1 **TITLE: The somatic genetic and epigenetic mutation rate in a wild long-lived**

2 **perennial *Populus trichocarpa***

3

4 **AUTHORS**

5

6 Brigitte T. Hofmeister[1], Johanna Denkena[2], Maria Colomé-Tatché[2,3,4], Yadollah

7 Shahryary[5], Rashmi Hazarika[5,6], Jane Grimwood[7,8], Sujan Mamidi[7], Jerry Jenkins[7], Paul

8 P. Grabowski[7], Avinash Sreedasyam[7], Shengqiang Shu[8], Kerrie Barry[8], Kathleen Lail[8],

9 Catherine Adam[8], Anna Lipzen[8], Rotem Sorek[9], Dave Kudrna[10], Jayson Talag[10], Rod

10 Wing[10], David W. Hall[11], Gerald A. Tuskan[12], Jeremy Schmutz[7,8], Frank Johannes[5,6,*],

11 Robert J. Schmitz[6,11,*]

12

13 [1]Institute of Bioinformatics, University of Georgia, Athens, GA, USA

14 [2]Institute of Computational Biology, Helmholtz Center Munich, German Research

15 Center for Environmental Health, Neuherberg, Germany

16 [3]European Research Institute for the Biology of Ageing, University of Groningen,

17 University Medical Centre Groningen, Groningen, The Netherlands

18 [4]TUM School of Life Sciences Weihenstephan, Technical University of Munich, Freising,

19 Germany

20 [5]Department of Plant Sciences, Technical University of Munich, Liesel-Beckmann-Str. 2,

21 Freising, Germany

22 [6]Institute for Advanced Study (IAS), Technical University of Munich, Lichtenbergstr. 2a,

23 Garching, Germany

24 7HudsonAlpha Institute of Biotechnology, Huntsville, Alabama, USA

25 8Department of Energy Joint Genome Institute, Walnut Creek, California, USA

26 9Department of Molecular Biology, Weizmann Institute of Science, Rehovot, Israel

27 10Arizona Genomics Institute, School of Plant Sciences, University of Arizona, Tucson,

28 AZ, USA

29 11Department of Genetics, University of Georgia, Athens, GA, USA

30 12The Center for Bioenergy Innovation, Oak Ridge National Laboratory, Oak Ridge, TN,

31 USA

32

33 *CORRESPONDING AUTHORS: Robert J. Schmitz, schmitz@uga.edu and Frank

34 Johannes, frank@johanneslab.org

35

37

38 **ABSTRACT**

39

40 **Background:** Plants can transmit somatic mutations and epimutations to offspring,

41 which in turn can affect fitness. Knowledge of the rate at which these variations arise is

42 necessary to understand how plant development contributes to local adaption in an eco-

43 evolutionary context, particularly in long-lived perennials.

44 **Results:** Here, we generated a new high-quality reference genome from the oldest

45 branch of a wild *Populus trichocarpa* tree with two dominant stems which have been

46 evolving independently for 330 years. By sampling multiple, age-estimated branches of

47   this tree, we used a multi-omics approach to quantify age-related somatic changes at

48   the genetic, epigenetic and transcriptional level. We show that the per-year somatic

49   mutation and epimutation rates are lower than in annuals and that transcriptional

50   variation is mainly independent of age divergence and cytosine methylation.

51   Furthermore, a detailed analysis of the somatic epimutation spectrum indicates that

52   transgenerationally heritable epimutations originate mainly from DNA methylation

53   maintenance errors during mitotic rather than during meiotic cell divisions.

54   **Conclusion:** Taken together, our study provides unprecedented insights into the origin

55   of nucleotide and functional variation in a long-lived perennial plant.

56

57   **BACKGROUND**

58

59   The significance of somatic mutations, i.e., variations in DNA sequence that occur after

60   fertilization, in long-lived plant and animal species have been a point of debate and

61   investigation for the past 30 years [1–4]. It has been hypothesized that the evolutionary

62   consequences of such mutations are likely even more profound in woody perennial

63   plants, where undifferentiated meristematic cells produce all above-ground and below-

64   ground structures. As meristems undergo constant cell division throughout the lifetime

65   of a plant, somatic mutations arising in meristems may result in genetic differences

66   being passed onto progeny cells [5–8]. The accumulation of somatic mutations can thus

67   lead to genetic and occasionally also phenotypic divergence among vegetative lineages

68   within the same individual. In trees, for instance, different branches have been shown to

69   differ in their responses to pest and pathogen attack, alternate reactions to drought

70    and/or nutrient availability, or dissimilar demands for photosynthate material, even

71    within the same individual [9]. Beyond the impact of point mutations and small

72    insertions/deletions on gene function, alterations in chromatin structure and DNA

73    methylation might also impact gene expression variation.

74

75    Phenotypic variation has been attributed to somatic mutations in several perennial

76    plants, including the derivation of Nectarines in peach [10] and the origin of modern

77    grape cultivars (*Vitis vinifera* L.) [11]. In *Populus tremuloides,* somatic mutations have

78    been hypothesized as the cause for variation in DNA markers among individual ramets

79    of a single genotype [12]. Initial attempts to demonstrate within-tree mosaicism using

80    genetic markers [13], showed at low-resolution that the degree of intra-tree variability

81    was positively correlated with the physical distance between sampled branches. More

82    recently, work in oak (*Quercus rubur*) has documented variation in DNA sequence

83    among an independent sampling of alternate branches from a single genotype [14, 15].

84    They estimated a fixed mutation rate of 4.2 - 5.2 x $10_{-8}$ substitutions per locus per

85    generation, which is only within one order of magnitude of the rate observed in the

86    herbaceous annual plant *Arabidopsis thaliana* [16]. These results are consistent with an

87    emerging hypothesis that the per-unit-time mutation rate of perennials is much lower

88    than in annuals to delay mutational meltdown [17, 18] and this lower rate is

89    accomplished by limiting the number of cell divisions between the meristem and the

90    new branch [19]. Additional recent studies have also revealed similar rates of

91    spontaneous mutations in a range of species including perennials [18]. Regardless of

92    the rate of mutation, the frequency of deleterious mutations in woody plants is high,

93    which is hypothesized to reduce survival of progeny resulting from inbreeding and favor

94    outcrossing as is observed in many forest trees [20, 21].

95

96    Similar to genetic mutations, phenotypic variation can be caused by epigenetic variation

97    such as stable changes in cytosine methylation or epimutations [22]. Cytosine

98    methylation is a covalent base modification that is inherited through both mitotic and

99    meiotic cell divisions in plants [23]. It occurs in three sequence contexts, CG, CHG, and

100   CHH (H = A, T, or C) and the pattern and distribution of methylation at these different

101   contexts is predictive of its function in genome regulation [24]. Spontaneous changes in

102   methylation independent of genetic changes can lead to phenotypic changes [25]. Well-

103   characterized examples in plants include the peloric phenotype in toadflax (*Linaria*

104   *vulgaris*), the colorless non-ripening phenotype in tomato (*Solanum lycopersicum*), and

105   the mantled phenotype in oil palm (*Elaeis guineensis*) [26–28].

106

107   Once established, epimutations can stably persist or be inherited across generations.

108   For example, the reversion rate from the colorless non-ripening epimutant allele to wild

109   type is about 1 in 1000 per generation in tomato [27]. Studies in *A. thaliana* mutation

110   accumulation lines have documented that the vast majority (91-99.998%) of methylated

111   regions in the genome are stably inherited across generations; only a small subset of

112   the methylome shows variation among mutation accumulation lines [29–31]. Estimates

113   in *A. thaliana* indicate that the spontaneous methylation gain and loss rates at CG sites

114   are $2.56 \times 10^{-4}$ and $6.30 \times 10^{-4}$ per generation per haploid methylome, respectively [32].

115   Despite the wealth of knowledge about transgenerational methylation inheritance, very

116    little is known about somatic epimutations, especially in long-lived perennial species.

117    Previous studies have been limited by resolution and time. Heer *et al.* observed no

118    global methylation changes and no consistent variation in gene body methylation

119    associated with growth conditions of Norway spruce [33]. Several studies have linked

120    stress conditions to differential methylation in perennials but did not look at the stability

121    of methylation after removing the stressor [34, 35]. One exception, Le Gac *et al.*,

122    identified environment-related differentially methylated regions in poplar, but only

123    examined stability across six months [36].

124

125    Detailed insights into the rate and spectrum of somatic mutations and epimutations are

126    necessary to understand how somatic development of long-lived perennials contribute

127    to population-level variation in an eco-evolutionary context. Here we generated a new

128    high-quality reference genome from the oldest branch of a wild *Populus trichocarpa* tree

129    with two dominant stems which have been evolving independently for approximately

130    330 years. By sampling multiple, age-estimated branches of this tree, we used a multi-

131    omics approach to quantify age-related somatic changes at the genetic, epigenetic and

132    transcriptional level. Our study provides the first quantitative insights into how nucleotide

133    and functional variation arise during the lifetime of a long-lived perennial plant.

134

135    **RESULTS**

136

137    **Experimental design for the discovery of somatic genetic and epigenetic variants**

138

139    A stand of trees was identified near Mount Hood, Oregon and vegetative samples were

140    collected from over 15 trees as part of an independent study. Of these trees, five were

141    chosen for subsequent analysis and five branches of each tree were identified (Fig. S1).

142    For each branch, the stem age was determined by coring the main stem at breast

143    height and where the branch meets the stem and the branch age was determined by

144    coring the base of the branch (Fig.1 and Fig. S2). Although 25 branches in total were

145    initially sampled, six were excluded from analysis because they were epicormic and age

146    estimates could not be determined. Two other branches had incomplete cores, but ages

147    could be estimated based on radial diameter.

148

149    From this, we were specifically interested in tree 13 and tree 14 (Fig. 1). Originally

150    identified as two separate genotypes, they are actually two main stems of a single basal

151    root system and trunk. Both tree 13 and tree 14 originated as stump sprouts off of an

152    older tree that was knocked down over 300 years ago. Attempts to determine the total

153    age were unsuccessful. However, statistical estimates based on molecular-clock

154    arguments and a regression analysis of diameter to age suggest that the tree is

155    approximately 330 years old (Shayary et al. 2019, co-submission).

156

157    Leaf samples were collected from eight age-estimated branches for multi-omics

158    analysis for tree 13 and tree 14. The oldest branch of tree 14 (branch 14.5) was used

159    for genome assembly of *Populus trichocarpa* var. *Stettler*. Genome resequencing was

160    performed for all branches to explore intra- and inter-tree genetic variation. PacBio,

7

161    MethylC-seq, and mRNA-seq libraries were constructed for the branches of tree 13 and

162    tree 14 to explore structural, methylation, and transcriptional variation, respectively.

163

164    **Genome assembly and annotation of *Populus trichocarpa* var. *Stettler***

165    We sequenced the *P. trichocarpa* var. *Stettler* using a whole-genome shotgun

166    sequencing strategy and standard sequencing protocols. Sequencing reads were

167    collected using Illumina and PacBio. The current release is based on PacBio reads

168    (average read length of 10,477 bp, average depth of 118.58x) assembled using the

169    MECAT CANU v.1.4 assembler [37] and subsequently polished using QUIVER [38]. A

170    set of 64,840 unique, non-repetitive, non-overlapping 1.0 kb sequences were identified

171    in the version 4.0 *P. trichocarpa* var. *Nisqually* assembly and were used to assemble

172    the chromosomes. The version 1 *Stettler* release contains 392.3 Mb of sequence with a

173    contig N50 of 7.5 Mb and 99.8% of the assembled sequence captured in the

174    chromosomes. Additionally, ~232.2 Mb of alternative haplotypes were identified.

175    Completeness of the final assembly was assessed using 35,172 annotated genes from

176    the version 4.0 *P. trichocarpa* var. *Nisqually* release (jgi.doe.gov). A total of 34,327

177    (97.72%) aligned to the primary *Stettler* assembly.

178    The annotation was performed using ~1.4 billion pairs of 2x150 stranded paired-end

179    Illumina RNA-seq GeneAtlas *P. trichocarpa* var. *Nisqually* reads, ~1.2 billion pairs of

180    2x100 paired-end Illumina RNA-seq *P. trichocarpa* var. *Nisqually* reads from Dr. Pankaj

181    Jaiswal, and ~430 million pairs of 2x75 stranded paired-end Illumina var. *Stettler* reads

182    using PERTRAN (Shu, unpublished) on the *P. trichocarp*a var. *Stettler* genome. About

183    ~3 million PacBio Iso-Seq circular consensus sequences were corrected and collapsed

184    by a genome-guided correction pipeline (Shu, unpublished) on the *P. trichocarpa* var.

185    *Stettler* genome to obtain ~0.5 million putative full-length transcripts. We annotated

186    34,700 protein-coding genes and 17,314 alternative splices for the final annotation.

187    Because of the extensive resources included in the annotation, 32,330 genes had full-

188    length transcript support.

189    **Identification and rate of somatic genetic variants**

190    Leaf samples from the five trees were sequenced to an average depth of ~87x (~60-

191    164x) using Illumina HiSeq. Roughly 88% of the high-quality reads map to the genome

192    and about 98.6% of the genome is covered by at least one read, and genome coverage

193    (~8-500x) used for SNP calling was about 97%. The initial number of SNPs per tree

194    (mutation on any branch) varied between 44,000 and 152,000, which is populated with

195    many false positives due to coverage, sequencing and alignment errors, etc. Applying

196    an additional filter requiring >20x coverage per position and requiring coverage in all

197    branches reduced the total amount genome space queried to ~40 Mb. Furthermore,

198    since most of the genome (99.9%) is homozygous at every base pair, a somatic

199    mutation will almost always result in a change from a homozygous to heterozygous site.

200    Restricting the analysis to sites that change from homozygous to heterozygous, we

201    identified 118 high-confidence SNPs in tree 13 and 143 high-confidence SNPs in tree

202    14 (Tables S1-2).

203    Over two-thirds of the SNPs in tree 13 and tree 14 were transition mutations, with C-G

204    to T-A mutations accounting for over 54% of the SNPs (Fig. 2a). Of the transversion

205    mutations C-G to G-C was the least common (3.8%) whereas C-G to A-T was most

206    common (10%). Nearly half of the SNPs (46%) occurred in transposable elements and

207    about 10% occur in promoter regions (Fig. 2b and Tables S1-S2). SNPs are significantly

208    enriched in TEs and depleted in promoter regions genome-wide (Chi-square, df = 3, $P <$

209    0.001)

210    To obtain an estimate of the rate of somatic point mutations from these SNP calls, we

211    developed *mutSOMA* (https://github.com/jlab-code/mutSOMA), a phylogeny-based

212    inference method that fully incorporates knowledge of the age-dated branching topology

213    of the tree (see Methods and Supplementary Text). Using this approach, we find that

214    the somatic point mutation rate in poplar is $1.33 \times 10^{-10}$ (95% CI: $1.53 \times 10^{-11}$ - $4.18 \times 10^{-10}$

215    $^{-10}$) per base per haploid genome per year (Table S3). Generation time can refer to two

216    measurements—time from seed to production of first seeds and the organism's lifespan.

217    In annual plants, these values can be considered the same; however, this is not the

218    case for perennials. Assuming 15 years from seed to production of first seeds [39], the

219    poplar seed-to-seed generation mutation rate would be approximately $1.99 \times 10^{-9}$. This

220    is slightly lower than the per-generation (seed-to-seed) mutation rate observed in the

221    annual *A. thaliana* ($7 \times 10^{-9}$) [16]. Next looking at the lifespan per-generation rate and

222    assuming a maximum age of 200 years [40], the lifespan per-generation rate is $2.66 \times$

223    $10^{-8}$. This estimate is slightly lower than the per-generation somatic mutation rate

224    recently reported in oak ($4.2 - 5.8 \times 10^{-8}$) [14].

225

226    To analyze structural variants (SV) between haplotypes and somatic SV mutations,

227    PacBio libraries were generated for the eight branches from tree 13 and tree 14 (Fig. 1).

228    For each branch, four PacBio cells were sequenced generating an average output of

229     3.05 million reads and 28.3 Gb per branch (Table S4). After aligning the PacBio output

230     to the *P. trichocarpa* var. *Stettler* genome, calling SVs larger than 20 bp, and filtering,

231     we identified ~10,466 deletions, ~6,702 insertions, 645 duplications, and three

232     inversions between the reference *Stettler* haplotype and the alternative haplotype

233     (Table S5). Upon manual inspection of read mapping for a representative subset of

234     SVs, 72.6% of SVs have strong support where multiple aligned reads support the SV

235     type and size (Table S6). Deletions and duplications are significantly enriched in

236     tandem repeat sequence and depleted in genic sequence (Kolmogorov-Smirnov two-

237     sample test, *P* value < 2.2 x $10_{-16}$). Furthermore, deletions generally have less genic

238     sequence and more tandem repeat sequence than do duplications (Fig. S3). Several of

239     the detected SVs are large, with 11 deletions and five duplications greater than 50 kb

240     (Table S5) with genic sequence content ranging from 0.0% to 23.7%. Comparisons of

241     the branches from tree 13 and tree 14 did not identify instances of somatic SV mutation.

242

243     **Identification and rate of somatic epigenetic variants**

244

245     To explore somatic epigenetic variation associated with changes in DNA methylation,

246     we generated whole-genome bisulfite sequencing libraries from the branch tips of tree

247     13 and tree 14 (Fig. 1). The average genome coverage for the samples was ~41.1x and

248     sequence summary statistics are located in Table S7. Genome-wide methylation levels

249     were similar across all samples with 36.61% mCG, 19.02% mCHG, and 2.07% mCHH%

250     (Fig. S4) [41], indicating that global methylation levels are relatively stable across

251     branches. Nonetheless, we observed significant age-dependent DNA methylation

252    divergence between branches in CG and CHG contexts, both at the level of individual

253    cytosines as well as at the level of regions, i.e. clusters of cytosines (Fig. 3a-b, Fig. S5,

254    and Table S8). These age-dependent divergence patterns indicate that spontaneous

255    methylation changes (i.e. epimutations) are cumulative across somatic development

256    and thus point to a shared meristematic origin (Shahryary et al. 2019, co-submission).

257

258    To obtain an estimate of somatic epimutation rates, we applied *AlphaBeta* (Shahryary et

259    al. 2019, co-submission). The method builds on our previous approach for estimating

260    'germline'-epimutation in mutation accumulation (MA) lines [32], except here we treat

261    the tree branching topology as an intra-organismal phylogeny and model mitotic instead

262    of meiotic inheritance. Focusing first on cytosine-level epimutations, we estimated that

263    at the genome-wide scale spontaneous methylation gains in contexts CG and CHG

264    occur at a rate of $1.8 \times 10^{-6}$ and $3.3 \times 10^{-7}$ per site per haploid genome per year,

265    respectively; whereas spontaneous methylation losses in these two sequence contexts

266    occur at a rate of $5.8 \times 10^{-6}$ and $4.1 \times 10^{-6}$ per site per haploid genome per year. Based

267    on these estimates, we extrapolate that the *seed-to-seed* per-generation epimutation

268    rate in poplar is about 10-5 and the *lifespan* per-generation rate is 10-4. Remarkably,

269    these estimates are very similar to those reported in *A. thaliana* MA lines [32]. The

270    observation that two species with such different life history traits and genome

271    architecture display very similar per-generation mutation and epimutation rates

272    suggests that the rates themselves are subject to strong evolutionary constraints.

273

274    In addition to global epimutation rates, we also estimated rates for different genomic

275    features (mRNA, promoters, intergenic, TEs). This analysis revealed highly significant

276    rate differences in the CG and CHG context between genomic features, with mRNAs

277    showing the highest and TEs the lowest combined rates (Fig. 3c-j). Interestingly, the

278    ordering of the magnitude of the mRNA, promoter, and intergenic rates is similar to that

279    previously observed in *A. thaliana* MA lines [32]. The differences in rates at local

280    genomic features likely reflect the distinct DNA methylation pathways that function on

281    these sequences (RNA-directed DNA methylation, CHROMOMETHYLASE3,

282    CHROMOMETHYLASE2, DNA METHYLTRANSFERASE1, etc.). For example, the high

283    rate of epimutation losses in mRNA relative to other features (Fig. 3g-h) could reflect the

284    activity of CMT3-mediated gene body DNA methylation [42, 43]. The observation that

285    the epimutation rates of these features is consistent between *A. thaliana* MA lines (>60

286    generations) and this long-lived perennial (within a single generation) seems to imply

287    that epimutations are not a result of biased reinforcement of DNA methylation during

288    sexual reproduction or environment/genetic variation, but instead a feature of DNA

289    methylation maintenance through mitotic cell divisions.

290

291    **Assessment of spontaneous differentially methylated regions**

292

293    Differentially methylated regions are functionally more relevant than individual cytosine-

294    level changes, as in certain cases they are linked to differential gene expression and

295    phenotypic variation [26–28, 44, 45]. To explore the extent of differentially methylated

296    regions (DMRs) that spontaneously arise in these trees we searched for all pairwise

297   DMRs between all branches. In total, we identified 10,909 DMRs that possessed

298   changes in all sequence contexts (CG, CHG and CHH - C-DMRs). Together they

299   constitute approximately 1.69 Mb of the total 167.4 Mb (~1%) of methylated sequences

300   in the *Stettler* genome and they reveal age-dependent accumulation (Fig. 4a). Most

301   DMRs occur in intergenic regions (56.7%), but a significant enrichment of DMRs were

302   detected within two kilobases from the transcriptional start site of genes compared to

303   methylated regions as a whole (Fig. 4b) (Fisher's exact test, one-sided, $P$ value <

304   0.001).

305   Given the heterozygous nature of wild *P. trichocarpa*, we explored allelic methylation

306   changes. After filtering for sufficient coverage and methylation change, we assigned the

307   pseudo-allele state of each branch at 4,488 DMRs. Possible states were homozygous

308   unmethylated, heterozygous, and homozygous methylated. In each sample, 43.0% of

309   DMRs, on average, were categorized as homozygous methylated (Fig. S6).

310   Interestingly, the youngest branches, 13.1 and 14.1 have about 10% more homozygous

311   methylated pseudo-alleles than the other branches (51.1% vs 41.7%). Next, we looked

312   at the number of changes of pseudo-allele states. This is expected as DMRs were

313   identified as having different methylation levels in the samples. On average, there are

314   3.02 state changes for each DMR with 94.4% of DMRs having one to five state changes

315   (Fig. 4c). These data suggest that many of these regions are metastable, a common

316   feature of epimutations in plants.

317

318   An example of a region with one state change are the tree specific DMRs (Fig. 4d). In

319   these regions, all branches of one tree are homozygous unmethylated and all branches

320 of the other tree are homozygous methylated. This suggesting the methylation state

321 change occurred shortly after the trees separated and remained stable throughout

322 subsequent mitotic divisions. In contrast, we also identified highly variable regions with

323 seven state changes, a change between each branch (Fig. 4e). Of the regions with two

324 state changes, 150 have branch-specific state changes. For example, in Fig. 4f

325 branches 13.1 to 13.3 are homozygous unmethylated, then it changes to homozygous

326 methylated for branch 13.5, and changes again to homozygous unmethylated for

327 branches 14.5 – 14.2. Similarly, in Fig. 4g, all branches except 14.5 are homozygous

328 methylated and 14.5 has spontaneously lost methylation.

329

330 We also used the identified C-DMRs (differential methylation in all cytosine sequence

331 contexts) to obtain region-level epimutation rates. To do this, we established control

332 regions ('non-DMR') with the same size distribution as observed for C-DMRs and used

333 the methylation levels of all cytosines in each (non-)DMR to calculate methylation levels

334 per region. Interestingly, this analysis shows that region-level epimutation rates are

335 comparable to epimutation rates of single cytosines. Even though there are far fewer

336 DMRs in comparison to epimutations at single cytosines, the similar rates are not too

337 unexpected if one considers that the total 'epimutable space' for regions in the genome

338 is much smaller than that for individual cytosines. In summary, these results might

339 suggest that the mechanisms which underlie spontaneous differential methylation are

340 the same for differential methylation in larger regions and at individual sites.

341

342 **Functional consequences of differential methylation on gene expression**

343

344  To assess if age-related cytosine methylation changes have functional consequences,

345  we performed mRNA-seq with three biological replicates for each branch of trees 13

346  and 14. On average, each library had over ~55 million reads and 96.8% mapping to the

347  *P. trichocarpa* var. *Stettler* genome (Table S9). We used DESeq2 to identify

348  differentially expressed genes (DEGs) pairwise between branches [46] and identified a

349  total of 2,937 genes. The *P. trichocarpa* var. *Stettler* genome has 34,700 annotated

350  genes, so this differential expression gene set is 8.46% of all genes and 10.5% of

351  expressed genes.

352

353  Since the somatic accumulation of spontaneous methylation changes could affect gene

354  expression, we asked if transcriptional divergence also increases as a function of tree

355  age. We found that in contrast to somatic mutations and epimutations, the divergence

356  between leaf transcriptomes is much more heterogeneous and displays only a weak

357  and non-significant accumulation trend (Fig. 5a). This observation suggests that the

358  accumulation of genetic and epigenetic changes are largely uncoupled from age-

359  dependent transcriptional changes in poplar, at least at the global scale.

360

361  However, this global analysis does not rule out that DNA methylation changes at

362  specific individual loci can have transcriptional consequences. To explore this in more

363  detail, we analyzed DMRs proximal to DEGs, and correlated the methylation level of the

364  DMR with the expression level of the gene. The correlation is positive when a higher

365  methylation level in the DMR is associated with higher expression of the gene.

366 Regardless of where the DMR was located relative to the gene, we observed positive

367 DMR-DEG correlations and negative DMR-DEG correlations. There was no bias for

368 direction of correlation and genomic feature type (Fig. 5b).

369

370 We further focused on four specific examples where DEG-DMR correlations were

371 statistically significant (Fig. S7). Of these four, three of the DMRs occurred within two

372 kilobases upstream of the transcription start site, and they have strong negative

373 correlations (Fig. 5c). The DMR located in the untranslated region of a gene encoding a

374 mitochondrial oxoglutarate/malate carrier protein was positively correlated with gene

375 expression (Fig. 5d), although it remains unclear if this relationship is causal.

376

377 Taken together, our transcriptome analysis indicates that gene expression changes in

378 this poplar tree are largely independent of methylation at both the global and local scale

379 except for a few rare examples. This observation is at least partly consistent with our

380 model-based analyses, which suggest that somatic epimutations in this tree accumulate

381 neutrally (Shahryary et al. 2019, co-submission).

382

383 **DISCUSSION**

384

385 Using a multi-omics approach, we were able to calculate the rates of somatic mutations

386 and epimutations in the long-lived perennial tree *P. trichocarpa*. Consistent with the per-

387 unit-time hypothesis, we find that the per-year genetic and epigenetic mutation rates in

388 poplar are lower than in A. thaliana, which is remarkable considering that the former

389    experienced hundreds of years of variable environmental conditions. This observation

390    supports the view that long-lived perennials may limit the number of meristematic cell

391    divisions during their lifetime and that they have evolved mechanisms to protect these

392    cell types from the persistent influence of environmental mutagens, such as UV-

393    radiation. Interestingly, in contrast to the observed differences in *per-year* mutation and

394    epimutation rates, our analysis reveals strong similarities in the *per-generation* rates

395    between these two species. This close similarity further suggests that the per-

396    generation rates of spontaneous genetic and epigenetic changes are under strong

397    evolution constraint, although it remains unclear from our experimental design how

398    many of these (epi)mutations will be successfully transferred to the next generation.

399

400    The results presented here are most certainly an underestimate of the actual rate. This

401    may be a result of the sampling biased used in this study, as we were only able to

402    sample surviving branches and identify mutations that occurred early enough that they

403    are present in the majority of the cells sampled in the tissues profiled. Perhaps variable

404    environmental conditions lower the epimutation rate by keeping the cells in sync, thus

405    few differences can be observed. Alternatively, meristematic cells that give rise to the

406    sampled tissues have highly reinforced and well-maintained DNA methylomes similar to

407    observations in embryonic tissue [47–51]. Either scenario would imply that most of the

408    identified epimutations are spontaneous in nature. Although the rate is different, the

409    ordering in feature-specific epimutation rates is the same between poplar and *A.*

410    *thaliana*, suggesting that this is a general pattern in plant genomes, which likely is

411    derived from maintenance of DNA methylation through mitotic cell divisions.

412

**CONCLUSION**

413

414

415     Taken together, our study provides unprecedented insights into the origin of nucleotide,

416     epigenetic, and functional variation in the long-lived perennial plant.

417

418

**METHODS**

419

420

**Sample collection and age estimation**

421

422

423     The trees used in this study were located at Hood River Ranger District [Horse Thief

424     Meadows area], Mt. Hood National Forest, 0.6 mi south of Nottingham Campground off

425     OR-35 at unmarked parking area, 500' west of East Fork Trail #650 across river, ca.

426     45.355313, -121.574284 (Fig. S1).

427

428     Cores were originally collected from the main stem and five branches from each of five

429     trees in April 2015 at breast height (~1.5 m) for standing tree age using a stainless-steel

430     increment borer (5 mm in diameter and up to 28 cm in length). Cores were mounted on

431     grooved wood trim, dried at room temperature, sanded and stained with 1%

432     phloroglucinol following the manufacturer's instructions (https://www.forestry-

433     suppliers.com/Documents/1568_msds.pdf). Annual growth rings were counted to

434     estimate age. For cores for which accurate estimates could not be made from the 2015

435 collection, additional collections were made in spring 2016. However, due to difficulty in

436 collecting by climbing, many of the cores did not reach the center of the stem or

437 branches (pith) and/or the samples suffered from heart rot. Combined with the difficulty

438 in demarcating rings in porous woods such as poplar *Populus* [52, 53], accurate

439 measures of tree age or branch age were challenging (Fig. S2).

440

441 Simultaneously with stem coring, leaf samples were collected from the tips of each of

442 the branches from the selected five trees. Branches 9.1, 9.5, 13.4, 14.1, 15.1, and 15.5

443 were too damaged to determine reasonable age estimates and were removed from

444 analysis. Branch 14.4 and the stems of 13.1 and 13.2 were estimated by simply

445 regressing the diameter of all branches and stems that could be aged by coring.

446

447 **Nuclei prep for DNA extraction**

448

449 Poplar leaves, that had been kept frozen at -80 °C, were gently ground with liquid

450 nitrogen and incubated with NIB buffer (10 mM Tris-HCL, PH8.0, 10 mM EDTA PH8.0,

451 100 mM KCL, 0.5 M sucrose, 4 mM spermidine, 1 mM spermine) on ice for 15 min.

452 After filtration through miracloth, Triton x-100 (Sigma) was added to tubes at a 1:20

453 ratio, placed on ice for 15 min, and centrifuged to collect nuclei. Nuclei were washed

454 with NIB buffer (containing Triton x-100) and re-suspended in a small amount of NIB

455 buffer (containing Triton x-100) then the volume of each tube was brought to 40 ml and

456 centrifuged again. After careful removal of all liquid, 10 ml of Qiagen G2 buffer was

457 added followed by gentle re-suspension of nuclei; then 30 ml G2 buffer with RNase A

458   (to final concentration of 50 mg/ml) was added. Tubes were incubated at 37 °C for 30

459   min. Proteinase K (Invitrogen), 30 mg, was added and tubes were incubated at 50 °C

460   for 2 h followed by centrifugation for 15 min at 8000 rpm, at 4 °C, and the liquid gently

461   poured to a new tube. After gentle extraction with Chloroform / isoamyl alcohol (24:1),

462   then centrifugation and transfer of the top phase to a fresh tube, HMW DNA was

463   precipitated by addition of 2/3 volume of iso-propanol and re-centrifugation to collect the

464   DNA. After DNA was washed with 70% ethanol, it was air dried for 20 min and dissolved

465   thoroughly in 1x TE.

466

467   **Whole-genome sequencing**

468

469   We sequenced *Populus trichocarpa* var. *Stettler* using a whole-genome shotgun

470   sequencing strategy and standard sequencing protocols. Sequencing reads were

471   collected using Illumina and PacBio. Both the Illumina and PacBio reads were

472   sequenced at the Department of Energy (DOE) Joint Genome Institute (JGI) in Walnut

473   Creek, California and the HudsonAlpha Institute in Huntsville, Alabama. Illumina reads

474   were sequenced using the Illumina HISeq platform, while the PacBio reads were

475   sequenced using the RS platform. One 400-bp insert 2x150 Illumina fragment library

476   was obtained for a total of ~349x coverage (Table S10). Prior to assembly, all Illumina

477   reads were screened for mitochondria, chloroplast, and phix contamination. Reads

478   composed of >95% simple sequence were removed. Illumina reads less than 75 bp

479   after trimming for adapter and quality (q < 20) were removed. The final Illumina read set

480   consists of 906,280,916 reads for a total of ~349x of high-quality Illumina bases. For the

481  PacBio sequencing, a total of 69 chips (P6C4 chemistry) were sequenced with a total

482  yield of 59.29 Gb (118.58x) with 56.2 Gb > 5 kb (Table S11), and post error correction a

483  total of 37.3 Gb (53.4x) was used in the assembly.

484

485  **Genome assembly and construction of pseudomolecule chromosomes**

486

487  The current release is version 1.0 release began by assembling the 37.3 Gb corrected

488  PacBio reads (53.4x sequence coverage) using the MECAT CANU v.1.4 assembler [37]

489  and subsequently polished using QUIVER v.2.3.3 [38]. This produced 3,693 scaffolds

490  (3,693 contigs), with a scaffold N50 of 1.9 Mb, 955 scaffolds larger than 100 kb, and a

491  total genome size of 693.8 Mb (Table S12). Alternative haplotypes were identified in the

492  initial assembly using an in-house Python pipeline, resulting in 2,972 contigs (232.3 Mb)

493  being labeled as alternative haplotypes, leaving 745 contigs (461.5 Mb) in the single

494  haplotype assembly. A set of 64,840 unique, non-repetitive, non-overlapping 1.0 kb

495  syntenic sequences from version 4.0 *P. trichocarpa* var. *Nisqually* assembly and aligned

496  to the MECAT CANU v.1.4 assembly and used to identify misjoins in the *P. trichocarpa*

497  var. *Stettler* assembly. A total of 22 misjoins were identified and broken. Scaffolds were

498  then oriented, ordered, and joined together into 19 chromosomes. A total of 117 joins

499  were made during this process, and the chromosome joins were padded with 10,000

500  Ns. Small adjacent alternative haplotypes were identified on the joined contig set.

501  Althap regions were collapsed using the longest common substring between the two

502  haplotypes. A total of 14 adjacent alternative haplotypes were collapsed.

503

504 The resulting assembly was then screened for contamination. Homozygous single

505 nucleotide polymorphisms (SNPs) and insertion/deletions (InDels) were corrected in the

506 release sequence using ~100x of Illumina reads (2x150, 400-bp insert) by aligning the

507 reads using bwa-0.7.17 mem [54] and identifying homozygous SNPs and InDels with

508 the GATK v3.6's UnifiedGenotyper tool [55]. A total of 206 homozygous SNPs and

509 11,220 homozygous InDels were corrected in the release. Heterozygous SNP/indel

510 phasing errors were corrected in the consensus using the 118.58x raw PacBio data. A

511 total of 66,124 (1.98%) of the heterozygous SNP/InDels were corrected. The final

512 version 1.0 improved release contains 391.2 Mb of sequence, consisting of 25 scaffolds

513 (128 contigs) with a contig N50 of 7.5 Mb and a total of 99.8% of assembled bases in

514 chromosomes. Plots of the *Nisqually* marker placements for the 19 chromosomes are

515 shown in Fig. S8.

516

517 **Genome annotation**

518

519 Transcript assemblies were made from ~1.4 billion pairs of 2x150 stranded paired-end

520 Illumina RNA-seq GeneAtlas *P. trichocarpa* Nisqually reads, ~1.2 billion pairs of 2x100

521 paired-end Illumina RNA-seq *P. trichocarpa* Nisqually reads from Dr. Pankaj Jaiswal,

522 and ~430M pairs of 2x75 stranded paired-end Illumina var. *Stettler* reads using

523 PERTRAN (Shu, unpublished) on *P. trichocarpa* var. *Stettler* genome. About ~3M

524 PacBio Iso-Seq circular consensus sequences were corrected and collapsed by

525 genome guided correction pipeline (Shu, unpublished) on *P. trichocarpa* var. *Stettler*

526 genome to obtain ~0.5 million putative full-length transcripts. 293,637 transcript

527    assemblies were constructed using PASA [56] from RNA-seq transcript assemblies

528    above. Loci were determined by transcript assembly alignments and/or EXONERATE

529    alignments of proteins from *A. thaliana*, soybean, peach, Kitaake rice, *Setaria viridis*,

530    tomato, cassava, grape and Swiss-Prot proteomes to repeat-soft-masked *P. trichocarpa*

531    var. *Stettler* genome using RepeatMasker [57] with up to 2-kb extension on both ends

532    unless extending into another locus on the same strand. Gene models were predicted

533    by homology-based predictors, FGENESH+[58], FGENESH_EST (similar to

534    FGENESH+, EST as splice site and intron input instead of protein/translated ORF), and

535    EXONERATE [59], PASA assembly ORFs (in-house homology constrained ORF finder)

536    and from AUGUSTUS via BRAKER1 [60]. The best scored predictions for each locus

537    are selected using multiple positive factors including EST and protein support, and one

538    negative factor: overlap with repeats. The selected gene predictions were improved by

539    PASA. Improvement includes adding UTRs, splicing correction, and adding alternative

540    transcripts. PASA-improved gene model proteins were subject to protein homology

541    analysis to above mentioned proteomes to obtain Cscore and protein coverage. Cscore

542    is a protein BLASTP score ratio to MBH (mutual best hit) BLASTP score and protein

543    coverage is highest percentage of protein aligned to the best of homologs. PASA-

544    improved transcripts were selected based on Cscore, protein coverage, EST coverage,

545    and its CDS overlapping with repeats. The transcripts were selected if its Cscore is

546    larger than or equal to 0.5 and protein coverage larger than or equal to 0.5, or it has

547    EST coverage, but its CDS overlapping with repeats is less than 20%. For gene models

548    whose CDS overlaps with repeats for more that 20%, its Cscore must be at least 0.9

549    and homology coverage at least 70% to be selected. The selected gene models were

550     subject to Pfam analysis and gene models whose protein is more than 30% in Pfam TE

551     domains were removed and weak gene models. Incomplete gene models, low

552     homology supported without fully transcriptome supported gene models and short single

553     exon (< 300-bp CDS) without protein domain nor good expression gene models were

554     manually filtered out.

555

556     **SNP calling methods**

557

558     Illumina HiSeq2500 paired-end (2×150) reads were mapped to the reference genome

559     using bwa-mem [54]. Picard toolkit was used to sort and index the bam files. GATK [55]

560     was used further to align regions around InDels. Samtools v1.9 [61] was used to create

561     a multi-sample mileup for each tree independently. Preliminary SNPs were called using

562     Varscan v2.4.0 [62] with a minimum coverage of 21.

563

564     At these SNPs, for each branch, we calculated the conditional probability of each

565     potential genotype (RR, RA, AA) given the read counts of each allele, following SeqEM

566     [63], using an estimated sequencing error rate of 0.01. We identified high-confidence

567     genotype calls as those with a conditional probability 10,000x greater than the

568     probabilities of the other possible genotypes. We identified potential somatic SNPs as

569     those with both a high-confidence homozygous and high-confidence heterozygous

570     genotype across the branches.

571

572 We notice that the default SNP calling parameters tend to overcall homozygous-

573 reference allele genotypes and that differences in sequencing depth can bias the

574 relative number of heterozygous SNPs detected. To overcome these issues, we re-

575 called genotypes using conditional probabilities using down sampled allele counts. To

576 do this, we first randomly selected a set number of sequencing reads for each library at

577 each potential somatic SNP so that all libraries have the same sequencing depth at all

578 SNPs. Using the down sampled reads, we calculate the relative conditional probability

579 of each genotypes by dividing the conditional probabilities by the sum of the conditional

580 probabilities of all three potential genotypes. These relative probabilities are then

581 multiplied by the dosage assigned to their respective genotype (0 for RR, 1 for RA, 2 for

582 AA), and the dosage genotype is the sum of these values across all 3 possible

583 genotypes. Discrete genotypes were assigned using the following dosage values: RR =

584 dosage < 0.1; RA = 0.9 < dosage < 1.1; AA = dosage > 1.9. Dosages outside those

585 ranges are assigned a NA discrete genotype. SNPs with an NA discrete genotype or

586 depth below the down sampling level in any branch of a tree were removed from further

587 analysis. We performed three replicates of this procedure for depths of 20, 25, 30, 35,

588 40, and 45 reads.

589

590 PacBio libraries for each branch were sequenced using the PacBio Sequel platform,

591 fastq files aligned to the *P. trichocarpa* var. Stettler14 reference genome using ngmlr

592 [64], and multi-sample mileup files generated using in Samtools v1.9 [61] to quantify the

593 allele counts at the potential somatic SNPs. We used a minimum per-sample sequence

594    depth of 20 reads and used an alternate-allele threshold of 0.1 to call a heterozygote

595    genotype in the PacBio data.

596

597    To identify high-confidence candidate somatic SNPs, we identified potential somatic

598    SNPs with the same genotypes across branches using both the Illumina-based PacBio-

599    based genotypes, only including SNPs with full data in all branches for both types of

600    genotypes. Of these, we only retained SNPs that are homozygous in a single branch or

601    have a single homozygous-to-heterozygous transition (and no reversion) going from the

602    lowest to highest branches.

603

604    **Estimating somatic nucleotide mutation rate**

605

606    Building on the analytical framework developed in van der Graaf et al. (2015) and

607    Shahryary et al. 2019 (co-submission), we developed *mutSOMA*

608    (https://github.com/jlab-code/mutSOMA), a statistical method for estimating genetic

609    mutation rates in long-lived perennials such as trees. The method treats the tree

610    branching structure as a pedigree of somatic lineages and uses the fact that these cell

611    lineages carry information about the mutational history of each branch. A detailed

612    mathematical description of the method is provided in Supplementary Text. But briefly,

613    starting from the .vcf* files from $S$ samples representing different branches of the tree,

614    we let $G_{ik}$ be the observed genotype at the $k$-th single nucleotide ($k$ = 1, …, $N$) in the $i$-th

615    sample, where $N$ is the effective genome size (i.e. the total number of bases with

616    sufficient coverage). With four possible nucleotides (A, C, T, G) , $G_{ik}$ can have 16

617    possible genotypes in a diploid genome, 4 homozygous (A|A, T|T, C|C, G|G) and 12

618    heterozygous (A|G, A|T, …, G|C). Using this coding, we calculate the genetic

619    divergence, $D$, between any two samples $i$ and $j$ as follows:

620

621
$$D_{ij} = \sum_{k=1}^{N} I(G_{ik}, G_{jk}) N^{-1},$$

622

623    where $I(G_{ik}, G_{jk})$ is an indicator function, such that, $I(G_{ik}, G_{jk})$ = 1 if the two samples

624    share no alleles at locus $k$, 0.5 if they share one, and 0 if they share both alleles. We

625    suppose that $D_{ij}$ is related to the developmental divergence time of samples $i$ and $j$

626    through a somatic mutation model $M_\Theta$. The divergence times can be calculated from the

627    coring data (Table S13). We model the genetic divergence using

628

629
$$D_{ij} = c + D_{ij}^{\bullet}(M_\Theta) + \epsilon_{ij},$$

630

631    where $\epsilon_{ij} \sim N(0, \sigma^2)$ is the normally distributed residual, $c$ is the intercept, and $D_{ij}^{\bullet}(M_\Theta)$

632    is the expected divergence as a function of mutation model $M$ with parameter vector $\Theta$.

633    Parameter vector $\Theta$ contains the unknown mutation rate $\delta$ and the unknown proportion

634    $\gamma$ heterozygote loci of the most recent common 'founder' cells of samples $i$ and $j$. The

635    theoretical derivation of $D_{ij}^{\bullet}(M_\Theta)$ and details regarding model estimation can be found in

636    Supplementary Text. The estimation of the residual variance in the model allows for the

637    fact that part of the observed genetic divergence between any two samples is driven

638 both by genotyping errors as well as by somatic genetic drift as meristematic cells pass

639 through bottlenecks in the generation of the lateral branches.

640

641 **Structural variant analysis methods**

642

643 For structural variant (SV) analysis, PacBio libraries were generated for four branches

644 from the tree 13 and four branches from tree 14 with four sequencing cells sequenced

645 per branch using the PacBio Sequel platform. PacBio fastq files were aligned to the *P.*

646 *trichocarpa* var. *Stettler* reference genome using ngmlr v.0.2.6 [64] using a value of 0.01

647 for the "-R" flag. SVs were discovered and called using pbsv (pbsv v2.2.0,

648 https://github.com/PacificBiosciences/pbsv). SV signatures were identified for each

649 sample using 'pbsv discover' using the '--tandem-repeats' flag and a tandem repeat

650 BED file generated using trf v4.09 [65] for the *P. trichocarpa* var. *Stettler* genome. SVs

651 were called jointly for all 8 branches using 'pbsv call'.  The output from joint SV calling

652 changes slightly depending on the order of the samples used for the input in 'pbsv call',

653 so four sets of SVs were generated using four different sample orders as input. We

654 used a custom R script [66] to filter the SV output from pbsv. We remove low-complexity

655 insertions or deletions with sequence containing > 80% of a mononucleotide 8-mer,

656 50% of a single type of binucleotide 8-mer, or 60% of two types of binucleotide 8-

657 mers.  We required a minimum distance of 1 kb between SVs. We removed SVs with

658 sequencing coverage of more than three standard deviations above the mean coverage

659 across a sample. After calling genotypes, any SVs with missing genotype data were

660 removed.

661

Genotypes were called based on the output from pbsv using a custom R script. We

required a minimum coverage of 10 reads in all sample and for one sample to have at

least 20 reads. We required a minimum penetrance (read ratio) of 0.25 and at least 2

reads containing the minor allele for a heterozygous genotype. We allowed a maximum

penetrance of 0.05 for homozygous genotypes. For each genotype, we assigned a

quality score based on the binomial distribution-related relative probability of the 3

genotype classes (RR, AR, AA) based on A:R read ratio, using an estimated

sequencing error of 0.032, and an estimated minimum allele penetrance of 0.35. For a

genotype with a score below 0.9 but with the same genotype at the SV as another

sample with a score above 0.98, the score was adjusted by multiplying by 1.67. Any

genotypes with adjusted scores below 0.9 were converted to NA. For deletions,

duplications, and insertions, 10 representatives in different size classes were randomly

selected and the mapping patterns of reads were visually inspected using IGV v2.5.3

[67] to assign scores indicating how well the visual mapping patterns support the SV

designation. Scores were defined by the following: "strong", multiple reads align to the

same locations in the reference genome that support the SV type and size; "moderate",

multiple reads align to the same reference location for one side of the SV but align to

different or multiple locations in the region for the other side of the SV; and "weak",

reads align to reference locations that indicate a different SV type or much different SV

size.

682

683     The percent of genic sequence and tandem repeat sequence in deletions and

684     duplications were calculated using the *P. trichocarpa* var. *Stettler* annotation and

685     tandem repeat BED from above, respectively. Genome-wide expectations were derived

686     by separating the genome into 10-kb windows and calculating the percent genic and

687     tandem repeat sequence in each window. The distribution of genic and tandem repeat

688     sequences in deletions and duplications were compared to genome-wide expectations

689     using the Kolmogorov-Smirnov two-sample test (one-sided, $N_{null}$ = 39,151, $N_{del}$ =

690     10,433, $N_{dup}$ = 630).

691

692     SVs showing variation between branches and identified in all 4 replicates are potential

693     instances of somatic SV mutations or loss-of-heterozygosity gene conversions, and the

694     mapping positions of sequencing reads were visually inspected with IGV [67] to confirm

695     the variation at these SVs.

696

697     **MethylC-seq sequencing and analysis**

698

699     A single MethylC-seq library was created for each branch from leaf tissue. Libraries

700     were prepared according to the protocol described in Urich *et al.* [68]. Libraries were

701     sequenced to 150-bp per read at the Georgia Genomics & Bioinformatics Core (GGBC)

702     on a NextSeq500 platform (Illumina). Average sequencing depth was ~41.1x among

703     samples (Table S7).

704

705 MethylC-seq reads were processed and aligned using Methylpy v1.3.2 [69]. Default

706 parameters were used expect for the following: clonal reads were removed, lambda

707 DNA was used as the unmethylated control, and binomial test was performed for all

708 cytosines with at least three mapped reads.

709

710 **Identification of Differentially Methylated Regions**

711

712 Identification of differentially methylated regions (DMRs) was performed using Methylpy

713 v1.3.2 [69]. All methylome samples were analyzed together to conduct an undirected

714 identification of DMRs across all samples in the CNN (N=A, C, G, T) context. Default

715 parameters were used. Only DMRs at least 40-bp long with at least three differentially

716 methylated cytosines (DMS) and five or more cytosines with at least one read were

717 retained. For each DMR, the weighted methylation level was computed as mC / (mC +

718 uC) where mC and uC are the number of reads supporting a methylated cytosine and

719 unmethylated cytosine, respectively [41].

720

721 To identify epigenetic variants in these samples, we used a one-sided z-test to test for a

722 significant difference in methylation level of DMRs pairwise between branches. For each

723 pair, only DMRs with at least 5% difference in methylation level were used, regardless

724 of underlying context. Resulting $P$ values were adjusted using Benjamini-Hochberg

725 correction (N = 383,600) with FDR = 0.05 [70] and DMRs are defined by adjusted $P$

726 value ≤ 0.05.

727

**Identification of Methylated Regions**

For each sample, an unmethylated methylome was generated by setting the number of methylated reads to zero while maintaining the total number of reads. Methylpy DMR identification program [69] was applied to each sample using the original methylome and unmethylated methylome with the same parameters as used for DMR identification. Regions less than 40 bp-long, fewer than three DMS, and fewer than five cytosines with at least one read were removed. Remaining regions from all samples were merged using BEDtools v2.27.1 [71].

**Assigning genomic features to DMRs**

A genomic feature map was created such that each base pair of the genome was assigned a single feature type (transposable element/repeat, promoter, untranslated region, coding sequence, and intron) based on the previously described annotation. Promoters were defined as 2 kb upstream of the transcription start site of protein-coding genes. At positions where multiple feature types could be applicable, such as a transposon in an intron or promoter overlapping with adjacent gene, priority was given to untranslated regions (highest), introns, coding sequences, promoter, and transposon (lowest). Positions without an assignment were considered intergenic. Genomic feature content of each DMR and methylated region was assigned proportionally based on the number of bases in each category.

**Identification of pseudo-allele methylation**

751

752

753 We aimed to categorize the DMRs into three pseudo-allele states: homozygous

754 methylated, heterozygous, and homozygous unmethylated. First, DMRs were filtered on

755 the following criteria: i) at least 25% change in weighted CG methylation level between

756 the highest and lowest methylation level of the samples; ii) at least one sample had a

757 CG methylation level of at least 75%; and iii) at least two "covered" CG positions. A

758 "covered" CG is defined as having at least one read for both symmetrical cytosines in all

759 samples. After filtering, 4,488 regions were used for analysis.

760

761 For each region in each sample, we next categorize the aligned reads overlapping the

762 region. If at least 35% of its "covered" CG sites are methylated, the read is categorized

763 as methylated. Otherwise it is an unmethylated read. Finally, we define the pseudo-

764 allele state by the portion of methylated reads; homozygous unmethylated: ≤ 25%,

765 heterozygous: > 25% and < 75%, and homozygous methylated: ≥ 75%.

766

767 The null distribution was created by randomly shuffling the filtered DMRs in the genome

768 such that each simulated region is the same length as the original and it has at least two

769 "covered" CGs. The above procedure was applied and number of epigenotype changes

770 was determined. This was repeated for a total of 10 times.

771

772 The following special classes of DMRs were identified: highly variable, single loss,

773 single gain, and tree specific. A DMR is highly variable if there were pseudo-allele

774    changes between all adjacent branches. A DMR is single loss if all but one branch was

775    homozygous methylated, and one was homozygous unmethylated. Similarly, a DMR is

776    single gain if all but one branch was homozygous unmethylated and one branch was

777    homozygous methylated. Finally, a DMR is "tree specific" if all tree 13 branches were

778    homozygous unmethylated and all tree 14 branches were homozygous methylated or

779    vice versa.

780

781    **Estimating somatic epimutation rate**

782

783    We previously developed a method for estimating 'germline' epimutation rates in *A.*

784    *thaliana* based on multi-generational methylation data from Mutation Accumulation lines

785    [32]. In a companion method paper to the present study (Shahryary et al. 2019, co-

786    submission), we have extended this approach to estimating somatic epimutation rates in

787    long-lived perennials such as trees using leaf methylomes and coring data as input.

788    This new inference method, which we call *AlphaBeta*, treats the tree branching structure

789    as a pedigree of somatic lineages using the fact that these cell lineages carry

790    information about the epimutational history of each branch. *AlphaBeta* is implemented

791    as a bioconductor R package

792    (http://bioconductor.org/packages/devel/bioc/html/AlphaBeta.html). Using this approach,

793    we estimate somatic epimutation rates for individual CG, CHG, and CHH sites

794    independently, but also for regions. For the region-level analysis, we first use the

795    differentially methylated regions (DMRs) identified above. Sampling from the distribution

796    of DMR sizes, we then split the remainder of the genome into regions, which we refer to

797    as "non-DMRs". Per sample, we aggregate the total number of methylated Cs and

798    unmethylated Cs in each region corresponding to a DMRs or a non-DMRs and used

799    these counts as input for *AlphaBeta*.

800

801    **mRNA-seq sequencing and analysis**

802

803    Total RNA was extracted from leaf tissue in each branch using the Direct-zol RNA

804    MiniPrep Plus kit (Zymo Research) with Invitrogen's Plant RNA Reagent. Total RNA

805    quality and quantity were assessed before library construction. Strand-specific RNA-seq

806    libraries were constructed using the TruSeq Stranded mRNA LT kit (Illumina) following

807    the manufacturer's instructions. For each sample, three independent libraries (technical

808    replicates) were constructed. Libraries were sequenced to paired-end 75-bp reads at

809    the GGBC on a NextSeq500 platform (Illumina). Summary statistics are included in the

810    Table S9.

811

812    For analysis, first, paired-end reads were trimmed using Trimmomatic v0.36 [72].

813    Trimming included removing TruSeq3 adapters, bases with quality score less than 10,

814    and any reads less than 50-bp long. Second, remaining reads were mapped to the

815    *Stettler* genome with HiSAT2 [73] using default parameters except to report alignments

816    for transcript assemblers (--dta). The HiSAT2 transcriptome index was created using

817    extracted splice sites and exons from the gene annotation as recommended. Last,

818    transcriptional abundances for genes in the reference annotation were computed for

819    each sample using StringTie v1.3.4d [74]. Default parameters were used except to limit

820    estimates to reference transcripts. TPM (transcripts per million) values were outputted

821    to represent transcriptional abundance.

822

823    **Identification of differentially expressed genes**

824

825    Differentially expressed genes (DEGs) were identified using DeSeq2 v1.22.2 [46]. The

826    count matrix was extracted from StringTie output files and the analysis was performed

827    using the protocol (ccb.jhu.edu/software/stringtie/index.shtml?t=manual#deseq).

828    Abundances for all samples were joined into one DESeq dataset with $\alpha = 0.01$. Gene

829    abundance was compared between all samples pairwise. In each pair, a gene was

830    considered differentially expressed if the adjusted $P$ value $\leq 0.01$ and the $\log_2$-fold

831    change $\geq 1$. Genes differentially expressed in any pair were included for subsequent

832    analysis.

833

834    **Overlap of DMRs and DEGs**

835

836    We identified DMRs which overlapped the promoter region (2 kb upstream of

837    transcription start site) and gene body of annotated genes. For each DMR-gene pair, we

838    computed the Pearson's product moment correlation coefficient between the weighted

839    methylation level of the DMR and average gene abundance among replicates in TPM.

840    Next, looking only at genes which were previously identified as differently expressed,

841    we performed a two-sided Pearson's correlation test for each DMR-DEG pair to test for

842    statistically significant correlations. Resulting $P$ values were multiple test corrected with

843   Benjamini-Hochberg correction (N = 382, FDR = 0.05) [70]. Adjusted *P* values ≤ 0.05

844   were considered significantly correlated.

845

846   **DECLARATIONS**

847

848   **Ethics approval and consent to participate**

849

850   Not applicable

851

852   **Consent for publication**

853

854   Not applicable

855

856   **Availability of data and materials**

857

858   Raw sequence data used for genome assembly, resequencing and identification of

859   structural variation of individual branches are available at NCBI SRA (PRJNA516415).

860   Raw sequence data for whole-genome bisulfite sequencing and mRNA-sequencing are

861   available in GEO under accession GSE132939.

862

863   Custom analysis scripts used in this study are available in the GitHub repository

864   https://github.com/schmitzlab/somatic-epigenetic-mutation-poplar.

865

**Competing interests**

The authors declare that they have no competing interests.

**Authors' Contributions**

888   RJS, FJ, GAT, RS and JS conceived and designed the experiments. JG, SS, KB, KL,

889   CA, AL, DK, JT, RW performed data generation. BTH, JD, MCT, YS, RH, SM, JJ, PPG,

890   FJ performed data analysis. BTH prepared the figures and manuscript. BTH, DWH,

891   GAT, FJ, and RJS wrote and revised the manuscript with input from all authors. All

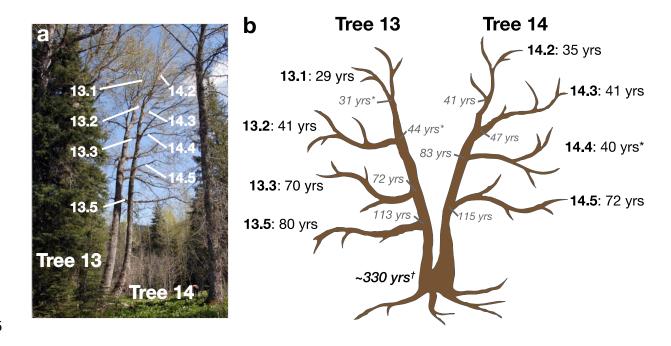892   authors read and approved the final manuscript.

893

894   **Acknowledgements**

895

903

904   **FIGURES**

905

**Fig. 1. Photograph and schematic drawing of Tree 13 and Tree 14.** This wild *P. trichocarpa*, located near Mt. Hood, Oregon, experienced a decapitation event ~300 years ago. Tree 14 re-sprouted from the stump and ~80-100 years later Tree 13 re-sprouted. (a) Leaf samples were collected from the labeled terminal branches. (b) Age was estimated for both the end of the branch (black font) and where it meets the main stem (gray italics). Ages with * indicate age was estimated using diameter; all other estimates were from core samples. Leaf samples of each branch was used to create genomic sequencing libraries, PacBio libraries, whole-genome bisulfite sequencing libraries, and mRNA-sequencing libraries.

**Fig. 2. Most somatic mutations are transitions and occur in non-genic regions.** (a) Distribution of reference to alternative allele observed in the high-confidence SNPs identified in Tree 13 and Tree 14. (b) Distribution of high-confidence SNPs separated by the genomic feature. Abbreviations: Pro, promoter; 2 kb upstream of TSS; TE, transposable elements and repeats; and IGR, intergenic regions.

42

**Fig. 3. Somatic epimutation rates for single sites, regions, and by genomic feature.** mCG (a) and mCHG (b) divergence by branch time divergence for single sites and regions; mCG (c) and mCHG (d) divergence by branch time divergence for

43

genomic features Pro (promoter; 2 kb upstream of TSS), mRNA, TE (transposable

elements), and IGR (intergenic regions); Estimated mCG (e) and mCHG (f) gain rates

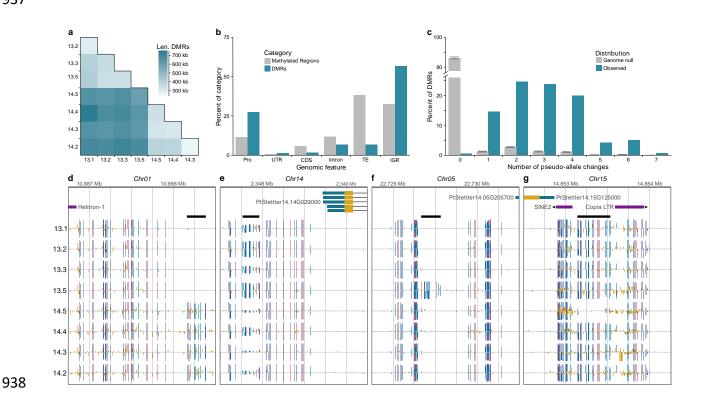by feature; Estimated mCG (g) and mCHG (h) loss rates by feature; Ratio of mCG (i)

and mCHG (j) loss to gain by feature. Error bars represent bootstrapped 95%

confidence intervals of the estimates. Abbreviations: Pro, promoter; 1.5 kb upstream of

TSS; TE, transposable elements and repeats; and IGR, intergenic regions.



**Fig. 4. Identification and quantification of somatic stability of differentially
methylated regions.** (a) Divergence of differentially methylated regions corresponds to
divergence in age. The darker color indicates combined length of the pairwise DMRs;
(b) The genome-wide distribution of identified DMRs in genomic features. Abbreviations:
TE, transposable elements and repeats; IGR, intergenic region; Pro, promoter region (2

945     kb upstream of transcription start site); UTR, untranslated regions; CDS, coding

946     sequence. Methylated regions were identified in as regions methylated in at least one

947     sample. (c) There are significantly more pseudo-allele changes between the branches

948     at DMRs (blue) compared to the genome-wide null (Wilcox rank sum, one-sided, *P*

949     value < 2 x 10$_{-16}$). Gray bars are the genome-wide null as mean +/- std. dev. across 10

950     simulations. (d) Browser screenshot of a tree specific DMR where all branches in tree

951     13 are homozygous unmethylated and all branches of tree 14 are homozygous

952     methylated. (e) Browser screenshot of a highly variable DMR where the pseudo allele

953     state changes between each branch. (f) Browser screenshot of a single gain DMR

954     where all branches except 13.5 are homozygous unmethylated and 13.5 gains

955     methylation. (g) Browser screenshot of a single loss DMR where all branches except

956     14.5 are homozygous methylated and 14.5 has lost methylation. For d-g, gene models

957     and transposable elements are shown at the top and methylome tracks are below.

958     Vertical bars indicate methylation at the position, where height corresponds to level and

959     color is context, red for CG, blue for CHG, and yellow for CHH. DMR is indicated by

960     thick black horizontal line.

961

962

963

**Fig. 5. Gene expression is largely independent from divergence age and nearby cytosine methylation except in a few examples.** a) Gene expression divergence is not significantly associated with divergence age. b) Distribution of positive and negative correlations for differentially expressed genes and overlapping/nearby DMRs. Positive

968     correlation occurs when the higher methylation level is associated with higher gene

969     expression among the samples. (c) A significantly negatively correlated, tree-specific

970     DMR and DEG where the DMR occurs in the promoter region of the gene (Pearson's

971     correlation test, two-sided, N = 8, adjusted *P* value = 0.0067). The higher methylation

972     levels in the DMR for tree 13 branches are associated with lower gene expression. (d) A

973     significantly positively correlated, single gain DMR and DEG where the DMR occurs in

974     the 5' untranslated region of the gene (Pearson's correlation test, two-sided, N = 8,

975     adjusted *P* = 0.0141). The higher methylation level in the DMR for branch 13.1 is

976     associated with greater gene expression. For c and d, gene expression, as transcripts

977     per million (TPM), is represented as points for the individual replicates and as bar for

978     mean among replicates. In the genome browser view, gene models and transposable

979     elements are shown at the top and methylome tracks are below. Vertical bars indicate

980     methylation at the position, where height corresponds to level and color is context, red

981     for CG, blue for CHG, and yellow for CHH. DMR is indicated by thick black horizontal

982     line.

983

984     **REFERENCES**

985

986     1.    Whitham TG, Slobodchikoff CN. Evolution by individuals, plant-herbivore

987           interactions, and mosaics of genetic variability: The adaptive significance of

988           somatic mutations in plants. Oecologia. 1981; 49: 287.

989     2.    Walbot V. On the life strategies of plants and animals. Trends in Genetics. 1985;

990           doi:10.1016/0168-9525(85)90071-X.

991    3.    Gill DE. Individual plants as genetic mosaics: Ecological organisms versus

992          evolutionary individuals. In: Crawley MJ, editor. Plant Ecology. Oxford: Blackwell

993          Scientific Publications; 1986. p. 321-343.

994    4.    Gill DE, Chao L, Perkins SL, Wolf JB. Genetic mosaicism in plants and clonal

995          animals. Annual review of Ecology and Systematics. 1995; 26: 423-444.

996    5.    Hadany L. A conflict between two evolutionary levels in trees. J Theor Biol. 2001

997          Feb 21; doi:10.1006/jtbi.2000.2236.

998    6.    Clarke E. Plant individuality and multilevel selection theory. The major transitions in

999          evolution revisited. 2011; 227-250.

1000   7.    Folse III HJ, Roughgarden J. Direct benefits of genetic mosaicism and

1001          intraorganismal selection: modeling coevolution between a long-lived tree and a

1002          short-lived herbivore. Evolution: International Journal of Organic Evolution. 2012;

1003          66: 1091-1113.

1004   8.    Tuskan GA, Groover AT, Schmutz J, DiFazio SP, Myburg A, Grattapaglia D et al.

1005          Hardwood Tree Genomics: Unlocking Woody Plant Biology. Frontiers in Plant

1006          Science. 2018; doi:10.3389/fpls.2018.01799.

1007   9.    Padovan A, Keszei A, Foley WJ, K√ºlheim C. Differences in gene expression

1008          within a striking phenotypic mosaic Eucalyptustree that varies in susceptibility to

1009          herbivory. BMC Plant Biology. 2013; doi:10.1186/1471-2229-13-29.

1010   10.   Wen I-C, Koch KE, Sherman WB. Comparing Fruit and Tree Characteristics of

1011          Two Peaches and Their Nectarine Mutants. J Amer Soc Hort Sci. 1995;

1012          doi:10.21273/JASHS.120.1.101.

1013    11.   Laucou V, Lacombe T, Dechesne F, Siret R, Bruno J-P, Dessup M et al. High

1014           throughput analysis of grape genetic diversity as a tool for germplasm collection

1015           management. Theor Appl Genet. 2011; doi:10.1007/s00122-010-1527-y.

1016    12.   Tuskan GA, Francis KE, Russ SL, Romme WH, Turner MG. RAPD markers reveal

1017           diversity within and among clonal and seedling stands of aspen in Yellowstone

1018           National Park, U.S.A. Can J For Res. 1996; doi:10.1139/x26-237.

1019    13.   Diwan D, Komazaki S, Suzuki M, Nemoto N, Aita T, Satake A et al. Systematic

1020           genome sequence differences among leaf cells within individual trees. BMC

1021           genomics. 2014; 15: 142.

1022    14.   Schmid-Siegert E, Sarkar N, Iseli C, Calderon S, Gouhier-Darimont C, Chrast J et

1023           al. Low number of fixed somatic mutations in a long-lived oak tree. Nat Plants.

1024           2017 Dec; doi:10.1038/s41477-017-0066-9.

1025    15.   Plomion C, Aury J-M, Amselem J, Leroy T, Murat F, Duplessis S et al. Oak

1026           genome reveals facets of long lifespan. Nat Plants. 2018 07; doi:10.1038/s41477-

1027           018-0172-3.

1028    16.   Ossowski S, Schneeberger K, Lucas-Lledö JI, Warthmann N, Clark RM, Shaw RG

1029           et al. The rate and molecular spectrum of spontaneous mutations in Arabidopsis

1030           thaliana. Science. 2010; 327: 92-94.

1031    17.   Groot EP, Laux T. Ageing: How Do Long-Lived Plants Escape Mutational

1032           Meltdown. Curr Biol. 2016 07 11; doi:10.1016/j.cub.2016.05.049.

1033    18.   Wang L, Ji Y, Hu Y, Hu H, Jia X, Jiang M et al. The architecture of intra-organism

1034           mutation rate variation in plants. PLoS biology. 2019; 17: e3000191.

1035  19.  Burian A, Barbier de Reuille P, Kuhlemeier C. Patterns of Stem Cell Divisions

1036       Contribute to Plant Longevity. Current Biology. 2016;

1037       doi:10.1016/j.cub.2016.03.067.

1038  20.  Klekowski EJ, Godfrey PJ. Ageing and mutation in plants. Nature. 1989;

1039       doi:10.1038/340389a0.

1040  21.  Bobiwash K, Schultz ST, Schoen DJ. Somatic deleterious mutation rate in a woody

1041       plant: estimation from phenotypic data. Heredity. 2013; 111: 338.

1042  22.  Schmitz RJ, Schultz MD, Urich MA, Nery JR, Pelizzola M, Libiger O et al. Patterns

1043       of population epigenomic diversity. Nature. 2013 Mar 14; doi:10.1038/nature11968.

1044  23.  Calarco JP, Borges F, Donoghue MT, Van Ex F, Jullien PE, Lopes T et al.

1045       Reprogramming of DNA methylation in pollen guides epigenetic inheritance via

1046       small RNA. Cell. 2012 Sep 28; doi:10.1016/j.cell.2012.09.001.

1047  24.  Law JA, Jacobsen SE. Establishing, maintaining and modifying DNA methylation

1048       patterns in plants and animals. Nat Rev Genet. 2010 Mar; doi:10.1038/nrg2719.

1049  25.  Johannes F, Schmitz RJ. Spontaneous epimutations in plants. New Phytol. 2019

1050       Feb; doi:10.1111/nph.15434.

1051  26.  Cubas P, Vincent C, Coen E. An epigenetic mutation responsible for natural

1052       variation in floral symmetry. Nature. 1999; 401: 157-161.

1053  27.  Manning K, Tör M, Poole M, Hong Y, Thompson AJ, King GJ et al. A naturally

1054       occurring epigenetic mutation in a gene encoding an SBP-box transcription factor

1055       inhibits tomato fruit ripening. Nat Genet. 2006 Aug; doi:10.1038/ng1841.

1056  28.  Ong-Abdullah M, Ordway JM, Jiang N, Ooi S-E, Kok S-Y, Sarpan N et al. Loss of

1057       Karma transposon methylation underlies the mantled somaclonal variant of oil

1058       palm. Nature. 2015; 525: 533-537.

1059  29.  Schmitz RJ, Schultz MD, Lewsey MG, O'Malley RC, Urich MA, Libiger O et al.

1060       Transgenerational epigenetic instability is a source of novel methylation variants.

1061       Science. 2011; doi:DOI: 10.1126/science.1212959.

1062  30.  Becker C, Hagmann J, Müller J, Koenig D, Stegle O, Borgwardt K et al.

1063       Spontaneous epigenetic variation in the Arabidopsis thaliana methylome. Nature.

1064       2011 Dec 8; doi:10.1038/nature10555.

1065  31.  Hofmeister BT, Lee K, Rohr NA, Hall DW, Schmitz RJ. Stable inheritance of DNA

1066       methylation allows creation of epigenotype maps and the study of epiallele

1067       inheritance patterns in the absence of genetic variation. Genome Biol. 2017 Aug

1068       16; doi:10.1186/s13059-017-1288-x.

1069  32.  van der Graaf A, Wardenaar R, Neumann DA, Taudt A, Shaw RG, Jansen RC et

1070       al. Rate, spectrum, and evolutionary dynamics of spontaneous epimutations. Proc

1071       Natl Acad Sci U S A. 2015 May 11; doi:10.1073/pnas.1424254112.

1072  33.  Heer K, Ullrich KK, Hiss M, Liepelt S, Schulze Brüning R, Zhou J et al. Detection of

1073       somatic epigenetic variation in Norway spruce via targeted bisulfite sequencing.

1074       Ecology and Evolution. 2018; doi:10.1002/ece3.4374.

1075  34.  Liang D, Zhang Z, Wu H, Huang C, Shuai P, Ye C-Y et al. Single-base-resolution

1076       methylomes of Populus trichocarpa reveal the association between DNA

1077       methylation and drought stress. BMC Genet. 2014; doi:10.1186/1471-2156-15-S1-

1078       S9.

1079   35.  Su Y, Bai X, Yang W, Wang W, Chen Z, Ma J et al. Single-base-resolution

1080          methylomes of Populus euphratica reveal the association between DNA

1081          methylation and salt stress. Tree Genetics & Genomes. 2018; doi:10.1007/s11295-

1082          018-1298-1.

1083   36.  Le Gac AL, Lafon-Placette C, Chauveau D, Segura V, Delaunay A, Fichot R et al.

1084          Winter-dormant shoot apical meristem in poplar trees shows environmental

1085          epigenetic memory. J Exp Bot. 2018 Sep 14; doi:10.1093/jxb/ery271.

1086   37.  Xiao CL, Chen Y, Xie SQ, Chen KN, Wang Y, Han Y et al. MECAT: fast mapping,

1087          error correction, and de novo assembly for single-molecule sequencing reads. Nat

1088          Methods. 2017 Nov; doi:10.1038/nmeth.4432.

1089   38.  Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C et al.

1090          Nonhybrid, finished microbial genome assemblies from long-read SMRT

1091          sequencing data. Nat Methods. 2013 Jun; doi:10.1038/nmeth.2474.

1092   39.  Ingvarsson PK. Multilocus patterns of nucleotide polymorphism and the

1093          demographic history of Populus tremula. Genetics. 2008;

1094          doi:10.1534/genetics.108.090431.

1095   40.  Rood SB, Polzin ML. Big old cottonwoods. Can J Bot. 2003; 81: 764-767.

1096   41.  Schultz MD, Schmitz RJ, Ecker JR. 'Leveling' the playing field for analyses of

1097          single-base resolution DNA methylomes. Trends Genet. 2012 Dec;

1098          doi:10.1016/j.tig.2012.10.012.

1099   42.  Bewick AJ, Ji L, Niederhuth CE, Willing E-M, Hofmeister BT, Shi X et al. On the

1100          origin and evolutionary consequences of gene body DNA methylation. PNAS.

1101          2016; doi:10.1073/pnas.1604666113.

1102    43.  Wendte JM, Zhang Y, Ji L, Shi X, Hazarika RR, Shahryary Y et al. Epimutations

1103        are associated with CHROMOMETHYLASE 3-induced de novo DNA methylation.

1104        eLife. 2019; doi:10.7554/elife.47891.001.

1105    44.  Melquist S, Luff B, Bender J. Arabidopsis PAI gene arrangements, cytosine

1106        methylation and expression. Genetics. 1999; 153: 401-413.

1107    45.  Silveira AB, Trontin C, Cortijo S, Barau J, Del Bem LE, Loudet O et al. Extensive

1108        natural epigenetic variation at a de novo originated gene. PLoS Genet. 2013 Apr;

1109        doi:10.1371/journal.pgen.1003437.

1110    46.  Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion

1111        for RNA-seq data with DESeq2. Genome Biology. 2014; doi:10.1186/s13059-014-

1112        0550-8.

1113    47.  Narsai R, Gouil Q, Secco D, Srivastava A, Karpievitch YV, Liew LC et al. Extensive

1114        transcriptomic and epigenomic remodelling occurs during Arabidopsis thaliana

1115        germination. Genome Biol. 2017 09 15; doi:10.1186/s13059-017-1302-3.

1116    48.  Kawakatsu T, Nery JR, Castanon R, Ecker JR. Dynamic DNA methylation

1117        reconfiguration during seed development and germination. Genome Biol. 2017 09

1118        15; doi:10.1186/s13059-017-1251-x.

1119    49.  Lin JY, Le BH, Chen M, Henry KF, Hur J, Hsieh TF et al. Similarity between

1120        soybean and Arabidopsis seed methylomes and loss of non-CG methylation does

1121        not affect seed development. Proc Natl Acad Sci U S A. 2017 11 07;

1122        doi:10.1073/pnas.1716758114.

1123    50.    Bouyer D, Kramdi A, Kassam M, Heese M, Schnittger A, Roudier F et al. DNA

1124           methylation dynamics during early plant life. Genome Biol. 2017 09 25;

1125           doi:10.1186/s13059-017-1313-0.

1126    51.    Ji L, Mathioni SM, Johnson S, Tucker D, Bewick AJ, Do Kim K et al. Genome-wide

1127           reinforcement of DNA methylation occurs during somatic embryogenesis in

1128           soybean. The Plant Cell. 2019; 31: 2315-2331.

1129    52.    Deflorio G, Hein S, Fink S, Spiecker H, Willis Mathew Robert Schwarze F. The

1130           application of wood decay fungi to enhance annual ring detection in three diffuse-

1131           porous hardwoods. Dendrochronologia. 2005; doi:10.1016/j.dendro.2005.02.002.

1132    53.    DeRose JR, Gardner RS. Technique to improve visualization of elusive tree-ring

1133           boundaries in aspen (Populus tremuloides). Tree-Ring Research. 2010; 66: 75-79.

1134    54.    Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-

1135           MEM. arXiv preprint arXiv:13033997. 2013;

1136    55.    McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A et al. The

1137           Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation

1138           DNA sequencing data. Genome research. 2010; 20: 1297-1303.

1139    56.    Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith RK, Hannick LI et al.

1140           Improving the Arabidopsis genome annotation using maximal transcript alignment

1141           assemblies. Nucleic Acids Res. 2003 Oct 01; doi:10.1093/nar/gkg770.

1142    57.    Smit AFA, Hubley R, Green P. RepeatMasker Open-4.0. 2013–2015.

1143           http://www.repeatmasker.org.

1144    58.    Salamov AA, Solovyev VV. Ab initio gene finding in Drosophila genomic DNA.

1145           Genome research. 2000; 10: 516-522.

1146    59.   Slater GSC, Birney E. Automated generation of heuristics for biological sequence

1147         comparison. BMC bioinformatics. 2005; 6: 31.

1148    60.   Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. BRAKER1:

1149         unsupervised RNA-Seq-based genome annotation with GeneMark-ET and

1150         AUGUSTUS. Bioinformatics. 2015; 32: 767-769.

1151    61.   Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N et al. The sequence

1152         alignment/map format and SAMtools. Bioinformatics. 2009;

1153         doi:10.1093/bioinformatics/btp352.

1154    62.   Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L et al. VarScan 2:

1155         somatic mutation and copy number alteration discovery in cancer by exome

1156         sequencing. Genome Res. 2012 Mar; doi:10.1101/gr.129684.111.

1157    63.   Martin ER, Kinnamon DD, Schmidt MA, Powell EH, Zuchner S, Morris RW.

1158         SeqEM: an adaptive genotype-calling approach for next-generation sequencing

1159         studies. Bioinformatics. 2010 Nov 15; doi:10.1093/bioinformatics/btq526.

1160    64.   Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, von Haeseler A et

1161         al. Accurate detection of complex structural variations using single-molecule

1162         sequencing. Nat Methods. 2018 06; doi:10.1038/s41592-018-0001-7.

1163    65.   Benson G. Tandem repeats finder: a program to analyze DNA sequences. Nucleic

1164         Acids Res. 1999 Jan 15; doi:10.1093/nar/27.2.573.

1165    66.   Team RC. R: A Language and Environment for Statistical Computing. Vienna,

1166         Austria: R Foundation for Statistical Computing; 2016. p.

1167    67.   Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G et al.

1168         Integrative genomics viewer. Nature biotechnology. 2011; 29: 24-26.

1169    68.    Urich MA, Nery JR, Lister R, Schmitz RJ, Ecker JR. MethylC-seq library

1170            preparation for base-resolution whole-genome bisulfite sequencing. Nat Protoc.

1171            2015 Mar; doi:10.1038/nprot.2014.114.

1172    69.    Schultz MD, He Y, Whitaker JW, Hariharan M, Mukamel EA, Leung D et al. Human

1173            body epigenome maps reveal noncanonical DNA methylation variation. Nature.

1174            2015 Jul 9; doi:10.1038/nature14465.

1175    70.    Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and

1176            powerful approach to multiple testing. Journal of the Royal statistical society: series

1177            B (Methodological). 1995; 57: 289-300.

1178    71.    Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic

1179            features. Bioinformatics. 2010; 26: 841-842.

1180    72.    Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina

1181            sequence data. Bioinformatics. 2014 Aug 1; doi:10.1093/bioinformatics/btu170.

1182    73.    Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory

1183            requirements. Nat Methods. 2015 Apr; doi:10.1038/nmeth.3317.

1184    74.    Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL.

1185            StringTie enables improved reconstruction of a transcriptome from RNA-seq reads.

1186            Nat Biotechnol. 2015 Mar; doi:10.1038/nbt.3122.

1187