



26 **ABSTRACT**

27 **Background :** In the plant sciences, results of laboratory studies often do not translate well  
28 to the field because lab growth conditions are very different from field conditions. To help  
29 close this lab-field gap, we developed a new strategy for studying the wiring of plant traits  
30 directly in the field, based on molecular profiling and phenotyping of individual plants of the  
31 same genetic background grown in the same field. This single-plant omics strategy leverages  
32 uncontrolled micro-environmental variation across the field and stochastic variation among  
33 the individual plants as information sources, rather than controlled perturbations. Here, we  
34 use single-plant omics on winter-type *Brassica napus* (rapeseed) plants to investigate to what  
35 extent rosette-stage gene expression profiles can be linked to the early and late phenotypes  
36 of individual field-grown plants.

37

38 **Results :** We find that rosette leaf gene expression in autumn has substantial predictive  
39 power for both autumnal leaf phenotypes and final yield in spring. Many of the top predictor  
40 genes are linked to developmental processes known to occur in autumn in winter-type *B.*  
41 *napus* accessions, such as the juvenile-to-adult and vegetative-to-reproductive phase  
42 transitions, indicating that the yield potential of winter-type *B. napus* is influenced by  
43 autumnal development.

44

45 **Conclusions :** Our results show that profiling individual plants under uncontrolled field  
46 conditions is a valid strategy for identifying genes and processes influencing crop yield in the  
47 field.

48

49 **KEYWORDS**

50 single-plant omics ; *Brassica napus* ; field trial ; machine learning ; transcriptome-based

51 phenotype prediction.

52

53 **BACKGROUND**

54 One of the major aims of molecular biology research is to unravel how genes influence  
55 phenotypes. This usually involves applying perturbations to the genome or growth  
56 environment of an organism of interest and analyzing the ensuing molecular and phenotypic  
57 responses. Generally, well-chosen perturbations are applied in a controlled experimental  
58 setting, and technical and biological replicates are performed to allow for sufficiently  
59 powerful analyses despite noise in the data. Noise in this context may refer to measurement  
60 errors, noise due to uncontrolled factors in the experimental setup, or noise due to cellular  
61 or environmental stochasticity. The main purpose of avoiding or averaging out such noise is  
62 to facilitate causal interpretation of the link between a perturbation and its molecular and  
63 phenotypic effects.

64

65 It is becoming increasingly clear however that data noise caused by uncontrolled  
66 experimental factors and even purely stochastic effects can be a valuable source of  
67 information, instead of merely a nuisance. Several studies have shown that stochastic gene  
68 expression noise in single cells can be used to infer regulatory influences (1-3). Gene networks  
69 are also increasingly inferred from single-cell gene expression datasets in which differences  
70 among cells are not purely due to stochastic effects in an otherwise homogeneous cell

71 population, but reflect additional uncontrolled heterogeneity among cells, e.g. in the  
72 temporal progression of a cell differentiation program (4-10).

73

74 In addition, several studies have investigated the information content of ‘noise’ datasets in  
75 which the profiled entities are multicellular individuals rather than single cells. Bhosale, Jewell  
76 et al. (11) found that gene expression noise among individual *Arabidopsis thaliana* plants  
77 grown under the same conditions harbored as much information on the function of genes as  
78 gene expression responses to controlled perturbations. The dataset analyzed by Bhosale,  
79 Jewell et al. (11) was however not ideal because it contained data on plants of three different  
80 accessions grown in six different labs (12), causing lab and accession effects that had to be  
81 removed computationally to uncover the individual plant noise of interest. Recently, a study  
82 on a cleaner *A. thaliana* seedling dataset confirmed that gene expression noise among  
83 individuals of the same background grown under the same lab conditions contains useful  
84 information on gene functions and regulatory relationships (13).

85

86 A common denominator in the aforementioned studies is that even under controlled  
87 conditions, each cell or individual is subject to a set of stochastic or other perturbations that  
88 escape experimental control, and that these uncontrolled perturbations, like any  
89 perturbations, generate responses that contain valuable information on the wiring of gene  
90 networks. Although most studies to date focused on the information content of noise under  
91 controlled lab conditions, there is no reason to believe that ‘noise’ datasets generated under  
92 less controlled conditions would be less valuable. On the contrary, studies performed in a  
93 more natural setting in which organisms are subject to uncontrolled perturbations may yield

94 information that cannot easily be recovered from experiments under controlled lab  
95 conditions.

96

97 In the plant sciences for instance, controlled growth conditions in a laboratory are generally  
98 very different from field conditions, in which plants are subject to a plethora of highly variable  
99 environmental cues that often have non-additive phenotypic effects (14-21). Results obtained  
100 in the laboratory therefore often translate poorly to the field (14, 22-26). Narrowing this lab-  
101 field gap is essential to speed up the development of new crop varieties and optimized  
102 agricultural practices, both of which are direly needed in view of the current challenges posed  
103 by world population growth, land use and climate change. One option to narrow the lab-field  
104 gap is to make lab conditions more field-like (22), but the decreased experimental control this  
105 implies challenges traditional experimental design practices to e.g. ensure reproducibility.  
106 Another option is to perform interventional experiments in the field rather than the lab, but  
107 controlled interventions in a field may be costly and the level of control that can be achieved  
108 is often limited (22). Observational ‘uncontrolled perturbation’ studies on the other hand can  
109 easily be set up in the field. Observational data come with their own array of challenges  
110 however, e.g. that many of the perturbations influencing the study subjects may remain  
111 unobserved and hence unknown, and that it is generally much more challenging to establish  
112 cause-effect relationships from observational data (27). Nevertheless, even purely  
113 correlational data generated in the field may help narrow the lab-field gap in plant sciences.

114

115 To assess the information content of plant molecular responses to uncontrolled perturbations  
116 occurring in a field environment, we previously generated transcriptome and metabolome  
117 data on the primary ear leaf of 60 individual *Zea mays* (maize) plants of the same genetic

118 background grown in the same field (28). Similar to what was found for lab-grown *A. thaliana*  
119 plants (11), the transcriptomes of the individual field-grown maize plants were found to  
120 contain as much information on maize gene function as transcriptomes profiling the response  
121 of maize plants to controlled perturbations in the lab. In addition, we found that the single-  
122 plant transcriptome and metabolome data had better-than-random predictive power for  
123 several phenotypes that were measured for the individual plants, and the prediction models  
124 also produced sensible candidate genes for these phenotypes (28). However, only a few  
125 phenotypes were measured in this study, and they were either closely associated with the  
126 material sampled for molecular profiling, not fully developed or both.

127

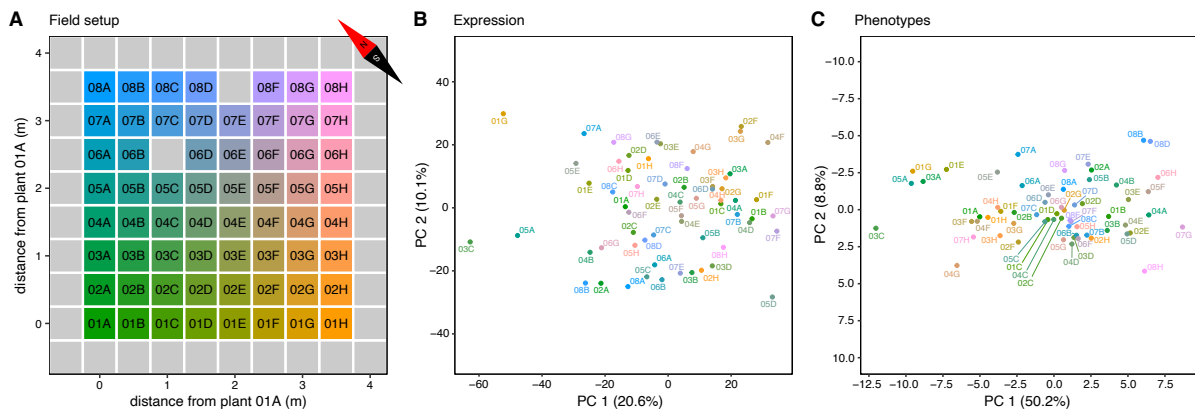
128 Here, we investigate in more detail how much phenotype information can be extracted from  
129 the transcriptomes of single plants subject to uncontrolled perturbations under field  
130 conditions. To this end, we profiled the rosette-stage leaf transcriptome of individual field-  
131 grown plants of the winter-type accession Darmor of *Brassica napus* (rapeseed), an important  
132 oilseed crop (29). Additionally, a wide range of phenotypes was measured for all plants  
133 throughout the growing season. We find that the autumnal leaf transcriptomes of the  
134 individual plants do not only have predictive power for autumnal leaf phenotypes but also for  
135 yield phenotypes measured more than 5 months later, such as silique count and total seed  
136 weight. Furthermore, we find that many of the genes that feature prominently in our  
137 predictive models are related to developmental processes known to occur in autumn in  
138 winter-type *Brassica napus*, in particular the juvenile-to-adult and vegetative-to-reproductive  
139 phase transitions. Our results suggest that micro-environmental variations across the field  
140 cause a gradual buildup of developmental differences among plants that ultimately result in  
141 yield differences at the end of the growing season.

## 142 RESULTS

### 143 Field trial, expression profiling and phenotyping

144 One hundred *Brassica napus* plants of the winter-type accession Darmor were grown in a field  
145 in a 10x10 equispaced grid pattern with 0.5 m distance between rows and columns (**Fig. 1**).  
146 On November 28, 2016, the eighth rosette leaf (leaf 8) of 62 non-border plants was harvested,  
147 and the harvested leaves were expression-profiled individually (see Methods and **Additional**  
148 **File 1: Table S1**). After leaf sampling, the plants were allowed to overwinter and set seed in  
149 spring. 62 phenotypes were recorded for all plants, ranging from rosette areas and individual  
150 leaf measurements in autumn to root and shoot measurements at harvest the following  
151 spring (**Additional File 1: Table S1**). Likely because of the low planting density, many of the  
152 plants developed one or more secondary inflorescence stems at ground level, which is not  
153 usually observed for *B. napus* grown under lab conditions or in the field at agronomically  
154 relevant planting densities. These secondary stems (further referred to as side stems) were  
155 harvested separately from the primary inflorescence stem with its cauline secondary  
156 inflorescences (further referred to as stem 1). Several yield phenotypes were measured for  
157 both stem 1 and the entire shoot (i.e. stem 1 plus side stems), including dry weight, seed  
158 weight, seed count and silique count. Cauline secondary inflorescence stems on stem 1 and  
159 tertiary inflorescence stems on the side stems (both further referred to as branches) were  
160 also counted, and branch counts are reported for both stem 1 and the entire shoot (the latter  
161 being the sum of branch counts on stem 1 and the side stems). Shoot growth phenotypes  
162 such as the time of maximum shoot growth, the maximum shoot growth rate and the end of  
163 shoot growth were derived from plant height time series data through curve fitting (see  
164 Methods). Several phenotypes were defined as ratios of other phenotypes, e.g. the ratio of

165 total seed weight to shoot dry weight and the ratio of the total number of seeds to the total  
166 number of siliques per plant.



167

168 **Fig. 1** Field trial layout and PCA plots for gene expression and phenotypes. **A** Plants were sown  
169 on a 10x10 equispaced grid with 0.5 m between rows and columns. Plant identifiers combine  
170 a number indicating the row (01-08) and a letter indicating the column (A-H) in which the  
171 plant was sown. Only plants with leaf 8 gene expression and phenotype profiles are labeled,  
172 border plants and grid positions at which no plants emerged are indicated by grey squares. **B**  
173 Plot of the first two principal components of the leaf 8 gene expression dataset, after  
174 normalization and RNA-seq batch correction (see Methods). **C** Plot of the first two principal  
175 components in the phenotype dataset. Individual plants in **B** and **C** are colored according to  
176 the color gradient in **A**, with similar coloring of plants indicating spatial proximity in the field.

177

## 178 Exploratory data analysis

179 Principal component analysis (PCA) suggests that there are no subpopulations of plants with  
180 distinct expression or phenotype profiles (**Fig. 1**). A few relative outliers are visible however,  
181 e.g. plant 04G in the phenotype PCA plot (**Fig. 1C**), a very small plant that yielded barely any  
182 seeds. Single-nucleotide polymorphism (SNP) analysis of the RNA-seq data (see Methods) did



183 not uncover signs of substantial genetic substructure in the plant population (**Additional File**  
184 **2: Fig. S1**).

185

186 Mapping of the field coordinates on the expression and phenotype PCA plots on the other  
187 hand suggests that there is spatial structure in the data (**Fig. 1**). However, the levels of only  
188 169 out of 76,808 transcripts and 1 out of 41 phenotypes (root system width) were found to  
189 be significantly spatially autocorrelated across the field (Moran's I, Benjamini-Hochberg (BH)  
190 adjusted permutation test  $p$ -values ( $q$ -values)  $< 0.05$ , **Additional File 3: Table S2, Additional**  
191 **File 2: Fig. S2**). In a previous study on a similar number of field-grown maize plants (28),  
192 14.17% of transcripts were found to be significantly spatially autocorrelated at  $q \leq 0.01$ , which  
193 is considerably more than the 0.22% recovered here at  $q \leq 0.05$ . This may be due to  
194 differences in the way Moran's I values and their significance were calculated in Cruz, De  
195 Meyer et al. (28) versus the present study (see Methods). To assess whether some functional  
196 classes of genes have on average a stronger or weaker spatial autocorrelation signal than  
197 other classes, regardless of the statistical significance of the Moran's I values, two-sided  
198 Mann-Whitney U (MWU) tests (30) were performed on the transcript list ranked in order of  
199 decreasing Moran's I value. Genes involved in e.g. photosynthesis, translation, the response  
200 to abiotic stimuli, response to cytokinin, regulation of circadian rhythm, photoperiodism and  
201 the vegetative to reproductive phase transition were found to have a significantly higher  
202 Moran's I on average than other genes (MWU  $q \leq 0.05$ , **Additional File 3: Table S2**). This  
203 suggests that there is spatial patterning in the data, but that its discovery may be hampered  
204 by a lack of statistical power due to the small size of the field trial.

205

206 Most continuous phenotypes and high-count discrete phenotypes (e.g. seed and silique  
207 counts) are at least approximately normally distributed (Anderson-Darling and Shapiro-Wilk  
208 normality tests,  $p > 0.01$ , **Additional File 4: Table S3**), with the exception of five ratio  
209 phenotypes (seeds per silique, seeds per silique stem 1, seed weight/dry weight stem 1, total  
210 seed weight/shoot dry weight and branches per stem), leaf count (74 DAS) and two shoot  
211 growth phenotypes (time of max shoot growth and end of shoot growth). Many of these  
212 phenotypes exhibit relative outliers that may influence normality testing results (**Additional**  
213 **File 2: Fig. S3**). When removing outliers (see Methods), four additional phenotypes (seeds per  
214 silique, seeds per silique stem 1, stem 1 seed weight/ stem 1 dry weight, total seed  
215 weight/shoot dry weight) were found to be approximately normally distributed (Anderson-  
216 Darling and Shapiro-Wilk normality test,  $p > 0.01$ , **Additional File 4: Table S3**).

217

218 Some phenotypes were found to be more variable across the field than others. Dry weight,  
219 seed and silique phenotypes at harvest are the most variable, with coefficients of variation  
220 (CVs) between 43.7% and 51.9% (**Additional File 4: Table S3**). Taproot length also has a high  
221 CV (42.8%). Plant height (278 DAS) and shoot growth parameters exhibit the lowest CV values  
222 ( $< 7\%$ ). Most ratio phenotypes also have relatively low CV values ( $\leq 20.3\%$ ), with the exception  
223 of siliques per branch (35.3%), siliques per branch stem 1 (35.0%) and branches per stem  
224 (33.6%). When removing outliers, the CV of some of these ratio phenotypes is further  
225 reduced, notably for seed weight stem 1/dry weight stem 1 (20.3%  $\rightarrow$  9.5%), total seed  
226 weight/shoot dry weight (18.1%  $\rightarrow$  8.9%), seeds per silique (19.5%  $\rightarrow$  14.7%) and seeds per  
227 silique stem 1 (19.4%  $\rightarrow$  14.6%). Leaf and branch phenotypes generally exhibit intermediate  
228 CVs. Whereas leaf 8 fresh weight (81 DAS), leaf 8 area (81 DAS), total branch count and rosette  
229 area (42 DAS) have a CV  $\geq 30\%$ , other leaf 8 and leaf 6 phenotypes and branch count stem 1

230 exhibit a CV in the range 19.1%-23.2%, and leaf 8 chlorophyll content (81 DAS) has a CV of  
231 only 12.5%.

232

233 Gene expression also exhibits substantial variability across the field. Ignoring genes expressed  
234 in less than 10 samples, the median gene has an expression CV of 34.2% (**Additional File 4:**  
235 **Table S3**). To investigate whether some classes of genes vary more in expression than others  
236 across the field, we ranked *B. napus* genes based on a normalized version of their expression  
237 CV (*normCV*, see Methods and **Additional File 4: Table S3**). MWU tests (30) were performed  
238 to assess whether any Gene Ontology (GO) biological processes are represented more at the  
239 top or bottom of the *normCV*-ranked gene list than expected by chance (**Additional File 4:**  
240 **Table S3**). As observed in earlier studies on populations of lab-grown *Arabidopsis thaliana*  
241 Col-0 plants (31) and field-grown *Zea mays* B104 plants (28), genes involved in photosynthesis  
242 and responses to biotic and abiotic stimuli were found to be on average more variably  
243 expressed than other genes, while genes involved in housekeeping functions related to  
244 protein, RNA and DNA metabolism were found to be on average more stably expressed across  
245 the field (**Additional File 4: Table S3**). To what extent high gene expression variability is due  
246 to either variability in the levels of external stimuli experienced by the individual plants or due  
247 to a higher intrinsic noisiness of a gene's expression levels (on the scale of entire leaves) is  
248 unclear. Some categories of genes with more variable expression across the field, such as  
249 genes involved in photosynthesis or response to abiotic stimuli, also exhibit higher Moran's I  
250 values on average, suggesting that their variability may be linked to external stimuli that are  
251 spatially patterned. On the other hand, most genes with highly variable expression do not  
252 exhibit strong spatial patterns (**Additional File 2: Fig. S4**), which indicates that their expression

253 variability may be caused by intrinsic stochastic factors, or alternatively by extrinsic factors  
254 that are not spatially autocorrelated at the field sampling resolution employed.

255

## 256 **Linking phenotypes to the leaf 8 expression profiles of single genes**

257 To assess how much information leaf 8 gene expression profiles contain on the phenotypes  
258 of individual plants, we used linear mixed-effect (lme) models to associate plant phenotypes  
259 with the autumnal leaf 8 expression profile of single genes, taking into account spatial  
260 autocorrelation effects (see Methods). Between 11,986 and 14,032 gene expression profiles,  
261 out of 76,808, were found to be significantly associated ( $q \leq 0.05$ ) with leaf 8 phenotypes such  
262 as leaf 8 length, width, area and fresh weight (**Table 1 , Additional File 5: Table S4**). That leaf  
263 8 phenotypes yield more associated genes than other phenotypes is not surprising, given that  
264 leaf 8 was used for gene expression profiling. Next to leaf 8 phenotypes, also other leaf and  
265 rosette phenotypes feature more associated genes than non-leaf phenotypes, except for leaf  
266 6 length (74 DAS). The gene sets associated with leaf phenotypes are generally significantly  
267 enriched (hypergeometric test,  $q \leq 0.05$ ) in genes involved in e.g. response to biotic and  
268 abiotic stimuli (salt), photosynthesis, circumnutation, cell wall biogenesis, amino acid  
269 metabolism and response to sulfate and nitrogen starvation (**Additional File 6: Table S5**).  
270 Additionally, leaf phenotype-related gene lists show significant enrichment, notably among  
271 transcription factors, in genes involved in dorsal/ventral, adaxial/abaxial and radial pattern  
272 formation, phloem, xylem and procambium histogenesis, and meristem development  
273 (**Additional File 6: Table S5**).

274

275 Interestingly, appreciable numbers of gene-phenotype associations were found as well for  
276 several phenotypes that are only distantly related in space and time to the leaf 8 material

277 profiled for RNA-seq. In particular seed, silique and shoot dry weight phenotypes yielded high  
278 numbers of associated genes, ranging from 1,859 genes for total shoot dry weight to 1,248  
279 genes for the silique count on stem 1 at harvest (**Table 1, Additional File 5: Table S4**). Many  
280 of the gene sets associated with these phenotypes are enriched in genes involved in nitrate  
281 assimilation, superoxide metabolism, circumnutation, circadian rhythm, response to biotic  
282 and abiotic stimuli (cold, salt, water deprivation), response to nutrient levels (nitrogen,  
283 sulphate and phosphate starvation), and, in particular among transcription factors, phosphate  
284 ion homeostasis, histone modification, regulation of the vegetative to reproductive phase  
285 transition and floral organ morphogenesis (**Additional File 6: Table S5**). 1,110 genes were  
286 found associated with the branch count on stem 1, with GO enrichments similar to those  
287 obtained for dry weight, silique and seed phenotypes (**Additional File 6: Table S5**). In  
288 contrast, the total branch count phenotype only yields a set of 89 associated genes (**Table 1,**  
289 **Additional File 5: Table S4**), which is however also strongly enriched in e.g. superoxide  
290 metabolism and salt stress genes. The fact that the total branch count is composed of cauline  
291 secondary inflorescence stems on stem 1 and tertiary inflorescence stems on the side stems  
292 may render this phenotype less relevant.

293

294 Phenotypes with very low CV such as leaf 8 chlorophyll content (81 DAS), the maximum shoot  
295 growth rate and end of shoot growth yielded no significantly associated genes, suggesting  
296 that the biological variation of these phenotypes is limited and that the observed variation  
297 may be dominated by technical noise (**Table 1, Additional File 5: Table S4**). The phenotype  
298 with the lowest CV on the other hand, the time of maximal shoot growth (CV=0.6%), features  
299 3,498 significant leaf 8 gene expression correlates. The associated gene set is strongly  
300 enriched in genes involved in e.g. cell wall biogenesis and response to biotic stimuli

301 **(Additional File 6: Table S5).** For plant height (278 DAS) (CV = 6.8%), 99 associated genes are  
 302 found with minor GO enrichments.

303

Phenotype	All genes			Transcription factors		
	# Significant	Most significant	<i>q</i>	# Significant	Most significant	<i>q</i>
leaf 8 length (76 DAS)	14032	BnaC07g39340D	4.35E-14	453	BnaA05g33840D	4.44E-09
leaf 8 width (76 DAS)	13695	BnaC07g39340D	1.62E-15	429	BnaA05g33840D	6.98E-09
leaf 8 width (81 DAS)	13605	BnaA02g18860D	1.37E-12	420	BnaA05g33840D	1.36E-08
leaf 8 fresh weight (81 DAS)	12989	BnaC07g39340D	5.79E-12	412	BnaAnng02740D	4.17E-08
leaf 8 area (81 DAS)	12569	BnaC07g39340D	1.29E-13	408	BnaAnng02740D	1.64E-08
leaf 8 length (81 DAS)	11986	BnaA01g14450D	6.39E-14	383	BnaA05g27750D	5.25E-08
leaf count (74 DAS)	10442	BnaC04g49060D	1.86E-07	313	BnaA05g33840D	2.16E-06
rosette area (42 DAS)	7196	BnaC09g39140D	7.14E-06	212	BnaA06g39930D	1.06E-04
leaf 6 width (74 DAS)	5386	BnaA09g04980D	2.05E-05	184	BnaA05g27750D	2.76E-04
time of max shoot growth	3498	BnaC06g28860D	7.29E-07	89	BnaC04g03950D	5.42E-06
total shoot dry weight	1859	BnaA05g29010D	1.37E-06	76	BnaCnng05590D	2.36E-04
total shoot dry weight (w/o seeds)	1802	BnaA05g29010D	8.88E-07	68	BnaAnng37500D	4.89E-04
dry weight stem 1	1612	BnaA05g29010D	2.79E-06	72	BnaC01g37260D	5.87E-05
dry weight stem 1 (w/o seeds)	1611	BnaA05g29010D	6.82E-06	75	BnaA08g12050D	3.13E-04
total seed weight	1598	BnaA06g35450D	1.02E-05	66	BnaCnng05590D	6.05E-05
total seed count	1545	BnaA06g35450D	9.10E-06	63	BnaCnng05590D	1.12E-04
seed weight stem 1	1539	BnaA05g29010D	1.65E-05	66	BnaC01g37260D	1.65E-05
total silique count	1520	BnaA06g35450D	3.92E-05	64	BnaCnng05590D	7.12E-04
seed count stem 1	1449	BnaC01g37260D	2.58E-05	67	BnaC01g37260D	2.58E-05
leaf 6 length (74 DAS)	1345	BnaA09g04980D	2.37E-04	34	BnaA02g18720D	8.39E-03
silique count stem 1	1248	BnaA05g29010D	6.32E-05	56	BnaC01g37260D	6.32E-05
branch count stem 1	1110	BnaA06g35450D	2.52E-04	39	BnaC01g37260D	1.95E-03
siliques per branch stem 1	593	BnaA01g17100D	2.24E-03	29	BnaC01g37260D	2.48E-03
total seed weight/shoot dry weight	458	BnaC09g50070D	2.24E-05	13	BnaC04g55440D	4.81E-04
seed weight stem 1/dry weight stem 1	280	BnaA03g50380D	8.90E-04	6	BnaC03g62970D	3.39E-03
branch count stem 1/length stem 1	240	BnaAnng11300D	8.22E-03	5	BnaC01g37260D	4.04E-02
seeds per silique stem 1	233	BnaC05g45470D	8.73E-04	9	BnaC04g03950D	4.42E-03
seeds per silique	112	BnaA08g07570D	9.01E-04	3	BnaC04g03950D	1.43E-02
plant height (278 DAS)	99	BnaA06g34140D	7.20E-04	5	BnaC01g37260D	2.48E-02
total branch count	89	BnaA06g35450D	2.39E-03	1	BnaCnng05590D	5.82E-03
root system width	4	BnaA01g06800D	7.53E-03	0	-	-
siliques per branch	3	BnaAnng39720D	2.22E-02	0	-	-
branches per stem	0	-	-	0	-	-
taproot length	0	-	-	0	-	-
leaf 8 chlorophyll content (81 DAS)	0	-	-	0	-	-
max shoot growth rate	0	-	-	0	-	-
end of shoot growth	0	-	-	0	-	-
rosette lesions (74 DAS)	0	-	-	0	-	-
leaf 6 lesions (74 DAS)	0	-	-	0	-	-
leaf 8 lesions (76 DAS)	0	-	-	0	-	-
stem count	0	-	-	0	-	-

304

305 **Table 1** Numbers of significant gene expression-phenotype associations. For any given  
 306 phenotype, results are reported on the complete gene set ( $n=76,808$ ; 'All genes' columns)  
 307 and on the set of transcription factors ( $n=2,521$ ; 'Transcription factors' columns). In both  
 308 cases, the results shown include (from left to right) the total number of significant gene  
 309 expression-phenotype associations ( $q \leq 0.05$ ), the most significant gene and its  $q$ -value.

310 No genes were found associated ( $q \leq 0.05$ ) with taproot length and only four with root system  
311 width, suggesting that autumnal leaf 8 gene expression may not contain a lot of information  
312 on root phenotypes. On the other hand, given the difficulty of recovering intact root systems  
313 from the soil, it is not excluded that root measurement errors may have influenced these  
314 results.

315

316 Ratio phenotypes exhibit between 0 and 593 associated genes. In particular the branches per  
317 stem and siliques per branch ratios do poorly (0 and 3 associated genes, respectively). Both  
318 involve the total branch count, which is itself only associated with 89 genes. Ratios involving  
319 the branch count on stem 1 on the other hand yield between 240 and 593 associated genes.  
320 One potential reason for ratio phenotypes having at most a few hundred gene associations is  
321 that ratios suffer from increased error levels due to the propagation of measurement errors  
322 from both the numerator and denominator. This may be particularly problematic for ratios of  
323 highly correlated variables such as the seeds per silique and seed weight/dry weight  
324 phenotypes (both for stem 1 and the entire shoot), which exhibit a low CV and likely have  
325 even lower true biological variation. No genes were found associated at  $q \leq 0.05$  with  
326 qualitative or low-count discrete phenotypes such as rosette lesions (74 DAS), leaf 6 lesions  
327 (74 DAS), leaf 8 lesions (76 DAS) and stem count (i.e. stem 1 plus the number of side stems).

328

329 **Leaf and final yield phenotypes of individual field-grown *B. napus* plants can**  
330 **be predicted to a considerable extent from their autumnal leaf 8**  
331 **transcriptome**

332 We built random forest (RF) and elastic net (enet) models to predict the phenotypes of  
333 individual plants from their autumnal leaf 8 transcriptome, using either all genes or only  
334 transcription factors (TFs) as potential features and using three different feature selection  
335 techniques (see Methods). For each combination of phenotype, model type (RF or enet),  
336 potential feature set (all genes or TFs) and feature selection technique, 9 repeat models were  
337 learned, each time using 10-fold cross-validation with different splits (see Methods), resulting  
338 per combination in a total of 90 test sets and 9 test set predictions per plant. The best model  
339 for a given phenotype and potential feature set was taken to be the one with the highest  
340 median test  $R^2$  value across all 90 test sets for continuous and high-count phenotypes (see  
341 Methods), or the highest median test accuracy for qualitative or low-count discrete  
342 phenotypes (**Table 2, Additional File 7: Table S6**).

343

344 Not surprisingly, leaf 8 phenotypes, which are most closely related in space and time to the  
345 material sampled for transcriptome profiling, are the most predictable. Except for the leaf 8  
346 chlorophyll content at sampling time (81 DAS), which features very poor prediction  
347 performance, the median test  $R^2$  scores for leaf 8 phenotypes range from 0.48 to 0.70 when  
348 using all genes as potential features. Other leaf-related phenotypes such as leaf 6 width (74  
349 DAS, median test  $R^2 = 0.38$ ), rosette area (42 DAS, median test  $R^2 = 0.23$ ) and leaf 6 length (74  
350 DAS, median test  $R^2 = 0.07$ ) are comparatively less predictable.

351

352 Surprisingly, many of the final seed, silique and shoot dry weight phenotypes are more  
353 predictable from the autumnal leaf 8 transcriptome than leaf 6 and rosette phenotypes, with  
354 seed weight on stem 1 rivaling the leaf 8 phenotypes in terms of median test  $R^2$  value (**Table**  
355 **2, Additional File 7: Table S6, Fig. 2, Additional File 2: Fig. S5**). All seed weight, seed and



356 silique count and shoot dry weight phenotypes have median test  $R^2$  values in the range 0.35  
357 - 0.51 for the 'all genes' models, which is in all cases higher than the 95<sup>th</sup> percentile of test  $R^2$   
358 values obtained from single train-test splits on 90 datasets in which the phenotype values  
359 were permuted (**Additional File 7: Table S6**). In other words, the model for the real data train-  
360 test split with median test  $R^2$  outperforms 95% of the models for comparable train-test splits  
361 on randomized data. Note that this serves only as an indication of model performance on real  
362 versus randomized data, not as a formal test assessing whether the median test  $R^2$  score on  
363 real data is significantly higher than expected at random. The latter would require the 9 times  
364 repeated 10-fold cross-validation setup used on the real data to be used on each of the  
365 permuted datasets as well (instead of the single train-test split per permutation used here),  
366 which is computationally prohibitive.

367

368 Interestingly, yield phenotypes measured for stem 1 are generally slightly more predictable  
369 than the corresponding phenotypes measured for the entire shoot, with median test  $R^2$  score  
370 differences between stem 1 and total shoot phenotypes in the range 0.02-0.09 for the 'all  
371 genes' models and 0.05-0.13 for the 'transcription factors' models. This suggests that gene  
372 expression levels in leaf 8 of the rosette may be more predictive for phenotypes of stem 1  
373 (i.e. the primary inflorescence stem and its cauline secondary inflorescences) than for  
374 phenotypes measured on the whole shoot (i.e. including the secondary inflorescence stems  
375 branching at ground level).

376

377 Root phenotypes, branching phenotypes, final plant height (278 DAS) and shoot growth  
378 phenotypes are generally poorly predictable (**Table 2**). Plant height and shoot growth  
379 phenotypes are likely poorly predictable because they show little variation across the field

380 **(Additional File 4: Table S3)**, increasing the risk that measurement error outweighs biological  
381 variation. Also taproot length and root system width may suffer from measurement errors.  
382 The total branch count and branch count stem 1 phenotypes on the other hand have a high  
383 CV and likely limited measurement error, suggesting that leaf 8 gene expression profiles may  
384 contain less information on these branching phenotypes than on leaf, seed, silique and dry  
385 weight phenotypes.

386

387 Most phenotypes calculated as ratios of other phenotypes are very poorly predictable, even  
388 if the constituent phenotypes have high prediction performance values. For instance, the  
389 median test  $R^2$  value for seeds per silique (total seed count divided by total silique count) is  
390 negative (-0.14), whereas both total seed count and total silique count have median test  $R^2$   
391 values  $\geq 0.38$ . In many cases however, the numerator and denominator phenotypes of a ratio  
392 are highly correlated, leading to a derived phenotype with a small range that may be  
393 dominated by noise propagated from measurement errors in the constituent phenotypes  
394 rather than biological variability. The number of siliques per branch on stem 1 and the entire  
395 shoot are notable exceptions with high CV values and reasonable prediction performance  
396 **(Table 2)**. The latter ratio phenotypes are highly correlated with the number of siliques on  
397 stem 1 (PCC = 0.92) and the entire shoot (PCC = 0.72), respectively, indicating that the number  
398 of siliques per branch is an important determinant of silique count, in addition to the number  
399 of branches (PCC between total branch count and total silique count = 0.87, PCC between  
400 branch count stem 1 and silique count stem 1 = 0.82).

401

402

403

	All genes				Transcription factors				Single gene			
	Feature sel.	Model type	Median test R2	Median pooled PCC	Feature sel.	Model type	Median test R2	Median pooled PCC	Top gene	Median test R2	Median pooled PCC	CV
<b>Continuous and high-count phenotypes</b>												
leaf 8 width (76 DAS)	median	enet	0.70 *	0.87	median	enet	0.64 *	0.84	BnaC04g39580D	0.67	0.83	2.32E-01
leaf 8 width (81 DAS)	median	enet	0.65 *	0.86	median	enet	0.65 *	0.84	BnaA02g18860D	0.62	0.83	2.32E-01
leaf 8 area (81 DAS)	median	enet	0.63 *	0.83	median	enet	0.58 *	0.81	BnaCmg33420D	0.60	0.81	3.70E-01
leaf 8 fresh weight (81 DAS)	median	enet	0.59 *	0.81	median	enet	0.53 *	0.78	BnaCmg33420D	0.58	0.79	3.88E-01
seed weight stem 1	spearman	enet	0.51 *	0.77	median	enet	0.53 *	0.78	BnaA05g29010D	0.42	0.67	4.51E-01
leaf 8 length (76 DAS)	spearman	rf	0.51 *	0.80	median	enet	0.53 *	0.79	BnaC07g39340D	0.58	0.80	2.21E-01
leaf 8 length (81 DAS)	spearman	rf	0.48 *	0.78	spearman	enet	0.51 *	0.78	BnaC01g17020D	0.52	0.81	2.11E-01
seed count stem 1	spearman	enet	0.47 *	0.74	median	enet	0.43 *	0.73	BnaA06g20870D	0.38	0.61	4.71E-01
silique count stem 1	median	enet	0.46 *	0.72	median	enet	0.45 *	0.72	BnaA05g29010D	0.37	0.66	4.37E-01
total seed count	spearman	enet	0.45 *	0.73	median	enet	0.38 *	0.71	BnaC03g60710D	0.39	0.61	4.78E-01
dry weight stem 1	spearman	enet	0.44 *	0.73	median	enet	0.39 *	0.70	BnaA05g29010D	0.40	0.70	4.83E-01
dry weight stem 1 (w/o seeds)	hsic-5000	enet	0.42 *	0.69	spearman	enet	0.35 *	0.64	BnaA05g29010D	0.39	0.69	5.19E-01
total seed weight	spearman	enet	0.42 *	0.74	median	enet	0.40 *	0.70	BnaA06g35450D	0.41	0.69	4.69E-01
total shoot dry weight	spearman	enet	0.41 *	0.71	median	enet	0.31 *	0.68	BnaA09g48720D	0.41	0.67	4.89E-01
leaf 6 width (74 DAS)	median	enet	0.38 *	0.68	hsic-5000	rf	0.07	0.51	BnaC03g15540D	0.35	0.61	1.91E-01
total silique count	median	enet	0.38 *	0.68	median	enet	0.36 *	0.68	BnaC04g21390D	0.40	0.63	4.56E-01
siliques per branch stem 1	hsic-5000	enet	0.37 *	0.67	hsic-5000	rf	0.14	0.51	BnaC04g21390D	0.25	0.60	3.50E-01
total shoot dry weight (w/o seeds)	spearman	enet	0.35 *	0.66	spearman	enet	0.29	0.62	BnaA06g05150D	0.40	0.69	5.13E-01
leaf count (74 DAS)	median	rf	0.24 *	0.66	median	rf	0.40 *	0.72	BnaA01g34700D	0.38	0.70	1.12E-01
rosette area (42 DAS)	median	enet	0.23 *	0.59	median	enet	0.36 *	0.68	BnaC05g30740D	0.33	0.64	3.00E-01
branch count stem 1	spearman	rf	0.19	0.56	median	enet	0.11 *	0.52	BnaA10g29560D	0.38	0.56	2.00E-01
siliques per branch	spearman	enet	0.16	0.49	spearman	enet	-0.01	0.36	BnaA08g09860D	0.12	0.53	3.53E-01
plant height (278 DAS)	median	enet	0.12 *	0.47	spearman	enet	0.16	0.51	BnaC07g25920D	0.34	0.63	6.81E-02
total branch count	median	rf	0.10	0.40	median	enet	0.17 *	0.57	BnaA09g48720D	0.26	0.59	3.42E-01
branch count stem 1/length stem 1	median	rf	0.07	0.39	median	rf	-0.06	0.23	BnaC05g15590D	0.22	0.55	1.58E-01
leaf 6 length (74 DAS)	median	enet	0.07 *	0.45	median	rf	0.06	0.43	BnaA09g04980D	0.31	0.64	1.96E-01
max shoot growth rate	median	rf	0.03	0.41	hsic-5000	rf	-0.02	0.34	BnaA10g21770D	0.15	0.50	6.75E-02
root system width	median	rf	0.01	0.36	median	rf	0.17 *	0.56	BnaC07g01150D	0.17	0.56	2.16E-01
time of max shoot growth	median	rf	-0.02	0.38	hsic-5000	rf	0.20	0.57	BnaA05g08250D	0.15	0.56	6.45E-03
taproot length	spearman	rf	-0.02	0.26	median	enet	-0.09	0.18	BnaA04g17830D	0.16	0.50	4.28E-01
branches per stem	median	rf	-0.09	0.26	hsic-5000	rf	-0.12	0.21	BnaC03g42190D	0.13	0.51	3.36E-01
leaf 8 chlorophyll content (81 DAS)	median	enet	-0.14	-0.35	median	enet	-0.12	-0.40	BnaA03g40350D	0.09	0.46	1.25E-01
seeds per silique	median	rf	-0.14	-0.18	median	enet	-0.17	-0.31	BnaC03g38990D	-0.09	0.22	1.95E-01
seeds per silique stem 1	median	enet	-0.15	-0.31	median	enet	-0.15	-0.08	BnaC01g44890D	-0.05	0.40	1.94E-01
seed weight stem 1/dry weight stem 1	median	enet	-0.18	-0.09	median	enet	-0.16	-0.39	BnaC09g50070D	-0.13	0.58	2.03E-01
total seed weight/shoot dry weight	median	enet	-0.19	-0.08	median	enet	-0.17	-0.12	BnaA02g15500D	-0.12	0.33	1.81E-01
end of shoot growth	hsic-5000	rf	-0.22	0.02	median	rf	-0.29	-0.17	BnaA05g09440D	0.04	0.45	1.02E-02
<b>Qualitative and low-count phenotypes</b>												
leaf 8 lesions (76 DAS)	hsic-5000	enet	0.67		hsic-5000	rf	0.67					
rosette lesions (74 DAS)	median	rf	0.50		spearman	enet	0.46					
leaf 6 lesions (74 DAS)	median	enet	0.33		spearman	enet	0.50					
stem count	median	enet	0.33		hsic-5000	enet	0.50					

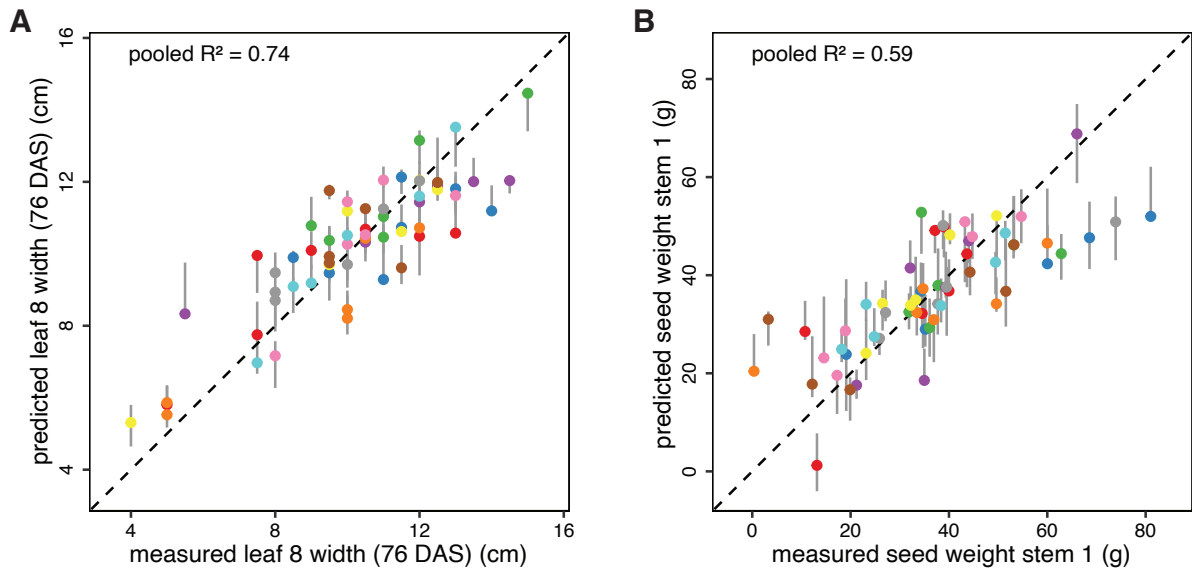
**Table 2** Best-performing multi-gene and single-gene models for each phenotype. Results are shown for models including all genes as potential features ('All genes' columns), models

404 including only TFs as potential features ('Transcription factors' columns) and models using a  
405 single gene as feature ('Single gene' columns). For the best all-gene and TF models for  
406 continuous or high-count discrete phenotypes, columns from left to right indicate the feature  
407 selection technique used (median = selection of features with median *rlog* gene expression >  
408 0, spearman = Spearman correlation, hsic-5000 = HSIC lasso, see Methods), the model type  
409 (enet = elastic net, rf = random forest), the median test  $R^2$  and the median pooled Pearson  
410 correlation coefficient (PCC, see Methods). Stars in the median test  $R^2$  column indicate that  
411 the median test  $R^2$  score on real data is higher than the 95<sup>th</sup> percentile of test  $R^2$  scores on  
412 permuted data (**Additional File 7: Table S6**). For qualitative and low-count phenotypes, the  
413 median test accuracy was used as a performance metric instead of the median test  $R^2$  (see  
414 Methods). Single-gene model columns include the best-performing gene and the  
415 corresponding median test  $R^2$  and median pooled PCC. All single-gene models are cross-  
416 validated lme models with spatial error structure. The CV column contains the coefficients of  
417 variation for the phenotypes.

418

419

420



421

422 **Fig. 2** Predictions versus observations for the best-scoring leaf and yield phenotypes. **A**

423 Predicted versus measured values for leaf 8 width (76 DAS), using the 'all genes' model with

424 the best median test  $R^2$  score (enet + median feature selection, **Table 2**). **B** Predicted versus

425 measured values for seed weight stem 1, using the 'all genes' model with the best median

426 test  $R^2$  score (enet + Spearman feature selection, **Table 2**). Vertical grey lines range from the

427 minimum to the maximum predicted value for a given plant across all model repeats, and

428 colored dots represent predictions for the repeat with the median pooled  $R^2$  score (i.e. the  $R^2$

429 score of the pooled test set predictions in the repeat concerned). Different marker colors

430 indicate the 10 different test sets in this repeat. Perfect predictions are located on the dashed

431 diagonal line in each panel.

432

433 To compare multi-gene models to single-gene models in terms of phenotype prediction

434 performance, we used the same repeated cross-validation setup as used for the multi-gene

435 models to calculate median test  $R^2$  scores and median pooled PCC values for single-gene

436 models (lme models with spatial structure, see previous section). Cross-validation scores were

437 calculated for each of the 100 genes most significantly associated with a given phenotype  
438 (lowest  $q$ -value for gene coefficient in lme model, **Additional File 7: Table S6**).

439

440 For leaf 8 phenotypes, the best multi-gene models generally have only slightly better median  
441 test  $R^2$  scores than the best single-gene models (**Table 2**). In other words, multi-gene models  
442 offer little benefit over single-gene models for quantitative prediction of leaf 8 phenotypes.

443 For leaf 8 length at 76 DAS and 81 DAS, the best single-gene models even outperform the  
444 multi-gene models. Single-gene models also outperform multi-gene models for several other  
445 phenotypes, sometimes with a wide margin, e.g. for plant height (278 DAS), branch count on  
446 stem 1, leaf count (74 DAS), leaf 6 length (74 DAS) and rosette area (42 DAS). This suggests  
447 that the multi-gene models are vulnerable to overfitting. In particular phenotypes with low  
448 single-gene model performance tend to exhibit a multi-gene model performance that is even  
449 lower, suggesting that the extent of multi-gene model overfitting is inversely correlated with  
450 the proportion of trait variance explained by single genes. An alternative explanation for the  
451 observation that the best single-gene models sometimes outperform the corresponding  
452 multi-gene models may be the ‘winner’s curse’ effect, also known as selection bias (32),  
453 whereby the apparently best-performing single-gene models may overestimate prediction  
454 performance.

455

456 Most of the phenotypes with comparatively high single-gene model performance scores  
457 however exhibit a modest increase of multi-gene model performance over the best single-  
458 gene model. Like most of the leaf 8-associated traits, total seed weight, total seed count and  
459 most of the shoot dry weight traits are modestly better predicted by multi-gene models than  
460 by single-gene models. Many of the seed and silique traits related to stem 1 on the other hand

461 (seed weight, seed count and silique count on stem 1, the number of siliques per branch on  
462 stem 1) are substantially better predicted by multi-gene models than by single-gene models.  
463 This indicates that several distinct gene expression patterns are likely relevant for quantitative  
464 prediction of stem 1 seed and silique traits.

465

466 For most ratio phenotypes, both the multi-gene and single-gene models have very poor  
467 prediction performance, in particular when the numerator and denominator phenotypes that  
468 make up the ratio are very highly correlated. In these cases, the denominator is essentially  
469 already a good predictor of the numerator. To assess whether any gene expression profiles  
470 contain additional information on the numerator given knowledge of the denominator, we  
471 used alternative single-gene models with a log link (see Methods) to predict the numerators  
472 of the seeds per silique ratio on stem 1 and the branches per stem ratio (seed count stem 1  
473 and total branch count, respectively) conditioned on their denominator (silique count stem 1  
474 and stem count, respectively). These models are not suited for making predictions in practice,  
475 given the need to know the denominator, but they may indicate whether prediction of the  
476 ratio based on gene expression is at all feasible and if so, which genes may be important. If  
477 no genes are found to be predictive for the numerator (and hence the ratio) conditioned on  
478 the denominator, then attempts to predict the ratio phenotype unconditionally are likely to  
479 be unsuccessful. For both seeds per silique stem 1 and branches per stem, the fitted  
480 coefficients and residuals look reasonable for the best predictor genes (**Additional File 2: Fig.**  
481 **S6, Additional File 2: Fig. S7**). The corresponding models succeed in suppressing a few of the  
482 more extreme residuals of the base model (without gene expression effect), without  
483 improving predictions for most other plants. However, no gene coefficients were found to be  
484 significantly different from zero for any phenotype after BH correction ( $q < 0.05$ ), neither in

485 models assuming constant error variance nor in models with heteroscedastic and/or spatially  
486 covarying error structures (see Methods). This indicates that the poor performance of the  
487 original multi-gene and single-gene models for these phenotypes is to be expected.

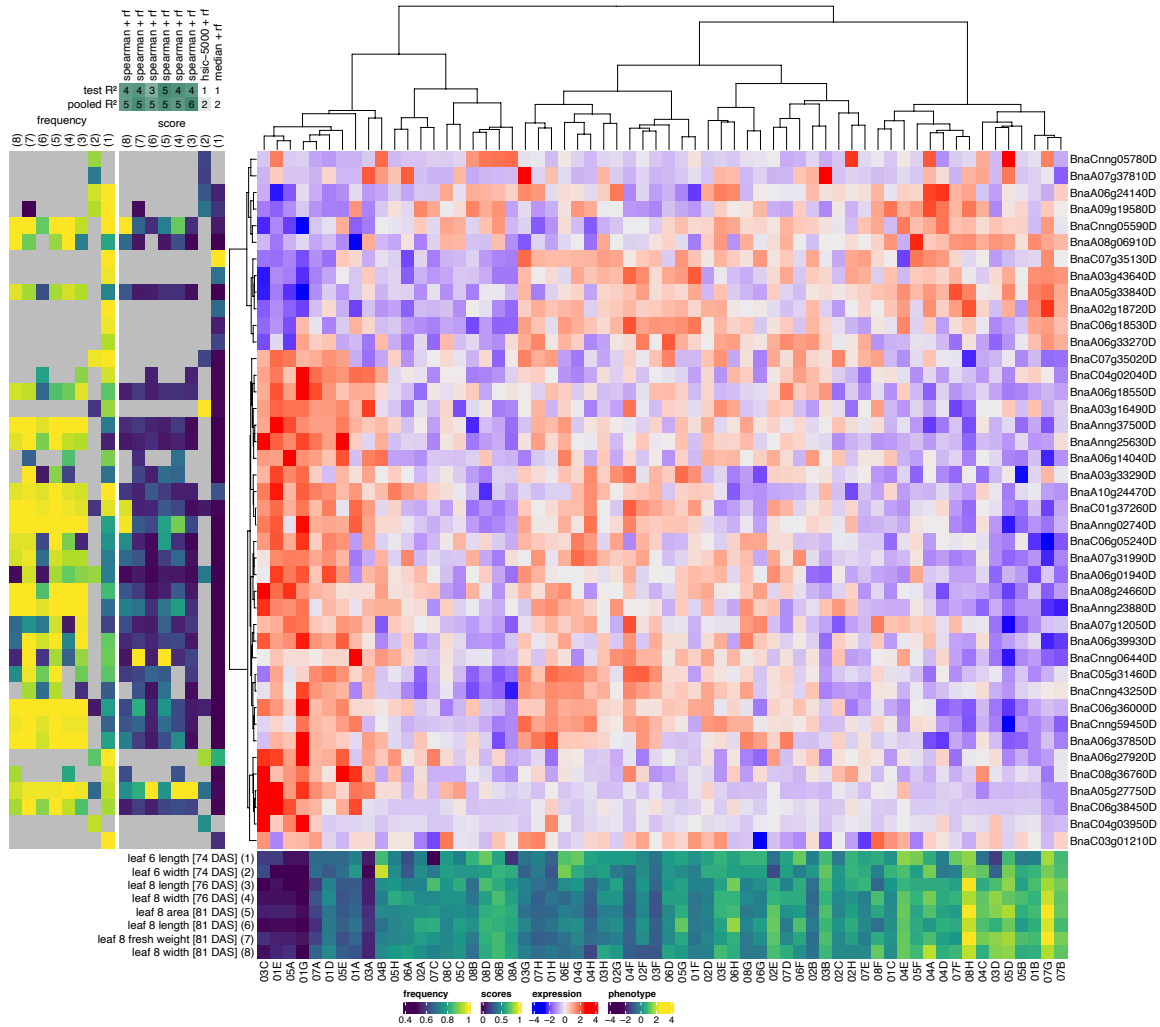
488

### 489 **Top predictors for leaf phenotypes**

490 The best multi-gene prediction performance scores were obtained for leaf 8 phenotypes. To  
491 assess whether the genes featuring most prominently in the multi-gene models for leaf  
492 phenotypes make biological sense, we focused on the top-10 predictor lists of the TF-based  
493 models for leaf 8 length and width (76 DAS and 81 DAS), fresh weight (81 DAS) and area (81  
494 DAS), and leaf 6 length and width (74 DAS) (**Additional File 7: Table S6**). As these leaf  
495 phenotypes are generally highly correlated (PCC between leaf 8 phenotype in the range  
496 [0.78,0.97], PCC between leaf 8 and leaf 6 phenotypes in range [0.45, 0.60]), many of the most  
497 important predictors (TFs) in the random forest and elastic net models are shared among  
498 phenotypes. We therefore grouped the top-10 predictor lists for the different phenotypes in  
499 two sets, one for the random forest (RF) models (**Fig. 3**) and one for the elastic net (enet)  
500 models (**Additional File 2: Fig. S8**). The rationale for looking at the TF models instead of the  
501 models built on all genes is that TFs are more likely than the average gene to have been  
502 functionally characterized to some extent, and are more likely to be causally involved in  
503 phenotype regulation (although it needs to be stressed that our analysis remains entirely  
504 correlational). Given the relative lack of experimentally determined gene functions in *B.*  
505 *napus*, most of the functional interpretation given below for *B. napus* genes is based on the  
506 experimentally determined functions of likely orthologs in *A. thaliana* (see Methods).

507





508

509 **Fig. 3** Top predictor genes in random forest models of leaf phenotypes. A clustered heatmap  
 510 of the z-scored gene expression profiles of the top genes for predicting leaf phenotypes is  
 511 shown centrally (blue-red color scale, Ward.D2 hierarchical clustering). The leaf phenotypes  
 512 concerned and their z-scored profiles across plants are shown at the bottom (dark blue-yellow  
 513 heatmap with plant identifiers at the bottom). For each of these phenotypes, the top-10 most  
 514 important genes (highest median gini importance across all 90 cross-validation splits) of the  
 515 RF model with the highest median test  $R^2$  score are included on the figure (gene identifiers  
 516 are shown at right). The mostly dark blue score panel to the left of the expression heatmap  
 517 shows the median gini importance scores of the selected genes in each of the selected  
 518 phenotype models, normalized to the maximum importance score per model to make the

519 color scales of the different models (columns) comparable. The mostly yellow frequency panel  
520 to the left of the score panel shows the frequencies at which genes were selected as features  
521 across all 90 cross-validation splits of a given model. Grey squares in the score and frequency  
522 panels indicate that a given gene was not selected as a feature in a given model. The  
523 phenotypes in the score and frequency panels are identified by numbers (1-8) on top of the  
524 panels, corresponding to the numbers associated with the phenotypes in the bottom  
525 phenotype panel. On top of the score panel, the feature selection techniques used in the best-  
526 scoring RF models for each phenotype are shown (median = selection of features with median  
527 *rlog* gene expression > 0, spearman = Spearman correlation, hsic-5000 = HSIC lasso, see  
528 Methods), as well as the corresponding test and pooled  $R^2$  scores rounded to the nearest 0.1  
529 and then multiplied by ten (e.g. a test  $R^2$  score of 0.38 would be denoted as 4).

530

531 Many of the top TF predictors for leaf phenotypes have *A. thaliana* orthologs with known  
532 functions in leaf development. One TF with high importance scores in both the RF and enet  
533 models is *BnaCnng05590D*, a putative ortholog of the homeodomain leucine zipper class I  
534 (*HD-ZIP I*) gene *ARABIDOPSIS THALIANA HOMEBOX 1* (*AtHB1/AT3G01470*). Both the RF and  
535 enet top predictor sets additionally contain *BnaA05g33840D*, another putative ortholog of  
536 *AtHB1*. Ectopic *AtHB1* overexpression in tobacco seedlings was previously shown to lead to  
537 de-etiolated phenotypes in the dark, including true leaf development (33). Mutation of an  
538 upstream open reading frame in the *AtHB1* 5' untranslated region that normally represses  
539 *AtHB1* translation was shown to lead to smaller, more serrated leaves, smaller rosettes, a  
540 delay of the vegetative-to-reproductive phase transition and siliques containing fewer seeds  
541 in *A. thaliana* (34). Similarly, *AtHB1* overexpression in a silencing-deficient *rdr6-12* mutant  
542 background resulted in plants with shorter and more serrated leaves (35).

543

544 The enet top predictor list for leaf phenotypes also contains another *HD-ZIP I* gene,  
545 *BnaC02g43700D*, which is putatively orthologous to *AtHB5* (*AT5G65310*) or *AtHB16*  
546 (*AT4G40060*). Similar to *AtHB1*, overexpression of *AtHB16* leads to smaller, more serrated  
547 leaves exhibiting reduced cell expansion, smaller rosettes and siliques containing fewer seeds  
548 (36). Additionally, *AtHB16* overexpression was reported to reduce the flowering time  
549 sensitivity to differences in photoperiod in *A. thaliana* (36).

550

551 Next to *HD-ZIP I* genes, the RF and enet top predictor lists contain several other *HD-ZIP* genes.  
552 *BnaA06g18550D* in the RF top predictor list is putatively orthologous to the *A. thaliana* gene  
553 *REVOLUTA* (*AtREV/AT5G60690*), which encodes a *HD-ZIP III* transcription factor known to  
554 regulate postembryonic meristem initiation (37) and several polarity-associated growth  
555 processes in *A. thaliana*, including abaxial-adaxial patterning in leaves (38). Loss-of-function  
556 *atrev-1* mutant plants were shown to exhibit overgrowth and deformation of rosette and  
557 cauline leaves after bolting (39). The RF top predictor list also contains two additional *HD-ZIP*  
558 *III* gene family members, *BnaC06g05240D* and *BnaA06g01940D*, that are putatively  
559 orthologous to *AtHB8* (*AT4G32880*) or *AtHB15* (*AT1G52150*). *AtHB8* and *AtHB15* are thought  
560 to have effects on postembryonic meristem initiation that are antagonistic to the effects of  
561 *AtREV* (40). On the other hand, gain-of-function mutations in *AtHB15*, like gain-of-function  
562 mutations in *AtREV*, have been shown to result in adaxialized leaves (41). Both *AtHB8* and  
563 *AtHB15* are thought to function prominently in vascular development, possibly  
564 antagonistically (41-44). Furthermore, the enet top predictor list includes *BnaC03g02700D*, a  
565 *HD-ZIP II* gene putatively orthologous to *AtHAT3* (*AT3G60390*), *AtHAT14* (*AT5G06710*),

566 *AtHB17* (AT2G01430) or *AtHB18* (AT1G70920). *AtHAT3* is known to be involved in leaf  
567 abaxial/adaxial patterning (45), and to be regulated by *AtREV* (46).

568

569 Both the RF and enet top predictor lists prominently feature putative orthologs of *A. thaliana*  
570 *WUSCHEL RELATED HOMEBOX* (*AtWOX*) genes: *BnaA05g27750D* (RF and enet) and  
571 *BnaC05g41930D* (enet). Both genes are putatively orthologous to *AtWOX5* (AT3G11260) or  
572 *AtWOX7* (AT5G05770). Next to roles in root development, *AtWOX5* was reported to act  
573 redundantly with *AtWOX1* and *AtWOX3* to control leaf shape by promoting lateral leaf growth  
574 (47). *AtWOX1*, 2, 3 and 5 were also shown to regulate the expression of *REVOLUTA* (*AtREV*)  
575 and other *HD-ZIP III* genes in the shoot apical meristem (48), and *AtWOX1* and *AtWOX3* are  
576 thought to regulate *HD-ZIP III* genes in lateral leaf regions, thereby contributing to the  
577 maintenance of adaxial/abaxial patterning at the margin of growing leaves (49).

578

579 Not all transcription factors in the RF and enet models are equally important for all leaf  
580 phenotypes. *BnaCnng06440D* (*AtMYB60/AT1G08810*) for instance has higher RF and (to a  
581 lesser extent) enet importance scores for leaf 8 area (81 DAS) and leaf 8 fresh weight (81 DAS)  
582 than for other leaf phenotypes. Its likely *A. thaliana* ortholog *AtMYB60* is involved in  
583 regulating stomatal opening, and its expression is downregulated under drought (50).  
584 *atmyb60-1* null mutant plants exhibit a constitutive reduction of stomatal opening and  
585 decreased transpirational water loss under drought (50). A second TF in the RF models with  
586 higher importance for leaf 8 area (81 DAS) and leaf 8 fresh weight (81 DAS) than for other leaf  
587 phenotypes is *BnaC06g36000D* (*AtHB33/AtZHD5/AT1G75240*). Its likely ortholog *AtHB33*  
588 codes for a zinc-finger homeodomain TF downregulated in response to abscisic acid (ABA),  
589 which e.g. induces stomatal closure (51). Constitutive *AtHB33* overexpression in *A. thaliana*

590 resulted in accelerated growth, larger leaves and larger epidermal cells (52). A third TF in the  
591 RF models with mildly higher importance for leaf 8 area (81 DAS) and leaf 8 fresh weight (81  
592 DAS) is *BnaCnn59450D*. Overexpression of its putative orthologs *AtSHN2* (*AT5G11190*) and  
593 *AtSHN3* (*AT5G25390*) in *A. thaliana* resulted in folded and twisted leaves, shiny green leaf  
594 surfaces with increased levels and altered composition of cuticular wax, increased cuticular  
595 permeability, larger pavement cells, reduced trichome number and stomatal density, and  
596 increased drought tolerance (53).

597

598 Both the RF and enet top-10 lists feature several orthologs of *A. thaliana* NUCLEAR FACTOR  
599 Y, SUBUNIT A (*AtNF-YA*) genes (putative *A. thaliana* orthologs in parentheses):  
600 *BnaAnng02740D* (*AtNF-YA2/10*, *AT3G05690/AT5G06510*, RF), *BnaA10g24470D* (*AtNF-*  
601 *YA2/10*, *AT3G05690/AT5G06510*, RF and enet), *BnaC06g33980D* (*AtNF-YA3/8*,  
602 *AT1G72830/AT1G17590*, enet) and *BnaC01g37260D* (*AtNF-YA5/6*, *AT1G54160/AT3G14020*,  
603 RF). All four genes are negatively correlated with leaf phenotypes in the field expression  
604 dataset. NF-Y transcription factor complexes are heterotrimers, consisting of A, B and C  
605 subunits, that function in various developmental programs and abiotic stress responses in  
606 plants (54). *AtNF-YA2* and *AtNF-YA10* were previously found to regulate leaf size in *A.*  
607 *thaliana*, with their overexpression promoting cell expansion (55). *AtNF-YA5* was found to  
608 promote drought resistance, with *atnf-ya5* knockout plants and *AtNF-YA5*-overexpressing  
609 plants displaying increased and reduced leaf water loss, respectively, relative to wild-type  
610 plants (56). *AtNF-YA8* was recently found to negatively regulate the juvenile-to-adult  
611 (vegetative) phase change by activating the transcription of *AtMIR156* genes (57).  
612 Overexpression of *AtNF-YA8* resulted in a delay of the juvenile-to-adult transition and thereby  
613 reduced leaf sizes and altered leaf shapes (57).

614

615 Interestingly, several of the top-TFs recovered in the multi-gene models for leaf phenotypes  
616 are linked to the regulation of flowering. Plant NF-Y complexes for instance are known to also  
617 function in the regulation of flowering time (54). Overexpression of the aforementioned *AtNF-*  
618 *YA8* gene was found to delay flowering time (57), and similar observations were made for  
619 other *AtNF-YA* genes such as *AtNF-YA1* (*AtHAP2A*, *AT5G12840*), *AtNF-YA2*, *AtNF-YA3*, *AtNF-*  
620 *YA4* (*At2g34720*), *AtNF-YA7* (*At1g30500*) and *AtNF-YA10* (58, 59). It has been suggested that  
621 the photoperiodic flowering regulator CONSTANS (*AtCO*) may compete with NF-YA subunits  
622 in the NF-Y complex to form an alternative complex activating *FLOWERING LOCUS T* (*FT*)  
623 expression in *A. thaliana*, thereby promoting flowering (58). Additionally, *AtNF-YA2* has been  
624 suggested to function as a negative regulator of flowering in an alternative, stress-mediated  
625 flowering pathway (60). On the other hand, *AtNF-YA2* was recently suggested to positively  
626 regulate flowering by directly influencing *AtFT* expression (61).

627

628 The *A. thaliana* orthologs of several of the aforementioned *HD-ZIP* genes  
629 (*BnaCnng05590D/AtHB1*, *BnaA05g33840D/AtHB1*, *BnaC02g43700D/AtHB16*,  
630 *BnaC06g05240D/AtHB15*, *BnaA06g01940D/AtHB15*) have also been linked to regulation of  
631 the juvenile-to-adult and/or vegetative-to-reproductive phase changes (34, 36, 41).  
632 Furthermore, both the enet and RF predictor lists contain *BnaA06g39930D*, a putative  
633 ortholog of *EARLY FLOWERING MYB PROTEIN* (*AtEFM/AT2G03500*) in *A. thaliana*. *AtEFM* is  
634 known to directly repress the expression of *FLOWERING LOCUS T* (*AtFT*, *AT1G65480*) in the  
635 leaf vasculature, and is thought to mediate the effects of temperature and light cues on the  
636 timing of the floral transition (62). The RF predictor list additionally contains *BnaC05g31460D*,  
637 a putative ortholog of *AtJMJD5* (*AtJMJD5*, *AT3G20810*), the protein product of which interacts

638 with AtEFM to repress AtFT (62). The RF and enet top predictor lists also contain  
639 *BnaA07g12050D*, a putative ortholog of the floral homeotic gene *APETALA2*  
640 (*AtAP2/AT4G36920*) or the related *euAPETALA2* gene *AtTOE3 (AT5G67180)*. Both AtAP2 and  
641 AtTOE3 are known to repress *AGAMOUS (AtAG)* expression during floral patterning (63).

642

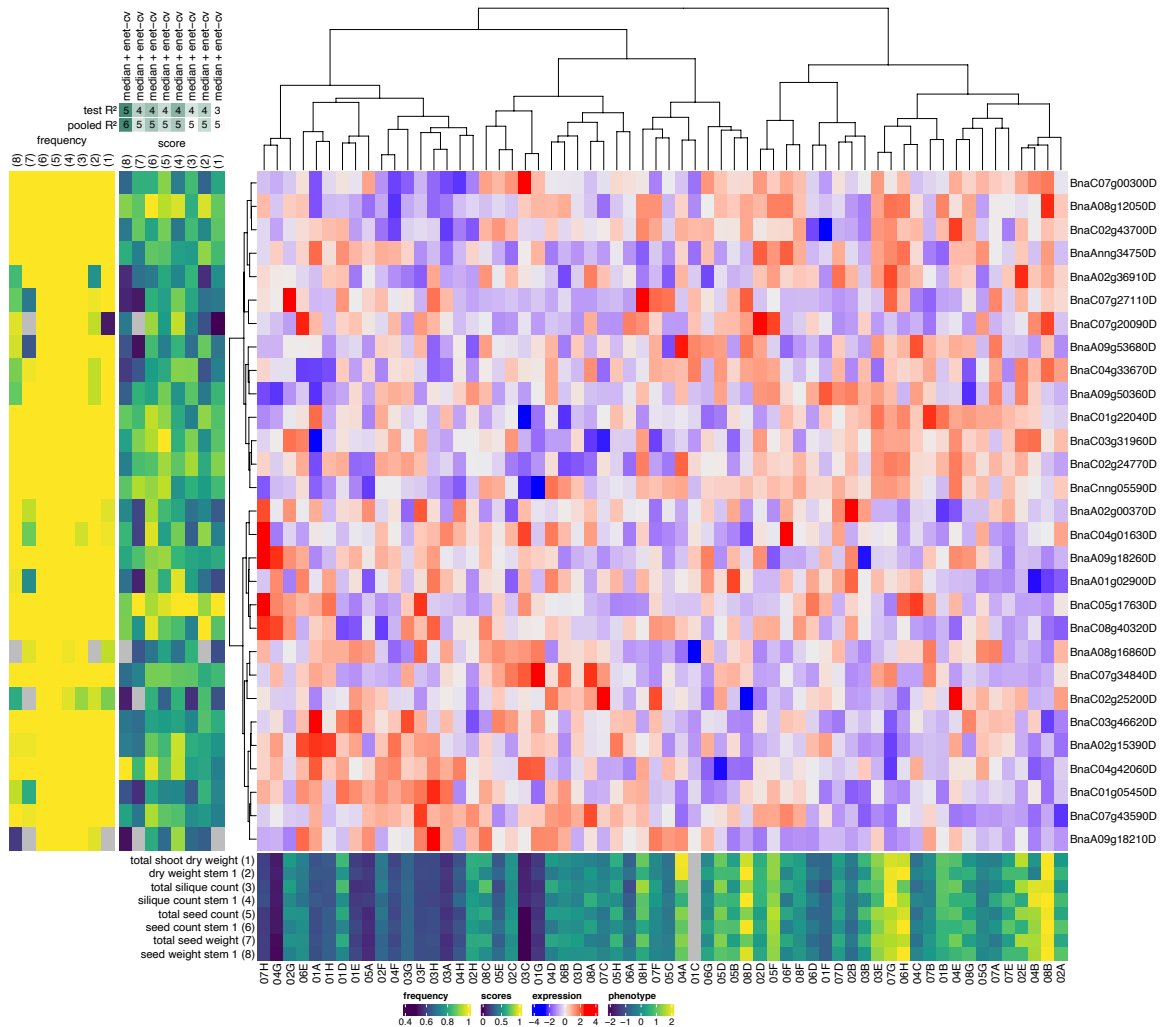
643 In summary, 15/42 and 11/35 transcription factors in the RF and enet lists of top leaf  
644 phenotype predictors, respectively, have putative *A. thaliana* orthologs linked to leaf  
645 development and patterning, the juvenile-to-adult phase change, the floral transition or  
646 drought response.

647

#### 648 **Top predictors for seed, silique and shoot dry weight phenotypes**

649 Next to leaf 8 phenotypes, also the seed, silique and shoot dry weight phenotypes (further  
650 referred to as ‘yield’ phenotypes) of the individual plants at harvest (late spring) could be  
651 predicted to a considerable extent from autumnal leaf 8 transcriptome data (see above).  
652 Similar to the leaf phenotypes, the yield phenotypes are highly correlated (PCC range [0.84-  
653 0.99]) and hence have a lot of high-scoring RF and enet predictors in common (**Additional File**  
654 **7: Table S6, Fig. 4, Additional File 2: Fig. S9**). Furthermore, these phenotypes are also  
655 significantly correlated with leaf phenotypes (PCC range [0.47, 0.74]), leading to a substantial  
656 overlap between the top-10 predictor lists of yield and leaf phenotypes.

657



658

659 **Fig. 4** Top predictor genes in elastic net models of yield phenotypes. A clustered heatmap of  
 660 the z-scored gene expression profiles of the top genes for predicting yield phenotypes is  
 661 shown centrally (blue-red color scale, Ward.D2 hierarchical clustering). The yield phenotypes  
 662 concerned and their z-scored profiles across plants are shown at the bottom (dark blue-yellow  
 663 heatmap with plant identifiers at the bottom). For each of these phenotypes, the top-10 most  
 664 important genes (highest median elastic net coefficients across all 90 cross-validation splits)  
 665 of the enet model with the highest median test  $R^2$  score are included on the figure (gene  
 666 identifiers are shown at right). The mostly green-blue score panel to the left of the expression  
 667 heatmap shows the median elastic net coefficients of the selected genes in each of the  
 668 selected phenotype models, normalized to the maximum coefficient per model to make the



669 color scales of the different models (columns) comparable. The mostly yellow frequency panel  
670 to the left of the score panel shows the frequencies at which genes were selected as features  
671 across all 90 cross-validation splits of a given model. Grey squares in the score and frequency  
672 panels indicate that a given gene was not selected as a feature in a given model. The  
673 phenotypes in the score and frequency panels are identified by numbers (1-8) on top of the  
674 panels, corresponding to the numbers associated with the phenotypes in the bottom  
675 phenotype panel. On top of the score panel, the feature selection techniques used in the best-  
676 scoring enet models for each phenotype are shown (median = selection of features with  
677 median *rlog* gene expression > 0, spearman = Spearman correlation, hsic-5000 = HSIC lasso,  
678 see Methods), as well as the corresponding test and pooled  $R^2$  scores rounded to the nearest  
679 0.1 and then multiplied by ten (e.g. a test  $R^2$  score of 0.38 would be denoted as 4).

680

681 In particular, virtually all TF genes in the leaf top-10 predictor lists with links to the juvenile-  
682 to-adult or vegetative-to-reproductive phase changes and flowering also feature prominently  
683 in the RF or enet top-10 predictor lists for yield phenotypes, including *BnaCnng05590D*  
684 (*AtHB1/AT3G01470*), *BnaC02g43700D* (*AtHB5/AT5G65310* or *AtHB16/AT4G40060*),  
685 *BnaA07g12050D* (*AtAP2/AT4G36920*), *BnaAnng02740D* (*AtNF-YA2/10*,  
686 *AT3G05690/AT5G06510*), *BnaA10g24470D* (*AtNF-YA2/10*, *AT3G05690/AT5G06510*),  
687 *BnaC06g33980D* (*AtNF-YA3/8*, *AT1G72830/AT1G17590*) and *BnaC01g37260D* (*AtNF-YA5/6*,  
688 *AT1G54160/AT3G14020*). Furthermore, like the top predictor lists for leaf phenotypes, the  
689 enet top predictor list for yield phenotypes contains a putative ortholog of the *A. thaliana*  
690 gene *EARLY FLOWERING MYB PROTEIN* (*AtEFM/AT2G03500*), but a different one  
691 (*BnaAnng34750D*).

692

693 Furthermore, many of the top predictor TF genes for yield phenotypes that are absent from  
694 the top-10 predictor lists for leaf phenotypes also have *A. thaliana* orthologs involved in  
695 processes related to the floral transition and flowering. In the combined set of top-10 enet  
696 predictors for shoot dry weight, seed and silique phenotypes (**Fig. 4**, n=29), five genes code  
697 for AGAMOUS-LIKE MADS-box transcription factors (best candidate *A. thaliana* orthologs and  
698 associated AGI codes are given between parentheses): *BnaC05g17630D*  
699 (*AtAGL104/AT1G22130*), *BnaA02g15390D* (*AtAGL12/AT1G71692*), *BnaA02g00370D*  
700 (*BnFLC.A2*, *AtFLC/AT5G10140*), *BnaA01g02900D* (*AtAGL16/AT3G57230*), and  
701 *BnaA09g53680D* (*AtAGL30/AT2G03060*). *BnFLC.A2* is orthologous to *A. thaliana* FLOWERING  
702 LOCUS C (*AtFLC*), a key repressor of the floral transition (64, 65). Two AGAMOUS-LIKE genes  
703 feature in the combined set of top-10 RF predictors for yield phenotypes (**Additional File 2:**  
704 **Fig. S9**, n=21): *BnaA02g15390D* (*AtAGL12/AT1G71692*) and *BnaA09g05500D*  
705 (*AtAGL8/AtFUL/FRUITFULL/AT5G60910*). *AtFUL* is thought to regulate the floral transition  
706 downstream of *AtFT* in the shoot apical meristem, partially redundantly with *AtSOC1*  
707 (*AtAGL20, AT2G45660*) (66, 67). *AtAP2* (APETALA2) and *AtFUL* are thought to form a bistable  
708 switch mechanism through mutual repression that regulates early stages of the floral  
709 transition at the shoot apical meristem (68). Negative regulation of *AtAP2* and several *AP2-*  
710 *LIKE* genes by *AtFUL* was also found to contribute to meristem arrest at the end of flowering  
711 (69). *ful* mutants were found to exhibit a delayed floral transition (68) and increased flower  
712 production, but decreased seed set (69).

713

714 The enet top predictor list also features *BnaA09g18260D*, a *HD-ZIP II* gene putatively  
715 orthologous to *JAIBA* (*AtJAB/AtHAT1/AT4G17460*) or *AtHAT2* (*AT5G47370*). *AtJAB* was shown  
716 to be involved in male and female reproductive development and floral meristem

717 determination in *A. thaliana*, and *jab* loss-of-function mutants exhibit an increased number  
718 of floral buds per inflorescence but a reduced number of seeds per silique (70), not unlike *ful*  
719 mutants. The enet top predictor list also contains two *HD-ZIP IV* genes, *BnaA09g50360D*  
720 (*AtHDG2/ AT1G05230*) and *BnaC03g31960D* (*AtANL2/ AT4G00730*). A combination of *hdg2*  
721 and *pdf2* null mutant alleles in *A. thaliana* was shown earlier to produce flowers with sepaloid  
722 petals and carpeloid stamens (71). *BnaC03g31960D* also features in the enet top predictor list  
723 for leaf phenotypes, but less prominently (**Additional File 2: Fig. S8**).

724

725 Furthermore, the gene *BnaA08g12050D* is ranked highly in both the enet and RF top predictor  
726 lists. The best candidate ortholog of this gene in *A. thaliana* is *AtMYB3R1* (*AT4G32730*), coding  
727 for a regulator of cell proliferation that acts in a module with *AtTSO1* to balance cell  
728 proliferation with differentiation in developing roots and shoots (72). Loss-of-function  
729 mutations in *AtMYB3R1* suppress all phenotypes of the *tso1-1* mutant, among others a lack  
730 of floral organ differentiation (72). *BnaA08g12050D* also features as a predictor for leaf 6  
731 length (74 DAS) and leaf 8 area (81 DAS) in **Additional File 2: Fig. S8**.

732

733 *BnaC07g27110D* and *BnaC01g22040D* in the enet predictor list are putative orthologs of  
734 *AtGATA16* (*AT5G49300*) and *AtGATA17*(*AT3G16870*) or *AtGATA17L*(*AT4G16141*),  
735 respectively. Evidence suggests these and other LLM-domain B-GATA transcription factors  
736 are involved (at least partially redundantly) in the regulation of flowering time, silique length,  
737 seed set and other developmental processes (73). The enet top predictor list also contains  
738 *BnaC04g33670D* and *BnaA08g16860D*, *BZIP* genes putatively orthologous to the *A. thaliana*  
739 genes *DRINK ME* (*AtDKM/AtBZIP30/AT2G21230*) and *DRINK ME-LIKE*  
740 (*AtDKML/AtBZIP29/AT4G38900*), respectively. *AtDKM* and *AtDKML* are negative regulators of

741 reproductive development and growth (74). *AtDKM* overexpression in *A. thaliana* results in  
742 smaller plants with fewer floral buds and shorter siliques, while a *dkm* mutant show the  
743 opposite phenotype (74). *dkml* mutant plants also exhibited increased silique length but  
744 slightly fewer floral buds than wild-type plants (74). *AtDKM* was shown to interact *in planta*  
745 with several regulators of meristem development, including *WUSCHEL* (*AtWU*), *HECATE1*  
746 (*AtHEC1*), the aforementioned *JAIBA* and *NGATHA1* (*AtNGA1*) (74). Interestingly, the RF top  
747 predictor list contains a putative ortholog of *NGATHA1* (*AtNGA1/AT2G46870*) or *NGATHA2*  
748 (*AtNGA2/AT3G61970*), namely *BnaA09g39540D*. *AtNGA1* and *AtNGA2* are known to be  
749 involved in gynoecium development and were recently shown to also have a function in  
750 regulating shoot apical meristem development (75). Another likely regulator of meristem  
751 development, *BnaC07g43590D*, is found in the enet predictor list. *BnaC07g43590D* is most  
752 likely an ortholog of *ARABIDOPSIS RESPONSE REGULATOR 10* (*AtARR10/AT4G31920*) or *12*  
753 (*AtARR12/AT2G25180*), both known to directly activate the expression of *WUSCHEL* and to  
754 play a role in shoot apical meristem regeneration and maintenance (76).

755

756 In summary, 17/29 and 11/21 TF genes in the enet and RF lists of top yield predictors,  
757 respectively, have putative *A. thaliana* orthologs linked to the juvenile-to-adult phase change,  
758 the floral transition, flowering or regulation of meristem development.

759

## 760 **Predicting final yield phenotypes from early growth phenotypes**

761 As a baseline to assess the prediction performance of the molecular models, we trained  
762 models predicting plant phenotypes in spring (mostly phenotypes at harvest) from single or  
763 multiple autumnal leaf and rosette phenotypes. For these single- and multi-phenotype

764 models, the same modeling approaches were used as for the single- or multi-gene models,  
765 respectively (see Methods).

766

767 Interestingly, many of the mature plant phenotypes can be predicted to a considerable extent  
768 from phenotypes measured earlier in the growing season (**Table 3, Additional File 8: Table**  
769 **S7**). In particular the models for phenotypes measured on the entire shoot (total seed, silique  
770 and branch count, total seed weight, total shoot dry weight) perform surprisingly well. For  
771 most of these phenotypes, the performance of the early-phenotype models is only slightly  
772 less than that of the best single-gene or multi-gene model, and the early-phenotype models  
773 for total seed weight and total branch count even outperform the molecular models (in the  
774 case of total branch count even substantially so). Also for branching phenotypes related to  
775 stem 1 (branch count stem 1, branch count stem 1/length stem 1), the best early-phenotype  
776 models feature high prediction performance scores. For other stem 1 phenotypes however  
777 (seed weight, seed count, silique count and siliques per branch on stem 1, stem 1 dry weight  
778 with and without seeds), the molecular models clearly outperform the early-phenotype  
779 models.

780

781 Most multi-phenotype models with appreciable prediction performance (median test  $R^2 >$   
782  $0.10$ ), both for whole-shoot and stem 1 phenotypes, feature leaf 8 area (81 DAS) as the top  
783 predictor (**Table 3**). Leaf 8 area (81 DAS) is generally also the most predictive early phenotype  
784 in the corresponding sets of single-phenotype models. The multi-phenotype models with the  
785 best prediction performance scores, i.e. those for whole-shoot phenotypes and stem 1  
786 branching phenotypes (but not the other stem 1 phenotypes), generally also feature rosette  
787 area (42 DAS) as a predictor of some importance (**Additional File 8: Table S7**). For total branch

788 count, branch count stem 1 and branch count stem 1/length stem 1, rosette area (42 DAS) is  
 789 even the top predictor in either the RF or enet model, or both (**Additional File 8: Table S7**).  
 790 Rosette area (42 DAS) itself is only moderately predictable from the leaf 8 molecular data,  
 791 which may explain why multi-phenotype models are better at predicting these branching  
 792 phenotypes than multi-gene models.  
 793

Mature plant phenotypes	All early phenotypes			Single early phenotypes		
	Model type	Median test R <sup>2</sup>	Median pooled PCC	Top phenotype	Median test R <sup>2</sup>	Median pooled PCC
seed weight stem 1	enet	0.33	0.59	leaf 8 area (81 DAS)	0.32	0.63
seed count stem 1	rf	0.25	0.57	leaf 8 area (81 DAS)	0.26	0.61
siliques count stem 1	enet	0.24	0.55	leaf 6 width (74 DAS)	0.22	0.53
total seed count	rf	0.40	0.66	leaf 8 area (81 DAS)	0.44	0.71
dry weight stem 1	enet	0.26	0.57	leaf 8 area (81 DAS)	0.35	0.62
dry weight stem 1 (w/o seeds)	enet	0.18	0.54	leaf 8 area (81 DAS)	0.30	0.59
total seed weight	enet	0.45	0.68	leaf 8 area (81 DAS)	0.46	0.72
total shoot dry weight	enet	0.38	0.67	leaf 8 area (81 DAS)	0.44	0.71
total siliques count	rf	0.36	0.63	leaf 8 area (81 DAS)	0.41	0.70
siliques per branch stem 1	enet	0.14	0.47	leaf 8 area (81 DAS)	0.07	0.48
total shoot dry weight (w/o seeds)	enet	0.29	0.63	leaf 8 area (81 DAS)	0.37	0.68
branch count stem 1	enet	0.35	0.64	leaf 8 area (81 DAS)	0.34	0.65
siliques per branch	enet	-0.04	0.32	leaf 6 width (74 DAS)	-0.07	0.36
plant height	enet	0.23	0.58	leaf 8 length (81 DAS)	0.25	0.61
total branch count	rf	0.40	0.69	rosette area (42 DAS)	0.38	0.65
branch count stem 1/length stem 1	rf	0.33	0.63	leaf 8 area (81 DAS)	0.22	0.56
max shoot growth rate	enet	0.04	0.40	leaf 8 length (81 DAS)	0.04	0.41
root system width	rf	0.04	0.42	leaf 8 length (81 DAS)	0.05	0.39
time of max shoot growth	enet	-0.01	0.53	leaf 8 width (81 DAS)	0.08	0.53
taproot length	rf	-0.02	0.33	leaf 8 width (81 DAS)	0.01	0.33
branches per stem	enet	-0.14	-0.22	leaf 8 lesions (76 DAS)	-0.15	0.14
seeds per siliques	enet	-0.17	-0.18	leaf 8 length (81 DAS)	-0.07	0.18
seeds per siliques stem 1	enet	-0.15	-0.06	leaf 8 length (81 DAS)	-0.04	0.23
seed weight stem 1/dry weight stem 1	enet	-0.15	-0.40	leaf 8 lesions (76 DAS)	-0.19	-0.19
total seed weight/shoot dry weight	enet	-0.16	-0.39	leaf 8 lesions (76 DAS)	-0.18	-0.06
794 end of shoot growth	enet	-0.15	0.20	leaf 8 width (81 DAS)	-0.12	0.30

795 **Table 3 Best-performing multi-phenotype and single-phenotype models for mature plant**  
 796 **phenotypes.** Results are shown for models including all early phenotypes as potential  
 797 features (multi-phenotype models) and models using a single early phenotype as feature  
 798 (single-phenotype models). For the best multi-phenotype models, columns from left to right  
 799 indicate the model type used (enet = elastic net, rf = random forest), the median test  $R^2$  and  
 800 the median pooled PCC (see Methods). Single-phenotype columns include the best-

801 performing early phenotype ('Top phenotype' column) and the corresponding median test  $R^2$   
802 and median pooled PCC. All single-phenotype models are cross-validated lme models with  
803 spatial error structure.

804

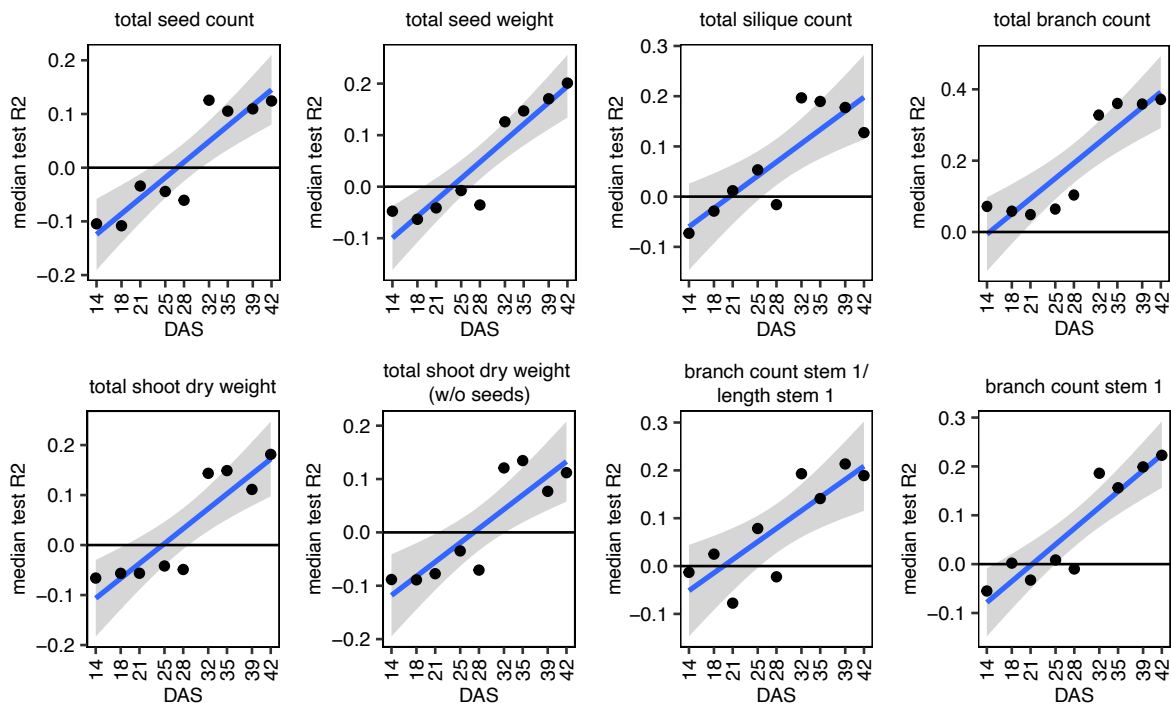
805 Our results indicate that the leaf 8 molecular data offer little benefit over early-phenotype  
806 measurements for quantitative prediction of mature phenotypes measured on the entire  
807 plant. On the other hand, the leaf 8 molecular data yields substantially better models than  
808 the early-phenotype data for most mature stem 1 phenotypes. Often, the multi-gene models  
809 for stem 1 phenotypes are also slightly better than the multi-gene models for the  
810 corresponding whole-plant phenotypes (see previous section). This suggests that the  
811 molecular makeup of the 8th rosette leaf at the time of sampling contained more information  
812 on the development of the primary flowering stem and its cauline secondary inflorescences  
813 than on the development of side stems at ground level. Early phenotypes on the other hand  
814 may contain more information on whole-plant yield phenotypes than on phenotypes  
815 specifically related to stem 1.

816

817 Given that even the earliest of the autumnal phenotypes considered thus far, the rosette area  
818 at 42 DAS, still has some predictive power for several yield phenotypes (median test  $R^2 > 0.10$   
819 for total branch count, seed count, seed weight and silique count, total dry weight with and  
820 without seeds, branch count stem 1 and branch count stem 1/length stem 1), we assessed  
821 whether earlier rosette areas (v2, see Methods) are also predictive for these phenotypes (**Fig.**  
822 **5, Additional File 9: Table S8**). Median test  $R^2$  scores were found to decrease when using  
823 earlier rosette areas as predictors, with rosette areas measured  $\leq 28$  DAS generally yielding  
824 low ( $< 0.10$ ) and in many cases negative median test  $R^2$  scores. When using the earliest

825 rosette area (14 DAS) as predictor, the median pooled  $R^2$  and PCC scores are however still in  
826 the ranges [0.05, 0.20] and [0.27, 0.45], respectively, indicating that even the earliest rosette  
827 area measurements contain some information on final yield phenotypes.

828



829

830 **Fig. 5** Predictive power of early rosette areas for yield phenotypes. In each subplot, median  
831 test  $R^2$  values are plotted for lme models predicting the given phenotype from early rosette  
832 areas v2 (14-42 DAS, x-axis). Only mature phenotypes that can be predicted from rosette area  
833 (42 DAS) with a median test  $R^2 > 0.1$  are shown. Blue lines are ordinary least-squares linear  
834 regressions, with shaded areas indicating 95% confidence intervals on the trendline. Most  
835 phenotypes exhibit a rather dichotomous median test  $R^2$  profile with rosette areas v2 from  
836 14 to 28 DAS yielding substantially lower median test  $R^2$  values than rosette areas v2 from  
837 32 to 42 DAS. Accordingly, linear model fits at 28 and 32 DAS are often poor.

838

839



## 840 DISCUSSION

841 In this study, we used machine learning models to predict the phenotypes of individual *B.*  
842 *napus* Darmor plants grown in the same field from rosette-stage leaf gene expression data.  
843 Our results show that many plant phenotypes can be predicted to a substantial extent from  
844 leaf 8 gene expression. Phenotypes closely related in time and space to the material sampled  
845 for RNA-seq, in particular leaf 8 phenotypes, generally feature good prediction performance,  
846 in accordance with results obtained earlier in a similar setup for maize (28). Interestingly  
847 however, also many of the phenotypes measured at the end of the growing season, ~5.5  
848 months after leaf sampling for RNA-seq, feature high prediction performance. In particular  
849 seed yield, silique and dry weight traits exhibit prediction performance scores in the same  
850 range as the autumnal leaf and rosette phenotypes.

851

852 Azodi et al. (77) predicted several agronomically relevant mature plant traits (plant height,  
853 grain yield and flowering time) in a population of maize inbred lines from genetic marker data,  
854 whole-seedling transcriptome data and combinations thereof. Their transcriptome-based  
855 models exhibited PCC scores between predicted and measured values in the range [0.50,  
856 0.61] for flowering time, [0.42, 0.51] for plant height and [0.47, 0.55] for 300 kernel weight  
857 (77). In the present study, the transcriptome-based models for mature plant traits in *B. napus*  
858 (ignoring ratio phenotypes) exhibit median pooled PCC scores in the range [0.57, 0.77] for  
859 seed phenotypes, [0.51, 0.74] for silique phenotypes, [0.56, 0.73] for shoot dry weight  
860 phenotypes, [0.40, 0.56] for branch count phenotypes, [0.40, 0.53] for plant height (278 DAS)  
861 and [0.07, 0.36] for root phenotypes (**Table 2, Additional File 7: Table S6**). Comparing the  
862 observed PCC ranges of both studies suggests that mature traits of individual plants of the  
863 same line grown in the same field are as predictable from early-stage transcriptome data as

864 average mature traits in a diversity panel. However, direct comparison of the PCC values  
865 across studies is complicated by differences in the phenotypes predicted, prediction and  
866 scoring methodology and factors affecting model training, scoring and overfitting potential  
867 such as the study population size (388 lines in the maize study versus 62 *B. napus* plants in  
868 the present study) and the number of potential model features (31,238 genes in the maize  
869 study versus 76,808 genes in the *B. napus* dataset). Also the species difference and the tissue  
870 and developmental time point sampled for RNA-seq (whole seedlings at the V1 stage for  
871 maize versus rosette leaf 8, 81 DAS, for *B. napus*) may impact how well a transcriptome can  
872 predict a given phenotype. The most comparable models are likely the whole-transcriptome-  
873 based random forest model for maize plant height, with a PCC of 0.42 (77), and the median-  
874 filter random forest model for the height of individual *B. napus* plants (without feature  
875 selection other than removing genes with *rlog* expression >0 in less than half of the samples,  
876 reducing the feature set to 55,166 genes), with a median pooled PCC of 0.43 (**Additional File**  
877 **7: Table S6**).

878

879 Given that the single-plant transcriptome data can quantitatively predict many plant  
880 phenotypes better than expected by chance, the top predictor genes may shed light on  
881 biological processes that impact phenotypes in the field. Many of the top predictors in the TF  
882 models for seed, silique and dry weight phenotypes for instance are known to function in the  
883 floral transition. From the perspective of our experimental setup, it makes sense that such  
884 genes are recovered, as it is known that the floral transition starts in autumn in winter-type  
885 *B. napus* accessions (78, 79), i.e. around the time that rosette leaves were harvested for RNA-  
886 seq in the present field trial, and is set in motion to a large extent by systemic signals  
887 emanating from leaves in Brassicaceae and other plant families (80-82).

888

889 Mechanistic interpretation of the correlational links between top predictor genes and  
890 phenotypes is however not straightforward. Putative orthologs of *AtHB1* and *AtHB16* are for  
891 instance found among the top predictors positively correlated with both leaf and yield  
892 phenotypes (**Fig. 4, Additional File 2: Fig. S9**), but upregulation of these genes in *A. thaliana*  
893 was previously found to lead to smaller and more serrated leaves (35, 36), to delay the  
894 vegetative-to-reproductive phase transition and to result in siliques bearing fewer seeds (34,  
895 36). Some top predictors that correlate negatively with yield phenotypes have putative *A.*  
896 *thaliana* orthologs that are thought to function primarily as negative regulators of the floral  
897 transition in leaves, e.g. *AtNF-YA* genes (58, 59), but others are putatively orthologous to a  
898 positive regulator of the floral transition, such as *AtFUL*. Other floral transition regulators  
899 recovered as predictors in our yield models, e.g. orthologs of *AtFLC* and *AtEFM*, do not by  
900 themselves exhibit a significant positive or negative correlation with yield phenotypes.

901

902 Most likely, the associations recovered between individual plant phenotypes and autumnal  
903 leaf gene expression patterns are due to developmental timing differences among the plants,  
904 rather than reflecting the effects of upregulation or downregulation of specific regulators. In  
905 the *A. thaliana* developmental gene expression atlas of Klepikova et al. (83), orthologs of  
906 predictors positively correlated with leaf size such as *AtHB1* and *AtHB16* are more highly  
907 expressed in mature *A. thaliana* leaves (at flowering), while orthologs of predictors negatively  
908 correlated with leaf size such as *AtREV*, *AtWOX5* and *AtHAT3* are more highly expressed in  
909 young leaves. This suggests that plants with low expression of *AtHB1/16* orthologs and high  
910 expression of *AtREV/AtWOX5/AtHAT3* orthologs had a more juvenile (and hence smaller) leaf  
911 8 at the autumnal sampling time point, which explains the observed gene expression-leaf

912 phenotype correlations. That autumnal leaf phenotypes and final yield phenotypes have  
913 several developmental predictors in common (e.g. *AtHB1*) and that the autumnal leaf  
914 phenotypes themselves are also predictive of yield indicates that the developmental  
915 differences in autumn impacted final yield. These differences were not limited to differences  
916 in leaf development, as evidenced by the fact that the predictor sets for both leaf and yield  
917 phenotypes also contain regulators of plant-wide developmental phase transitions occurring  
918 in autumn (juvenile-to-adult, vegetative-to-reproductive).

919

920 In summary, our results indicate that the yield potential of the individual plants was already  
921 determined to a large extent by their developmental state at the time of leaf sampling in  
922 autumn. Mendham and Scott (84) previously found that the size of winter-type *B. napus*  
923 plants at the time of inflorescence initiation affects their yield potential, in the context of an  
924 experiment assessing sowing date effects on yield. Our results show that even when sown on  
925 the same date in the same field, individual winter-type *B. napus* plants of the same line display  
926 developmental differences in autumn that correlate with yield differences in spring. Even if  
927 only part of the variability in e.g. total seed weight (CV = 46.9%) observed in our trial is due to  
928 autumnal effects on plant growth and development, the gains of mitigating such effects could  
929 be substantial.

930

931 The question remains however what could have caused the developmental differences  
932 among plants in the present field trial. One potential cause is differences in seed germination  
933 and seedling emergence across the field. In wheat, it was established previously that relative  
934 differences in seedling emergence date are strongly correlated with differences in final yield  
935 (85). Next to seed quality, many environmental factors are known to impact the timing of

936 seed germination and seedling emergence, including soil structure (86), soil temperature (87),  
937 sowing depth (85, 87), soil water potential, oxygenation and light quality (88), and soil  
938 nutrients (89). The seedling emergence date was not recorded in the present field trial, but  
939 the closest proxy that was measured, namely rosette area at 14 DAS, was found to be a bad  
940 predictor for yield, indicating that variation in seed germination and seedling emergence  
941 across the field did not by themselves have a major impact on yield. The observation that  
942 later rosette areas are progressively better at predicting yield rather suggests that  
943 developmental differences among plants accumulated over time.

944

945 The observation that genes involved in the regulation of circadian rhythm, photoperiodism  
946 and the vegetative-to-reproductive phase transition are on average more spatially  
947 autocorrelated in the autumnal gene expression dataset than the average gene suggests that  
948 spatially patterned micro-environmental factors may be linked to the variability of  
949 developmental gene expression in autumn, and ultimately yield variability in spring. That the  
950 phenotypes are influenced by environmental factors is also suggested by the observation that  
951 the sets of genes associated with leaf and yield phenotypes are heavily enriched in genes  
952 involved in responses to abiotic and biotic stimuli and nutrient levels (**Additional File 6: Table**  
953 **S5**). The finding that developmental processes feature more prominently in the TF-based  
954 phenotype prediction models than responses to environmental stimuli indicates that micro-  
955 environmental variations among plants in the present field trial may have influenced plant  
956 phenotypes mainly by influencing development. More work is needed however to establish  
957 whether and how micro-environmental variability impacts the growth and development of  
958 individual plants in the same field. To address this, additional field trials need to be performed  
959 in which, next to the gene expression and phenotypes of individual plants, also a range of

960 environmental parameters is measured on the single-plant level (e.g. soil structure and  
961 chemistry).

962

963 Additional single-plant field trials are also needed to assess to what extent the predictive  
964 models, gene-phenotype and process-phenotype associations learned from the present field  
965 trial generalize to other soils and meteorological conditions, other time points or tissues  
966 sampled for RNA-seq, and other cultivars. Given the developmental nature of many of the top  
967 predictors in the current models, it is likely that our current prediction models, based on leaf  
968 gene expression data for a single field trial at a single time point, will not perform well when  
969 applied on follow-up field trials, even when using the same cultivar in a similar field under  
970 roughly the same climate conditions. Differences in weather conditions and other  
971 environmental factors across trials may for instance influence the timing of developmental  
972 phase transitions, making it all but impossible to sample the exact same developmental time  
973 point in follow-up trials. If leaf gene expression were to be profiled at a slightly earlier or later  
974 developmental time point than in the present trial, the current top predictors may no longer  
975 be adequate phenotype proxies and other genes that function earlier or later in e.g. the floral  
976 transition may become relevant instead. The construction of robust prediction models will  
977 therefore likely require single-plant data generated under a wide variety of field conditions  
978 and sampling schemes. We want to emphasize however that quantitative prediction of single-  
979 plant phenotypes is not the primary goal we envision for single-plant omics experiments.  
980 Rather, the primary aim is to identify which biological processes, environmental factors and  
981 associated genes may influence plant phenotypes in the field. In this respect, any additional  
982 genes and processes identified in follow-up trials would add to our overall knowledge on how  
983 rapeseed plants grow in a field.

984

985 It is worth pointing out that the dataset generated in this study may also serve other purposes  
986 than gene-phenotype association. Earlier, we have shown that field-generated single-plant  
987 transcriptomics data can also be used efficiently to predict the function of genes (28). Given  
988 the complex genome duplication history of *B. napus* (90), the combination of gene function  
989 prediction and gene-phenotype association may be particularly useful to shed light on which  
990 *B. napus* genes in a (long) list of paralogs are most likely functionally orthologous to a given  
991 *A. thaliana* gene, and how paralogs have diverged in function. This knowledge may in turn be  
992 useful in the context of genetic engineering and breeding efforts to optimize yield and stress  
993 tolerance in *B. napus*.

994

## 995 **CONCLUSIONS**

996 We have shown that individual *B. napus* plants of the same background grown in the same  
997 field exhibit considerable variation in gene expression and phenotypes, and that the plants'  
998 autumnal gene expression profiles have predictive power for their yield in spring. Many of the  
999 top yield predictor genes are associated with developmental processes occurring in autumn  
1000 in winter-type *B. napus*, such as the juvenile-to-adult and floral transitions. Together, this  
1001 indicates that autumnal development has a major influence on the yield potential of winter-  
1002 type *B. napus* plants. In summary, our data show that profiling individual plants under  
1003 uncontrolled field conditions is a valid strategy for identifying genes and processes influencing  
1004 crop yield in the field.

1005

1006

1007

## 1008    **METHODS**

### 1009    **Field trial setup**

1010    Seeds from the winter-type *Brassica napus* accession Darmor (BnASSYST-120) were sown in  
1011    a field in Merelbeke, Belgium (50°58'24.9"N 3°46'49.1"E) on September 8, 2016. Three seeds  
1012    were sown at ~2 cm depth at each of 100 points arranged in a 10x10 grid with 0.5 m spacing  
1013    within and between rows. Seedlings were thinned out to leave one seedling growing at each  
1014    grid point. Early- and late-emerging seedlings were pruned preferentially (based on visual  
1015    assessment) to make the remaining seedling population as homogeneous as possible. At two  
1016    points, no seedlings emerged.

1017    Plots of *Miscanthus sinensis*, *M. sacchariflorus* and *Miscanthus* hybrids were grown to the  
1018    northeast and southeast of the *B. napus* field trial, and maize was grown to the northwest, at  
1019    distances >5 m. The field plot was surrounded by chicken wire and covered by netting to keep  
1020    out birds and large herbivores. The netting was removed in spring when plants grew taller  
1021    than ~1 m. Additionally, perimeter fencing was used to protect the field trial and the mobile  
1022    weather station on site (see **Additional File 1: Table S1** for weather station data).

1023    After germination, individual plant images were taken twice a week between September 22  
1024    and October 20, 2016 (9 time points) to assess the projected leaf area of the growing rosettes.  
1025    Nadir images were taken using a D90 camera (Nikon Inc., USA) equipped with a 35 mm lens  
1026    (AF-S DX Nikkor 35 mm F1.8G, Nikon Inc., USA) set at iso 200, f/8. The shutter speed could  
1027    vary to allow for a proper exposure, determined by the camera. The camera/tripod was  
1028    positioned away from the sun to avoid shadows in the images taken. For each time point a  
1029    grey calibration card (Novoflex grey card 15 x 20 cm, NOVOFLEX Präzisionstechnik GmbH,



1030 Germany) was used to correct the white balance. This card was also used as reference to  
1031 convert pixels to areas in  $\text{cm}^2$  (see below). The ground sampling difference (GSD) was 0.015  
1032 cm/pixel.

1033 At 74 DAS, the length and width of leaf 6 (counting upward from the first true leaf) were  
1034 measured non-destructively, leaf 6 lesion and total rosette lesion severity were scored and  
1035 the number of fully emerged rosette leaves (area  $> \sim 2 \text{ cm}^2$ ) was recorded. At 76 DAS, leaf 8  
1036 length and width were measured and leaf 8 lesions were scored. The width of leaf blades was  
1037 measured at the widest point. Leaf lengths were measured from the leaf tip to the point  
1038 where the petiole first lacked conspicuous laminar tissue (looking from the leaf tip toward the  
1039 base). Lesion severity was scored qualitatively on a scale from 0 (lesions cover at most five  
1040 percent of the leaf blade or rosette) to 2 (more than half of the leaf or rosette eaten).

1041 At 81 DAS, on November 28, 2016, the eighth rosette leaves of 62 non-border plants (i.e. the  
1042 plants at all non-border locations where seedlings emerged) were harvested for RNA-  
1043 sequencing in a time span of  $\sim 1$  hour (13:25-14:27). Leaves were cut off where the petiole  
1044 first lacked conspicuous laminar tissue (looking from the leaf tip toward the base) and washed  
1045 with DEPC-treated and sterilized water. The chlorophyll content of each leaf was measured  
1046 at four different positions on the leaf with a CCM-200 chlorophyll content meter (Opti-  
1047 Sciences, Inc., Hudson, USA), and the average of these measurements was used in the  
1048 analyses. Leaves were then photographed twice against a white background with a piece of  
1049 millimeter paper to assess the image scale and perspective, a ruler, and color and greyscale  
1050 references, the second time covered with a glass plate to flatten them. Next, the midvein of  
1051 every leaf was cut out using scissors, and the residual leaf material was folded into a pre-  
1052 weighed 50 ml tube. The filled 50 ml tube and the midvein were weighed together to measure

1053 leaf fresh weight, after which the tube was stored in liquid nitrogen on the field. The entire  
1054 leaf processing pipeline, from cutting a leaf to storing it in liquid nitrogen, was completed for  
1055 each leaf in less than 5 minutes.

1056 After leaf sampling, the plants were left to overwinter and set seed in spring. After bolting,  
1057 plant height was measured from ground level to the top of the primary flowering stem at 13  
1058 time points between 189 and 231 DAS (**Additional File 1: Table S1**). One of the plants sampled  
1059 in autumn for RNA-seq, 01C, did not survive until the end of the growth season. The remaining  
1060 61 non-border plants were harvested on June 13, 2017 (278 DAS), at which time ~50% of  
1061 seeds had started changing color from green to black but no significant pod shattering or seed  
1062 predation had occurred. Final plant height at 278 DAS was measured on the field, from ground  
1063 level to the top of the primary flowering stem. Afterwards, shoots were cut off and the root  
1064 systems were dug up. Taproot length was measured from ground level to the deepest root  
1065 tip. Root system width was measured perpendicular to the taproot at the root system's widest  
1066 point.

1067 For each harvested plant, the primary flowering stem plus its cauline secondary  
1068 inflorescences (stem 1) and the secondary inflorescence stems branching at ground level (side  
1069 stems) were dried in two separate bags in a well-ventilated, dry attic. The number of branches  
1070 and siliques per stem, the total shoot dry weight and the dry weight of stem 1 were measured  
1071 on dried plants. Seeds were recovered manually from the dried-out pods for stem 1 and the  
1072 side stems separately, and separated from dust and small pod debris using a customized seed  
1073 aspirator with vibration channel (Baumann Saatzuchtbedarf GmbH, Waldenburg, Germany).  
1074 The resulting seed batches for stem 1 and the side stems were weighed and counted using an

1075 elmor C3 seed counter (elmor AG, Schwyz, Switzerland). Seed counts and weights are  
1076 reported for stem 1 and the entire plant (i.e. the sum of stem 1 and the side stems).

## 1077 **Determination of shoot growth parameters**

1078 Shoot growth parameters (time of maximum shoot growth  $t_m$ , maximum shoot growth rate  
1079 and the end of shoot growth  $t_e$ ) were derived by fitting a beta-sigmoid growth curve to the  
1080 time series of 14 plant height measurements between 189 and 278 DAS (91) :

$$\begin{aligned} 1081 \quad h(t) &= h_0 + (h_{max} - h_0) * \left(1 + \frac{t_e - t}{t_e - t_m}\right) * \left(\frac{t}{t_e}\right)^{\frac{t_e}{t_e - t_m}} & t < t_e \\ 1082 \quad h(t) &= h_{max} & t \geq t_e \end{aligned} \quad (\text{Eq. 1})$$

1083

1084 With  $h(t)$  the plant height at plant age  $t$ ,  $h_0$  and  $h_{max}$  the initial and final plant height at  $t =$   
1085 0 and  $t = t_e$ , respectively,  $t_e$  the plant age at the end of growth and  $t_m$  the plant age at the  
1086 moment of maximal growth. Before curve fitting, the time points in day of year (DOY) at which  
1087 the plant heights were measured were translated to plant ages  $t$  in growing degree days  
1088 (GDD), i.e.  $t(i) = \sum_{j=0}^{j=i} \max(T_j - T_b, 0)$  with  $i$  the time point in DOY,  $T_j$  the average air  
1089 temperature at  $j$  DOY (**Additional File 1: Table S1**) and  $T_b = 5$  °C a base temperature below  
1090 which no growth is assumed to occur (79, 92). Optimization of the parameters  $h_0$ ,  $h_{max}$ ,  $t_e$   
1091 and  $t_m$  was done with the *nls* function in R using the 'port' algorithm. The maximum shoot  
1092 growth rate was obtained by calculating the derivative of  $h(t)$  (**Eq. 1**) at  $t_m$ . After curve fitting,  
1093 the values obtained for  $t_m$  and  $t_e$  were converted back from GDD to DOY and subsequently  
1094 to DAS.

## 1095 **Image-based phenotyping**

1096 Leaf 8 areas (81 DAS) were estimated by segmenting the flattened leaf images taken at the  
1097 time of leaf harvest. The millimeter grid scale on each image was used to correct for  
1098 perspective distortion and to create a uniform spatial resolution across the entire image of  
1099 100 pixels per cm. Images were cropped to remove the grid scale and sample label.  
1100 Segmentation was done by training a U-Net convolutional neural network (93) on a small  
1101 dataset of 25 images for which random patches of foreground and background were  
1102 annotated using VGG Image Annotator (via) v:2.0.7 (94). Random cropping, resizing, rotating  
1103 (by multiples of 90 degrees), mirroring, color-jittering and gaussian blurring were applied to  
1104 artificially increase the training dataset size. The training was done using the Adam optimizer  
1105 (95) in Pytorch v:1.7.1 (96) with default settings. The pixel-wise cross-entropy loss was back-  
1106 propagated only for annotated regions of each image. The learning rate was initially set to 1e-  
1107 3 and was automatically halved as soon as the minimal training loss stagnated for more than  
1108 3 epochs. The network was trained for 16 epochs. The trained network was validated by  
1109 visually evaluating it on unseen images, and then applied to all flattened leaf images.

1110 Leaf 8 length and width at 81 DAS were measured on the flattened leaf images using ImageJ  
1111 v:1.50 (97). For measuring leaf 8 length, the midvein was traced from the leaf tip to the cutting  
1112 point (i.e. where the petiole first lacked conspicuous laminar tissue) using the ImageJ  
1113 segmented line tool. Leaf 8 width was measured at the widest point.

1114 For measuring the projected area of the rosettes photographed at 42 DAS (i.e. the rosette  
1115 imaging date closest to leaf sampling), a dedicated script was developed using the image  
1116 analysis software HALCON (version 13.0.1.1, MVTec Software GmbH, Germany). First, the  
1117 images were cropped to remove parts of adjacent plants visible on the pictures. To remove  
1118 noise, both a gentle Gaussian filter and a median filter were applied. Each RGB image was

1119 then converted to the HSV color space, where the Hue channel was used to select the green  
1120 plant parts using a threshold range for the green pixels (34-80) defined based on trial and  
1121 error. Care was taken to also include the petioles. After this, a 'closing\_circle' operator was  
1122 used and remaining small lesions (due to insect damage) were filled up using the 'fill\_up'  
1123 operator. Only the largest segmented area was taken into account, to differentiate between  
1124 the plant of interest and small weeds nearby.

1125 The HALCON segmentation strategy worked well for the rosette images taken at 42 DAS, but  
1126 regularly produced segmentation errors for images of smaller rosettes taken closer to the  
1127 sowing date. An alternative segmentation approach was therefore used on rosette images  
1128 taken at 14, 18, 21, 25, 28, 32, 35 and 39 DAS (and 42 DAS as control). The main difficulty for  
1129 the earlier time points is distinguishing small rosettes from weeds and other distracting  
1130 objects occurring on the field. This requires an algorithm with a larger field of view than what  
1131 a HALCON script or standard U-net (see above) can provide. Instead, a standard pre-trained  
1132 DenseNet M161 (98) was taken and augmented with additional bilinear upsampling layers  
1133 after each 'dense' layer of the original algorithm. That is, the last feature layer of DenseNet  
1134 was upsampled with bilinear interpolation and a weighted sum was made with the higher  
1135 resolution 'dense' features. This was repeated for each dense layer up to the original input  
1136 resolution. The network was trained for 175 epochs (final mean epoch loss = 0.01) on 54 hand-  
1137 labeled images (6 images per time point) using stochastic gradient descent (SGD) with  
1138 momentum (learning rate = 0.001 and momentum = 0.99). The learning rate was divided by  
1139 10 each time the train loss plateaued for more than 4 epochs. Image rotations, mirroring and  
1140 HSV augmentations were used to augment the training data. The trained model was used to  
1141 segment all rosette images. After segmentation, a post-processing step was performed to

1142 remove segmented parts of *B. napus* plants adjacent to the plant of interest and remaining  
1143 weeds, using scikit-image v: 0.19.2 (99). Only the connected component closest to the  
1144 centroid of the image and other components within a 25-pixel distance of this central  
1145 component (e.g. leaves of which the stalk was segmented incorrectly because of a lower  
1146 chlorophyll content) were associated with the plant of interest. Connected components with  
1147 an area less than 10,000 pixels were filtered out to eliminate small weeds. This approach was  
1148 evaluated visually for all segmentations and proved to work well for most plants.  
1149 Segmentations with missing plant parts or weeds that weren't filtered out by this post-  
1150 processing step were manually corrected. A grey calibration card (Novoflex grey card 15 x 20  
1151 cm, NOVOFLEX Präzisionstechnik GmbH, Germany) was used as a reference to convert pixels  
1152 to areas in cm<sup>2</sup>. The projected rosette areas at 42 DAS estimated by this segmentation  
1153 approach exhibit a Pearson correlation of 0.997 with the areas estimated by the  
1154 aforementioned HALCON script.

## 1155 **RNA sequencing**

1156 The frozen leaf samples for the 62 harvested non-border plants were grinded, and total RNA  
1157 was extracted using the guanidinium thiocyanate-phenol-chloroform extraction method  
1158 using TRI-reagent (Thermo Fisher Scientific) followed by DNA digestion using the RQ1 RNase-  
1159 free DNase kit (Promega). ds cDNA was prepared using the Maxima H Minus Double-Stranded  
1160 cDNA Synthesis Kit (#K2561, Thermo Fisher Scientific) to a concentration of ~17-38 ng/ul in  
1161 10mM Tris-Cl buffer (pH 8.5) at a minimum volume of 30ul. (~0.6 - 1.1 ug total). ds cDNA  
1162 samples were sent to the University of Missouri Genomics Technology Core, where library  
1163 preparation was performed (average insert size of 500 bp) using the Illumina TruSeq DNA  
1164 PCR-Free Library Prep Kit according to the protocol described in (100). 250 bp paired-end

1165 sequencing was performed at the Tufts University Genomics Core on an Illumina HiSeq 2500  
1166 machine in Rapid Run mode. The samples were sequenced in 3 batches (**Additional File 1:**  
1167 **Table S1**).

1168 The raw RNA-seq data was processed using a custom Galaxy pipeline (101) implementing the  
1169 following steps. First, the fastq files were quality-checked using FastQC (v:0.5.1) (102). Next,  
1170 Trimmomatic (v:0.32.1) (103) was used to remove adapters, read fragments with average  
1171 quality below 20 and trimmed reads shorter than 125 base pairs. The trimmed and filtered  
1172 reads were mapped to the *Brassica napus* Darmor-bzh reference genome v:5  
1173 (<https://www.genoscope.cns.fr/brassicanapus/data/>) (90) using HISAT2 v:2.0.5 (104) with  
1174 default values for all parameters. Only the uniquely mapping reads or (in the case of multiple  
1175 mappings) the best secondary alignment were kept for the following analyses. The mapping  
1176 files were quantified using HTSeq v:0.6.1p1 (105) with the option 'Intersection-union', using  
1177 the genome annotation of the *Brassica napus* Darmor-bzh reference genome v:5  
1178 (<https://www.genoscope.cns.fr/brassicanapus/data/>). No filtering steps were performed  
1179 during preprocessing except for removing genes that were not expressed in any samples.  
1180 Counts were normalized across samples and batches using a modified regularized log (*rlog*)  
1181 model of the DESeq2 (106) package in R. Counts are still modeled in the same way as in the  
1182 original *rlog* implementation, that is :

1183

$$1184 \quad k_{ij} \sim NB(\mu_{ij}, \alpha_i)$$

$$1185 \quad \mu_{ij} = s_j \times q_{ij} \quad (\text{Eq. 2})$$

$$1186 \quad \log_2(q_{ij}) = \mathbf{x}_j \cdot \boldsymbol{\beta}_i$$

1187

1188 Where  $k_{ij} \in \mathbb{N}^+$  is the count of gene  $i$  in sample  $j$ , which is assumed to be sampled from a  
1189 negative binomial distribution (NB) with estimated mean  $\mu_{ij} \in \mathbb{R}^+$  and estimated dispersion  
1190 of the  $i$ th gene  $\alpha_i$ .  $\mu_{ij}$  is taken as the expected count  $q_{ij}$  for a ‘typical’ library size (i.e. with a  
1191 size factor  $s_j = 1$ ), scaled by a library size normalization factor  $s_j$  for sample  $j$ . Note that  $q_{ij}$   
1192 still contains batch effects :  $\mathbf{x}_j \in \mathbb{R}^p$  is a vector of  $p = 65$  predictors for sample  $j$ , including  
1193 an intercept, 2 dummy variables for the smallest sequencing batches (1 and 3) that capture  
1194 batch effects relative to the largest sequencing batch (2, the effects of which are absorbed in  
1195 the intercept) and dummy variables for each of the 62 plants that were sampled.  $\boldsymbol{\beta}_i \in \mathbb{R}^p$   
1196 contains the estimated coefficients for those predictors for gene  $i$ . As in (106), an empirical  
1197 Bayes shrinkage procedure is used to estimate  $\boldsymbol{\beta}_i$ , using a flat prior for the intercept  $\beta_{i0}$  and  
1198 the sequencing batch coefficients, and a zero-centered normal prior for each plant coefficient  
1199  $\beta_{ip_j}$  (with  $p_j$  the index of the plant corresponding to sample  $j$ ), with prior variance estimated  
1200 using quantile matching as described in Love et al. (106). There are only two differences  
1201 compared to Love et al. (106): the first is the addition of two batch coefficients as fixed effects  
1202 in the design matrix, and the second is that log-fold changes used in the prior random effect  
1203 variance computation are estimated relative to the mean of each batch instead of to the mean  
1204 of all samples. Once the model is estimated,  $rlog$  counts are computed as in Love et al. (106),  
1205 that is:

1206

$$1207 \quad rlog_{ij} \equiv \beta_{i0} + \beta_{ip_j} \quad (\text{Eq. 3})$$

1208

1209 Note that all samples  $j$  belonging to the same plant (technical repeats) have the same value  
1210 for  $\beta_{ip_j}$ . The modified  $rlog$  transformation removes library size effects and batch effects,



1211 unites technical repeats into one estimate and log-transforms the data (reducing  
1212 heteroscedasticity) in a single step. In addition, using random effects for each plant allows  
1213 pooling information from technical repeats while simultaneously basing variance estimates  
1214 on all samples (including samples without technical repeats). This method therefore makes  
1215 maximal use of the available data. The resulting data is show in **Additional File 2: Fig. S10**.

## 1216 **SNP detection and population structure analysis**

1217 Trimmed and filtered RNA-seq reads were aligned to the *Brassica napus* Darmor-bzh  
1218 reference genome v:5 (<https://www.genoscope.cns.fr/brassicanapus/data/>) (90) using  
1219 HISAT2 v:2.0.5 (104) with default values for all parameters. Genomic variants were detected  
1220 for each plant using NGSEP v:3.3.2 (107) on the aligned reads. For downstream analyses, we  
1221 focused on biallelic SNPs with a minimum genotype quality of 40 and called in at least 49  
1222 samples (80% of the population). Missing calls were imputed using Beagle v:5.1 (108) using  
1223 default parameters, and only SNPs with minor allele frequency (MAF)  $\geq 0.05$  after imputation  
1224 were kept, resulting in a dataset of 23,188 SNPs.

1225

1226 A neighbor-joining tree was made based on the SNP dataset with TASSEL v:5.2.60 (109), using  
1227 1-IBS (identity by state) as the distance measure while setting the distance from an individual  
1228 to itself to zero. The tree was rendered using the polar tree layout in FigTree v:1.4.3 (110).

1229

## 1230 **Spatial autocorrelation analysis**

1231 Moran's I was calculated for each gene (phenotype) as  $I = \frac{n}{w} \frac{(\mathbf{x}-\bar{\mathbf{x}})^T \mathbf{C}(\mathbf{x}-\bar{\mathbf{x}})}{\|(\mathbf{x}-\bar{\mathbf{x}})\|^2}$ . Where  $\mathbf{x}$  is a  
1232 column vector of *rlog* gene expression (phenotype) values,  $n$  is the number of samples and  $w$

1233 is the sum of elements of the connectivity matrix  $\mathbf{C}$ . For  $\mathbf{C}$  a binary  $n \times n$  ‘queens’  
1234 connectivity matrix was chosen, meaning that neighboring horizontal, vertical and diagonal  
1235 plants are seen as connected. Note that  $\mathbf{C}$  can differ from one phenotype to the next since  
1236 not all phenotypes were available for all samples. For each gene (phenotype), the Moran’s I  
1237 was recalculated on  $10^5$  random permutations of  $\mathbf{x}$  to obtain an empirical null distribution,  
1238 which was then compared to the real Moran’s I to obtain a  $p$ -value. Finally,  $p$ -values were  
1239 corrected across all genes (phenotypes) using the Benjamini-Hochberg (BH) procedure (111).  
1240 All calculations were done using the PySAL python library (112).

1241

## 1242 **Variance analysis**

1243 Principal component analysis (PCA) was done on various normalized versions of the gene  
1244 expression count matrix and on the phenotype dataset (including qualitative phenotypes such  
1245 as leaf 6 lesion severity (74 DAS) but excluding the plant height and rosette area time series  
1246 except for the final time points, i.e. plant height (278 DAS) and rosette area (42 DAS)), using  
1247 the ‘prcomp’ function in the R stats package on the centered gene expression datasets and  
1248 the ‘ppca’ method in `pcaMethods v:1.88.0` (113) on the z-scored phenotype dataset.  
1249 Phenotype distributions were plotted using the ‘histogram’ function in Matlab R2018b with  
1250 probability normalization option. Shapiro-Wilk and Anderson-Darling tests were performed  
1251 using the ‘normalitytest’ script (114) and ‘adtest’ functions in Matlab R2018b, respectively.  
1252 Outliers were defined as values more than three scaled median absolute deviations (MAD)  
1253 away from the median, as is default in the Matlab R2018b ‘isOutlier’ function. Outliers were  
1254 only removed for the purpose of calculating their effect on the phenotypes’ normality and  
1255 coefficient of variation (CV), all other analyses used the complete phenotype dataset.

1256

1257 Normalized coefficients of variation (*normCVs*) for gene expression profiles were computed  
1258 on batch and library size corrected data (without *rlog* transform). Normalized counts were  
1259 obtained as  $x_{ij} = k_{ij}/(\beta_{ib_j} \times s_j)$  where  $\beta_{ib_j}$  is the batch effect for gene  $i$  in sample  $j$  as  
1260 estimated in the *rlog* calculation (see above). Since batch 2 is absorbed in the intercept,  $\beta_{ib_j} =$   
1261 1 for samples of batch 2. Contrary to the *rlog* transform, this method does not collapse  
1262 technical repeats, and they were instead collapsed by averaging (as in **Additional File 2: Fig.**  
1263 **S10**, panel B, but without the  $\log_2$ -transform). From here on, variance analysis followed the  
1264 same procedure as described in Cruz, De Meyer et al. (28). Briefly, a trendline was fitted to  
1265 the  $CV^2$  versus mean expression relationship (omitting genes expressed in <10 samples) using  
1266 a generalized linear model (GLM) of the gamma family with identity link of the form  $CV^2(\mathbf{x}) =$   
1267  $a/\bar{x} + b$ , with fitting parameters  $a$  and  $b$  (115) (**Additional File 2: Fig. S11**). Code from the  
1268 M3Drop R package (116) was used for this purpose. A normalized CV accounting for the  
1269 observed mean-variance relationship was then calculated as  $normCV(\mathbf{x}) = \log_2(CV^2(\mathbf{x})/$   
1270  $trend(\bar{x}))$  where  $trend(\bar{x})$  is the fitted value at the mean of  $\mathbf{x}$ .

1271

## 1272 **GO enrichment analysis**

1273 A Gene Ontology (GO) annotation for *Brassica napus* was generated using the TRAPID v.2.0  
1274 platform (117) with default parameters on April 16, 2020. Transcript sequences parsed from  
1275 the *B. napus* Darmor-bzh reference genome annotation v:5 (90) using the gffread v.0.9.6  
1276 utility (118) were used as input for TRAPID, and PLAZA 4.5 dicots (119) was used as the  
1277 reference database. GO enrichment  $p$ -values were calculated with hypergeometric tests and  
1278 adjusted for multiple testing ( $q$ -values) using the BH procedure (111), either using custom R  
1279 scripts or using BiNGO v:3.0.3 (120). GO categories gravitating toward the top or bottom of

1280 gene lists ranked in order of decreasing Moran's I or normalized CV were detected using two-  
1281 sided Mann-Whitney U tests (with genes belonging to the GO category of interest classified  
1282 as group 1 and other genes as group 2), as implemented in the 'wilcox.test' function in the R  
1283 stats package v:4.0.5, followed by BH *p*-value adjustment.

1284

## 1285 **Ortholog inference**

1286 Putative *A. thaliana* orthologs of *B. napus* genes were identified in two steps. First, putative  
1287 orthologs of *B. napus* genes were identified in *B. rapa* and *B. oleracea* (source of the A and C  
1288 subgenomes of *B. napus*, respectively), based on best similarity hits returned by TRAPID v.2.0  
1289 (117) and on the syntenic relationships reported in Chalhoub et al. (90) and Sun et al. (121).  
1290 Second, putative *A. thaliana* orthologs of the identified *B. rapa* and *B. oleracea* genes were  
1291 retrieved from PLAZA 4.5 dicots (119), which provides orthology inferences integrating four  
1292 different lines of evidence : orthogroup inference within gene families using OrthoFinder  
1293 (122), orthology inference using gene tree-species tree reconciliation, orthology inference  
1294 from best DIAMOND (123) hits and their inparalogs, and positional orthology inference  
1295 through collinearity analysis (124). The most likely *A. thaliana* orthologs of a given *B. napus*  
1296 gene were taken to be the putative orthologs that are most strongly supported across both  
1297 inference steps.

1298

## 1299 **Single-feature phenotype prediction models**

1300 **Single-gene models.** Linear mixed-effects (lme) models (125) were fitted to predict a  
1301 phenotype given the expression profile of a single gene. That is, given a phenotype vector *y*

1302 and a vector  $\mathbf{x}$  of a given gene's z-scored expression values across the field, we fit the  
1303 following model:

1304

$$\begin{aligned} \mathbf{y} &= \beta_0 + \beta_1 \mathbf{x} + \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon} &\sim \mathcal{N}(0, \Sigma) \end{aligned} \tag{Eq. 4}$$

1307

1308 where  $\beta_0$  is the intercept (average phenotype value),  $\beta_1$  the gene effect coefficient, and  $\boldsymbol{\varepsilon}$  the  
1309 residual error which is assumed to follow a multivariate normal distribution with a Gaussian  
1310 covariance structure  $\Sigma$  given by:

$$\Sigma_{ij} = \sigma_{\varepsilon}^2 \times \left( v \times I_{ij} + (1 - v) \times \exp \left[ - \left( \frac{d_{ij}}{r} \right)^2 \right] \right) \tag{Eq. 5}$$

1312 where  $d_{ij}$  is the physical distance between plant  $i$  and  $j$  on the field,  $\sigma_{\varepsilon}^2$  is the overall residual  
1313 phenotype variance, the nugget  $v$  (between 0 and 1) determines the proportion of the  
1314 residual variance that is independently and identically distributed (*iid*) as opposed to  
1315 governed by spatial autocorrelation, the range  $r$  determines how fast the residual phenotype  
1316 correlation between plants drops when the distance between them increases, and  $I$  is an  
1317 identity matrix. The same model form was used to predict final yield phenotypes, e.g. total  
1318 seed weight, as a function of one of the phenotypes measured early in the growing season,  
1319 e.g. leaf 8 area (81 DAS). All parameters ( $\beta_0, \beta_1, \sigma_{\varepsilon}, n, r$ ) are estimated from the data by  
1320 Restricted Maximum Likelihood (ReML) estimation, implemented in the nlme package (126)  
1321 in R. In some cases the lme model didn't converge and a regular linear model (lm) was used  
1322 instead.  $p$ -values for the  $\beta_1$  coefficients were determined using Wald tests and adjusted for  
1323 multiple testing using the BH procedure (111).

1324

1325 For each of the 100 genes with the lowest BH-adjusted  $\beta_1$   $p$ -value for a given phenotype, a 9-  
1326 times repeated 10-fold cross-validation scheme was used to assess the gene's predictive  
1327 power (see section on multi-gene models for details). The median test  $R^2$  score across all 90  
1328 splits was used as a measure of prediction performance.

1329

1330 **Single-phenotype models.** The same linear mixed-effects (lme) modeling and cross-validation  
1331 strategy as used for the single-gene models was also used also to model spring phenotypes  
1332 as a function of autumnal leaf or rosette phenotypes. Leaf 6 and leaf 8 phenotypes and the  
1333 rosette area at 42 DAS were used as features for predicting all spring phenotypes. In a  
1334 separate analysis, also earlier rosette areas (14-42 DAS) were used as features, in order to  
1335 assess how the predictive power of the projected rosette area for yield phenotypes evolves  
1336 over time.

1337

1338 **Alternative single-gene models for ratio phenotypes.** For seeds per silique (on stem 1 or the  
1339 entire plant), the following alternative log-link model was fitted using the nlme package (126)  
1340 in R :

1341

$$1342 \quad \ln(E(\mathbf{n} \oslash \mathbf{d})) = \beta_0 + \beta_1 \mathbf{x} \quad (\text{Eq. 6})$$

1343

1344 where  $\oslash$  stands for the element-wise division of the numerator  $\mathbf{n}$ , a vector containing the  
1345 seed count stem 1 for all plants, by the denominator  $\mathbf{d}$ , a vector containing the silique count  
1346 stem 1 for all plants.  $\mathbf{x}$  is the expression profile of a given gene across plants. The numerator  
1347 is assumed to follow a normal distribution given the denominator  $\mathbf{d}$  and the gene expression  
1348 profile  $\mathbf{x}$  :

1349 
$$\mathbf{n} \sim \mathcal{N}(\mathbf{d} \cdot \exp(\beta_0 + \beta_1 \mathbf{x}), \Sigma) \quad (\text{Eq. 7})$$

1350

1351 Various error models  $\Sigma$  were tried out. For each gene,  $\Sigma$  is either a constant  $\sigma^2$  across all  
1352 plants estimated from the data, a spatially covarying error structure (using a Gaussian  
1353 covariance structure as for the other single-gene models, see above), a heteroscedastic error  
1354 structure with the error variance increasing linearly with the estimate, or a both spatially  
1355 covarying and heteroscedastic error structure. The parameters  $\beta_0, \beta_1, \sigma^2$  (and optionally the  
1356 nugget and range for spatial models) were estimated using the ‘ngls’ function in nlme.  $p$ -  
1357 values for the gene expression coefficients  $\beta_1$  were determined using Wald tests and adjusted  
1358 for multiple testing using the BH procedure (111).

1359

1360 A similar model was used for the branches per stem phenotype :

1361

1362 
$$\ln(E((c + \mathbf{n}) \oslash \mathbf{d})) = \beta_0 + \beta_1 \mathbf{x} \quad (\text{Eq. 8})$$

1363

1364 where  $\mathbf{n}$  is a vector containing the total branch count for all plants,  $\mathbf{d}$  is a vector containing  
1365 the stem count for all plants, and  $c$  is an extra offset introduced to account for the amount of  
1366 branches per stem decreasing with increasing numbers of stems on a plant.

1367

### 1368 **Multi-feature phenotype prediction models**

1369 **Multi-gene models.** Predictive models were made for each phenotype based on z-scored *rlog*  
1370 gene expression data, using either all genes or only transcription factors as potential features.  
1371 Random forest (127) and elastic net (128) models were constructed with scikit-learn v:0.23.2  
1372 (129) using a 10-fold cross-validation scheme. Model learning on the training data in each

1373 cross-validation split was done in two steps. First a feature selection model was used to select  
1374 promising features, and then a random forest or elastic net model was built on the selected  
1375 features. Three methods were used as alternatives for feature selection. The first feature  
1376 selection technique used was HSIC lasso (130) as implemented in the pyHSICLasso package  
1377 (131), which generally selected at most 200 genes. The second feature selection technique  
1378 was a filter selecting gene expression profiles exhibiting a significant Spearman correlation  
1379 with the phenotype of interest (BH-adjusted  $q \leq 0.01$  ; if no features survived this filter, the  
1380 threshold was set at  $p \leq 0.001$ ). The third feature selection technique was a filter selecting  
1381 genes with *rlog* gene expression  $> 0$  in at least half of the samples (median *rlog* gene  
1382 expression  $> 0$ ). Elastic net models were built using a fourfold inner cross-validation loop to  
1383 estimate the model hyperparameters. For random forest models, 1000 trees were estimated  
1384 ( $n\_estimators = 1000$ ) using bootstrapping (`bootstrap=True`), and  $\sqrt{n}$  features (with  $n$  the  
1385 total number of features) were considered when looking for the best split (`max_features =`  
1386 `"auto"`). The hyperparameters `'max_depth'` (the maximum number of nodes) and  
1387 `'min_samples_leaf'` (the minimal number of samples at each leaf node) were optimized using  
1388 a grid search with possible values (1, 2, 5, 10, 20, 50) and (1, 2, 5) for `'max_depth'` and  
1389 `'min_samples_leaf'`, respectively. Optimal hyperparameters were selected based on  
1390 generalization scores on out-of-bag (oob) samples (`oob_score=True`).

1391

1392 For each combination of phenotype, machine learning method and feature selection  
1393 technique, 9 repeats of the aforementioned 10-fold cross-validation scheme were performed,  
1394 giving rise to 90 train-test data splits in total. For each split, an out-of-bag (oob)  $R^2$  score was  
1395 computed from the predicted and observed phenotype values in the test set, and the median  
1396 oob  $R^2$  across all 90 splits (= median test  $R^2$ ) is reported as a measure of model prediction



1397 performance. Alternative  $R^2$  values and Pearson correlation (PCC) values were computed  
1398 based on the combined set of test predictions across all 10 splits of a cross-validation repeat.  
1399 The medians of those  $R^2$  and PCC values across the 9 cross-validation repeats for a given  
1400 model are reported as the median pooled  $R^2$  and median pooled PCC score of the model,  
1401 respectively.

1402

1403 For both elastic net and random forest models, genes of potential interest for a given  
1404 phenotype were ranked based on their median importance across the 90 cross-validation  
1405 splits of the model version with the highest median test  $R^2$  score (the difference between  
1406 model versions being the use of different feature selection techniques). For random forest  
1407 models, the gini importance of a gene was used as its importance score. For elastic net  
1408 models, the absolute value of a gene's estimated model coefficient was used.

1409

1410 ***Models on permuted datasets.*** For all continuous and high-count phenotypes and for both  
1411 the 'all genes' and 'transcription factors' feature sets, models were trained and tested on 90  
1412 datasets in which the phenotype values were permuted, using the same machine learning  
1413 method and feature selection technique as for the model with the best median test  $R^2$  score  
1414 on real data for the given phenotype and feature set. For each phenotype and feature set,  
1415 one model was trained per permuted dataset, using a single 90-10 train-test split mimicking  
1416 one fold of the cross validation setup used on real data.

1417

1418 ***Multi-phenotype models.*** For all phenotypes measured in spring, additional predictive  
1419 models were made based on z-scored data for 14 leaf and rosette phenotypes measured in  
1420 the preceding autumn. We used the same modeling approach as for the expression-based

1421 models (random forest and elastic net, 9 repeats of 10-fold nested cross-validation), except  
1422 that the feature selection step of the expression-based modeling protocol was skipped given  
1423 the low number of potential model features. In this respect, using elastic net models instead  
1424 of a simple linear regression framework is technically also unnecessary, but elastic nets were  
1425 used nevertheless to maximize comparability of the early phenotype-based and expression-  
1426 based modeling results.

1427

## 1428 **DECLARATIONS**

### 1429 **Ethics approval and consent to participate**

1430 Not applicable

1431

### 1432 **Consent for publication**

1433 Not applicable

1434

### 1435 **Availability of data and materials**

1436 The raw RNA-seq data generated in this study is available at ArrayExpress, experiment E-  
1437 MTAB-11904 (<https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-11904>). Data  
1438 analysis scripts are available from Zenodo (<https://zenodo.org/record/7072001>) and GitHub  
1439 ([https://github.com/MMichaelVdV/Brassica\\_segmentation](https://github.com/MMichaelVdV/Brassica_segmentation)).

1440

### 1441 **Competing interests**

1442 The authors declare that they have no competing interests.

1443

1444 **Funding**

1445 SDM is a fellow of the Research Foundation-Flanders (FWO, grant 1146319N). FWO had no  
1446 role in the design of the study, the collection, analysis, and interpretation of data or the  
1447 writing of the manuscript.

1448

1449 **Author contributions**

1450 SM conceived and supervised the study; TDS, PL, IR-R and SM designed the field trial ; HN, DI  
1451 and IR-R provided resources ; TDS and SM set up the field trial ; SDM, DC, TDS, KB and SM  
1452 phenotyped plants during the growing season ; SDM, DC, TDS, PL, JDB, KB, TVH, HN, IR-R and  
1453 SM sampled leaves for RNA-seq ; JDB prepared samples for RNA-seq ; SDM, DC, KB, HS and  
1454 SM harvested and phenotyped plants at the end of the field trial ; SDM, DC, TDS, PL, MVdV,  
1455 SH and SM analyzed data ; SDM, DC and SM wrote the manuscript with input from the other  
1456 authors. All authors read and approved the final manuscript.

1457

1458 **Acknowledgements**

1459 The authors thank Benjamin Wittkop and Rod Snowdon for providing *B. napus* Darmor seeds,  
1460 Luc van Gysegem, Thomas Vanderstocken and Katleen Sucaet for their help with setting up  
1461 and maintenance of the field trial, Dorota Herman and Kirin Demuynck for assistance with  
1462 leaf sampling, and Chris Pires for advice on sample prep for RNA sequencing.

1463

1464

## 1465 REFERENCES

- 1466 1. Stewart-Ornstein J, Weissman JS, El-Samad H. Cellular noise regulons underlie  
1467 fluctuations in *Saccharomyces cerevisiae*. *Mol Cell*. 2012;45(4):483-93.
- 1468 2. Dunlop MJ, Cox RS, 3rd, Levine JH, Murray RM, Elowitz MB. Regulatory activity  
1469 revealed by dynamic correlations in gene expression noise. *Nat Genet*.  
1470 2008;40(12):1493-8.
- 1471 3. Munsky B, Neuert G, van Oudenaarden A. Using gene expression noise to understand  
1472 gene regulation. *Science*. 2012;336(6078):183-7.
- 1473 4. Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, et al. The dynamics and  
1474 regulators of cell fate decisions are revealed by pseudotemporal ordering of single  
1475 cells. *Nat Biotechnol*. 2014;32(4):381-6.
- 1476 5. Moignard V, Woodhouse S, Haghverdi L, Lilly AJ, Tanaka Y, Wilkinson AC, et al.  
1477 Decoding the regulatory network of early blood development from single-cell gene  
1478 expression measurements. *Nat Biotechnol*. 2015;33(3):269-76.
- 1479 6. Moignard V, Macaulay IC, Swiers G, Buettner F, Schutte J, Calero-Nieto FJ, et al.  
1480 Characterization of transcriptional networks in blood stem and progenitor cells using  
1481 high-throughput single-cell gene expression analysis. *Nat Cell Biol*. 2013;15(4):363-72.
- 1482 7. Aibar S, González-Blas CB, Moerman T, Huynh-Thu VA, Imrichova H, Hulselmans G, et  
1483 al. SCENIC: single-cell regulatory network inference and clustering. *Nat Methods*.  
1484 2017;14(11):1083-6.
- 1485 8. van Dijk D, Sharma R, Nainys J, Yim K, Kathail P, Carr AJ, et al. Recovering gene  
1486 interactions from single-cell data using data diffusion. *Cell*. 2018;174(3):716-29 e27.
- 1487 9. Wagner A, Regev A, Yosef N. Revealing the vectors of cellular identity with single-cell  
1488 genomics. *Nat Biotechnol*. 2016;34(11):1145-60.
- 1489 10. Pratapa A, Jaliyal AP, Law JN, Bharadwaj A, Murali TM. Benchmarking algorithms for  
1490 gene regulatory network inference from single-cell transcriptomic data. *Nat Methods*.  
1491 2020;17(2):147-54.
- 1492 11. Bhosale R, Jewell JB, Hollunder J, Koo AJK, Vuylsteke M, Michoel T, et al. Predicting  
1493 gene function from uncontrolled expression variation among individual wild-type  
1494 *Arabidopsis* plants. *Plant Cell*. 2013;25(8):2865-77.
- 1495 12. Massonnet C, Vile D, Fabre J, Hannah MA, Caldana C, Lisec J, et al. Probing the  
1496 reproducibility of leaf growth and molecular phenotypes: a comparison of three  
1497 *Arabidopsis* accessions cultivated in ten laboratories. *Plant Physiol*. 2010;152(4):2142-  
1498 57.
- 1499 13. Cortijo S, Bhattarai M, Locke JCW, Ahnert SE. Co-expression networks from gene  
1500 expression variability between genetically identical seedlings can reveal novel  
1501 regulatory relationships. *Front Plant Sci*. 2020;11:599464.
- 1502 14. Atkinson NJ, Urwin PE. The interaction of plant biotic and abiotic stresses: from genes  
1503 to the field. *J Exp Bot*. 2012;63(10):3523-43.
- 1504 15. Barah P, Naika MBN, Jayavelu ND, Sowdhamini R, Shameer K, Bones AM.  
1505 Transcriptional regulatory networks in *Arabidopsis thaliana* during single and  
1506 combined stresses. *Nucleic Acids Res*. 2016;44(7):3147-64.
- 1507 16. Cabello JV, Lodeyro AF, Zurbriggen MD. Novel perspectives for the engineering of  
1508 abiotic stress tolerance in plants. *Curr Opin Biotech*. 2014;26:62-70.

- 1509 17. Davila Olivas NH, Kruijer W, Gort G, Wijnen CL, van Loon JJ, Dicke M. Genome-wide  
1510 association analysis reveals distinct genetic architectures for single and combined  
1511 stress responses in *Arabidopsis thaliana*. *New Phytol.* 2017;213(2):838-51.
- 1512 18. Johnson SM, Lim FL, Finkler A, Fromm H, Slabas AR, Knight MR. Transcriptomic analysis  
1513 of *Sorghum bicolor* responding to combined heat and drought stress. *BMC Genom.*  
1514 2014;15:456.
- 1515 19. Rasmussen S, Barah P, Suarez-Rodriguez MC, Bressendorff S, Friis P, Costantino P, et  
1516 al. Transcriptome responses to combinations of stresses in *Arabidopsis*. *Plant Physiol.*  
1517 2013;161(4):1783-94.
- 1518 20. Suzuki N, Rivero RM, Shulaev V, Blumwald E, Mittler R. Abiotic and biotic stress  
1519 combinations. *New Phytol.* 2014;203(1):32-43.
- 1520 21. Thoen MP, Davila Olivas NH, Kloth KJ, Coolen S, Huang PP, Aarts MG, et al. Genetic  
1521 architecture of plant stress resistance: multi-trait genome-wide association mapping.  
1522 *New Phytol.* 2017;213(3):1346-62.
- 1523 22. Poorter H, Fiorani F, Pieruschka R, Wojciechowski T, van der Putten WH, Kleyer M, et  
1524 al. Pampered inside, pestered outside? Differences and similarities between plants  
1525 growing in controlled conditions and in the field. *New Phytol.* 2016;212(4):838-55.
- 1526 23. Mittler R. Abiotic stress, the field environment and stress combination. *Trends Plant*  
1527 *Sci.* 2006;11(1):15-9.
- 1528 24. Nelissen H, Moloney M, Inzé D. Translational research: from pot to plot. *Plant*  
1529 *Biotechnol J.* 2014;12(3):277-85.
- 1530 25. Nelissen H, Sprenger H, Demuyneck K, De Block J, Van Hautegeem T, De Vliegheer A, et al.  
1531 From laboratory to field: yield stability and shade avoidance genes are massively  
1532 differentially expressed in the field. *Plant Biotechnol J.* 2020;18(5):1112-4.
- 1533 26. Oh SJ, Kim YS, Kwon CW, Park HK, Jeong JS, Kim JK. Overexpression of the  
1534 transcription factor *AP37* in rice improves grain yield under drought conditions. *Plant*  
1535 *Physiol.* 2009;150(3):1368-79.
- 1536 27. Maathuis MH, Colombo D, Kalisch M, Bühlmann P. Predicting causal effects in large-  
1537 scale systems from observational data. *Nat Methods.* 2010;7(4):247-8.
- 1538 28. Cruz DF, De Meyer S, Ampe J, Sprenger H, Herman D, Van Hautegeem T, et al. Using  
1539 single-plant-omics in the field to link maize genes to functions and phenotypes. *Mol*  
1540 *Syst Biol.* 2020;16(12):e9667.
- 1541 29. FAO. World Food and Agriculture - Statistical Yearbook 2021. Rome; 2021.
- 1542 30. Mann HB, Whitney DR. On a test of whether one of two random variables is  
1543 stochastically larger than the other. *Ann Math Statist.* 1947;18(1):50-60.
- 1544 31. Cortijo S, Aydin Z, Ahnert S, Locke JC. Widespread inter-individual gene expression  
1545 variability in *Arabidopsis thaliana*. *Mol Syst Biol.* 2019;15(1):e8591.
- 1546 32. Efron B. Tweedie's formula and selection bias. *J Am Stat Assoc.* 2011;106(496):1602-  
1547 14.
- 1548 33. Aoyama T, Dong CH, Wu Y, Carabelli M, Sessa G, Ruberti I, et al. Ectopic expression of  
1549 the *Arabidopsis* transcriptional activator *Athb-1* alters leaf cell fate in tobacco. *Plant*  
1550 *Cell.* 1995;7(11):1773-85.
- 1551 34. Ribone PA, Capella M, Arce AL, Chan RL. A uORF represses the transcription factor  
1552 *AtHB1* in aerial tissues to avoid a deleterious phenotype. *Plant Physiol.*  
1553 2017;175(3):1238-53.

- 1554 35. Miguel VN, Manavella PA, Chan RL, Capella MA. The AtHB1 transcription factor  
1555 controls the miR164-*CUC2* regulatory node to modulate leaf development. *Plant Cell*  
1556 *Physiol.* 2020;61(3):659-70.
- 1557 36. Wang Y, Henriksson E, Söderman E, Henriksson KN, Sundberg E, Engström P. The  
1558 *Arabidopsis* homeobox gene, *ATHB16*, regulates leaf development and the sensitivity  
1559 to photoperiod in *Arabidopsis*. *Dev Biol.* 2003;264(1):228-39.
- 1560 37. Otsuga D, DeGuzman B, Prigge MJ, Drews GN, Clark SE. *REVOLUTA* regulates meristem  
1561 initiation at lateral positions. *Plant J.* 2001;25(2):223-36.
- 1562 38. Byrne ME. Shoot meristem function and leaf polarity: the role of class III HD-ZIP genes.  
1563 *Plos Genet.* 2006;2(6):e89.
- 1564 39. Talbert PB, Adler HT, Parks DW, Comai L. The *REVOLUTA* gene is necessary for apical  
1565 meristem development and for limiting cell divisions in the leaves and stems of  
1566 *Arabidopsis thaliana*. *Development.* 1995;121(9):2723-35.
- 1567 40. Prigge MJ, Otsuga D, Alonso JM, Ecker JR, Drews GN, Clark SE. Class III homeodomain-  
1568 leucine zipper gene family members have overlapping, antagonistic, and distinct roles  
1569 in *Arabidopsis* development. *Plant Cell.* 2005;17(1):61-76.
- 1570 41. Ochando I, Jover-Gil S, Ripoll JJ, Candela H, Vera A, Ponce MR, et al. Mutations in the  
1571 microRNA complementarity site of the *INCURVATA4* gene perturb meristem function  
1572 and adaxialize lateral organs in *Arabidopsis*. *Plant Physiol.* 2006;141(2):607-19.
- 1573 42. Baima S, Possenti M, Matteucci A, Wisman E, Altamura MM, Ruberti I, et al. The  
1574 *Arabidopsis* ATHB-8 HD-zip protein acts as a differentiation-promoting transcription  
1575 factor of the vascular meristems. *Plant Physiol.* 2001;126(2):643-55.
- 1576 43. Gardiner J, Donner TJ, Scarpella E. Simultaneous activation of *SHR* and *ATHB8*  
1577 expression defines switch to preprocambial cell state in *Arabidopsis* leaf development.  
1578 *Dev Dyn.* 2011;240(1):261-70.
- 1579 44. Kim J, Jung JH, Reyes JL, Kim YS, Kim SY, Chung KS, et al. microRNA-directed cleavage  
1580 of *ATHB15* mRNA regulates vascular development in *Arabidopsis* inflorescence stems.  
1581 *Plant J.* 2005;42(1):84-94.
- 1582 45. Bou-Torrent J, Salla-Martret M, Brandt R, Musielak T, Palauqui JC, Martínez-García JF,  
1583 et al. ATHB4 and HAT3, two class II HD-ZIP transcription factors, control leaf  
1584 development in *Arabidopsis*. *Plant Signal Behav.* 2012;7(11):1382-7.
- 1585 46. Brandt R, Salla-Martret M, Bou-Torrent J, Musielak T, Stahl M, Lanz C, et al. Genome-  
1586 wide binding-site analysis of *REVOLUTA* reveals a link between leaf patterning and  
1587 light-mediated growth responses. *Plant J.* 2012;72(1):31-42.
- 1588 47. Zhang Z, Runions A, Mentink RA, Kierzkowski D, Karady M, Hashemi B, et al. A  
1589 *WOX/auxin* biosynthesis module controls growth to shape leaf form. *Curr Biol.*  
1590 2020;30(24):4857-68 e6.
- 1591 48. Zhang Z, Tucker E, Hermann M, Laux T. A molecular framework for the embryonic  
1592 initiation of shoot meristem stem cells. *Dev Cell.* 2017;40(3):264-77 e4.
- 1593 49. Nakata M, Matsumoto N, Tsugeki R, Rikirsch E, Laux T, Okada K. Roles of the middle  
1594 domain-specific *WUSCHEL-RELATED HOMEODOMAIN* genes in early development of leaves  
1595 in *Arabidopsis*. *Plant Cell.* 2012;24(2):519-35.
- 1596 50. Cominelli E, Galbiati M, Vavasseur A, Conti L, Sala T, Vuylsteke M, et al. A guard-cell-  
1597 specific MYB transcription factor regulates stomatal movements and plant drought  
1598 tolerance. *Curr Biol.* 2005;15(13):1196-200.

- 1599 51. Wang L, Hua D, He J, Duan Y, Chen Z, Hong X, et al. *Auxin Response Factor2 (ARF2)* and  
1600 its regulated homeodomain gene *HB33* mediate abscisic acid response in *Arabidopsis*.  
1601 Plos Genet. 2011;7(7):e1002172.
- 1602 52. Hong SY, Kim OK, Kim SG, Yang MS, Park CM. Nuclear import and DNA binding of the  
1603 ZHD5 transcription factor is modulated by a competitive peptide inhibitor in  
1604 *Arabidopsis*. J Biol Chem. 2011;286(2):1659-68.
- 1605 53. Aharoni A, Dixit S, Jetter R, Thoenes E, van Arkel G, Pereira A. The SHINE clade of AP2  
1606 domain transcription factors activates wax biosynthesis, alters cuticle properties, and  
1607 confers drought tolerance when overexpressed in *Arabidopsis*. Plant Cell.  
1608 2004;16(9):2463-80.
- 1609 54. Zhao H, Wu D, Kong F, Lin K, Zhang H, Li G. The *Arabidopsis thaliana* Nuclear Factor Y  
1610 transcription factors. Front Plant Sci. 2016;7:2045.
- 1611 55. Zhang M, Hu X, Zhu M, Xu M, Wang L. Transcription factors NF-YA2 and NF-YA10  
1612 regulate leaf growth via auxin signaling in *Arabidopsis*. Sci Rep. 2017;7(1):1395.
- 1613 56. Li WX, Oono Y, Zhu J, He XJ, Wu JM, Iida K, et al. The *Arabidopsis* NFYA5 transcription  
1614 factor is regulated transcriptionally and posttranscriptionally to promote drought  
1615 resistance. Plant Cell. 2008;20(8):2238-51.
- 1616 57. Zhao H, Lin K, Ma L, Chen Q, Gan S, Li G. *Arabidopsis* NUCLEAR FACTOR Y A8 inhibits  
1617 the juvenile-to-adult transition by activating transcription of *MIR156s*. J Exp Bot.  
1618 2020;71(16):4890-902.
- 1619 58. Wenkel S, Turck F, Singer K, Gissot L, Le Gourrierc J, Samach A, et al. CONSTANS and  
1620 the CCAAT box binding complex share a functionally important domain and interact to  
1621 regulate flowering of *Arabidopsis*. Plant Cell. 2006;18(11):2971-84.
- 1622 59. Leyva-González MA, Ibarra-Laclette E, Cruz-Ramírez A, Herrera-Estrella L. Functional  
1623 and transcriptome analysis reveals an acclimatization strategy for abiotic stress  
1624 tolerance mediated by *Arabidopsis* NF-YA family members. Plos One.  
1625 2012;7(10):e48138.
- 1626 60. Xu MY, Zhang L, Li WW, Hu XL, Wang MB, Fan YL, et al. Stress-induced early flowering  
1627 is mediated by miR169 in *Arabidopsis thaliana*. J Exp Bot. 2014;65(1):89-101.
- 1628 61. Siriwardana CL, Gnesutta N, Kumimoto RW, Jones DS, Myers ZA, Mantovani R, et al.  
1629 NUCLEAR FACTOR Y, subunit A (NF-YA) proteins positively regulate flowering and act  
1630 through *FLOWERING LOCUS T*. Plos Genet. 2016;12(12):e1006496.
- 1631 62. Yan Y, Shen L, Chen Y, Bao S, Thong Z, Yu H. A MYB-domain protein EFM mediates  
1632 flowering responses to environmental cues in *Arabidopsis*. Dev Cell. 2014;30(4):437-  
1633 48.
- 1634 63. Jung JH, Lee S, Yun J, Lee M, Park CM. The miR172 target TOE3 represses *AGAMOUS*  
1635 expression during *Arabidopsis* floral patterning. Plant Sci. 2014;215-216:29-38.
- 1636 64. Schiessl S. Regulation and subfunctionalization of flowering time genes in the  
1637 allotetraploid oil crop *Brassica napus*. Front Plant Sci. 2020;11:605155.
- 1638 65. Tudor EH, Jones DM, He Z, Bancroft I, Trick M, Wells R, et al. QTL-seq identifies  
1639 *BnaFT.A02* and *BnaFLC.A02* as candidates for variation in vernalization requirement  
1640 and response in winter oilseed rape (*Brassica napus*). Plant Biotechnol J.  
1641 2020;18(12):2466-81.
- 1642 66. Torti S, Fornara F, Vincent C, Andrés F, Nordström K, Göbel U, et al. Analysis of the  
1643 *Arabidopsis* shoot meristem transcriptome during floral transition identifies distinct  
1644 regulatory patterns and a leucine-rich repeat protein that promotes flowering. Plant  
1645 Cell. 2012;24(2):444-62.

- 1646 67. Melzer S, Lens F, Gennen J, Vanneste S, Rohde A, Beeckman T. Flowering-time genes  
1647 modulate meristem determinacy and growth form in *Arabidopsis thaliana*. Nat Genet.  
1648 2008;40(12):1489-92.
- 1649 68. O'Maoiléidigh DS, van Driel AD, Singh A, Sang Q, Le Bec N, Vincent C, et al. Systematic  
1650 analyses of the *MIR172* family members of *Arabidopsis* define their distinct roles in  
1651 regulation of *APETALA2* during floral transition. PLoS Biol. 2021;19(2):e3001043.
- 1652 69. Balanzà V, Martínez-Fernández I, Sato S, Yanofsky MF, Kaufmann K, Angenent GC, et  
1653 al. Genetic control of meristem arrest and life span in *Arabidopsis* by a *FRUITFULL*-  
1654 *APETALA2* pathway. Nat Commun. 2018;9(1):565.
- 1655 70. Zúñiga-Mayo VM, Marsch-Martínez N, de Folter S. JAIBA, a class-II HD-ZIP transcription  
1656 factor involved in the regulation of meristematic activity, and important for correct  
1657 gynoecium and fruit development in *Arabidopsis*. Plant J. 2012;71(2):314-26.
- 1658 71. Kamata N, Okada H, Komeda Y, Takahashi T. Mutations in epidermis-specific HD-ZIP IV  
1659 genes affect floral organ identity in *Arabidopsis thaliana*. Plant J. 2013;75(3):430-40.
- 1660 72. Wang W, Sijacic P, Xu P, Lian H, Liu Z. *Arabidopsis* TSO1 and MYB3R1 form a regulatory  
1661 module to coordinate cell proliferation with differentiation in shoot and root. Proc  
1662 Natl Acad Sci U S A. 2018;115(13):E3045-E54.
- 1663 73. Ranftl QL, Bastakis E, Klermund C, Schwechheimer C. LLM-domain containing B-GATA  
1664 factors control different aspects of cytokinin-regulated development in *Arabidopsis*  
1665 *thaliana*. Plant Physiol. 2016;170(4):2295-311.
- 1666 74. Lozano-Sotomayor P, Chávez Montes RA, Silvestre-Vañó M, Herrera-Ubaldo H, Greco  
1667 R, Pablo-Villa J, et al. Altered expression of the bZIP transcription factor DRINK ME  
1668 affects growth and reproductive development in *Arabidopsis thaliana*. Plant J.  
1669 2016;88(3):437-51.
- 1670 75. Gailloch C, Jamge S, van der Wal F, Angenent G, Immink R, Lohmann JU. A molecular  
1671 network for functional versatility of HECATE transcription factors. Plant J.  
1672 2018;95(1):57-70.
- 1673 76. Meng WJ, Cheng ZJ, Sang YL, Zhang MM, Rong XF, Wang ZW, et al. Type-B  
1674 ARABIDOPSIS RESPONSE REGULATORS specify the shoot stem cell niche by dual  
1675 regulation of *WUSCHEL*. Plant Cell. 2017;29(6):1357-72.
- 1676 77. Azodi CB, Pardo J, VanBuren R, de Los Campos G, Shiu SH. Transcriptome-based  
1677 prediction of complex traits in maize. Plant Cell. 2020;32(1):139-51.
- 1678 78. O'Neill CM, Lu X, Calderwood A, Tudor EH, Robinson P, Wells R, et al. Vernalization  
1679 and floral transition in autumn drive winter annual life history in oilseed rape. Curr  
1680 Biol. 2019;29(24):4300-6 e2.
- 1681 79. Diepenbrock W. Yield analysis of winter oilseed rape (*Brassica napus* L.): a review.  
1682 Field Crop Res. 2000;67:35-49.
- 1683 80. Searle I, He Y, Turck F, Vincent C, Fornara F, Kröber S, et al. The transcription factor  
1684 FLC confers a flowering response to vernalization by repressing meristem competence  
1685 and systemic signaling in *Arabidopsis*. Genes & development. 2006;20(7):898-912.
- 1686 81. Turnbull C. Long-distance regulation of flowering time. J Exp Bot. 2011;62(13):4399-  
1687 413.
- 1688 82. Lang A. Physiology of flower initiation. Differentiation and Development Encyclopedia  
1689 of Plant Physiology. 15. Berlin, Heidelberg: Springer; 1965. p. 1380–536.
- 1690 83. Klepikova AV, Kasianov AS, Gerasimov ES, Logacheva MD, Penin AA. A high resolution  
1691 map of the *Arabidopsis thaliana* developmental transcriptome based on RNA-seq  
1692 profiling. Plant J. 2016;88(6):1058-70.



- 1693 84. Mendham NJ, Scott RK. Limiting effect of plant size at inflorescence initiation on  
1694 subsequent growth and yield of oilseed rape (*Brassica napus*). J Agr Sci.  
1695 1975;84(Jun):487-502.
- 1696 85. Gan YT, Stobbe EH, Moes J. Relative date of wheat seedling emergence and its impact  
1697 on grain yield. Crop Sci. 1992;32(5):1275-81.
- 1698 86. Liu T, Chen W, Li F, Wu W, Sun C, Ding J, et al. Characterization of the 3D structure of a  
1699 cultivated land surface and its influence on wheat seedlings growth using Kinect. Sci  
1700 Rep. 2017;7(1):3927.
- 1701 87. Soltani A, Robertson MJ, Torabi B, Yousefi-Daz M, Sarparast R. Modelling seedling  
1702 emergence in chickpea as influenced by temperature and sowing depth. Agr Forest  
1703 Meteorol. 2006;138(1-4):156-67.
- 1704 88. Forcella F, Benecch Arnold, R.L., Sanchez, R., Ghera, C.M. Modeling seedling  
1705 emergence. Field Crop Res. 2000;67(2):123-39.
- 1706 89. Varma V, Iyengar SB, Sankaran M. Effects of nutrient addition and soil drainage on  
1707 germination of N-fixing and non-N-fixing tropical dry forest tree species. Plant Ecol.  
1708 2016;217(8):1043-54.
- 1709 90. Chalhoub B, Denoeud F, Liu SY, Parkin IAP, Tang HB, Wang XY, et al. Early allopolyploid  
1710 evolution in the post-Neolithic *Brassica napus* oilseed genome. Science.  
1711 2014;345(6199):950-3.
- 1712 91. Yin X, Goudriaan J, Lantinga EA, Vos J, Spiertz HJ. A flexible sigmoid function of  
1713 determinate growth. Ann Bot. 2003;91(3):361-71.
- 1714 92. Habekotté B. A model of the phenological development of winter oilseed rape  
1715 (*Brassica napus* L.). Field Crop Res. 1997;54(2-3):127-36.
- 1716 93. Ronneberger O, Fischer P, Brox T. U-Net: convolutional networks for biomedical image  
1717 segmentation. arXiv. 2015;doi:10.48550/arXiv.1505.04597.
- 1718 94. Dutta A, Zisserman A. The VIA annotation software for images, audio and video.  
1719 Proceedings of the 27th ACM International Conference on Multimedia.  
1720 2019;doi:10.1145/3343031.3350535.
- 1721 95. Kingma DP, Ba J. Adam: a method for stochastic optimization. arXiv.  
1722 2014;doi:10.48550/arXiv.1412.6980.
- 1723 96. Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: an  
1724 imperative style, high-performance deep learning library. In: H. Wallach HL, A.  
1725 Beygelzimer, F. d'Alché-Buc, E. Fox and R. Garnett, editor. Advances in Neural  
1726 Information Processing Systems 32: Curran Associates, Inc.; 2019. p. 8024-35.
- 1727 97. Schneider CA, Rasband WS, Eliceiri KW. NIH Image to ImageJ: 25 years of image  
1728 analysis. Nat Methods. 2012;9(7):671-5.
- 1729 98. Huang G, Liu Z, Weinberger KQ. Densely connected convolutional networks. arXiv.  
1730 2016;doi:10.48550/arXiv.1608.06993.
- 1731 99. van der Walt S, Schönberger JL, Nunez-Iglesias J, Boulogne F, Warner JD, Yager N, et al.  
1732 scikit-image: image processing in Python. PeerJ. 2014;2:e453.
- 1733 100. Bivens NJ, Zhou M. RNA-Seq library construction methods for transcriptome analysis.  
1734 Curr Protoc Plant Biol. 2016;1(1):197-215.
- 1735 101. Goecks J, Nekrutenko A, Taylor J. Galaxy: a comprehensive approach for supporting  
1736 accessible, reproducible, and transparent computational research in the life sciences.  
1737 Genome Biol. 2010;11(8):R86.
- 1738 102. Andrews S. FastQC: a quality control tool for high throughput sequence data.  
1739 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>. 2010.

- 1740 103. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence  
1741 data. *Bioinformatics*. 2014;30(15):2114-20.
- 1742 104. Pertea M, Kim D, Pertea GM, Leek JT, Salzberg SL. Transcript-level expression analysis  
1743 of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat Protoc*.  
1744 2016;11(9):1650-67.
- 1745 105. Anders S, Pyl PT, Huber W. HTSeq—a Python framework to work with high-throughput  
1746 sequencing data. *Bioinformatics*. 2015;31(2):166-9.
- 1747 106. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for  
1748 RNA-seq data with DESeq2. *Genome Biol*. 2014;15(12):550.
- 1749 107. Tello D, Gil J, Loaiza CD, Riascos JJ, Cardozo N, Duitama J. NGSEP3: accurate variant  
1750 calling across species and sequencing protocols. *Bioinformatics*. 2019;35(22):4716-23.
- 1751 108. Browning BL, Zhou Y, Browning SR. A one-penny imputed genome from next-  
1752 generation reference panels. *Am J Hum Genet*. 2018;103(3):338-48.
- 1753 109. Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES. TASSEL:  
1754 software for association mapping of complex traits in diverse samples. *Bioinformatics*.  
1755 2007;23(19):2633-5.
- 1756 110. Rambaut A. FigTree. <http://treebioedacuk/software/figtree/>. 2016.
- 1757 111. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful  
1758 approach to multiple testing. *J R Stat Soc Ser B Methodol*. 1995;57(1):289-300.
- 1759 112. Rey SJ, Anselin L. PySAL: A Python library of spatial analytical methods. *Handbook of*  
1760 *applied spatial analysis*: Springer; 2010. p. 175-93.
- 1761 113. Stacklies W, Redestig H, Scholz M, Walther D, Selbig J. pcaMethods—a bioconductor  
1762 package providing PCA methods for incomplete data. *Bioinformatics*. 2007;23(9):1164-  
1763 7.
- 1764 114. Öner M, Kocakoç ID. JMASM 49: A compilation of some popular goodness of fit tests  
1765 for normal distribution: their algorithms and MATLAB codes (MATLAB). *J Mod Appl*  
1766 *Stat Meth*. 2017;16(2):547-75.
- 1767 115. Brennecke P, Anders S, Kim JK, Kolodziejczyk AA, Zhang X, Proserpio V, et al.  
1768 Accounting for technical noise in single-cell RNA-seq experiments. *Nat Methods*.  
1769 2013;10(11):1093-5.
- 1770 116. Andrews TS, Hemberg M. M3Drop: dropout-based feature selection for scRNASeq.  
1771 *Bioinformatics*. 2019;35(16):2865-7.
- 1772 117. Bucchini F, Del Cortona A, Kreft Ł, Botzki A, Van Bel M, Vandepoele K. TRAPID 2.0: a  
1773 web application for taxonomic and functional analysis of de novo transcriptomes.  
1774 *Nucleic Acids Res*. 2021;49(17):e101.
- 1775 118. Pertea G, Pertea M. GFF Utilities: GffRead and GffCompare [version 2; peer review: 3  
1776 approved]. *F1000Research*. 2020;9:304.
- 1777 119. Van Bel M, Diels T, Vancaester E, Kreft Ł, Botzki A, Van de Peer Y, et al. PLAZA 4.0: an  
1778 integrative resource for functional, evolutionary and comparative plant genomics.  
1779 *Nucleic Acids Res*. 2018;46(D1):D1190-D6.
- 1780 120. Maere S, Heymans K, Kuiper M. BiNGO: a Cytoscape plugin to assess  
1781 overrepresentation of gene ontology categories in biological networks. *Bioinformatics*.  
1782 2005;21(16):3448-9.
- 1783 121. Sun F, Fan G, Hu Q, Zhou Y, Guan M, Tong C, et al. The high-quality genome of *Brassica*  
1784 *napus* cultivar 'ZS11' reveals the introgression history in semi-winter morphotype.  
1785 *Plant J*. 2017;92(3):452-68.

- 1786 122. Emms DM, Kelly S. OrthoFinder: phylogenetic orthology inference for comparative  
1787 genomics. *Genome Biol.* 2019;20(1):238.
- 1788 123. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND.  
1789 *Nat Methods.* 2015;12(1):59-60.
- 1790 124. Proost S, Fostier J, De Witte D, Dhoedt B, Demeester P, Van de Peer Y, et al. i-ADHoRe  
1791 3.0--fast and sensitive detection of genomic homology in extremely large data sets.  
1792 *Nucleic Acids Res.* 2012;40(2):e11.
- 1793 125. Pinheiro JC, Bates DM. Linear mixed-effects models: basic concepts and examples.  
1794 *Mixed-effects models in S and S-Plus.* New York, NY: Springer; 2000. p. 3-56.
- 1795 126. Pinheiro J, Bates D, DebRoy S, Sarkar D, R Core Team. nlme: Linear and nonlinear  
1796 mixed effects models. 2019.
- 1797 127. Breiman L. Random forests. *Mach Learn.* 2001;45(1):5-32.
- 1798 128. Zou H, Hastie T. Regularization and variable selection via the elastic net. *J Roy Stat Soc*  
1799 *B.* 2005;67:301-20.
- 1800 129. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn:  
1801 machine learning in Python. *J Mach Learn Res.* 2011;12:2825-30.
- 1802 130. Yamada M, Jitkrittum W, Sigal L, Xing EP, Sugiyama M. High-dimensional feature  
1803 selection by feature-wise kernelized Lasso. *Neural Comput.* 2014;26(1):185-207.
- 1804 131. Climente-González H, Azencott CA, Kaski S, Yamada M. Block HSIC Lasso: model-free  
1805 biomarker detection for ultra-high dimensional data. *Bioinformatics.*  
1806 2019;35(14):i427-i35.

1807