1    Recent Y chromosome divergence despite ancient origin of dioecy in poplars (*Populus*)

2

3    Armando Geraldes[a#], Charles A. Hefer[a#], Arnaud Capron[a], Natalia Kolosova[a], Felix

4    Martinez-Nuñez[a], Raju Y. Soolanayakanahally[b], Brian Stanton[c], Robert D. Guy[d], Shawn

5    D. Mansfield[e], Carl J. Douglas[a] and Quentin C. B. Cronk[a]

6    [#]Equal contribution

7    Departments of [a]Botany, [d]Forest and Conservation Sciences and [e]Wood Science,

8    University of British Columbia, Vancouver, BC V6T 1Z4, Canada

9    [b]Agroforestry Development Centre, Agriculture and Agri-Food Canada, Indian Head, SK

10   S0G 2K0, Canada

11   [c]Greenwood Resources, Portland, OR 97201, USA

12

13   Corresponding author:  Armando Geraldes

14   6270 University Boulevard – Botany Department, UBC

15   Vancouver, BC V6T 1Z4 Canada

16   email: geraldes_at_mail_dot_ubc_dot_ca

17

18   Keywords: GWAS, SNP, XY-sex-determining-system, gender, methylation

19

20   Supplementary information: All supplementary files will be made available upon request.

21   Data Accessibility: All sequence data will be deposited on Genbank and SRA and all

22   other data files on Dryad.

23

**Abstract**

All species of the genus *Populus* (poplar, aspen) are dioecious, suggesting an ancient origin of this trait. Theory suggests that non-recombining sex-linked regions should quickly spread, eventually becoming heteromorphic chromosomes. In contrast, we show using whole genome scans that the sex-associated region in *P. trichocarpa* is small and much younger than the age of the genus. This indicates that sex-determination is highly labile in poplar, consistent with recent evidence of "turnover" of sex determination regions in animals. We performed whole genome resequencing of 52 *Populus trichocarpa* (black cottonwood) and 34 *P. balsamifera* (balsam poplar) individuals of known sex. Genome-wide association studies (GWAS) in these unstructured populations identified 650 SNPs significantly associated with sex. We estimate the size of the sex-linked region to be ~100 Kbp. All significant SNPs were in strong linkage disequilibrium despite the fact that they were mapped to six different chromosomes (plus 3 unmapped scaffolds) in version 2.2 of the reference genome. We show that this is likely due to genome misassembly. The segregation pattern of sex associated SNPs revealed this to be an XY sex determining system. Estimated divergence times of X and Y haplotype sequences (6-7 MYA) are much more recent than the divergence of *P. trichocarpa* (poplar) and *P. tremuloides* (aspen). Consistent with this, in *P. tremuloides* we found no XY haplotype divergence within the *P. trichocarpa* sex-determining region. These two species therefore have a different genomic architecture of sex, suggestive of at least one turnover event in the recent past.

46      **Introduction**

47          The separation of male and female sexual function into different individuals

48      (dioecy) is an efficient way to ensure that sexual reproduction results in the

49      recombination of genetic information from different individuals and is common in

50      eukaryotes, occurring in 94% of animals [1] but only in about 6% of flowering plant

51      species [1, 2]. Dioecy usually evolves from a cosexual ancestral state and involves at

52      least two mutations. In one model, the pathway to XY systems involves one recessive

53      mutation that suppresses male function ($M^F$ -> $M^s$) and a dominant mutation that

54      suppresses female function ($F^f$ -> $F^S$) [3], where the Y chromosome harbors the alleles

55      $M^F$ and $F^S$ and the X chromosome the alleles $M^s$ and $F^f$. Recombination suppression

56      between these loci on the Y chromosome likely evolves under the action of natural

57      selection because recombination generates unfit sterile individuals [4]. With time,

58      recombination suppression may extend to the rest of the chromosome via the

59      accumulation of sexually antagonistic mutations on the Y [5], leading to the degeneration

60      of the heterogametic sex chromosome (the Y or the W) via Muller's ratchet, background

61      selection and hitchhiking [6]. Under this view, old sex chromosomes are structurally and

62      genetically divergent. The mammalian Y chromosome, having evolved ~170 MYA [7],

63      is one such case of a degenerate Y chromosome that retains only a small fraction of the

64      genes thought to be present in the autosomal pair from which the Y arose [8].

65          Studying old and degenerate Y chromosomes allows only for retrospective

66      insights into their evolutionary origins. In some groups, sex chromosomes may be young

67      and therefore provide windows into the initial stages of their evolution (e.g., [9, 10]). In

68      plants, dioecy evolved independently in several clades allowing for a comparative

3

69    approach that may reveal commonalities and peculiarities among independent origins of

70    sex chromosomes [11]. Despite recent progress in the use of genomic resources to

71    unravel the genetic basis of dioecy in plants such as papaya and white campion, the

72    nature of sex-determining regions and sex-determining genes in plants remains elusive

73    [12].

74         *Populus* species (poplars, cottonwoods and aspens) present an excellent

75    opportunity to study the evolution of sex chromosomes. *Populus* and *Salix*, sister genera

76    in the Salicaceae, are composed exclusively of dioecious species (with reports of rare

77    cosexual genotypes, e.g. [13]), consistent with a single ancient origin of dioecy in this

78    group around 65 MYA [14]. The cytological evidence (reviewed in [15]) for the

79    existence of heteromorphic sex chromosomes is mixed, but in general there is no strong

80    evidence for their existence (or for different chromosome counts in males and females),

81    and the nature of the sex-determining region in *Populus* has remained elusive. Previous

82    genetic mapping studies have mapped the sex-determining region to the proximal

83    telomeric end of chromosome 19 in poplars and cottonwoods (*Populus* sections

84    Tacamahaca and Aigeiros, [16, 17]) or to a pericentromeric region in aspens (*Populus*

85    section Populus, [18-20]). Some studies have proposed that females are the

86    heterogametic sex (ZW system, [16, 19]) while other evidence suggests that males are

87    (XY system, [17, 18, 20-22]). Recently, markers associated with sex were described for

88    aspens, corresponding to the presence of the gene *TOZ19* on the Y chromosome of *P.*

89    *tremula* and *P. tremuloides* and its absence from the X chromosome [22]. Here we use a

90    genome-wide association approach (GWAS) to determine the genomic architecture of sex

91    in two species of poplar.

4

92

## Results

Genome-wide association analysis (GWAS)

We performed a simple case control GWAS between allele frequency at

3,656,736 loci with MAF>0.1 (minor allele frequency) and GR>0.9 (genotyping rate,

Table S1) and sex (male vs. female) of 34 female and 18 male *P. trichocarpa* individuals

(hereafter T52 association population, SI). After Bonferroni correction we recovered 623

single nucleotide polymorphisms (SNPs) significantly associated with sex ($\alpha<0.05$; Fig. 1

and Table S2). Across all significant SNPs and accessions, females were homozygous at

99.9% of the genotypes and males were heterozygous at 94.0% of the genotypes, a

pattern consistent with an XY sex determining system (Table 1).

A similar analysis for 1,140,437 SNPs (Table S1) and sex (18 female and 16 male

individuals, SI) in *P. balsamifera* (hereafter B34) recovered no SNPs statistically

associated with sex (Fig. 1 and Table S2). Inspection of the results of the two analyses

revealed that for 72.6% (452/623) of the significantly associated SNPs in *P. trichocarpa*,

no data was available in *P. balsamifera* (i.e. SNPs had GR<0.9 and/or MAF<0.1). For the

remaining SNPs, the vast majority (157/171) showed a similar pattern to that of SNPs

significantly associated with sex in *P. trichocarpa*, i.e. females were homozygous and

males heterozygous (with less than 10% of accessions deviating from this pattern) and

the observed uncorrected p-values range was $2.07 \times 10^{-4} - 1.22 \times 10^{-6}$ (Fig. 1 and Table S2).

Finally, we created a third association population consisting of 36 females and 32

males where, in each sex, equal numbers of accessions were *P. trichocarpa* and *P.*

*balsamifera* (hereafter BT68, SI). In this population there were 1,782,995 SNPs with

5

115    MAF>0.1 and GR>0.9 (Table S1) and 303 SNPs were significantly associated with sex

116    (α<0.05; Fig. 1 and Table S2), of which only 27 were not significant in the analysis with

117    *P. trichocarpa* alone (T52, Table S2). Across all significant SNPs and accessions in

118    BT68, females were homozygous at 99.6% of the genotypes and males were

119    heterozygous at 94.1% of the genotypes, a pattern again consistent with an XY sex

120    determining system (Table 1).

121          In all three cases, Q-Q plots (Fig. S1) did not reveal an inflation of observed p-

122    values with regards to the expected distribution of p-values, except for the extreme

123    observed p-values. This is as expected given that T52 and B34 are unstructured

124    populations and the population structure observed in BT68 did not co-vary with the

125    phenotype (SI).

126

127    <u>Genomic distribution of sex-associated SNPs</u>

128          Surprisingly, SNPs significantly associated with sex were located in 10 different

129    regions of v2.2 of the *P. trichocarpa* reference genome assembly. The majority of SNPs

130    associated with sex in T52 were located in the proximal end of chromosome 19 (hereafter

131    Chr19P, 62.12%, 387/623) and the distal end of the same chromosome (hereafter

132    Chr19D, 14.60%, 91/623). Remaining SNPs were located on chromosomes 1, 4, 5, 8, 9

133    and scaffolds 261, 1817 and 2325 (Fig. 1 and Table 2). Despite being distributed across

134    different genomic regions, pairwise estimates of linkage disequilibrium between SNPs

135    associated with sex were very high (all significant SNPs, average $r^2$=0.93, range 0.46-1;

136    significant SNPs in putatively different genomic regions, average $r^2$=0.90, range 0.46-1;

137    significant SNPs in the same genomic region, average $r^2$=0.97, range 0.51-1; Fig. S2).

6

138

139    Evidence for a single sex-specific locus in *P. trichocarpa*

140         Previous QTL experiments in *Populus* mapped sex to a single location suggesting

141    inheritance as a single genetic locus [15]. Furthermore, sex always mapped to

142    chromosome 19, albeit to different positions on the chromosome in different

143    crosses/species [15]. Thus, the presence of SNPs associated with sex in different

144    genomic regions in our GWAS might be due to problems with the assembly of the

145    reference genome v2.2. To address this we repeated the GWAS for the T52 population

146    after read mapping and SNP calling to genome assemblies v1.0 and v3.0. In both cases,

147    sex again mapped to multiple regions although the details of the locations differed among

148    assemblies (Fig. S3 and Tables S3-S4). The sex-determining region was therefore highly

149    unstable with respect to assembly version.

150         We also performed cross mapping of sex-linked regions among assemblies with a

151    BlastN (E-value cutoff $10^{-10}$, best 10 hits kept) search of the regions containing

152    significant SNPs in v2.2 against assemblies v1.0 and v3.0 (Table S5). The resulting

153    alignments indicated that the Chr09 and Chr19D sex-linked regions have similar

154    locations in all three assemblies. All other sex-linked regions mapped to a different

155    location in at least one of the three assemblies, e.g., the sex-linked region in Chr19P (v2.2

156    and v1.0) has moved to the distal end of Chr18 in v3.0.

157         Finally, we queried (BlastN, E-value cut-off $10^{-6}$, best hit kept) the BAC end

158    sequences of a male *P. trichocarpa* library [23] to look for BAC clones nearby our

159    sex-linked regions in which the two ends mapped to different locations in v2.2 and

160    we identified 17 such BACs (SI). For three of these BACs, both ends were sex-linked.

161　One BAC-end sequence from clone POR18-C06 maps to Chr19D and the other end

162　maps to Chr01. One BAC-end sequence from clone POR02-A02 maps to Chr19P and

163　the other end maps to scaffold 2325. Finally, one BAC-end sequence from clone

164　POR07-E07 maps to Chr19P and the other end maps to Chr08. These results suggest

165　that in this male, the sex-linked regions in Chr19D and Chr01, as well as in Chr08,

166　Chr19P and scaffold 2325 are physically linked.

167　　　The above evidence, taken together, strongly suggests that assembly problems are

168　sufficient to explain the genomic distribution of the sex-associated markers.

169

170　Sex-linked regions in other accessions and species

171　　　We developed two PCR-RFLP assays for rapid genotyping of accessions in two

172　of the regions with SNPs significantly associated with sex (Chr09 and Chr19P).

173　Application of these assays to 8 samples of each sex and species used in the GWAS

174　revealed full agreement between WGS and PCR-RFLP inferred genotypes (Fig. S4) and

175　confirmed that one male of each species, BELA18-5 and AP2446, appears to be

176　recombinant; i.e. both of these males are homozygous for the majority of SNPs in

177　Chr19P, but are heterozygous for significant SNPs in the other sex-linked regions (Fig.

178　S2). Application of these assays to *P. trichocarpa* and *P. balsamifera* accessions of each

179　sex that were not used in the GWAS showed that these SNPs are linked with sex in

180　independent accessions (Fig. S4). Finally, we used these assays to determine whether

181　these SNPs are also linked to sex in other species. All 16 *P. deltoides* and 16 *P. nigra*

182　accessions of known sex assayed were homozygous (XX) in females and heterozygous

183　(XY) in males (Fig. S4). This indicates that the *P. trichocarpa/P. balsamifera* sex-linked

8

184  markers are conserved in these species. However, for one female and three male *P.*

185  *tremuloides* accessions no differences between sexes were observed (Fig. S4), suggesting

186  that in aspens these regions are not sex-linked.

187

188  Phylogeny of X and Y alleles

189       We performed allele-specific amplification and sequencing of X and Y alleles in

190  two regions associated with sex (gene POPTR_0019s00240 on Chr19P and gene

191  POPTR_0009s08410 on Chr09) using several males from each of four species: *P.*

192  *trichocarpa*, *P. balsamifera*, *P. deltoides* and *P. nigra* (hereafter referred to as

193  "cottonwoods"). We also included sequences cloned from *P. tremuloides* (hereafter

194  referred to as "aspen"), the reference genome sequence, as well as the genome sequence

195  of the paralog of each gene that resulted from the Salicoid whole genome duplication

196  (WGD) event [24]. Maximum likelihood phylogenies of each region (Fig. 2 and Fig. S5)

197  show that both X and Y chromosome alleles from all four cottonwood species group by

198  gametolog (i.e., X or Y) and not by species, indicating that X and Y chromosome alleles

199  began to diverge before species did. Note that because for one of the amplicons in

200  Chr19P we failed to amplify the X gametolog of *P. nigra, P. nigra* alleles are not shown

201  in the concatenated phylogeny of Chr19P (Fig. 2); nevertheless phylogenies of the other

202  two amplicons in Chr19P show unequivocally that *P. nigra* alleles cluster by gametolog

203  (Fig. S5). The placement of aspen alleles with respect to X and Y alleles from

204  cottonwoods is uncertain. For the region in Chr09, they cluster with cottonwood

205  sequences from the X gametolog, but with low bootstrap support, while for the region in

206  Chr19P they appear basal to the X and Y clades (Fig. 2).

207

208    <u>Divergence at X and Y regions</u>

209    The phylogenies in Fig. 2 clearly suggest that recombination between the X and Y

210    regions identified here ceased, and their divergence in cottonwoods started, after the split

211    between cottonwoods and aspens. The amount of divergence at silent sites ($K_s$), between

212    the X and Y clade (Chr09 $K_s$=0.0224 and Chr19P $K_s$=0.0163) was only slightly lower

213    than $K_S$ between all XY cottonwood alleles and aspen (Chr09 $K_s$=0.0638 and Chr19P

214    $K_s$=0.0186), and both were roughly one tenth the $K_s$ between the XY clade and the

215    paralog from the Salicoid WGD (Chr09 $K_s$=0.2027 and Chr19P $K_s$=0.1774; Table 3).

216    Assuming the timing of the WGD to be 65 MYA [24], then XY divergence for Chr9

217    would be approx. 7.2 MYA and Chr19 divergence approx. 6.0 MYA.

218    For both regions, the ratio of non-synonymous substitutions per non-synonymous

219    site to synonymous substitutions per synonymous site ($K_a/K_s$) is higher for the Y lineage

220    than for the X lineage (Table 3). This pattern is consistent with an accumulation of

221    deleterious mutations following recombination suppression. The fact that this difference

222    is larger when divergence is measured to aspen (Chr09 X $K_a/K_s$=0.560, Y $K_a/K_s$=0.870

223    and Chr19P X $K_a/K_s$=0.247 and Y $K_a/K_s$=0.463), than when divergence is measured to

224    the Salicoid paralog (Chr09 X $K_a/K_s$=0.582, Y $K_a/K_s$=0.737 and Chr19P X $K_a/K_s$=0.188

225    and Y $K_a/K_s$=0.208), suggests that the increase in $K_a/K_s$ in the Y lineage is recent.

226    Furthermore, despite its recent origin our data suggest that the Y-haplotype is already

227    becoming non-functional as we observe frame-shift insertions/deletions in Y sequences

228    of POPTR_0009s08410.

229

230 Size and composition of the sex-linked region.

231   The 650 sex-associated SNPs, if concatenated, cover a total genomic region of

232 ~100 Kbp. Thus given the evidence above that a single region is involved, that region is

233 extremely small. We considered if there might be large missing tracts of Y sequence that

234 were not detected by our read-mapping protocol. *De novo* assembly of unmapped reads

235 from male accessions revealed four male-specific contigs that are candidates for such Y

236 sequences (SI). However these are short (longest contig is 2514 bp) and BlastN searches

237 reveal that they either are repetitive in nature or have significant similarity to the sex-

238 linked regions identified by GWAS. There is thus no present evidence that the sex-locus

239 in *P. trichocarpa* is significantly larger than reported here. The 13 genes in the sex-linked

240 region (Table 2) cover a range of functional classes, including DNA methylation,

241 hormone regulation, ion transport and plant defense.

242

243 **Discussion**

244 XY sex-determining system

245   The identification of 650 sex-specific SNPs heterozygous in males and

246 homozygous in females by GWAS unequivocally shows that an XY system is involved in

247 sex-determination in *P. trichocarpa/P. balsamifera*. The findings were fully and

248 independently supported by PCR/RFLP-assays for two representative SNPs that

249 distinguish X and Y alleles carried out on *P. trichocarpa*, *P. balsamifera*, *P. deltoides*,

250 and *P. nigra* individuals of known sex not included in the GWAS. This finding of an XY

251 system in cottonwoods (*Populus* sections Tacamahaca and Aigeiros) is further supported

252 by previous reports of an XY system in *P. nigra* of section Aigeiros [17] and aspens of

11

253     section Populus [20-22] but is at odds with previous suggestions that a ZW (female

254     heterogamy) system of sex determination may function in *P. trichocarpa* [16].

255             The previous suggestion that *P. trichocarpa* has a ZW system was based on

256     inferences from a cross of *P. deltoides* x (*P. nigra* x *P. deltoides*) and was not supported

257     by sex-specific markers [16]. Our results run counter to those inferences, but it is

258     conceivable that a ZW system, with a highly divergent W chromosome that is not

259     represented in the *P. trichocarpa* reference sequence [24], could produce the observed

260     pattern of homozygosity in females and heterozygosity in males at SNPs significantly

261     associated with sex, as the W sequence would be absent in males and divergent enough

262     from the Z that reads from the W chromosome do not map to the reference sequence.

263     Thus, in females, apparent homozygosity would in fact be due to hemizygosity. Several

264     observations contradict this hypothesis: a) we observed heterozygous positions in females

265     at sex-linked regions intermingled with SNPs significantly associated with sex (SI), b)

266     Sanger sequencing of females for sex-linked regions revealed heterozygous positions

267     (SI), c) qPCR of two sex-linked regions (Chr19P and Chr19D) revealed a 1:1 ratio of

268     amplification of autosomal to sex-linked regions in both sexes (SI), d) WGS coverage is

269     approximately similar in males and females at sex-linked regions (SI) and e) *de novo*

270     assembly of female specific regions did not reveal unassembled regions unique to

271     females (SI). Given strong direct evidence for an XY system from sex-linked markers,

272     and absence of evidence for hemizygosity in females, we now argue that the ZW

273     hypothesis can be discounted.

274

275     Genomic architecture of the sex locus

12

276    The 623 sex-specific SNP markers identified by GWAS in T52 are in nearly

277    complete genetic linkage (Fig. S2). The majority of these markers map to Chr19P

278    confirming previous studies that implicate this region as the location of the sex locus [16,

279    17]. However, remarkably, we found that sex-linked markers in apparent genetic linkage

280    map to multiple physical locations in the three *P. trichocarpa* genome assemblies (Fig. 1

281    and Tables S2-S4). Our data do not support the existence of a multi-locus system of sex

282    determination in *P. trichocarpa*, but instead suggest that a single genetic region controls

283    dioecy and that the genome assembly is a work in progress with some contigs from

284    Chr19P having been misassembled into other genomic regions. Sex determining regions

285    and sex chromosomes are notoriously difficult to assemble [25]. Further refinement of

286    the assembly regarding the sex locus may require complementary methods.

287

288    <u>The age of the cottonwood sex locus and evolution of dioecy in *Populus*</u>

289    We find the same sex-linked markers in *P. trichocarpa* and *P. balsamifera*

290    (*Populus* section Tacamahaca) as well as in *P. nigra* and *P. deltoides* (*Populus* section

291    Aigeiros). The sex locus therefore predates the divergence of these species. Sequence

292    analysis of sex regions of these species suggests an approximate date of 6-7 MYA (late

293    Tertiary) for the divergence of X and Y. The fossil record indicates that aspens (*Populus*

294    sect. Populus) likely diverged from cottonwoods long before this, given that middle to

295    late Oligocene (~25 MYA) fossils of section Populus from Alaska have been reported

296    [26]. Consistent with this is the fact that the polymorphic loci we identified do not

297    provide sex-specific markers in aspens. Furthermore, sex linked markers have recently

298    been identified in the pericentromeric region of Chr19 in aspens [22]. Genes in this

13

299 region are not homologous to sex-linked genes identified in our study, and SNPs in this

300 region do not segregate with sex in our mapping populations (SI); hence aspens and

301 cottonwoods likely have independent sex determining mechanisms.

302        If there were a single origin of dioecy in this group, it is problematic that there are

303 apparently distinct sex-determining loci in *Populus*. One plausible explanation is that

304 there has been at least one sex-determination mechanism "turnover" since the divergence

305 of poplars and aspens. The labile nature of sex determining regions is well known, with

306 many examples of "turnover" of sex determining regions from diverse groups [27].

307 Mapping of sex-linked regions in other *Populus* species as well as in the sister genus

308 *Salix* (willows) would provide further insight into the dynamics of sex-linked region

309 turnover in the Salicaceae.

310

311 <u>The size of the sex locus</u>

312        One remarkable feature of the sex locus described here is its compactness.

313 Concatenating all the regions with sex specific markers leads to a total estimated size for

314 the sex-determining region of ~100 Kbp. This small size is consistent with the difficulties

315 encountered in finding sex-specific markers in the Salicaceae (reviewed in [15]).

316 However, there are good reasons for supposing that a non-recombining region at a sex

317 locus will rapidly expand, eventually to encompass an entire chromosome [28]. Such

318 expansion is empirically well documented in other plant systems [29] and is driven by

319 sexual conflict making it advantageous for more and more genes to be captured by the

320 non-recombining regions. Even the 6-7 MYA date we estimate for the divergence of X

321 and Y alleles would likely be sufficient for expansion to encompass a considerable

14

322    portion of a chromosome. Therefore the apparent remarkably small size of the *P.*

323    *trichocarpa* sex locus requires explanation.

324        One possibility is that the actual size of the cottonwood sex-determining locus is

325    larger than it appears due to large tandem duplications and transposable element

326    insertions in the Y. Yet, our *de novo* assembly of male-specific unmapped reads revealed

327    only four small male-specific contigs (average length 1877 bp, SI) and these have either

328    Blast hits to the sex-linked regions identified with GWAS (SI) or consist mostly of low

329    complexity repetitive sequence. We were unable to retrieve further male-specific contigs,

330    specifically, regions of higher divergence to the female reference sequence that may be

331    indicative of older divergence strata as observed in other animal [8] and plant species

332    [29]. Future investigations might reveal larger Y-specific regions. Alternatively, it is

333    possible that features unique to trees dampen the expansion of sex determining regions.

334    For instance, sexual conflict may be minimal in trees as carbon investment in

335    reproduction is a relatively small annual cost compared to the massive storage of carbon

336    in wood, a tissue with no obvious secondary sexual characteristics.

337

338    Functional insights into sex-determination in cottonwoods

339        The sex-linked specific region in *P. trichocarpa* contains 13 genes (Table 2).

340    However it is too early to say which, if any, of these genes are the master-regulators of

341    sex. The reference genome is from a female (XX) individual [24] and, as suggested

342    above, further work is required to fully characterize the Y chromosome. Furthermore,

343    many of the genes in this region have poorly defined functions. Nevertheless, there are at

344    least two plausible candidate genes. One, a poplar ortholog of the *Arabidopsis thaliana*

15

345     [Arabidopsis] cytokinin pathway-associated *ARABIDOPSIS RESPONSE REGULATOR*

346     *17* (*ARR17*), is implicated in phytohormone signaling and the other, the poplar ortholog

347     of Arabidopsis *METHYLTRANSFERASE 1*, (*MET1*), is involved in DNA methylation.

348          Phytohormone signaling is involved in other plant sex determination systems,

349     such as the ethylene pathway in cucumber [30], and it is possible that cytokinin

350     signaling, mediated by *ARR17* is used in poplar. DNA methylation has been implicated in

351     sex determination in other plant systems, e.g. *Silene latifolia* [31]. In an

352     andromonoecious clone of *P. tomentosa* expression of the poplar orthologue of *MET1*

353     was significantly higher in all stages of female flower development [32]. In Arabidopsis,

354     *MET1* is required for maintenance of epigenetic memory [33] and is involved in

355     reproductive development including the control of floral homeotic genes such as

356     *AGAMOUS*, *APETALA3* and *SUPERMAN* [34, 35].

357          Due to the Salicoid WGD [24] there are two paralogs of genes in the sex locus

358     region such as *ARR17* and *MET1*, relative to Arabidopsis. Neofunctionalization, in which

359     one copy has evolved a specific sex-determining function while the other copy retains the

360     ancestral function, is therefore possible. The WGD may thus be important in the

361     evolution of dioecy in this group. Functional differences between the paralog in the sex-

362     specific region and an autosomal sister paralog could reveal pathways involved in sex

363     determination.

364

365     **Materials and Methods**

366          Tree sex was determined by visual inspection of flowers. DNA from *Populus*

367     *trichocarpa* and *P. balsamifera* association populations was extracted from leaves and

16

368    sequenced (100bp paired-end reads) on an Illumina HiSeq at the Genome Sciences

369    Centre, Vancouver, BC to either 15x or 30x coverage (SI). Sequence data generated

370    ranged from 31-241 million reads. All sequences are deposited at the NCBI short read

371    archive under SRA XXX. Illumina reads were aligned to reference *P. trichocarpa*

372    genome assemblies v1.0, v2.2 and v3.0 (http://www.phytozome.net) using BWA version

373    0.6.1 [36] with a 4 bp misalignment threshold, disallowing insertions or deletions within

374    5bp of the end of the sequence (aln –n 0.04 -i 5), maximum insert size of 500 bp (sampe -

375    a 500), and default values for the remaining parameters. Paired-end mate information was

376    synced using Picard-tools FixMateInformation (http://picard.sourceforge.net/). Local re-

377    alignment was performed on identified regions with high SNP entropy, using a window-

378    size of 10 bp, and a mismatch fraction of 0.15 for base qualities to identify mismatched

379    regions using GATK version 1.5 [36].  Indel re-alignments were restricted to regions

380    with a maximum insert size of 3 Kbp, and the maximum positional change of an indel set

381    to 200 bp. Variant calls were made using the duplicate-marked alignment files and the

382    UnifiedGenotyper from GATK emitting variant with a minimum phred-scaled confidence

383    threshold of 30. We used vcftools [37] to filter out any variants where coverage was <5X

384    and where more than two bases were segregating.

385        We performed a standard case/control GWAS between allele frequencies and sex

386    phenotype using Plink v1.07 [38]. We report associations at $\alpha<0.05$ after Bonferroni

387    correction for multiple testing. Analysis of population structure in the three association

388    populations is given in SI.

389        PCR-RFLP (polymerase chain reaction followed by restriction fragment length

390    polymorphism) genotyping assays in two regions associated with sex were developed as

391    follows: mpileup files were converted into fasta files by generating calls at each base of

392    the reference whenever coverage at the position in each individual was higher than six

393    and whenever heterozygote genotypes were present by requiring that each allele had

394    coverage of at least three. All other positions were considered missing data. The fasta

395    sequences were used to design PCR primers in regions conserved across all accessions to

396    amplify two short fragments on the sex-linked regions that mapped to Chr09 and Chr19P.

397    PCR primers, amplicons and protocol details are in SI. The Chr09 amplicon was digested

398    with *Bsl*I (New England Biolabs, Ipswich, MA) and *Cla*I (New England Biolabs,

399    Ipswich, MA); the Chr19 amplicon was digested with *TspR*I (New England Biolabs,

400    Ipswich, MA); see SI for details. The same assays were used in *P. deltoides, P. nigra* and

401    *P. tremuloides* accessions (SI).

402        To generate haplotypic Sanger sequences from selected male accessions (SI),

403    allele-specific primers [39] were designed for three regions of the gene

404    POPTR_0019s00240 on Chr19P and for one region of the gene POPTR_0009s08410 on

405    Chr09 (SI). Each allele-specific primer was used with the common primer to generate an

406    allele-specific PCR fragment that was subsequently cloned and Sanger sequenced. PCR

407    protocol, amplicon and cloning details are in SI. Chromatograms from Sanger sequencing

408    were visually inspected, trimmed, and aligned with BioEdit [40]. Sequences were aligned

409    to the closest *P. trichocarpa* paralog (resulting from the Salicoid WGD); for Chr09 the

410    paralog is POPTR_0001s29310, and for Chr19 the paralog is POPTR_0004s14140) and

411    neighbor joining maximum likelihood trees for each amplicon were estimated in MEGA

412    v5.03 [41] using the Tamura-Nei model and complete deletion of all sites with missing

413    data and gaps. Levels of divergence were calculated for synonymous sites ($K_s$) only and

414    for replacement sites only ($K_a$) in DNAsp v5 [42]. Sequence data is deposited in NCBI

415    under accession numbers XXX.

416

425

426    **References**

427    1. Jarne P, Auld JR (2006) Animals mix it up too: the distribution of self-fertilization
428        among hermaphroditic animals. Evolution 60:1816–1824.

429    2. Renner SS, Ricklefs RE (1995) Dioecy and its correlates in the flowering plants.
430        American Journal of Botany 596–606.

431    3. Charlesworth B (1991) The evolution of sex chromosomes. Science 251:1030–1033.

432    4. Bull JJ (1983) Evolution of sex determining mechanisms. Benjamin-Cummings
433        Publishing Company.

434    5. Rice WR (1984) Sex chromosomes and the evolution of sexual dimorphism. Evolution
435        38:735–742.

436    6. Bachtrog D (2013) Y-chromosome evolution: emerging insights into processes of Y-
437        chromosome degeneration. Nat Rev Genet 14:113–124.

438    7. Graves JAM (2006) Sex chromosome specialization and degeneration in mammals.
439        Cell 124:901–914.

440    8. Skaletsky H, Kuroda-Kawaguchi T, Minx PJ, et al. (2003) The male-specific region of

441     the human Y chromosome is a mosaic of discrete sequence classes. Nature 423:825–
442     837.

443  9. Vicoso B, Emerson JJ, Zektser Y, et al. (2013) Comparative sex chromosome
444     genomics in snakes: differentiation, evolutionary strata, and lack of global dosage
445     compensation. Plos Biol 11:e1001643.

446  10. Yoshida K, Makino T, Yamaguchi K, et al. (2014) Sex chromosome turnover
447     contributes to genomic divergence between incipient stickleback species. PLoS
448     Genet 10:e1004223.

449  11. Ming R, Bendahmane A, Renner SS (2011) Sex chromosomes in land plants. Annu
450     Rev Plant Biol 62:485–514

451  12. Zhang J, Boualem A, Bendahmane A, Ming R (2014) Genomics of sex determination.
452     Current Opinion in Plant Biology 18:110–116.

453  13. Stettler RF (1971) Variation in sex expression of black cottonwood and related
454     hybrids. Silvae Genet 20:42–46.

455  14. Cronk QCB (2005) Plant eco-devo: the potential of poplar as a model organism. New
456     Phytologist 166:39–48.

457  15. Tuskan GA, DiFazio SP, Faivre-Rampant P, et al. (2012) The obscure events
458     contributing to the evolution of an incipient sex chromosome in *Populus*: a
459     retrospective working hypothesis. Tree Genet Genomes 8:559–571.

460  16. Yin T, DiFazio SP, Gunter LE, et al. (2008) Genome structure and emerging evidence
461     of an incipient sex chromosome in *Populus*. Genome Research 18:422–430.

462  17. Gaudet M, Jorge V, Paolucci I, et al. (2008) Genetic linkage maps of *Populus nigra*
463     L. including AFLPs, SSRs, SNPs, and sex trait. Tree Genet Genomes 4:25–36.

464  18. Pakull B, Groppe K, Meyer M, et al. (2009) Genetic linkage mapping in aspen
465     (*Populus tremula* L. and *Populus tremuloides* Michx.). Tree Genet Genomes 5:505–
466     515.

467  19. Paolucci I, Gaudet M, Jorge V, et al. (2010) Genetic linkage maps of *Populus alba* L.
468     and comparative mapping analysis of sex determination across *Populus* species. Tree
469     Genet Genomes 6:863–875.

470  20. Pakull B, Groppe K, Mecucci F, et al. (2011) Genetic mapping of linkage group XIX
471     and identification of sex-linked SSR markers in a *Populus tremula × Populus*
472     *tremuloides* cross. Can J For Res 41:245–253.

473  21. Kersten B, Pakull B, Groppe K, et al. (2014) The sex-linked region in *Populus*
474     *tremuloides* Turesson 141 corresponds to a pericentromeric region of about two
475     million base pairs on *P. trichocarpa* chromosome 19. Plant Biol (Stuttg) 16:411–418.

476   22. Pakull B, Kersten B, Lüneburg J, Fladung M (2014) A simple PCR-based marker to
477        determine sex in aspen. Plant Biol (Stuttg). doi: 10.1111/plb.12217

478   23. Rampant PF, Lesur I, Boussardon C, et al. (2011) Analysis of BAC end sequences in
479        oak, a keystone forest tree species, providing insight into the composition of its
480        genome. BMC Genomics 12:292.

481   24. Tuskan GA, DiFazio SP, Jansson S, et al. (2006) The genome of black cottonwood,
482        *Populus trichocarpa* (Torr. & Gray). Science 313:1596–1604.

483   25. Li G, Davis BW, Raudsepp T, et al. (2013) Comparative analysis of mammalian Y
484        chromosomes illuminates ancestral structure and lineage-specific evolution. Genome
485        Research 23:1486–1495.

486   26. Wolfe JA, Tanai T (1980) The Miocene Seldovia Point Flora from the Kenai Group,
487        Alaska. United States Government Printing Office, Washington DC.

488   27. Stöck M, Horn A, Grossen C, et al. (2011) Ever-young sex chromosomes in european
489        tree frogs. Plos Biol 9:e1001062.

490   28. Charlesworth D (2013) Plant sex chromosome evolution. Journal of Experimental
491        Botany 64:405–420.

492   29. Bergero R, Forrest A, Kamau E, Charlesworth D (2007) Evolutionary strata on the X
493        chromosomes of the dioecious plant *Silene latifolia*: evidence from new sex-linked
494        genes. Genetics 175:1945–1954.

495   30. Boualem A, Fergany M, Fernandez R, et al. (2008) A conserved mutation in an
496        ethylene biosynthesis enzyme leads to andromonoecy in melons. Science 321:836–
497        838.

498   31. Janousek B, Siroký J, Vyskot B (1996) Epigenetic control of sexual phenotype in a
499        dioecious plant, *Melandrium album*. Mol Gen Genet 250:483–490.

500   32. Song Y, Ma K, Ci D, et al. (2013) Sexual dimorphic floral development in dioecious
501        plants revealed by transcriptome, phytohormone, and DNA methylation analysis in
502        *Populus tomentosa*. Plant molecular biology 83:559–576.

503   33. Blevins T, Pontvianne F, Cocklin R, et al. (2014) A two-step process for epigenetic
504        inheritance in Arabidopsis. Molecular Cell 54:30–42.

505   34. Jacobsen SE, Sakai H, Finnegan EJ, et al. (2000) Ectopic hypermethylation of flower-
506        specific genes in Arabidopsis. Current Biology 10:179–186.

507   35. Finnegan EJ, Peacock WJ, Dennis ES (1996) Reduced DNA methylation in
508        Arabidopsis thaliana results in abnormal plant development. Proc Natl Acad Sci USA
509        93:8449–8454.

21

510   36. Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler
511       transform. Bioinformatics 25:1754–1760.

512   37. DePristo MA, Banks E, Poplin R, et al. (2011) A framework for variation discovery
513       and genotyping using next-generation DNA sequencing data. Nat Genet 43:491–498.

514   38. Danecek P, Auton A, Abecasis G, et al. (2011) The variant call format and VCFtools.
515       Bioinformatics 27:2156–2158.

516   39. Purcell S, Neale B, Todd-Brown K, et al. (2007) PLINK: a tool set for whole-genome
517       association and population-based linkage analyses. The American Journal of Human
518       Genetics 81:559–575.

519   40. Ye S, Dhillon S, Ke XY, et al. (2001) An efficient procedure for genotyping single
520       nucleotide polymorphisms. Nucleic Acids Research 29:art. no.–e88.

521   41. Hall TA (1999) BioEdit: a user friendly biological sequence alignment editor and
522       analyses program for Windows 95/98/NT. Nucleic Acids Symposium Series 41:95–
523       98.

524   42. Tamura K, Peterson D, Peterson N, et al. (2011) MEGA5: molecular evolutionary
525       genetics analysis using maximum likelihood, evolutionary distance, and maximum
526       parsimony methods. Mol Biol Evol 28:2731–2739.

527   43. Librado P, Rozas J (2009) DnaSP v5: a software for comprehensive analysis of DNA
528       polymorphism data. Bioinformatics 25:1451–1452.

Table 1. Number of loci associated with sex (and percent observed genotypes) in T52 and BT68.

| | | T52 | | | BT68 | |
| --- | --- | --- | --- | --- | --- | --- |
| | Loci | Genotypes | | Loci | Genotypes | |
| | | Females | Males | | Females | Males |
| Total | 623 | 20849 | 10796 | 303 | 10439 | 9220 |
| % YY | | 0.0 | 0.9 | | 0.0 | 0.2 |
| % XY | | 0.1 | 94.0 | | 0.4 | 94.1 |
| % XX | | 99.9 | 5.1 | | 99.6 | 5.7 |

Table 2- All regions significantly associated with sex in T52 (v2.2 of the genome annotation).

| Chr[1] | Range[2] | Length (bp)[3] | T52[4] | Genes v2.2 (POPTR_)[5] | Arabidopsis ortholog[6] | Arabidopsis name or description [function] |
|---|---|---|---|---|---|---|
| 1 | 3391740..3395419 | 3680 | 12 | 0001s04290 | AT4G25650 | ACD1-LIKE [plastid function] |
| 4 | 15911439..15911672 | 234 | 2 | | | |
| 5 | 4647049..4647851 | 803 | 4 | | | |
| 8 | 5205..5205 | 1 | 1 | | | |
| 9 | 7708230..7709240 | 1011 | 36 | 0009s08410 | AT1G58290 | ATHEMA1 [chlorophyll biosynthesis] |
| 19 | 4440..67615 | 63176 | 387 | 0019s00210 | AT5G26360 | TCP-1/cpn60 chaperonin family protein [protein folding] |
| | | | | 0019s00220 | AT5G49890 | ATCLC-C; [transmembrane chloride transport] |
| | | | | 0019s00230 | AT5G49890 | ATCLC-C; [transmembrane chloride transport] |
| | | | | 0019s00240 | AT5G49160 | MET1; [cytosine methyltransferase] |
| | | | | 0019s00250 | AT1G12210 | RFL1; [defense response] |
| | | | | 0019s00260 | AT5G47260 | NB-ARC protein [defense response] |
| 19 | 15953054..15958519 | 5466 | 91 | 0019s15410 | AT3G56380 | ARR17; [cytokinin-mediated signaling pathway] |
| | | | | 0019s15415 | AT1G11300 | EGM1; [protein kinase] |
| 261 | 160..24417 | 24258 | 62 | 0261s00200 | AT5G26360 | TCP-1/cpn60 chaperonin family protein [protein folding] |
| | | | | 0261s00210 | NA | |
| | | | | 0261s00220 | AT5G49890 | ATCLC-C; [transmembrane chloride transport] |
| 1817 | 162..561 | 400 | 7 | | | |
| 2325 | 155..2156 | 2002 | 21 | | | |

[1]Chromosome/scaffold. [2]Position of the first and last significant SNPs. [3]Distance in base pairs between the first and last significant SNPs. [4]Number of significant SNPs. [5]Arabidopsis ortholog retrieved from v2.2 annotation.

Table 3- Divergence estimates at two regions associated with sex.

| | | X-Y[1] | XY[2] | | X[3] | | Y[4] | |
|---|---|---|---|---|---|---|---|---|
| | | | *P. tremuloides* | Salicoid Paralog | *P. tremuloides* | Salicoid Paralog | *P. tremuloides* | Salicoid Paralog |
| Chr09 | $K_s$[5] | 0.0224 | 0.0638 | 0.2027 | 0.0547 | 0.2174 | 0.0683 | 0.1889 |
| | $K_a/K_s$[6] | 1.606 | 0.761 | 0.667 | 0.56 | 0.582 | 0.87 | 0.737 |
| Chr19 | $K_s$[5] | 0.0163 | 0.0186 | 0.1774 | 0.0227 | 0.182 | 0.0144 | 0.1728 |
| | $K_a/K_s$[6] | 0.295 | 0.331 | 0.198 | 0.247 | 0.188 | 0.463 | 0.208 |

[1]Divergence between the X and Y clades from Fig. 2. [2]Divergence between all cottonwood sequences and *P. tremuloides*/Salicoid paralog. [3]Divergence between all sequences from the X lineage in cottonwoods and *P. tremuloides*/Salicoid paralog. [4]Divergence between all sequences from the Y lineage in cottonwoods and *P. tremuloides*/Salicoid paralog. [5]Synonymous substitutions at synonymous sites included in the estimation of K. [6]The ratio of non-synonymous substitutions per non-synonymous site to synonymous substitutions per synonymous site.
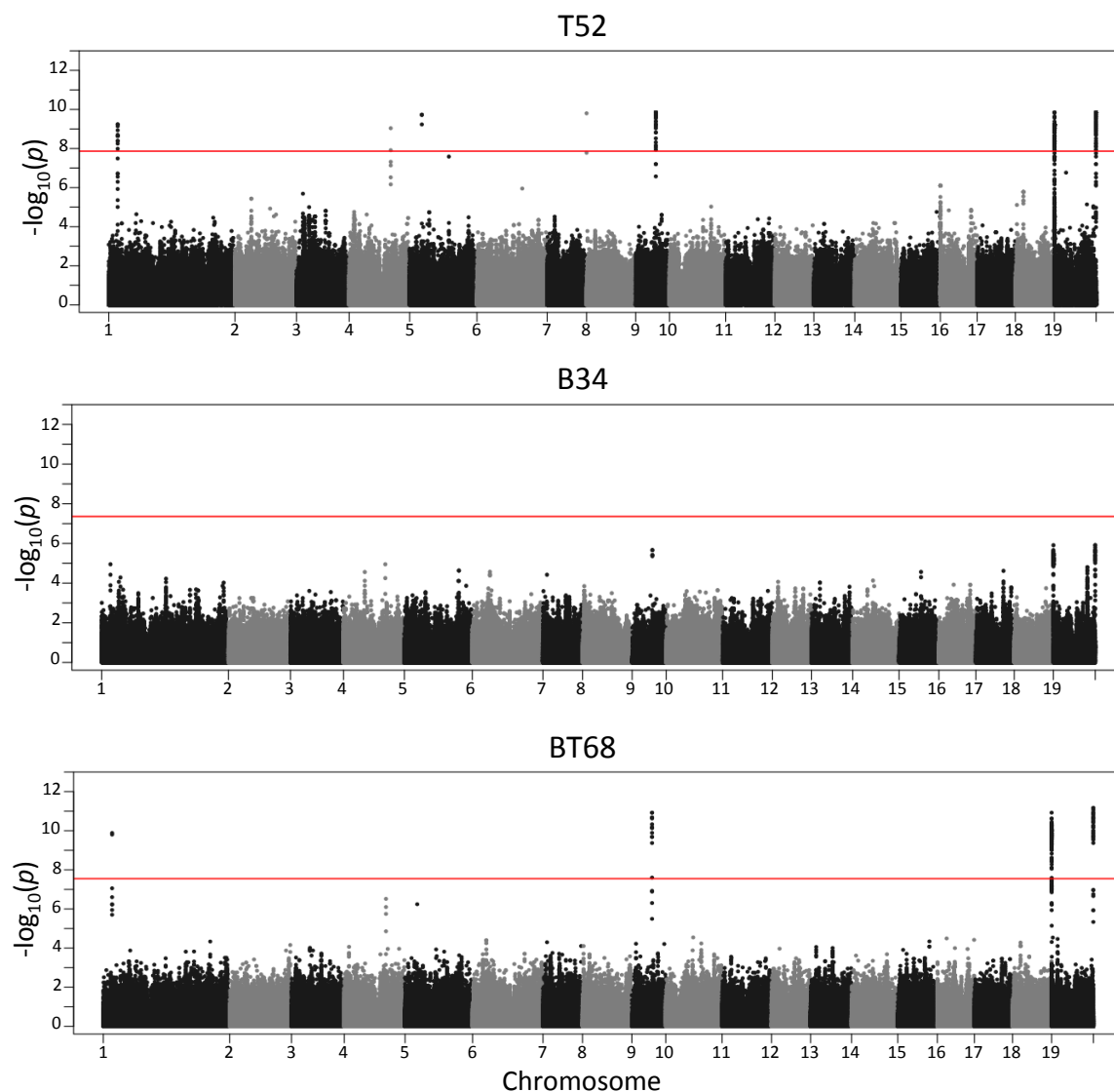
Fig. 1- Manhattan plots depicting the GWAS results for association between allele frequency (v2.2 of the reference genome) and sex in three populations: 34 female and 18 male *P. trichocarpa* accessions (T52), 18 female and 16 male *P. balsamifera* accessions (B34) and 36 female and 32 male accessions, where half the samples of each sex are *P. trichocarpa* and the other half are *P. balsamifera* (BT68). SNPs mapped to unassembled scaffolds are not represented. The horizontal line indicates the $-\log_{10}(p$ value) corresponding to $\alpha<0.05$ after Bonferroni correction for multiple testing.
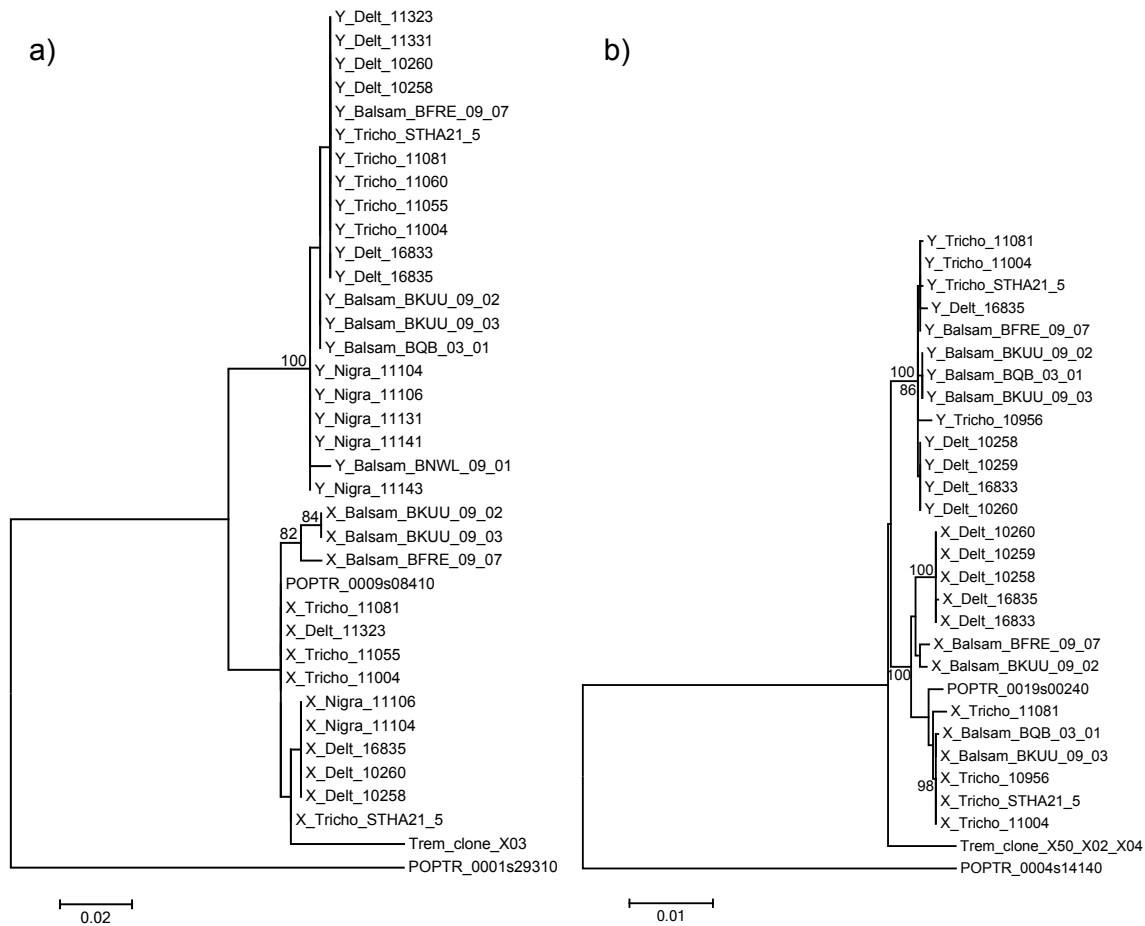
Fig. 2- Neighbor joining maximum likelihood phylogenies of regions significantly associated with sex. For each phylogeny, only male accessions were used. Y chromosome alleles are indicated with a Y and X chromosome alleles with an X at the beginning of the sequence name, followed by species and accession identifiers. Only one random *P. tremuloides* allele is depicted. Phylogenies including all *P. tremuloides* sequences are available in Fig. S5. Each phylogeny also includes the *P. trichocarpa* reference sequence from genome assembly v2.2 (POPTR_0009s08410 and POPTR_00019s00240) and the reference sequence from genome assembly v2.2 of the paralog from the Salicoid WGD (POPTR_0001s29310 and POPTR_0004s14140). Only bootstrap values higher than 80 are shown. a) Phylogeny of Chr09 region (Amplicon1:Chr09:7690067, SI) and b) Concatenated phylogeny of Chr19 (Amplicon1:Chr19:40024, Amplicon2:Chr19:41515 and Amplicon3:Chr19:44107, SI).

27