# 1  Taming the massive genome of Scots pine with PiSy50k,
# 2  a new genotyping array for conifer research

3  Chedly Kastally[a]*, Alina K. Niskanen[a]*, Annika Perry[b], Sonja T. Kujala[c], Komlan Avia[d], Sandra

4  Cervantes[a], Matti Haapanen[e], Robert Kesälahti[a], Timo A. Kumpula[a], Tiina M. Mattila[a,f], Dario I.

5  Ojeda[a,g], Jaakko S. Tyrmi[a], Witold Wachowiak[h], Stephen Cavers[b], Katri Kärkkäinen[c], Outi

6  Savolainen[a], Tanja Pyhäjärvi[a,i]**

7  [a]Department of Ecology and Genetics, University of Oulu, 90014 University of Oulu, Finland

8  [b]UK Centre for Ecology & Hydrology, Bush Estate, Penicuik, Midlothian, UK, EH26 0QB, UK.

9  [c]Natural Resources Institute Finland (Luke), Paavo Havaksen tie 3, 90570 Oulu, Finland

10  [d]Université de Strasbourg, INRAE, SVQV UMR-A 1131, F-68000, Colmar, France

11  [e]Natural Resources Institute Finland (Luke), Latokartanonkaari 9, FI-00790 Helsinki, Finland

12  [f]Department of Organismal Biology, EBC, Uppsala University, Uppsala, Sweden

13  [g]Norwegian Institute of Bioeconomy Research, Ås, Norway

14  [h]Institute of Environmental Biology, Faculty of Biology, Adam Mickiewicz University in Poznań,

15  Uniwersytetu Poznańskiego 6, 61-614 Poznań, Poland

16  [i]Department of Forest Sciences, University of Helsinki, 00014 University of Helsinki, Finland

17  *Contributed equally

18  **Corresponding author: Tanja Pyhäjärvi, Department of Forest Sciences, University of Helsinki,

19  FIN-00014 University of Helsinki, Finland, tanja.pyhajarvi@helsinki.fi

# Summary

Scots pine (*Pinus sylvestris*) is the most widespread coniferous tree in the boreal forests of Eurasia and has major economic and ecological importance. However, its large and repetitive genome presents a challenge for conducting genome-wide analyses such as association studies and genomic selection. We present a new 50K SNP genotyping array for Scots pine research, breeding programs, and other applications. To select the SNP set, we first genotyped 480 Scots pine samples on a 407 540 SNP screening array, and identified 47 712 high-quality SNPs for the final array (called 'PiSy50k'). Here, we provide details of the design and testing, as well as allele frequency estimates from the discovery panel, functional annotation, tissue-specific expression patterns, and expression level information for the SNPs or corresponding genes, when available. We validated the performance of the PiSy50k array using samples from breeding populations from Finland and Scotland. Overall, 39 678 (83.2%) SNPs showed low error rates (mean = 0.92%). Relatedness estimates based on array genotypes were consistent with the expected pedigrees, and the amount of Mendelian error was negligible. In addition, array genotypes successfully discriminate Scots pine populations from different geographic origins. The PiSy50k array will be a valuable tool for future genetic studies and forestry applications.

# Significance statement

38    Scots pine is an evolutionary, economically and ecologically impressive coniferous species but

39    its gigantic genome has limited studying e.g. the genetic basis of its functional trait variation. We

40    have developed a genotyping array that facilitates Scots pine genetic research and linking its

41    trait variation to genetic polymorphisms and gene expression levels across the genome.

# Introduction

43  Scots pine (*Pinus sylvestris*) is one of the world's most widely distributed conifers (Durrant *et al.*,

44  2016) and is dominant in forests across 145 million hectares in Northern Eurasia (Mason and

45  Aía, 2000; Mullin *et al*., 2011; Pyhäjärvi *et al.*, 2020). The species is an important source of

46  timber and other wood-based products (CABI, 2013) and boreal forests, of which Scots pine is

47  an essential part, are a significant carbon sink (Pan *et al.*, 2011). In addition to traditional timber,

48  pulp, paper, and energy production, more diverse uses for Scots pine biomass are currently

49  being developed (e.g., Agbor *et al.*, 2011; Rusanen *et al.*, 2019). The combination of large

50  biomass volumes, the species capability of adapting to varying marginal environments (Durrant

51  *et al.*, 2016), and modern genomic tools provide new possibilities for improving the desired

52  economic and ecological properties.

53  Breeding activities of Scots pine are centralized in Fennoscandia and the Baltic region, Sweden

54  and Finland having the most advanced breeding programs (Haapanen *et al.*, 2015). A first cycle

55  of selection and breeding was completed in the UK in the late 20th century (Lee, 2002), and

56  there is currently substantial interest in further improvement of the species, to reduce national

57  dependency on exotics. The genetic gains from breeding are delivered by seed from seed

58  orchards, comprising copies of field-tested plus trees (outstanding selections from wild stands).

59  Orchard-reproduced stock has been predicted to yield 20-25% improvement in per unit area

60  wood production above unimproved seed lots (Rosvall *et al.*, 2001, Haapanen *et al.*, 2016;

61  Jansson *et al.*, 2017). Forest tree breeding programs traditionally operate on large numbers of

62  individuals. Cost-effective genotyping platforms are therefore essential in incorporating

63  genomics to tree breeding schemes in the extent that is now true for cattle and crop breeding

64  (Grattapaglia *et al.*, 2018, Meuwissen *et al.*, 2016; Voss-Fels *et al.*, 2019).

65 Genotyping arrays are efficient and easy in comparison to other cost efficient sequencing

66 methods such as genotyping-by-sequencing (Pavan *et al.*, 2020). They are more reproducible

67 across studies, have less missing data and, importantly, require less bioinformatic pre-

68 processing (e.g., Darrier *et al.*, 2019). For forest tree species, SNP arrays are available for

69 walnut (Marrano *et al.*, 2019), Norway spruce (Bernhardsson *et al.*, 2020) and several eucalypt

70 species (Silva-Junior *et al.*, 2015). They have been used to build linkage maps (Silva-Junior and

71 Grattapaglia, 2015), develop genomic selection (GS) models (Tan *et al.*, 2017) and in genome-

72 wide association studies (GWAS) (Bernard *et al.*, 2020).

73 We foresee four primary applications for a new Scots pine SNP genotyping array:

74 1) Genomic selection

75 Genomic selection aims to predict the breeding value of an individual based on its genotypes,

76 where markers are assumed to be in linkage disequilibrium (LD) with the causal variation

77 (Meuwissen et al., 2001). In a set of individuals with both genotype and phenotype data (training

78 population), genomic prediction models are first generated and tested, leading to a prediction

79 equation. Genomic estimated breeding values (GEBV) can then be calculated from this

80 equation for individuals with genotype data only (e.g., Wray *et al.*, 2019). GS in trees shows

81 promising results (Isik 2014) and good predictive ability has been achieved with a few thousand

82 of SNPs (e.g., Bartholomé *et al.*, 2016; Calleja-Rodriguez *et al.*, 2020; Cappa *et al.*, 2019; Chen

83 *et al.*, 2018; Grattapaglia *et al.*, 2018; Lenz *et al.*, 2017; Resende *et al.,* 2012; but see

84 Thistlethwaite *et al.,* 2020).

85 GS has potential to increase genetic gains per unit of time when the breeding cycle can be

86 shortened, i.e. when reproductive maturity is reached soon after prediction of GEBV. There are

87 significant biological constraints to achieve this in Scots pine that reaches sexual maturity at 8-

88 20 years of age (Sarvas 1964). Nevertheless, genomic markers can provide other benefits by

89    reducing the phenotyping costs and achieving higher selection intensities in situations when a

90    large number of selection candidates are more easily genotyped than phenotyped (Calleja-

91    Rodriguez *et al.*, 2020; Grattapaglia *et al.*, 2018; Voss-Fels *et al.*, 2019). The operational

92    viability of such measures is obviously dependent on the costs of genotyping.

93    2) Pedigree construction

94    Genotyping data can be used to confirm and reconstruct pedigrees, identify labeling and

95    grafting errors, and estimate genomic relatedness among individuals. Realized genomic

96    relationships are potentially very useful for Scots pine breeding programs, as they allow more

97    accurate genomic prediction of breeding values. Genomic relationships can also help to bridge

98    unconnected progeny-testing series in a multi-environment genetic evaluation. Pedigree

99    reconstruction and parentage analysis using markers also opens opportunities for implementing

100   less costly breeding strategies, such as polymix breeding (Isik 2014).

101   3) Genome-wide association studies

102   Many of the most valued characteristics of Scots pine and other conifers are complex traits,

103   controlled by many genes. GWAS offers a way to detect the loci responsible governing the

104   variation, improving our understanding of the genetic architecture and biological mechanisms

105   behind these traits (Burghardt *et al.*, 2017; González-Martínez *et al.*, 2007 ; Neale and

106   Savolainen, 2004; Yeaman *et al.*, 2016). Large sample sizes are crucial for detecting the

107   associations, since polygenic traits are mostly controlled by numerous small effect

108   polymorphisms (Tam *et al.*, 2019; Yang *et al.*, 2010). A genome-wide SNP array is a convenient

109   tool for quickly genotyping many samples. Use of a common genotyping platform will allow for

110   comparison across studies.

111   4) Genetic mapping

112    High resolution genetic maps inform about the linkage relationships. They are important tools in

113    quantitative trait locus mapping (Lander and Thompson 1990). Combined with physical maps or

114    partial genomic information, they allow analysing of the recombination rate landscape of the

115    genome. To achieve high resolution, large numbers of progeny need to be genotyped, for which

116    SNP arrays are a cost efficient and powerful solution. When SNPs are anchored to scaffolds of

117    a genome assembly, maps derived from SNP array genotyping can be used to improve the

118    scaffolding of reference genomes, by linking together or re-ordering contigs (Fierst, 2015).

119    In addition, other potential applications for a SNP genotyping array include monitoring genetic

120    diversity, tracing geographic origin, estimating population structure, demographic inference,

121    identifying segregation distortion and identifying large structural variants.

122    SNP arrays are valuable universal tools for genetic fingerprinting and evaluation of diversity, but

123    they also have limitations. For instance, SNPs are typically accumulated close to or within

124    coding regions, because data are easier to obtain during SNP discovery using RNA-seq or

125    exome-targeted approaches than with whole genome sequences. Further, coding regions are

126    often of high interest and favored in array design. Also, as arrays only score preassigned SNPs

127    with a minimum minor allele frequency (MAF) threshold often applied, there is always an

128    ascertainment bias. This bias affects analyses performed on new individuals using the same set

129    of markers in two ways (McTavish and Hillis, 2015). First, loci with rare alleles in the discovery

130    population will not be scored. This may cause a bias in diversity estimation in favor of those with

131    common alleles. Second, at the population level, allele frequencies, and thus diversity, in

132    samples genetically close to the discovery panel will be biased upward compared to samples

133    from a distant lineage. Ascertainment bias thus is especially problematic for inferences requiring

134    information on rare alleles and not suitable for identifying new genetic variants. However, in

135    many analyses, the ascertainment bias can be taken into account if the original SNP discovery

136    panel and the array design is known (Clark *et al.*, 2005).

137 Here, we present the Axiom PiSy50k (Thermo Fisher Scientific), a new genotyping array for

138 Scots pine. We describe the different SNP sources and discovery panels and the selection

139 process used during the array design. The final array combines a set of high-performing SNPs

140 from a previously developed Axiom_PineGAP trans-specific SNP array of *Pinus* (Perry *et al.*,

141 2020) and a new set of curated SNPs originating from exome capture, RNA-seq, PacBio and

142 candidate gene studies (Table 1). We provide a detailed description of SNP discovery,

143 screening, filtering, evaluation of ascertainment bias, error rates, and the metadata we collected

144 during the design, such as gene expression and copy-number variation. We also explore the

145 array's capability to discriminate populations and reconstruct pedigrees.

# Results and discussion

146

## Array design

147

148 SNP choice and array design had four main stages: collection, filtering, *in silico* evaluation and

149 screening array evaluation (Figure 1). We first collected SNPs from eight data sets that differed

150 in sample size, sampling design, source material (RNA or DNA, tissue) and sequencing

151 technology (Sanger sequencing, PacBio, Illumina-seq). We filtered these initial data, tailoring

152 our approach to each data source's specific characteristics. We removed markers likely to be in

153 paralogous areas of the genome. Paralogy is a common problem for conifer species, which

154 have large genomes with a lot of repetitive elements (Neale *et al.*, 2014). Partly, this was done

155 by checking haplotypes from seed megagametophyte tissue, where observed heterozygosity

156 indicates false SNPs generated by paralogy. After the initial filtering, Thermo Fisher Scientific

157 conducted an *in silico* evaluation of 1.3 million SNPs and from these, we selected 407 540

158 SNPs of high interest and strong predicted performance.

8

## Performance of the screening array

160    We evaluated the performance of the screening array by genotyping a natural population

161    sample of 470 trees, six megagametophytes and four diploid embryos from full-sib crosses, all

162    from Finland. SNPs were assigned to six classes: Poly High Resolution (PHR, three well-

163    separated genotypesclusters), No Minor Homozygote (NMH, two well-separated genotypes

164    clusters, homozygous and heterozygous), Mono High Resolution (MHR, one homozygous

165    genotype cluster), Call Rate Below Threshold (CRBT), Off-Target Variant (OTV, more than

166    three clusters) and Others. When choosing SNPs for the PiSy50k array based on the screening

167    array, we considered conversion types PHR, NMH and MHR as successful. Of 407 540 SNPs in

168    the screening array, 245 149 (60.2%) converted successfully and 157 325 (38.6%) were

169    polymorphic (Table S1, Figure 2). The success rate varied among sources from 10% to 50%,

170    with lowest and highest rates in the LUKE candidate and UOULU candidate derived SNPs

171    respectively (Table S1, Figure 2). The latter set had already gone through several rounds of

172    verification and thus its higher conversion rate was not surprising. The genotyping success rate

173    at sample level was high; 476 (99%) samples had a call rate above the 97% threshold in the

174    conversion classes PHR, NMH, and MHR.

175    To assess the effects of ascertainment bias throughout the PiSy50k design, we evaluated its

176    effects on the screening array by investigating the minor allele frequency (MAF) distribution and

177    the genetic structure in the sample. The MAF distribution of the screening array is characterized

178    by a deficit of intermediate frequency alleles (MAF values between 0.15 and 0.5) compared to

179    the distribution expected based on the standard neutral model (SNM) (Figure 3A). This is not

180    surprising, as previous studies on Scots pine's genetic diversity across Europe have

181    demonstrated an overall deficit of intermediate alleles and excess of rare alleles in natural

182    populations of this species compared to the SNM (Tyrmi *et al.*, 2020; Pyhäjärvi *et al.*, 2020 and

9

183    references therein). However, the pattern of rare alleles in the screening set differs from the one

184    in earlier studies. We observed an excess of rare allele classes (MAF between 0.007 and 0.15,

185    Figure S1), but a deficit in the extremely rare classes (MAF below 0.007, Figure S1), as

186    expected from ascertainment bias.

187    In addition, ascertainment bias influenced the estimates of genetic structure among samples.

188    Principal component (PC) analyses of the screening array genotypes of UOULU RNAseq and

189    UOULU exomeFEB2019 sets clearly separate trees included in the discovery panel from the

190    rest of the samples (Figure S2 a and c). The ascertainment bias was more subtle in the other

191    sources, even when samples from the discovery panel were genotyped (Figure S2 e). This

192    difference was due to the larger size of the other discovery panels (Table 1). The effect of

193    ascertainment bias was particularly severe when the exact discovery panel samples or their

194    close relatives were included (Figure S2 a and c). For most applications and datasets not

195    related to the discovery panels, these effects on genetic structure are unlikely to be so extreme,

196    but we recommend that users of the array carefully consider sample origin when performing

197    analyses.

198    Finally, from the remaining 75 629 SNPs, we excluded SNPs with heterozygous calls in

199    megagametophyte haploid samples (but allowing one error in SNPs from three high priority

200    sources, see Table 1) or with more than one Mendelian error. We also pruned SNPs in high LD

201    ($r^2 > 0.9$), keeping the SNPs with the higher minor allele frequency (MAF) from each such pair.

202    From the remaining loci, we first retained all SNPs from high priority sources and favored SNPs

203    with higher MAF in the remaining set. SNPs in a highly outcrossing wind pollinated natural

204    population of Scots pine are expected to be in Hardy-Weinberg (HW) equilibrium and we used

205    deviation from HW ($p$-values < 0.001) to identify and filter out potentially paralogous and other

206    error prone SNPs. As expected, the markers selected for the PiSy50k array deviated less from

207    the HW expectations and showed less extreme heterozygosities compared to all screening

10

208    array markers before selection (Figure S3). The final PiSy50k array includes 47 712 SNPs.

## Performance of the PiSy50k array

210    The 47 712 SNPs in the final PiSy50k array were in 31 657 contigs (average of 1.51 SNPs per

211    contig). Of the eight data sources, markers from RNA-seq origin were the most prevalent (44%;

212    Table S2). The majority of markers have been used in previous studies and come associated

213    with various information depending on the source, including functional annotation, gene

214    expression at the tissue level, and allele frequency estimates in up to 20 European populations

215    (Supporting Data S1).

216    Altogether, 1 619 markers derived from ProCoGen haploid (1 544) and diploid sources (75)

217    were located on one of the 4 226 scaffolds mapped on *P. taeda* linkage map (Westbrook *et al.*,

218    2015; Figure 4, Table S3). There was an average of 134 SNPs per linkage group (LG) and they

219    were homogeneously distributed among LGs. Even though the majority of SNPs do not have a

220    known position on the map yet, the quick genotyping of large numbers of progeny with the

221    PiSy50k array could be used to improve the genetic map of Scots pine and help anchor

222    genomic reads, scaffolds and SNPs at the chromosome scale in the future.

223    We evaluated the performance of the PiSy50k array by genotyping 2 688 samples from Finland

224    (2178, including 14 controls), Scotland (496), Australia (3), and Estonia (11). Of these, 2 308

225    samples had call rates above 97% (85.9% of samples), the recommended threshold for Axiom

226    genotyping arrays. In total, 40 405 (84.69 %) markers were successfully converted of which 39

227    678 markers were polymorphic (Table S4).

228    Of the 21 control samples, three needle and six megagametophyte samples passed the 97%

229    CR threshold (Table S5). Of the six megagametophyte samples, one replicated pair was

230    recovered. Based on the five control samples retained (three needles and one

11

231   megagametophyte pair), the error rates were relatively low (mean 0.83%). The error rate in the

232   subset of SNPs shared with the Axiom_PineGAP suggests a similar, or slightly lower, error rate

233   in the PiSy50k (mean 0.52% compared to 0.64% in the Axiom_PineGAP). Overall, these values

234   are close to those obtained in other arrays, e.g. 0.8% in the walnut genotyping array (Marrano

235   *et al.*, 2019), 0.1% in Affymetrix GeneChip Human mapping 50k Array (Saunders *et al.*, 2007),

236   or ranging between 0.03% and 0.05% in the Axiom Apple480K genotyping array (Bianco *et al.*,

237   2016).

238   Of the 930 markers with errors among pairs (including both needle and megagametophyte

239   controls), the majority (N = 916) were not shared among controls. This suggests that the error

240   probably occurred during the genotype call for a single sample only, as opposed to the marker

241   itself being unreliable. There are 14 markers for which errors were observed among both

242   megagametophyte and needle controls and they are indicated in Supporting Data S2.

243   Comparison of markers shared between the PiSy50k and Axiom_PineGAP arrays (N = 7592)

244   using the needle control present on both arrays also showed low error rates (mean 0.55%,

245   Table S5) indicating cross-array reproducibility, which allows data obtained by the two arrays to

246   be combined.

247   To confirm that the variants at the selected SNPs in the PiSy50k array are indeed allelic (not

248   paralog), we assessed the heterozygosity levels of the megagametophyte samples. The two

249   megagametophyte replicates have very low heterozygosity levels (mean 0.89%) compared to

250   the needle replicates (mean 29.30%), suggesting a low level of errors due to paralogy. Of the 40

251   405 converted markers, 38 906 were homozygous in both replicates, 1 060 were 'no call' in at

252   least one replicate, 165 were heterozygous in both replicates and 274 were homozygous in one

253   replicate and heterozygous in the other. The SNPs that were heterozygous in the

254   megagametophyte samples are indicated in Supporting Data S2.

255    To evaluate the potential of the PiSy50k array for pedigree reconstruction and assess the

256    proportion of Mendelian errors in the array, we analyzed the pairwise relatedness of the full-sib

257    progeny and their parents in a subset of 135 trios across 10 families of our sample. By plotting

258    the kinship coefficient (K, (Manichaikul *et al.*, 2010)) against the proportion of sites where

259    individuals share no allele (IBS0), we identified four distinct groups (Figure 5a): (1) known

260    parent-offspring pairs (mean +/- standard deviations: K = 0.245 +/- 0.004, IBS0 = 0.001 +/- 2e-

261    04), (2) full-sibs (K = 0.246 +/- 0.027, IBS0 = 0.015 +/- 4e-03), (3) half-sibs (K = 0.120 +/- 0.018,

262    IBS0 = 0.030 +/- 4e-03), and finally (4) the remaining unrelated pairs (K = -0.002 +/- 0.009; IBS0

263    = 0.059 +/- 2e-03). We separated parent-offspring pairs from full-sibs, which have expected K

264    values close to 0.25, using the IBS0 statistic (equal or close to 0 between a parent and an

265    offspring but with higher values between siblings (Manichaikul *et al.*, 2010)). Within each family,

266    the K estimates were around the expected value of 0.25, while between families it was close to

267    0, except for progeny pairs between families 5 and 31, and families 14 and 20, which shared a

268    common parent and had a K estimate around 0.125, as expected for half-sibs (Figure 5). The

269    pedigree relationships identified with PiSy50k matched those expected from the crossing

270    design, demonstrating the array's power to resolve relatedness structure and reconstruct

271    pedigrees, a critical feature for a multitude of applications in tree breeding and genetics: GWAS,

272    GS, breeding program management and seed production.

273    To further assess the error rate in the PiSy50k data, we evaluated the number of Mendelian

274    errors (ME) within each family. We examined all 40 405 SNPs in 135 trios and identified 16 040

275    errors across 5 837 loci (mean error rate per locus = 0.29%; Figure S4a). More than 98% of all

276    SNPs had a ME below 5%. Across families, we identified an average of 1 604 errors per family,

277    majority in different SNPs across families (Figure S4b: 4277 SNPs with an error only in a single

278    family and 1110 in at least two). These values are in line with the ME measured in other arrays

279    (Bernhardsson *et al.*, 2020; Silva-Junior *et al.*, 2015).

## Genetic diversity

281   To explore the power of genotypes from the PiSy50k array to discriminate trees from different

282   geographic origins, we ran a principal component analysis (PCA) using a subset of 120 samples

283   from different localities in Scotland and Finland (Figure 6). The first two PCs separated two main

284   groups consistent with the two countries of origin. We then ran PCAs using only samples from

285   each country. Although no distinct groups appeared in those analyses, some differentiation was

286   found between samples from different geographic origins in Scotland (Figure 6b) – a level of

287   geographic resolution not previously possible. In the Finnish subset, variation was more

288   homogeneous with less geographic structure (Figure 6c), although samples from Northern

289   origins were located slightly apart from samples from Southern and Central origins.

290   To assess the effects of ascertainment bias on the MAF distribution in the PiSy50k array, we

291   compared the frequency distributions obtained from the array to a previously published exome

292   capture dataset (Tyrmi *et al.*, 2020) (Figure 3B). We observed a similar but stronger effect of

293   ascertainment on the MAF estimated with the PiSy50k array genotyping results than with the

294   screening array results. Indeed, in the PiSy50k results, the distribution reaches a maximum at

295   frequency 0.13, with decreasing frequencies of lower MAF values, as opposed to the screening

296   array where the peak is at the lowest allele frequency class. This could be explained by the

297   more stringent filtering of SNPs with low allele frequencies when selecting markers for the final

298   PiSy50k set, whereas there was no intentional allele frequency filtering from the source data to

299   the screening set. In addition, the discovery process naturally has an inherent filter for allele

300   frequency, which is the sample size of the discovery panel.

301   In summary, PiSy50k is a novel genotyping tool for Scots pine, an economically important and

302   widely distributed conifer. It greatly improves the genotyping capacity for the species, which will

303   facilitate wide application of modern breeding tools and supports the development of a new,

304   forest-based bioeconomy. The metadata provided connects the genotyping data to functional

305 properties via annotations and tissue-specific expression patterns. Low error rates indicate high

306 reproducibility even across the previous Scots pine array Axiom_PineGAP (Perry *et al.*, 2020),

307 hence new datasets will be back-compatible and all new work will add value to our knowledge of

308 the species.

# Experimental procedure

309

## Selection of SNPs for initial screening

310

311 *ProCoGen haploid and diploid sets*.

312 The ProCoGen haploid and diploid sets were generated with two exome-capture experiments

313 both based on the same bait set used by Tyrmi *et al.* (2020). A total of 177 trees collected

314 across Europe, from Spain to northern Finland, were genotyped using DNA extracted from

315 megagametophyte tissue (haploid set, 109 samples, 12 populations) or needles (diploid set, 68

316 samples, 8 populations). Bait design, DNA extraction, library preparation, and sequencing steps

317 followed the procedure described in (Tyrmi *et al.*, 2020). We processed the sequences

318 generated to identify SNPs following the same method described in Tyrmi *et al.* (2020) for the

319 haploid set, but applied a few adjustments for the diploid set: we used BWA (Li, 2013) for

320 mapping reads and used samtools v0.9 (command *mpileup*, default parameters) (Li *et al.*, 2009)

321 for variant calling. To filter potential paralogs, we  removed loci with heterozygous calls in the

322 haploid set or significantly departing from the HWE in the diploid set (PLINK v1.90b5.2 (Chang

323 *et al.*, 2015) command --hardy, at alpha = 0.05). During this procedure, we excluded one

324 haploid sample with an exceptionally high proportion of heterozygous calls. Finally, we excluded

325 all SNPs within 50 bp distance of these markers. We retained 248 591 and 32 649 SNPs in the

326 haploid and diploid sets, respectively.

## *UOULU exomeFEB2019*

328    We used 95 504 SNPs identified in exome capture of a family originating from Punkaharju ISS,

329    in southeast Finland: a cross between Maternal tree 463 and paternal tree 485 (Kesälahti *et al.*,

330    In Prep). The material sampled consisted of: needles of both parental trees, one

331    megagametophyte of the paternal tree, two megagametophyte of the maternal tree from open-

332    pollinated seeds, and, from two seeds of the cross progeny, two embryos and a

333    megagametophyte were sampled. We excluded positions with depth below 4 per genotype. We

334    removed twenty-five base pairs both upstream and downstream from each heterozygous site

335    found in haploid megagametophyte as potential areas with paralog or mapping issues.

## *UOULU RNA-seq*

337    The UOULU RNA-seq set refers to markers derived from RNA-seq data (Ojeda *et al.*, 2019)

338    originating from five tissues (needle, phloem, vegetative bud, embryo and megagametophyte) of

339    six unrelated individuals of Scots pine (but 18 haploid genomes when accounting for diploidy

340    and paternal contribution in embryos) collected from Punkaharju ISS. We considered 1 349 291

341    SNPs obtained by mapping RNA-seq reads to the Scots pine reference transcriptome

342    (https://a3s.fi/pinus_sylvestris_transcriptome_public_data/Trinity_CD-HIT.fa). From this initial

343    set, we first excluded markers identified in contigs associated with potential contaminants (fungi

344    or microbes) (Cervantes *et al.*; Ojeda *et al.*, 2019)

345    (https://a3s.fi/pinus_sylvestris_transcriptome_public_data/Trinity_guided_gene_level_info.txt).

346    Second, we removed heterozygous SNPs in haploid samples. Finally, we compared the

347    genotypes called in megagametophyte, embryo and diploid tissues collected from the same tree

348    to identify and exclude loci with Mendelian errors. In total, we retained 736 827 SNPs.

349    For the *UOULU RNA-seq* set, we provide information about the predicted multi-copy status,

350    orthologous genes identified in *P. taeda* (Zimin *et al.*, 2014) and *P. lambertiana* (Stevens *et al.*,

351    2016) based on blastn results (see details in Ojeda *et al.*, 2019), and expression levels and

352    tissue-specificity in five tissues (Cervantes *et al.*). This information is available in Supporting

353    Data S1.

## *UOULU candidate*

355    The UOULU candidate set contains SNPs reported in multiple publications and genetic

356    databases on various candidate genes of Scots pine. This set includes SNP markers used in

357    Kujala *et al.* (2017), and additional SNPs from phenology related genes (Kujala and Savolainen

358    2012; Palme *et al.,* 2008; Pyhäjärvi *et al.,* 2007, Wachowiak *et al.,* 2009), stress and phenology

359    related genes (Avia *et al.*, 2014), polyamine genes (Vuosku *et al.*, 2018, 2019), genes from

360    comparative resequencing projects (Wegrzyn *et al.*, 2008; Grivet *et al.,* 2017), and markers

361    identified in sequences from the Evoltree EST database (www.evoltree.eu). Additionally, for a

362    subset of those markers, we have collected allele frequency estimates from two genotyping

363    assay experiments on 426 Scots pine trees (data unpublished). These SNPs, referred to as

364    UOULU candidate VIP in the metadata, were given higher priority during the array manufacture,

365    in both the screening and PiSy50k arrays, by increasing their probeset counts and, this way,

366    improving their call rates during the genotyping.

## *LUKE candidate*

368    The LUKE candidate set comprises SNPs extracted from candidate genes related to phenology

369    (e.g. Bouché *et al.*, 2016) and genes of the primary and secondary metabolism pathways active

370    during heartwood formation (Lim *et al.*, 2016). DNA libraries targeting these candidate genes

371    were produced from one individual of Southern Finnish origin and sequenced using a PacBio

372    sequencer (Kujala *et al.,* in prep). We used the long PacBio sequences as a reference to map

17

373    short reads from exome captures of megagametophyte samples of Scots pine collected across

374    Europe (Tyrmi *et al.,* (2020) excluding samples from Baza, Spain) with BWA mem (Li, 2013).

375    Since a preliminary variant calling based on this initial mapping resulted in a large number of

376    errors (heterozygous calls in haploid samples), we isolated short reads mapping to individual

377    PacBio contigs and re-assembled them with MIRA (Chevreux, 2007) for each individual. We

378    then aligned the resulting individual re-assemblies to each other with cap3 (Huang and Madan,

379    1999), and called variants using bcftools (commands *mpileup* and *call*). In addition, some SNPs

380    were identified and included solely as being polymorphic within the reference individual.

381    *UKCEH sets 1 and 2*

382    We used SNPs collected during the Axiom_PineGAP (Thermo Fisher Scientific) array design

383    (Perry *et al.*, 2020) and from the comparative transcriptomics of four pine species (*P. sylvestris*,

384    *P. mugo*, *P. uncinata* and *P. uliginosa*) by Wachowiak *et al.* (2015). Briefly, we identified 196

385    636 polymorphic positions from transcriptomes, candidate gene sequences and markers from

386    previous population genetic studies on the four above mentioned pine species. From these, we

387    retained two distinct sets: (1) UKCEH1, comprised of 20 795 successfully converted SNPs from

388    the Axiom_PineGAP array, and (2) UKCEH2, a set of 175 841 SNPs including 29 034 SNPs

389    from the Axiom_PineGAP array which were not successfully converted, 31 897 SNPs that

390    passed the initial filtering during the design but were not included in the final array and 114 910

391    SNPs identified by Wachowiak *et al.* (2015) which were polymorphic in Scots pine but not

392    included in the Axiom_PineGAP array design.

# 393    SNP scoring for inclusion in the screening array

394    For each retained site, we built 71-mer probes by extracting up to 35 bp up- and downstream

395    from the source references. We submitted 1 317 798 probes to Microarray Research Services

396   Laboratory (Thermo Fisher Scientific), Santa Clara, US, for scoring (Table S1). During this step,

397   probes' score were downgraded if: they contained polymorphic sites within 35 bp distance of the

398   focal marker (interfering polymorphism), they were mapped to highly repetitive regions of the

399   genome (using TrinityCD-HIT.fasta.gz and Pita 1.01

400   (https://treegenesdb.org/FTP/Genomes/Pita/v1.01/genome/Pita.1_01.fa.gz) as references for

401   RNA and DNA based probes respectively), or were highly similar to other probes. Each marker

402   was given a classification: 'Recommended', 'Neutral', 'Not recommended' or 'Not possible'.

403

404   Based on Thermo Fisher Scientific's evaluation and the available metadata on each data

405   source, we established the following priority groups (in order of priority): (1) the 20 795 high

406   quality SNPs from the Axiom_PineGAP array, (2) all recommended or neutral markers identified

407   by Thermo Fisher Scientific, (3) UOULU candidate markers, (4) LUKE candidate markers, (5)

408   from the 'not recommended' set in the ProCoGen haploid set, SNPs of high interest identified in

409   (Tyrmi *et al.*, 2020), (6) SNPs with less than 50% of missing data in the discovery panel from the

410   'not recommended' set in the ProCoGen sets, and finally, (7) we relaxed the filtering criterion

411   used by Thermo Fisher Scientific and selected the best markers in the remaining set. More

412   specifically, we relaxed the wobble count filter threshold (number of polymorphic sites on the

413   same 71-mer) from < 4 to < 6, based on the assumption that a high proportion of the variable

414   sites are associated with rare alleles, and thus interfering polymorphism should have lower

415   impact on the probe performance in the case of Scots pine. During the screening array

416   manufacture, out of the 428 516 SNPs retained, a total of 407 540 markers were fitted on the

417   array.

# Screening set genotyping

419   The screening set of 407 540 SNPs was used to confirm the normal segregation of

420    polymorphism in a larger sample from a natural population, to identify potential deviations from

421    HW equilibrium, indications of paralog mapping — such as heterozygote sites in haploid

422    samples, deviations from Mendelian segregation, and identification of loci in strong LD with

423    each other. To this end, we used the screening array to genotype 480 samples of Scots pine

424    from Punkaharju ISS population, including: 470 diploid needle samples from adult trees, six

425    haploid megagametophytes and four diploid embryos. Two families, "463 x 485" and "320 x

426    251", with two parents and two offspring (embryos) from each were used to estimate Mendelian

427    error rate.

428    DNA was extracted from dry needles and fresh megagametophytes using E.Z.N.A.® SP Plant

429    DNA Kit (Omega Bio-tek, Inc.). Genotyping and array manufacturing for the screening set was

430    performed by Thermo Fisher Scientific at Santa Clara, US. Genotype calling was performed by

431    Thermo Fisher Scientific (Applied Biosystems™ Axiom™ Genotyping Services) following the

432    Axiom Best Practices Workflow (Axiom Genotyping Solution Data Analysis Guide). In short,

433    genotype clusters were defined using samples with quality control call rate (QC CR) >= 0.97

434    and dish quality control rate (dQC) >= 0.82. The markers were classified into five conversion

435    categories: PolyHighResolution (PHR), NoMinorHom (NMH), MonoHighResolution (MHR),

436    CallRateBelowThreshold (CRBT), Off-Target Variant (OTV), and Other. We retained markers

437    only from classes PHR and NMH with call rate (CR) >= 0.97 in the subsequent analyses of the

438    screening array and for inclusion on the PiSy50k array.

439    During the array design of both screening and PiSy50k arrays, identical SNPs discovered

440    independently across different sources were identified and merged. To keep track of as much

441    information as possible for those markers, we recorded their common presence and IDs in

442    different sources but eventually assigned a single authoritative origin.

## Selection of markers for the PiSy50k array

443

444 For the PiSy50k array, we filtered the markers based on their performance on the screening

445 array prioritizing markers in candidate genes of interest or markers that performed well in the

446 Axiom_PineGAP array (Perry *et al.*, 2020). These markers were within the Axiom Best Practices

447 Workflow default quality thresholds (see above). For each marker with conversion type PHR or

448 NMH, we estimated MAF and tested departure from HW equilibrium (exact test) for 466

449 individuals, excluding the haploid megagametophyte samples, the offspring samples and four

450 samples with QC CR < 0.97 using PLINK version 1.9 (Purcell *et al.*, 2007). We estimated the

451 number of Mendelian errors in PLINK using the family data.

452 We excluded markers deviating from HW equilibrium ($p < 0.001$) and markers with more than

453 one Mendelian error. Markers from the candidate gene sources (LUKE candidate and UOULU

454 candidate) were selected using a lenient inclusion threshold of MAF >= 0.01 and marker call

455 rate > 0.90, which also included markers from the Thermo Fisher Scientific conversion type "call

456 rate below threshold". We filtered SNPs from the Axiom_PineGAP array first to include markers

457 with MAF >= 0.05. To increase the number of well performing markers, we also included

458 markers with MAF >= 0.05 in previously genotyped European samples (Perry *et al.* 2020).

459 To avoid markers in paralogous genomic regions, we excluded markers with heterozygous call

460 in the haploid megagametophyte samples except in three high priority sources (UKCEH1, LUKE

461 candidate and UOULU candidate) for which we allowed at most one, erroneous, heterozygous

462 call per marker. We further granted 381 markers of high interest from sources UOULU

463 candidate (335) and UOULU RNA-seq (23) a higher probeset count in the array to increase their

464 call rate. Finally, to remove the excess from the retained set, we excluded markers from the low

465 priority sources with lowest MAF (MAF after filtering >= 0.08). The final number of markers for

466 PiSy50k was 47 712 (Figure 1). The distribution of the markers by source is shown in Table S2.

21

467 To inspect how SNP selection for the PiSy50k array affected HW deviation compared to the

468 screening array on average, we plotted the observed *p*-alues from the exact HW tests against

469 the expected *p*-values based on the null distribution in a cumulative Q-Q plot before and after

470 SNP selection. We compared the observed *p*-values of 10 000 random loci against 100 samples

471 drawn from the null distribution using HardyWeinberg package (Graffelman, 2015) in R (The R

472 Project for Statistical Computing)(version 3.6.3). We also illustrated the distribution of genotypes

473 with respect to HW expectations in ternary plots showing genotypes before and after the

474 PiSy50k SNP choice.

475 To assess the effects of ascertainment bias on the screening array, we ran two analyses. First,

476 we plotted the MAF distribution for loci with conversion types PHR or NMH (n loci without

477 missing data = 56 693, n individuals = 466) against the expected MAF assuming a standard

478 neutral model (Tajima, 1989). Second, we looked at the effects of ascertainment bias on the

479 inference of genetic structure by conducting PCAs using the R package pcadapt (Privé *et al.*,

480 2020). We performed PCAs using SNPs separately from each source and retained the results

481 from two sets where we observed the strongest effects of ascertainment bias, from sources

482 UOULU RNA-seq and UOULU exomeFEB2019, and one in which the effects were minimal, the

483 ProCoGen haploid sources. To further illustrate the root cause of the observed biases, we

484 performed those PCAs with and without the individuals present in the original discovery panel

485 and driving the patterns observed.

## Linkage map position of PiSy50k markers

487 To assess whether markers from the PiSy50k array are homogeneously distributed across all

488 chromosomes, we positioned them on a genetic map produced for *P. taeda* by Westbrook *et al.*

489 (2015) comprised of 12 linkage groups (LG) and to which contigs from *P. taeda* reference

490 genome Pita v1.01 have been mapped. We included all PiSy50k SNPs previously mapped to

22

491  one of the contigs or scaffolds from the same reference genome (data sources ProCoGen

492  haploid and diploid). When a given SNP was outside the aligned segment of the reference

493  contig, we used the closest position effectively aligned on the genetic map from the same contig

494  as a reference point to infer the position of the focal SNP on the map, assuming that the

495  physical distances covered by single contigs from the Pita v1.01 reference genome to be

496  negligible compared to the size of each individual LG.

## PiSy50k array genotyping

497

498  We tested the PiSy50k array performance by genotyping 2 688 samples (across seven plates).

499  The 2688 samples consisted of 317 Finnish plus trees, 1847 full-sib offspring from the Finnish

500  breeding population, 489 Scottish samples, three Austrian samples, 11 Estonian samples and

501  21 controls. The needle control was a single tree from Scotland, UK, and was included on each

502  genotyping plate; this sample had also been genotyped on the Axiom_PineGAP array. In

503  addition, seven haploid megagametophyte samples were genotyped twice, such that each

504  sample was genotyped on two different random plates. Other samples were randomized over

505  the plates such that the different geographic locations and sample categories (plus trees and

506  offspring) were spread on all plates to avoid plate effects that may bias genotyping results of a

507  specific sample category.

508  The arrays were manufactured by Thermo Fisher Scientific (Waltham, MA, US) and genotyping

509  was conducted by University of Bristol Genomics Facility (Bristol, UK). Needle samples (n = 2

510  674, including 7 controls) were dried and stored in bags with silica gel. For megagametophyte

511  samples (7 control samples included twice each), germination was initiated by placing the seeds

512  on a moist filter paper inside a petri dish for 24 hours at room temperature. Seeds were then

513  dissected under a microscope to separate megagametophyte from the embryo tissue. The DNA

514  from Finnish and Estonian samples was extracted using E.Z.N.A.® SP Plant DNA Kit (Omega

23

515 Bio-tek, Inc.). DNA of Scottish needles were extracted using a Qiagen DNeasy Plant kit and

516 checked visually on a 1 % agarose gel. DNA was quantified with a Qubit spectrophotometer.

517 We performed the genotype call using Axiom Analysis Suite (version 5.1.1.1) following the

518 Axiom Best Practices Workflow with default parameters concordantly to the screening array

519 genotype calling, except for the plate QC threshold for average call rate for passing samples,

520 which we set to 0.97. We retained the markers in the PHR and NMH conversion classes for

521 analyses.

# Evaluation of the PiSy50k array performance

## Error rate and heterozygosity in haploid samples

524 We genotyped 21 control samples to estimate error rates for each array: one needle and two

525 megagametophyte controls per plate, with replicate megagametophyte pairs arranged over

526 sequential plates. We estimated the error rates as the proportion of calls which did not match

527 among pairs of controls across plates (excluding calls where one or both were missing). We

528 also measured the heterozygosity in megagametophyte samples to assess probe specificity and

529 identify putative paralogous markers in the PiSy50k array.

## Pedigree inference and mendelian error rate

531 We used a subset of 153 samples from 10 crosses, including 18 parents and their 135 offspring,

532 to estimate the coefficients of kinship (K) and the proportion of sites where individuals share no

533 allele (IBS0) between all pairs using converted SNPs (40 405) with KING v2.2.5 (options --

534 related --degree 2) (Manichaikul *et al.,* 2010). We estimated the Mendelian error rate within

535 each family independently using PLINK v1.90b5.2 (option --mendel).

## Population clustering and ascertainment bias.

537 To evaluate the power of the PiSy50k in discriminating samples from different origins, we used

538 a subset of 120 plus-tree samples: 30 samples from Scotland, grouped in four geographic

539 areas, and 30 samples from Southern, Central and Northern Finland each. We assessed the

540 genetic structure by performing three PCAs: using all 120 samples, the 90 Finnish samples or

541 the 30 Scottish samples separately. We used the function *prcomp* (core R, with scaling and

542 centering options enabled) after replacing missing data for a given genotype by the locus' allele

543 frequency. Finally, to assess the effect of ascertainment bias on the MAF generated with

544 PiSy50k, we compared the MAF distribution of the Finnish subset of 90 plus trees to the one

545 obtained using exome capture data of Scots pine trees published in Tyrmi *et al.* 2020. From the

546 published vcf file, we extracted the data of 42 megagametophyte samples from four Finnish

547 populations (Inari, Kälviä, Kolari and Punkaharju). We then replaced genotypes with depth

548 below 5 with missing data and kept only loci with a minimum call rate of 50%. Finally, to have

549 comparable MAF distributions, we downsampled both distributions to a sample size of 30.

# Author Contributions

550

551   Design of the study: AKN, AP, CK, KK, MH, OS, StC, STK, TP. Field and laboratory work: AKN,

552   AP, SaC, STK, TAK, TP, RK, StC. Computational analyses: AKN, AP, CK, DIO, JST, KA, STK,

553   TMM, TP, WW. Initial draft of the manuscript: AKN, CK, TP. Final manuscript: All authors.

# Acknowledgements

554

# Conflict of interest

564

565   The authors declare no conflict of interest.

# Supporting information legends

566

567 The following material is included in the supporting information:

568 • Supporting Figures S1 to S4

569 • Supporting Tables S1 to S5

570 • Supporting Methods S1: "Additional steps/details in selecting markers from screening
571 array to PiSy50k array".

572 • Supporting Data S1 and S2

## Legends:

573

574 **Figure S1.** Minor allele frequencies for the Intensive Study Site Punkaharju (southeast Finland)

575 population (N=466) and 56 693 SNPs without missing data in the screening array. The red line

576 illustrates the expected neutral MAF (Tajima, 1989). Note that this figure is identical to Figure 3

577 but is represented with a logarithmic scale on both the x- and y-axes.

578 **Figure S2.** Principal component analysis on the screening array data illustrating the

579 ascertainment bias on the observed genetic structure. (a, c, e) Analysis including samples used

580 in SNP discovery panels of each SNP source, discovery individuals are highlighted and labelled,

581 except in e) for clarity. (b, d, f) Analysis excluding samples used in SNP discovery. SNP

582 sources: (a, b) UOULU RNA-seq (48 357 SNPs), (c, d) UOULU exomeFEB2019 (6 137 SNPs)

583 and (e, f) ProCoGen haploid (23 204 SNPs).

584 **Figure S3**. Hardy-Weinberg equilibrium (HW) test results for the screening array data before

585    filtering (a,b) and for the selected set for the PiSy50k (c,d). (a,c) Q-Q plots comparing the *p*

586    values expected based on the null distribution against the observed *p* values from the exact HW

587    tests of 10 000 random SNPs on the screening array before (a) and after (c) selecting markers

588    for the PiSy50k array. The green line indicates the expected under HW. (b,d) Ternary plots

589    showing the genotype frequencies of 10 000 random SNPs on the screening array before (b)

590    and after (d) selecting markers for the PiSy50k array. Blue and red dots are markers

591    respectively following or deviating significantly from the HW expectations (Chi-square test at

592    alpha level 0.001).

593    **Figure S4**. Mendelian errors (ME) of the PiSy50k identified in 40 405 SNPs genotyped in 135

594    trios (10 crosses). (a) Distribution of ME across loci, the red line indicates the mean error rate

595    across loci (0.29%). (b) ME across families (bars at 0 and 1 indicate the number of SNPs with

596    no ME and with ME in only one family).

597    **Table S1.** Conversion type for markers from each data set in the screening array based on

598    individuals with call rate 97% or above. We included the markers with the PHR and NMH

599    conversion types (in bold) in the selection of markers for the PiSy50k array. PHR = Poly High

600    Resolution, NMH = No Minor Homozygote, MHR = Mono High Resolution, CRBT = Call Rate

601    Below Threshold, OTV = Off-Target Variant. Values in parenthesis are the proportion (per cent)

602    of each conversion type in each data set.

603    **Table S2.** Number and proportions of markers from each source at different steps of the

604    PiSy50k array design.

605    **Table S3.** Distribution of PiSy50k markers on *P. taeda* linkage groups (Westbrook *et al.,* 2015).

606    **Table S4.** Conversion type for markers from each data set in the PiSy50k array based on

607    individuals with call rate 97% or above. We included the markers with the PHR and NMH

608    conversion types (in bold) in further analyses. Count and proportion (%) of each conversion type

609    is given within each data set.

610    **Table S5.** Evaluation of the PiSy50k array for the control samples with call rate above 97%.

611    Values before the forward slash indicate estimates obtained from the full PiSy50k array (40 405

612    SNPs). Values after the forward slash indicate estimates obtained from the subset of SNPs and

613    the needle sample also genotyped by the Axiom_PineGAP array (7 592 SNPs). CR: call rate;

614    Het: heterozygosity. Mean pairwise error rate estimated as percentage of calls among control

615    pairs that were different (excluding markers which had missing data in at least one of the pairs).

616    **Methods S1.** Additional steps/details in selecting markers from screening array to PiSy50k

617    array.

618    **Data S1.** The metadata for markers included on the PiSy50k array.

619    **Data S2.** Shared errors across controls identified in the error evaluation of the PiSy50k, see the

620    main text.

# References

622  Agbor, V.B., Cicek, N., Sparling, R., Berlin, A., and Levin, D.B. (2011) *Biomass pretreatment:*
623  *fundamentals toward application*. *Biotechnol. Adv.*, **29**, 675–685.

624  Ahtikoski, A., Ojansuu, R., Haapanen, M., Hynynen, J., and Kärkkäinen, K. (2012) *Financial*
625  *performance of using genetically improved regeneration material of Scots pine (Pinus*
626  *sylvestris L.) in Finland*. *New Forests*, **43**, 335–348.

627  Avia, K., Kärkkäinen, K., Lagercrantz, U., and Savolainen, O. (2014) *Association of*
628  *FLOWERING LOCUS T/TERMINAL FLOWER 1-like gene FTL2 expression with growth*
629  *rhythm in Scots pine (Pinus sylvestris)*. *New Phytol.*, **204**, 159–170.

630  Bartholomé, J., Van Heerwaarden, J., Isik, F., Boury, C., Vidal, M., Plomion, C., and Bouffier, L.
631  (2016) *Performance of genomic prediction within and across generations in maritime pine*.
632  *BMC Genomics*, **17**, 604.

633  Bernard, A., Marrano, A., Donkpegan, A., Brown, P.J., Leslie, C.A., Neale, D.B., et al. (2020)
634  *Association and linkage mapping to unravel genetic architecture of phenological traits and*
635  *lateral bearing in Persian walnut (Juglans regia L.)*. *BMC Genomics*, **21**, 203.

636  Bernhardsson, C., Zan, Y., Chen, Z., Ingvarsson, P.K., and Wu, H.X. (2020) *Development of a*
637  *highly efficient 50K SNP genotyping array for the large and complex genome of Norway*
638  *spruce (Picea abies L. Karst) by whole genome re-sequencing and its transferability to*
639  *other spruce species*. *Mol. Ecol. Resour.*

640  Bianco, L., Cestaro, A., Linsmith, G., Muranty, H., Denancé, C., Théron, A., et al. (2016)
641  *Development and validation of the Axiom(®) Apple480K SNP genotyping array*. *Plant J.*,
642  **86**, 62–74.

643  Bouché, F., Lobet, G., Tocquin, P., and Périlleux, C. (2016) *FLOR-ID: an interactive database*
644  *of flowering-time gene networks in Arabidopsis thaliana*. *Nucleic Acids Res.*, **44**, D1167–71.

645  Burghardt, L.T., Young, N.D., and Tiffin, P. (2017) *A Guide to Genome-Wide Association*
646  *Mapping in Plants*. *Curr Protoc Plant Biol*, **2**, 22–38.

647  CABI (2013) *The CABI Encyclopedia of Forest Trees*. CABI.

648  Calleja-Rodriguez, A., Pan, J., Funda, T., Chen, Z., Baison, J., Isik, F., et al. (2020) *Evaluation*
649  *of the efficiency of genomic versus pedigree predictions for growth and wood quality traits*
650  *in Scots pine*. *BMC Genomics*, **21**, 796.

651  Cappa, E.P., de Lima, B.M., da Silva-Junior, O.B., Garcia, C.C., Mansfield, S.D., and
652  Grattapaglia, D. (2019) *Improving genomic prediction of growth and wood traits in*
653  *Eucalyptus using phenotypes from non-genotyped trees by single-step GBLUP*. *Plant Sci.*,
654  **284**, 9–15.

655  Cervantes, S., Vuosku, J., Paczesniak, D., and Pyhäjärvi, T. *Atlas of tissue-specific and tissue-*
656  *preferential gene expression in ecologically and economically significant conifer Pinus*
657  *sylvestris*.

658  Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J. (2015) *Second-*
659  *generation PLINK: rising to the challenge of larger and richer datasets*. *Gigascience*, **4**, 7.

660  Chen, Z.-Q., Baison, J., Pan, J., Karlsson, B., Andersson, B., Westin, J., et al. (2018) *Accuracy*
661  *of genomic selection for growth and wood quality traits in two control-pollinated progeny*
662  *trials using exome capture as the genotyping platform in Norway spruce*. *BMC Genomics*,
663  **19**, 946.

664  Chevreux, B. (2007) *MIRA: an automated genome and EST assembler*.

665  Clark, A.G., Hubisz, M.J., Bustamante, C.D., Williamson, S.H., and Nielsen, R. (2005)
666  *Ascertainment bias in studies of human genome-wide polymorphism*. *Genome Res.*, **15**,
667  1496–1502.

668  Darrier, B., Russell, J., Milner, S.G., Hedley, P.E., Shaw, P.D., Macaulay, M., et al. (2019) *A*

669         *Comparison of Mainstream Genotyping Platforms for the Evaluation and Use of Barley*
670         *Genetic Resources*. *Front. Plant Sci.*, **10**, 544.
671 Fierst, J.L. (2015) *Using linkage maps to correct and scaffold de novo genome assemblies:*
672         *methods, challenges, and computational tools*. *Front. Genet.*, **6**, 220.
673 González-Martínez, S.C., Wheeler, N.C., Ersoz, E., Nelson, C.D., and Neale, D.B. (2007).
674         *Association Genetics in Pinus taeda L. I. Wood Property Traits*. *Genetics* **175**, 399–409.
675 Graffelman, J. (2015) *Exploring Diallelic Genetic Markers: TheHardyWeinbergPackage*. *Journal*
676         *of Statistical Software*, **64**.
677 Grattapaglia, D., Silva-Junior, O.B., Resende, R.T., Cappa, E.P., Müller, B.S.F., Tan, B., et al.
678         (2018) *Quantitative Genetics and Genomics Converge to Accelerate Forest Tree Breeding*.
679         *Front. Plant Sci.*, **9**, 1693.
680 Grivet, D., Avia, K., Vaattovaara, A., Eckert, A.J., Neale, D.B., Savolainen, O., and González-
681         Martínez, S.C. (2017) *High rate of adaptive evolution in two widespread European pines*.
682         *Mol. Ecol.*, **26**, 6857–6870.
683 Haapanen M., Jansson G., Nielsen U.B., Steffenrem A., Stener L.G. (2015). *The status of tree*
684         *breeding and its potential for improving biomass production – a review of breeding activities*
685         *and genetic gains in Scandinavia and Finland.* Skogforsk, Uppsala. 56 p.
686         http://www.skogforsk.se/contentassets/9d9c6eeaef374a2283b2716edd8d552e/the-status-
687         of-tree-breeding-low.pdf.
688 Haapanen, M., Hynynen, J., Ruotsalainen, S., Siipilehto, J., and Kilpeläinen, M.-L. (2016)
689         *Realised and projected gains in growth, quality and simulated yield of genetically improved*
690         *Scots pine in southern Finland*. *European Journal of Forest Research*, **135**, 997–1009.
691 Houston Durrant, T., De Rigo, D., and Caudullo, G. (2016) *Pinus sylvestris in Europe:*
692         *distribution, habitat, usage and threats*. *European Atlas of Forest Tree Species.*
693         *Luxembourg: Publications Office of the European Union*, e016b94.
694 Huang, X. and Madan, A. (1999) *CAP3: A DNA sequence assembly program*. *Genome Res.*, **9**,
695         868–877.
696 Isik (2014) *Genomic selection in forest tree breeding: the concept and an outlook to the future.*
697         *New Forests,* 45, 379–401.
698 Jansson, G., Hansen, J.K., Haapanen, M., Kvaalen, H., and Steffenrem, A. (2017) *The genetic*
699         *and economic gains from forest tree breeding programmes in Scandinavia and Finland*.
700         *Scand. J. For. Res.*, **32**, 273–286.
701 Kujala, S.T. and Savolainen, O. (2012) *Sequence variation patterns along a latitudinal cline in*
702         *Scots pine (Pinus sylvestris): signs of clinal adaptation? Tree Genet. Genomes*, **8**, 1451–
703         1467.
704 Kujala, S.T., Knürr, T., Kärkkäinen, K., Neale, D.B., Sillanpää, M.J., and Savolainen, O. (2017)
705         *Genetic heterogeneity underlying variation in a locally adaptive clinal trait in Pinus*
706         *sylvestris revealed by a Bayesian multipopulation analysis*. *Heredity,* **118**, 413–423.
707 Lande, R. and Thompson, R. (1990) *Efficiency of marker-assisted selection in the improvement*
708         *of quantitative traits*. *Genetics*, **124**, 743–756.
709 Lenz, P.R.N., Beaulieu, J., Mansfield, S.D., Clément, S., Desponts, M., and Bousquet, J. (2017)
710         *Factors affecting the accuracy of genomic selection for growth and wood quality traits in an*
711         *advanced-breeding population of black spruce (Picea mariana)*. *BMC Genomics*, **18**.
712 Lee, S.J. (2002). Selection of parents for the Scots pine breeding population in Britain. Forestry:
713         An International Journal of Forest Research, **75**, 293–303.
714 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.,
715         Durbin, R., and 1000 Genome Project Data Processing Subgroup (2009). *The Sequence*
716         *Alignment/Map format and SAMtools. Bioinformatics.* **25**, 2078–2079.
717 Li, H. (2013) *Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM*.
718         *arXiv [q-bio.GN]*.
719 Lim, K.-J., Paasela, T., Harju, A., Venäläinen, M., Paulin, L., Auvinen, P., et al. (2016)

720    *Developmental Changes in Scots Pine Transcriptome during Heartwood Formation*. *Plant*
721    *Physiol.*, **172**, 1403–1417.
722  Manichaikul, A., Mychaleckyj, J.C., Rich, S.S., Daly, K., Sale, M., and Chen, W.-M. (2010)
723    *Robust relationship inference in genome-wide association studies*. *Bioinformatics*, **26**,
724    2867–2873.
725  Marrano, A., Martínez-García, P.J., Bianco, L., Sideli, G.M., Di Pierro, E.A., Leslie, C.A., et al.
726    (2019) *A new genomic tool for walnut (Juglans regia L.): development and validation of the*
727    *high-density Axiom^{TM} J. regia 700K SNP genotyping array*. *Plant Biotechnol. J.*, **17**, 1027–
728    1036.
729  Mason, W.L., and Alía, R. (2000). Current and future status of Scots Pine (Pinus sylvestris L.)
730    forests in Europe. Forest Systems 9(1), pp.317-335.
731  McTavish, E.J. and Hillis, D.M. (2015) *How do SNP ascertainment schemes and population*
732    *demographics affect inferences about population history? BMC Genomics*, **16**, 266.
733  Meuwissen, T., Hayes, B., and Goddard, M. (2016) *Genomic selection: A paradigm shift in*
734    *animal breeding*. *Anim Fron*, **6**, 6–14.
735  Meuwissen, T.H., Hayes, B.J., and Goddard, M.E. (2001) *Prediction of total genetic value using*
736    *genome-wide dense marker maps*. *Genetics*, **157**, 1819–1829.
737  Mullin, T.J., Andersson, B., Bastien, J.-C., Beaulieu, J., Burdon, R., Dvorak, W., King, J., Kondo,
738    T., Krakowski, J., Lee, S., et al. (2011). *Economic importance, breeding objectives and*
739    *achievements. Genetics, Genomics and Breeding of Conifers,* 40–127.
740  Neale, D.B., and Savolainen, O. (2004). *Association genetics of complex traits in conifers.*
741    *Trends in Plant Science.* **9**, 325–330.
742  Neale, D.B., Wegrzyn, J.L., Stevens, K.A., Zimin, A.V., Puiu, D., Crepeau, M.W., et al. (2014)
743    *Decoding the massive genome of loblolly pine using haploid DNA and novel assembly*
744    *strategies*. *Genome Biol.*, **15**, R59.
745  Nikkanen, T., Karvinen, K., Koski, V., Rusanen, M., and Yrjänä-Ketola, L. (1999) *Kuusen ja*
746    *männyn siemenviljelykset ja niiden käyttöalueet*. Metsäntutkimuslaitos.
747  Ojeda, D.I., Mattila, T.M., Ruttink, T., Kujala, S.T., Kärkkäinen, K., Verta, J.-P., and Pyhäjärvi, T.
748    (2019) *Utilization of Tissue Ploidy Level Variation in de Novo Transcriptome Assembly of*
749    *Pinus sylvestris*. *G3,* **9**, 3409–3421.
750  Oróstica, K.Y. and Verdugo, R.A. (2016) *chromPlot: visualization of genomic data in*
751    *chromosomal context*. *Bioinformatics*, **32**, 2366–2368.
752  Palmé, A.E., Wright, M., and Savolainen, O. (2008) *Patterns of divergence among conifer ESTs*
753    *and polymorphism in Pinus sylvestris identify putative selective sweeps*. *Mol. Biol. Evol.*,
754    **25**, 2567–2577.
755  Pan, Y., Birdsey, R.A., Fang, J., Houghton, R., Kauppi, P.E., Kurz, W.A., et al. (2011) *A large*
756    *and persistent carbon sink in the world's forests*. *Science*, **333**, 988–993.
757  Pavan, S., Delvento, C., Ricciardi, L., Lotti, C., Ciani, E., and D'Agostino, N. (2020)
758    *Recommendations for Choosing the Genotyping Method and Best Practices for Quality*
759    *Control in Crop Genome-Wide Association Studies*. *Front. Genet.*, **11**, 447.
760  Perry, A., Wachowiak, W., Downing, A., Talbot, R., and Cavers, S. (2020) *Development of a*
761    *SNP array for population genomic studies in four European pine species*. *Mol. Ecol.*
762    *Resour.*
763  Privé, F., Luu, K., Vilhjálmsson, B.J., and Blum, M.G.B. (2020) *Performing Highly Efficient*
764    *Genome Scans for Local Adaptation with R Package pcadapt Version 4*. *Mol. Biol. Evol.*,
765    **37**, 2153–2154.
766  Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A.R., Bender, D., et al. (2007)
767    *PLINK: a tool set for whole-genome association and population-based linkage analyses*.
768    *Am. J. Hum. Genet.*, **81**, 559–575.
769  Pyhäjärvi, T., García-Gil, M.R., Knürr, T., Mikkonen, M., Wachowiak, W., and Savolainen, O.
770    (2007) *Demographic history has influenced nucleotide diversity in European Pinus*

771       *sylvestris populations*. *Genetics*, **177**, 1713–1724.

772   Pyhäjärvi, T., Kujala, S.T., and Savolainen, O. (2020) *275 years of forestry meets genomics in*
773       *Pinus sylvestris*. *Evol. Appl.*, **13**, 11–30.

774   Resende, M.F.R., Muñoz, P., Resende, M.D.V., Garrick, D.J., Fernando, R.L., Davis, J.M.,
775       Jokela, E.J., Martin, T.A., Peter, G.F., and Kirst, M. (2012). *Accuracy of Genomic Selection*
776       *Methods in a Standard Data Set of Loblolly Pine (Pinus taeda L.). Genetics,* **190**, 1503–
777       1510.

778   Rosvall O, Jansson G, Andersson B, Ericsson T, Karlsson B, Sonesson J, Stener L (2001)
779       Genetiska vinster i nuvarande och framtida fröplantager och klonblandningar [Summary:
780       Genetic gains from present and future seed orchards and clone mixes]. Redogörelse nr 1,
781       Skogforsk: 1–41.

782   Rusanen, A., Lappalainen, K., Kärkkäinen, J., Tuuttila, T., Mikola, M., and Lassi, U. (2019)
783       *Selective hemicellulose hydrolysis of Scots pine sawdust*. *Biomass Conversion and*
784       *Biorefinery*, **9**, 283–291.

785   Sarvas, R. (1964): Havupuut. WSOY, Porvoo–Helsinki. 518 p. (In Finnish)

786   Saunders, I.W., Brohede, J., and Hannan, G.N. (2007) *Estimating genotyping error rates from*
787       *Mendelian errors in SNP array genotypes and their impact on inference*. *Genomics*, **90**,
788       291–296.

789   Silva-Junior, O.B., Faria, D.A., and Grattapaglia, D. (2015) *A flexible multi-species genome-*
790       *wide 60K SNP chip developed from pooled resequencing of 240 Eucalyptus tree genomes*
791       *across 12 species*. *New Phytol.*, **206**, 1527–1540.

792   Silva-Junior, O.B. and Grattapaglia, D. (2015) *Genome-wide patterns of recombination, linkage*
793       *disequilibrium and nucleotide diversity from pooled resequencing and single nucleotide*
794       *polymorphism genotyping unlock the evolutionary history of Eucalyptus grandis*. *New*
795       *Phytol.*, **208**, 830–845.

796   Stevens, K.A., Wegrzyn, J.L., Zimin, A., Puiu, D., Crepeau, M., Cardeno, C., et al. (2016)
797       *Sequence of the Sugar Pine Megagenome*. *Genetics*, **204**, 1613–1626.

798   Tajima, F. (1989) *Statistical method for testing the neutral mutation hypothesis by DNA*
799       *polymorphism*. *Genetics*, **123**, 585–595.

800   Tam, V., Patel, N., Turcotte, M., Bossé, Y., Paré, G., and Meyre, D. (2019) *Benefits and*
801       *limitations of genome-wide association studies*. *Nat. Rev. Genet.*, **20**, 467–484.

802   Tan, B., Grattapaglia, D., Martins, G.S., Ferreira, K.Z., Sundberg, B., and Ingvarsson, P.K.
803       (2017) *Evaluating the accuracy of genomic prediction of growth and wood traits in two*
804       *Eucalyptus species and their F1 hybrids*. *BMC Plant Biol.*, **17**, 110.

805   R Core Team (2018). R: A language and environment for statistical computing. R Foundation for
806       Statistical Computing, Vienna, Austria.

807   Thistlethwaite, F.R., El-Dien, O.G., Ratcliffe, B., Klápště, J., Porth, I., Chen, C., Stoehr, M.U.,
808       Ingvarsson, P.K., and El-Kassaby, Y.A. (2020). Linkage disequilibrium vs. pedigree:
809       Genomic selection prediction accuracy in conifer species. PLOS ONE, *15*, e0232201.

810   Tyrmi, J.S., Vuosku, J., Acosta, J.J., Li, Z., Sterck, L., Cervera, M.T., et al. (2020) *Genomics of*
811       *Clinal Local Adaptation in Pinus sylvestris Under Continuous Environmental and Spatial*
812       *Genetic Setting*. *G3: Genes|Genomes|Genetics*, **10**, 2683–2696.

813   Voss-Fels, K.P., Cooper, M., and Hayes, B.J. (2019) *Accelerating crop genetic gains with*
814       *genomic selection*. *Theor. Appl. Genet.*, **132**, 669–686.

815   Vuosku, J., Karppinen, K., Muilu-Mäkelä, R., Kusano, T., Sagor, G.H.M., Avia, K., et al. (2018)
816       *Scots pine aminopropyltransferases shed new light on evolution of the polyamine*
817       *biosynthesis pathway in seed plants*. *Ann. Bot.*, **121**, 1243–1256.

818   Vuosku, J., Muilu-Mäkelä, R., Avia, K., Suokas, M., Kestilä, J., Läärä, E., et al. (2019)
819       *Thermospermine Synthase (ACL5) and Diamine Oxidase (DAO) Expression Is Needed for*
820       *Zygotic Embryogenesis and Vascular Development in Scots Pine*. *Front. Plant Sci.*, **10**,
821       1600.

822 Wachowiak, W., Balk, P.A., and Savolainen, O. (2009) *Search for nucleotide diversity patterns*
823     *of local adaptation in dehydrins and other cold-related candidate genes in Scots pine*
824     *(Pinus sylvestris L.)*. *Tree Genet. Genomes*, **5**, 117.
825 Wachowiak, W., Trivedi, U., Perry, A., and Cavers, S. (2015) *Comparative transcriptomics of a*
826     *complex of four European pine species*. *BMC Genomics*, **16**, 234.
827 Wegrzyn, J.L., Lee, J.M., Tearse, B.R., and Neale, D.B. (2008) *TreeGenes: A forest tree*
828     *genome database*. *Int. J. Plant Genomics*, **2008**, 412875.
829 Westbrook, J.W., Chhatre, V.E., Wu, L.-S., Chamala, S., Neves, L.G., Muñoz, P., et al. (2015) *A*
830     *Consensus Genetic Map for Pinus taeda and Pinus elliottii and Extent of Linkage*
831     *Disequilibrium in Two Genotype-Phenotype Discovery Populations of Pinus taeda*. *G3,* **5**,
832     1685–1694.
833 Wray, N.R., Kemper, K.E., Hayes, B.J., Goddard, M.E., and Visscher, P.M. (2019) *Complex*
834     *Trait Prediction from Genome Data: Contrasting EBV in Livestock to PRS in Humans:*
835     *Genomic Prediction*. *Genetics*, **211**, 1131–1141.
836 Yang, J., Benyamin, B., McEvoy, B.P., Gordon, S., Henders, A.K., Nyholt, D.R., et al. (2010)
837     *Common SNPs explain a large proportion of the heritability for human height*. *Nat. Genet.*,
838     **42**, 565–569.
839 Yeaman, S., Hodgins, K.A., Lotterhos, K.E., Suren, H., Nadeau, S., Degner, J.C., Nurkowski,
840     K.A., Smets, P., Wang, T., Gray, L.K., et al. (2016). Convergent local adaptation to climate
841     in distantly related conifers. Science, **353**, 1431–1433.
842 Zimin, A., Stevens, K.A., Crepeau, M.W., Holtz-Morris, A., Koriabine, M., Marçais, G., et al.
843     (2014) *Sequencing and assembly of the 22-gb loblolly pine genome*. *Genetics*, **196**, 875–
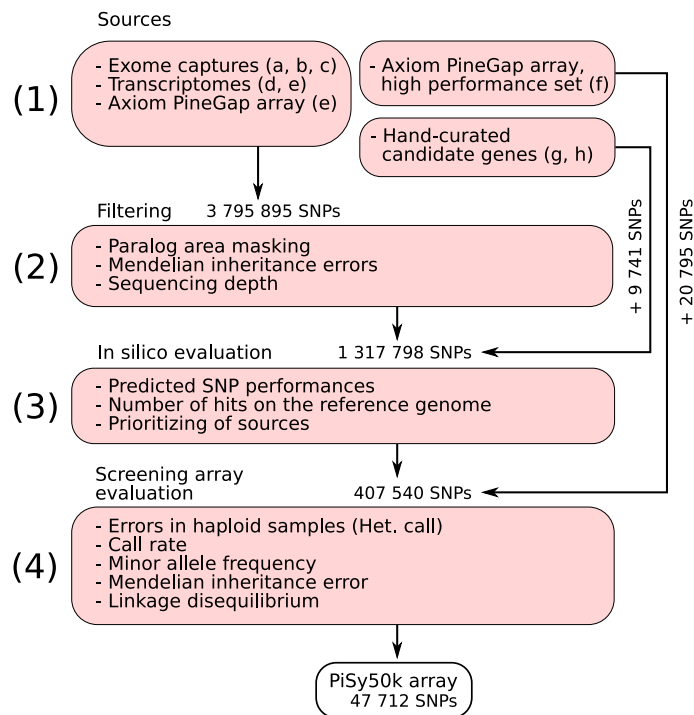844     890.

# Tables

845

846    **Table 1**. Sources of SNPs used in the design of PiSy50k array (M = megagametophyte, N =

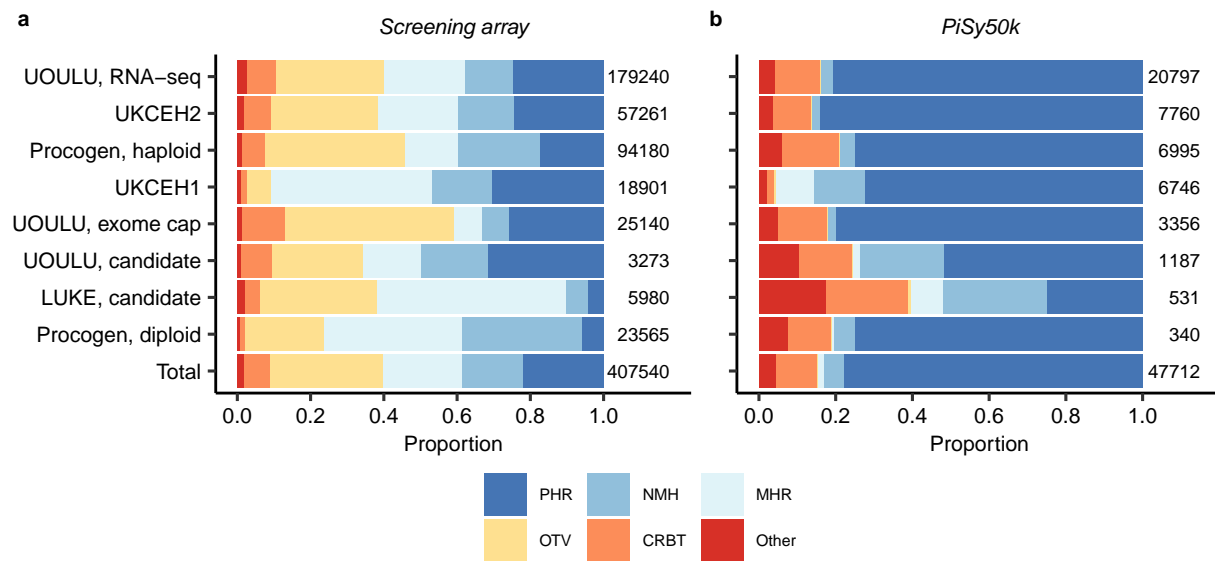847    needle, E = embryo; ISS Punkaharju = Intensive Study Site Punkaharju in southeast Finland).

| Data ID | Source tissue | Ascertainment size | Sampling area | DNA/RNA | Method | Reference |
|---|---|---|---|---|---|---|
| a. ProCoGen haploid | M | 109 haploids | Europe | DNA | Exome capture, Illumina | Tyrmi *et al.*, 2020 |
| b. ProCoGen diploid | N | 68 diploids | Europe | DNA | Exome capture, Illumina | Kastally et al. In prep. |
| c. UOULU exomeFEB2019 | NEM | 2 diploids | ISS Punkaharju | DNA | Exome capture, Illumina | Kesälahti et al. In prep. |
| d. UOULU RNA-seq | NEM | 18 lineages | ISS Punkaharju | RNA | Transcriptome | Ojeda *et al.*, 2019 |
| e. UKCEH1† | N | 17 diploids | Europe | RNA | SNP array Axiom_PineGAP (best set) | Perry *et al.*, 2020 |
| f. UKCEH2 | N | 17 diploids | Europe | RNA | SNP array Axiom_PineGAP; Transcriptomes of 4 Pine species | Perry *et al.*, 2020; Wachowiak *et al.*, 2015 |
| g. UOULU candidate† | M | 12-119 haploids | Europe | DNA | Sanger sequencing, Illumina sequencing | Avia *et al.*, 2014; Grivet *et al.*, 2017; Kujala & Savolainen 2012; Kujala *et al.*, 2017; Palmé *et al.*, 2008; Pyhäjärvi *et al.*, 2007;Vuosku *et al.*, 2018, 2019; Wachowiak *et al.*, 2009; Wegrzyn *et al.*, 2008, Evoltree EST database (http://www.evoltree.eu) |
| h. LUKE candidate† | M | 2-102 haploids | Europe | DNA | Sequence capture, Pacific Bioscience, Illumina | Kujala et al. In prep, Tyrmi *et al.*, 2020 |

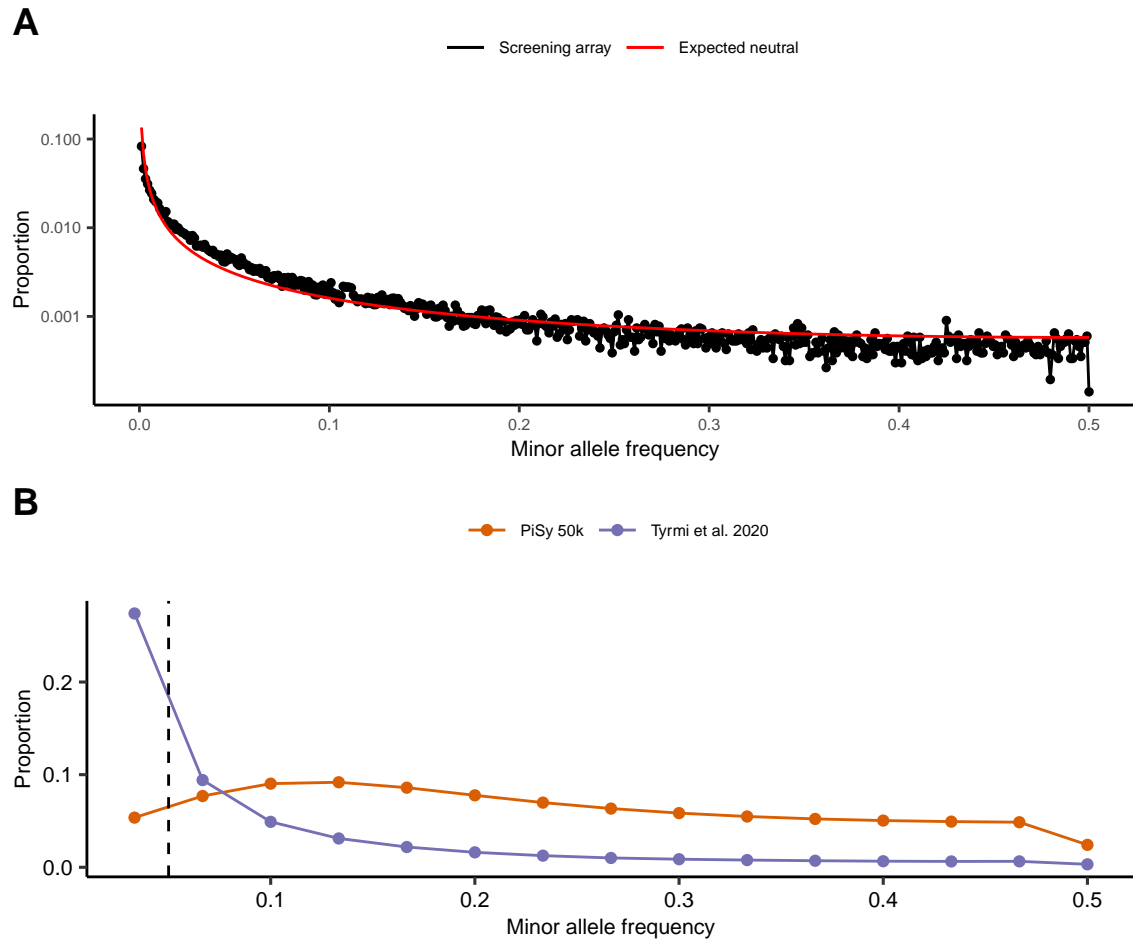848    †High priority sources, favored during the array design.
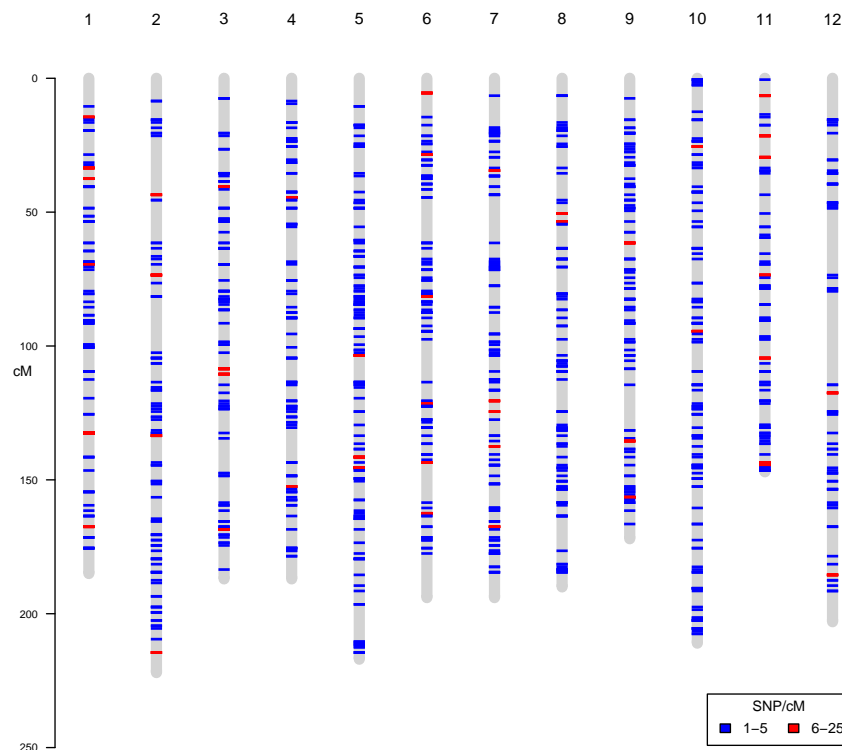
# Figures



**Figure 1.** Flow chart of the PiSy50k array design. We proceeded in four steps: (1) the collection of SNPs from 8 sources (Table 1; a: ProCoGen Haploid, b: ProCoGen Diploid, c: UOULU exomeFEB2019, d: UOULU RNA-seq, e: UKCEH2, f: UKCEH1, g: UOULU candidate and h: LUKE PacBio); (2) filtering to remove SNPs from paralogous genomic areas, SNPs with low sequencing depth or Mendelian errors; (3) evaluation to retain the best set of 407 540 markers (screening set) and (4) filtering based on the screening array performance to select the 47 712 markers retained in the PiSy50k array.
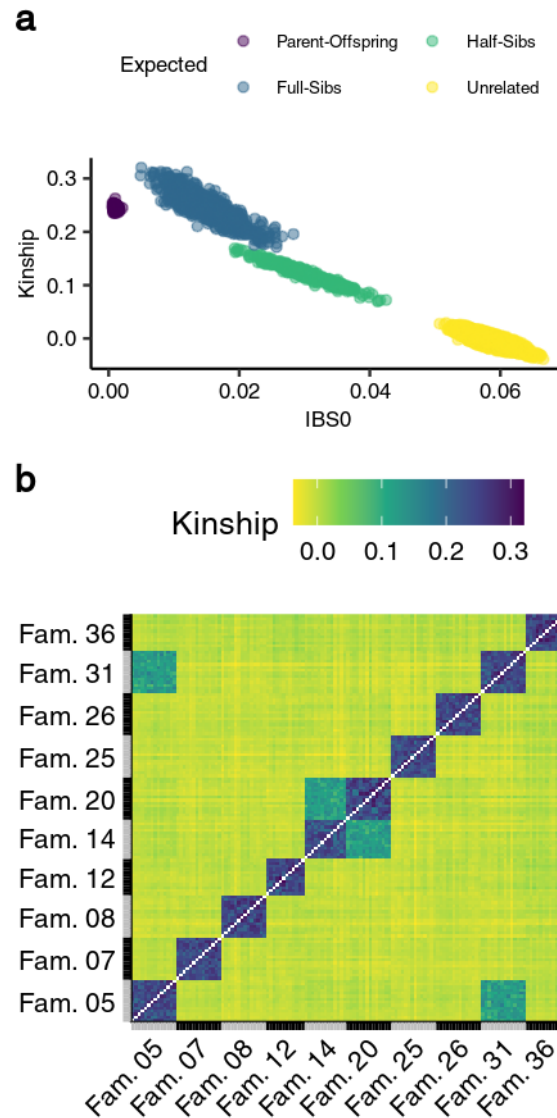
**Figure 2.** The proportions of conversion types of each marker source in (a) the screening array and (b) the PiSy50k array. PHR = Poly High Resolution, NMH = No Minor Homozygote, MHR = Mono High Resolution, CRBT = Call Rate Below Threshold, OTV = Off-Target Variant. Number right to the bar indicates the total number of SNPs per marker source.

**Figure 3.** Minor allele frequency (MAF) spectra of the screening and PiSy50k arrays. (A) MAF for the screening population sample (N = 466) and 56 693 SNPs (conversion type PHR and NMR) without missing data in the screening array. The red line illustrates the expected neutral MAF (Tajima, 1989). Note the log scale on the y-axis. (B) MAF based on the PiSy50k array including 38 302 SNPs genotyped in 90 plus-trees across three Finnish breeding populations (red line) and 42 exome captures of Scots pine trees sampled in four natural populations of Finland (Tyrmi et al., 2020). To be comparable, we downsampled both distributions to 30 samples. The vertical dashed line marks the filter threshold of 0.05 used during the array design and below which SNPs were partly excluded. As expected, there is a deficiency of rare alleles in the data obtained from the PiSy50k, as a result of ascertainment bias.
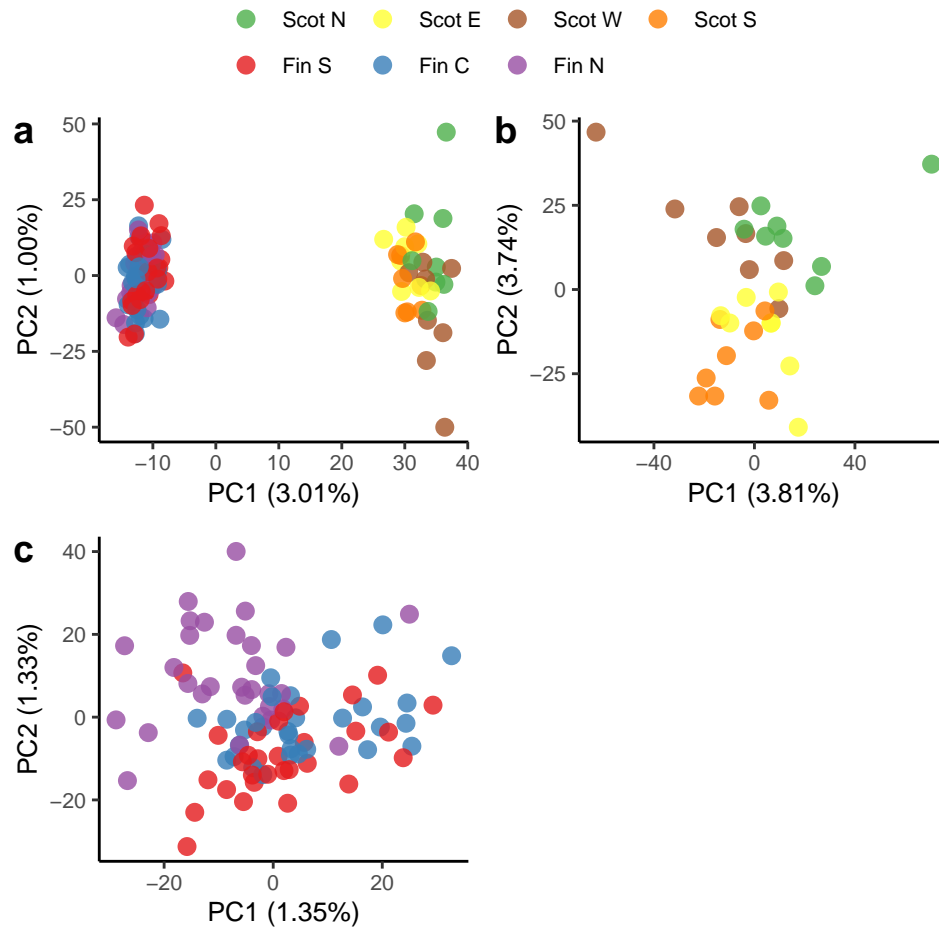
**Figure 4.** Position and density of 1 619 SNPs from the PiSy50k array on the *P. taeda* linkage map (Westbrook et al., 2015). The vertical grey lines represent the 12 linkage groups in *P. taeda*, while horizontal colored lines indicate the marker positions and density. This plot was made with the R package chromPlot (v 1.12.0) (Oróstica and Verdugo, 2016).

**Figure 5.** Relatedness analyses of 10 families (including 18 parents and 135 offspring) using the PiSy50k array. (a) Kinship coefficients (Manichaikul et al., 2010) and proportion of sites where individuals share no allele (IBS0) between all pairs and using 39 678 SNPs (PHR + NMH). Expected relationships between pairs are outlined: parent-offspring in purple, full sibs in blue, half sibs in green, and unrelated pairs in yellow. (b) Heat map of the kinship coefficients between all pairs of the 135 offspring.

**Figure 6.** Principal Component Analysis (PCA) using 39 678 polymorphic SNPs from the PiSy50k array genotyped in 120 trees from seven areas in Finland (90) and Scotland y (30). PCA including (a) all 120 samples from Finland and Scotland, (b) 30 samples collected across 21 localities grouped in four geographical areas of Scotland, or (c) 90 samples from Southern, Central and Northern Finland (30 samples each). Scot N, E, W and S: Northern, Eastern, Western and Southern Scotland. Fin S, C and N: Southern, Central and Northern Finland.