional species (Tuskan et al., 2018). Our design strategy leveraged this abundance of existing

83    genomic information to quantify polymorphism and the distribution of insertion- deletion events

84    within and among species in order to maximize capture efficiency. Additionally, because we

85    consistently target regions of exons we are able to characterize the nucleotide-site degeneracy

86    with these data to quantify population genomic summary statistics. We demonstrate the utility of

87    this resource for *Populus* and *Salix* species by presenting a fully resolved phylogenetic tree for

88    six species and an outgroup, and by estimating the distribution of nucleotide diversity within

89    species for our targeted genes.

90    **Methods**

91    *Probe Design*

92    Our goal was to identify regions that could be efficiently captured using RNA bait hybridization

93    for diverse species across the family Salicaceae. The family Salicaceae is thought to have

94    diverged from other clades approximately 92.5 Mya (Zhang et al., 2018b). Our primary focus

95    was on the genera *Populus* and *Salix*, which diverged approximately 48 Mya, and the species

96    *Idesia polycarpa* Maxim, which diverged from other clades approximately 56 Mya, which we

97    use as an outgroup (Zhang et al., 2018b). Although we were interested in using these probes for

98    phylogenetics with both *Populus* and *Salix* species, we focused on maximizing capture efficiency

99    for the species in *Salix*, because the phylogeny for *Populus* is already much better resolved than

100   that for *Salix* (Trybush et al., 2008; Wang et al., 2014, 2020; Percy et al., 2014; Liu et al., 2017).

101   For this reason, the capture baits were designed to target regions in *Salix pupurea* that also would

102    have high capture efficiency across the Salicaceae. The efficiency of RNA bait binding, and thus

103    capture efficiency, is reduced as target regions diverge due to sequence polymorphism (Lemmon

104    and Lemmon, 2013). To improve capture efficiency, we quantified sequence polymorphism

105    among whole-genome resequencing data from a diverse array of *Populus* and *Salix* species

106    (Table S1). The whole-genome short reads of the *Populus* and *Salix* species were aligned to the

107    *Populus trichocarpa* genome assembly version 3 (Tuskan et al., 2006) using bwa mem v. 0.7.12

108    with default parameters (Li, 2013). We used the *P. trichocarpa* genome as our initial reference

109    because it was the most polished and annotated genome in genus. Variable sites and insertion-

110    deletion mutations (indels) were identified using samtools mpileup (Li, 2011), and read depth for

111    the variant calls was quantified using vcftools (Danecek et al., 2011). Custom Python scripts

112    were used to identify variant and indel frequencies for all exons in the *P. trichocarpa* genome

113    annotation (scripts available at https://github.com/BrianSanderson/phylo-seq-cap; Sanderson,

114    2020).

115        Orthologs for our candidate loci in the in the *Salix purpurea* 94006 genome assembly

116    version 1 (*Salix purpurea* v1.0, DOE-JGI; Carlson et al., 2017; Zhou et al., 2018) were identified

117    using a list of orthologs between the *P. trichocarpa* and *S. purpurea* prepared using a tree-based

118    approach by Phytozome v 12 (Goodstein et al., 2012). We further screened candidate regions to

119    exclude high-similarity duplicated regions by accepting only loci with single BLAST (Camacho

120    et al., 2009) hits against the highly contiguous assembly of *S. purpurea* 94006 version 5 (Zhou et

121    al., 2020), which is less fragmented than the *S. purpurea* 94006 version 1. Genes from the

122    Salicoid whole-genome duplication were identified using MCScanX (Wang et al., 2012), using

123    default parameters and selected those segments for which the average $K_S$ value for paralogous

124    genes was between 0.2 and 0.8. Genes for which at least 600 base pairs (bp) of exon sequence

125    contained 2-12% polymorphism and fewer than two indels were selected for probe design by

126    Arbor Biosciences (Ann Arbor, MI, USA). Probes were designed with 50% overlap across the

127    targeted regions, so that each nucleotide position would potentially be captured by two probes.

128    Finally, to ensure that loci with high divergence across the family would be captured, we

129    identified targets with less than 95% identity (based on BLAST results) between *S. purpurea* and

130    *P. trichocarpa* and designed supplementary probes from orthologs of these genes in the *Idesia*

131    *polycarpa* genome.

132    *Library Preparation and Sequence Capture*

133    Libraries for two individuals from each *Populus balsamifera* L., *P. tremula* L., *P. mexicana*

134    Wesmael., *Salix nigra* Marshall, *S. exigua* Nutt., and *S. phlebophylla* Andersson (Table S3) were

135    prepared using the NEBNext Ultra II DNA Prep Kit following the published protocol for this kit

136    (New England Biolabs, Ipswitch, MA, USA), and quantified using an Agilent Bioanalyzer 2100

137    DNA 1000 kit (Agilent Technologies, Santa Clara, CA, USA). Libraries were pooled at

138    equimolar concentrations into two pools of six prior to probe hybridization following the Arbor

139    Biosciences myBaits protocol v 3.0.1 and Hale et al. (2020). The hybridized samples were

140    subsequently pooled at equimolar ratios and sequenced at the Texas Tech Center for

141    Biotechnology and Genomics using a MiSeq with the Micro chemistry and 150 bp paired-end

142    reads (Illumina, Inc., San Diego, CA, USA).

143    *Analysis of Sequence Capture Data*

144    Read data was trimmed for primer sequences and low quality scores using Trimmomatic v. 0.36

145    (Bolger et al., 2014). The trimmed read data, as well as the whole-genome reads for *I. polycarpa*,

146    were assembled into gene sequences using the HybPiper pipeline (Johnson et al., 2016). We

147    estimated the depth of read coverage across all targeted genes as well as at off-target sites in R.

148    The assembled amino acid sequences were aligned with mafft v. 7.310 with the parameters –

149    localpair and –maxiterate 1000 (Katoh and Standley, 2013), converted into codon-aligned

150    nucleotide alignments with pal2nal v. 14 (Suyama et al., 2006), and trimmed for quality and

151    large gaps with trimal v. 1.4.rev15 with the parameter -gt 0.5 (Capella-Gutiérrez et al., 2009).

152         HybPiper provides warnings for genes that have multiple competing assemblies that are

153    within 80% of the length of the target region, because the alternate alignments may indicate that

154    those genes have paralogous copies in the genome. We estimated phylogenetic relationships

155    using the full set of gene sequences recovered from our sequence capture data, as well as a

156    restricted set of putatively single copy genes, based on our *a priori* list of paralogs between *S.*

157    *purpurea* and *P. tricocarpa*, and supplemented by the list of paralog warnings from HybPiper.

158         We estimated gene trees using RAxML v. 8.2.10, specifying a GTR$\Gamma$ model of sequence

159    evolution (Stamatakis, 2014). A set of 250 bootstrap replicates was generated for each gene tree.

160    We used ASTRAL-III to infer the species tree from the RAxML gene trees (Zhang et al., 2018a;

161    Rabiee et al., 2019). Because all nodes are weighted equally during quartet decomposition in

162    ASTRAL-III, we used sumtrees in the Python package DendroPy v. 4.4.0 to collapse nodes with

163    less than 33% bootstrap support values prior to species tree estimation (Sukumaran and Holder,

164    2010). A set of 100 multilocus bootstrap replicates was generated for the species tree. We used

165    phyparts to determine the extent of congruence among gene trees for each node in the species

166    tree (Smith et al., 2015). Cladograms representing the gene tree congruence and alternate

167    topologies were plotted with the scripts phypartspiecharts.py and minority_report.py, written by

168    Matt Johnson (scripts available at https://github.com/mossmatters/phyloscripts).

169         Finally, we used custom Python scripts to quantify nucleotide diversity at synonymous

170    and non-synonymous sites between the individuals of the same species, as well as correlations in

171    values of per-site nucleotide diversity between all species. The scripts described above as well as

172    the full details of these analyses including are available in Jupyter notebooks at

173    https://github.com/BrianSanderson/phylo-seq-cap (Sanderson, 2020).

174    **Results**

175    *Sequence capture efficiency*

176    The final capture kit targets 972 genes covered by 12,951 probes based on the *S. purpurea*

177    reference, and an additional 7049 (redundant) probes based on the *I. polycarpa* genome that

178    target genes with the highest divergence between *S. purpurea* and *P. trichocarpa* identified by

179    Phytozome. This included an average of $680 \pm 309$ (mean $\pm$ sd) probes on each *S. purpurea*

180    chromosome (Table S2), with an average of $1098 \pm 489$ (mean $\pm$ sd) bp of exon sequence per

181    gene. Of the 972 target genes, 593 are putatively single copy based on our identification of

182    paralogs in the *S. purpurea* genome assembly, 142 represent pairs of paralogs from the shared

183    Salicoid whole-genome duplication (i.e. 71 pairs of genes), and 237 are genes that have known

184    paralogs for which we were not able to design targets in this kit (i.e. each of these genes has one

185    or more paralogs in the *S. purpurea* genome that is not targeted by probes). We included a total

186    of 1219 genes in the target file used to assemble the capture data, which includes the 972

187    targeted genes as well as paralogous copies for which probes were not designed. Because the

188    issues of paralogy become more complex when we add species other than *S. purpurea* and *P.*

189    *trichocarpa*, we advise using the HybPiper warnings of multiple competing long assemblies to

190    assess paralogy in novel species following guidance here

191    https://github.com/mossmatters/HybPiper/wiki/Paralogs. The sequences of the capture probes as

192    well as the target reference file are accessible at https://github.com/BrianSanderson/phylo-seq-

193    cap (Sanderson, 2020). The sequence capture kit is available from Arbor Biosciences

194    (Ref#170424-30 "Salicaceae").

195         Sequence capture efficiency was high among the libraries. We recovered 805,820 $\pm$

196    178,482 reads (mean $\pm$ sd) from our *Populus* and *Salix* target capture libraries, of which 86.7 $\pm$

197    1.15% (mean $\pm$ sd) mapped to the target sequence reference (Table 1). An average of 94.48 $\pm$

198    1.37% of targeted exon sequences were covered by $\geq$ 10 reads. The average read depth was

199    44.65 $\pm$ 1.61 for on-target sites, and 14.48 $\pm$ 2.10 for off-target sites (Table S4).

200    *Phylogenetics*

201    The species tree estimated with putatively single copy genes correctly paired all individuals of

202    the same species and revealed a fully resolved phylogeny for the *Populus* and *Salix* species with

203    100% multilocus bootstrap support for all nodes (Fig. 1A). At least 85% of gene trees support the

204    topology of the species tree (Fig. 1B), with the exceptions of the bipartition that separates *P.*

205    *balsamifera* and *P. tremula*, and the bipartition that separates *S. phlebophylla* from the other

206    *Salix* species, which had dominant alternate topologies that were supported by a large number of

207    gene trees (Figs. S1 and S2). The topology of the species tree estimated with the full set of genes

208    and known paralogs was nearly identical to the tree estimated with only the putatively single

209    copy genes. The major difference between these trees was evident in the bipartition separating *P.*

210    *balsamifera* and *P. tremula*, where there were a large number of alternative topologies supported

211    by small numbers of gene trees (the top 3 were supported by 13, 11, and 10 gene trees; Fig. S3).

212    *Population genomics*

213    Patterns of nucleotide diversity, measured as Nei's $\pi$ (Nei and Li, 1979), varied among species,

214    with the greatest variation at synonymous sites (Fig. S4; Table S5). *P. tremula* had the highest

215    average values of $\pi$ at both synonymous and non-synonymous sites (Fig. 2). The values of $\pi$

216　among species were highly correlated for species within genus and exhibited lower correlations

217　between genera (Fig. 3).

218　**Discussion**

219　The decreasing cost of obtaining genomic and transcriptomic sequence data holds great promise

220　for unlocking our understanding of phylogenetic relationships and population genetic patterns

221　within and among complex taxonomic groups. However, assembling complete genomes is still

222　not a trivial task, and there exist relatively few polished plant reference genomes onto which

223　genome skimming data can be mapped. Many methods have been developed to reduce the

224　sequencing and analytical burdens associated with obtaining genome data. We believe that

225　targeted sequence capture is one of the most promising contemporary methods of inexpensively

226　generating genomic information.

227　　　The efficiency of our targeted sequence capture array was extremely high, which yielded

228　abundant phylogenetic information for six species of *Populus* and *Salix*. Overall, the phylogeny

229　was fully resolved and conformed to our general understanding the relationships among the taxa

230　(Wu et al., 2015; Wang et al., 2020). One strength of the sequence capture approach is that it

231　provides sufficiently long contiguous segments of gene sequences to assemble gene trees and it

232　can overcome the problems introduced by concatenation of multiple gene regions with divergent

233　histories (Edwards et al., 2007; Edwards, 2009). The super tree approach also allowed for the

234　identification of alternative evolutionary histories that are supported by different regions of the

235　genome, as often occurs during historical hybridization and introgression (Zhang et al., 2018a;

236　Rabiee et al., 2019). Our species tree identified three alternative gene tree relationships among

237　the three *Populus* species (Fig. S1). Previous studies have provided evidence of historical

238　introgression among these species, including a history of chloroplast capture and hybridization

239    between *P. mexicana* and species in the section *Tacamahaca* (including *P. balsamifera*; Wang et

240    al., 2014, 2020; Liu et al., 2017). The second most supported alternative topology that we

241    recovered placed *P. mexicana* and *P. tremula* as sister taxa, a pattern that does not support this

242    hypothesis, likely due to incomplete lineage sorting (Wang et al., 2020). *P. tremula* likely has a

243    greater long-term effective population size than *P. balsamifera* (Wang et al., 2016), and so

244    coalescence times may be shorter on average in *P. balsamifera*. Among the *Salix* species, we

245    identified three alternative gene tree relationships between the *S. phlebophylla* and *S. exigua*

246    individuals, which may reflect the histories of rapid speciation and hybridization that have long

247    vexed attempts at phylogenetic reconstruction in the genus *Salix* (Fig. S2; Trybush et al., 2008;

248    Percy et al., 2014). Both of these patterns in *Populus* and *Salix* may be better understood once

249    additional taxa are added to this phylogeny.

250        We have also shown that this sequence capture design can be applied to address questions

251    related to population genomics in the Salicaceae. Many of the advantages of targeted sequence

252    capture over competing methods are of particular relevance for population genomics studies,

253    including specific knowledge of loci being sequenced, the ability to differentiate among

254    synonymous, non-synonymous, intronic, and intergenic loci, and the ability to collect data on the

255    same set of loci across different experiments, either within species or across species, for

256    comparative studies. In particular, synonymous sites, especially four-fold synonymous sites, are

257    among the fastest evolving regions of the genome and the sites within genic regions least

258    influenced by selection (Wright and Andolfatto, 2008), and are thus among the best regions for

259    estimating patterns of historical demography. Our estimates of nucleotide diversity are similar to

260    those that have been previously reported for *P. balsamifera* and *P. tremula* using Sanger

261    sequencing data (Ingvarsson, 2005; Olson et al., 2010) and whole-genome sequencing data

262    (Wang et al., 2016). The high estimates of diversity in *S. phlebophylla* compared to the other two

263    *Salix* species is curious and may result from a history with relatively little migration due to the

264    absence of glaciation over a large portion of its Beringian distribution (Hultén, 1937).

265        The current study is based on a small sample size per species (n = 2), and so our ability to

266    account for population structure or robustly perform population genomic inferences with these

267    data is limited. Additionally, a potential limitation for using this sequence capture array for

268    comparative population genomics is that we screened loci for a range of among-species

269    variability between 2-12%, which excludes loci that exhibit extremely high or low values of

270    nucleotide diversity. This may bias estimates of nucleotide diversity arising from these probes

271    toward greater evenness. The ability to identify synonymous sites, which are the closest to

272    neutral among all classes of sites (Wright and Andolfatto, 2008), should partially address this

273    bias. Another feature of sequence capture data is the recovery of "off-target" sequences that

274    result from the fact that the insert size of libraries is larger than the 120 bp bait length, and so

275    regions upstream and downstream of the target will be sequenced as well. These regions may

276    include intronic and intergenic regions, as well as exonic sequences that deviate from the

277    constraints we used for our design. The results we report here only incorporate the "on-target"

278    sites that we sequenced, but HybPiper implements methods to assemble intronic sequences as

279    well. However, the potential effects of hitchhiking selection on synonymous site variation will

280    likely remain apparent.

281        We also found that it was straightforward to integrate the targeted sequence capture data

282    with whole-genome sequence data using the HybPiper pipeline by simply including the FASTQ

283    files from whole-genome reads in the pipeline. This strategy was used to successfully incorporate

284    whole-genome sequencing data from *Idesia polycarpa*, to act as our outgroup. The proportion of

285    gene coverage as well as the read depth for the *I. polycarpa* data was similar to the sequence

286    capture libraries (Table 1).

287         A whole genome duplication occurred prior to the divergence of *Salix* and *Populus*, and

288    there are at least 8000 known paralog pairs in the *P. trichocarpa* reference genome (Tuskan et

289    al., 2006). Genes with paralogous copies in the genome can complicate gene assemblies, because

290    sequence data from both copies may alternately align to the same target sequence. We identified

291    paralogous sequences in the *S. purpurea* genome assembly using MCScanX, and used that

292    information to assist in the design the sequence capture array. The final array includes 593

293    putatively single copy genes, 142 pairs of paralogs, and 237 genes which have paralogs but for

294    which we were not able to include both paralogs in the kit due our selection criteria. The target

295    reference file we used to map the sequence capture data thus includes 1219 genes including the

296    single copy and known paralogs from *S. purpurea*. In addition to this, HybPiper provides

297    warnings for genes that have multiple competing alignments that cover the majority of the target

298    sequence, which may indicate the presence of multiple paralogous copies in the genome

299    (Johnson et al., 2016). This will be particularly useful because the genes that have maintained

300    paralogous copies are likely to differ among species throughout the diversification of willows.

301    We estimated evolutionary relationships using both the full set of 1219 single copy and known

302    paralog genes, as well as a limited set of just single copy genes that did not report paralog

303    warnings. The results from both analyses were nearly the same, but this will likely not be true for

304    a more complex phylogenetic analysis that includes more than six species and an outgroup. For

305    those more complex phylogenetic analyses, the ability to compare trees constructed with single

306    copy genes with those using paralogous copies may provide crucial information for reconciling

307    evolutionary relationships.

308    This sequence capture array will provide the community with an excellent resource to

309    consistently sequence a set of variable regions of the genome for phylogenetic and population

310    genomic investigations in the Salicaceae. The rate of read mapping and coverage of target genes

311    was remarkably consistent across both genera, despite the fact that the taxa were selected to

312    maximize sampling of phylogenetic diversity within each genus. The Salicaceae are important

313    plants in the northern hemisphere both ecologically and economically and have been the subjects

314    of numerous population genetics and genomics investigations of speciation, hybridization,

315    introgression, selection, and historical population size and migration. This resource will allow

316    phylogenetic and comparative population genomic studies to assess the same loci across different

317    studies, which will allow us to build a worldwide diversity database and facilitate more precise

318    comparative research questions. Our results demonstrate that the rate of gene capture is

319    extremely high, such that it would be unnecessary to filter data and determine appropriate

320    overlapping genotype thresholds, as is necessary with random genome partitioning methods such

321    as RAD-seq.

322    **Acknowledgements**

331    Agriculture and Agri-Food. Agriculture and Agri-Food Canada retains complete ownership of

332    the resources presented here.

333    **Author Contributions**

334    S.P.D. and M.S.O. conceived the study. S.P.D., Q.C.C., T.M., and M.S.O. secured funding to

335    support the project. B.J.S. and S.P.D. designed the sequence capture array. Q.C.C. and T.M.

336    provided whole genome sequence data. B.J.S and M.S.O. prepared and sequenced the DNA

337    samples, analyzed the data, interpreted the results, and wrote the manuscript. All authors edited

338    drafts of the manuscript.

339    **Data Accessibility**

340    Accession numbers for all sequence data used to design the sequence capture array are presented

341    in Table S1. The raw reads of targeted sequence capture data from the six species of *Populus* and

342    *Salix* are available on the NCBI sequence read archive under the BioProject accession number

343    PRJNA627181. The raw reads of the *Idesia polycara* whole genome sequence data are available

344    in the Genome Warehouse of the Beijing Institute of Genomics (BIG), under the accession

345    number PRJCA002959. The sequences of the probes that were designed, all of the custom

346    Python scripts that were used for this study, and the full details of analyses summarized in

347    notebooks are available at https://github.com/BrianSanderson/phylo-seq-cap (Sanderson, 2020).

348    **Literature Cited**

349    BOLGER, A.M., M. LOHSE, and B. USADEL. 2014. Trimmomatic: A flexible trimmer for Illumina

350    sequence data. *Bioinformatics* 30: 2114–2120.

351    CAMACHO, C., G. COULOURIS, V. AVAGYAN, N. MA, J. PAPADOPOULOS, K. BEALER, and T.L.

352    MADDEN. 2009. BLAST+: Architecture and applications. *BMC Bioinformatics* 10: 421–421.

353     CAPELLA-GUTIÉRREZ, S., J.M. SILLA-MARTÍNEZ, and T. GABALDÓN. 2009. trimAl: A tool for

354     automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25: 1972–

355     1973.

356     CARLSON, C.H., Y. CHOI, A.P. CHAN, M.J. SERAPIGLIA, C.D. TOWN, and L.B. SMART. 2017.

357     Dominance and Sexual Dimorphism Pervade the Salix purpurea L. Transcriptome. *Genome*

358     *Biology and Evolution* 9: 2377–2394.

359     DANECEK, P., A. AUTON, G. ABECASIS, C.A. ALBERS, E. BANKS, M.A. DEPRISTO, R.E.

360     HANDSAKER, ET AL. 2011. The variant call format and VCFtools. *Bioinformatics* 27: 2156–2158.

361     DE QUEIROZ, A., and J. GATESY. 2007. The supermatrix approach to systematics. *Trends in*

362     *Ecology and Evolution* 22: 34–41.

363     DICKMANN, D.I., and J. KUZOVKINA. 2014. Poplars and Willows of the World, With Emphasis

364     on Silviculturally Important Species. *In* J. Isebrands, and J. Richardson [eds.], Poplars and

365     Willows Trees for Society and the Environment, 8–91. The Food and Agriculture Organization

366     of the United Nations, Rome, Italy.

367     DOE-JGI. Salix purpurea version 1.

368     EDWARDS, S.V. 2009. Is a new and general theory of molecular systematics emerging? *Evolution*

369     63: 1–19.

370     EDWARDS, S.V., L. LIU, and D.K. PEARL. 2007. High-resolution species trees without

371     concatenation. *Proceedings of the National Academy of Sciences* 104: 5936–5941.

372   GNIRKE, A., A. MELNIKOV, J.R. MAGUIRE, P. ROGOV, E.M. LEPROUST, W. BROCKMAN, T.J.

373   FENNELL, ET AL. 2009. Solution hybrid selection with ultra-long oligonucleotides for massively

374   parallel targeted sequencing. *Nature Biotechnology* 27: 182–189.

375   GOODSTEIN, D.M., S. SHU, R. HOWSON, R. NEUPANE, R.D. HAYES, J. FAZO, T. MITROS, ET AL.

376   2012. Phytozome: A comparative platform for green plant genomics. *Nucleic Acids Research* 40:

377   D1178–D1186.

378   GROVER, C.E., A. SALMON, and J.F. WENDEL. 2012. Targeted sequence capture as a powerful

379   tool for evolutionary analysis. *American Journal of Botany* 99: 312–319.

380   HALE, H., E.M. GARDNER, J. VIRUEL, L. POKORNY, and M.G. JOHNSON. 2020. Strategies for

381   reducing per-sample costs in target capture sequencing for phylogenomics and population

382   genomics in plants: Low-cost Hyb-Seq. *Applications in Plant Sciences* e11337.

383   HULTÉN, E. 1937. Outline of the history of arctic and boreal biota during the Quarternary period :

384   Their evolution during and after the glacial period as indicated by the equiformal progressive

385   areas of present plant species. Bokförlags aktiebolaget Thule, Stockholm, Sweden.

386   INGVARSSON, P.K. 2005. Nucleotide polymorphism and linkage disequilibrium within and

387   among natural populations of European aspen (Populus tremula L., Salicaceae). *Genetics* 169:

388   945–953.

389   JOHNSON, M.G., E.M. GARDNER, Y. LIU, R. MEDINA, B. GOFFINET, A.J. SHAW, N.J.C. ZEREGA,

390   and N.J. WICKETT. 2016. HybPiper: Extracting Coding Sequence and Introns for Phylogenetics

391   from High-Throughput Sequencing Reads Using Target Enrichment. *Applications in Plant*

392   *Sciences* 4: 1600016–1600016.

393    JONES, M.R., and J.M. GOOD. 2016. Targeted capture in evolutionary and ecological genomics.

394    *Molecular Ecology* 25: 185–202.

395    KATOH, K., and D.M. STANDLEY. 2013. MAFFT Multiple Sequence Alignment Software

396    Version 7: Improvements in Performance and Usability. *Molecular Biology and Evolution* 30:

397    772–780.

398    LEMMON, E.M., and A.R. LEMMON. 2013. High-Throughput Genomic Data in Systematics and

399    Phylogenetics. *Annual Review of Ecology, Evolution, and Systematics* 44: 99–121.

400    LI, H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.

401    00: 1–3.

402    LI, H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping

403    and population genetical parameter estimation from sequencing data. *Bioinformatics* 27: 2987–

404    2993.

405    LINNAEUS, C. 1753. Species plantarum. Impensis Laurentii Salvii, Stockholm.

406    LIU, X., Z. WANG, W. SHAO, Z. YE, and J. ZHANG. 2017. Phylogenetic and Taxonomic Status

407    Analyses of the Abaso Section from Multiple Nuclear Genes and Plastid Fragments Reveal New

408    Insights into the North America Origin of Populus (Salicaceae). *Frontiers in Plant Science* 7: 1–

409    9.

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

411    ET AL. 2010. Target-enrichment strategies for next-generation sequencing. *Nature Methods* 7:

412    111–118.

413    MCKAIN, M.R., M.G. JOHNSON, S. URIBE-CONVERS, D. EATON, and Y. YANG. 2018. Practical

414    considerations for plant phylogenomics. *Applications in Plant Sciences* 6: e1038–e1038.

415    NEI, M., and W.H. LI. 1979. Mathematical model for studying genetic variation in terms of

416    restriction endonucleases. *Proceedings of the National Academy of Sciences* 76: 5269–5273.

417    OLSON, M.S., A.L. ROBERTSON, N. TAKEBAYASHI, S. SILIM, W.R. SCHROEDER, and P. TIFFIN.

418    2010. Nucleotide diversity and linkage disequilibrium in balsam poplar (Populus balsamifera).

419    *New Phytologist* 186: 526–536.

420    PERCY, D.M., G.W. ARGUS, Q.C. CRONK, A.J. FAZEKAS, P.R. KESANAKURTI, K.S. BURGESS,

421    B.C. HUSBAND, ET AL. 2014. Understanding the spectacular failure of DNA barcoding in willows

422    ( Salix ): Does this result from a trans-specific selective sweep? *Molecular Ecology* 23: 4737–

423    4756.

424    RABIEE, M., E. SAYYARI, and S. MIRARAB. 2019. Multi-allele species reconstruction using

425    ASTRAL. *Molecular Phylogenetics and Evolution* 130: 286–296.

426    RAUP, H.M. 1959. The willows of boreal Western America. *Contributions from the Gray*

427    *Herbarium of Harvard University* 185: 3–95.

428    RICHARDS, S. 2018. Full disclosure: Genome assembly is still hard. *PLOS Biology* 16:

429    e2005894–e2005894.

430    SANDERSON, B.J. 2020. BrianSanderson/phylo-seq-cap: Publication (Version v1.0). Zenodo.

431    http://doi.org/10.5281/zenodo.3979562.

432    SANDERSON, M.J., A. PURVIS, and C. HENZE. 1998. Phylogenetic supertrees: Assembling the

433    trees of life. *Trends in Ecology and Evolution* 13: 105–109.

434    SKVORTSOV, A.K. 1999. Willows of Russia and Adjacent Countries. I. N. (. Kadis, G. R. Argus,

435    and A. G. Zinovjev [eds.], University of Joensuu, Joensuu.

436    SMITH, S.A., M.J. MOORE, J.W. BROWN, and Y. YANG. 2015. Analysis of phylogenomic datasets

437    reveals conflict, concordance, and gene duplications with examples from animals and plants.

438    *BMC Evolutionary Biology* 15: 150–150.

439    STAMATAKIS, A. 2014. RAxML version 8: A tool for phylogenetic analysis and post-analysis of

440    large phylogenies. *Bioinformatics* 30: 1312–1313.

441    SUKUMARAN, J., and M.T. HOLDER. 2010. DendroPy: A Python library for phylogenetic

442    computing. *Bioinformatics* 26: 1569–1571.

443    SUYAMA, M., D. TORRENTS, and P. BORK. 2006. PAL2NAL: Robust conversion of protein

444    sequence alignments into the corresponding codon alignments. *Nucleic Acids Research* 34:

445    W609–W612.

446    TRYBUSH, S., Š. JAHODOVÁ, W. MACALPINE, and A. KARP. 2008. A genetic study of a Salix

447    germplasm resource reveals new insights into relationships among subgenera, sections, and

448    species. *Bioenergy Research* 1: 67–79.

449    TUSKAN, G.A., S.P. DIFAZIO, S. JANSSON, J. BOHLMANN, I. GRIGORIEV, U. HELLSTEN, N.

450    PUTNAM, ET AL. 2006. The Genome of Black Cottonwood, Populus trichocarpa (Torr. & Gray).

451    *Science* 313: 1596–1604.

452    TUSKAN, G.A., A.T. GROOVER, J. SCHMUTZ, S.P. DIFAZIO, A. MYBURG, D. GRATTAPAGLIA, L.B.

453    SMART, ET AL. 2018. Hardwood Tree Genomics: Unlocking Woody Plant Biology. *Frontiers in*

454    *Plant Science* 9: 1799–1799.

455    WAGNER, N.D., S. GRAMLICH, and E. HÖRANDL. 2018. RAD sequencing resolved phylogenetic

456    relationships in European shrub willows (Salix L. Subg. Chamaetia and subg. Vetrix) and

457    revealed multiple evolution of dwarf shrubs. *Ecology and Evolution* 8: 8243–8255.

458    WANG, J., N.R. STREET, D.G. SCOFIELD, and P.K. INGVARSSON. 2016. Natural Selection and

459    Recombination Rate Variation Shape Nucleotide Polymorphism Across the Genomes of Three

460    Related Populus Species. *Genetics* 202: 1185–1200.

461    WANG, M., L. ZHANG, Z. ZHANG, M. LI, D. WANG, X. ZHANG, Z. XI, ET AL. 2020.

462    Phylogenomics of the genus *Populus* reveals extensive interspecific gene flow and balancing

463    selection. *New Phytologist* 225: 1370–1382.

464    WANG, Y., H. TANG, J.D. DEBARRY, X. TAN, J. LI, X. WANG, T.-H. LEE, ET AL. 2012. MCScanX:

465    A toolkit for detection and evolutionary analysis of gene synteny and collinearity. *Nucleic Acids*

466    *Research* 40: e49–e49.

467    WANG, Z., S. DU, S. DAYANANDAN, D. WANG, Y. ZENG, and J. ZHANG. 2014. Phylogeny

468    reconstruction and hybrid analysis of populus (Salicaceae) based on nucleotide sequences of

469    multiple single-copy nuclear genes and plastid fragments. *PLoS ONE* 9:.

470    WRIGHT, S.I., and P. ANDOLFATTO. 2008. The Impact of Natural Selection on the Genome:

471    Emerging Patterns in Drosophila and Arabidopsis. *Annual Review of Ecology, Evolution, and*

472    *Systematics* 39: 193–213.

473    WU, J., T. NYMAN, D.-C. WANG, G.W. ARGUS, Y.-P. YANG, and J.-H. CHEN. 2015. Phylogeny of

474    Salix subgenus Salix s.l. (Salicaceae): Delimitation, biogeography, and reticulate evolution.

475    *BMC Evolutionary Biology* 15: 31.

476    ZHANG, C., M. RABIEE, E. SAYYARI, and S. MIRARAB. 2018a. ASTRAL-III: Polynomial time

477    species tree reconstruction from partially resolved gene trees. *BMC Bioinformatics* 19: 153–153.

478    ZHANG, L., Z. XI, M. WANG, X. GUO, and T. MA. 2018b. Plastome phylogeny and lineage

479    diversification of Salicaceae with focus on poplars and willows. *Ecology and Evolution* 8: 7817–

480    7823.

481    ZHOU, R., D. MACAYA-SANZ, C.H. CARLSON, J. SCHMUTZ, J.W. JENKINS, D. KUDRNA, A.

482    SHARMA, ET AL. 2020. A willow sex chromosome reveals convergent evolution of complex

483    palindromic repeats. *Genome Biology* 21: 38–38.

484    ZHOU, R., D. MACAYA-SANZ, E. RODGERS-MELNICK, C.H. CARLSON, F.E. GOUKER, L.M.

485    EVANS, J. SCHMUTZ, ET AL. 2018. Characterization of a large sex determination region in Salix

486    purpurea L. (Salicaceae). *Molecular Genetics and Genomics* 293: 1437–1452.

487

488    **Tables**

| Name | Num. Reads | Reads Mapped | Prop. Mapped | Genes Mapped | Genes with 25% Seq | Genes with 50% Seq | Genes with 75% Seq | Genes with 100% Seq |
|---|---|---|---|---|---|---|---|---|
| I_polycarpa_WGS-2 | 223470714 | 1653494 | 0.007 | 971 | 970 | 966 | 944 | 123 |
| P_balsamifera_MGR-01 | 614093 | 523321 | 0.852 | 972 | 971 | 960 | 884 | 122 |
| P_balsamifera_MGR-04 | 769303 | 659712 | 0.858 | 972 | 972 | 965 | 915 | 145 |
| P_mexicana_PM3 | 843032 | 739728 | 0.878 | 972 | 972 | 964 | 917 | 140 |
| P_mexicana_PM5 | 880962 | 768927 | 0.873 | 972 | 972 | 967 | 913 | 142 |
| P_tremula_R01-01 | 749220 | 638002 | 0.852 | 971 | 970 | 960 | 907 | 134 |
| P_tremula_R04-01 | 634625 | 539805 | 0.851 | 971 | 969 | 956 | 876 | 122 |
| S_exigua_SE002 | 1139616 | 998698 | 0.876 | 969 | 969 | 966 | 937 | 229 |
| S_exigua_SE053 | 843120 | 741938 | 0.88 | 969 | 969 | 964 | 928 | 195 |
| S_nigra_SG037 | 1166615 | 1028635 | 0.882 | 971 | 971 | 967 | 932 | 205 |
| S_nigra_SG051 | 602649 | 524993 | 0.871 | 971 | 970 | 961 | 903 | 136 |
| S_phlebophylla_SP15M | 753628 | 651791 | 0.865 | 972 | 972 | 967 | 939 | 204 |
| S_phlebophylla_SP7F | 672975 | 581147 | 0.864 | 972 | 972 | 967 | 925 | 203 |

489

490    **Table 1.** Coverage summary statistics for sequence capture read data. For each library, values

491    represent the number of reads in the sequenced library, the number of those reads that mapped to

492    the reference file for the targeted genes, the proportion of mapped reads, the number of targeted

493    genes (out of 972) that had read data mapped to them, and the number of genes that had 25%,

494    50%, 75%, and 100% of the targeted sequences covered with > 10X reads. Footnote: the

495    I_polycarpa_WGS-2 data is from whole-genome sequencing data, rather than targeted sequence

496    capture, and thus the low percent of read mapping reflects the lack of target enrichment

497    (although the read coverage across targets was comparable to the sequence capture libraries,
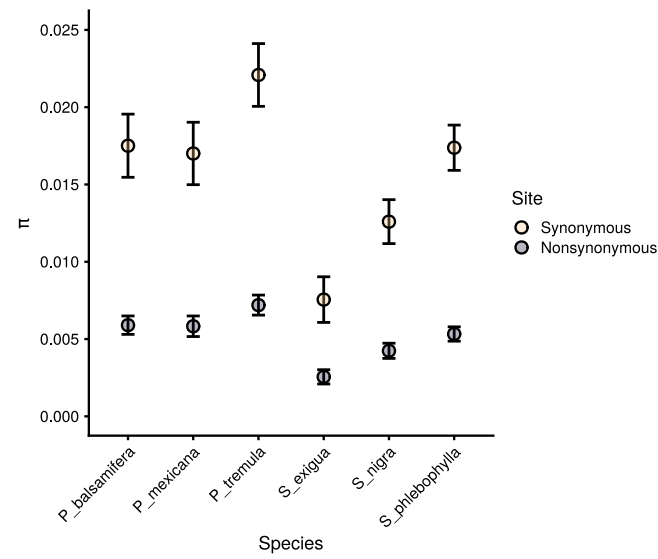
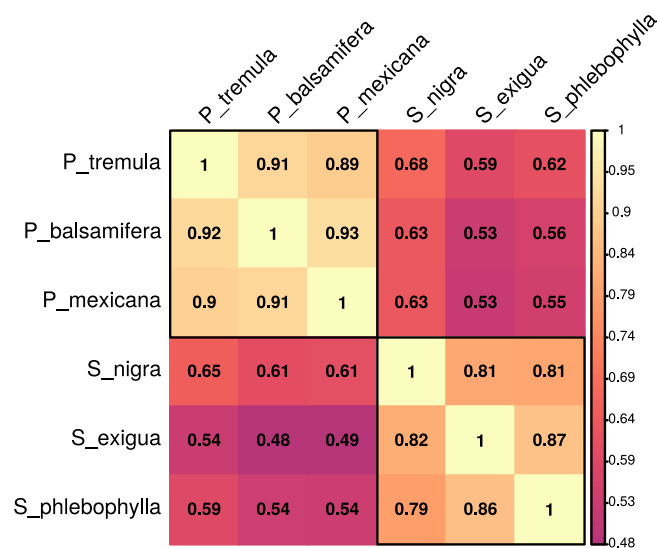498    Table S3).

499 **Figures**



500

501 **Figure 1.** Species trees estimated for the 432 putatively single copy genes that did not have

502 paralog warnings reported by HybPiper. **A)** Species tree generated by ASTRAL-III for the gene

503 trees. Node values represent bootstrap support from 100 multilocus bootstrap replicates in

504 ASTRAL-III. Branch lengths represent coalescent units. **B)** Cladogram showing the congruence

505 of gene trees for all nodes in the ASTRAL-III species tree. The numbers above each node

506 represent the number of gene trees that support the displayed bipartition, and numbers below the

507 node represent the number of gene trees that support all alternate bipartitions. Purple wedges

508 represent the proportion of gene trees that support the displayed bipartition. Blue wedges

509 represent the proportion of gene trees that support a single alternative bipartition (see Figs S1 &

510 S2). Green wedges represent the proportion of gene trees that have multiple conflicting

511 bipartitions. Yellow wedges represent the proportion of gene trees that have no supported

512 bipartition. Plotting code and its interpretation were provided by Matt Johnson (for more detail,

513 see:

514 https://github.com/mossmatters/MJPythonNotebooks/blob/master/PhyParts_PieCharts.ipynb)

515

**Figure 2.** Means and 95% confidence intervals of values of nucleotide diversity (Nei's $\pi$) within

each species at synonymous (yellow) and nonsynonymous (purple) sites.

518

**Figure 3.** Pairwise correlation (Pearson's r) of values of Nei's $\pi$ between all species. Values

520    above the diagonal represent the correlation of $\pi$ at synonymous sites, values below the diagonal

521    represent non-synonymous sites. Black boxes represent within-genus comparisons.

522 **Supporting Information**

523 **Table S1.** Coverage summary statistics for whole-genome reads used to design sequence capture

524 array. For each library, values represent the name of the sequenced individual, the number of

525 reads in the sequenced library, the number of reads that mapped to the *Populus trichocarpa* v3

526 reference genome, the proportion of reads that mapped to the reference genome, and the mean

527 and standard deviation of read depth.

528 **Table S2.** Distribution of probes across the *Salix purpurea* genome.

529 **Table S3.** Collection details for *Populus* and *Salix* species.

530 **Table S4.** Summary of read depth at on- and off-target sites. For each library values represent

531 the 5%, 25%, 50%, 75%, and 99% quantiles of read depth, the maximum number of reads

532 mapped to a site, and the mean and standard deviation of read depth.

533 **Table S5.** Nucleotide diversity expressed as Nei's $\pi$ for nonsynonymous and synonymous sites.

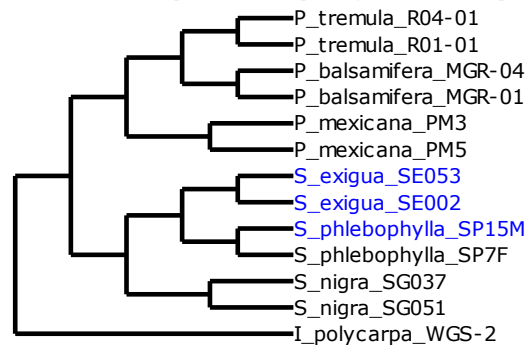534    **Supplemental Figures**



535

536    **Figure S1**. Alternate bipartitions for the three species of *Populus*, based on gene tree
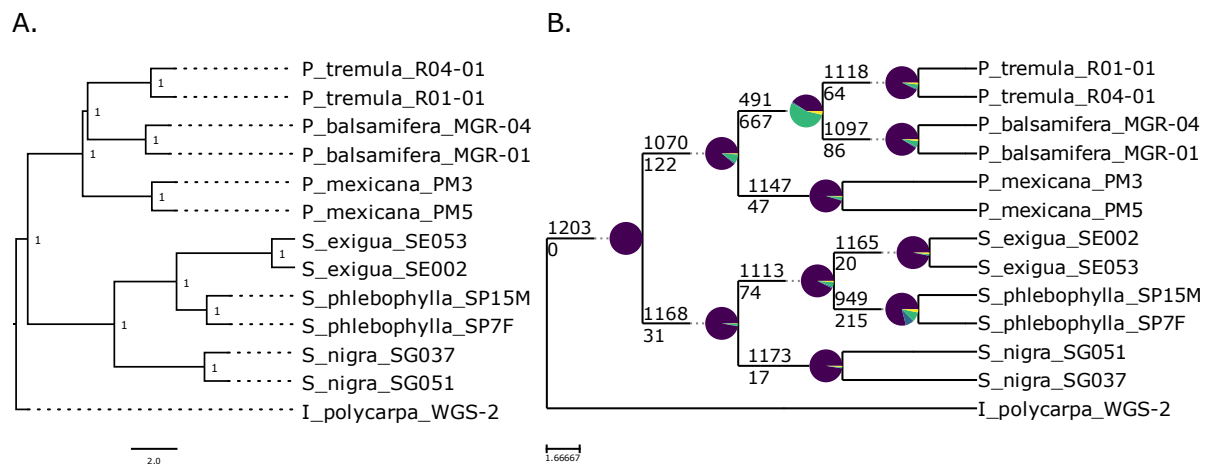
537    concordance. The cladogram in all three panels is that of the ASTRAL-III species tree (Figure

538    1), and the blue color represents the bipartition supported by the indicated number of gene trees

539    in each panel. **A)** 158 gene trees support the displayed ASTRAL-III species tree topology. **B)**

540    114 gene trees support a bipartition that places *P. tremula* and *P. mexicana* together. **C)** 94 gene

541    trees support a bipartition that places *P. balsamifera* and *P. mexicana* together.

**Figure S2**. Alternate bipartitions for the three species of *Salix*, based on gene tree concordance.
The cladogram in all three panels is that of the ASTRAL-III species tree (Figure 1), and the blue
color represents the bipartition supported by the indicated number of gene trees in each panel. **A)**
327 gene trees support the displayed ASTRAL-III species tree topology. **B)** 44 gene trees and **C)**
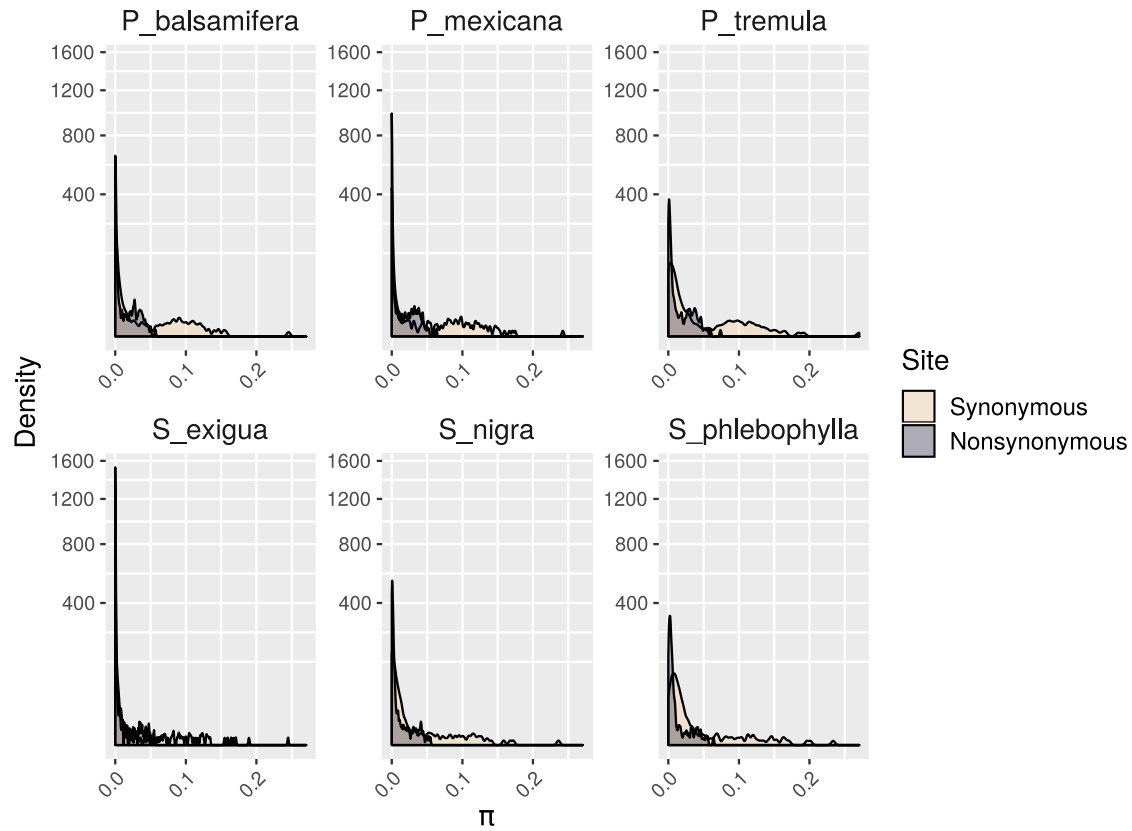39 gene trees place one of the *S. phlebophylla* individuals within *S. exigua*.

**Figure S3**. Species trees estimated for all genes and known paralogs. **A)** Species tree generated by ASTRAL-III for the gene trees. Node values represent bootstrap support from 100 multilocus bootstrap replicates in ASTRAL-III. Branch lengths represent coalescent units. **B)** Cladogram showing the congruence of gene trees for all nodes in the ASTRAL-III species tree. The numbers above each node represent the number of gene trees that support the displayed bipartition, and numbers below the node represent the number of gene trees that support all alternate bipartitions. Purple wedges represent the proportion of gene trees that support the displayed bipartition. Blue wedges represent the proportion of gene trees that support a single alternative bipartition. Green wedges represent the proportion of gene trees that have multiple conflicting bipartitions. Yellow wedges represent the proportion of gene trees that have no supported bipartition. Plotting code and its interpretation were provided by Matt Johnson (for more detail, see: https://github.com/mossmatters/MJPythonNotebooks/blob/master/PhyParts_PieCharts.ipynb)

561

**Figure S4**. Distributions of values of nucleotide diversity (Nei's $\pi$) within each species at synonymous (yellow) and nonsynonymous (purple) sites.

562

563

564