1    **Pangenome of white lupin provides insights into the**

2    **diversity of the species**

3

4    Bárbara Hufnagel[1]*, barbara.hufnagel@supagro.fr

5    Alexandre Soriano[1], alexandre.soriano@supagro.fr

6    Jemma Taylor[2], j.taylor2@kew.org

7    Fanchon Divol[1], fanchon.divol@supagro.fr

8    Magdalena Kroc[3], mkro@igr.poznan.pl

9    Heather Sanders[4], heather.sanders@secure-harvests.com

10    Likawent Yeheyis[5], likawenty@yahoo.com

11    Matthew Nelson[2,6], matthew.nelson@csiro.au

12    Benjamin Péret[1]*, benjamin.peret@supagro.fr

13

14    *1*    *BPMP, Univ Montpellier, CNRS, INRAE, Institut Agro, Montpellier, France*

15    *2*    *Royal Botanic Gardens, Kew, UK*

16    *3*    *Institute of Plant Genetics Polish Academy of Sciences, Poznan, Poland*

17    *4*    *Secure Harvests, Bradford on Avon, UK*

18    *5*    *Amhara Agricultural Research Institute, Bahir Dar, Ethiopia*

19    *6*    *CSIRO, Perth, Australia*

20    *    *Corresponding authors*

21

22

23

24    **ABSTRACT**

25 **Background:** White lupin is an old crop with renewed interest due to its seed high

26 protein content and high nutritional value. Despite a long domestication history in the

27 Mediterranean basin, modern breeding efforts have been fairly scarce. Recent

28 sequencing of its genome has provided tools for further description of genetic

29 resources but detailed characterization is still missing.

30 **Results:** Here, we report the genome sequencing of several accessions that were

31 used to establish a white lupin pangenome. We defined core genes that are present

32 in all individuals and variable genes that are absent in some and may represent a

33 gene pool for stress adaptation. We believe that the identification of novel genes,

34 together with a more comprehensive reference sequence, represents a significant

35 improvement of the white lupin genetic resources. As an example, we used this

36 pangenome to identify selection footprints and to provide a candidate gene for one of

37 the main QTLs associated with late flowering in Ethiopian lupin types. A 686

38 nucleotide deletion was identified in exon 3 of the *LaFTa1* (*Lupinus albus Flowering*

39 *Time a1*) gene that suggests a molecular origin for this trait of importance, defining

40 the need for vernalization in some lupins.

41 **Conclusions:** The white lupin pangenome provides a novel genetic resource to

42 better understand how domestication has shaped the genomic variability amongst

43 this crop. It will be of major importance for breeders to select new breeding traits and

44 incorporate them into new, more efficient and robust cultivars in order to face a

45 growing demand for plant protein sources, notably in Europe.

46

47 **Keywords:** White lupin, pangenome, flowering time, domestication, plant diversity.

48 **BACKGROUND**

49  White lupin (*Lupinus albus* L.) is a pulse whose domestication started about

50 3000 - 4000 years ago in the Mediterranean region [1]. It is cultivated for its seeds

51 that contain high levels of proteins and are used both for food and feed [2]. The wild

52 forms of the species can only be found in the Balkan region and evidence of its

53 earliest use as a green manure and grain crop come from that same region [3]. Early

54 Greek farmers selected larger seeds and white flowers, and presumably soft-

55 seededness (water permeable seeds) was the earliest domestication trait. Greek and

56 Roman literature suggests that seed indehiscence (*i.e.* resistance to pod shattering)

57 had not yet been incorporated by the first century A.D. [4].

58  Wild collections and landraces of white lupin contain high levels of

59 quinolizidine alkaloids that accumulate in the seed, resulting in a bitter taste and

60 possible toxicity. Lysine-derived alkaloids are characteristic of the Genistoids [5–7], a

61 monophyletic basal clade belonging to the Fabaceae family. Traditionally these bitter

62 compounds are removed from white lupin seeds by soaking in water, a practice that

63 is still carried out today across the Mediterranean and Nile regions [1]. However, this

64 is uneconomic on a broad-scale, which motivated the identification of low alkaloid

65 mutants in Germany in the 1930s, aided by advances in chemistry [4]. Modern

66 cultivars of white lupin incorporate low alkaloid genes, hence the term 'sweet' lupins.

67  Breeding efforts have rarely been intensive or sustained over long periods. As

68 a result, white lupin yields remain low and highly variable, in comparison to similar

69 pulses like soybean for which important breeding efforts have been made

70 internationally. Although white lupin cultivation represents a promising crop for

71 Europe, in a political context aiming towards plant protein independence, the lack of

72 well characterised genetic resources has been hampering a fast deployment of white

73 lupin as an alternative crop to soybean imports. The recent sequencing of white lupin

74   genome [8,9] demonstrated a resurgence of interest for this "old" crop. We believe

75   that white lupin intragenomic diversity might reflect the early traces of its slow and

76   sporadic domestication history.

77          Here we report a pangenome for white lupin that reveals important aspects of

78   the species diversity, single nucleotide polymorphisms (SNPs) and gene presence–

79   absence variations (PAVs). We construct a species pangenome consisting of 'core'

80   genes that are present in all individuals and 'variable' (soft-core or shell) genes that

81   are absent in some individuals [10,11]. Building on this comprehensive dataset, we

82   were able to identify a deletion in the QTL region associated with late flowering in

83   Ethiopian white lupins. The deleted gene is a homolog of the FT (Flowering Time)

84   gene, suggesting that this deletion is at the origin of the need for vernalization in

85   these accessions. Our analyses provide new perspectives on white lupin intra-

86   species diversity and domestication history.

87

88   **RESULTS**

89   ***De novo* assembly and pangenome construction**

90          We gathered a set of 39 white lupin accessions, including 25 modern cultivars,

91   10 landraces and 4 wild accessions from 17 countries (Supplementary Table 1).

92   Genome sequence of 15 out of these accessions was available from a previous

93   report [8], whereas 24 accessions have been sequenced within this study to obtain

94   broader species representation. Short-read sequences have been assembled *de*

95   *novo* for each accession (28.5x mean depth, 150 bp pair-end, Supplementary Table

96   2).

97          The *de novo* assembly for each accession produced a total of 14.9 Gb of

98   contigs longer than 500 base pairs (bp) with an N50 value (the minimum contig

99    length needed to cover 50% of the assembly) of 24,475 bp. These *de novo*

100   assemblies showed a mean complete BUSCOs score of 96.3%, a value similar to the

101   AMIGA reference genome (97.7%). Assembly completeness assessed by BUSCO

102   was higher than 91.7%, for all accessions and in case of three accessions (Kiev,

103   P27174 and Magnus) the score was similar to the reference genome (Fig. 1a).

104        The pangenome was built using the iterative mapping and assembly

105   approach, in a similar strategy used to generate the *Brassica oleracea* [10] and

106   tomato [12] pangenomes. The assembly of *L. albus* reference genome based on

107   AMIGA accession is 450,972,408 bp size with 38,258 predicted protein-coding genes

108   [8]. All *de novo* assembled contigs were compared with the reference genome to

109   identify previously unknown sequences. A total of 270 Mb of nonreference sequence

110   with identity <90% to the reference genome was obtained. After pangenome

111   construction and removal of contaminants and overly repetitive sequences, we

112   assembled an additional 3,663 scaffolds, with a length greater than 2,000 bp, for a

113   total length of 11,733,253 bp. Using a threshold of a minimum 10x coverage, we

114   identified 178 newly predicted protein-coding genes, among which 61 could be

115   annotated with gene ontology (GO) terms or Pfam domains (Supplementary Dataset

116   1). The white lupin pangenome, including reference and nonreference genome

117   sequences, had a total size of 462,705,661 bp and contained 38,446 protein-coding

118   genes. The total size of the constructed pangenome is compatible with nuclear DNA

119   content estimates based on flow cytometry [13] which suggests that it represents the

120   complete genome sequence of the species.  We added to the White Lupin Genome

121   portal (www.whitelupin.fr) dedicated user-friendly tools for the exploitation of the

122   pangenome, such as a BLAST tool for individual accessions, download of specific

123   regions of accessions and a genome browser mapping all the variants.

124

## Core and variable genes

126       The presence or absence of each protein-coding gene was predicted for each

127 of the 39 accessions based on the mapping of reads from each accession to the

128 pangenome assembly using SGSGeneLoss [14]. Likewise to other plants

129 pangenomes [10,12,15–18], we categorized genes in the white lupin pangenome

130 according to their presence frequencies, using Markov clustering in the

131 GET_HOMOLOGUES-EST pipeline [19]. The majority of the genes, 32,068 (78.5%),

132 are core genes shared by all the 39 accessions; 6,046 soft-core (14.8%), being

133 absent in at least one accession; and 8,776 (21.4%) are shell, present in 2-38

134 accessions (Fig. 1b). The size of the pangenome expanded with each additional

135 accession to 38,443 genes, and extrapolation leads to a predicted pangenome size

136 of 40,844 +/-289 genes (Figure 1b).

137

## Single-nucleotide polymorphism detection and annotation

139       To capture and broadly characterize white lupin diversity we applied a strict

140 SNP identification pipeline, using GATK 4.1.0.0. A total of 9,442,876 raw SNPs were

141 identified, 806,740 of which were recognized in the newly assembled pangenome

142 scaffolds.  After filtering, 3,527,872 SNPs were retained in the 39 accessions,

143 corresponding to a rate of 1 variant every 127 bp (Supplementary Figure S1). The

144 majority (85.8%) of the high-quality variants are SNPs (3,027,761) and the other

145 501,111 variants detected are insertions and deletions (Fig. 1c – blue). Most variants

146 (59.3%) are distributed on intergenic regions, 7.1% are within introns and only 1.9%

147 (96,576) of the variants are located in exons (Fig. 1c – red). From the variants

148 present in the CDS region 4,725 showed potentially large effects by causing start

149    codon changes, premature stop codons or elongated transcripts, and 50,478 are

150    considered to produce a moderate effect by leading to codon changes in annotated

151    genes. The frequency of these missense SNPs in the core gene set was one each

152    4.26 kb, which was lower than the variable gene set, with a rate of one for 1.84 kb.

153    The rest of the variants lead to synonymous changes in proteins (low effect variants)

154    or modifiers, causing changes outside the coding regions (Fig. 1c − green).

155    Collectively, this comprehensive dataset of the genome variation of white lupin

156    provides a resource for biology and breeding of this species.

157

**Population structure**

159    To establish a phylogenetic benchmark for the analysis of the pangenome, we

160    built a consensus maximum likelihood tree (Fig. 2a) to infer the phylogenetic

161    relationships for these *L. albus* accessions using the complete set of 3.5 M SNPs

162    described above. This phylogenetic tree clustering supported six clades, which

163    exhibited distinctive geographic origin and distinctive botanical features. In the Type

164    1 are grouped accessions with early flowering traits, including the Chilean

165    agrogeotypes, and German and French accessions used in breeding programs. This

166    group also included the widely used cv. Kiev Mutant, which was generated by

167    mutagenesis techniques with the intention to induce early flowering, and the

168    accessions that are derived forms of it (Primorsky and Dieta, [3]). Type 2 is also

169    composed by accessions with early flowering, anumber of which have characteristics

170    of Polish agroecotypes described by Kurlovich [3] and are adapted to grow in Eastern

171    Europe. One of the most representative accessions of this group is the cv. Kalina [3],

172    an old cultivar created in the Polish breeding program sharing similar genetic

173    background with the broadly used cultivar Start. Interestingly, Start is reported to

7

174    carry different early-flowering genes than Kiev Mutant [20]. Type 2 also comprises

175    two landraces with from Syria and Israel/Palestine. Type 3 encompasses autumn-

176    sown genotypes with strong vernalisation requirement and dwarf phenotype from the

177    French breeding program, and the Algerian landrace ALB01. Algerian landraces are

178    also reported to have a strong need of vernalization [3]. Type 4 comprises landraces

179    from Iberian and Apennine Peninsula together with the described thermoneutral

180    cultivars (*i.e.* Neutra, [2]). Type 5 is composed only by Ethiopian landraces and Wild

181    group is composed by the four "*graecus*-type" accessions of the panel, all presenting

182    small black-speckled seeds and non-domesticated traits (hard seeds and shattering

183    pods).

184         We examined genetic structure by performing a Bayesian model-based

185    clustering analysis and found that the six population groups matched the maximum-

186    likelihood tree (Fig. 2b). This presented evidence of significant admixture in some

187    lines and a weak population structure, a pattern already seen in other studies of *L.*

188    *albus* [21]. This weak population structure is also seen through the population-

189    differentiation statistic ($F_{ST}$). The $F_{ST}$ value between all six groups were 0.27,

190    however, $F_{ST}$ between Type 1 and Type 2 are low as 0.086, and Type 4 and Wild

191    have an $F_{ST}$ of 0.092. Indeed, regarding the Bayesian model, in scenarios dividing

192    the accessions in 4 or 5 sub-populations (Fig 2b, K=4 and K=5), accessions from

193    Type 4 are merged with the Wild group. On the other hand, Type 5 showed a strong

194    differentiation from the other groups, with $F_{ST}$ values ranging from 0.34 to 0.46, with

195    Type 4 and Type 3, respectively, which is corroborating with previous studies [22].

196    Principal component analysis reinforced the similarity among some groups (Fig. 2c).

197    The two first principal components explain 65.9% of genotypic variance and it is

198    highlighting the overlap among certain groups, in particular, Type 1 and Type 2.

199        Differentiation of genetic diversity between the 6 groups was investigated

200    further through analysis of decay of linkage disequilibrium (LD, Fig. 2d). The decay of

201    LD with physical distance between SNPs to half of the maximum values occurred at

202    3.85 Kb ($r^2 = 0.38$), consistent with a high level of diversity and partially outcrossing

203    mode of reproduction in this species [23]. Type 4 group also showed a fast LD decay

204    of 5.7 Kb ($r^2 = 0.40$) and Type 1-3 groups have an average LD decay of 10.5 Kb.

205    Wild group showed a slower LD decay (38.1 Kb, $r^2 = 0.39$) when comparing with the

206    other white lupin groups, presumably an effect of the small number of wild

207    accessions in the analysis. Nevertheless, these LD decay levels can still be

208    considered fast compared with other plant species, for example rice (~75–150 Kb,

209    [24]), soybean (~340-580 Kb, [25]) or wheat (~7-12.4 Mb, [26]), also self-pollinated

210    crops. The Type 5 group (Ethiopian landraces) only reached half of its LD decay after

211    1.5 Mb, reinforcing the high similarity of its accessions and a possible genetic

212    isolation of this group [21]. The average nucleotide diversity π per site [27] showed

213    that diversity was five times lower in Type 5 group (π = 0.068) compared to the

214    general nucleotide diversity (π = 0.372). While the Wild group, although is also

215    composed of only four accessions, showed a nucleotide diversity π = 0.402.

216

**Protein-coding genes presence and absence characterization**

218        Presence and absence variants (PAVs) are an important type of structural

219    variation and play an important role in shaping genomes, therefore contributing to

220    phenotypic diversity [28]. The construction of a white lupin pangenome allowed

221    identification of 1195 PAVs, representing protein-coding genes that are absent in at

222    least one of the accessions, being 1132 genes from the reference genome and 63

223    from the newly identified genes (Supplementary Dataset 2-3). We further examined if

224     the phylogenetic groups have an influence in the number of PAVs and if the PAVs

225     are homogeneous within the groups (Fig. 3a-c). The wild accessions have a

226     significatively higher number of newly identified genes, with the accessions

227     GRAECUS and GR38 only missing 4 of them. The four wild accessions share 157

228     out of the 178 new-identified genes in the pangenome (Fig. 3a).

229          The number of missing genes within individual genomes ranges from 45

230     (AMIGA – Type 1) to 348 genes (GRC5262B – Wild). Each group shares a median of

231     31 common lost genes amongst all its accessions and a total of 103 genes are

232     absent in at least one accession of each group (Fig. 3b). There are 137 genes that

233     have been exclusively lost within accessions of the Wild group, however only 30

234     genes are shared among all the *graecus* accessions. On the other hand, genomes of

235     Ethiopian landraces (Type 5) share a total of 118 common missing genes, amongst

236     39 are unique for this group. Remarkably, for this group there is a concentration of

237     lost genes on Chr17. This includes a set of 9 tandem duplicated genes covering a

238     region of 120 Kb (Supplementary Fig. 2). They are annotated as "Putative ferric-

239     chelate reductase (NADH)" homologs of *Arabidopsis* gene *FRO2*, known for its role

240     of iron uptake by the roots under stress condition [29].

241          Checking the position of the PAVs on the chromosomes we could identify

242     some peculiarity regarding the PAVs within the groups. For example, on Chr13 there

243     is a concentration of PAVs in the region of 5-10 Mb that are missing from most

244     accessions of Type 2-5 and Wild, but are present in the genomes of most Type 1

245     members. Similar pattern happens in the 3.6-6.4 Mb region of Chr04. Chr23 has the

246     highest number of PAVs (78), a common feature of all the groups.

247          Functional analysis of PAVs suggests enrichment of GO terms as "integral

248     component of membrane" (GO:0016021) and "oxidation-reduction process"

249    (GO:0055114) (Fig. 3d, supplementary Fig. 3 and Supplementary Data 2). These

250    suggest an enrichment of genes and gene families coding for membrane receptors

251    proteins or membrane transporters. Other GO terms suggest that some of the genes

252    may be involved in cell wall remodeling ("cell wall" - GO:0005618; "cell wall

253    organization" - GO:0071555). Genes with these functions are frequently linked to

254    biotic and abiotic stress responses [30,31]. PAV genes related to abiotic and biotic

255    stress responses have been observed in several plant species [15,17,18,32–35] and

256    these may reflect the evolution for adaptive traits related for each agroecotype.

257    Moreover, the presence/absence of these stress-response related genes may also

258    be partially due to whole-genome triplication event on white lupin genome [8], which

259    caused an overlapping roles in various loci.

260

261    **Footprints of selection and alleles identification in candidate genes**

262    To demonstrate the power of white lupin pangenome to address basic

263    research questions, we used it to detect possible footprints of selection and to

264    identify alleles in candidate genes underlying major QTLs. Firstly, to examine

265    potential selective signals during white lupin domestication and breeding, we

266    scanned white lupin genome searching for regions with marked reductions in

267    nucleotide diversity (Fig. 4).

268    The domestication and breeding efforts in white lupin have focused in

269    searching for accessions with reduced seed alkaloid content, reduced time to flower

270    as well as excessive indeterminate branching. Therefore, we combined Type 1 and

271    Type 2 accessions, that are spring types and went to a more intense breeding

272    process and compared them with Type 3 and Type 4 accessions, that are winter

273    types (Fig. 4a). A selective sweep affecting only the spring white lupin accessions

274    would be expected to leave a typical low-polymorphism and high-divergence signal

275    around the region of the selected genes. We measured the sweep on the nucleotide

276    diversity (π value [27]), by comparing the two groups (πWinter/πSpring) over 250-kb

277    windows. We identified 167 putative selection sweeps associated to the breeding of

278    the spring accessions (πWinter/πSpring > 2.101). We observed that some of the

279    peaks co-localized with previously reported white lupin QTLs for flowering time and

280    alkaloid content [36,37]. The same pattern was observed when checking for the

281    divergence of the gene pool between these two groups along the chromosomes (Fst,

282    Supplementary Fig. 4a).

283    Interestingly, other peaks with higher sweeps of diversity are present,

284    indicating that other genomic regions may be implicated with these traits and may

285    carry other important genes of these pathways. Furthermore, they highlight specific

286    genomic regions of spring accessions that have been selected during domestication

287    and breeding. For instance, we checked for orthologs of domestication genes from

288    the close relative narrow-leafed lupin and found that the gene Lalb_Chr12g0203121,

289    a homolog of a candidate gene for the reduced pod shattering locus *tardus*

290    (Lup002448, [38]) , is co-localized with a sweep peak on Chr12.

291    The reported white lupin QTLs were identified in a recombinant inbred line

292    (RIL) mapping population derived from the cross between Kiev Mutant (Type 1) and

293    the Ethiopian landrace P27174 (Type 5). Thereupon, we checked the sweep of

294    diversity between all accessions compared to Ethiopian accessions, T5

295    (πGeneral/πT5, Fig. 4b) and identified 84 sweep peaks (πGeneral/πT5 > 83.97).  A

296    similar trend of co-localization of the QTL peaks were observed, with steep peaks

297    around the QTL regions. Interestingly, the region corresponding the QTL *pauper* did

298    not show a peak, being far below the statistical significance threshold. This indicates

299    that the two groups have similar level of nucleotide diversity in this region

300    (Supplementary Fig. 4b). It can be explained by the above-mentioned similarities

301    among accessions of group T5 and that many of the modern accessions carry the

302    low alkaloid alleles for the pauper region. However, although this region showed a

303    similar nucleotide diversity between the two groups, it presented a high genetic

304    variance, with a median FST of 0.94 for the region (Fig. 4c).

305          In another approach to demonstrate the power of white lupin pangenome, we

306    used its assembly to identify a candidate gene underlying a major QTL and describe

307    the associated allelic diversity. Chromosome 2 is the location of an important QTL

308    associated with early flowering white lupins. We used the protein sequences of

309    *Lupinus angustifolius* that have been previously mapped in syntenic regions of the

310    these QTLs [39] to perform an homology search against the pangenome. we

311    identified the gene *LaFTa1* (Lalb_Chr02g0156991), a homolog of the gene *LanFTa1*

312    (Lup21189) mapped on this QTL region. The white lupin *LaFTa1*

313    (Lalb_Chr02g0156991) was annotated as "Putative phosphatidylethanolamine-

314    binding protein" (PEBP) in the reference genome. The FT proteins belonging to the

315    PEBP family are the key control points of the flowering time in plants. The *LaFTa1*

316    gene presented a deletion of 686 on the third intron that is present only on Type 5

317    accessions, that have late flowering phenotypes (Fig. 5a-b and Supplementary Fig.

318    5). Indeed, one of the parents of this QTL mapping population belongs to this group

319    (P27174). It is reported that changes in FT promoter and introns can alter FT

320    expression in response to photoperiod and vernalization, and consequently, induce

321    flowering [40]. This suggests that the identified *LaFTa1* is the gene underlying this

322    QTL and that this deletion on the intron of Type 5 accessions may be contributing for

323    the late flowering pattern of this group.

324

325    **DISCUSSION**

326         A pangenome is a complete set of genes for a species, including core genes

327    which are present in all individuals, and variable genes which are absent in one or

328    more individuals [41]. We generated a *de novo* assembly for 38 white lupin

329    accessions and, taking advantage of a good reference assembly for the species [8],

330    we constructed a *L. albus* pangenome by iteratively and randomly sampling these

331    sequenced accessions. This dataset is representative of the diversity of the species,

332    containing wild accessions, landraces and cultivars of white lupin from across their

333    respective distributions. As a result, we estimate that this white lupin pangenome

334    assembly effectively encompasses the complete sequence for the genome of the

335    species, with 462,7Mb sequence and containing 38,446 protein-coding genes. The

336    finding that 21.5% of genes in the pangenome exhibit varying degrees of genic

337    presence/absence variants (PAVs) highlights the diverse genetic feature of white

338    lupin and the significant improvement of the reference genome, by including genomic

339    information of other accessions and discovery of new genes. Remarkably, the white

340    lupin pangenome showed a high content of core genes (78.5%), as compared with

341    other plant species as tomato (74.2%, [34]), *Arabidopsis thaliana* (70%,[19]),  bread

342    wheat (64%, [42]), sesame (58%, [16]) and wild soybean (49%, [43]), which might be

343    a reflection of its domestication history and modest breeding efforts to date.

344         The domestication of white lupin started during Bronze Age [4], and the

345    ancestral history of this species is different than other major crops such as rice,

346    maize, sorghum, tomato, and soybeans, which are more ancient [44]. The early

347    cultivated forms have the same Mediterranean distribution that its wild ancestor types

348    (*graecus*), which led to small adaptation or selection differences. *L. albus*

349  domestication was slow with potentially centuries between acquisition of each

350  domestication trait, which may explain why there is not a more pronounced genetic

351  differentiation between wild, landrace and cultivated types [45]. This is echoed in the

352  lack of population structure presented within these accessions and in the low LD

353  extent, which generally reduce the diversity and change allele frequencies either to

354  fixation or intermediate frequencies  [46]. Despite being a largely self-pollinating crop

355  (with an out-crossing rate reported as 8–10 %, [23]), white lupin showed a

356  remarkably low LD extent (< 4kb), even lower than the wild population of its relative,

357  narrow-leafed lupin, that showed a decay of LD after 19.01 Kb [47]. One distinction

358  between these closely related species is that narrow-leafed lupin is almost

359  exclusively self-pollinating and so the modest levels of outcrossing in white lupin may

360  be a key factor governing the differences in LD between these two species. Having a

361  low LD and weak population structure together mean that association mapping is

362  likely to be particularly powerful in white lupin, in contrast to the more highly

363  structured and high LD species narrow-leafed lupin, where association studies have

364  so far proved rather weak [47,48].

365     Type 5 accessions, from Ethiopia, are the only group which shown a strong

366  genetic differentiation from the others, with $F_{ST}$ values higher than 0.3. Such a distinct

367  separation is an evidence that the Ethiopian accessions have evolved in isolation and

368  the genetic differences are probably due to ancient founder effects. The differences

369  of Type 5 group are also highlighted by the PAVs. Together with the Wild group,

370  Ethiopian landraces carry most of the new identified genes and also miss a large

371  number of genes of the reference genome (Fig. 3). Moreover, it is a highly

372  homogeneous group, with all accessions sharing a large number of these lost genes.

373  The loss of these genes is probably an adaptive response for the local environment.

374    For instance, the loss of the nine tandem duplicated homologue *AtFRO2* on Chr17

375    might be an adaptive response to highland Ethiopian soils that are iron-rich [49]. A

376    more detailed look into the PAVs among the different groups may be useful to better

377    understand their specificities.

378        Our analysis brings a high resolution to the within-species diversity. Using the

379    pangenome dataset, we performed genome-wide comparisons of the assemblies,

380    enabling the characterization of more than 3 million complex variants, including many

381    large-effect coding variants which should be helpful in pinpointing causal variations in

382    QTLs for important traits and in future genome-wide association studies. In particular,

383    our study demonstrated that 4,725 genes were found to contain important coding

384    variation in at least one accession and might have important biological functions

385    underlying the variation of complex traits.

386        We wanted to demonstrate how a pangenome can be a useful tool to identify

387    allelic differences that are responsible for phenotypic variation. By performing a

388    genome wide analysis, we detected that nucleotide diversity were quite variable

389    across the genome. The efforts of breeding in white lupin have been focused in

390    combining of domestication traits such as soft and white seeds and reduced pod-

391    shattering, which were already available from ancient times, with that of reduced

392    alkaloids, increased yield and the reduction of flowering time and excessive

393    branching [45]. Looking for differences in nucleotide diversity across the genome

394    amongst breeding accessions and comparing with landraces/wild accessions, we

395    could detect some peaks of sweep of diversity. In these peaks there is an important

396    decrease of nucleotide diversity within the breeding lines and they represent marks of

397    selection (Fig. 4). In these regions were also detect a high divergence between the

398    two gene pools (Fst). However, although the sweeps of diversity co-localize with

399   some identified QTLs for flowering time and low alkaloid content, there are other

400   higher peaks along the chromosomes. These regions should be explored in order to

401   find genes underlying phenotypic traits that have been selected directly or indirectly

402   during domestication and breeding of white lupin. For instance, white lupin is known

403   for thriving in soils with low nutrient availability by producing specialized root

404   structures called cluster roots [50]. In a previous work, we demonstrated that the

405   breeding accessions have an earlier establishment of the root system through lateral

406   and cluster root formation that was indirectly selected [8]. By looking closer in these

407   chromosome regions with low nucleotide diversity and high genetic differentiation we

408   might be able to find genes with important roles in the root architecture of white lupin.

409   Hence, integration of the information from studies of gene function and the high

410   density of variants described in this pangenome can provide a complementary

411   approach to forward genetic studies and can contribute to develop the research and

412   breeding of white lupin.

413

414   **CONCLUSION**

415       In summary, the white lupin pangenome comprises a wealth of information on

416   genetic variation that has yet to be fully exploited by researchers and breeders.

417   Although a there is large collection of white lupin accessions available in genebanks

418   worldwide, they barely have been explored and genetically characterized. This

419   pangenome represents a comprehensive and important resource to facilitate the

420   exploration of white lupin as a legume model for future functional studies and

421   molecular breeding.

422

423

424 **METHODS**

425 **Genome sequences of white lupin accessions**

426 We retrieved the genome sequencing data of 15 white lupin accessions that

427 were published previously [8], including 11 modern cultivars, 1 landrace and 2 wild

428 relatives. They were sequenced using Illumina technology using paired-end 2 × 150

429 bp short-reads with average sequencing depth of 45.99×. It included Illumina genome

430 data of 64.47x depth for the reference cultivar "AMIGA". Genome sequences of

431 additional 24 accessions were generated here, including 12 modern cultivars, 9

432 landraces and 2 wild relatives. Young leaves of individual plants were used to extract

433 genomic DNA of each accession using the QIAGEN DNeasy Plant Mini kit following

434 the supplier's recommendations. The accessions were sequenced using Illumina

435 technology using paired-end 2 × 150 bp short-reads (Macrogen, South Korea). It was

436 generated a total of 196.85 Gb of data with average sequencing depth of 19.1x.

437 (Supplementary Table 2).

438

439 **De novo genome assembly and pangenome construction**

440 Reads were processed to trim adapters and low-quality sequences using

441 Cutadapt 1.15 [51] with parameters '--pair-filter=any -q20,20 -m 35' and the forward

442 and reverse Illumina TruSeq Adapters. The final high-quality cleaned Illumina reads

443 from each sample were *de novo* assembled using Spades 3.13.0 [52] with k-mer size

444 of 21,33,55,77,99,121. The assembled contigs were then aligned to the white lupin

445 reference genome [8] (GenBank accession no.: WOCE00000000,

446 http://www.whitelupin.fr.), using the steps 7 and 8 of the EUPAN Pipeline [53], in

447 order to extract contigs that were not aligning to the reference. Then, redundancy in

448 the extracted contigs has been reduced using CD-hit 4.8.1 with default parameters.

449    The resulting contigs were then search against the NCBI nt nucleotide database

450    using blastn 2.10 [54]. Sequences with best hits from outside the Eudicots, or

451    covered by known plant mitochondrial or chloroplast genomes, were possible

452    contaminations and were therefore removed.

453

454    **Annotation of the white lupin pangenome**

455         A custom repeat library was constructed by screening the pangenome and the

456    white         lupin         reference         genome         using         RepeatModeler

457    (http://www.repeatmasker.org/RepeatModeler/), and used to screen the nonreference

458    genome       to       identify       repeat       sequences       using       RepeatMasker

459    (http://www.repeatmasker. org/). Contigs with more than 98% of repetitive sequences

460    were removed from the annotation pipeline. Protein-coding genes were predicted

461    from nonreference genome using MAKER2 [55]. *Ab initio* gene prediction was

462    performed using Augustus [56] and SNAP [57]. Augustus [58] has been previously

463    trained for white lupin as described in the documentation, and SNAP was trained for

464    two rounds based on already assembled transcriptome of white lupin, as described in

465    maker2 documentation. In addition, protein sequences of white lupin, *Medicago*

466    *truncatula* and the Viridiplantae subset of Swissprot were used as evidence. Finally,

467    gene predictions based on *ab initio* approaches, and transcript and protein evidence

468    were integrated using the MAKER2 pipeline. A set of high-confidence gene models

469    supported by transcript and/or protein evidence were generated by MAKER2. In

470    order to remove possible remaining contamination, all high confidence maker

471    generated protein sequences were aligned against the nr databses, and sequences

472    with best hits from outside Eudicots or with best hit inside chloroplastic and

473    mitochondrial sequences were removed. Genes that matched white lupin reference

474    sequences were also removed the same way.

475       In parallel, contigs with a length superior to 2Kb from the whole assembly of

476    the 39 lupin accession were annotated using the Egnep 1.5.1 pipeline [59].

477    RepeatMasker was used to detect and remove contigs constitute by more than 98%

478    of known repeat sequences based on the previously built white lupin repetitive

479    element sequences database. The white lupin transcriptome [8]was used as ESTs

480    evidence, using a minimum identity percentage of 95%, along with the proteome of

481    white lupin, *Medicago truncatula*, and the Viridiplantae subset of the swissprot

482    database, with weight of 0.4, 0.3 and 0.3 respectively. Resulting predicted proteins

483    were search against REXdb and repbase in order to remove possible transposable

484    elements. The resulting genes prediction were again scan with repeat-masker, and

485    genes composed of more than 90% of detected repetitive sequences were removed

486    from further analyses in order to control false positive.

487

488    **Gene presence/absence variation and pangenome modeling**

489       Reads were processed to trim adapters and low-quality sequences using

490    Cutadapt 1.15 [51] with parameters '--pair-filter=any -q20,20 -m 35' and the forward

491    and reverse Illumina TruSeq adapters. Resulting high quality reads were then aligned

492    to the pangenome using BWA-MEM [60] with default parameters. Picard tools was

493    used to remove possible PCR and optical duplicates, and reads considered as not

494    properly paired were removed using samtools view. The presence or absence of

495    each gene in each accession was determined using SGSGeneLoss [14] . In brief, for

496    a given gene in a given accession, if less than 10% of its exon regions were covered

497    by at least five reads (minCov = 5, lostCutoff = 0.1), this gene was treated as absent

498    in that accession, otherwise it was considered present. The parameters used for the

499    new set of gene discovered in the pangenome were different:  minCov = 10 and

500    lostCutoff = 0.8. For more precise pangenome studies taking into account all the

501    genes discovered in all the different varieties, GET_HOMOLOGUES_EST  was used

502    on the whole CDS and proteome of the whole 39 varieties with parameters "-R

503    123545 -P -M -c -z -A -t 2" to detect clusters of genes shared by at least two

504    varieties.

505

506    **SNP discovery and annotation**

507    Cutadapt [61] was used to remove Illumina Truseq adapters from the

508    sequencing data and to remove bases with a quality score lower than 30, in both 5'

509    and 3' end of the reads. Reads with a length lower than 35 were discarded. We then

510    used BWA-MEM version 0.7.17 [60] to map the resequencing reads from all 39

511    genotypes to the white lupin reference genome. PCR and Optical duplicates were

512    detected and removed using Picard Tools. After that, GATK 4 HaplotypeCaller tool

513    was used in emit-ref-confidence GVCF mode to produce one gvcf file per sample.

514    These files were merged using GATK Combine GVCFs. Finally, GATK

515    GenotypeGVCFs was used to produce a vcf file containing variants from all the 39

516    samples. This identified a total of 9,442,876 SNPs/indel. After filtering for minimum

517    allele frequency of 0.15 and heterozygosity frequency of 0–0.2, 3,527,872 SNPs

518    were retained for further analysis.

519

520    **Evolutionary analysis**

521    A maximum-likelihood phylogenetic tree was constructed based on 3,121,673

522    parsimony-informative SNPs with 1,000 bootstraps using IQ-TREE [62] using

523  ModelFinder [63] option. Then, a phylogenetic tree was prepared using the iTOL v

524  4.3 [64].

525  Population structure based on the same set of SNPs was investigated using

526  STRUCTURE [65]. Thirty independent runs for each K from 1 to 15 were performed

527  with an admixture model at 50,000 Markov chain Monte Carlo (MCMC) iterations and

528  a 10,000 burn-in period. Principal component analysis using this SNP dataset was

529  performed using the function "princomp" in R (http://www.R-project.org/).  The linkage

530  disequilibrium (LD) pattern was computed using PopLDdecay v3.40 [66]. LD decay

531  was measured on the basis of the $r^2$ value and the corresponding distance between

532  two given SNPs.

533

534  **Selective sweep analyses**

535  To detect genomic regions affected by domestication we used the same set of

536  3,121,673 SNPs using Tassel [67]. The level of genetic diversity (π) was measured

537  with a window size of 2000 SNPs and a step size of the same length, generating

538  windows of approximately 250-kb. Genome regions affected by selection or

539  domestication should have substantially lower diversity in spring white lupin (Types 1

540  and 2, πSpring) than the diversity in winter accession (Type 3 and 4, πWinter) and

541  Ethiopian accessions (πT5). Windows with the top 10% highest ratios of

542  πWinter/πSpring (≥2.101) or πGeneral/πCA (≥83.969) were selected as candidate

543  selection and domestication sweeps. The PopGenome package [68] in R with its

544  sliding window method was used to calculate the interpopulation differentiation, FST.

545  Using a set of 40 k random good-quality SNPs evenly distributed along the 25

546  chromosomes, we calculated nonoverlapping sliding-windows of 10 SNPs each.

547

548 **REFERENCES**

549 1. Taylor JL, De Angelis G, Nelson MN. How Have Narrow-Leafed Lupin Genomic

550 Resources Enhanced Our Understanding of Lupin Domestication? Springer, Cham;

551 2020. p. 95–108.

552 2. Wolko B, Clements JC, Naganowska B, Nelson M, Hua'an Y. Lupinus. Kole C,

553 editor. Wild Crop Relat. Genomic Breed. Resour. Legum. Crop. Forages. Berlin,

554 Heidelberg: Springer Berlin Heidelberg; 2011.

555 3. Kurlovich BS. Lupins: Geography, Classification, Genetic Resources and

556 Breeding. Publishing House "Intan"; 2002.

557 4. Gladstones JS. Distribution, origin. taxonomy, history and importance. Lupins as

558 Crop Plants Biol Prod Uti-lization Gladstones JS, Atkins C, Hamblin J(eds) CAB Int

559 Oxon, New York. 1998. p. 1–39.

560 5. Kinghorn AD, Hussain RA, Robbins EF, Balandrin MF, Stirton CH, Evans S V.

561 Alkaloid distribution in seeds of Ormosia, Pericopsis and Haplormosia.

562 Phytochemistry. 1988;27:439–44.

563 6. van Wyk B-E. The value of chemosystematics in clarifying relationships in the

564 genistoid tribes of papilionoid legumes. Biochem Syst Ecol. 2003;31:875–84.

565 7. Wink M, Mohamed GIA. Evolution of chemical defense traits in the Leguminosae:

566 mapping of distribution patterns of secondary metabolites on a molecular phylogeny

567 inferred from nucleotide sequences of the rbcL gene. Biochem Syst Ecol.

568 2003;31:897–917.

569 8. Hufnagel B, Marques A, Soriano A, Marquès L, Divol F, Doumas P, et al. High-

570 quality genome sequence of white lupin provides insight into soil exploration and

571 seed quality. Nat Commun. Springer US; 2020;11:492.

572 9. Xu W, Zhang Q, Yuan W, Xu F, Muhammad Aslam M, Miao R, et al. The genome

573 evolution and low-phosphorus adaptation in white lupin. Nat Commun. Springer US;

574 2020;11:1069.

575 10. Golicz AA, Bayer PE, Barker GC, Edger PP, Kim H, Martinez PA, et al. The

576 pangenome of an agronomically important crop plant Brassica oleracea. Nat

577 Commun. Nature Publishing Group; 2016;7:13390.

578 11. Vernikos G, Medini D, Riley DR, Tettelin H. Ten years of pan-genome analyses.

579 Curr Opin Microbiol. 2015;23:148–54.

580 12. Gao L, Gonda I, Sun H, Ma Q, Bao K, Tieman DM, et al. The tomato pan-

581 genome uncovers new genes and a rare allele regulating fruit flavor. Nat Genet.

582 Springer US; 2019;51:1044–51.

583 13. Naganowska B, Wolko B, Śliwińska E, Kaczmarek Z. Nuclear DNA content

584 variation and species relationships in the genus Lupinus (Fabaceae). Ann Bot.

585 2003;92:349–55.

586 14. Golicz AA, Martinez PA, Zander M, Patel DA, Van De Wouw AP, Visendi P, et al.

587 Gene loss in the fungal canola pathogen Leptosphaeria maculans. Funct Integr

588 Genomics. 2015;15:189–96.

589 15. Gordon SP, Contreras-Moreira B, Woods DP, Des Marais DL, Burgess D, Shu S,

590 et al. Extensive gene content variation in the Brachypodium distachyon pan-genome

591 correlates with population structure. Nat Commun. Springer US; 2017;8:2184.

592 16. Yu J, Golicz AA, Lu K, Dossa K, Zhang Y, Chen J, et al. Insight into the evolution

593 and functional characteristics of the pan-genome assembly from sesame landraces

594 and modern cultivars. Plant Biotechnol J. 2019;17:881–92.

595 17. Zhao J, Bayer PE, Ruperao P, Saxena RK, Khan AW, Golicz AA, et al. Trait

596 associations in the pangenome of pigeon pea ( Cajanus cajan ). Plant Biotechnol J.

597 2020;pbi.13354.

598    18. Montenegro JD, Golicz AA, Bayer PE, Hurgobin B, Lee H, Chan C-KK, et al. The

599    pangenome of hexaploid bread wheat. Plant J. 2017;90:1007–13.

600    19. Contreras-Moreira B, Cantalapiedra CP, García-Pereira MJ, Gordon SP, Vogel

601    JP, Igartua E, et al. Analysis of Plant Pan-Genomes and Transcriptomes with

602    GET_HOMOLOGUES-EST, a Clustering Solution for Sequences of the Same

603    Species. Front Plant Sci. 2017;8:1–16.

604    20. Adhikari K, Buirchell B, Yan G, Sweetingham M. Two complementary dominant

605    genes control flowering time in albus lupin (Lupinus albus L.). Plant Breed.

606    2011;130:496–9.

607    21. Raman R, Cowley RB, Raman H, Luckett DJ. Analyses Using SSR and DArT

608    Molecular    Markers    Reveal    that    Ethiopian    Accessions    of    White    Lupin

609    (&amp;lt;i&amp;gt;Lupinus    albus&amp;lt;/i&amp;gt;    L.)    Represent    a    Unique

610    Genepool. Open J Genet. 2014;04:87–98.

611    22. Raman R, Cowley RB, Raman H, Luckett DJ. Analyses Using SSR and DArT

612    molecular markers reveal that Ethiopian accessions of white lupin (Lupinus albus L.)

613    represent a unique genepool. Open J Genet. 2014;4:87–98.

614    23. Green A, Brown A, Oram R. Determination of outcrossing rate in a breeding

615    population of Lupinus albus L. (White Lupin). Plant Breed. 1980;84:181–91.

616    24. Mather K a, Caicedo AL, Polato NR, Olsen KM, McCouch S, Purugganan MD.

617    The    extent    of    linkage    disequilibrium    in    rice    (Oryza    sativa    L.).    Genetics.

618    2007;177:2223–32.

619    25. Hyten DL, Choi I-Y, Song Q, Shoemaker RC, Nelson RL, Costa JM, et al. Highly

620    Variable Patterns of Linkage Disequilibrium in Multiple Soybean Populations.

621    Genetics. 2007;175:1937–44.

622    26. Molero G, Joynson R, Pinera-Chavez FJ, Gardiner L, Rivera-Amado C, Hall A, et

623    al. Elucidating the genetic basis of biomass accumulation and radiation use efficiency

624    in spring wheat and its role in yield potential. Plant Biotechnol J. 2019;17:1276–88.

625    27. Tajima F. Evolutionary relationship of DNA sequences in finite populations.

626    Genetics. 1983;105:437–60.

627    28. Marroni F, Pinosio S, Morgante M. Structural variation and genome complexity: is

628    dispensable really dispensable? Curr Opin Plant Biol. Elsevier Current Trends;

629    2014;18:31–6.

630    29. Connolly EL, Campbell NH, Grotz N, Prichard CL, Guerinot M Lou.

631    Overexpression of the FRO2 Ferric Chelate Reductase Confers Tolerance to Growth

632    on Low Iron and Uncovers Posttranscriptional Control. Plant Physiol.

633    2003;133:1102–10.

634    30. Osakabe Y, Yamaguchi-Shinozaki K, Shinozaki K, Tran LSP. Sensing the

635    environment: Key roles of membrane-localized kinases in plant perception and

636    response to abiotic stress. J Exp Bot. 2013;64:445–58.

637    31. Novaković L, Guo T, Bacic A, Sampathkumar A, Johnson KL. Hitting the Wall-

638    Sensing and Signaling Pathways Involved in Plant Cell Wall Remodeling in

639    Response to Abiotic Stress. Plants (Basel, Switzerland). MDPI; 2018;7:89.

640    32. Bayer PE, Golicz AA, Tirnaz S, Chan CKK, Edwards D, Batley J. Variation in

641    abundance of predicted resistance genes in the Brassica oleracea pangenome. Plant

642    Biotechnol J. 2019;17:789–800.

643    33. Zhou P, Silverstein KAT, Ramaraj T, Guhlin J, Denny R, Liu J, et al. Exploring

644    structural variation and gene family architecture with De Novo assemblies of 15

645    Medicago genomes. BMC Genomics. BMC Genomics; 2017;18:1–14.

646    34. Gao L, Gonda I, Sun H, Ma Q, Bao K, Tieman DM, et al. The tomato pan-

647    genome uncovers new genes and a rare allele regulating fruit flavor.

648   35. Shen X, Liu ZQ, Mocoeur A, Xia Y, Jing HC. PAV markers in Sorghum bicolour:

649   genome pattern, affected genes and pathways, and genetic linkage map

650   construction. Theor Appl Genet. 2015;128:623–37.

651   36. Phan HTT, Ellwood SR, Adhikari K, Nelson MN, Oliver RP. The first genetic and

652   comparative map of white lupin (Lupinus albus L.): Identification of QTLs for

653   anthracnose resistance and flowering time, and a locus for alkaloid content. DNA

654   Res. 2007;14:59–70.

655   37. Książkiewicz M, Nazzicari N, Yang H, Nelson MN, Renshaw D, Rychel S, et al. A

656   high-density consensus linkage map of white lupin highlights synteny with narrow-

657   leafed lupin and provides markers tagging key agronomic traits. Sci Rep.

658   2017;7:15335.

659   38. Plewiński P, Książkiewicz M, Rychel-Bielska S, Rudy E, Wolko B. Candidate

660   domestication-related genes revealed by expression quantitative trait loci mapping of

661   narrow-leafed lupin (Lupinus angustifolius L.). Int J Mol Sci. 2019;20:1–24.

662   39. Rychel S, Książkiewicz M, Tomaszewska M, Bielski W, Wolko B. FLOWERING

663   LOCUS T, GIGANTEA, SEPALLATA, and FRIGIDA homologs are candidate genes

664   involved in white lupin (Lupinus albus L.) early flowering. Mol Breed. 2019;39:43.

665   40. Andrés F, Coupland G. The genetic basis of flowering responses to seasonal

666   cues. Nat. Rev. Genet. 2012.

667   41. Golicz AA, Batley J, Edwards D. Towards plant pangenomics. Plant Biotechnol J.

668   2016;14:1099–105.

669   42. Montenegro JD, Golicz AA, Bayer PE, Hurgobin B, Lee HT, Chan CKK, et al. The

670   pangenome of hexaploid bread wheat. Plant J. 2017;90:1007–13.

671   43. Li YH, Zhou G, Ma J, Jiang W, Jin LG, Zhang Z, et al. De novo assembly of

672   soybean wild relatives for pan-genome analysis of diversity and agronomic traits. Nat

673    Biotechnol. 2014;32:1045–52.

674    44. Diamond J. Evolution, consequences and future of plant and animal

675    domestication. Nature. Nature Publishing Group; 2002. p. 700–7.

676    45. Wolko B, Clements JC, Naganowska B, Nelson MN, Yang H. Lupinus. Wild Crop

677    Relat Genomic Breed Resour. Berlin, Heidelberg: Springer Berlin Heidelberg; 2011.

678    p. 153–206.

679    46. Hamblin MT, Jannink J-L. Factors Affecting the Power of Haplotype Markers in

680    Association Studies. Plant Genome J. 2011;4:145.

681    47. Mousavi-Derazmahalleh M, Nevado B, Bayer PE, Filatov DA, Hane JK, Edwards

682    D, et al. The western Mediterranean region provided the founder population of

683    domesticated narrow-leafed lupin. Theor Appl Genet. Springer Berlin Heidelberg;

684    2018;131:2543–54.

685    48. Mousavi-Derazmahalleh M, Bayer PE, Nevado B, Hurgobin B, Filatov D, Kilian A,

686    et al. Exploring the genetic and adaptive diversity of a pan-Mediterranean crop wild

687    relative: narrow-leafed lupin. Theor Appl Genet. Springer Berlin Heidelberg;

688    2018;131:887–901.

689    49. Elias E. Soils of the Ethiopian Highlands: Geomorphology and Properties. 2016.

690    50. Lambers H, Bishop JG, Hopper SD, Laliberté E, Zúñiga-Feest A. Phosphorus-

691    mobilization ecosystem engineering: the roles of cluster roots and carboxylate

692    exudation in young P-limited ecosystems. Ann Bot. 2012;110:329–48.

693    51. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina

694    sequence data. Bioinformatics. 2014/04/01. Oxford University Press; 2014;30:2114–

695    20.

696    52. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al.

697    SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell

698     Sequencing. J Comput Biol. 2012;19:455–77.

699     53. Hu Z, Sun C, Lu KC, Chu X, Zhao Y, Lu J, et al. EUPAN enables pan-genome

700     studies of a large number of eukaryotic genomes. Bioinformatics. 2017;33:2408–9.

701     54. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al.

702     BLAST+: architecture and applications. BMC Bioinformatics. BioMed Central;

703     2009;10:421.

704     55. Yandell M, Holt C. MAKER2: an annotation pipeline and genome-database

705     management tool for second-generation genome projects. BMC Bioinformatics.

706     2011;12:491.

707     56. Stanke M, Morgenstern B. AUGUSTUS: a web server for gene prediction in

708     eukaryotes that allows user-defined constraints. Nucleic Acids Res. 2005;33:W465–

709     7.

710     57. Korf I. Gene finding in novel genomes. BMC Bioinformatics. 2004;5.

711     58. Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntenically

712     mapped cDNA alignments to improve de novo gene finding. Bioinformatics.

713     2008;24:637–44.

714     59. Sallet E, Gouzy J, Schiex T. EuGene-PP: A next-generation automated

715     annotation pipeline for prokaryotic genomes. Bioinformatics. 2014;30:2659–61.

716     60. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler

717     transform. Bioinformatics. 2010;26:589–95.

718     61. Martin M. Cutadapt removes adapter sequences from high-throughput

719     sequencing reads. EMBnet.journal. 2011;17:10.

720     62. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: A Fast and

721     Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. Mol

722     Biol Evol. 2015;32:268–74.

723   63. Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermiin LS.

724   ModelFinder: fast model selection for accurate phylogenetic estimates. Nat Methods.

725   Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights

726   Reserved.; 2017;14:587–9.

727   64. Letunic I, Bork P. Interactive tree of life (iTOL) v3: an online tool for the display

728   and annotation of phylogenetic and other trees. Nucleic Acids Res. 2016;44:W242–5.

729   65. Hubisz MJ, Falush D, Stephens M, Pritchard JK. Inferring weak population

730   structure with the assistance of sample group information. Mol Ecol Resour.

731   2009;9:1322–32.

732   66. Zhang C, Dong SS, Xu JY, He WM, Yang TL. PopLDdecay: A fast and effective

733   tool for linkage disequilibrium decay analysis based on variant call format files.

734   Bioinformatics. 2019;35:1786–8.

735   67. Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES.

736   TASSEL: software for association mapping of complex traits in diverse samples.

737   Bioinformatics. 2007;23:2633–5.

738   68. Pfeifer B, Wittelsbürger U, Ramos-Onsins SE, Lercher MJ. PopGenome: An

739   Efficient Swiss Army Knife for Population Genomic Analyses in R. Mol Biol Evol.

740   2014;31:1929–36.

741   **ACKNOWLEDGEMENTS**

742

743   **DECLARATIONS**

744   **Ethics approval and consent to participate**

745   Not applicable

746

747   **Consent for publication**

748     Not applicable

749

**Competing interests**

751     The authors declare that they have no competing interests.

752

**Funding**

754     This project has received funding from the European Research Council (ERC) under

755     the European Union's Horizon 2020 research and innovation program (Starting Grant

756     LUPINROOTS - grant agreement No 637420 to B.P.) and from the Innovate UK

757     project 133048 (Ethiopian Lupins for Food and Feed) to H.S.

758

**Availability of data and materials**

760     The detailed methods and datasets supporting the conclusions of this report are

761     included within the article and its additional files. All deep sequencing data reported

762     in this paper have been submitted to the NCBI. The datasets generated and

763     analyzed during the current study are available from the corresponding author upon

764     request. Full genomic and raw sequence data are publicly available for download on

765     the White Lupin genome portal [www.whitelupin.fr/pangenome] that contains a

766     Genome Browser, Expression tools and a Sequence retriever dedicated to the

767     pangenome. The pangenome project and raw data has been deposited at

768     DDBJ/ENA/GenBank under the accession PRJNA608889.

769

**Authors' contributions**

771     A.S. developed bioinformatic resources and performed pangenome assembly. J.T.

772     and F.D. performed DNA extraction and experiments. M.N., H.S., L.Y. and M.K.

773    provided genetic material. B.H. performed data analysis. B.H., M.K., M.N. and B.P.

774    designed experiments and wrote the article.

775

776    **Corresponding authors**

777    Correspondence to Benjamin Péret (benjamin.peret@supagro.fr) and Bárbara
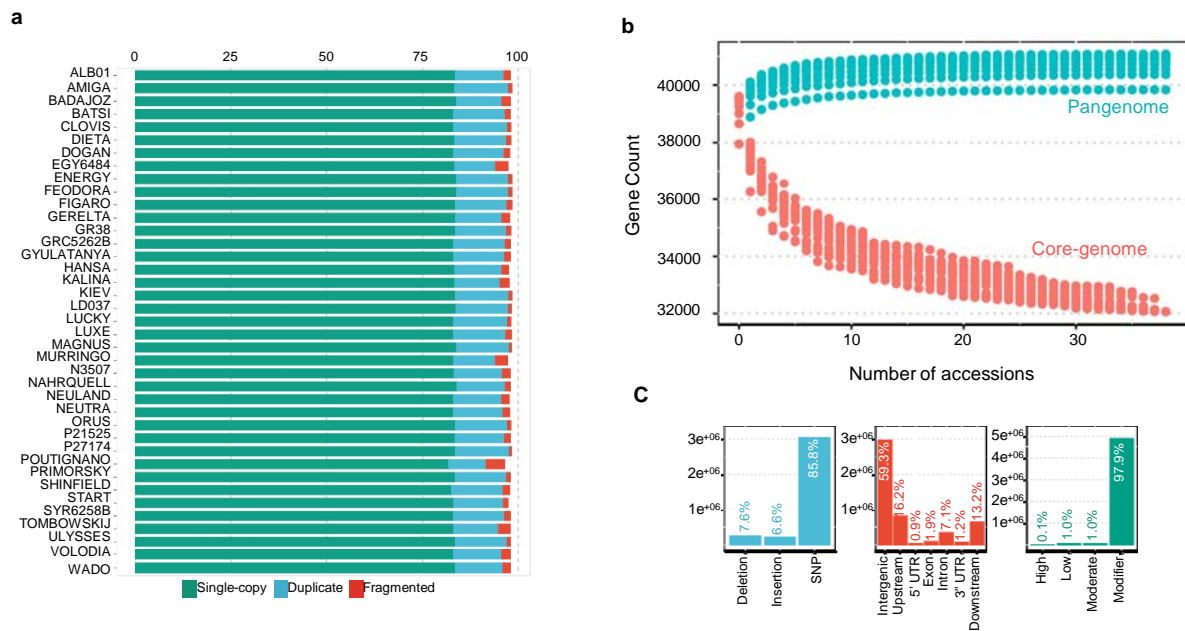
778    Hufnagel (barbara.hufnagel@supagro.fr).

**Figure 1. Pangenome of *L. albus*. (a)** BUSCO percent completeness of all assemblies. All of the assemblies of this study have BUSCO completeness higher than 91.7%. **(b)** Pangenome modeling **(c)** Distribution of variants along white lupin pangenome. Types of variations identified (blue); positioning of the variants in the genome in relation to the gene structures (red); impact of the variants (green).

**Figure 2. Phylogeny and population structure of 39 accessions of *L. albus*. (a)** Maximum likelihood phylogenetic tree of white lupin constructed based on 3.5 M SNPs. The accessions are divided in 6 idiotypes. **(b)** Model-based clustering analysis with different numbers of ancestral kinships (k=4, 5 and 6). The y axis quantifies cluster membership and the x axis list the different accessions. The positions of these accessions on the x axis are consistent with those in the phylogenetic tree. **(c)** Principal component analysis based on 3.5 M SNPs. The ellipses are discriminating the accessions of each idiotype groups. **(d)** Genome-wide average LD decay estimated from different white lupin group. The decay of LD with physical distance between SNPs to half of the maximum values occurred at 3.85 kb ($r^2$ =0.38) considering all accessions.
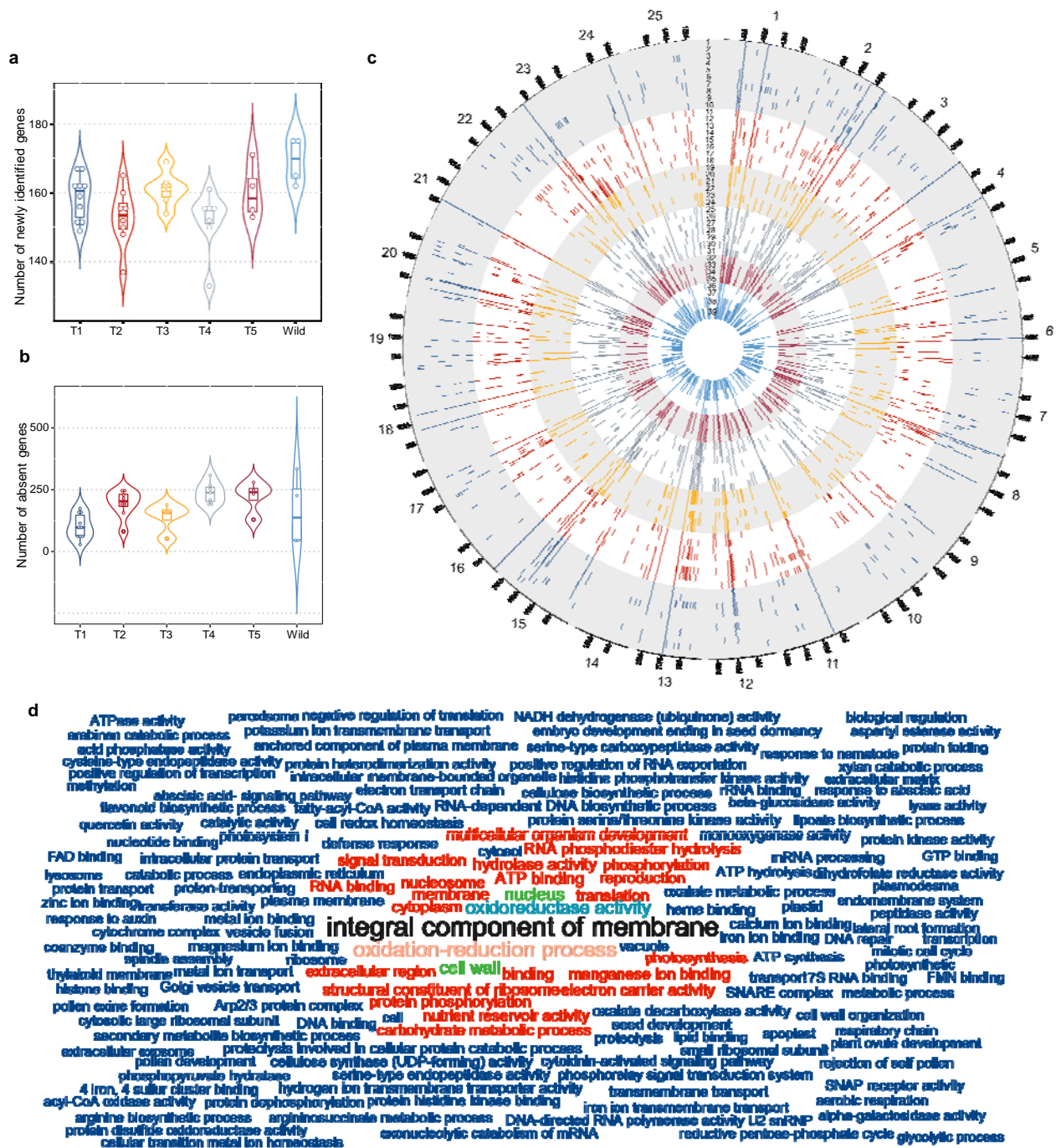
**Figure 3. PAV of coding gene in *L. albus*. (a)** Number of newly identified genes by phylogenetic group. **(b)** Number of absent genes by phylogenetic groups. **(c)** Positioning of absent genes in the 25 white lupin chromosomes in each one of the 39 accessions. Order of accessions from outer to inner track: 1-AMIGA, 2-FEODORA, 3-FIGARO, 4-ENERGY, 5-KIEV MUTANT, 6-HANSA, 7-P21525, 8-PRIMORSKY, 9-DIETA, 10-VOLODIA, 11-START, 12-N3507, 13-TOMBOWSKIJ, 14-KALINA, 15-SYR6258B, 16-LUCKY, 17-MURRINGO, 18-SHINFIELD, 19-ALB01, 20-LUXE, 21-ULYSSE, 22-MAGNUS, 23-CLOVIS, 24-ORUS, 25-NAHRQUELL, 26-GYUNLATANYA, 27-NEULAND, 28-NEUTRA, 29-BADAJOZ, 30-EGY6484B, 31-POUTIGANO, 32-P27174, 33-GERELTA, 34-DOGAN, 35-WADO, 36-GR38, 37-GRAECUS, 38-BATSI, 39-GRC5262B. The accessions' colors reflect the 6 idiotypes. **(d)** Functional enrichment analysis of the variable genome. Graphical representation of enriched biological process (GOs). Size of the words and colors are proportional to their representativeness in the gene pool.
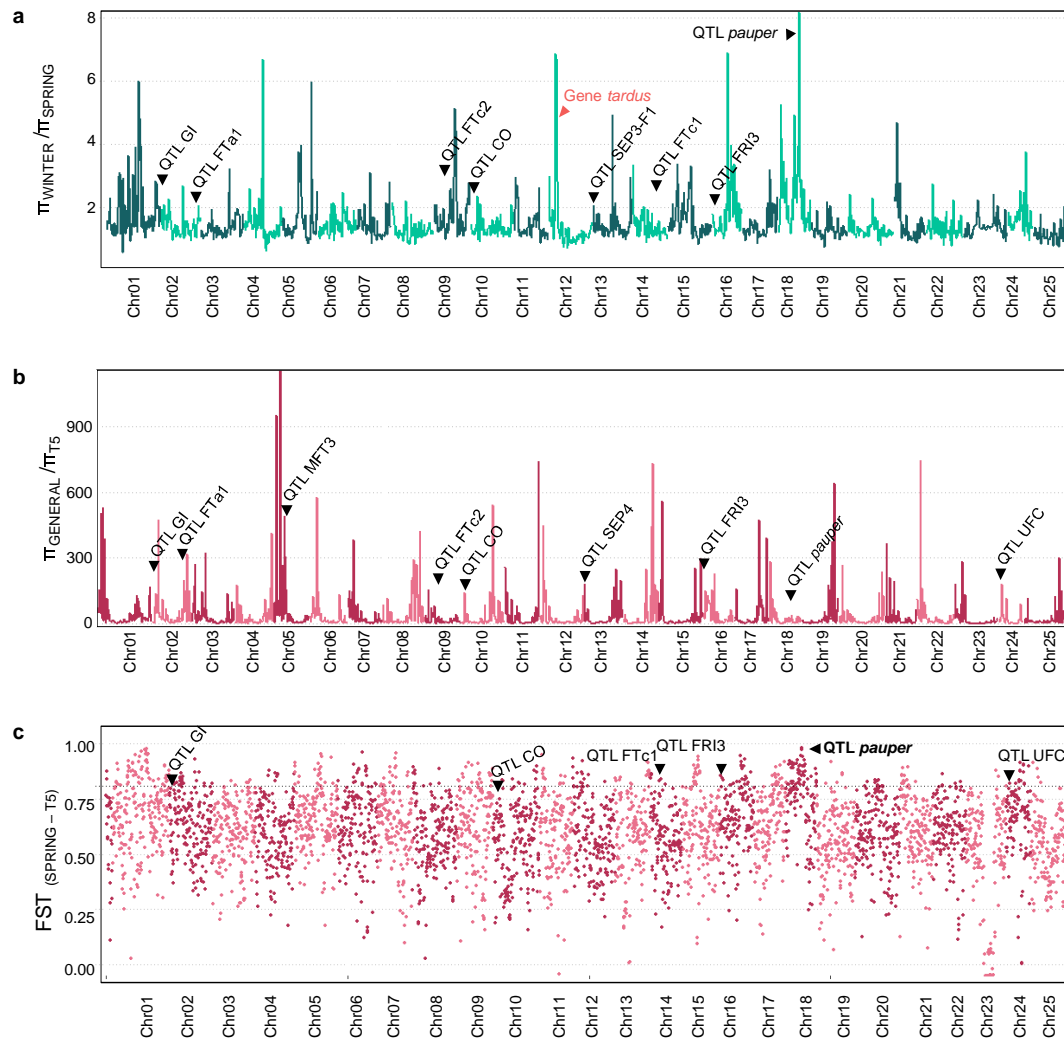
**Figure 4. Footprints of selection in the white lupin genome.** Nucleotide diversity (π) comparison between **(a)** Winter (T3 and T4) and Spring accessions (T1 and T2) and **(b)** between all accessions (General) and Ethiopian accessions (T5). QTLs previously reported and a *L. angustifolius* domestication gene (red) that overlapped with selective sweeps are marked. **(c)** Fst-based genome-wide analysis of population differentiation estimated between Spring (T1 and T2) and Ethiopian (T5) accessions. Black horizontal dashed line marks the .90 percentile of distribution of Fst estimated (Fst = 0.81).
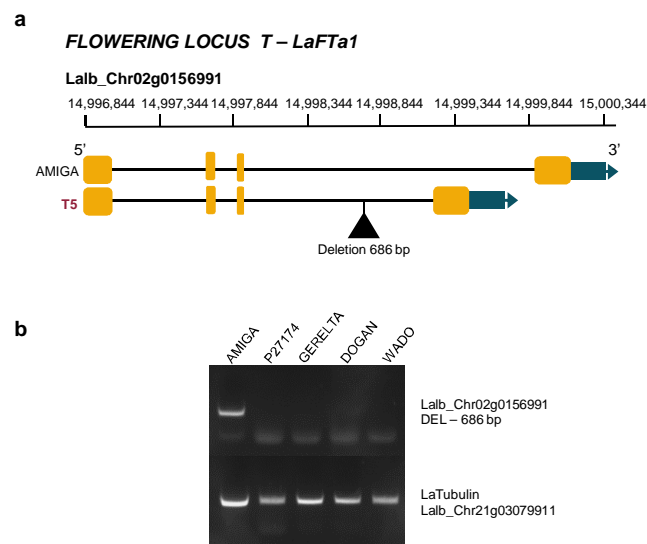
**Figure 5. Identification and allele variation of candidate gene *FLOWERING LOCUS T.* (a)** Candidate gene located on chromosome 2. Type 5 accessions, originated from Ethiopia, have a deletion of 686 bp in the third intron. **(b)** Confirmation of the deletion in the third intron of Type 5 accession by PCR. Gene *LaTubulin* was used as positive control.