

 Open access • Posted Content • DOI:10.1101/2020.09.23.310268

## **Genetic variation, environment and demography intersect to shape Arabidopsis defense metabolite variation across Europe — Source link**

Ella Katz, Clement Bagaza, Samuel Holden, Ruthie Angelovici ...+2 more authors

**Institutions:** University of California, Berkeley, University of Missouri, University of Copenhagen

**Published on:** 24 Sep 2020 - bioRxiv (Cold Spring Harbor Laboratory)

**Topics:** Population

Related papers:

- [Explaining intraspecific diversity in plant secondary metabolites in an ecological context.](#)
- [The scales and signatures of climate adaptation by the Arabidopsis transcriptome](#)
- [Yeast growth responses to environmental perturbations are associated with rewiring of large epistatic networks](#)
- [Evolution in heterogeneous environments and the potential of maintenance of genetic variation in traits of adaptive significance.](#)
- [Evolutionary insights from studies of geographic variation: Contemporary variation and looking to the future.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/genetic-variation-environment-and-demography-intersect-to-2qsuacl9qu>

1 **Genetic variation, environment and demography intersect to shape**  
2 ***Arabidopsis* defense metabolite variation across Europe**

3 Ella Katz<sup>1</sup>, Clement Bagaza<sup>2</sup>, Samuel Holden<sup>2</sup>, Ruthie Angelovici<sup>2</sup>, Daniel J. Kliebenstein<sup>1,3x</sup>

4 <sup>1</sup>Department of Plant Sciences, University of California, Davis, One Shields Avenue, Davis, CA,  
5 95616, USA

6 <sup>2</sup>Division of Biological Sciences, Interdisciplinary Plant Group, Christopher S. Bond Life  
7 Sciences Center, University of Missouri, Columbia, Missouri 65211

8 <sup>3</sup>DynaMo Center of Excellence, University of Copenhagen, Thorvaldsensvej 40, DK-1871,  
9 Frederiksberg C, Denmark

10 <sup>x</sup>Corresponding Author: [Kliebenstein@ucdavis.edu](mailto:Kliebenstein@ucdavis.edu)

11

12 **Abstract**

13 Plants face a variety of challenges within their ever-changing environment. Diverse metabolites  
14 are central to the plants ability to overcome these challenges. Understanding the environmental  
15 and genetic factors influencing the variation in specialized metabolites is the key to understand  
16 how plants survive and develop under changing environments. Here we measure the variation in  
17 specialized metabolites across a population of 797 natural *Arabidopsis thaliana* accessions. We  
18 show a combination of geography, environmental parameters, demography, and different genetic  
19 processes that creates a specific pattern in their accumulation and distribution. By identifying and  
20 tracking causal polymorphisms at multiple loci controlling metabolites variation we show that  
21 each locus displays extensive allelic heterogeneity with signatures of both parallel and  
22 convergent evolutionary processes. These loci combine epistatically and show differing  
23 relationships to environmental parameters leading to different distributions. This provides a  
24 detailed perspective about the complexity of the forces and mechanisms that shape the  
25 accumulation and distribution of a family of specialized metabolites critical for plant fitness.

26

27

28

## 29 **Introduction**

30 The biotic and abiotic components of a plant's habitat/environment are continuously changing.  
31 This creates a complex system to which a plant must develop adaptation strategies to ensure  
32 survival and reproduction. Metabolites are frequent keys to these strategies, involving the  
33 production and accumulation of different metabolites from signaling hormones, primary  
34 metabolites and a wide array of multi-functional specialized metabolites (Erb & Kliebenstein,  
35 2020; Hanower & Brzozowska, 1975; Hayat et al., 2012; Kim et al., 2012; D J Kliebenstein,  
36 2004; Malcolm, 1994; Thakur & Rai, 1982; Wolters & Jürgens, 2009; Yang, Lin, & Kao, 2000).  
37 The complete suite of these metabolites helps to determine the plants survival and development.  
38 A complication in the plants ability to create an optimal blend of metabolite-based strategies is  
39 the fact that individual specialized metabolites can have contrasting effects in a complex  
40 environment. For example, individual specialized metabolites can provide defense against some  
41 attackers while simultaneously causing sensitivity to other biotic attackers or abiotic stresses  
42 (Agrawal, 2000; Bialy, Oleszek, Lewis, & Fenwick, 1990; Erb & Kliebenstein, 2020; Futuyma  
43 & Agrawal, 2009; Hu et al., 2018; Lankau, 2007; Opitz & Müller, 2009; Uremis, Arslan,  
44 Sangun, Uygur, & Isler, 2009; Züst & Agrawal, 2017). This creates offsetting ecological benefits  
45 and costs for individual metabolites that when summed across all the metabolites means that  
46 there are complex selective pressures driving the differentiation of metabolic profiles within a  
47 species and shaping genetic variation within and between populations depending on the diverse  
48 challenges faced (Fan, Leong, & Last, 2019; R. Kerwin et al., 2015; Malcolm, 1994; Sønderby,  
49 Geu-Flores, & Halkier, 2010; Szakiel, Pączkowski, & Henry, 2011; Wentzell & Kliebenstein,  
50 2008; Züst et al., 2012).

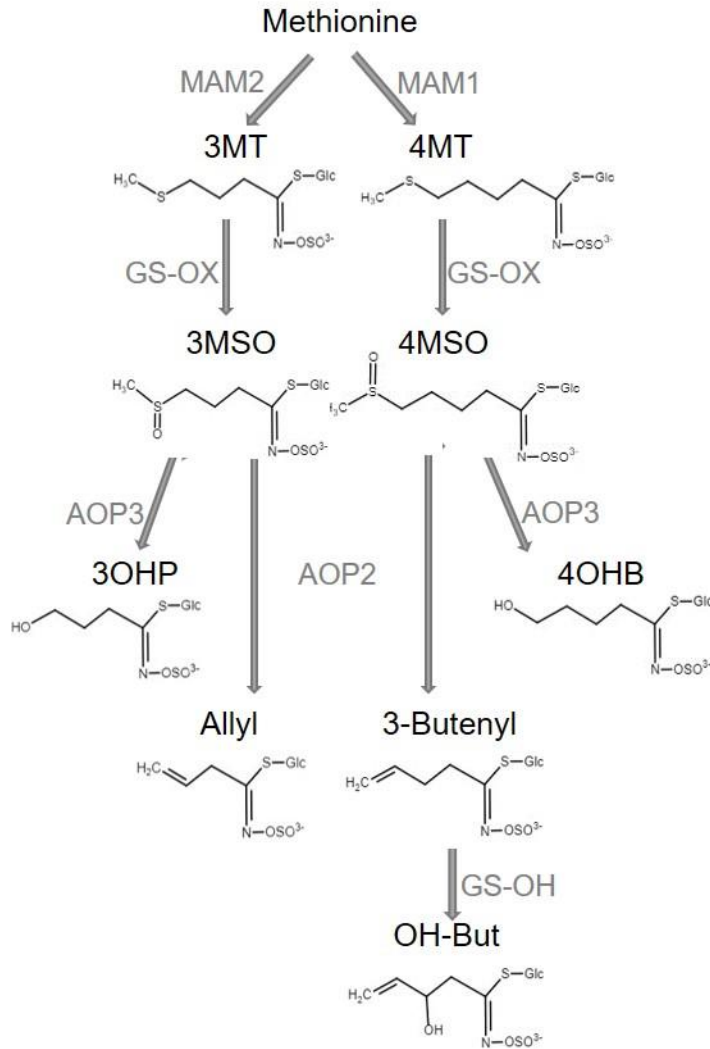
51 Recent decades have seen significant advances in the identification of the genetic variation  
52 creating this metabolic variation. A common theme developing from these studies is that the  
53 metabolic variation within and between species is the result of structural variation at the enzymes  
54 responsible for the chemical structures (Chan, Rowe, Corwin, Joseph, & Kliebenstein, 2011;  
55 Chan, Rowe, & Kliebenstein, 2010; Fan et al., 2019; Kroymann, Donnerhacke, Schnabelrauch,  
56 & Mitchell-Olds, 2003; Moore et al., 2019; Schillmiller, Pichersky, & Last, 2012). These  
57 structural variants and the resulting chemical variation strongly influence plant fitness in  
58 response to a broad range of biotic interactions including at least herbivores, but also other plant

59 species and other members of the same plant species (Bednarek & Osbourn, 2009; Brachi et al.,  
60 2015; R. E. Kerwin et al., 2017; R. Kerwin et al., 2015; Lankau & Kliebenstein, 2009; Lankau &  
61 Strauss, 2007, 2008; Lankau, 2007). Most mechanistic studies of natural variation in specialized  
62 metabolism have focused on apparent biallelic phenotypic variation linked to loss-of-function  
63 variants. However, it is not clear if biallelic genetic causation is true when extended to a large  
64 collection of individuals from wide-ranging populations within a species. If selective pressures  
65 are sufficiently strong and non-linear, it is possible to have repeated and independent generation  
66 of structural variants creating the same metabolic variation. This raises the possibility for  
67 chemical variation within a species to show hallmarks of parallel evolution, wherein  
68 phenotypically similar variants independently arise from the same genetic background. Equally it  
69 may be possible to find within-species convergent evolution, where different allele with identical  
70 metabolic consequences arise from independent genetic backgrounds through different  
71 mechanisms. Because these genetic processes are occurring simultaneously with neutral  
72 demographic processes like migration, there is a need to better understand how the intersection  
73 of environmental pressure, demography and genomic complexity gives rise to the pattern of  
74 metabolic variation across a plant species.

75 To better understand how genomic variation, demography and environmental pressures shape the  
76 variation of specialized metabolism within a species, we used the model Glucosinolates (GSLs)  
77 pathway. GSLs are a diverse class of specialized metabolites produced in the order Brassicales,  
78 including the model plant *Arabidopsis* (*Arabidopsis thaliana*), that show extensive variation  
79 between and within species across the order (Bakker, Traw, Toomajian, Kreitman, & Bergelson,  
80 2008; Benderoth et al., 2006; Brachi et al., 2015; Chan et al., 2010; Daxenbichler et al., 1991;  
81 Halkier & Gershenzon, 2006; R. Kerwin et al., 2015; D J Kliebenstein, Gershenzon, & Mitchell-  
82 Olds, 2001; D J Kliebenstein, Kroymann, et al., 2001; D J Kliebenstein, Lambrix, Reichelt,  
83 Gershenzon, & Mitchell-Olds, 2001; James E. Rodman, Kruckeberg, & Al-Shehbaz, 1981;  
84 James Eric Rodman, 1980; Sønderby et al., 2010; Wright, Lauga, & Charlesworth, 2002). GSLs  
85 consist of a common core structure with a highly diverse side chain that determines the GSLs  
86 biological activity in defense, growth, development and abiotic stress resistance (Beekwilder et  
87 al., 2008; Hansen et al., 2008; Hasegawa, Yamada, Kosemura, Yamamura, & Hasegawa, 2000;  
88 Katz et al., 2020; Katz, Nisani, Sela, Behar, & Chamovitz, 2015; Malinovsky et al., 2017;  
89 Salehin et al., 2019; Yamada et al., 2003). The *Arabidopsis*-GSL system is an optimal model to

90 study the species wide processes driving specialized metabolite variation because the identity of  
91 the whole biosynthetic pathway is known, including the major causal loci for natural variation  
92 (Benderoth et al., 2006; Brachi et al., 2015; Chan et al., 2011, 2010; Hansen, Kliebenstein, &  
93 Halkier, 2007; D J Kliebenstein, Gershenzon, et al., 2001; D. Kliebenstein, Pedersen, Barker, &  
94 Mitchell-Olds, 2002; Daniel J Kliebenstein, Figuth, & Mitchell-Olds, 2002; Kroymann &  
95 Mitchell-Olds, 2005; Pfalz, Vogel, Mitchell-Olds, & Kroymann, 2007; S nderby et al., 2010;  
96 Wentzell et al., 2007). These major loci, have been proven to influence Arabidopsis fitness and  
97 can be linked to herbivore pressure (Brachi et al., 2015; Hansen et al., 2008; Jander, Cui, Nhan,  
98 Pierce, & Ausubel, 2001; R. E. Kerwin et al., 2017; R. Kerwin et al., 2015; Z st et al., 2012).  
99 Beyond the major causal loci, there is also evidence from genome wide association studies for  
100 highly polygenic variation in the genetic background that further contributes to modulating GSL  
101 variation (Chan et al., 2011). The public availability of over 1000 widely distributed accessions  
102 with genomic sequences provides the ability to phenotype GSL variation across a large spatial  
103 scale and query the distribution and relationship of causal haplotypes at the major GSL causal  
104 loci.

105 In Arabidopsis and other Brassicas, the main GSLs are Methionine-derived, Aliphatic, GSLs.  
106 Genetic variation in Aliphatic GSLs structure is controlled by natural variation at three loci, GS-  
107 Elong, GS-AOP and GS-OH with these three loci combining to create a dominant Aliphatic GSL  
108 chemotype. In addition to these expressed loci, there is a large suite of loci that can modify these  
109 dominant patterns (Brachi et al., 2015; Chan et al., 2011, 2010). GS-Elong differentially  
110 elongates the Methionine side chain by structural variation influencing the expression of  
111 divergent methylthioalkylmalate synthase enzymes (MAM) that add carbons to the side chain  
112 (Abrahams, Pires, & Schranz, 2020). In Arabidopsis, MAM2 catalyzes the addition of two  
113 carbons to the side chain, creating GSLs with 3 carbon side chains. MAM1 catalyzes the addition  
114 of three carbons to make GSLs with 4 carbon side chains (Figure 1). MAM3 (also known as  
115 MAM-L) catalyzes the addition of up to 6 carbons (D J Kliebenstein, Lambrix, et al., 2001;  
116 Kroymann et al., 2003; Mithen, Clarke, Lister, & Dean, 1995). The core pathway leads to the  
117 creation of the methylthio GSL (MT). Then, the MT will be converted to a methylsulfinyl  
118 (MSO) with a matching number of carbons (Giamoustaris & Mithen, 1996; Hansen et al., 2007).  
119 Structural variation at the GS-AOP locus leads to differential modification of the MSO by  
120 differential expression of a family of 2-oxoacid-dependent dioxygenases (2ODD). The AOP2



121

122 **Figure 1: Aliphatic GSL biosynthesis pathway.** Short names and structures of the GSLs are in  
 123 black. Genes encoding the causal enzyme for each reaction (Arrow) are in grey. GS-OX is a gene  
 124 family of five or more genes. OH-But= 2-OH-3-Butenyl.

125

126 enzyme removes the MSO moiety leaving an alkenyl sidechain, while AOP3 leaves a hydroxyl  
 127 moiety. Previous work has suggested three alleles of GS-AOP: AOP3 expressing, AOP2  
 128 expressing and a null allele (i.e. Col-0 and similar accessions) with nonfunctional copies of  
 129 AOP2 and AOP3 leading to MSO accumulation, the AOP substrate (Figure 1) (Chan et al., 2010;  
 130 D J Kliebenstein, Kroymann, et al., 2001; D J Kliebenstein, Lambrix, et al., 2001; Mithen et al.,  
 131 1995). The 4C alkenyl side-chain can be further modified by adding a hydroxyl group at the 2C  
 132 via the GS-OH 2-ODD (Figure 1) (Hansen et al., 2008). In spite of the evolutionary distance,  
 133 independent variation at the same three loci influence the structural diversity in Aliphatic-GSLs

134 within Brassica, Streptanthus and Arabidopsis (D J Kliebenstein & Cacho, 2016; Lankau &  
135 Kliebenstein, 2009). For example the C3 MAM in Arabidopsis and Brassica represent two  
136 independent lineages as are the MAMs responsible for C4 GSLs, in fact the MAM locus contains  
137 at least three independent lineages that recreate the same length variation (Abrahams et al.,  
138 2020). This indicates repeated evolution across species, but it is not clear how frequently these  
139 loci are changing within a single species or how ecological or demographic processes may shape  
140 within-species variation at these loci.

141 In this work we described GSL variation in seeds of a collection of 797 *Arabidopsis thaliana*  
142 natural accessions collected from locations across Europe. The amounts of GSLs can vary across  
143 different tissues and life stages, but there is a strong correlation in the type of Aliphatic GSL  
144 produced across tissues because the major causal effect loci are not plastic due to structural  
145 variation (Brown, Tokuhisa, Reichelt, & Gershenzon, 2003; D J Kliebenstein, Gershenzon, et al.,  
146 2001; D J Kliebenstein, Kroymann, et al., 2001; Petersen, Chen, Hansen, Olsen, & Halkier,  
147 2002). Thus, the seeds chemotype is the same as the leaves. Further, seeds have the highest level  
148 of GSLs in Arabidopsis and they are stable at room temperature until germination, hence making  
149 seeds a perfect tissue to survey variation. Further, GSLs are known to be important for seed  
150 defenses against herbivores and pathogens (Raybould & Moyes, 2001). By measuring GSLs in  
151 seeds, we identified that the distribution of GSLs and the causal alleles are influenced by a  
152 diverse set of factors with a primary contribution by geography, environmental parameters and  
153 their interaction. We describe here the complex genetic architecture of the three main causal loci  
154 responsible for the GSL composition, show how it effects the actual phenotype, and how it  
155 evolved. Interestingly, the combination of these elements reveal that several evolutionary  
156 mechanisms are involved in shaping the GSL variation and distribution, and the integration of  
157 them result in the pattern described here.

158

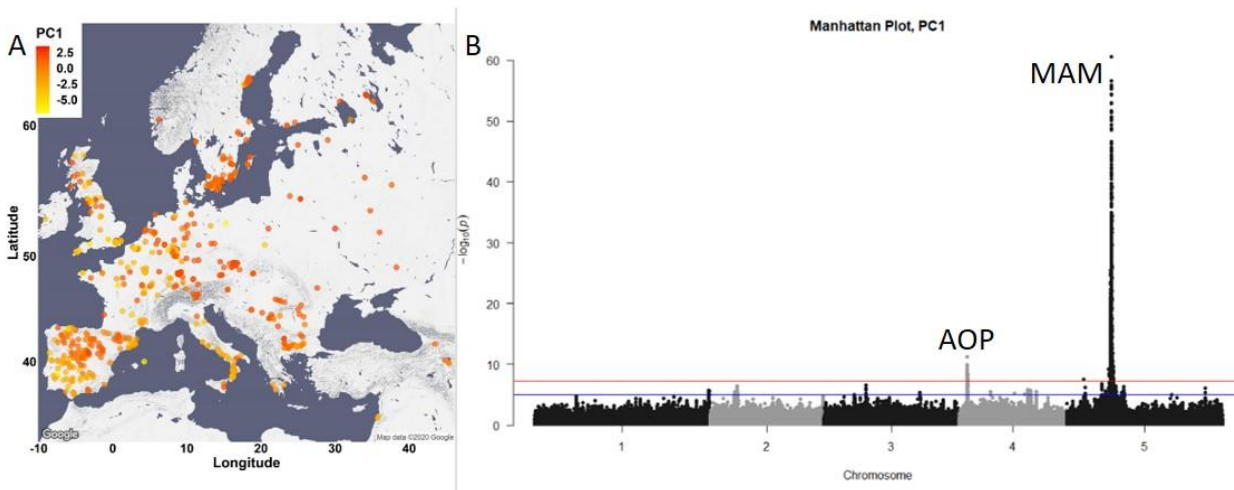
## 159 **Results**

### 160 **GSL variation across Europe**

161 To investigate the genetic, environmental and demographic parameters influencing the  
162 distribution of Arabidopsis GSL chemotypes, we measured GSLs from seeds of a collection of



163 797 *Arabidopsis thaliana* natural accessions (The 1001 Genomes Consortium, 2016). These  
164 *Arabidopsis* accessions were collected from different geographical locations, mainly in and  
165 around Europe. 23 different GSLs were detected and quantified identifying a wide diversity in  
166 composition and amount among the natural accessions with a median heritability of 83%,  
167 ranging from 34% to 93% (Supplemental Table 1). To summarize the GSL variation among the  
168 accessions we performed principal component analyses (PCA) on the accumulation of all the  
169 individual GSLs across the accessions as an unbiased first step. The first two PCs only captured  
170 33% of the total variation with PC1 describing GSLs with 4 and 7 carbons and PC2 mainly  
171 capturing GSLs with 8 carbons in their side chain (supp. Figure 1). Previous work using a  
172 collection of predominantly central European accessions had suggested a simple continental  
173 gradient chain-elongation variation from the south-west to the north-east (Brachi et al., 2015;  
174 Züst et al., 2012). To assess if this was still apparent in this larger collection, we plotted the  
175 accessions based on their geographical locations, and colored them based on their PC1 and PC2  
176 scores that are linked to chain elongation variation (Figure 2A and supp. Figure 2A,  
177 respectively). This larger collection shows that there is not a single gradient shaping GSL  
178 diversity across Europe (Figure 2A). Instead the extended sampling of accessions around the  
179 Mediterranean in this collection shows that the SW to NE pattern reiterates within the Iberian  
180 Peninsula.



181

182 **Figure 2: GSL variation across Europe is dominated by two loci.** A. The accessions are  
183 plotted on the map based on their collection site, and colored based on their PC1 score. B.  
184 Manhattan plot of GWAS analyses using PC1. Horizontal lines represent 5% significance  
185 thresholds using Bonferroni (red) and permutations (blue).



186 To test which of the major causal loci are detectable in this collection and to identify new  
187 genomic regions that are associated with the observed GSL variation, we performed genome  
188 wide association (GWA, with EMMAX algorithms) analyses using the PC1 and PC2 values.  
189 This collection of natural accessions presents a dense variant map and is 3x larger than previous  
190 GSL GWA mapping populations. In spite of the large population size, both PC1 and PC2 based  
191 analyses identified the same two major peaks covering two of the known causal genes  
192 controlling GSL diversity (Figure 2B for PC1 GWA analyses, supp. Figure 2B for PC2 GWA  
193 analyses) (Brachi et al., 2015; Chan et al., 2011, 2010). The largest peak in both cases, is the GS-  
194 Elong locus on chromosome 5, containing the MAM1 (AT5G23010), MAM2 and MAM3  
195 (AT5G23020) genes. The peak on chromosome 4 is the GS-AOP locus containing the AOP2 and  
196 AOP3 genes (AT4G03060 and AT4G03050, respectively). Previous F2, QTL and molecular  
197 experiments have shown that the genes within GS-AOP and GS-Elong loci are the causal genes  
198 for GSL variation within these regions (Benderoth et al., 2006; Brachi et al., 2015; Chan et al.,  
199 2011, 2010; D J Kliebenstein, Gershenzon, et al., 2001; D. Kliebenstein et al., 2002; Daniel J  
200 Kliebenstein et al., 2002; Kroymann & Mitchell-Olds, 2005; Pfalz et al., 2007; Wentzell et al.,  
201 2007). Surprisingly, none of the 8 other known natural variants within the GSL biosynthetic  
202 pathway were identified by GWA including three that were found with 96 accessions and three  
203 that were found with 595 accessions using PC1 and 2 (Brachi et al., 2015; Chan et al., 2011,  
204 2010; Daniel J. Kliebenstein, 2009). It is possible that the extended sampling of accessions may  
205 have created genomic and demographic issues that influenced this high false-negative error rate  
206 where ~80% of validated natural variants found using multiple RIL populations were missed.

207

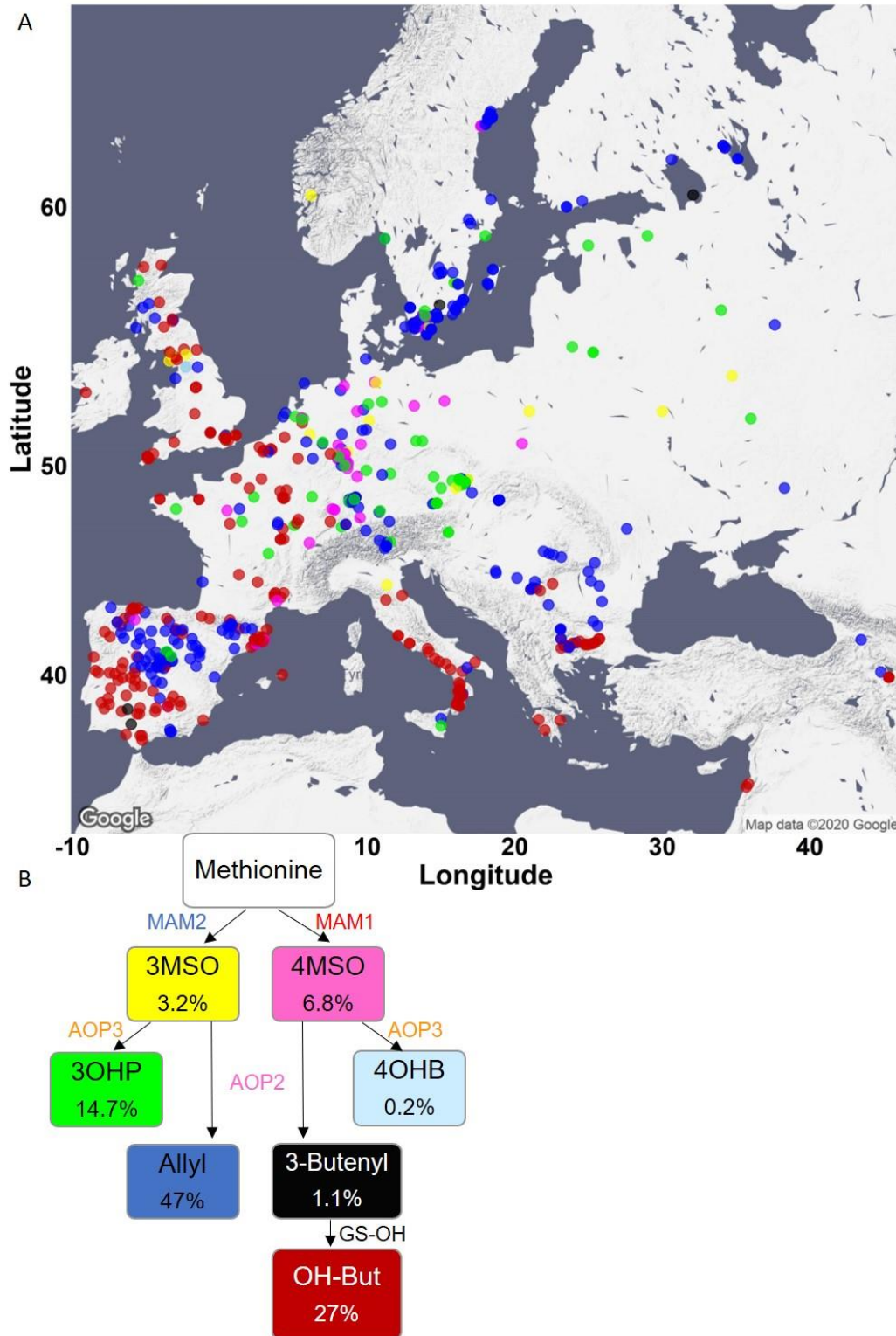
## 208 **Complex GSL Chemotypic Variation**

209 One potential complicating factor is that GSL chemotypic variation is best described as a discrete  
210 multimodal distribution involving the epistatic interaction of multiple genes which PCA's linear  
211 decomposition cannot accurately capture (Figure 1). To test if PCA was inaccurately describing  
212 GSL chemotypic variation, we directly called the specific GSL chemotypes in each accession.  
213 Using Arabidopsis QTL mapping populations and GWA, we have shown that the GS-AOP,  
214 Elong and OH loci determine seven discrete chemotypes, 3MSO, 4MSO, 3OHP, 4OHB, Allyl,  
215 3-Butenyl, 2-OH-3-Butenyl, that can be readily assigned from GSLs phenotypic data (Brachi et

216 al., 2015; Chan et al., 2011, 2010; D J Kliebenstein, Gershenzon, et al., 2001). Using accessions  
217 with previously known chemotypes and genotypes, we developed a phenotypic classification  
218 scheme to assign the chemotype for each accession (Figure 3, for details see methods and supp.  
219 Figures 3-5, for structures see Figure 1 and supp. Table 1). Since the Aliphatic GSLs  
220 composition in the seeds reliably indicate the GSL structural composition in the other plant's life  
221 stages and tissues, assigning a chemotype for each accession based on the seeds composition is  
222 expected to be highly stable across tissues of the same accession (Brown et al., 2003; Chan et al.,  
223 2011, 2010; D J Kliebenstein, Gershenzon, et al., 2001; D J Kliebenstein, Kroymann, et al.,  
224 2001). Most accessions were classified as 2-OH-3-Butenyl (27%) or Allyl (47%) with lower  
225 frequencies for the other chemotypes. Mapping the chemotypes on Europe showed that the PCA  
226 decomposition was missing substantial information on GSL chemotype variation (Figure 3).  
227 Instead of a continuous distribution across Europe, the chemotype classifications revealed  
228 specific geographic patterns. Central and parts of northern Europe were characterized by a high  
229 variability involving the co-occurrence of individuals from all chemotypes. In contrast, southern  
230 Europe, including the Iberian Peninsula, Italy and the Balkan, has two predominant chemotypes,  
231 Allyl or 2-OH-3-Butenyl, that are separated by a sharp geographic partitioning (Figure 3, and  
232 supp. Figure 6). The few accessions in southern Europe belonging to other chemotypes were all  
233 accessions previously identified as having genomes identical to accessions in central Europe,  
234 suggesting that they are likely stock center seed contaminations (The 1001 Genomes  
235 Consortium, 2016). Uniquely, Swedish accessions displayed a striking presence of almost solely  
236 Allyl chemotypes that was not mirrored on the eastern coast of the Baltic Sea (Finnish,  
237 Lithuanian, Latvian or Estonian accessions). Directly assigning GSL variation by discrete  
238 chemotypes provided a more detailed image not revealed by PCA decomposition. Further, the  
239 different chemotypic to geographic patterns suggests that there may be different pressures  
240 shaping GSL variation particularly when comparing central and southern Europe.

241 Figure 3: **Phenotypic classification based on GSL content.** A. Using the GSL accumulation,  
242 each accession was classified to one of seven aliphatic short chained GSL chemotypes based on  
243 the enzyme functions as follows: MAM2, AOP null: classified as 3MSO dominant, colored in  
244 yellow. MAM1, AOP null: classified as 4MSO dominant, colored in pink. MAM2, AOP3:  
245 classified as 3OHP dominant, colored in green. MAM1, AOP3: classified as 4OHB dominant,  
246 colored in light blue. MAM2, AOP2: classified as Allyl dominant, colored in blue. MAM1,  
247 AOP2, GS-OH non-functional: classified as 3-Butenyl dominant, colored in black. MAM1,  
248 AOP2, GS-OH functional: classified as 2-OH-3-Butenyl dominant, colored in red. The

249 accessions were plotted on a map based on their collection sites and colored based on their  
 250 dominant chemotype. B. The coloring scheme with functional GSL enzymes in the aliphatic  
 251 GSL pathway is shown with the percentage of accessions in each chemotypes (out of the total  
 252 797 accessions) shown in each box.



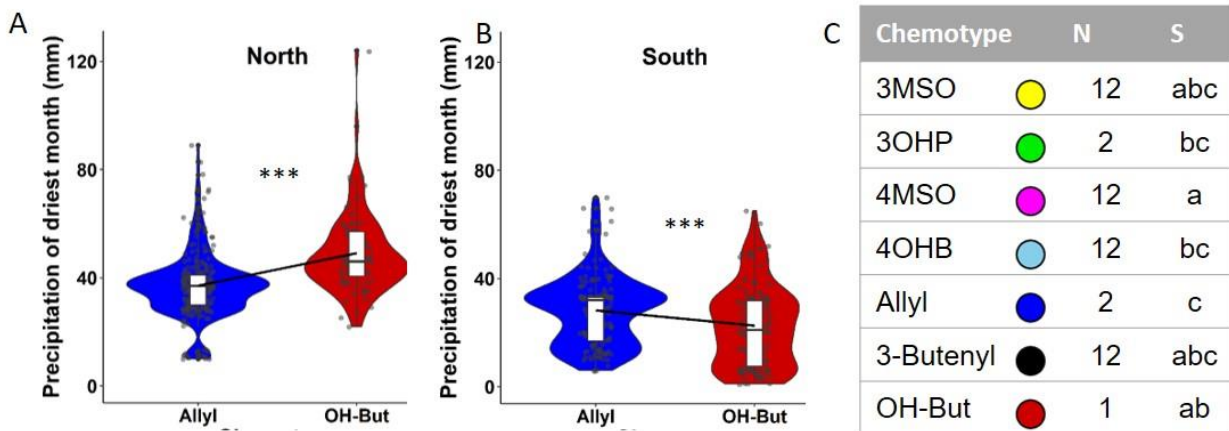
253

## 254 **Geography and environmental parameters affect GSL variation**

255 Because GSL chemotypes may be more reflective of local environment, we proceed to test if  
256 they are associated with weather parameters and landscape conditions. Further, given the  
257 difference in chemotype occurrence in central and southern Europe we hypothesized that these  
258 environmental connections may change between central and southern Europe. For these tests, we  
259 chose environmental parameters that capture a majority of the environmental variance and by  
260 that may describe the type of ecosystem (Ferrero-Serrano & Assmann, 2019). We assigned each  
261 accession the environmental value based on its location. These environmental parameters include  
262 geographic proximity (distance to the coast), precipitation descriptors (precipitation of wettest  
263 and driest month) and temperature descriptors (maximal temperature of warmest month and  
264 minimal temperature of coldest month) capture major abiotic pressures as well as provide  
265 information about the type of ecosystem in which each accession exists. We ran a  
266 multivariate analysis of variance (MANOVA) for each geographic area separately (north and  
267 central vs south, as shown in supp. Figure 6). This showed significant difference in how the GSL  
268 chemotypes associated to the environmental parameters across Europe. This was best illustrated  
269 by the two dominant chemotypes, Allyl and 2-OH-3-Butenyl, showing opposing relationships to  
270 the precipitation in the driest month. In Northern and Central Europe, the Allyl chemotype is  
271 more associated with lower precipitation in the driest month, while accessions with 2-OH-3-  
272 Butenyl as the dominant chemotype are associated with higher precipitation in the driest month.  
273 In Southern European accessions, this association is inverted (Figure 4A,B). This suggests that  
274 the relationship of GSL chemotype to environmental parameters vary across geographic regions  
275 of Europe rather than fitting a simple linear model.

276 As the two main chemotypes in the collection differ by the length of the carbon chain (C3 for  
277 Allyl, C4 for 2-OH-3-Butenyl), we created a linear model to further check the interaction  
278 between each environmental condition to geography in respect to the carbon chain length. Most  
279 of the environmental parameters significantly interacted with geography, meaning that the  
280 relationship of environment to GSL alleles change across geographic areas (supp. Figure 7, for  
281 details on the models see methods). Conducting this analysis for each of the geographic areas  
282 separately highlighted this by showing that these parameters have different effects on the carbon  
283 chain length in each of the areas (supp. Figure 7). This was true when the model was run with or

284 without ancestral population state being included in the model (The 1001 Genomes Consortium,  
285 2016).



286

287 **Figure 4: Environmental conditions differentially associate with GSLs across geographic**  
 288 **location.** A. B. The association of the two major chemotypes allyl and 2-OH-3-Butenyl to  
 289 precipitation values of the driest month. Significance was tested by t-Test,  $P = 0.00000258$  for  
 290 the North (Slope= 0.01),  $P = 0.0005521$  for the South (Slope= -0.007). C. MANOVA was  
 291 performed for the south and north as indicated in methods section, followed by pairwise  
 292 comparisons of least-squares means. Numbers indicate chemotypes with significant differences  
 293 in the North and letters indicate chemotypes with significant differences in the south. OH-But=  
 294 2-OH-3-Butenyl.

295

296 Using a random forest machine learning approach provided similar results with different  
 297 environment to chemotype relationships in the north and south (sup. Figure 8), supporting the  
 298 hypothesis that the GSL chemotype to environment relationships change across regions within  
 299 Europe.

300

### 301 **The genetic architecture of GSL variation**

302 The presence of different GSL chemotype to environmental relationships across Europe raises  
 303 the question of how these chemotypes are generated. Are these chemotypes from locally derived  
 304 alleles or obtained by the intermixing of widely distributed causal alleles. Further, if there are  
 305 multiple alleles, do they display within species convergent or parallel signatures. We focus on  
 306 the GS-AOP, GS-Elong, and GS-OH loci, the causal genes creating Arabidopsis GSL

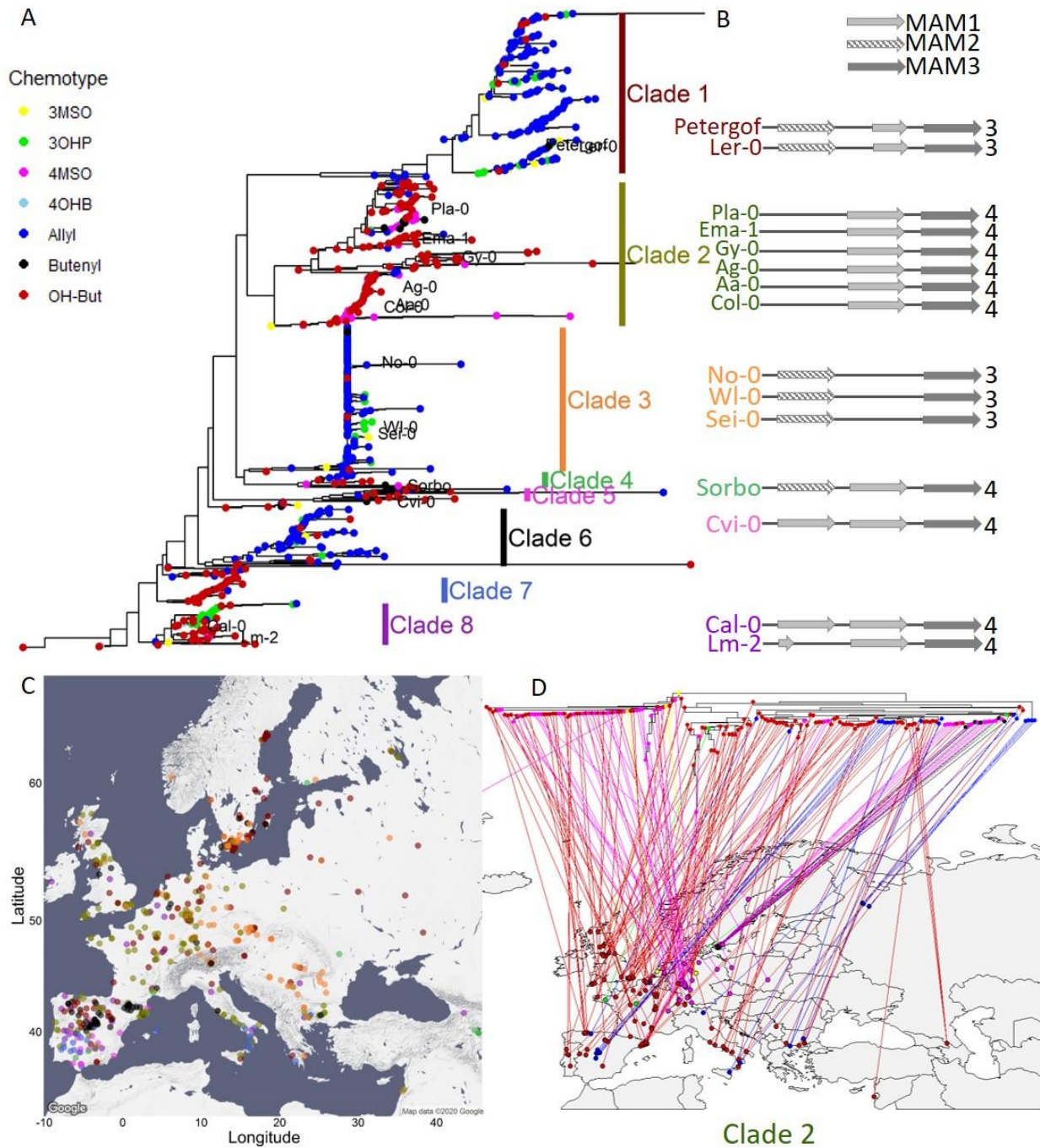


307 chemotypes, and use the available genomic sequences in all of these accessions to investigate the  
308 allelic variation in these genes to map the allelic distribution and test the potential for convergent  
309 and/or parallel evolution within each locus.

310 GS-Elong: Because the variation in the GS-Elong locus is caused by complex structural variation  
311 in MAM1 and MAM2 that is not resolvable using the available data from short-read genomic  
312 sequence, we used the MAM3 sequence within this locus to ascertain the genomic relationship of  
313 accessions at the causal GS-Elong locus (Kroymann et al., 2003). We aligned the MAM3  
314 sequence from each of the accessions, rooted the tree with the *Arabidopsis lyrata* orthologue  
315 (MAMc), and colored the tree tips based on the accessions dominant chemotype.

316 The accessions were distributed across eight distinctive clades with each clade clustering  
317 accessions having either a C3 or C4 status (Figure 5A). The clades C3/C4 status altered across  
318 the tree with three of the clades C3 dominant (MAM2 expressed), and five clades being C4  
319 dominant (MAM1 expressed). Further supporting the use of MAM3 is that the accession  
320 assignments to these clades agree with available bacterial artificial chromosome-based  
321 sequencing of the GS-Elong region from 15 accessions (Figure 5B). While there are multiple  
322 functional alleles for both C3 and C4 chemotypes, the genomic sequence and phylogeny does not  
323 appear consistent with a simple parallel evolution model where one allele/population is the basis  
324 for the independent derivation of all alternative alleles. This is illustrated by the difference in the  
325 genomic arrangement of Clade 4 and 5 which both create C4 GSLs. Clade 4 has a copy of  
326 MAM2 and MAM1 while Clade 5 has two copies of MAM1 (Figure 5). It appears that Clade 4 is  
327 the basis for two independent C3 alleles via separate deletions of MAM1 (Clades 1 and 3) and a  
328 separate C4 allele via a deletion of MAM2 (Clade 2, Figure 5). Unfortunately, no long-read  
329 sequencing is available in accessions from Clade 6 or 7 and locus-specific de novo alignment of  
330 short-read sequences in these accessions was not able to resolve the regions complexity. Filling  
331 in these clades would be necessary to better understand convergent/parallel events giving rise to  
332 GSL chemotypes.

333 Interestingly, the most basal clade has no copy of MAM2, raising the question of where MAM2  
334 arose (Figure 5). This suggests that true ancestral state(s) of this locus is not represented in this  
335 collection and would need to be searched for in other populations of *Arabidopsis thaliana*, if it  
336 exists in extant populations.



337

338 **Figure 5: MAM3 phylogeny.** A. MAM3 phylogeny of *Arabidopsis thaliana* accessions, rooted  
 339 by *Arabidopsis lyrata* MAMc, that is not shown because of distance. Tree tips are colored based  
 340 on the accession chemotype. The named accessions indicate that GSL-Elong region of these  
 341 accessions was previously sequenced (Kroymann et. al. 2003). B. The genomic structure of the  
 342 GSL-Elong regions in the previously sequenced accessions is shown based on Kroymann et. al.  
 343 2003. The accession names are colored based on their clades. The color of the name of the  
 344 accession indicates the clade it belongs. Bright grey arrows represents MAM1 sequences, dashed



345 arrows represents MAM2 sequences. Dark grey arrows represent MAM3 sequences. The number  
346 to the right of the genomic cartoon represents the number of carbons in the side chain. C.  
347 Collection sites of the accessions, colored by their clade classification (from section A). D. Clade  
348 2 reflection on the map.

349

350 Using this phylogeny, we investigated the presence of the different GS-Elong haplotypes across  
351 Europe to ask if each region has a specific allele/clade or if the alleles are distributed across the  
352 continent. Specifically, we were interested if the strong C3/C4 partitioning in southern Europe  
353 was driven by the creation of local alleles or if this partitioning might contain a wide range of  
354 alleles. If the latter is true, this would argue for a selective pressure shaping this C3/C4 divide.  
355 We plotted the accessions on the map and colored them based on their GS-Elong clade (Figure  
356 5C). This showed that the strong C3/C4 partition in the Iberian Peninsula contains haplotypes  
357 from all the GS-Elong clades except Clade 3 and is not caused by local alleles. This suggests that  
358 the strong geographic partitioning of the C3/C4 chemotypes in Iberia may be driven by selective  
359 pressure causing the partitioning of the chemotypes rather than neutral demographic processes.

360 Shifting focus to all of Europe showed that while most clades were widely distributed across  
361 Europe there were a couple over-arching patterns (Figure 5C, and supp. Figure 9). GS-Elong  
362 clades 1 and 6 follow a pattern that fits with alleles located within the Iberian glacial refugia that  
363 then moved north. In contrast the absence of clade 3 from Iberia is more parsimonious with a  
364 glacial refugia in the Balkans followed by a northward movement wherein it mixed with the  
365 other clades. Other clades never moved north and are exclusive to the south as shown by clades 5  
366 and 7. While these are both C4 clades, other C4 clades like clades 2 and 8 were able to move  
367 north (Figure 5D, and supp. Figure 9, respectively). This suggests that there are either differences  
368 in their GSL chemotype influencing their distribution or there are neighboring genes known to be  
369 under selection in Arabidopsis like FLC (AT5G10140) that may have influenced their  
370 distribution. In combination, this suggests that a complex demography is involved in shaping the  
371 chemotypes identity with some regions, Iberia, showing evidence of local selection while other  
372 regions, central Europe, possibly showing a blend requiring further work to delineate (supp.  
373 Figure 9).

374 GS-AOP: Side chain modification of the core MSO GSL is determined by the GS-AOP locus.  
375 Most of the accessions contain a copy of AOP2 and a copy of AOP3, but only one of them will

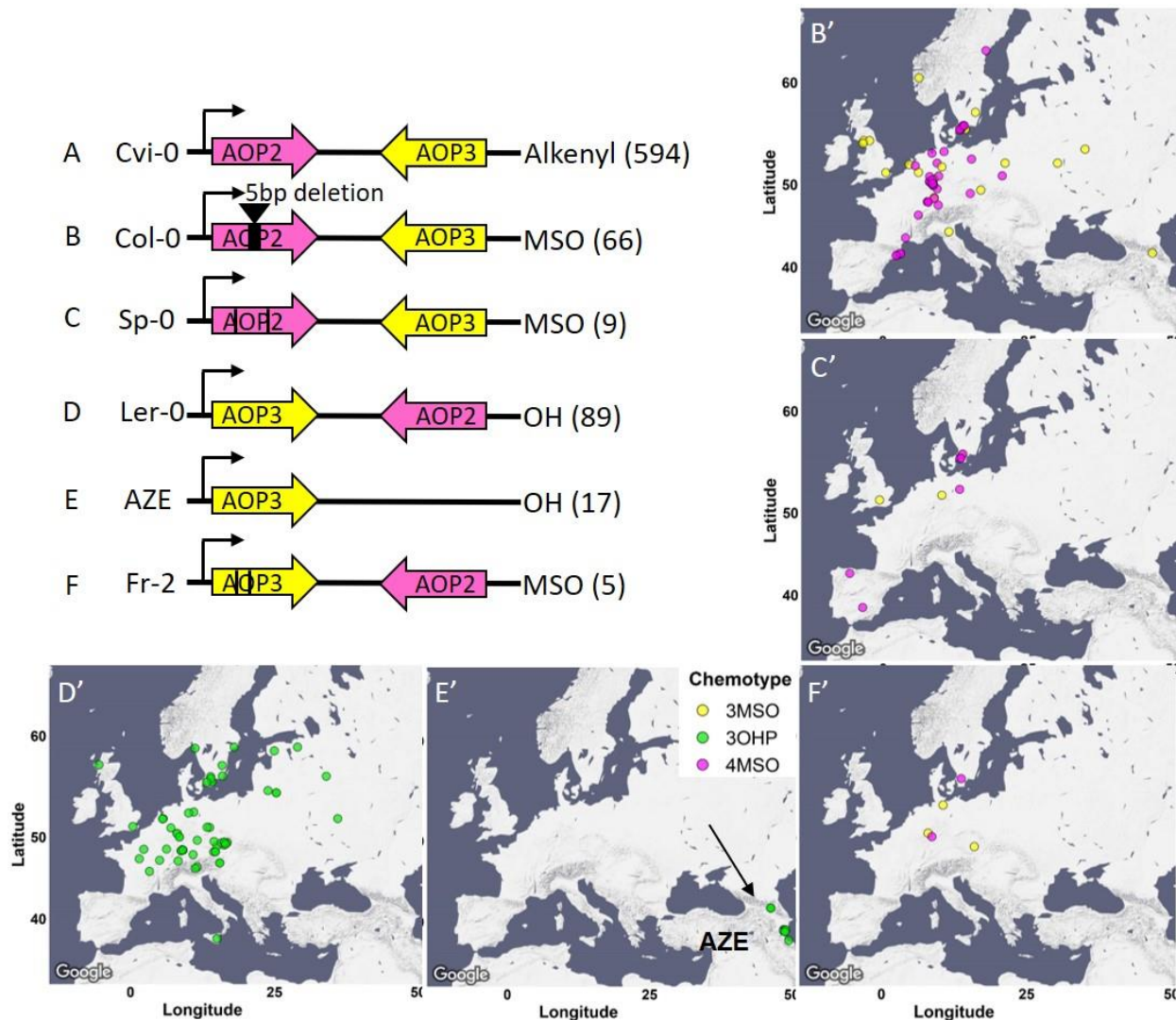
376 be functionally expressed (Chan et al., 2010), while in some cases both will be nonfunctional. To  
377 better understand the demography and evolution of the GS-AOP locus, we separately aligned the  
378 AOP2 and AOP3 sequences, rooted each tree with the *Arabidopsis lyrata* orthologue, and  
379 colored the trees tips based on the accessions dominant chemotype.

380 The phylogenetic trees shared a very similar topology, yielding a clear separation between  
381 alkenyl (AOP2 expressed) and hydroxyalkyl (AOP3 expressed) accessions. Alkenyl expressing  
382 accessions like Cvi-0 with an expressed copy of the AOP2 enzyme formed a single contiguous  
383 cluster (Figure 6A). In contrast, hydroxyalkyl accessions clustered into two separate groups with  
384 one group of 3OHP dominant accessions partitioning from the rest of the accessions at the most  
385 basal split in the tree (supp. Figure 10). This haplotype is marked by having an inversion  
386 swapping the AOP2 and AOP3 promoters as shown in bacterial artificial chromosome  
387 sequencing of the Ler-0 accession (Figure 6D) (Chan et al., 2010). The tree also identified a  
388 second group of 3OHP dominant accessions located among the alkenyl accessions. Analyzing  
389 the sequences of these accessions reveals that this small group of 3OHP accessions have a  
390 complete deletion of AOP2 and contain only AOP3 (Figure 6E). Thus, there are at least two  
391 independent transitions from Alkenyl to Hydroxyalkyl GSLs within *Arabidopsis*, neither of  
392 which are related to the Alkenyl to Hydroxyalkyl conversion within *Arabidopsis lyrata*.

393 The null accessions (MSO dominant chemotypes) were identifiable in all the major clades on the  
394 tree (supp. Figure 10, middle column of heatmap) suggesting that there are independent LOF  
395 mutations that abolish either AOP2 or AOP3. Deeper examination of the sequences of these  
396 accessions identified three convergent LOF alleles leading to the MSO chemotype. Most of the  
397 null accessions harbor a 5 bps deletion in their AOP2 sequence, that causes a frameshift  
398 mutation. This mutation arose within the Alkenyl haplotype and was first reported in the Col-0  
399 reference genome (Figure 6B) (D J Kliebenstein, Lambrix, et al., 2001). In addition, there are  
400 additional independent LOF events arising in both the alkenyl haplotype (e.g. Sp-0, Figure 6C),  
401 and within the Ler-0 inversion haplotype (e.g. Fr-2, Figure 6F). Thus, GS-AOP has repeated  
402 LOF alleles arising within all the major AOP haplotypes suggesting convergent evolution of the  
403 MSO chemotype out of both the Alkenyl and Hydroxyalkyl chemotypes.

404 Using the combined chemotype/genotype assignments at GS-AOP, we investigated the  
405 distribution of the alleles across Europe. The Alkenyl haplotype is spread across the entire

406 continent. In contrast, the hydroxyalkyl haplotypes are more local. The Ler-like 3OHP haplotype  
 407 is present in only central and north Europe (Figure 6D), while the other 3OHP haplotype,  
 408 possessing only AOP3, is limited to Azerbaijan, along the Caspian Sea (Figure 6E). In contrast to  
 409 the distinct hydroxyalkyl locations, the distribution of the independent LOF null haplotypes  
 410 overlaps with all of them being located within central and north Europe (Figure 6B, C and F).  
 411 The fact that these independently derived LOF alleles are all contiguous suggests that there may  
 412 be a benefit to these alleles specific to Central Europe.



413

414 **Figure 6: AOP Genomic structure.** The genomic structure and causality of the major  
 415 AOP2/AOP3 haplotypes are illustrated. Pink arrows show the AOP2 gene while yellow arrows  
 416 represent AOP3. The black arrows represent the direction of transcription from the AOP2  
 417 promoter as defined in the Col-0 reference genome. Its position does not change in any of the  
 418 regions. The black lines in Sp-0 and Fr-2 presence the position of independent variants creating





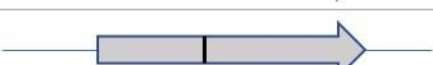




419 premature stop codons. The GSL chemotype for each haplotype is listed to the right with the  
420 number of the accessions in brackets. The maps show the geographic distribution of the  
421 accessions from each structure.

422

423 GS-OH: The final major determinant of natural variation in Arabidopsis GSL chemotype is the  
424 GS-OH enzyme that adds a hydroxyl group to the 2 carbon on 3-butenyl GSL to create 2-OH-3-  
425 butenyl GSL. Previous work had suggested two GS-OH alleles measurable in the seed, a  
426 functional allele in almost all accessions and a non-functional allele caused by active site  
427 mutations represented by the Cvi-0 accession (Hansen et al., 2008). Because of functional  
428 epistasis, we can only obtain functional phenotypic information from accessions that accumulate  
429 the GS-OH substrate, 3-butenyl GLS. This identified 11 accessions with a non-functional GS-  
430 OH. Surveying these 11 accessions in the polymorph database identified multiple independent  
431 LOF events. One of these 11 accessions have the Cvi active site mutations, two accessions have  
432 a shared nonsense SNP that introduce premature stop codons, and two accessions have a  
433 complete loss of this gene (Table 1). We could not identify the causal LOF allele in the other six  
434 accessions due to sequence quality in the databases. All of these independent GS-OH LOF  
435 alleles are phylogenetically positioned within groups of accessions that largely do not  
436 accumulate 3-butenyl GLS, e.g. 3 carbon or non-alkenyl accessions suggesting that the  
437 functional epistasis may be influencing the generation of these alleles. Thus, we searched the  
438 accessions that do not accumulate 3-butenyl GLS and have effectively hidden the GS-OH  
439 function for these GS-OH LOF events (Supp table 2). In each case, the LOF allele is more  
440 frequent in the non 4 carbon-alkenyl accessions than expected by random chance. This suggests  
441 that there is a bias against 3-butenyl GSL synthesis as the LOF alleles are more frequent when  
442 the GS-OH gene is hidden by functional epistasis. This agrees with the fact that the 3-butenyl  
443 chemotype is the most sensitive to generalist lepidopteran herbivory (Hansen et al., 2008). Thus,  
444 these mutations may represent ongoing pseudogenization of the GS-OH gene when it is  
445 functionally hidden by epistasis at the GS-AOP and GS-Elong loci. These LOF events would  
446 then only be displayed upon rare admixture with 2-OH-3-Butenyl accessions.

447

448

Accession	Type of mutation	Allele structure
Sorbo, Pien	polymorphism at SNP10831302	
Cvi-0	Active site mutation	
IP-Mot-0, IP-Tri-0	Gene deletion	
T670	Independent mutation	
FlyA 3	Independent mutation	
Ting-1	Independent mutation	
T880	Independent mutation	
T710	Independent mutation	
T850	Independent mutation	

449

450 Table 1: **GS-OH structure**. The structures of GS-OH in the 3-Butenyl accessions are illustrated.  
 451 These mutations create premature stop codons.

452

## 453 Discussion

454 Understanding the genetic, demographic and environmental factors that shape variation within a  
 455 trait in a population is key to understanding trait evolution. In this work we used Aliphatic GSLs  
 456 in seeds of *Arabidopsis thaliana* to query how genetics, geography, environment and  
 457 demography intersect to shape chemotypic variation across Europe. We found that  
 458 environmental conditions, together with geography affect the presence and distribution of  
 459 chemotypes within the accessions. This was demonstrated by specific traits that were associated  
 460 with specific environmental conditions, and this association was shifted across the continent.  
 461 Comparing the associations of traits to specific environmental conditions in central Europe  
 462 versus the south revealed different, sometimes even inverse, behaviors. For example, In the  
 463 Iberian Peninsula, 2-OH-3-Butenyl was positively associated with potential drought while in  
 464 Central Europe, it was the opposite GS-Elong allele showing association. This showed that



465 chemotypic variation across Europe is created by a blend of all these processes that differ at the  
466 individual loci and required the simultaneous analysis of genotype and phenotype to fully  
467 interpret.

468 In contrast to the bimodal distribution in the Aliphatic GSL traits, each of the three major  
469 Aliphatic GSL loci showed allelic heterogeneity with multiple independent structural variants  
470 that recreated the same phenotypic variation. The GS-AOP locus had numerous events with the  
471 AOP3 variant of GS-AOP arising via at least two independent events and the Null allele being  
472 generated at least 15 independent times. The GS-AOP null alleles convergently arose from all  
473 the different functional haplotypes. It is less clear if the independent GS-AOP AOP3 alleles  
474 should be classified as convergent or parallel due to a lack of clarity in what is the ancestral state.  
475 Similar to the independent GS-AOP null alleles, there were numerous independent GS-OH LOF  
476 variants with at least 9 independent events. While these are parallel GS-OH LOF events because  
477 they came from a single functional GS-OH group, their ability to accumulate depends on the  
478 epistatic silencing of GS-OH by the GS-AOP and GS-Elong loci. The GS-Elong locus also had  
479 an extensive level of allelic heterogeneity hallmarked by a shifting expression of the MAM1 or  
480 MAM2 gene, again with hallmarks of both parallel and convergent processes. Interestingly, at  
481 both the GS-AOP and GS-Elong loci, one gain-of-function event (e.g. AOP3 in GS-AOP locus,  
482 and MAM2 in GS-Elong locus) is concurrently linked to a loss-of-function of the other gene at  
483 the locus (AOP2 and MAM1, respectively). These structural variants are shaped such that the  
484 chemotypes show distinct separations without any intermediate phenotypes. This allelic  
485 heterogeneity is in contrast to previous work on other biotic interactions genes like pathogen  
486 resistance gene-for-gene loci that typically have two moderate frequency stable alleles creating  
487 the phenotypic variation within the species (Atwell et al., 2010; Corrión & Day, 2001;  
488 MacQueen, Sun, & Bergelson, 2016). In other cases alleles of genes involved in biotic defense  
489 can present more complex patterns, e.g. natural variation in the immune gene *ACCELERATED*  
490 *CELL DEATH 6* (*ACD6*) is caused by a rare allele causing an extreme lesion phenotype. It is not  
491 yet clear what selective pressures influence *ACD6* genetic variation (Todesco et al., 2010; Zhu et  
492 al., 2018). The contrast where Aliphatic GSL loci have high levels of allelic heterogeneity for  
493 independent and recurrent LOF and GOF events while other resistance genes have more stable  
494 biallelic variation suggests that there are different selective regimes influencing these loci.

495 Further work is needed to assess the range of allelic heterogeneity in loci controlling resistance  
496 to diverse biotic traits within the environment.

497

498 The allelic heterogeneity at these loci illustrates the benefit of simultaneously tracking the  
499 phenotype and genotype when working to understand the distribution of trait variation. For  
500 example, the Iberian Peninsula and the Mediterranean had low variability in Aliphatic GSL  
501 chemotype with the chemotypes not overlapping while central/north Europe had high Aliphatic  
502 GSL diversity with the chemotypes overlapping. At first glance, this contrasts with previous  
503 work showing that the Iberian Peninsula and the Mediterranean are more genetically diverse.  
504 However, this discrepancy was caused by one of the causal loci. Specifically, the GS-AOP locus  
505 is largely fixed as the Alkenyl allele in Iberia/Mediterranean with the alternative GS-AOP alleles  
506 enriched in central Europe. In contrast to GS-AOP, Iberia and the Mediterranean were highly  
507 genetically diverse for the GS-Elong locus and appear to contain all the variation in GS-Elong  
508 found throughout Europe (The 1001 Genomes Consortium, 2016). Thus, the chemotypic  
509 divergence from genomic variation expectations was driven by just the GS-AOP locus. This  
510 indicates that the high level of chemotypic variation in central Europe is a blend of alleles that  
511 moved from the south (GS-Elong) and alleles that possibly arose locally (GS-AOP, both nulls  
512 and AOP3). Further, the chemotypes found in any one region appear to be created by a  
513 combination of alleles moving across the continent, local generation of new polymorphisms and  
514 local selective pressures that shape the chemotypes distribution across the landscape.

515

516 One difficulty in interpreting the evolutionary processes, e.g. parallel v convergent, especially for  
517 structural variants illustrated by all the three loci is the complication in properly identifying the  
518 ancestral state of the population. While this could typically be done by relying on shared loci  
519 with sister species, this is not possible in this case as *Arabidopsis lyrata* and *halleri* have genetic  
520 variation at GS-Elong and GS-AOP creating the exact same phenotypes. Further, neither of these  
521 sister species have yet been found to have a functional GS-OH (Heidel, Clauss, Kroymann,  
522 Savolainen, & Mitchell-Olds, 2006; Ramos-Onsins, Stranger, Mitchell-Olds, & Aguadé, 2004;  
523 Windsor et al., 2005). The MSO chemotypes could be viewed as convergent evolution within a  
524 species, as the MSO phenotype independently re-occurred multiple times in the AOP2 and AOP3  
525 genetic backgrounds. However, it is not clear how to classify the different AOP3 (AZE v Ler)



526 types as it not clear if the AOP2 or AOP3/Ler haplotype is ancestral within *Arabidopsis thaliana*.  
527 Another option to calling ancestral state is deep sampling in the species but even with these  
528 accessions, we do not appear to have reached the necessary threshold. For example previous  
529 work at the GS-Elong locus had suggested that the Sorbo accession, collected from Tajikistan,  
530 was the most likely ancestral state as it had a copy of both MAM1 and MAM2 (Kroymann et al.,  
531 2003). However, the phylogeny with this larger collection of accessions suggested that Sorbo is  
532 not ancestral. Further, a recent phylogeny of MAM genes across the Brassicales suggests that  
533 MAM2 is an *Arabidopsis thaliana* specific gene with an undefined origin (Abrahams et al.,  
534 2020). This suggests that to get a better understanding of the ancestral state to define  
535 evolutionary processes, especially for loci with allelic heterogeneity and structural variants, we  
536 need to broaden our phylogenetic context by deeper sampling within and between species.  
537

538 Another complication caused by the allelic heterogeneity and differential selective pressures  
539 displayed within this system is that we were unable to detect a number of known and validated  
540 natural variants that are causal within this population. Specifically, the GWAS with this  
541 collection of 797 accessions was unable to find 80% of the known causal loci including one of  
542 the three major effect loci, GS-OH. Maximizing the number of genotypes and the SNP marker  
543 density was unable to overcome the complications imposed by the complex pressures shaping  
544 the distribution of these traits. In this system, the optimal path to identifying the causal  
545 polymorphisms has instead been a small number of Recombinant Inbred Line populations  
546 derived from randomly chosen parents. In complex adaptive systems, the optimal solution to  
547 identifying causal variants is likely a blend of structured mapping populations and then  
548 translating the causal genes from this system to the GWAS results and tracking the causal loci  
549 directly.

550 In this work we combined different approaches to uncover some of the parameters shaping the  
551 Aliphatic GSL content across Europe. Widening the size of the population will enable us to  
552 deepen our understanding on the evolutionary mechanisms shaping a phenotype in a population.  
553  
554  
555

## 556 **Methods**

### 557 **Plant materiel:**

558 Seeds for 1135 *Arabidopsis* (*Arabidopsis thaliana*) genotypes were obtained from the 1001  
559 genomes catalog of *Arabidopsis thaliana* genetic variation (<https://1001genomes.org/>). All  
560 *Arabidopsis* genotypes were grown at 22°C/24°C (day/night) under long-day conditions (16-h of  
561 light/8-h of dark). Two independent experiments were performed, each of them included the full  
562 set of genotypes. In the analyses only accessions from Europe and around Europe were included  
563 (Figure 2A), resulting in an analysis of 797 accessions. List of the accessions can be found in  
564 supp. Table 1.

### 565 **GSL extractions and analyses:**

566 GSLs were measured as previously described (D J Kliebenstein, Gershenzon, et al., 2001; D J  
567 Kliebenstein, Kroymann, et al., 2001; D J Kliebenstein, Lambrix, et al., 2001). Briefly, ~3mg of  
568 seeds were harvested in 200 µL of 90% methanol. Samples were homogenized for 3 min in a  
569 paint shaker, centrifuged, and the supernatants were transferred to a 96-well filter plate with  
570 DEAE sephadex. The filter plate with DEAE sephadex was washed with water, 90% methanol,  
571 and water again. The sephadex bound GSLs were eluted after an overnight incubation with  
572 110µL of sulfatase. Individual desulfo-GSLs within each sample were separated and detected by  
573 HPLC-DAD, identified, quantified by comparison to standard curves from purified compounds,  
574 and further normalized to the weight. List of GSLs and their structure are in supplementary table  
575 1. Row GSLs data are in supplementary table 1B.

### 576 **Statistics, heritability, and data visualization:**

577 Statistical analyses were conducted using R software (<https://www.R-project.org/>) with the  
578 RStudio interface (<http://www.rstudio.com/>). For each independent GLS, a linear model  
579 followed by ANOVA was utilized to analyze the effect of accession, replicate, and location in  
580 the experiment plate upon the measured GLS amount. Broad-sense heritability (supplementary  
581 table 1C) for the different metabolites was estimated from this model by taking the variance due  
582 to accession and dividing it by the total variance. Estimated marginal means (emmeans) for each  
583 accession were calculated for each metabolite from the same model using the package emmeans  
584 (“CRAN - Package emmeans,” n.d.) (supplementary table 1D). Principal component analyses

585 were done with FactoMineR and factoextra packages (Abdi & Williams, 2010). Data analyses  
586 and visualization was done using R software with tidyverse (Wickham et al., 2019) and ggplot2  
587 (Kahle & Wickham, 2013) packages.

588 Principal component analyses were done with FactoMineR and factoextra packages (Abdi &  
589 Williams, 2010).

590 Maps were generated using ggmap package (“[https://journal.r-project.org/archive/2013-1/kahle-  
591 wickham.pdf](https://journal.r-project.org/archive/2013-1/kahle-wickham.pdf),” n.d.).

### 592 **Phenotypic classification based on GSL content:**

593 For each accession the expressed enzyme in each of the following families was determined based  
594 on the content (presence and amounts) of short chained Aliphatic GSLs:

595 MAM enzymes: the total amount of 3 carbons GSLs and 4 carbons GSLs was calculated for each  
596 accession. 3 carbons GSLs include 3MT, 3MSO, 3OHP and Allyl GSL. 4 carbons GSLs include  
597 4MT, 4MSO, 4OHB, 3-butenyl and 2-OH-3-butenyl GSL (for structures and details see supp.  
598 Table 1). Accessions that the majority of Aliphatic short chained GSL contained 3 carbons in  
599 their side chains classified as MAM2 expressed (supp. Figure 3). Accessions that the majority of  
600 Aliphatic short chained GSL contained 4 carbons in their side chains classified as MAM1  
601 expressed (supp. Figure 3). The accessions were plotted on a map based on their original  
602 collection sites (supp. Figure 3).

603 AOP enzymes: the relative amount of alkenyl GSL, alkyl GSL and MSO GSL were calculated in  
604 respect to the total short chained Aliphatic GSL as follows:

$$605 \text{ Alkenyl GSL (AOP2 expressed)} = \frac{\text{Allyl} + 2\text{-OH-3-butenyl} + 3\text{-butenyl}}{\text{Total short chained GSL}}$$

$$606 \text{ Alkyl GSL (AOP3 expressed)} = \frac{3\text{OHP} + 4\text{OHB}}{\text{Total short chained GSL}}$$

$$607 \text{ MSO GSL (AOP null)} = \frac{3\text{MSO} + 4\text{MSO}}{\text{Total short chained GSL}}$$

608 The expressed AOP enzyme was determined based on those ratios: accessions with majority  
609 alkenyl GSL were classified as AOP2 expressed. Accessions with majority of alkyl GSL were  
610 classified as AOP3 expressed. Accessions with majority of MSO GSL were classified as AOP

611 null. The accessions were plotted on a map based on their original collection sites (supp. Figure  
612 4).

613 GS-OH enzyme: the ratio between 2-OH-3-butenyl GSL to 3-butenyl GSL was calculated only  
614 for MAM1 expressed accessions (accessions that the majority of GSLs contain 4 carbons in their  
615 side chain). Accessions with high amounts of 2-OH-3-butenyl GSL were classified as GS-OH  
616 functional. Accessions with high amounts of 3-butenyl GSL were classified as GS-OH non-  
617 functional. The accessions were plotted on a map based on their original collection sites (supp.  
618 Figure 5).

619 Each accession was classified to one of seven Aliphatic short chained GSLs based on the  
620 combination of the dominancy of the enzymes as follows: MAM2, AOP null: classified as  
621 3MSO dominant. MAM1, AOP null: classified as 4MSO dominant. MAM2, AOP3: classified as  
622 3OHP dominant. MAM1, AOP3: classified as 4OHB dominant. MAM2, AOP2: classified as  
623 Allyl dominant. MAM1, AOP2, GS-OH non-functional: classified as 3-Butenyl dominant.  
624 MAM1, AOP2, GS-OH functional: classified as 2-OH-3-Butenyl dominant. The accessions were  
625 plotted on a map based on their original collection sites and colored based on their dominant  
626 chemotype (Figure 3).

#### 627 **Environmental data:**

628 Environmental data was obtained from the 1001 genomes website (<https://1001genomes.org/>, for  
629 geographical data) and from the Arabidopsis CLIMtools  
630 (<http://www.personal.psu.edu/sma3/CLIMtools.html>, (Ferrero-Serrano & Assmann, 2019)) for  
631 environmental data. We used the five variables that captured a majority of the variance in this  
632 dataset including maximal temperature of warmest month (WC2\_BIO5), minimal temperature of  
633 coldest month (WC2\_BIO6), precipitation of wettest month (WC2\_BIO13), precipitation of  
634 driest month (WC2\_BIO14), and distance to the coast (in Km).

#### 635 **Environmental MANOVA:**

636 Linear models to test the effect of geographical and environmental parameters (supp. Figure 1, 8)  
637 were conducted using dplyr package (“CRAN - Package dplyr,” n.d.) and included the following  
638 parameters:

639 Supp. Figure 1- linear models for collection sites: PC score ~ Latitude + Longitude + Latitude \*  
640 Longitude.

641 Supp. Figure 7 - for all the data: C length (C3 or C4) ~ Genomic group + Geography (north  
642 versus south) +Max temperature of warmest month+ Min temperature of coldest month+  
643 Precipitation of wettest month+ Precipitation of driest month+ Distance to the coast + Geography  
644 \*Genomic group + Geography \* Max temperature of warmest month + Geography \* Min  
645 temperature of coldest month+ Geography \* Precipitation of driest month+ Geography \*  
646 Precipitation of wettest month + Geography \*Distance to the coast.

647 For the north and the south: C length (C3 or C4) ~ Genomic group + Geography (north versus  
648 south) +Max temperature of warmest month+ Min temperature of coldest month+ Precipitation  
649 of wettest month+ Precipitation of driest month+ Distance to the coast.

650 Multivariate analysis of variance (MANOVA) models to check the effect of environmental  
651 variables on the chemotype identity for each collection (Figure 4 and supp. Figure 7) included  
652 the following parameters:

653 Max temperature of warmest month+ Min temperature of coldest month+ Precipitation of wettest  
654 month+ Precipitation of driest month+ Distance to the coast~ Chemotype.

655 **Random Forest analyses:** Random forest analyses was conducted using the “randomForest”  
656 and “ElemStatLearn” packages in Rstudio (“CRAN - Package ElemStatLearn,” n.d., “CRAN: R  
657 News,” n.d.; Liaw & Wiener, 2002). In these analyzes we used the environmental parameters  
658 and genomic group data to predict the chemotype identity, after excluding the low frequencies  
659 chemotypes (4OHB and 3-Butenyl from all of them, 3MSO from the south).

660 **Genome wide association studies:**

661 The phenotypes for GWAS were each accession value for PC1 and 2. GWAS was implemented  
662 with the easyGWAS tool (Grimm et al., 2017) using the EMMAX algorithms (Kang et al.,  
663 2010)and a minor allele frequency (MAF) cutoff of 5%. The results were visualized as  
664 manhattan plots using the qqman package in R (Turner, 2014).

665 **Phylogeny:**

666 Genomic sequences from the accessions for MAM3 – AT5G23020, AOP2 – Chr4, 1351568 until  
667 1354216, AOP3 - AT4G03050.2, and GS-OH – AT2G25450 were obtained using the  
668 Pseudogenomes tool ([https://tools.1001genomes.org/pseudogenomes/#select\\_strains](https://tools.1001genomes.org/pseudogenomes/#select_strains)).

669 Multiple sequence alignment was done with the msa package in R, using the ClustalW,  
670 ClustalOmega, and Muscle algorithms (Bodenhofer, Bonatesta, Horejš-Kainrath, & Hochreiter,  
671 2015). Phylogenetic trees were generated with the ggtree package in R (Yu, 2020). Each tree was  
672 rooted by the genes matching *Arabidopsis layrata*'s functional orthologue or closest homologue.

673

## 674 **Acknowledgments**

675 This work was supported by the National Science Foundation, Directorate for Biological  
676 Sciences, Division of Molecular and Cellular Biosciences (grant no. MCB 1906486 to D.J.K)  
677 and Division of Integrative Organismal Systems (grant no. IOS 1655810 to D.J.K, and IOS  
678 1754201 to R.A), and by the United States-Israel Binational Agricultural Research and  
679 Development Fund (to D.J.K. and E.K., grant no. FI-560-2017). We thank Dr. Allison  
680 Gaudinier (Department of Plant and Microbial Biology, University of California, Berkeley), Dr.  
681 Tobias Züst (Institute of Plant Sciences, University of Bern), Dr. Christopher W. Wheat  
682 (Department of Zoology, Stockholm University) and Dr. Daniel Runcie (Department of Plant  
683 Sciences, University of California Davis.) for critical reading of the manuscript.

684

685

686

687

688

689

690

691

## 692 Bibliography

- 693 Abdi, H., & Williams, L. J. (2010). Principal component analysis. *Wiley Interdisciplinary*  
694 *Reviews: Computational Statistics*, 2(4), 433–459. <https://doi.org/10.1002/wics.101>
- 695 Abrahams, R. S., Pires, J. C., & Schranz, M. E. (2020). Genomic origin and diversification of the  
696 glucosinolate MAM locus. *Frontiers in Plant Science*, 11, 711.  
697 <https://doi.org/10.3389/fpls.2020.00711>
- 698 Agrawal, A. A. (2000). Overcompensation of plants in response to herbivory and the by-product  
699 benefits of mutualism. *Trends in Plant Science*, 5(7), 309–313.  
700 [https://doi.org/10.1016/S1360-1385\(00\)01679-4](https://doi.org/10.1016/S1360-1385(00)01679-4)
- 701 Atwell, S., Huang, Y. S., Vilhjálmsson, B. J., Willems, G., Horton, M., Li, Y., ... Nordborg, M.  
702 (2010). Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana* inbred  
703 lines. *Nature*, 465(7298), 627–631. <https://doi.org/10.1038/nature08800>
- 704 Bakker, E. G., Traw, M. B., Toomajian, C., Kreitman, M., & Bergelson, J. (2008). Low levels of  
705 polymorphism in genes that control the activation of defense response in *Arabidopsis*  
706 *thaliana*. *Genetics*, 178(4), 2031–2043. <https://doi.org/10.1534/genetics.107.083279>
- 707 Bednarek, P., & Osbourn, A. (2009). Plant-microbe interactions: chemical diversity in plant  
708 defense. *Science*, 324(5928), 746–748. <https://doi.org/10.1126/science.1171661>
- 709 Beekwilder, J., van Leeuwen, W., van Dam, N. M., Bertossi, M., Grandi, V., Mizzi, L., ... Bovy,  
710 A. (2008). The impact of the absence of aliphatic glucosinolates on insect herbivory in  
711 *Arabidopsis*. *Plos One*, 3(4), e2068. <https://doi.org/10.1371/journal.pone.0002068>
- 712 Benderoth, M., Textor, S., Windsor, A. J., Mitchell-Olds, T., Gershenzon, J., & Kroymann, J.  
713 (2006). Positive selection driving diversification in plant secondary metabolism.  
714 *Proceedings of the National Academy of Sciences of the United States of America*,  
715 103(24), 9118–9123. <https://doi.org/10.1073/pnas.0601738103>
- 716 Bialy, Z., Oleszek, W., Lewis, J., & Fenwick, G. R. (1990). Allelopathic potential of  
717 glucosinolates (mustard oil glycosides) and their degradation products against wheat.  
718 *Plant and Soil*, 129(2), 277–281. <https://doi.org/10.1007/BF00032423>
- 719 Bodenhofer, U., Bonatesta, E., Horejš-Kainrath, C., & Hochreiter, S. (2015). msa: an R package  
720 for multiple sequence alignment. *Bioinformatics*, 31(24), 3997–3999.  
721 <https://doi.org/10.1093/bioinformatics/btv494>
- 722 Brachi, B., Meyer, C. G., Villoutreix, R., Platt, A., Morton, T. C., Roux, F., & Bergelson, J.  
723 (2015). Coselected genes determine adaptive variation in herbivore resistance throughout  
724 the native range of *Arabidopsis thaliana*. *Proceedings of the National Academy of*  
725 *Sciences of the United States of America*, 112(13), 4032–4037.  
726 <https://doi.org/10.1073/pnas.1421416112>



- 727 Brown, P. D., Tokuhisa, J. G., Reichelt, M., & Gershenzon, J. (2003). Variation of glucosinolate  
728 accumulation among different organs and developmental stages of *Arabidopsis thaliana*.  
729 *Phytochemistry*, 62(3), 471–481. [https://doi.org/10.1016/S0031-9422\(02\)00549-6](https://doi.org/10.1016/S0031-9422(02)00549-6)
- 730 Chan, E. K. F., Rowe, H. C., Corwin, J. A., Joseph, B., & Kliebenstein, D. J. (2011). Combining  
731 genome-wide association mapping and transcriptional networks to identify novel genes  
732 controlling glucosinolates in *Arabidopsis thaliana*. *PLoS Biology*, 9(8), e1001125.  
733 <https://doi.org/10.1371/journal.pbio.1001125>
- 734 Chan, E. K. F., Rowe, H. C., & Kliebenstein, D. J. (2010). Understanding the evolution of  
735 defense metabolites in *Arabidopsis thaliana* using genome-wide association mapping.  
736 *Genetics*, 185(3), 991–1007. <https://doi.org/10.1534/genetics.109.108522>
- 737 Corrion, A., & Day, B. (2001). Pathogen Resistance Signalling in Plants. In John Wiley & Sons  
738 Ltd (Ed.), *eLS* (pp. 1–14). Chichester, UK: John Wiley & Sons, Ltd.  
739 <https://doi.org/10.1002/9780470015902.a0020119.pub2>
- 740 CRAN - Package dplyr. (n.d.). Retrieved June 16, 2020, from [https://CRAN.R-](https://CRAN.R-project.org/package=dplyr)  
741 [project.org/package=dplyr](https://CRAN.R-project.org/package=dplyr)
- 742 CRAN - Package ElemStatLearn. (n.d.). Retrieved June 16, 2020, from [https://CRAN.R-](https://CRAN.R-project.org/package=ElemStatLearn)  
743 [project.org/package=ElemStatLearn](https://CRAN.R-project.org/package=ElemStatLearn)
- 744 CRAN - Package emmeans. (n.d.). Retrieved June 16, 2020, from [https://CRAN.R-](https://CRAN.R-project.org/package=emmeans)  
745 [project.org/package=emmeans](https://CRAN.R-project.org/package=emmeans)
- 746 CRAN: R News. (n.d.). Retrieved June 16, 2020, from <https://CRAN.R-project.org/doc/Rnews/>
- 747 Daxenbichler, M. E., Spencer, G. F., Carlson, D. G., Rose, G. B., Brinker, A. M., & Powell, R.  
748 G. (1991). Glucosinolate composition of seeds from 297 species of wild plants.  
749 *Phytochemistry*, 30(8), 2623–2638. [https://doi.org/10.1016/0031-9422\(91\)85112-D](https://doi.org/10.1016/0031-9422(91)85112-D)
- 750 Erb, M., & Kliebenstein, D. J. (2020). Plant secondary metabolites as defenses, regulators, and  
751 primary metabolites: the blurred functional trichotomy. *Plant Physiology*.  
752 <https://doi.org/10.1104/pp.20.00433>
- 753 Fan, P., Leong, B. J., & Last, R. L. (2019). Tip of the trichome: evolution of acylsugar metabolic  
754 diversity in Solanaceae. *Current Opinion in Plant Biology*, 49, 8–16.  
755 <https://doi.org/10.1016/j.pbi.2019.03.005>
- 756 Ferrero-Serrano, Á., & Assmann, S. M. (2019). Phenotypic and genome-wide association with  
757 the local environment of *Arabidopsis*. *Nature Ecology & Evolution*, 3(2), 274–285.  
758 <https://doi.org/10.1038/s41559-018-0754-5>
- 759 Futuyma, D. J., & Agrawal, A. A. (2009). Macroevolution and the biological diversity of plants  
760 and herbivores. *Proceedings of the National Academy of Sciences of the United States of*  
761 *America*, 106(43), 18054–18061. <https://doi.org/10.1073/pnas.0904106106>

- 762 Giamoustaris, A., & Mithen, R. (1996). Genetics of aliphatic glucosinolates. IV. Side-chain  
763 modification in *Brassica oleracea*. *TAG. Theoretical and Applied Genetics. Theoretische*  
764 *Und Angewandte Genetik*, 93(5-6), 1006–1010. <https://doi.org/10.1007/BF00224105>
- 765 Grimm, D. G., Roqueiro, D., Salomé, P. A., Kleeberger, S., Greshake, B., Zhu, W., ...  
766 Borgwardt, K. M. (2017). easyGWAS: A Cloud-Based Platform for Comparing the  
767 Results of Genome-Wide Association Studies. *The Plant Cell*, 29(1), 5–19.  
768 <https://doi.org/10.1105/tpc.16.00551>
- 769 Halkier, B. A., & Gershenzon, J. (2006). Biology and biochemistry of glucosinolates. *Annual*  
770 *Review of Plant Biology*, 57, 303–333.  
771 <https://doi.org/10.1146/annurev.arplant.57.032905.105228>
- 772 Hanower, P., & Brzozowska, J. (1975). Influence d'un choc osmotique sur la composition des  
773 feuilles de cotonnier en acides amines libres. *Phytochemistry*, 14(8), 1691–1694.  
774 [https://doi.org/10.1016/0031-9422\(75\)85275-7](https://doi.org/10.1016/0031-9422(75)85275-7)
- 775 Hansen, B. G., Kerwin, R. E., Ober, J. A., Lambrix, V. M., Mitchell-Olds, T., Gershenzon, J., ...  
776 Kliebenstein, D. J. (2008). A novel 2-oxoacid-dependent dioxygenase involved in the  
777 formation of the goiterogenic 2-hydroxybut-3-enyl glucosinolate and generalist insect  
778 resistance in *Arabidopsis*. *Plant Physiology*, 148(4), 2096–2108.  
779 <https://doi.org/10.1104/pp.108.129981>
- 780 Hansen, B. G., Kliebenstein, D. J., & Halkier, B. A. (2007). Identification of a flavin-  
781 monooxygenase as the S-oxygenating enzyme in aliphatic glucosinolate biosynthesis in  
782 *Arabidopsis*. *The Plant Journal: For Cell and Molecular Biology*, 50(5), 902–910.  
783 <https://doi.org/10.1111/j.1365-313X.2007.03101.x>
- 784 Hasegawa, T., Yamada, K., Kosemura, S., Yamamura, S., & Hasegawa, K. (2000). Phototropic  
785 stimulation induces the conversion of glucosinolate to phototropism-regulating  
786 substances of radish hypocotyls. *Phytochemistry*, 54(3), 275–279.  
787 [https://doi.org/10.1016/S0031-9422\(00\)00080-7](https://doi.org/10.1016/S0031-9422(00)00080-7)
- 788 Hayat, S., Hayat, Q., Alyemini, M. N., Wani, A. S., Pichtel, J., & Ahmad, A. (2012). Role of  
789 proline under changing environments: a review. *Plant Signaling & Behavior*, 7(11),  
790 1456–1466. <https://doi.org/10.4161/psb.21949>
- 791 Heidel, A. J., Clauss, M. J., Kroymann, J., Savolainen, O., & Mitchell-Olds, T. (2006). Natural  
792 variation in MAM within and between populations of *Arabidopsis lyrata* determines  
793 glucosinolate phenotype. *Genetics*, 173(3), 1629–1636.  
794 <https://doi.org/10.1534/genetics.106.056986>
- 795 <https://journal.r-project.org/archive/2013-1/kahle-wickham.pdf>. (n.d.). Retrieved June 16, 2020,  
796 from <https://journal.r-project.org/archive/2013-1/kahle-wickham.pdf>
- 797 Hu, L., Mateo, P., Ye, M., Zhang, X., Berset, J. D., Handrick, V., ... Erb, M. (2018). Plant iron  
798 acquisition strategy exploited by an insect herbivore. *Science*, 361(6403), 694–697.  
799 <https://doi.org/10.1126/science.aat4082>

- 800 Jander, G., Cui, J., Nhan, B., Pierce, N. E., & Ausubel, F. M. (2001). The TASTY locus on  
801 chromosome 1 of *Arabidopsis* affects feeding of the insect herbivore *Trichoplusia ni*.  
802 *Plant Physiology*, *126*(2), 890–898. <https://doi.org/10.1104/pp.126.2.890>
- 803 Kahle, D., & Wickham, H. (2013). ggmap: Spatial Visualization with ggplot2. *The R Journal*,  
804 *5*(1), 144. <https://doi.org/10.32614/RJ-2013-014>
- 805 Kang, H. M., Sul, J. H., Service, S. K., Zaitlen, N. A., Kong, S.-Y., Freimer, N. B., ... Eskin, E.  
806 (2010). Variance component model to account for sample structure in genome-wide  
807 association studies. *Nature Genetics*, *42*(4), 348–354. <https://doi.org/10.1038/ng.548>
- 808 Katz, E., Bagchi, R., Jeschke, V., Rasmussen, A. R. M., Hopper, A., Burow, M., ... Kliebenstein,  
809 D. J. (2020). Diverse allyl glucosinolate catabolites independently influence root growth  
810 and development. *Plant Physiology*, *183*(3), 1376–1390.  
811 <https://doi.org/10.1104/pp.20.00170>
- 812 Katz, E., Nisani, S., Sela, M., Behar, H., & Chamovitz, D. A. (2015). The effect of indole-3-  
813 carbinol on PIN1 and PIN2 in *Arabidopsis* roots. *Plant Signaling & Behavior*, *10*(9),  
814 e1062200. <https://doi.org/10.1080/15592324.2015.1062200>
- 815 Kerwin, R. E., Feusier, J., Muok, A., Lin, C., Larson, B., Copeland, D., ... Kliebenstein, D. J.  
816 (2017). Epistasis × environment interactions among *Arabidopsis thaliana* glucosinolate  
817 genes impact complex traits and fitness in the field. *The New Phytologist*, *215*(3), 1249–  
818 1263. <https://doi.org/10.1111/nph.14646>
- 819 Kerwin, R., Feusier, J., Corwin, J., Rubin, M., Lin, C., Muok, A., ... Kliebenstein, D. J. (2015).  
820 Natural genetic variation in *Arabidopsis thaliana* defense metabolism genes modulates  
821 field fitness. *eLife*, *4*. <https://doi.org/10.7554/eLife.05604>
- 822 Kim, J., Kang, K., Gonzales-Vigil, E., Shi, F., Jones, A. D., Barry, C. S., & Last, R. L. (2012).  
823 Striking natural diversity in glandular trichome acylsugar composition is shaped by  
824 variation at the Acyltransferase2 locus in the wild tomato *Solanum habrochaites*. *Plant*  
825 *Physiology*, *160*(4), 1854–1870. <https://doi.org/10.1104/pp.112.204735>
- 826 Kliebenstein, D J. (2004). Secondary metabolites and plant/environment interactions: a view  
827 through *Arabidopsis thaliana* tinted glasses. *Plant, Cell & Environment*, *27*(6), 675–684.  
828 <https://doi.org/10.1111/j.1365-3040.2004.01180.x>
- 829 Kliebenstein, D J, & Cacho, N. I. (2016). Nonlinear Selection and a Blend of Convergent,  
830 Divergent and Parallel Evolution Shapes Natural Variation in Glucosinolates. In *In S.*  
831 *Kopriva (Ed.), Glucosinolates*. Elsevier. <https://doi.org/10.1016/bs.abr.2016.06.002>
- 832 Kliebenstein, D J, Gershenzon, J., & Mitchell-Olds, T. (2001). Comparative quantitative trait loci  
833 mapping of aliphatic, indolic and benzylic glucosinolate production in *Arabidopsis*  
834 *thaliana* leaves and seeds. *Genetics*, *159*(1), 359–370.
- 835 Kliebenstein, D J, Kroymann, J., Brown, P., Figuth, A., Pedersen, D., Gershenzon, J., &  
836 Mitchell-Olds, T. (2001). Genetic control of natural variation in *Arabidopsis*

- 837 glucosinolate accumulation. *Plant Physiology*, 126(2), 811–825.  
838 <https://doi.org/10.1104/pp.126.2.811>
- 839 Kliebenstein, D J, Lambrix, V. M., Reichelt, M., Gershenzon, J., & Mitchell-Olds, T. (2001).  
840 Gene duplication in the diversification of secondary metabolism: tandem 2-oxoglutarate-  
841 dependent dioxygenases control glucosinolate biosynthesis in Arabidopsis. *The Plant*  
842 *Cell*, 13(3), 681–693. <https://doi.org/10.1105/tpc.13.3.681>
- 843 Kliebenstein, D., Pedersen, D., Barker, B., & Mitchell-Olds, T. (2002). Comparative analysis of  
844 quantitative trait loci controlling glucosinolates, myrosinase and insect resistance in  
845 Arabidopsis thaliana. *Genetics*, 161(1), 325–332.
- 846 Kliebenstein, Daniel J, Figuth, A., & Mitchell-Olds, T. (2002). Genetic architecture of plastic  
847 methyl jasmonate responses in Arabidopsis thaliana. *Genetics*, 161(4), 1685–1696.
- 848 Kliebenstein, Daniel J. (2009). A quantitative genetics and ecological model system:  
849 understanding the aliphatic glucosinolate biosynthetic network via QTLs. *Phytochemistry*  
850 *Reviews : Proceedings of the Phytochemical Society of Europe*, 8(1), 243–254.  
851 <https://doi.org/10.1007/s11101-008-9102-8>
- 852 Kroymann, J., Donnerhacke, S., Schnabelrauch, D., & Mitchell-Olds, T. (2003). Evolutionary  
853 dynamics of an Arabidopsis insect resistance quantitative trait locus. *Proceedings of the*  
854 *National Academy of Sciences of the United States of America*, 100 Suppl 2, 14587–  
855 14592. <https://doi.org/10.1073/pnas.1734046100>
- 856 Kroymann, J., & Mitchell-Olds, T. (2005). Epistasis and balanced polymorphism influencing  
857 complex trait variation. *Nature*, 435(7038), 95–98. <https://doi.org/10.1038/nature03480>
- 858 Lankau, R. A. (2007). Specialist and generalist herbivores exert opposing selection on a  
859 chemical defense. *The New Phytologist*, 175(1), 176–184. <https://doi.org/10.1111/j.1469-8137.2007.02090.x>
- 861 Lankau, R. A., & Kliebenstein, D. J. (2009). Competition, herbivory and genetics interact to  
862 determine the accumulation and fitness consequences of a defence metabolite. *Journal of*  
863 *Ecology*, 97(1), 78–88. <https://doi.org/10.1111/j.1365-2745.2008.01448.x>
- 864 Lankau, R. A., & Strauss, S. Y. (2007). Mutual feedbacks maintain both genetic and species  
865 diversity in a plant community. *Science*, 317(5844), 1561–1563.  
866 <https://doi.org/10.1126/science.1147455>
- 867 Lankau, R. A., & Strauss, S. Y. (2008). Community complexity drives patterns of natural  
868 selection on a chemical defense of Brassica nigra. *The American Naturalist*, 171(2), 150–  
869 161. <https://doi.org/10.1086/524959>
- 870 Liaw, A., & Wiener, M. (2002). Classification and Regression by randomForest. *R News* 2.
- 871 MacQueen, A., Sun, X., & Bergelson, J. (2016). Genetic architecture and pleiotropy shape costs  
872 of Rps2-mediated resistance in Arabidopsis thaliana. *Nature Plants*, 2, 16110.  
873 <https://doi.org/10.1038/nplants.2016.110>

- 874 Malcolm, S. B. (1994). Milkweeds, monarch butterflies and the ecological significance of  
875 cardenolides. *Chemoecology*, 5-6(3-4), 101–117. <https://doi.org/10.1007/BF01240595>
- 876 Malinovsky, F. G., Thomsen, M.-L. F., Nintemann, S. J., Jagd, L. M., Bourguine, B., Burow, M.,  
877 & Kliebenstein, D. J. (2017). An evolutionarily young defense metabolite influences the  
878 root growth of plants via the ancient TOR signaling pathway. *eLife*, 6.  
879 <https://doi.org/10.7554/eLife.29353>
- 880 Mithen, R., Clarke, J., Lister, C., & Dean, C. (1995). Genetics of aliphatic glucosinolates. III.  
881 Side chain structure of aliphatic glucosinolates in *Arabidopsis thaliana*. *Heredity*, 74(2),  
882 210–215. <https://doi.org/10.1038/hdy.1995.29>
- 883 Moore, B. M., Wang, P., Fan, P., Leong, B., Schenck, C. A., Lloyd, J. P., ... Shiu, S.-H. (2019).  
884 Robust predictions of specialized metabolism genes through machine learning.  
885 *Proceedings of the National Academy of Sciences of the United States of America*,  
886 116(6), 2344–2353. <https://doi.org/10.1073/pnas.1817074116>
- 887 Opitz, S. E. W., & Müller, C. (2009). Plant chemistry and insect sequestration. *Chemoecology*,  
888 19(3), 117–154. <https://doi.org/10.1007/s00049-009-0018-6>
- 889 Petersen, B. L., Chen, S., Hansen, C. H., Olsen, C. E., & Halkier, B. A. (2002). Composition and  
890 content of glucosinolates in developing *Arabidopsis thaliana*. *Planta*, 214(4), 562–571.  
891 <https://doi.org/10.1007/s004250100659>
- 892 Pfalz, M., Vogel, H., Mitchell-Olds, T., & Kroymann, J. (2007). Mapping of QTL for resistance  
893 against the crucifer specialist herbivore *Pieris brassicae* in a new *Arabidopsis* inbred line  
894 population, Da(1)-12 x Ei-2. *Plos One*, 2(6), e578.  
895 <https://doi.org/10.1371/journal.pone.0000578>
- 896 Ramos-Onsins, S. E., Stranger, B. E., Mitchell-Olds, T., & Aguadé, M. (2004). Multilocus  
897 Analysis of Variation and Speciation in the Closely Related Species *Arabidopsis halleri*  
898 and *A. lyrata*. *Genetics*, 166(1), 373–388. <https://doi.org/10.1534/genetics.166.1.373>
- 899 Raybould, A. F., & Moyes, C. L. (2001). The ecological genetics of aliphatic glucosinolates.  
900 *Heredity*, 87(Pt 4), 383–391. <https://doi.org/10.1046/j.1365-2540.2001.00954.x>
- 901 Rodman, James E., Kruckeberg, A. R., & Al-Shehbaz, I. A. (1981). Chemotaxonomic diversity  
902 and complexity in seed glucosinolates of caulanthus and streptanthus (cruciferae).  
903 *Systematic Botany*, 6(3), 197. <https://doi.org/10.2307/2418282>
- 904 Rodman, James Eric. (1980). Population variation and hybridization in sea-rockets (cakile,  
905 cruciferae): seed glucosinolate characters. *American Journal of Botany*, 67(8), 1145–  
906 1159. <https://doi.org/10.1002/j.1537-2197.1980.tb07748.x>
- 907 Salehin, M., Li, B., Tang, M., Katz, E., Song, L., Ecker, J. R., ... Estelle, M. (2019). Auxin-  
908 sensitive Aux/IAA proteins mediate drought tolerance in *Arabidopsis* by regulating  
909 glucosinolate levels. *Nature Communications*, 10(1), 4021.  
910 <https://doi.org/10.1038/s41467-019-12002-1>



- 911 Schilmiller, A. L., Pichersky, E., & Last, R. L. (2012). Taming the hydra of specialized  
912 metabolism: how systems biology and comparative approaches are revolutionizing plant  
913 biochemistry. *Current Opinion in Plant Biology*, 15(3), 338–344.  
914 <https://doi.org/10.1016/j.pbi.2011.12.005>
- 915 Sønderby, I. E., Geu-Flores, F., & Halkier, B. A. (2010). Biosynthesis of glucosinolates--gene  
916 discovery and beyond. *Trends in Plant Science*, 15(5), 283–290.  
917 <https://doi.org/10.1016/j.tplants.2010.02.005>
- 918 Szakiel, A., Pączkowski, C., & Henry, M. (2011). Influence of environmental biotic factors on  
919 the content of saponins in plants. *Phytochemistry Reviews : Proceedings of the*  
920 *Phytochemical Society of Europe*, 10(4), 493–502. [https://doi.org/10.1007/s11101-010-](https://doi.org/10.1007/s11101-010-9164-2)  
921 [9164-2](https://doi.org/10.1007/s11101-010-9164-2)
- 922 Thakur, P., & Rai, V. (1982). Dynamics of amino acid accumulation of two differentially  
923 drought resistant *Zea mays* cultivars in response to osmotic stress. *Environmental and*  
924 *Experimental Botany*, 22(2), 221–226. [https://doi.org/10.1016/0098-8472\(82\)90042-9](https://doi.org/10.1016/0098-8472(82)90042-9)
- 925 The 1001 Genomes Consortium. (2016). 1,135 Genomes reveal the global pattern of  
926 polymorphism in *Arabidopsis thaliana*. *Cell*, 166(2), 481–491.  
927 <https://doi.org/10.1016/j.cell.2016.05.063>
- 928 Todesco, M., Balasubramanian, S., Hu, T. T., Traw, M. B., Horton, M., Epple, P., ... Weigel, D.  
929 (2010). Natural allelic variation underlying a major fitness trade-off in *Arabidopsis*  
930 *thaliana*. *Nature*, 465(7298), 632–636. <https://doi.org/10.1038/nature09083>
- 931 Turner, S. D. (2014). qqman: an R package for visualizing GWAS results using Q-Q and  
932 manhattan plots. *BioRxiv*. <https://doi.org/10.1101/005165>
- 933 Uremis, I., Arslan, M., Sangun, M. K., Uygur, V., & Isler, N. (2009). Allelopathic potential of  
934 rapeseed cultivars on germination and seedling growth of weeds. *Asian Journal of*  
935 *Chemistry*, 21(3), 2170–2184.
- 936 Wentzell, A. M., & Kliebenstein, D. J. (2008). Genotype, age, tissue, and environment regulate  
937 the structural outcome of glucosinolate activation. *Plant Physiology*, 147(1), 415–428.  
938 <https://doi.org/10.1104/pp.107.115279>
- 939 Wentzell, A. M., Rowe, H. C., Hansen, B. G., Ticconi, C., Halkier, B. A., & Kliebenstein, D. J.  
940 (2007). Linking metabolic QTLs with network and cis-eQTLs controlling biosynthetic  
941 pathways. *PLoS Genetics*, 3(9), 1687–1701.  
942 <https://doi.org/10.1371/journal.pgen.0030162>
- 943 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L., François, R., ... Yutani, H.  
944 (2019). Welcome to the tidyverse. *The Journal of Open Source Software*, 4(43), 1686.  
945 <https://doi.org/10.21105/joss.01686>
- 946 Windsor, A. J., Reichelt, M., Figuth, A., Svatos, A., Kroymann, J., Kliebenstein, D. J., ...  
947 Mitchell-Olds, T. (2005). Geographic and evolutionary diversification of glucosinolates

- 948 among near relatives of *Arabidopsis thaliana* (Brassicaceae). *Phytochemistry*, 66(11),  
949 1321–1333. <https://doi.org/10.1016/j.phytochem.2005.04.016>
- 950 Wolters, H., & Jürgens, G. (2009). Survival of the flexible: hormonal growth control and  
951 adaptation in plant development. *Nature Reviews. Genetics*, 10(5), 305–317.  
952 <https://doi.org/10.1038/nrg2558>
- 953 Wright, S. I., Lauga, B., & Charlesworth, D. (2002). Rates and patterns of molecular evolution in  
954 inbred and outbred *Arabidopsis*. *Molecular Biology and Evolution*, 19(9), 1407–1420.  
955 <https://doi.org/10.1093/oxfordjournals.molbev.a004204>
- 956 Yamada, K., Hasegawa, T., Minami, E., Shibuya, N., Kosemura, S., Yamamura, S., &  
957 Hasegawa, K. (2003). Induction of myrosinase gene expression and myrosinase activity  
958 in radish hypocotyls by phototropic stimulation. *Journal of Plant Physiology*, 160(3),  
959 255–259. <https://doi.org/10.1078/0176-1617-00950>
- 960 Yang, C. W., Lin, C. C., & Kao, C. H. (2000). Proline, ornithine, arginine and glutamic acid  
961 contents in detached rice leaves. *Biologia Plantarum*, 43(2), 305–307.  
962 <https://doi.org/10.1023/A:1002733117506>
- 963 Yu, G. (2020). Using ggtree to Visualize Data on Tree-Like Structures. *Current Protocols in*  
964 *Bioinformatics*, 69(1), e96. <https://doi.org/10.1002/cpbi.96>
- 965 Zhu, W., Zaidem, M., Van deaana 1Weyer, A.-L., Gutaker, R. M., Chae, E., Kim, S.-T., ...  
966 Weigel, D. (2018). Modulation of ACD6 dependent hyperimmunity by natural alleles of  
967 an *Arabidopsis thaliana* NLR resistance gene. *PLoS Genetics*, 14(9), e1007628.  
968 <https://doi.org/10.1371/journal.pgen.1007628>
- 969 Züst, T., & Agrawal, A. A. (2017). Trade-Offs Between Plant Growth and Defense Against  
970 Insect Herbivory: An Emerging Mechanistic Synthesis. *Annual Review of Plant Biology*,  
971 68, 513–534. <https://doi.org/10.1146/annurev-arplant-042916-040856>
- 972 Züst, T., Heinricher, C., Grossniklaus, U., Harrington, R., Kliebenstein, D. J., & Turnbull, L. A.  
973 (2012). Natural enemies drive geographic variation in plant defenses. *Science*, 338(6103),  
974 116–119. <https://doi.org/10.1126/science.1226397>
- 975