1
2
3
4 **The *p*-coumaroyl arabinoxylan transferase *HvAT10* underlies natural variation in whole-grain cell**
5 **wall phenolic acids in cultivated barley**
6

7 Kelly Houston[1§], Amy Learmonth[2§], Ali Saleh Hassan[3,4§], Jelle Lahnstein[3], Mark Looseley[5], Alan Little[3],
8 Robbie Waugh[1,2,3]*, Rachel A Burton[3]*, Claire Halpin[2]*

9 § Authors contributing equally

10 *Corresponding authors: Robbie Waugh, Rachel A Burton, Claire Halpin

11

12 [1]Cell and Molecular Sciences, The James Hutton Institute, Errol Road Invergowrie, Dundee, DD2 5DA,

13 Scotland, UK.

14 [2]Division of Plant Sciences, School of Life Sciences, University of Dundee at The James Hutton

15 Institute, Invergowrie, Dundee, DD2 5DA, Scotland, UK.

16 [3]School of Agriculture, Food and Wine, University of Adelaide, Urrbrae SA 5064

17 **Current affiliation if different to those above**
18 Ali Saleh Hassan;
19 [4]CSIRO Agriculture and Food, Waite campus, Wine Innovation West Bld, Hartley Grove, Urrbrae, SA
20 5064,
21
22 Mark Looseley;
23 [5]Xelect Ltd, Horizon House, Abbey Walk, St Andrews, Fife, KY16 9LB, Scotland, UK
24

25

26 Communicating Author:          Claire Halpin

27 e-mail:                        c.halpin@dundee.ac.uk

28 tel:                           +44 (0) 1382 568775

29

30

31 **Keywords:** barley, p-coumarate, ferulate, GWAS

32

33 **Phenolic acids in cereal grains have important health-promoting properties and influence**

34 **digestibility for industrial or agricultural uses. Here we identify alleles of a single BAHD *p*-**

35 **coumaroyl arabinoxylan transferase gene, *HvAT10,* as responsible for the natural variation in cell**

36 **wall-esterified *p*-coumaric and ferulic acid in whole grain of a collection of cultivated two-row**

37 **spring barley genotypes. We show that *HvAT10* is rendered non-functional by a premature stop**

38 **codon mutation in approximately half of the genotypes in our mapping panel. The causal**

39 **mutation is virtually absent in wild and landrace germplasm suggesting an important function for**

40 **grain arabinoxylan *p*-coumaroylation pre-domestication that is dispensable in modern agriculture.**

41 **Intriguingly, we detected detrimental impacts of the mutated locus on barley grain quality traits.**

42 **We propose that *HvAT10* could be a focus for future grain quality improvement or for**

43 **manipulating phenolic acid content of wholegrain food products.**

44 Phenolic acids in the cell walls of cereals limit digestibility[1] when grain or biomass is used for animal

45 feed or processed to biofuels and chemicals. They are also important dietary antioxidant, anti-

46 inflammatory and anti-carcinogenic compounds and contribute to beer flavour and aroma[2,3]. The

47 hydroxycinnamates, *p*-coumarate and ferulate (*p*CA and FA respectively), are the major phenolic

48 acids in grasses. Both occur as decorations ester-linked to cell wall arabinoxylan. Lignin also has

49 esterified *p*CA decorations but FA in lignin is incorporated directly into the growing polymer by ether

50 linkages[4]. Besides its role as a lignin monomer, FA in the cell wall acts to cross-link arabinoxylans to

51 each other and to lignin, and it is this cross-linking that may impede digestibility. The role of *p*CA in

52 cell walls is less clear. In lignin, it may promote polymerisation of sinapyl alcohol monolignols[5] and

53 act as a termination unit[4], but there are no clear theories about its role when attached to

54 arabinoxylans. Given the importance of *p*CA and FA to plant health and the uses of cereal crops,

55 there has been much recent interest in identifying genes that can be manipulated in transgenic

56 plants to influence phenolic acid content[6-14]. Given current GM legislation in some countries it

57 would be more appropriate for crop improvement to identify genes and alleles determining natural

58 variation in *p*CA and FA that could be exploited immediately in contemporary plant breeding.

59 We quantified cell wall-esterified *p*CA and FA in the wholegrain of a replicated GWAS panel of 211

60 elite 2-row spring barley cultivars grown in a field polytunnel. We observed a 6-fold variation for

61 esterified *p*CA (54 µg/g - 327 µg/g) and a greater than 2-fold variation in esterified FA (277 µg/g -748

62 µg/g) (Supplementary Fig. 1a,b, Supplementary Data 1, 2) with no correlation between FA and *p*CA

63 levels ($R^2$ = 0.04). A GWAS of this data using 43,834 SNP markers identified a single highly significant

64 association for grain esterified *p*CA on chromosome 7H (-log10(p)=13.9; Fig. 1a, Supplementary Data

65 3) and a co-locating peak for FA just below statistical significance (-log10(p)=3.9; Fig. 1b,

66   Supplementary Data 3). Given the closeness of FA and *p*CA on the phenylpropanoid pathway we

67   also conducted a GWAS using FA:*p*CA concentration ratios which provides internal data

68   normalisation, reducing inherent variability in single compound measurements[15]. Mapping

69   log[FA:*p*CA] values increased both the strength and significance of association with the locus (-

70   log10(p)=19.4; Fig. 1c, Supplementary Fig. 2c, Supplementary Data 3), confirming a level of

71   dependency between esterified FA and esterified *p*CA concentrations. GWAS on similar data from a

72   semi-independent set of 128 greenhouse-grown barley genotypes identified the same associations

73   (Supplementary Fig. 2a-c, Supplementary Data 3).

74   The entire region above the adjusted false discovery rate (FDR) threshold for the log[FA:*p*CA] values

75   spanned a 65.7MB segment of chromosome 7H (459,131,547bp - 524,825,783bp) containing 347

76   high-confidence gene models. We surveyed this region for genes involved in phenolic acid or cell

77   wall biosynthesis. This revealed several candidates including two cinnamyl alcohol dehydrogenases

78   (*CAD*s), a caffeate-O-methyltransferase (*HvCOMT1*[16]) and three BAHD acyltransferases.

79   Interrogation of an RNA-seq dataset for 16 barley tissues[17] revealed that five of these six candidates

80   exhibited moderate to low levels of expression across all surveyed tissues (Fig. 2a). However, the

81   *BAHD* gene HORVU7Hr1G085100 stood out as being highly expressed in the hull lemma and palea

82   where 80% of grain *p*CA is found[18] (Fig. 2a, b). We then consulted a database of variant calls from a

83   barley RNA-seq dataset that included 118 of our GWAS genotypes[19]. We observed no SNP variation

84   in two of the candidate genes. Three had one SNP each; *COMT1* (HORVU7Hr1G082280) had a

85   synonymous SNP, one *CAD* (HORVU7Hr1G079380) a SNP in the 3' UTR and one *BAHD*

86   (HORVU7Hr1G085390) a non-synonymous but rare SNP. None appeared likely to impair gene

87   function. However, the *BAHD* HORVU7Hr1G085100 had 3 SNPs including one causing a premature

88   stop codon leading to loss of a third of the protein sequence. BLASTp of the predicted full-length

89   HORVU7Hr1G085100 protein sequence revealed it was 79% identical to rice *OsAT10*

90   (LOC_Os06g39390.1), a gene functionally characterised as a *p*-coumaroyl CoA arabinoxylan

91   transferase[7]. Critically, overexpression of *OsAT10* in rice dramatically increases cell wall-esterified

92   *p*CA levels in leaves while concomitantly reducing the levels of esterified FA[7]. A maximum likelihood

93   phylogenetic tree of *BAHD* gene sequences confirmed HORVU7Hr1G085100 as *HvAT10* (Fig. 2c) and

94   another of our candidates, HORVU7Hr1G085390, as a possible *HvAT10* paralog with negligible

95   expression in the tissues surveyed (Fig. 2a). The third *BAHD*, HORVU7Hr1G085060, is likely an

96   AT8[7,13].

97   To more accurately document polymorphisms in *HvAT10*, we PCR-sequenced the gene from 52

98   genotypes of the GWAS panel (Supplementary Data 1). Two nonsynonymous SNPs, one in each of

99     *HvAT10*'s two exons (Fig. 3a), were in complete linkage disequilibrium across the 52 lines. A G/A

100     SNP at 430bp translates to either a valine or isoleucine, substituting one non-polar, neutral amino

101     acid for another, so unlikely to affect function. By contrast, a C/A SNP at 929bp produces either

102     serine in the full length protein, or a premature stop codon that truncates the protein by 124 amino

103     acids, removing the BAHD family conserved DFGWG motif (DVDYG in barley and other grasses)

104     thought to be essential for catalysis[20-22] (Fig. 3b). The *at10^STOP* mutation is therefore predicted to

105     knock-out gene function. We designed a diagnostic Kompetitive Allele Specific PCR (KASP) assay to

106     distinguish the two *HvAT10* alleles and genotyped all 212 cultivars in our GWAS population

107     (Supplementary Table 1). Consistent with the hypothesis that *at10^STOP* is the causal variant

108     underlying the log[FA:*p*CA] GWAS peak, no SNP scored higher than the KASP diagnostic when

109     included in the GWAS although one, JHI-Hv50k-2016-488774, in complete LD scored equally highly.

110     *HvAT10* had a minor allele frequency of 0.48 and appears to significantly influence levels of both *p*CA

111     (*p*= 4.30e-19) and FA in grain (p=1.80e-11) with the median for *at10^STOP* genotypes being 28% lower

112     for *p*CA (Fig. 3c) and 14% higher for FA (Fig. 3d) than those with the wildtype allele. Comparing the

113     median log[FA:*p*CA] for *at10^STOP* cultivars (0.58) to the wildtype cultivar group (0.37) showed an even

114     higher significant difference between the groups (*p*= 7.56e-50) (Fig. 3e, Supplementary Fig. 3).

115     In contrast to our initial observation on the whole population, plotting grain esterified *p*CA against

116     FA (Fig. 3f) within each allele group now reveals positive correlations, suggesting that although flux

117     into phenolic acid biosynthesis may differ between cultivars, it co-ordinately affects both phenolic

118     acids. The *at10^STOP* genotypes show approximately one-third less *p*CA than wildtype genotypes

119     reflecting a deficiency of *p*CA on arabinoxylan in cultivars that lack a functional *p*-coumaroyl CoA

120     arabinoxylan transferase. Nevertheless, two-thirds of cell wall esterified *p*CA remains since most *p*CA

121     is associated with lignin[23,24] through the action of other BAHD genes. The influence of *at10^STOP* on FA

122     is evidenced by considering the 27 cultivars with grain esterified FA above 600 µg/g; 23 of these

123     have the *at10^STOP* allele (Fig. 3f; Supplementary Data 1). This effect on FA might occur in several

124     ways: *p*CA that cannot be esterified onto arabinoxylan could be methoxylated to produce FA thereby

125     increasing FA pools for transfer onto arabinoxylan, or alternatively, *p*CA and FA may compete for

126     transfer onto a shared acceptor (likely UDP-arabinose[12]) before incorporation into arabinoxylan such

127     that loss of *p*CA transfer by *HvAT10* leaves more free acceptor for FA transfer. Either mechanism

128     could explain how *at10^STOP* can indirectly increase grain cell wall esterified ferulate. An inverse

129     interaction between levels of *p*CA and FA on arabinoxylan was also seen in transgenic rice[7],

130     switchgrass[26], and *Setaria viridis*[13] where BAHD expression was manipulated.

131   Intrigued by the prevalence of $at10^{STOP}$ in 50% of our elite barley genepool we were curious about

132   whether this had any ecological, evolutionary, or performance-related significance.  To explore, we

133   PCR-sequenced a collection of 114 georeferenced barley landraces and 76 wild barley (*Hordeum*

134   *spontaneum*) genotypes[26] across the $at10^{STOP}$ polymorphism (Supplementary Data 1).  We found

135   $at10^{STOP}$ to be extremely rare, present in three of 114 landraces and absent in all 76 wild genotypes

136   (Supplementary Fig. 4a, Supplementary Data 1).  The three landraces show a clear pattern of identity

137   by descent, clustering in the same clade of the dendrogram (Supplementary Fig. 4a).  We interpret

138   these data as suggesting strong selection against the premature stop codon in wild germplasm and

139   that $at10^{STOP}$ was a post-domestication mutation that under cultivation has no pronounced negative

140   effects on fitness.

141   Several possibilities could explain enrichment of $at10^{STOP}$ in the cultivated genepool.  To explore, we

142   first  calculated  genome  wide  $F_{ST}$  by  locus  using  two  groups  based  on  the  *HvAT10*  allele.

143   HORVU7Hr1G084140 (a Serine/threonine-protein kinase not expressed in the lemma or palea) also

144   had an $F_{ST}$ of 1.0, and three other genes an $F_{ST}$ above 0.875 (Supplementary Fig. 5a,b Supplementary

145   Data 4).  Based on their functional annotations and gene expression patterns (Supplementary Data

146   4, Supplementary Fig. 5c) we observed no obvious reason for these to be under strong selection and

147   responsible for enhancing the frequency of $at10^{STOP}$ via extended LD.

148   Next, due to the exclusive expression of *HvAT10* in the lemma and palea, we measured a series of

149   grain morphometric traits across our panel.  We found that, on average, grain from the $at10^{STOP}$

150   genotypes had significantly reduced grain width compared to cultivars with the wildtype allele (Table

151   1) suggesting a potential role for arabinoxylan-esterified phenolic acids in modifying grain shape.  Xu

152   *et al*[27] previously identified a QTL hotspot on chromosome 7H for traits including grain area, and

153   grain width.  The eight 9K iSelect markers defining this QTL can be positioned on the current physical

154   map at 482-500MB on 7H, corresponding to the location of *HvAT10*.  Wang *et al*[28] also identified a

155   QTL for grain length:width, grain perimeter, and grain roundness at the same location.

156   Prompted by these observations and the prevalence of registered UK barley varieties in our panel,

157   we then explored grain parameters recorded in an extensive historical dataset from the UK's

158   National and Recommended Lists trials 1988-2016[29].  Different grain quality phenotypes were

159   available for up to 106 of our cultivars.  Group comparisons of WT and $at10^{STOP}$genotypes revealed

160   surprising differences for hot water extract, diastatic power, germinative energy in 4ml, and wort

161   viscosity (Table 1).  In all cases, the group of $at10^{STOP}$ cultivars had poorer quality, offering no

162   evidence of positive selection during breeding.  The variation associated with the *HvAT10* locus is

163   however highly significant and of potential interest for optimising grain quality traits (Table 1).

164     Finally, to understand more about the origin of the *at10*<sup>STOP</sup> in elite germplasm, we investigated its

165     occurrence in the pedigree of our GWAS population. The earliest cultivar with *at10*<sup>STOP</sup> is *cv.* Kenia

166     (cross between the Swedish landrace Gull and Danish landrace Binder) released in 1931 and

167     subsequently introduced into NW European breeding programmes.  Despite smaller grain and

168     slightly poorer malting properties compared to its contemporary UK varieties, it established a long-

169     standing position as a parent for further crop improvement due to its short stiff straw, earliness and

170     high yield[30].  Several decades later, *at10*<sup>STOP</sup>-containing derivatives of Kenia, such as *cv.* Delta

171     (National list 1959), were still being used as parents in our pedigree chart.

172     Taken together, we conclude that the continued prevalence of Kenia-derived germplasm may go

173     some way to explaining the frequency of the *at10*<sup>STOP</sup> allele in our population.  While this may simply

174     be a straightforward genetic legacy of historical barley breeding, our data suggests that purging this

175     mutation could assist the development of superior quality barley varieties. Conversely, much

176     research has focussed on the beneficial bioactivity of ferulate in the diet and the *at10*<sup>STOP</sup> allele could

177     enable breeding for increased ferulate in wholegrain products.

178     [2047 words]

179     **References**

180       1.   Hatfield RD, Ralph J., & Grabber JH. Cell wall structural foundations: Molecular basis for

181          improving forage digestibilities. *Crop Sci.* **39**, 27-37 (1999).

182       2.   Calinoiu LF & VodnarDC. Whole grains and phenolic acids: a review on bioactivity,

183          functionality, health benefits and bioavailability. *Nutrients* **10**, 1615 (2018).

184       3.   Lentz M. The Impact of simple phenolic compounds on beer aroma and flavour.

185          *Fermentation* **4**, 20 (2018).

186       4.   Hatfield RD, Rancour DM. & Marita JM. Grass cell walls: a story of cross-linking. *Front. Plant*

187          *Sci.* **7**, 2056 (2017).

188       5.   Ralph J, et al. Peroxidase-dependent cross-linking reactions of *p*-hydroxycinnamates in plant

189          cell walls. *Phytochem. Rev.* **3,** 79–96. (2004).

190       6.   Withers S, Lu FC, Kim H, Zhu YM, Ralph J, Wilkerson CG. Identification of grass-specific

191          enzyme that acylates monolignols with *p*-coumarate. *J Biol Chem*. **287**:8347–8355. (2012).

192       7.   Bartley,LE. et al. Overexpression of a BAHD acyltransferase, OsAt10, alters rice cell wall

193          hydroxycinnamic acid content and saccharification. *Plant Physiol*. **161**, 1615-1633 (2013).

194       8.   Petrik DL, et al., *p*-Coumaroyl-CoA: monolignol transferase (PMT) acts specifically in the

195          lignin biosynthetic pathway in *Brachypodium distachyon*. *Plant J.***77**:713–726. (2014).

196   9.  Marita JM, Hatfield RD, Rancour DM, Frost KE. Identification and suppression of the p-
197       coumaroyl CoA: hydroxycinnamyl alcohol transferase in *Zea mays* L. *Plant J.* **78**:850–864.
198       (2014).

199   10. Sibout R, Le Bris P, Legee F, Cezard L, Renault H, Lapierre C. Structural redesigning
200       Arabidopsis lignins into alkali-soluble lignins through the expression of p-coumaroyl-coA:
201       monolignol transferase PMT. *Plant Physiol.*;**170**:1358–1366. (2016).

202   11. Karlen SD, et al.  Monolignol ferulate conjugates are naturally incorporated into plant lignins.
203       *Sci Adv.***2:**1–9. (2016).

204   12. Buanafina MMD, Fescemyer HW, Sharma M, Shearer EA. Functional testing of a PF02458
205       homologue of putative rice arabinoxylan feruloyl transferase genes in *Brachypodium*
206       *distachyon*. *Planta*.**243**:659–674. (2016).

207   13. De Souza WR et al. Suppression of a single BAHD gene in *Setaria viridis* causes large, stable
208       decreases in cell wall feruloylation and increases biomass digestibility. *New Phytol*.**218**:81–
209       93. (2018).

210   14. Mota et al. Suppression of a BAHD acyltransferase decreases *p*-coumaroyl on arabinoxylan
211       and improves biomass digestibility in the model grass *Setaria viridis*. *Plant J.*
212       https://doi.org/10.1111/tpj.15046 (2020).

213   15. Petersen A, et al. On the hypothesis-free testing of metabolite ratios in genome-wide and
214       metabolome-wide association studies. *BMC Bioinformatics* **13,** 120 (2012).

215   16. Daly P et al. RNAi-suppression of barley caffeic acid O-methyltransferase modifies lignin
216       despite redundancy in the gene family. *Plant Biotechnol J* **17**; 3, 594-607 (2019).

217   17. Colmsee C et al. BARLEX - the Barley Draft Genome Explorer. *Mol Plant*. **8**(6):964-966.(2015).

218   18. Barron C et al. Assessment of biochemical markers identified in wheat for monitoring barley
219       grain tissue. *J Cereal Sci.* **74,** 11-18. (2017).

220   19. Rapazote-Flores P, et al., BaRTv1.0: an improved barley reference transcript dataset to
221       determine accurate changes in the barley transcriptome using RNA-seq.  *BMC Genomics*
222       **11**;20(1):968. (2019).

223   20. Ma X, Koepke J, Panjikar S, Fritzsch G, Stöckigt J. Crystal structure of vinorine synthase, the
224       first representative of the BAHD superfamily. *J Biol Chem*. Apr **8**;280(14):13576-83.(2005).

225   21. D'Auria ,JC. Acyltransferases in plants: a good time to be BAHD. *Curr  Opin Plant Biol.*
226       Jun;**9**(3):331-40. (2006)

227   22. Morales-Quintana, L., Alejandra Moya-Leon, M., Herrera, R. Computational study enlightens
228       the structural role of the alcohol acyltransferase DFGWG motif. *J. Mol. Model.* **21**, 216 (2015)

229   23. Ralph, J. Hydroxycinnamates in lignification. *Phytochem Rev*. **9**, 65-83 2010.

230    24. Lapierre, C., Voxeur, A., Karlen, SD, Helm, RF. & Ralph J. Evaluation of Feruloylated and
231        p-Coumaroylated Arabinosyl Units in Grass Arabinoxylans by Acidolysis in Dioxane/Methanol
232        *J. Agric. Food Chem.* **66**, 5418–5424 (2018).
233    25. Li GT, et al, Overexpression of a rice BAHD acyltransferase in switchgrass (panicum virgatum
234        L.) enhances saccharification. *BMC Biotechnol* **18**:54 (2018).

235    26. Russell J., et al. Exome sequencing of geographically diverse barley landraces and wild
236        relatives gives insights into environmental adaptation. *Nat Genet* **48,** 1024–1030 (2016)
237    27. Xu X, et al. Genome-Wide Association Analysis of Grain Yield-Associated Traits in a Pan-
238        European Barley Cultivar Collection. *The Plant Genome*, **11**: 1-11 170073. (2018).

239    28. Wang Q., et al. Dissecting the Genetic Basis of Grain Size and Weight in Barley (*Hordeum*
240        *vulgare* L.) by QTL and Comparative Genetic Analyses. *Front. Plant Sci*. **10,** 469. (2019).
241    29. Looseley M, et al. Association mapping of malting quality traits in UK spring and winter
242        barley cultivar collections. *Theor Appl Genet* **133**; 2567-2582. (2020)
243    30. Bell GDH. Barley breeding and related researches. *J Institute Brewing* **57**; 4, 247-260 (1951).
244

245    **Methods**

246    **Plant material and growth conditions**

247    Two populations of 2-row spring type barley were used to carry out the GWAS[31].   The first

248    population includes 211 elite lines grown in a polytunnel under field conditions in Dundee, Scotland.

249    For each line, 5 whole grains were ground to a fine powder using a ball mill (Mixer Mill MM400;

250    Retsch Haan Germany) and stored in dry conditions until the HPLC analysis. The second population

251    which was used for verification of the results of the analysis of the first subpopulation includes 128

252    elite lines grown in a glasshouse compartment in a mix of clay-loam and cocopeat (50:50 v/v) at

253    daytime and night time temperatures of 22°C and 15°C respectively in The Plant Accelerator,

254    Adelaide, Australia. As described previously, the collection of germplasm these populations are

255    sampled from has minimum population structure while maintaining as much genetic diversity as

256    possible[32]. Mature grains were stored until phenolic acid content analysis.

257    **Genotyping of SNP markers**

258    All lines were genotyped using the 50K iSelect SNP genotyping platform described previously[33]. Prior

259    to marker-trait association analysis, all markers with a minimum allele frequency of <5% and

260    markers with missing data >5% were excluded from the analysis.

261    **Phenotyping for cell wall-bound phenolic acids**

262     A ~ 20 mg amount of wholegrain barley was used per sample. *Trans*-ferulic and *trans*-*p*-coumaric

263     acid standards were purchased from SIGMA Aldrich (Castle Hill NSW, Australia). Standards were

264     prepared at 62.5 µm, 250 µm and 1000 µm by dissolving the appropriate amount of powder in 50%

265     methanol. Extraction of cell wall esterified phenolic acids was carried out following the methods

266     described by [34,35] with the following modifications. Samples were washed twice with 500 µl 80%

267     ethyl alcohol, with shaking for 10 minutes at room temperature to remove free phenolic acids. To

268     release total cell wall esterified phenolic acids, alkaline treatment was carried out by adding 600 µl

269     2M NaOH to the pellet. Samples were incubated on a rotary rack under nitrogen for 20 h in the dark

270     at room temperature. Samples were centrifuged at 15000 x g for 15 minutes at room temperature,

271     after which the supernatant was collected, acidified by adding 110 µl concentrated HCL and

272     extracted three times with 1 mL ethyl acetate. Following each extraction, samples were centrifuged

273     at 5000 x g for 7 minutes and the organic solution was collected. Extracts were combined,

274     evaporated to dryness in a rotary evaporator and dissolved in 100 µl of 50 % methanol prior to

275     injecting 40 µl into the HPLC column. For each sample two technical replicates were applied.

276

277     **HPLC conditions**

278     An Agilent Technologies 1260 Infinity HPLC equipped with a Diode Array detector was used. Samples

279     were analysed on an Agilent Poroshell 120 SB-C18 3.0x100mm 2.7- micron column kept at 30 C˚.

280     Eluents were A (0.5mM trifluoroacetic acid) and B (0.5mM trifluoroacetic acid, 40% methanol, 40%

281     acetonitrile, 10% water). Starting conditions were 85% A and 15% B. Flow rate was 0. 7 mL/min.

282     Eluting gradients were as follow; min 0-10: 15% to 55% B, min 11-12: column washed with 100% B,

283     min 13 back to the starting condition (85% A and 15% B). Detection was carried out at 280 nm and

284     spectral data was collected from 200 to 400 nm when required. Ferulic and *p*-coumaric acid peaks

285     were identified by comparing retention times and spectra to their corresponding standards. The

286     area under the peaks was quantified at 280 nm for *trans* forms.

287

288     **GWAS analysis of grain alkaline extractable pCA and FA and FA:pCA ratio**

289     Marker- trait association analysis was carried out using R 2.15.3 (www.R-project.org) and performed

290     with a compressed mixed linear model[36] implemented in the GAPIT R package[37]. For phenotype

291     values, the mean values of the barley wholegrain total alkaline extractable *trans*-ferulic and *trans*- *p*-

292     coumaric acid (w/w) were used. To identify genes within intervals associated with our trait we used

293     [17]. We also used the ratio of FA:pCA as a trait in our GWAS analysis. The ratio between the two

294     compounds was log transformed i.e. log(FA:pCA) to provide a more normally distributed dataset.

295     When using ratios in GWAS, a significant increase in the *p*-gain statistic[15] (a comparison between the

296  lowest -log10(p) values of the individual compounds and the -log10(p) value of the ratio) indicates

297  that ratios carry more information than the corresponding metabolite concentrations alone.  A

298  significant p-gain identifies a biologically meaningful association between the individual compounds.

299  We used B/(2*α) to derive a critical value of $3.42 \times 10^5$ for the FDR-adjusted p-gain, where α is the

300  level of significance (0.05) and B the number of tested metabolite pairs[15]. Therefore, as we tested

301  two traits our threshold was $2 \times 10^1$ and our p-gain was above this threshold.

302  To identify local blocks of LD, facilitating a more precise delimitation of QTL regions Linkage

303  disequilibrium (LD) was calculated across the genome between pairs of markers using a sliding

304  window of 500 markers and a threshold of $R^2 < 0.2$ using Tassel v 5 [38].  We anchored markers that

305  passed FDR and represented initial borders of the QTL on 7H to the physical map and then expanded

306  this region using local LD derived from genome wide LD analysis as described above. When the

307  GWAS had not resulted in an association that passed the FDR we used the arbitrary threshold of -

308  LOG10(P) to define the initial border.  The SNP with the highest LOD score was used to represent the

309  QTL. After identification and Sanger sequencing of the candidate gene *HvAT10* the GWAS was

310  repeated including the allele present at the S309Stop as an additional marker.

311  **Bioinformatics and gene identification**

312  We used BARLEX[17] to identify gene models present with the QTL defined by our analysis and their

313  expression profile based on RNAseq data in 16 different tissues/ developmental stages

314

315  **Phylogenetic analysis of barley BAHD acyltransferases**

316  Coding sequences of all BAHD acyltransferases with the PFAM domain PF02458 from rice, barley and

317  *Brachypodium* were downloaded from the Ensembl Plants database (http://plants.ensembl.org/).

318  Sequences were aligned using the MUSCLE alignment function[39]available in the Geneious 9.1.4

319  (https://www.geneious.com). The translation alignment option was used. A neighbour-joining tree

320  was produced from the alignment. Barley genes within group A and B clades were identified,

321  realigned with their rice and *Brachypodium* orthologs and a maximum likelihood tree was produced

322  from the translation alignment of the sequences. The following settings were applied: substitution

323  model: General-Time-Reversible (GTR), branch support: bootstrap, number of bootstrap: 1000.

324

325  **Resequencing and genotyping of *HvAT10* in the main and supplemental set**

326  Aligning the translation of AK376450 to Os06g39390 allowed the identification of the putative

327  genomic sequence of *HvAT10.* We designed four pairs of primers, details of sequences and reaction

328  conditions are in Supplementary Data 5, to amplify the full length CDS using reaction volumes,

329  reagents, and conditions as described in[40]. To facilitate quick and efficient genotyping of large

330  numbers of cultivars we subsequently designed a KASP genotyping assay to a SNP at 430bp in

331  *HvAT10* (Supplementary Data 5). Reactions were performed in an 8.1 μL reaction volume, with 3 μL

332  H2O, 1 μL DNA (20ng/μl), 4 μL KASP genotyping master mix, and 0.11 μL of the KASP assay.

333  Box plots to demonstrate the contribution of the SNP at 436bp in *HvAT10* to variation in grain pCA

334  and FA content were produced using R 2.15.3 (www.R-project.org).  To test for identity by descent of

335  the *HvAT10* allele within the set of accessions using for the GWAS a dendrogram was constructed

336  using maximum likelihood using the genotypic data from the 9k-select array[32] in MEGA7[41] with

337  default settings except for including bootstrapping and visualised in FigTree (v.1.4.4)

338  http://tree.bio.ed.ac.uk/software/figtree/.

**Characterisation of diversity of *HvAT10*  in *H. spontaneum* from the fertile crescent and barley**

340  **landraces.**

341  DNA was extracted as described above from 76 *H. spontaneum* and 114 barley landraces from[26]. The

342  S309Stop SNP was PCR amplified and Sanger sequenced with primer pair 5 using conditions

343  described above. A dendrogram was constructed using maximum likelihood using 4000 exome

344  capture derived SNPs from[26] in MEGA7[41] with default settings except for including bootstrapping and

345  visualised in FigTree ( v.1.4.4) http://tree.bio.ed.ac.uk/software/figtree/.

346

**Genome wide $F_{ST}$ analysis**

348  The fixation index ($F_{ST}$) is a measure of genetic differentiation between groups of individuals.

349  Genome wide $F_{ST}$ was calculated by locus using GenAlEx 6.502[42,43] after dividing the accessions into

350  two populations based on their HvAT10 allele using all informative 50K iSelect markers.

351

**Phenotypic analysis of cultivars with wildtype vs *at10^STOP* allele**

353  We characterised mature grain morphology using from plants grown in a polytunnel under field

354  conditions in Dundee, Scotland as described above, over two years (2010 and 2011). Grain area,

355  width and length were quantified using the MARVIN Seed Analyzer (GTA Sensorik GmbH, 2013).

356  BLUPs calculated from this data using R 2.15.3 (www.R-project.org) were used in subsequent

357  comparisons between allelic groups.

358

**Data availability**

360  All sequences of *HvAT10* generated in this study are available from NCBI, accession numbers are

361  provided in **Supplementary Table 1**.

362

**References (for methods section)**

31. Oakey H, et al. Identification of crop cultivars with consistently high lignocellulosic sugar release requires the use of appropriate statistical design and modelling. *Biotechnol Biofuels* **6,** 185 (2013).

32. Comadran et al. Natural variation in a homolog of Antirrhinum CENTRORADIALIS contributed to spring growth habit and environmental adaptation in cultivated barley. *Nat Genet*. **44**:1388–1392(2012)

33. Bayer MM, et al. Development and Evaluation of a Barley 50k iSelect SNP Array. *Front Plant Sci* **8**(1792) (2017).

34. Hernanz D, et al. Hydroxycinnamic Acids and Ferulic Acid Dehydrodimers in Barley and Processed Barley. *Journal of Agricultural and Food Chemistry* **49**(10):4884-4888 (2001).

35. Irakli MN, Samanidou VF, Biliaderis CG, Papadoyannis IN: Development and validation of an HPLC-method for determination of free and bound phenolic acids in cereals after solid-phase extraction. *Food chemistry* **134**(3):1624-1632 (2012).

36. Zhang Z, et al: Mixed linear model approach adapted for genome-wide association studies. *Nat Genet*, **42**(4):355-360 (2010).

37. Lipka AE, et al. GAPIT: genome association and prediction integrated tool. *Bioinformatics* **28**(18):2397-2399 (2012).

38. Bradbury PJ, et al. TASSEL: software for association mapping of complex traits in diverse samples. *Bioinformatics* **23**(19):2633-2635 (2007).

39. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res*.**32**(5):1792-1797 (2004).

40. Houston K, et al. Analysis of the barley bract suppression gene *Trd1*. *Theor Appl Genet*, **125**(1):33-45 (2012).

41. Kumar S, Stecher G, Tamura K. MEGA7: molecular evolutionary genetics analysis version 7.0 for bigger datasets. *Mol Biol Evol* . **33**(7):1870–1874(2016).

42. Peakall, R. and Smouse P.E. GENALEX 6: genetic analysis in Excel. Population genetic software for teaching and research. *Mol Ecol Notes*. **6**, 288-295 (2006).

43. Peakall, R. and Smouse P.E. GenAlEx 6.5: genetic analysis in Excel. Population genetic software for teaching and research-an update. *Bioinformatics* **28**, 2537-2539 (2012).

401

**Author contributions**

403    RW, KH, RB, CH, Alan Little, designed experiments. KH, ASH, JL, Amy Learmonth, carried out

404    experiments. KH, ASH, Amy Learmonth, ML**,** Alan Little, JL analysed data.  The manuscript was

405    written by CH, KH, RW, RB, Amy Learmonth, ASH with contributions from all other authors.

406
**Ethics declarations**

408    The authors declare no competing interests.

**Tables**

410

411    **Table 1 - T-test results for comparisons between *HvAT10* alleles.** For grain area, length and width

412    data are available in Supplementary Data 1 and analysis was carried out using BLUPS derived from 2

413    – years' worth of samples. Data used for comparison of hot water extract, germinative energy,

414    fermentable extract, diastatic power, wort viscosity and friability between *HvAT10* alleles are

415    published [29].

416

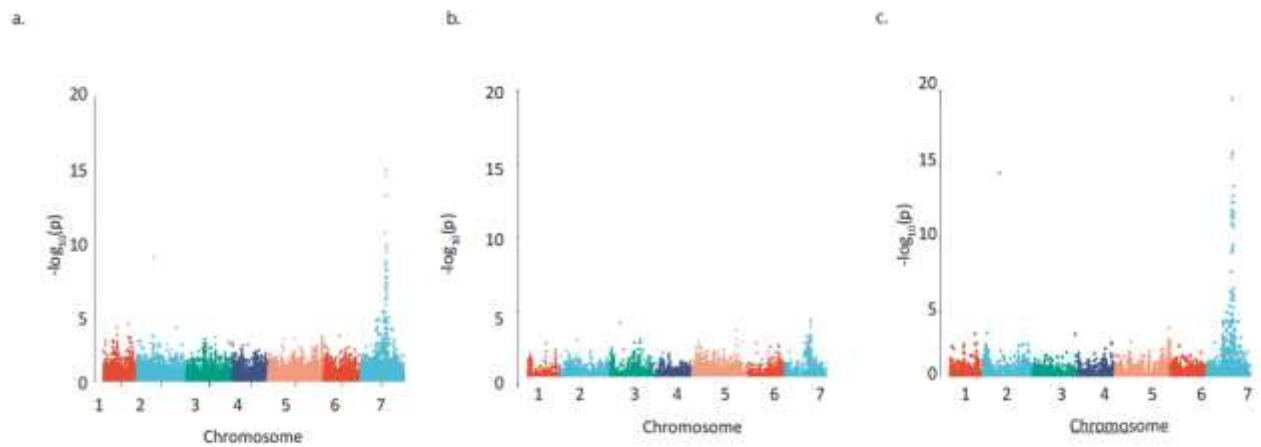| | Grain Area | Grain Length (mm) | Grain Width (mm) | Hot water extract (l°/Kg) | Germinative energy 4 ml (%) | Fermentable extract | Diastatic power IoB | Wort viscosity (mPa/s) | Friability (%) |
|---|---|---|---|---|---|---|---|---|---|
| *at10 STOP* | 27.67 | 9.02 | 3.92 | 309.44 | 97.34 | 70.72 | 100.77 | 1.5 | 87.97 |
| WT | 28.16 | 9.09 | 3.97 | 311.23 | 97.4 | 70.88 | 106.07 | 1.48 | 89.95 |
| *p* value | 0.0214* | 0.152 | 0.0009*** | 0.0005*** | 0.0206* | 0.0203* | 0.0282* | 0.0009*** | 0.0123* |

417
418

419 **Figures**

420

421 **Figure 1. Detecting regions of the barley genome associated with grain phenolic acid content using**
422 **a collection of 211 spring 2-row barleys.** Manhattan plots of the GWAS of the phenolic acid content
423 of wholegrain 2-row spring barley indicating regions of the genome associated with grain **a**. *p*-
424 coumaric acid, **b**. ferulic acid content, **c.** using the a ratio of these two phenolic acids calculated by
425 log[FA:p-Coumaric acid]. The –Log 10 (P-value) is shown on the Y axis, and the X axis shows the 7
426 barley chromosomes. The FDR threshold = –log 10(P)=6.02, plots use numerical order of markers on
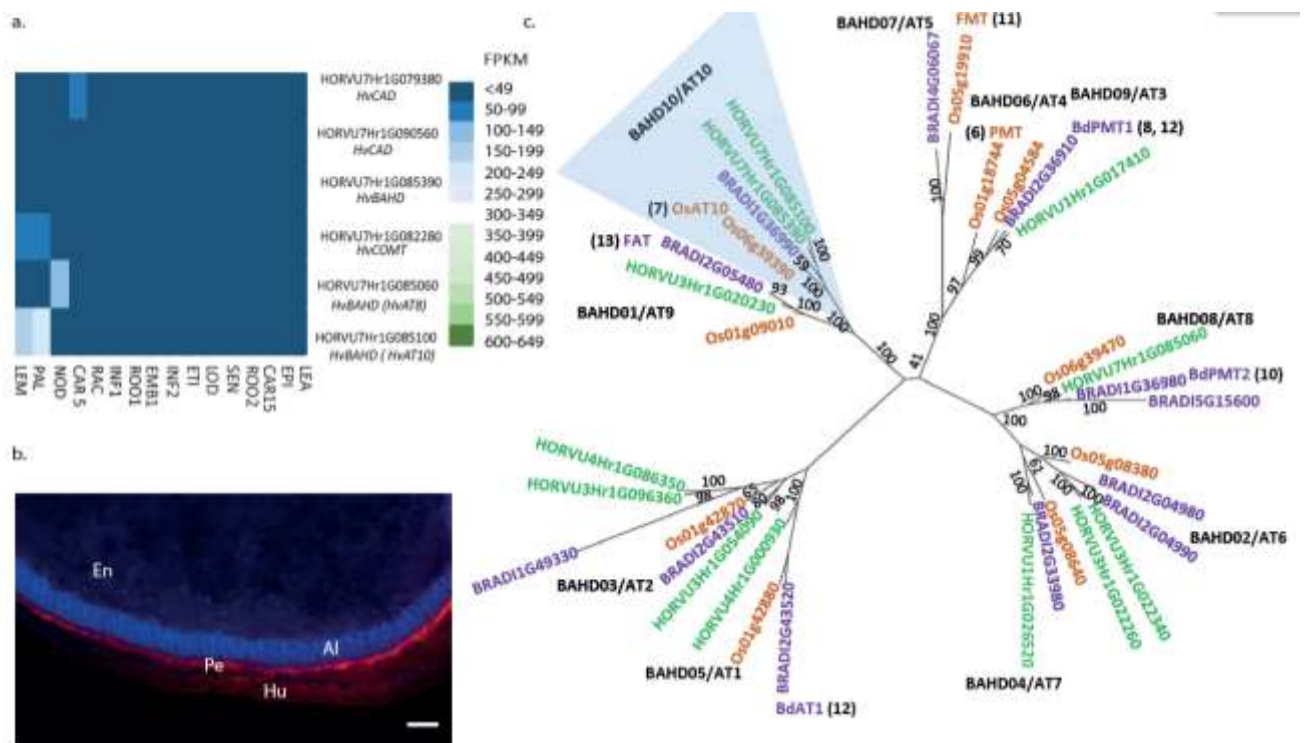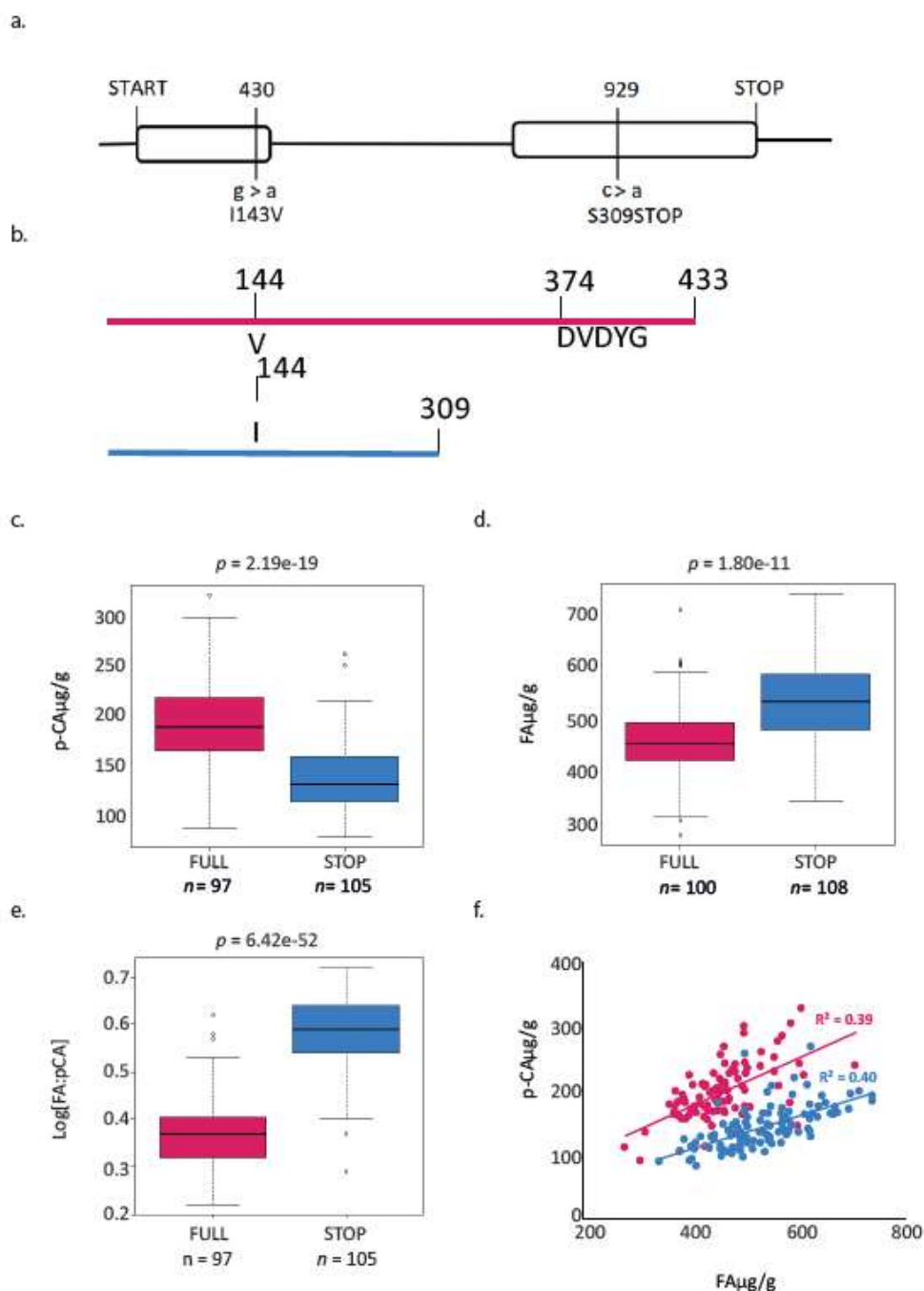427 the physical map.

428



429
430
431

432 **Figure 2. Putative candidates contributing to variation in grain *p*-coumaric (pCA) and Ferulic (FA)**
433 **content from GWAS. a**. Expression pattern for candidate genes under GWAS peak on 7H for grain *p*-
434 coumaric and ferulic acid content in 16 different tissues/ developmental stages. Values are FPKM
435 and a scale bar is provided. This expression data is derived from the publicly available RNAseq
436 dataset BARLEX, https://apex.ipk-gatersleben.de/apex/f?p=284:39 **b.** Phenolic acid autofluorescence
437 in whole grain sections. En: Endosperm, Al: Aleurone, Pe: Pericarp and Hu: Husk. Scale bar = 100µm
438 **c.** Phylogenetic tree of the *BAHD* acyltransferases**. A maximum-likelihood tree of the translation
439 alignment of the coding sequences of group A and B *BAHD* genes from barley, rice and
440 *Brachypodium*. Bootstrap support for branches is provided. Horvu numbers represent the barley
441 gene models in green, BRADI represents *Brachypodium* in purple, and *Os* represents the rice genes in
442 orange. The clade including LOC_OS06g39390 and HORVU7Hr1G085060 is highlighted in blue,
443 *OsAT10* is indicated in green and the closest barley orthologue is marked in red. Where function of a
444 gene model has been assigned the relevant reference is provided. Black text in bold indicates branch
445 names, both BAHD and AT[13].
446



447
448

449 **Figure 3. Gene and protein models for *HvAT10*. a.** Gene model for *HvAT10* including location, and
450 effect of SNPs detected from resequencing this gene in the 211 barley cultivars which have been
451 assayed for *p*-coumaric and ferulic acid. The numbering above the gene model represent locations in
452 the CDS which vary between these cultivars. The SNP, and the resulting change in the particular
453 amino acid are indicated underneath the gene model. The full length of the gene is 2117bp (with a
454 CDS of 1302bp) which translates to a protein of 435 amino acids as indicated. Protein model for
455 translation of HvAT10. **b.** a Full length protein and **c.** when the premature stop codon is present this
456 results in a truncated protein. Box plots demonstrate the effect of the SNP at 929bp within *HvAT10*
457 where the grain of the 211 barley cultivars were quantified for **d.** *p*-coumaric acid levels and **e.**
458 ferulic acid levels. **f.** Correlation between pCA and FA content based on *HvAT10* allele using 211
459 lines. The allele which results in full length version of HvAT10 are in pink, and the allele leading to a
460 premature stop codon are coloured blue.
461



462

**Supplementary Information**

**Supplementary Data 1. Phenolic acid and genetic data for all cultivars included in this study.** *p*-coumaric and ferulic acid content, KASP data and NCBI number for those lines that where sequenced for *HvAT10* is included.

**Supplementary Data 2. Summary of number of accessions used for each GWAS**. A total of 211 accessions were included in the main dataset but data for both phenolic acids is not available for all lines, therefore the number of individuals included in different analysis varies. Includes number of accessions for the GWAS presented in the main and supplementary analysis for both individual trait and the ratio analysis. Number of individuals with each allele of *HvAT10* based on genotyping of A430G is also included.

**Supplementary Data 3. Details of QTL identified on 7H for all analysis carried out.** Physical location, LOD score, and 50k iSelect marker with the highest LOD score are provided. * indicates that analysis passed the FDR threshold of -log10(p)=6.1.

**Supplementary Data 4. Gene models containing SNPs that have an $F_{ST}$>0.875 when $F_{ST}$ analysis carried out based on *HvAT10* allele.** This table includes 50k iSelect marker name, the chromosome the marker is located on, gene model and annotation based on Morex v1 Gene Models (2016).

**Supplementary Data 5. Details of primers and genotyping assays used in this study.** This includes details of primers for Sanger sequencing and KASP genotyping assay sequence for *HvAT10*.

**Supplementary Figure 1. Phenolic acid content of wholegrain flour from 211 2-row spring barleys linea. a**. *p*CA and **b**. ferulic acid content. Values represent the mean for FA and pCA expressed as w/w. Error bars represent standard deviation of the replicates.

**Supplementary Figure 2. Manhattan plots of the GWAS of the phenolic acid content of wholegrain flour from 128 2-row spring barley lines indicating regions of the genome associated with grain phenolic acid content.** Manhattan plots of the GWAS of the phenolic acid content of wholegrain 2-row spring barley indicating regions of the genome associated with grain **a**. *p*-coumaric acid, **b**. ferulic acid content and **c.** log[FA:p-Coumaric acid]. The –Log 10 (P-value) is shown on the Y axis, and the X axis shows the 7 barley chromosomes. FDR threshold = −log 10(P)=6.02, plots use numerical order of markers on the physical map.

**Supplementary Figure 3. Distribution of ratio between two phenolic acids quantified in the grain of 211 spring 2 row barleys lines and used to carry out GWAS**. The ratio was calculated as log[FA:p-Coumaric acid]. Accessions containing the allele which results in a full length version of *HvAT10* are in pink, and accessions containing the allele leading to a premature stop codon are coloured blue.

**Supplementary Figure 4. Distribution of the *HvAT10* premature stop codon in *H. vulgare* landraces and cultivated barley lines. a**. A dendrogram of 114 *H. vulgare* landraces constructed using a selection of SNPs with a genome-wide distribution with maximum likelihood methods. **b.** A dendrogram of cultivated barley germplasm using a selection of SNPs with a genome-wide distribution using maximum likelihood methods. Accessions containing the allele which results in full length version of HvAT10 are in pink, and accessions containing the allele leading to a premature stop codon are coloured blue.

**Supplementary Figure 4. $F_{ST}$ analysis based on *HvAT10*. a**. Plot displaying genome wide $F_{ST}$ with $F_{ST}$ index provided on the Y axis, an $F_{ST}$ of 1 indicating a complete fixation of each allele within the two subpopulations determined by their allele of *HvAT10*. **b**. Just $F_{ST}$ of markers at 7H. Red box indicates

514    location of the centromere. Two SNPs whose location overlap on this plot, including one in *HvAT10*,
515    have an $F_{ST}$ of 1.0.  Note shape of peak appears different in **a.** and **b.** due to the difference in scale of
516    the plots. **c**. RNAseq data for genes with $F_{ST}>0.875$ from 16 different tissues/ developmental stages.
517    Values are FPKM and a scale bar is provided.  This expression data is derived from the publicly
518    available RNAseq dataset BARLEX, https://apex.ipk-gatersleben.de/apex/f?p=284:39