

1 **Title**

2 The *Eruca sativa* genome and transcriptome: A targeted analysis of sulfur metabolism and
3 glucosinolate biosynthesis pre and postharvest

4
5 **Authors**

6 Luke Bell ^{1*}, Martin Chadwick ², Manik Puranik ², Richard Tudor ³, Lisa Methven ², Sue
7 Kennedy ³, Carol Wagstaff ²

8
9 **Affiliations**

10 ¹ School of Agriculture, Policy & Development, PO Box 237, University of Reading,
11 Whiteknights, Reading, Berkshire. RG6 6AR. UK.

12
13 ² School of Chemistry Food & Pharmacy, PO Box 226, University of Reading,
14 Whiteknights, Reading, Berkshire. RG6 6AP. UK.

15
16 ³ Elsoms Seeds Ltd., Pinchbeck Road, Spalding, Lincolnshire. PE11 1QG. UK.

17
18 * luke.bell@reading.ac.uk

19
20 **Abstract**

21 Rocket (*Eruca sativa*) is a source of health-related metabolites called glucosinolates
22 (GSLs) and isothiocyanates (ITCs) but little is known of the genetic and transcriptomic
23 mechanisms responsible for regulating pre and postharvest accumulations. We present the first *de*
24 *novo* reference genome assembly and annotation, with ontogenic and postharvest transcriptome
25 data relating to sulfur assimilation, transport, and utilization. Diverse gene expression patterns
26 related to sulfur metabolism and GSL biosynthesis are present between inbred lines of rocket. A
27 clear pattern of differential expression determines GSL abundance and the formation of hydrolysis
28 products. One breeding line sustained GSL accumulation and hydrolysis product formation
29 throughout storage. Copies of *MYB28*, *SLIM1*, *SDII* and *ESMI* orthologs have increased and
30 differential expression postharvest, and are associated with GSLs and hydrolysis product

31 formation. Two glucosinolate transporter gene orthologs (*GTR2*) were found to be associated with
32 increased GSL accumulations.

33

34 **Introduction**

35 Sulfur (S) is a critical macronutrient that plants require for growth and development ¹.
36 Sulfate (SO₄²⁻) is utilized as a primary means of synthesizing numerous S-containing metabolites,
37 such as amino acids (cysteine and methionine), glutathione (GSH), and glucosinolates (GSLs) ².
38 GSL compounds are present in species of the order Brassicales, and are abundant in many
39 vegetables and condiments worldwide, such as rapeseed (*Brassica napus*), Chinese cabbage
40 (*Brassica rapa*), cabbage (*Brassica oleracea* var. *capitata*), and broccoli (*B. oleracea* var. *italica*)
41 ³. GSLs are also found in the leafy vegetable *Eruca sativa* (“salad” rocket), which has gained
42 significant popularity amongst consumers over the last ten years ⁴. Rocket salad is known for its
43 distinctive flavour, aroma, and pungency, and can be eaten raw without the need for cooking ⁵.

44 GSLs are synthesised as part of plant defense mechanisms against pests and diseases ⁶, and
45 can also act as important S storage molecules ¹. Compounds such as glucosativin (4-mercaptobutyl
46 GSL; GSV) and glucorucolamine (4-cystein-S-yl)butyl GSL; GRL) are unique to the genera *Eruca*
47 and *Diplotaxis* (‘wild’ rocket) ⁷. GSV can exist in a dimer form (dimeric 4-mercaptobutyl GSL;
48 DMB), and diglucothiobeinin (4-(β-D-glucopyranosyldisulfanyl)butyl GSL; DGTB) is a unique
49 GSL dimer of these species ⁸.

50 Aliphatic GSLs are regulated by *MYB28*, *MYB29*, and *MYB76* transcription factors (TFs),
51 and indolic GSLs by *MYB34*, *MYB51*, and *MYB122* ². These MYBs are in turn regulated by basic
52 helix-loop-helix (bHLH) transcription factors such as *MYC2*, which are involved in plant defense
53 response ⁹. Other transcriptional regulators, such as *SLIM1* (*SULFUR LIMITATION 1*) and *SDII*
54 (*SULFUR DEFICIENCY INDUCED 1*) also interact with MYB transcription factors to regulate
55 the use and efficiency of sulfur within the plant. As GSLs contain significant amounts of sulfur
56 (up to 30% of total plant S-content) the synthesis and catalysis of these compounds is crucial in
57 times of stress (Figure 1) ^{10,11}.

58 GSLs themselves are not bioactive, and are hydrolysed by myrosinase enzymes (TGGs)
59 when tissue damage takes place. They form numerous breakdown products including
60 isothiocyanates (ITCs), which are of foremost interest for their anticarcinogenic effects in humans
61 ¹². Salad rocket produces the ITC sulforaphane (SF; a breakdown product of glucoraphanin; 4-

62 methylsulfinylbutyl GSL, GRA), which has been well documented for its potent anticarcinogenic
63 properties¹³. SF is abundant in broccoli, however its hydrolysis from GRA is often inhibited or
64 prevented due to high cooking temperatures employed by consumers, which denatures myrosinase
65 at temperatures >65°C¹⁴.

66 A previous study by Bell et al.¹⁵ observed that both GSL and ITC concentrations increased
67 significantly in rocket salad post-processing, but that this varied according to cultivar. The authors
68 proposed that in response to the harvesting and washing process, stress responses within leaf
69 tissues were initiated, leading to the increase in synthesis of GSLs and subsequent hydrolysis into
70 ITCs. Sugar content, by comparison, showed little dynamic change and little reduction in the same
71 samples, which could have implications for sensory perceptions and consumer acceptance⁵.

72 We present a *de novo* *E. sativa* reference genome sequence, and report on the specific
73 effects harvest, wash treatment, and postharvest storage have on GSL biosynthesis and sulfur
74 metabolism gene expression through RNA sequencing (RNAseq) in three elite inbred lines. We
75 also present evidence of transcriptomic changes between first and second cuts of rocket plants, and
76 how this in turn leads to elevated concentrations of both GSLs and ITCs. We hypothesised that
77 each rocket line would vary in their ability to retain and synthesize GSLs post washing and during
78 shelf life cold storage, as well as vary in their relative abundances between first and second cuts.

79

80 **Results**

81 ***E. sativa* genome assembly and annotation**

82 *De novo* reference genome sequence was produced by interleaving Illumina MiSeq and
83 HiSeq2500 sequence data (Illumina Inc., San Diego, CA, USA). An elite breeding line (designated
84 **C**) was selected for PCR-free paired-end sequencing and long mate-pair end sequencing and
85 assembled into 49,933 contigs (≥ 500 bp). The resulting assembly was ~851 Mb in size
86 (Supplementary Table S1).

87 Transposable elements (TEs) within the *E. sativa* genome comprise 66.3% of its content.
88 The majority of TEs are long terminal repeat (LTR) retrotransposons (37.3%), with long
89 interspersed nuclear elements (LINEs; 3.3%) and short interspersed nuclear elements (SINEs;
90 0.3%) having lower relative abundance. 18.2% of all TEs identified were of unknown
91 classification (Supplementary Table S2).

92 A total of 45,438 protein-coding genes were identified within the assembly, with an
93 average length of 1,889.6 bp, and an average of 4.76 exons per gene. This genome size is smaller
94 than that predicted for radish (*Raphanus sativus*), and larger than *Arabidopsis lyrata*
95 (Supplementary Table S3), and is consistent with what is known of Brassicales phylogeny¹⁶.
96 98.3% of predicted genes were found to have homology with other plant species (Figure 2b,
97 Supplementary Table S4).

98

99 **RNAseq analysis of *E. sativa* plants**

100 RNA sequencing and bioinformatics was conducted on 18 plants from three elite inbred
101 lines designated **A**, **B**, and **C**; giving a total of 54 plant samples. Time points corresponded to three
102 harvest times ('early harvest' at 22 days after sowing, **EH**; 'harvest' at 30 days after sowing, **H**;
103 'second cut', **SC**; leaves harvested from the same **H** plants 43 days after sowing), and three
104 consecutive postharvest time points (harvested at 30 days after sowing and designated: 'pre-wash',
105 **PW**; 'day 0' of shelf life, 1 day post wash, **D0**; and 'day 7' of shelf life, **D7**). See Supplementary
106 Figure S1 for a schematic of the experimental design.

107 After sample QC and clean-up over 2.6 billion clean paired-end reads were produced,
108 averaging ~49 million reads per sample. Q20 (<1% error rate) averaged 96.3%, Q30 (<0.1% error
109 rate) averaged 90.8%, and GC content ranged from 44.5% to 47.4%.

110

111 **Global differential gene expression**

112 The total numbers of differentially expressed genes (DEGs) for line **A**, **B** and **C** are
113 presented in Figure 2c, 2d and 2e, respectively. Few significant DEGs were observed between **EH**
114 and **H** samples for lines **A** and **B** (<333; Figure 2c and 2d), whereas they were observed at a higher
115 rate in **C** (2,234; Figure 2e). This indicates a high level of plasticity of **C** across growth stages.

116 This trend was reversed at **PW**, where 180 DEGs were observed compared to **H** in **C**, and
117 1,343 were observed in **A**. During shelf life (**D0** and **D7**) **C** expressed a greater number of DEGs
118 compared to **H**, than **A** or **B** (2,340 at **D0** and 3,075 at **D7**, respectively). By contrast, DEGs at **SC**
119 were much less variable between the three lines (330 – 676) indicating a greater degree of
120 uniformity of expression. In terms of DEGs different from **H** across all sample points, **A**, **B** and **C**
121 had 480, 308, and 270, respectively.

122 Figure 3 presents the numbers of DEGs at each sample point between cultivars, along with
123 the degree of spatial variation and separation using Principal Component Analysis (PCA). The
124 close clustering observed is indicative of the high degree of reproducibility of expression seen
125 between each biological replicate for each line tested (Supplementary Figure S2). It also
126 emphasizes the differences in global expression patterns at any given ontological growth stage or
127 postharvest time point between lines **A**, **B** and **C**. DEGs ranged from 17,167 at **PW** (37.78% of
128 total genes; Figure 3b) to 22,482 at **EH** (49.48% of total genes; Figure 3a).

129 PCA of each of the DEGs (Figure 3) was used to scrutinize genes responsible for the
130 highest degree of spatial separation according to factor scores. This analysis yielded 1,568 genes.
131 Of particular note are several genes related to sulfur metabolism and GSL biosynthesis: *SLIM1a*,
132 *MYB28b*, *GSTF11*, *IGMT4a*, *TGG1c*, *NIT2c*, *ESM1b*, *ESM1d*, *GTR2a*, *GTR2d*, *DHAR3*, *SiRb*,
133 *SULTR2;2*, and *SULTR3;5a*.

134

135 **Sulfur and phytochemical composition of *E. sativa***

136 **Sulfur content of *E. sativa***

137 Total sulfur content for each of the breeding lines is presented in Figure 4a. No significant
138 differences were observed between lines and sample time points ($P = 0.434$). As will be discussed
139 in the subsequent sections this observation is of significance.

140

141 **Glucosinolate profiles and contents of *E. sativa***

142 For each of the cultivars between the first (**H**) and second cuts (**SC**), an increase in total
143 GSL concentrations was observed due to elevations of GSV (**A**, a 1.4-fold increase, $P<0.0001$; **B**,
144 a 1.6-fold increase, $P<0.0001$; **C**, a 1.8-fold increase, $P<0.0001$) and GRA (**B**, a 2.6-fold increase,
145 $P<0.0001$; **C**, a 1.8-fold increase, $P<0.0001$; Supplementary Data File S1). Line **C** produced the
146 highest total concentrations of GSLs in **SC** (a 1.7-fold increase; $P<0.0001$), and line **B** also saw
147 significant elevations compared to **H** (a 1.6-fold increase; $P<0.0001$).

148 Line **A** contained the greatest GSL concentrations compared to **B** and **C**, until **D7** where
149 content declined significantly (a 0.6-fold decrease compared to **D0**, $P<0.0001$; Supplementary
150 Data File S1). **C** by comparison contained high concentrations of GSLs during shelf life, peaking
151 at **D0**, with a non-significant decrease at **D7** (0.3-fold reduction). This line did not demonstrate the
152 same decline in GSLs towards the end of shelf life as in the other two, and displays a propensity
153 for maintaining GSLs for longer into the shelf life period. We hypothesised that this was due to
154 fundamental differences in expression of GSL biosynthesis genes pre and postharvest.

155

156 **Glucosinolate hydrolysis product profiles and contents of *E. sativa***

157 Glucosinolate hydrolysis product (GHP) concentrations are presented in Figure 4c (see
158 Supplementary Data File S1 for ANOVA and Tukey's HSD significances). As with previous
159 studies of rocket¹⁷, three main GHPs were detected: sativin (a 1,3-thiazepane-2-thione; hydrolysis
160 product of GSV; SAT), erucin (ITC of glucoerucin; GER), and SF. The fluctuations in total GHP
161 concentration mirror those observed for GSLs, however the increases between **H** and **SC** are much
162 less pronounced, with no significant differences between cuts.

163 As with GSLs, line **B** displayed the lowest concentrations of GHPs, whereas the differences
164 between lines **A** and **C** are less apparent. The trend of reduction of GHPs over shelf life is also
165 visible for lines **A** and **B**, though only significant in **B** (a 0.9-fold reduction, $P<0.0001$).
166 Concentrations remained higher in line **C** (1.2 mg.g⁻¹ dw, a 0.6-fold reduction from **D0**).

167

168 **Monosaccharide profiles and contents of *E. sativa***

169 Monosaccharides are important in terms of sensory attributes and the masking of bitter and
170 pungent sensory attributes in rocket⁵ altering consumer perception and preference. Glucose is also
171 known to influence stress responses and interact with MYB transcription factors¹⁸ (Figure 1b).

172 The concentrations of sugars observed in *E. sativa* lines are presented in Figure 4d (see
173 Supplementary Data File S1 for ANOVA and Tukey's HSD significances).

174 Unlike previous reports¹⁵ the changes in sugar concentrations in this study were dynamic
175 across each of the respective time points. Both lines **A** and **B** contained low total concentrations
176 compared to line **C**. Line **B** displayed consistent concentrations, with no significant differences
177 observed. **A** showed a similar trend to GSL and GHP concentrations by declining at the end of
178 shelf life (**D7**; a 0.5-fold decrease from **D0**, $P < 0.0001$).

179 Line **C** is distinct from the others in terms of its sugar profile and the relative differences
180 between sample points. Concentrations increased postharvest (**D0** and **D7**; a 1.4 and 1.2-fold
181 increase relative to **H**, respectively), perhaps owing to a breakdown of stored carbohydrate to
182 facilitate respiration. Line **C** sugar content consists primarily of glucose, whereas **B** tended to have
183 greater concentrations of galactose, and **A** was composed of similar amounts of each
184 monosaccharide.

185

186 **Sulfur assimilation and glucosinolate biosynthesis pathway gene expression analysis**

187 **Sulfate assimilation gene expression**

188 Figure 5a presents differential gene expression within the sulfate assimilation pathway of
189 *E. sativa*. All significances quoted hereafter were at the $P < 0.001$ significance level. In the primary
190 stages of sulfur metabolism, assimilated sulfate is activated via adenylation to adenosine-5'-
191 phosphosulfate (APS), catalyzed by ATP sulfurylase (ATPS)¹⁹. In *E. sativa* four ATPS-encoding
192 genes were identified: *APS1a*, *APS1b*, *APS2*, and *APS3* (Figure 1a). Very few significant DEGs
193 were observed between sample points for each respective rocket line (see Supplementary Data File
194 S2 for full values and statistics of each sample comparison). However, between **H**, **SC**, and **PW**,
195 each respective line did show significant differential expression of ATPS genes.

196 In the second stage of the pathway, APS is reduced to sulfite by adenosine-5'-
197 phosphosulfate reductase (APR) genes²⁰. Four APRs were identified (*APR1a*, *APR1b*, *APR2a*,
198 and *APR2b*) as well as six ARP-like genes (*APRL4*, *APRL5a*, *APRL5b*, *APRL5c*, *APRL7a*, and
199 *APRL7b*). *APR1a* and *APR2a* showed significant differential expression across multiple samples
200 and time points (Figure 4b). Line **B** displayed low relative expression of these genes compared to
201 **A** and **C**. Line **C** exhibits significantly higher expression postharvest compared to **H**; 2.2 log₂-fold
202 (**D0**) and 2.7 log₂-fold (**D7**) increases of *APR1a*, and 0.9 log₂-fold (**D0**) and 1.1 log₂-fold (**D7**)

203 increases of *APR2a* were observed. We hypothesise that this may be indicative of a greater ability
204 to remobilize sulfate via APS genes to facilitate and maintain secondary metabolite biosynthesis
205 for longer into shelf life.

206 Two copies of genes encoding sulfite reductase (SIR; *SiRa* and *SiRb*) were identified. *SiRa*
207 showed significantly higher levels of expression in line **C** (Figure 5b). Line **C** had no significant
208 change in activity of this gene relative to time point **H**, however both lines **A** and **B** had
209 significantly lower expression postharvest (Figure 5b).

210

211 **Sulfur metabolism transcription, regulation, and transport gene expression**

212 Three copies of *SDII* (*SDIIa*, *SDIIb*, and *SDIIc*) and three copies of *SLIMI* (*SULFUR*
213 *LIMITATION 1*, aka *ETHYLENE INDUCED 3*; *SLIMIa*, *SLIMIb*, and *SLIMIc*) were identified
214 within the genome annotation. These genes are thought to play critical roles in the management
215 and use-efficiency of sulfur in plants, and have been linked with optimization of GSL biosynthesis
216 under S-limited conditions in *A. thaliana*¹⁰.

217 *SDIIa* and *SDIIc* were differentially expressed between each line (Figure 5b,
218 Supplementary Data File S2), with **C** having the highest levels of expression postharvest. It might
219 be expected that each line would see a similar trend of expression over the course of shelf life, as
220 additional sulfur is not obtainable; however only line **C** displayed this (Figure 5b). As shown in
221 Figure 4a, sulfur content was not significantly different between any experimental stages or
222 between breeding lines. Expression of MYB28 orthologs were not negatively associated with
223 expression of *SDII* gene copies. Previous research has shown that the SDI1 protein binds to
224 MYB28, inactivating expression and reducing GSL biosynthesis²¹. In *E. sativa* the opposite
225 appears to be true, with significant positive correlations between respective expression of two
226 MYB28 copies and *MYB29* with SDI1 copies (*SDIIa* and *MYB29*, $r = 0.72$; *SDIIb* and *MYB28a*,
227 $r = 0.507$; *SDIIc* and *MYB28c*, $r = 0.459$; Figure 6). At **D7**, both **A** and **B** had significantly lower
228 expression levels compared with **H** (a 2.2 and 3.7 log₂-fold reduction of *SDIIa*, respectively; and
229 a 3.9 and 3.2 log₂-fold significant reduction of *SDIIc*, respectively).

230 A similar pattern was observed for *SLIMIb*, where expression in **C** was significantly higher
231 at time points **D0** and **D7** relative to **H** (a 1.4 and 1.1 log₂-fold significant increase, respectively).
232 Expression of *SLIMIc* by comparison was not significantly different for each respective plant line
233 between time points, but there were clear and significant differences in expression between lines

234 (Supplementary Data File S2). Line **C** had highest expression of this gene, followed by **A**; with **B**
235 having significantly lower expression (Figure 5b). Previous studies have shown that *SLIM1* down
236 regulates *APK* gene expression and *GSL* biosynthesis as a way of conserving sulfur for primary
237 metabolism²¹. Our data suggest that this is only the case between *SLIM1a* and *APK3* ($r = -0.597$,
238 $P < 0.001$; Figure 6). *SLIM1a* expression was positively (and significantly) correlated with *APK*
239 expression ($r = 0.521$), and *SLIM1b* and *SLIM1c* with *APK4* ($r = 0.575$ and 0.698 , respectively;
240 Figure 6). This suggests *E. sativa* has a complex and interacting network of sulfur metabolism
241 genes, where functions may not be analogous to those found in *A. thaliana*.

242 16 sulfur transport (*SULTR*) genes were identified; of note were *SULTR1;2a*, *SULTR2;1a*,
243 *SULTR2;1b*, *SULTR4;1a*, and *SULTR4;2*. *SULTR1;2a* has been associated with the uptake of
244 environmental sulfate in root tissues (Supplementary Data File S2), but low levels of expression
245 were detected in leaf tissues. Postharvest, line **C** had differential expression of this gene compared
246 to **A** and **B** in **D7** samples (Figure 4b). This was more pronounced for *SULTR2;1a* and *SULTR2;1b*,
247 and both **A** and **B** had significant reductions in expression at **D0** and **D7** relative to **H**. **SC** samples
248 showed significant increases relative to **H**, with the exception of *SULTR2;1a* in **B**.

249 *SULTR4;1* and *SULTR4;2* genes also had distinct patterns of expression between lines.
250 *SULTR4;1a* saw significant increases in expression in postharvest samples relative to **H** (Figure
251 5b). Line **A** had higher expression of *SULTR4;2* during growth before declining significantly post-
252 wash (**D0**). The opposite trend was seen in **C**, where gene expression peaked at **D7**. These data are
253 suggestive of more active intra-leaf sulfur transport in line **C** postharvest, and may be associated
254 with the higher expression of *APR*, *SIR*, *SDII* and *SLIM1* genes to facilitate more efficient S
255 utilisation during this period.

256

257 **Glutathione synthesis**

258 With the exception of *GSH2b*, glutathione synthase genes were most highly expressed in
259 rocket line **B**, with significant increases observed postharvest (Figure 5b). Line **A** and **C** were
260 unchanged between sample points for these genes, but had a marked difference in expression for
261 *GSH2b* relative to each other. **B** had negligible levels of *GSH2b* expression.

262 As both glutathione and secondary S-containing metabolites, such as *GSLs*, have been
263 associated with antioxidant responses²¹ the differences observed between each of the lines in terms
264 of both *GSL* concentrations and glutathione-related gene expression, may be indicative of different

265 adaptive metabolic strategies for dealing with oxidative stress postharvest. Lines **A** and **C** favor
266 secondary sulfur metabolism and the synthesis of GSLs, and **B** favors primary sulfur metabolism
267 and glutathione synthesis.

268

269 **Glucosinolate-related transcription factors**

270 *MYC2a* and *MYC2c* were highly expressed in line **A**, and had uniform patterns of relative
271 expression. **SC** had the highest expression values for this line, suggesting a general response to
272 mechanical wounding and stress, however this was not significantly different from **H**. The only
273 significant difference for *MYC2c* between **H** and **SC** was in line **B** (a 0.7 log₂-fold increase; Figure
274 5c).

275 *MYB28a* and *MYB28b* display high degrees of differential expression between each rocket
276 line. While **A** has high expression of *MYB28a* in samples **EH**, **H**, **SC**, and **PW**, it has by
277 comparison lower expression of *MYB28b* compared to **C** (Figure 5c). **C** on the other hand has
278 relatively high expression for both of these TFs, and displays significantly higher expression
279 postharvest, up to and including **D7**. Combined with what is known about these TFs in other
280 Brassicaceae species, it is likely that the differences in GSL concentrations observed postharvest
281 are linked to the differential expression of *MYB28a* and *MYB28b* between the respective lines.

282

283 **Glucosinolate biosynthesis**

284 Rocket contains two genes encoding BCAT4, and two genes encoding BCAT3; converting
285 2-oxo acids to homomethionine and dihomomethionine. *BCAT3-1a* displayed no significant
286 variation between lines during growth, but saw significant increases for all (compared to **H**) at **D0**
287 and **D7** (Supplementary Data File S2). The most marked and significant increase was in **B**. It is
288 unclear how this ‘preference’ for BCAT3 activity over BCAT4 is regulated or affects the synthesis
289 pathway, but the relative and respective activity of these genes is correlated with GSL content.

290 Only orthologs of *MAM1* were identified, and each of the three copies had differing
291 expression patterns (Figure 5d). Line **A** displayed higher relative expression of *MAM1a*, whereas
292 **C** had greater expression for *MAM1b* and *MAM1c*. **B** however maintained low expression for all
293 three of these genes. **A** had reduced expression activity during shelf life, whereas in **C**, levels were
294 significantly higher compared to **H** (Figure 5d).

295 One *CYP79F1* homolog was found in rocket, with no expression found for a corresponding
296 *CYP79F2* gene. The lack of a *CYP79F2* homolog in rocket may be suggestive of a loss of function,
297 and/or redundancy with other enzymes. Of note for *CYP79F1* expression was the significant
298 differences observed between **EH** and **H**, indicating that earlier harvests of rocket leaves may have
299 a reduced ability for GSL biosynthesis compared with later ones and second cuts (**SC**). Expression
300 was significantly greater in **C** during shelf life. In the conversion of aldoximes to nitrile oxides,
301 *CYP83A1* expression was higher in **A** and **C** than **B**, with line **C** having significantly higher
302 expression in shelf life samples (Figure 5d).

303

304 **Glucosinolate hydrolysis**

305 11 *TGG1* (myrosinase) orthologs, and three *TGG2* (Supplementary Data File S2) were
306 identified within the annotation. Some of these genes appear to have differential expression
307 according to ontogeny and shelf life, with some copies expressed at **EH** with none during
308 postharvest (e.g. *TGG1h*, *TGG2a*, and *TGG2c*; Figure 5e). Others however display the inverse of
309 this, with increased relative expression postharvest (*TGG1a*; Supplementary Data File S2). It is
310 known that myrosinases are functionally redundant, however it has also been noted that their
311 activity and specificity is linked with developmental processes, and may explain some of the high
312 levels of expression observed at **EH**.

313 An explanation for the lack of nitrile GHPs in rocket may be that the high expression of
314 *NSP5* is inhibited by the five *EPITHIOSPECIFIER MODIFIER 1 (ESM1)* orthologs found in the
315 rocket genome. These proteins are known to inhibit the action of NSPs and promote ITC formation.
316 Expression was significantly greater in line **C** for *ESM1b* (Figure 5e) at all sample points, and fits
317 with the observed pattern of sustained GHP formation postharvest. The lower activities in **A** and
318 **B** did not however correspond to a reciprocal decrease in the relative concentrations of GHPs, and
319 neither were nitrile concentrations at anything above trace levels. Much further work is needed to
320 explain the genetic regulation of GHP formation in rocket and the high expression of *NSP5*.

321

322 **Glucosinolate transporters**

323 Eight GSL transporter genes were identified in the rocket annotation; four *GTR1* and four
324 *GTR2* homologs. These genes are involved in leaf distribution and long-range phloem GSL
325 distribution, respectively. Expression of *GTR2a* and *GTR2d* were significantly correlated with

326 GSL abundance and GHP formation in the analysed leaf tissues. Of particular note is that **B** had
327 no expression of *GTR2d* at any of the sample points, indicating that the gene may be non-
328 functional, and contribute to the overall low concentrations of GSLs observed. The exact site(s) of
329 GSL synthesis has not been established in *E. sativa*, however it is known that transport from the
330 root tissues to leaves occurs in other Brassicaceae. If this transport system is impaired in **B**, this
331 would explain the significantly lower abundance of GSLs observed in leaves (Figure 4b). Coupled
332 with the high expression of glutathione-related genes and similar sulfur content of **B** compared to
333 lines **A** and **C**, the inactivity of this gene copy may have significant effects on leaf sulfur transport,
334 metabolism, and antioxidant response. The lower GSL content in leaves may therefore be
335 compensated by increased glutathione synthesis.

336

337 **Principal component analysis of sulfur and glucosinolate metabolism genes**

338 Hereafter, only correlations significant at the $P < 0.001$ level are presented and discussed.
339 *SULTR4;1a* was significantly correlated with GRA concentrations ($r = 0.577$), which is associated
340 with shelf life samples for lines **A** and **C** (Figure 6b, cluster **II**). Figure 6a and 6b show a distinct
341 separation between ontogenic and shelf life samples along PC1. The increased expression of sulfur
342 transport genes such as this postharvest may provide some explanation as to why GSL
343 concentrations increase in the initial stages shelf life (**PW**), as S may be re-mobilized to facilitate
344 biosynthesis. Efficient transport and storage of sulfur pre-harvest may also facilitate better
345 retention and decreased degradation of GSLs postharvest. This can be seen in Figure 6b (**V**) where
346 *SULTR3;1a* and *SULTR3;2* are associated with pre-harvest expression.

347 SF and SAT concentrations were significantly correlated with *APR2a* gene expression
348 (Figure 6c **I** and **II**; $r = 0.58$, SF; $r = 0.464$, SAT) and associated in particular with **A** ontogenic
349 samples and **PW**. *APR2* is known to contribute to sulfur accumulation and homeostasis, as well as
350 facilitating cysteine synthesis, and is associated with increased myrosinase activity and GSL
351 recycling. Line **A** (on average) contained the highest ontogenic concentrations of GRA, SF, GSV,
352 and SAT (Figure 4b and 4c); this is supported by a significant correlation and association with
353 GSL-related transcription factors *MYB28a* ($r = 0.486$, SF; $r = 0.53$, SAT), *MYC2a* ($r = 0.596$, SF;
354 $r = 0.626$, SAT) and *MYC2c* ($r = 0.584$, SF; $r = 0.583$, GSV; $r = 0.634$, SAT; Figure 6c **I** and **II**),
355 as well as a drought tolerance-related gene *SAL1b* ($r = 0.595$, GRA; $r = 0.547$, SF; $r = 0.499$, SAT;
356 Figure 6a **II**, **III** and **IV**, 6c **II**). **A** was also associated with increased activity of *MAM1a* (Figure

357 6c **II**), facilitating greater GRA biosynthesis through chain elongation. It may be that lines **A** and
358 **C** have increased relative GRA concentrations at **EH** and **H**, but preferentially express *MYB28c*
359 and *MYB28b*, respectively. It is unknown if the function of each MYB28 TF are redundant in
360 rocket, but these data would suggest that there is some clear overlap of function, though the
361 expression of *MYB28b* is associated with increased GSL biosynthesis postharvest (Figure 5c).

362 The lower relative expression in line **B** for many of these genes is consistent with its lower
363 GSL and hydrolysis product concentrations, irrespective of sample point. GRA/SF, GSV/SAT,
364 and GER concentrations were significantly and negatively correlated with *SPERMIDINE*
365 *SYNTHASE 1c* (*SPDS1c*; $r = -0.622$, GRA; $r = -0.614$, SF; $r = -0.6$, GSV; $r = -0.454$, SAT; $r = -$
366 0.604 , GER), which had a high degree of co-separation in all **B** samples (Figure 6c **III**). This
367 association may be related to increased primary S metabolism and reduced partitioning of
368 methionine for secondary S metabolites (Figure 1a).

369

370 **Discussion**

371 ***E. sativa* has a distinct and complex glucosinolate pathway**

372 Gene orthologs have undergone numerous duplications in *Eruca*, and it is not clear what
373 the function(s) of these numerous copies may be. It may be the result of segmental duplications
374 within the genome, such as those observed in the *Brassica A* genome²², and future, more detailed
375 studies of the *Eruca* genome structure may reveal the nature and number of any such events. Genes
376 such as SOTs, FMO_{GS-OXS}, myrosinases (TGGs), *ESM1s*, and GSL transporters all have several
377 copies, and it has yet to be determined if these perform the same function as in *Arabidopsis*, or
378 have evolved new ones that are responsible for the unique GSL profile of rocket. Examples of this
379 include increases in copy numbers compared to *B. rapa* and *B. oleracea*. In these two species two
380 copies of *SOT18* have been identified²³, whereas *E. sativa* has seven. *B. rapa* has two copies of
381 FMO_{GS-OX} genes, and salad rocket has at least ten. The related *Diplotaxis tenuifolia* (“wild”
382 rocket) transcriptome has been reported to contain three copies of *MYB28*²⁴, and is consistent with
383 the hypothesis that duplication occurred after their divergence with a common ancestor within the
384 *B. oleracea* lineage.

385 Both *Arabidopsis* and *B. rapa* have four myrosinase gene copies, and *B. oleracea* has six
386²³. Our data indicate that *Eruca* has at least 14 copies, as well as two copies encoding PEN2

387 myrosinase. There has evidently been a massive diversification and duplication of these enzymes
388 in rocket, but it has yet to be established if this is reflected in functionality and spatial expression.

389

390 **Glutathione metabolism competes with glucosinolate biosynthesis for sulfur**

391 As shown in Figure 4a, the content of sulfur between the three tested breeding lines was
392 not significantly different. In light of the observed differences in gene expression and GSL
393 accumulations, we theorize that primary and secondary sulfur metabolism pathways ‘compete’ for
394 assimilated environmental sulfur. As content was not significantly different postharvest (**PW, D0,**
395 **and D7**) compared to pre-harvest first cut (**H**) in any of the breeding lines, the degree of
396 remobilization and ability to synthesize/recycle GSLs is under strict genetic control. The exact
397 reasons why each line differs in this respect are unclear, but as shown in Figure 4b, the amount of
398 total sulfur assimilated during growth is not reflected in the postharvest concentrations of GSLs.
399 Line **B** displays this: it contains significantly no more or less sulfur than lines **A** or **C**, yet
400 synthesizes far fewer GSLs.

401 We hypothesise that the lack of expression of *GTR2d* impairs GSL transport from major
402 sites of synthesis (possibly in the roots) thereby forcing leaf tissues to cope with oxidative stress
403 via glutathione synthesis and preferential shunting of sulfur into the primary metabolism pathway.
404 This is evidenced by the significantly higher expression of glutathione synthetase (*GSH2a*) and
405 glutamate-cysteine ligase genes (*GSH1a* and *GSH1c*). Lines **A** and **C** by contrast may cope with
406 oxidative stress postharvest by sustaining GSL biosynthesis. As **C** contained greater
407 concentrations of monosaccharides, we theorize that this facilitates GSL metabolism for longer,
408 reducing oxidative stress, and delaying the onset of senescence.

409

410 **Second cut rocket has greater uniformity of gene expression and glucosinolate content than** 411 **first cut**

412 Anecdotal evidence has suggested that the second cut (**SC**) is more consistent in terms of
413 taste and flavour than the first (**H**). For the first time we present transcriptomic and phytochemical
414 evidence to support this assertion. All cultivars saw increases in GSL abundance in the second cut,
415 and the shift is best exemplified by line **B**. First cut (**H**) samples were low in expression of
416 numerous GSL biosynthesis related genes that later saw significant increases at **SC**. This had the

417 result of making its expression profile more similar to that of **A** and **C**, along with comparable
418 concentrations of GSLs.

419 When the respective monosaccharide profiles are taken into account, it can be seen that
420 there is a general decrease in sugars at **SC**. This then leads to a shift in the ratio between
421 monosaccharides and GHPs, making the **SC** samples more pungent. Although flavour was not
422 tested in this study, our previous work has highlighted the importance of GHP:sugar ratios in
423 determining pungency²⁵.

424

425 **Materials and Methods**

426 **Plant material for genome sequencing**

427 Three elite inbred lines of salad rocket were produced through self-pollination for five
428 generations at Elsoms Seeds Ltd. (Spalding, UK) from 2010-2016, giving an estimated inbreeding
429 coefficient of 0.969²⁶. Each line was derived from germplasm accessions obtained from the
430 Leibniz-Institut für Pflanzengenetik und Kulturpflanzenforschung (IPK Gatersleben, Germany).
431 For reasons of commercial sensitivity these lines (**A**, **B**, and **C**) and their lineage will not be
432 identified.

433 For genome sequencing, plants of each line were grown under controlled growth room
434 conditions, and leaf samples had DNA extracted and sent to the Earlham Institute for QC analysis.
435 Samples were quantified using a Qubit fluorometer and dsDNA assay kit (ThermoFisher
436 Scientific, Loughborough, UK) and assessed for quality using NanoDrop (ThermoFisher
437 Scientific; according to 260/280 and 260/230 ratios).

438

439 **Genome sequence library preparation and assembly**

440 DNA sequencing and assembly was performed as a service by the Earlham Institute
441 (Norwich, UK). *De novo* genome sequencing and assembly was performed using PCR free paired-
442 end (PE) and LMP sequencing. After DNA sample QC, line **C** was selected for sequencing and
443 reference genome assembly. One PCR free PE library was constructed from gDNA, and sequenced
444 on one lane of an Illumina HiSeq2500 in rapid run-mode (v2) using 250 bp PE reads. LMP
445 sequencing was also conducted using one set of Nextera libraries (Illumina) from gDNA, and
446 sequenced on one lane of an Illumina MiSeq with 250 bp PE reads. After data QC and assembly
447 of the high coverage PE library, LMP libraries were mapped to determine their suitability for

448 assembly improvement. Three additional libraries were selected and re-sequenced to a higher
449 depth of coverage on a single lane of an Illumina HiSeq2500 in rapid run-mode, to again yield 250
450 bp PE reads.

451
452 **Genome sequencing bioinformatics**
453 FASTQ files were converted to BAM format using PicardTools (v1.84,
454 <http://broadinstitute.github.io/picard/>; FastqToSam option) and then assembled using DISCOVAR
455 *de novo* sequence assembler (build revision 52488) ²⁷. All LMP libraries were processed using
456 NextClip ²⁸ to analyse and create a high quality read subset for scaffolding the DISCOVAR-
457 assembled sequences. SOAP ²⁹ and SSPACE ³⁰ were used to scaffold the DISCOVAR assembly
458 using data from three of the NextClip-processed LMP read libraries.

459
460 **Genome annotation**
461 Annotation was performed by Novogene Co. Ltd. (Hong Kong). A homology and *de novo*-
462 based approach was taken in order to identify TEs. The homology-based approach used known
463 repetitive sequence databases: RepBase ³¹, RepeatProteinMask, and RepeatMasker
464 (<http://www.repeatmasker.org/>). *De novo* repeat libraries were created using LTR_FINDER ³²,
465 RepeatScout (<http://www.repeatmasker.org/>), and RepeatModeler
466 (<http://www.repeatmasker.org/RepeatModeler.html>).

467 An integrated approach was taken to compute consensus gene structures, such as cDNA,
468 proteins in related species, and *de novo* predictions. The homology-based approach used the
469 related genomes of *A. lyrata*, *A. thaliana*, *B. napus*, *Boechera stricta*, *Capsella rubella*, and *R.*
470 *sativus* (Supplementary Table S3) to compare against *E. sativa* to find homologous sequences, and
471 predict gene structures (using BLAST and genewise) ^{33,34}. *Ab initio* statistical models were also
472 used to predict genes and their intron-exon structures; e.g. Augustus ³⁵, GlimmerHMM ³⁶, and
473 SNAP (<http://homepage.mac.com/iankorf/>). EVidenceModeler (EVM) ³⁷ software was then used
474 to combine *ab initio* predictions, protein and transcript alignments, and RNAseq data into weighted
475 consensus gene structures. Lastly, PASA was used to update the consensus predictions by adding
476 UTR annotations and models for alternative splicing isoforms. All predicted proteins were
477 functionally annotated using alignments to SwissProt, TrEMBL ³⁸, KEGG ³⁹, and InterPro ⁴⁰
478 (Figure 2b).

479 The full reference genome sequence and annotation can be found in the European
480 Nucleotide Archive (Assembly accession no: GCA_902460325; Study ID: PRJEB34051; Sample
481 ID: ERS3673677; Annotation accession number ERZ1066251).

482

483 **Plant material growth and collection for RNA, elemental, and phytochemical analyses**

484 For RNAseq analyses seeds were sown in a random order in seedling compost, and raised
485 under controlled environment conditions in plastic trays inside a Weiss-Technik Fitotron cabinet
486 (Weiss-Technik UK Ltd., Loughborough, UK). Daytime temperature was set to 20 °C, and
487 nighttime temperatures to 14 °C (long day cycle; 16 h light, 8 h dark). Light intensity was set at
488 200 $\mu\text{mol m}^{-2} \text{s}^{-1}$. During a one-hour period of ‘dawn’ and ‘dusk’, light and temperature changes
489 were ramped on a gradient. Humidity was ambient. After ten days of growth, seedlings were
490 transplanted to one-litre pots in standard peat-based compost.

491 Postharvest (post sample **H**), leaves were stored for two days in a cold store (4 °C)¹⁵.
492 Samples for **D0** and **D7** were washed individually in mildly chlorinated water (sodium
493 hypochlorite; 30 ppm⁴¹) for two minutes, then rinsed for one minute with distilled water (all at
494 >14 °C to avoid cold-shock). Leaves were dried of excess moisture for one minute using a kitchen
495 salad spinner, then placed in fresh bags, sealed, and stored overnight at 4 °C. Shelf life leaves were
496 stored in the cold and dark (4 °C) for seven days (**D7**) – typical of the use-by date of commercially
497 bagged leaves.

498 All samples were taken between the hours of 1 – 3 pm to mitigate diurnal fluctuations in
499 phytochemical content and gene expression⁴². Immediately after each of the aforementioned
500 samples was taken, leaves were frozen using liquid nitrogen and ground into a fine powder using
501 a pestle and mortar. Samples were stored at -80 °C in tubes and lyophilized prior to chemical
502 analysis. A subset of non-lyophilized sample was kept aside for RNA extractions.

503

504 **RNA extraction and quality control**

505 RNA for RNAseq and qRT-PCR analyses was extracted using RNeasy Plant Mini Kit
506 (Qiagen, Manchester, UK) according to the manufacturer ‘Plants and Fungi’ procedure. As part of
507 the protocol, an on-column DNase digestion was incorporated according to the RNase-Free DNase
508 Set (Qiagen) protocol. Samples were checked for degradation and contamination prior to
509 sequencing using agarose gel electrophoresis (1%, TAE buffer), Qubit, and NanoPhotometer

510 (Implen, CA, USA) methods. Briefly, ≥ 2 μg of total RNA was obtained for each sample at a
511 minimum concentration of ≥ 50 $\text{ng}\cdot\mu\text{L}^{-1}$. RNA integrity was determined and evaluated using an
512 Agilent 2100 Bioanalyzer⁴³.

513

514 **RNAseq library preparation and sequencing**

515 After QC procedures, sequencing libraries were prepared using NEBNext Ultra RNA
516 Library Prep Kit for Illumina (NEB, MA, USA) following the manufacturer's instructions, and
517 index codes were added to attribute sequences to each sample. mRNA was purified from total
518 RNA by using poly-T oligo-attached magnetic beads. Fragmentation was carried out using divalent
519 cations under elevated temperature in NEBNext First Strand Synthesis Reaction Buffer (5x). First
520 strand cDNA was synthesized using random hexamer primer and M-MuLV Reverse Transcriptase
521 (RNase H-). Second strand cDNA synthesis was subsequently performed using DNA Polymerase
522 I and RNase H. Remaining overhangs were converted into blunt ends via exonuclease/polymerase
523 activities. After adenylation of 3' ends of DNA fragments, NEBNext Adaptor with hairpin loop
524 structure were ligated to prepare for hybridization. In order to select cDNA fragments of 150 –
525 200 bp in length preferentially, the library fragments were purified with an AMPure XP system
526 (Beckman Coulter, MA, USA). 3 μl USER Enzyme (NEB) was used with size-selected, adaptor-
527 ligated cDNA at 37 °C for 15 min, followed by 5 min at 95 °C before PCR. PCR was performed
528 with Phusion High-Fidelity DNA polymerase, Universal PCR primers and Index (X) Primer.
529 Finally, PCR products were purified (AMPure XP system) and library quality was assessed on
530 using an Agilent 2100 Bioanalyzer system.

531 The clustering of the index-coded samples was performed on a cBot Cluster Generation
532 System using HiSeq PE Cluster Kit cBot-HS (Illumina) according to the manufacturer's
533 instructions. After cluster generation, the library preparations were sequenced on an Illumina
534 HiSeq platform and 125 bp/150 bp paired-end reads were generated.

535

536 **RNAseq bioinformatics**

537 Raw data (raw reads) of FASTQ format were firstly processed through Novogene Co. Ltd.
538 perl scripts. Clean reads were obtained by removing reads containing adapter, reads containing
539 ploy-N, and low quality reads from the raw data. At the same time, Q20, Q30 and GC content of
540 the clean data were calculated.

541 An index of the reference genome was built using Bowtie (v2.2.3), and PE clean reads
542 were aligned to the reference genome using TopHat (v2.0.12)^{44–46}. TopHat was selected as the
543 mapping tool as it can generate a database of splice junctions based on the gene model annotation
544 file and thus a better mapping result is achieved than other non-splice mapping tools.

545 HTSeq (v0.6.1) was used to count the read numbers mapped to each gene⁴⁷. FPKM
546 (Fragments Per Kilobase of transcript sequence per Millions base pairs sequenced) of each gene
547 was calculated based on the length of the gene and reads count mapped to each gene. Differential
548 expression analysis of each sample point/inbred line (three biological replicates) was performed
549 using the DESeq R package (1.18.0)^{46,48}. The resulting *P*-values were adjusted using Benjamini
550 and Hochberg's approach for controlling the false discovery rate. Genes with an adjusted *P*-value
551 (<0.05) were assigned as being significantly differentially expressed.

552

553 **RNAseq validation by qRT-PCR**

554 Independent RNA extractions were conducted for qRT-PCR validation, and quality
555 checked according to the same protocols and instrumentation as for RNAseq. cDNA synthesis was
556 conducted using qPCRBIO cDNA Synthesis Kit (PCR Biosystems Ltd., London, UK) according
557 to the manufacturer instructions. cDNA was then diluted 10x prior to analysis. All 54 biological
558 samples were tested in triplicate.

559 PCR primers were designed using PRIMER3 (<http://bioinfo.ut.ee/primer3/>) using default
560 settings. Ten genes related to GSL biosynthesis and transcription were selected at random for the
561 validation analysis (*BCAT4*, *CYP83B1*, *MYB122-1a*, *MYB51a*, *SOT16*, *SUR1*, *TGG1b*, *TGG1d*,
562 *TGG1j*, and *UGT74B1*), with *ACT11* used as a reference gene⁴⁹. Gene sequences of *E. sativa* were
563 obtained using NovoFinder (Novogene Co. Ltd.), and primer annealing sites were designed to span
564 intron-intron boundaries where possible (see Supplementary Table S5).

565 Analysis was performed using the 2^{-ΔΔCt} method⁵⁰ on a Roche LightCycler 480 Instrument
566 and the Advanced Relative Quantification protocol (v1.5.1). Primer efficiencies were determined
567 by analyzing each primer set with log-fold dilutions of cDNA (Supplementary Table S5). 2x
568 qPCRBIO SyGreen Blue Mix Lo-ROX (PCR Biosystems Ltd.) was used to prepare a master mix
569 for all reactions. Reaction volumes totaled 10 μL, and the PCR method used was as per the
570 manufacturer recommendations.

571 Data were normalized and expressed as the log₂-fold change relative to *ACT11*. RNAseq
572 data for each of the tested genes were similarly converted for direct comparison of the two
573 methodologies (Supplementary Figure S3). An ANOVA test found no significant differences
574 between the two data sets for each of the respective genes.

575

576 **Intact glucosinolate extraction and analysis by LC-MS**

577 Intact GSLs were extracted according to the protocol used by Bell et al.⁸. Immediately
578 before LC-MS analysis, samples were diluted with 4 mL of HPLC-grade water. Samples were
579 analyzed in a random sequence with standards and QC samples. External standards of sinigrin
580 (SIN; >99% TLC), GRA (99.86%, HPLC), glucoalyssin (GAL) (98.8%, HPLC), 4OHB (96.19%,
581 HPLC) and GER (99.68%, HPLC) were prepared for quantification of GSL compounds. SIN was
582 used to quantify DGTB, GSV, and DMB, as no standards are available for these compounds.
583 4OHB was used to quantify the indole GSLs 4MOB and neoglucobrassicin (NGB). All standards
584 with the exception of SIN (Sigma Merck, Gillingham, UK) were purchased from PhytoPlan
585 (Heidelberg, Germany). Recovery of extracted GSLs was calculated by spiking six random
586 samples (in duplicate) with SIN (0.06 μM). The average recovery was 104.8%, indicating excellent
587 preservation of GSLs throughout the extraction process. Limits of detection (LOD) and
588 quantification (LOQ) were established for the method by running serial dilutions of SIN (LOD =
589 2.14 mg.kg⁻¹; LOQ = 6.48 mg.kg⁻¹).

590 LC-MS analysis was performed in the negative ion mode on an Agilent 1260 Infinity Series
591 LC system (Agilent, Stockport, UK) equipped with a binary pump, degasser, auto-sampler, column
592 heater and diode array detector, coupled to an Agilent 6120 Series single quadrupole mass
593 spectrometer. Separation of samples was achieved on a Gemini 3 μm C₁₈ 110Å (150 x 4.6 mm)
594 column (with Security Guard column, C₁₈; 4mm x 3mm; Phenomenex, Macclesfield, UK). GSLs
595 were separated during a 40 min chromatographic run, with a 5 min post-run sequence. Mobile
596 phases consisted of ammonium formate (0.1%; A) and acetonitrile (B) with the following gradient
597 timetable: (i) 0 min (A-B, 95:5, v/v); (ii) 0-13 mins (A-B, 95:5, v/v); (iii) 13-22 mins (A-B, 40:60,
598 v/v); (iv) 22-30 mins (A-B, 40:60, v/v); 30-35 mins (A-B, 95:5, v/v); (v) 35-40 mins (A-B, 95:5,
599 v/v). The flow rate was optimized for the system at 0.4 mL min⁻¹, with a column temperature of
600 30 °C, and 20 μl of sample injected into the system. Quantification was conducted using a diode
601 array detector at a wavelength of 229 nm.

602 MS analysis settings were as follows: Atmospheric pressure electrospray ionization was
603 carried out in negative ion mode (scan range m/z 100–1500 Da). Nebulizer pressure was set at 50
604 psi, gas-drying temperature at 350 °C, and capillary voltage at 2,000 V. Compounds were
605 identified using their primary ion mass $[M-H]^-$, and comparison to authentic standards^{51,52}. Data
606 were analyzed using Agilent OpenLAB CDS ChemStation Edition for LC-MS (vA.02.10). GSL
607 concentrations from each time point were averaged over three biological replicates with two
608 technical replicates of each ($n = 6$). This approach was also conducted for GHP and
609 monosaccharide content.

610

611 **Glucosinolate hydrolysis product extraction and analysis by GC-MS**

612 GHPs were extracted according to the protocol presented by Ku et al.⁵³ with the following
613 modification: samples were hydrolysed in d.H₂O for three hours at 30 °C before extraction with
614 dichloromethane (DCM) for 21 hours. This duration was optimized for maximum yields of GHPs
615 by comparison of extractions for three hours incubation at 30 °C with immediate DCM extraction,
616 and three, nine, and 21 hours post incubation with DCM. GC-MS analysis and GHP identification
617 was conducted according to the method presented by Bell et al.¹⁵. Concentrations of all GHPs
618 were calculated as equivalents of SF standard (Sigma).

619

620 **Monosaccharide extraction and analysis by HPLC**

621 Free monosaccharides were extracted according to the method presented by Bell et al.⁵,
622 with the exception that 0.2 g of lyophilized leaf powder was extracted. Extracts were analyzed on
623 an Agilent 1100 series HPLC system equipped with a binary pump, degasser, and auto-sampler,
624 with an external column heater (50 °C). A Bio-Rad Aminex HPX-87H (300 x 7.8 mm, 9 µm
625 particle size) column with a Micro-Guard Cation H guard column (Bio-Rad, Watford, UK) was
626 used to achieve separation with an isocratic gradient of 5 mM sulfuric acid, and a flow rate of 0.6
627 mL per min. A Polymer Laboratories ERC-7515 refractive index detector (Church Stretton, UK)
628 was used to detect monosaccharides. Compounds were quantified using authentic standards and
629 analyzed with Agilent ChemStation software (Santa Clara, CA, USA).

630

631 **Sulfur content analysis by ICP-OES**

632 Lyophilized samples were weighed into acid washed glass boiling tubes, pre-digested in 70% nitric
633 acid for 24 hours, before being heated to 90 °C for two hours using a heat block. Once cooled,
634 these were filtered through a 0.45 µM syringe filter, and diluted give an acid concentration of 3%.
635 These samples were analysed using inductively coupled plasma-optical emission spectroscopy
636 (ICP-OES) (Perkin Elmer Optima 7300 DV). Sulphur content was determined using the radial
637 signal at 181.975 nm.

638

639 **Statistical analyses**

640 All statistical analyses (not included in bioinformatics sections) were performed using XL
641 Stat (Addinsoft, Paris, France). Shapiro-Wilk normality tests were conducted for all variables and
642 fit a normal distribution. ANOVA with post-hoc Tukey's Honest Significant Difference (HSD)
643 tests were performed to generate multiple pairwise comparisons between sampling points for each
644 cultivar (i.e. **H** vs. **D7** for cultivar **B**) and between cultivars at each respective time point (i.e. **A**
645 vs. **B** for time point **H**) for phytochemical data (Supplementary Data File S1). All PCAs were
646 performed using Pearson correlation coefficient analysis, *n*-1 standardization, Varimax rotation,
647 and Kaiser Normalization. Phytochemical data were regressed onto the gene expression data as
648 supplementary variables for the targeted analysis.

649

650 **Supplementary Materials**

651

652 Fig. S1. RNAseq experimental design and sampling diagram.

653 Fig. S2. RNAseq sample replicate Pearson correlation analysis.

654 Fig. S3. RNAseq gene expression validation by qRT-PCR.

655 Table S1. Summary of genome assembly and annotation statistics.

656 Table S2. Reference genome sequence composition.

657 Table S3. Predicted protein-coding genes within the *E. sativa* reference genome.

658 Table S4. Numbers of genes with homology or functional assignment within the *E. sativa*
659 genome annotation.

660 Table S5. qRT-PCR primer sequences and amplification efficiencies.

661 Data file S1. Phytochemical Analyses of Variance (ANOVA) pairwise comparisons and RNAseq
662 gene expression correlation analysis.

663 Data file S2. Sulfur metabolism and glucosinolate biosynthesis gene expression values and
664 annotation descriptions.

665

666 **References and Notes**

- 667 1. Kopriva, S. *et al.* Editorial: *Frontiers of Sulfur Metabolism in Plant Growth, Development, and Stress*
668 *Response. Frontiers in Plant Science* **6**, (2016).
- 669 2. Frerigmann, H. & Gigolashvili, T. Update on the role of R2R3-MYBs in the regulation of glucosinolates
670 upon sulfur deficiency. *Front. Plant Sci.* **5**, 626 (2014).
- 671 3. Yan, X. F. & Chen, S. X. Regulation of plant glucosinolate metabolism. *Planta* **226**, 1343–1352 (2007).
- 672 4. Bell, L. & Wagstaff, C. Glucosinolates, Myrosinase Hydrolysis Products, and Flavonols Found in Rocket
673 (*Eruca sativa* and *Diplotaxis tenuifolia*). *J. Agric. Food Chem.* **62**, 4481–92 (2014).
- 674 5. Bell, L., Methven, L., Signore, A., Jose Oruna-Concha, M. & Wagstaff, C. Analysis of Seven Salad Rocket
675 (*Eruca sativa*) Accessions: The Relationships Between Sensory Attributes and Volatile and Non-volatile
676 Compounds. *Food Chem.* **218**, 181–191 (2017).
- 677 6. Winde, I. & Wittstock, U. Insect herbivore counteradaptations to the plant glucosinolate-myrosinase system.
678 *Phytochemistry* **72**, 1566–1575 (2011).
- 679 7. Kim, S. J. *et al.* Structural elucidation of 4-(cystein-S-yl)butyl glucosinolate from the leaves of *Eruca sativa*.
680 *Biosci. Biotechnol. Biochem.* **71**, 114–121 (2007).
- 681 8. Bell, L., Oruna-Concha, M. J. & Wagstaff, C. Identification and quantification of glucosinolate and flavonol
682 compounds in rocket salad (*Eruca sativa*, *Eruca vesicaria* and *Diplotaxis tenuifolia*) by LC-MS: highlighting
683 the potential for improving nutritional value of rocket crops. *Food Chem.* **172**, 852–861 (2015).
- 684 9. Kazan, K. & Manners, J. M. MYC2: The Master in Action. *Mol. Plant* **6**, 686–703 (2013).
- 685 10. Aarabi, F. *et al.* Sulfur deficiency-induced repressor proteins optimize glucosinolate biosynthesis in plants.
686 *Sci. Adv.* **2**, (2016).
- 687 11. Chan, K. X., Wirtz, M., Phua, S. Y., Estavillo, G. M. & Pogson, B. J. Balancing metabolites in drought: the
688 sulfur assimilation conundrum. *Trends Plant Sci.* **18**, 18–29 (2013).
- 689 12. Satyan, K. S. *et al.* Phenethyl isothiocyanate (PEITC) inhibits growth of ovarian cancer cells by inducing
690 apoptosis: Role of caspase and MAPK activation. *Gynecol. Oncol.* **103**, 261–270 (2006).
- 691 13. Herr, I., Buechler, M. W. & Büchler, M. W. Dietary constituents of broccoli and other cruciferous
692 vegetables: Implications for prevention and therapy of cancer. *Cancer Treat. Rev.* **36**, 377–383 (2010).
- 693 14. Rungapamestry, V., Duncan, A. J., Fuller, Z. & Ratcliffe, B. Effect of cooking brassica vegetables on the
694 subsequent hydrolysis and metabolic fate of glucosinolates. *Proc. Nutr. Soc.* **66**, 69–81 (2007).
- 695 15. Bell, L., Yahya, H. N., Oloyede, O. O., Methven, L. & Wagstaff, C. Changes In Rocket Salad
696 Phytochemicals Within The Commercial Supply Chain: Glucosinolates, Isothiocyanates, Amino Acids And

- 697 Bacterial Load Increase Significantly After Processing. *Food Chem.* **221**, 521–534 (2017).
- 698 16. Arias, T. & Pires, J. C. A fully resolved chloroplast phylogeny of the brassica crops and wild relatives
699 (Brassicaceae: Brassicaceae): Novel clades and potential taxonomic implications. *Taxon* **61**, 980–988 (2012).
- 700 17. Fechner, J. *et al.* The major glucosinolate hydrolysis product in rocket (*Eruca sativa* L.), sativin, is 1,3-
701 thiazepane-2-thione: Elucidation of structure, bioactivity, and stability compared to other rocket
702 isothiocyanates. *Food Chem.* **261**, 57–65 (2018).
- 703 18. Gigolashvili, T., Berger, B. & Flügge, U.-I. Specific and coordinated control of indolic and aliphatic
704 glucosinolate biosynthesis by R2R3-MYB transcription factors in *Arabidopsis thaliana*. *Phytochemistry*
705 *Reviews* **8**, 3–13 (2009).
- 706 19. Samuilov, S. *et al.* Knock-down of phosphoserine phosphatase gene effects rather N- than S-metabolism in
707 *Arabidopsis thaliana*. *Front. Plant Sci.* **9**, 1830 (2018).
- 708 20. Wang, M., Jia, Y., Xu, Z. & Xia, Z. Impairment of Sulfite Reductase Decreases Oxidative Stress Tolerance
709 in *Arabidopsis thaliana*. *Front. Plant Sci.* **7**, 1843 (2016).
- 710 21. Chan, K. X., Phua, S. Y. & Van Breusegem, F. Secondary sulfur metabolism in cellular signalling and
711 oxidative stress responses. *J. Exp. Bot.* **70**, 4237–4250 (2019).
- 712 22. Jiang, C. *et al.* Structural and functional comparative mapping between the Brassica A genomes in
713 allotetraploid *Brassica napus* and diploid *Brassica rapa*. *Theor. Appl. Genet.* **123**, 927–941 (2011).
- 714 23. Liu, S. *et al.* The *Brassica oleracea* genome reveals the asymmetrical evolution of polyploid genomes. *Nat.*
715 *Commun.* **5**, 499–507 (2014).
- 716 24. Cavaiuolo, M. *et al.* Gene expression analysis of rocket salad under pre-harvest and postharvest stresses: A
717 transcriptomic resource for *Diplotaxis tenuifolia*. *PLoS One* **12**, e0178119 (2017).
- 718 25. Bell, L., Methven, L. & Wagstaff, C. The influence of phytochemical composition and resulting sensory
719 attributes on preference for salad rocket (*Eruca sativa*) accessions by consumers of varying TAS2R38
720 diplotype. *Food Chem.* **222**, 6–17 (2017).
- 721 26. Falconer, D. S. & Mackay, T. F. C. *Introduction to quantitative genetics*. (Longman, 1996).
- 722 27. Weisenfeld, N. I. *et al.* Comprehensive variation discovery in single human genomes. *Nat. Genet.* **46**, 1350–
723 1355 (2014).
- 724 28. Leggett, R. M., Clavijo, B. J., Clissold, L., Clark, M. D. & Caccamo, M. NextClip: an analysis and read
725 preparation tool for Nextera Long Mate Pair libraries. *Bioinformatics* **30**, 566–568 (2014).
- 726 29. Li, R., Li, Y., Kristiansen, K. & Wang, J. SOAP: short oligonucleotide alignment program. *Bioinformatics*
727 **24**, 713–714 (2008).
- 728 30. Boetzer, M., Henkel, C. V., Jansen, H. J., Butler, D. & Pirovano, W. Scaffolding pre-assembled contigs
729 using SSPACE. *Bioinformatics* **27**, 578–579 (2011).
- 730 31. Jurka, J. *et al.* Repbase Update, a database of eukaryotic repetitive elements. *Cytogenet. Genome Res.* **110**,
731 462–7 (2005).
- 732 32. Xu, Z. & Wang, H. LTR_FINDER: an efficient tool for the prediction of full-length LTR retrotransposons.
733 *Nucleic Acids Res.* **35**, W265-8 (2007).

- 734 33. Kent, W. J. BLAT--the BLAST-like alignment tool. *Genome Res.* **12**, 656–64 (2002).
- 735 34. Birney, E., Clamp, M. & Durbin, R. GeneWise and Genomewise. *Genome Res.* **14**, 988–95 (2004).
- 736 35. Stanke, M. *et al.* AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Res.* **34**, W435-9
737 (2006).
- 738 36. Majoros, W. H., Pertea, M. & Salzberg, S. L. TigrScan and GlimmerHMM: two open source ab initio
739 eukaryotic gene-finders. *Bioinformatics* **20**, 2878–9 (2004).
- 740 37. Haas, B. J. *et al.* Automated eukaryotic gene structure annotation using EVidenceModeler and the Program
741 to Assemble Spliced Alignments. *Genome Biol.* **9**, R7 (2008).
- 742 38. Bairoch, A. & Apweiler, R. The SWISS-PROT protein sequence database and its supplement TrEMBL in
743 2000. *Nucleic Acids Res.* **28**, 45–8 (2000).
- 744 39. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **28**, 27–30
745 (2000).
- 746 40. Zdobnov, E. M. & Apweiler, R. InterProScan--an integration platform for the signature-recognition methods
747 in InterPro. *Bioinformatics* **17**, 847–8 (2001).
- 748 41. Suslow, T. Chlorination In The Production And Postharvest Handling Of Fresh Fruits And Vegetables. in
749 *Fruit and Vegetable Processing* (ed. McLaren, D.) 2–15 (Food Processing Center at the University of
750 Nebraska, 2000).
- 751 42. Huseby, S. *et al.* Diurnal and light regulation of sulphur assimilation and glucosinolate biosynthesis in
752 *Arabidopsis*. *J. Exp. Bot.* **64**, 1039–48 (2013).
- 753 43. Fleige, S. & Pfaffl, M. W. RNA integrity and the effect on the real-time qRT-PCR performance. *Mol.*
754 *Aspects Med.* **27**, 126–139 (2006).
- 755 44. Langmead, B. *et al.* Ultrafast and memory-efficient alignment of short DNA sequences to the human
756 genome. *Genome Biol.* **10**, R25 (2009).
- 757 45. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359
758 (2012).
- 759 46. Anders, S. *et al.* Differential expression analysis for sequence count data. *Genome Biol.* **11**, R106 (2010).
- 760 47. Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and
761 isoform switching during cell differentiation. *Nat. Biotechnol.* **28**, 511–515 (2010).
- 762 48. Wang, L., Feng, Z., Wang, X., Wang, X. & Zhang, X. DEGseq: an R package for identifying differentially
763 expressed genes from RNA-seq data. *Bioinformatics* **26**, 136–8 (2010).
- 764 49. Hu, R., Fan, C., Li, H., Zhang, Q. & Fu, Y.-F. Evaluation of putative reference genes for gene expression
765 normalization in soybean by quantitative real-time RT-PCR. *BMC Mol. Biol.* **10**, 93 (2009).
- 766 50. Livak, K. J. & Schmittgen, T. D. Analysis of Relative Gene Expression Data Using Real-Time Quantitative
767 PCR and the 2- $\Delta\Delta$ CT Method. *Methods* **25**, 402–408 (2001).
- 768 51. Lelario, F., Bianco, G., Bufo, S. A. & Cataldi, T. R. I. Establishing the occurrence of major and minor
769 glucosinolates in Brassicaceae by LC-ESI-hybrid linear ion-trap and Fourier-transform ion cyclotron
770 resonance mass spectrometry. *Phytochemistry* **73**, 74–83 (2012).

- 771 52. Cataldi, T. R. I., Rubino, A., Lelario, F. & Bufo, S. A. Naturally occurring glucosinolates in plant extracts of
772 rocket salad (*Eruca sativa* L.) identified by liquid chromatography coupled with negative ion electrospray
773 ionization and quadrupole ion-trap mass spectrometry. *Rapid Commun. Mass Spectrom.* **21**, 2374–2388
774 (2007).
- 775 53. Ku, K.-M., Kim, M. J., Jeffery, E. H., Kang, Y.-H. & Juvik, J. A. Profiles of Glucosinolates, their
776 Hydrolysis Products, and Quinone Reductase Inducing Activity from 39 Arugula (*Eruca sativa* Mill.)
777 Accessions. *J. Agric. Food Chem.* **64**, 6524–6532 (2016).

778

779 **Acknowledgments:**

780

781 **General:**

782 The authors would like to thank: the Vegetable Plant Breeding Team at Elsoms Seeds Ltd.;
783 members of the Genomics Pipelines Group in the BBSRC National Capability in Genomics and
784 Single Cell (BB/CCG1720/1) at Earlham Institute; Dr. Yunan Lin and Irene Wei for project and
785 technical support for genome annotation, RNAseq, and bioinformatics at Novogene Co. Ltd.;
786 Matthew Richardson for maintenance of controlled environment facilities at the University of
787 Reading Controlled Environment Laboratory; and Dr. Marcia Boura for advice on qRT-PCR.

788

789 **Funding:**

790 Dr. Luke Bell was supported by a BBSRC Case Award (BB/J012629/1) in partnership with Elsoms
791 Seeds Ltd. (Spalding, UK) and Bakkavor Group Ltd. (Spalding, UK) for *de novo* genome
792 sequencing and assembly. Dr. Luke Bell, Dr. Martin Chadwick, and Manik Puranik were
793 supported by a BBSRC LINK award (BB/N01894X/1) for all other work.

794

795 **Author contributions:**

796 LB and CW conceived and designed the experiment and analyses. RT and SK produced the
797 breeding line seed material for genome sequencing and the RNAseq experiment. LB grew plants
798 in controlled environment, performed RNA extractions and quality control, qRT-PCR validation,
799 glucosinolate analysis by LC-MS, and hydrolysis product analysis by GC-MS. MP performed
800 sugar analysis by HPLC. MC performed sulfur content analysis by ICP-OES. LB performed
801 ANOVAs, Pearson's correlation analyses, and Principal Component Analyses of all

802 phytochemical and gene expression data. LB wrote the paper, with contributions from MC, RT,
803 SK, and CW. Funding was obtained by LB, LM, and CW.

804

805 **Competing interests:**

806 The authors declare no conflicts of interest.

807

808 **Data and materials availability:**

809 Full reference genome sequence and annotation can be found in the European Nucleotide Archive.

810 Assembly accession no: GCA_902460325; Study ID: PRJEB34051; Sample ID: ERS3673677;

811 Annotation accession number ERZ1066251.

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

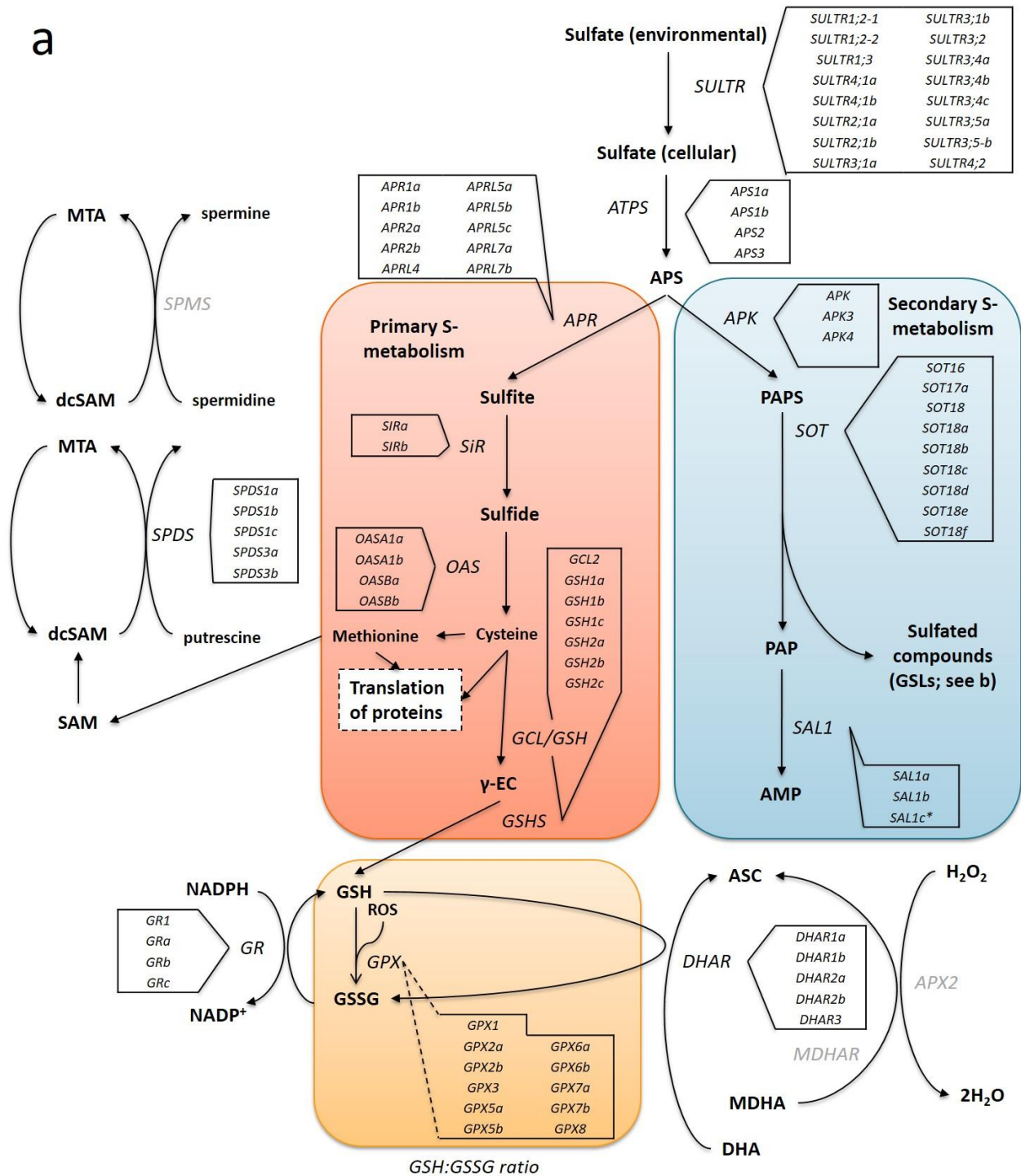
840

841

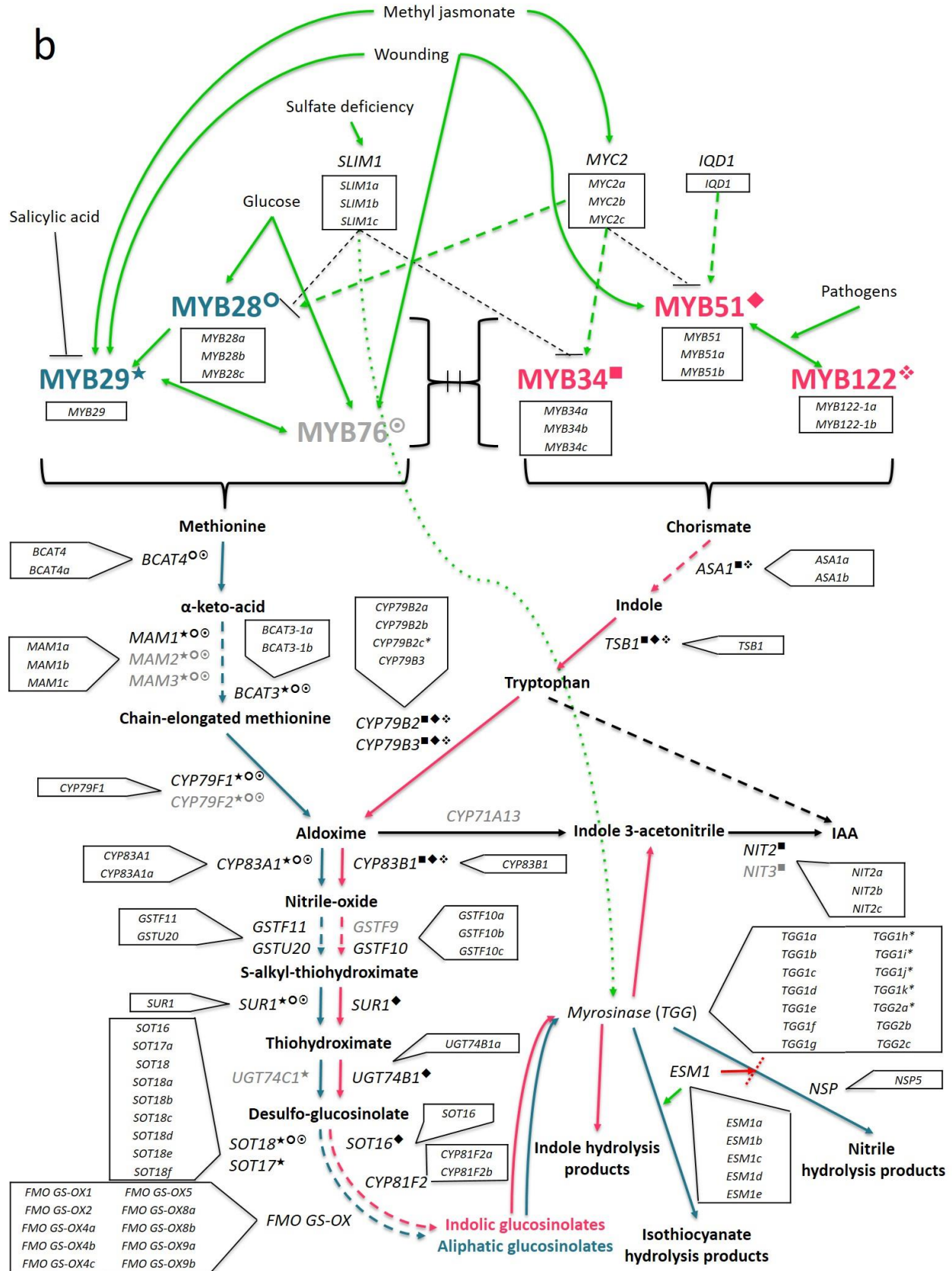
842

843 **Figures and Tables**
844

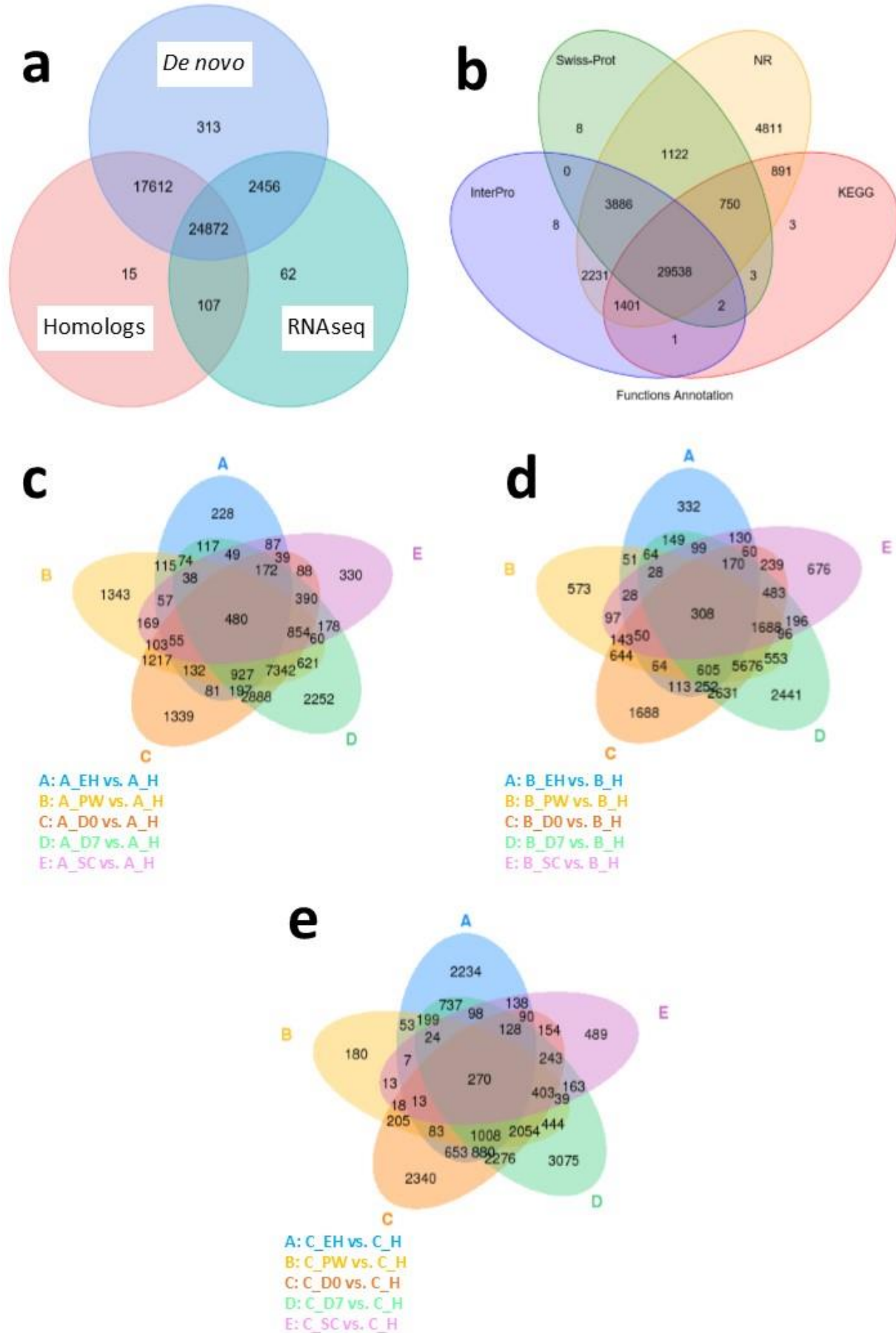
a



845



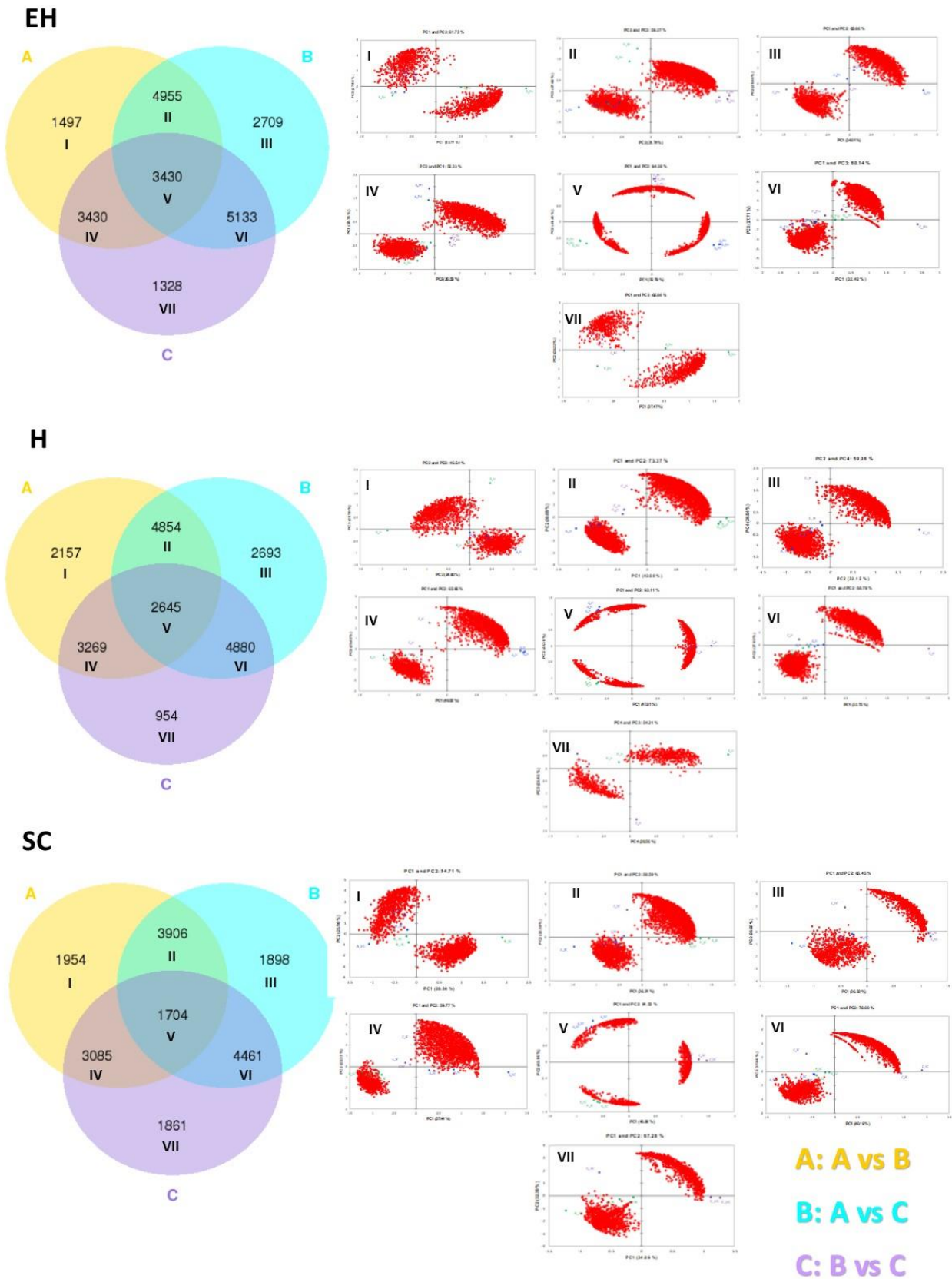
847
848 **Fig. 1. Sulfur metabolism and glucosinolate biosynthesis pathways.**
849 The primary (red box) and secondary (blue box) sulfur metabolism pathways (a) of *Arabidopsis*
850 *thaliana* with identified homologous genes within the *Eruca sativa* genome annotation (see boxed
851 insets) adapted from Chan et al. ¹¹. Environmental sulfur is assimilated and integrated into key
852 amino acids (cysteine and methionine) and enzymes. Sulfur metabolism is also intrinsically linked
853 with oxidative stress via glutathione synthesis. Under stress conditions 5'-phosphoadenosine-3'-
854 phosphate (PAP), glutathione disulfide (GSSG), and reduced glutathione (GSH) direct sulfate
855 towards GSH production. The GSH:GSSG redox state ratio is also known to influence sulfur
856 assimilation rates (orange box). *SOT* (*sulfotransferase*) genes link secondary sulfur metabolism
857 with the final sulfation step of GSL biosynthesis, and it is thought that *SAL1* plays an important
858 role in regulating the activity of these genes through interaction with PAP. GSL biosynthesis (b;
859 adapted from Gigolashvili et al. ¹⁷) is initiated by a complex and interacting network of abiotic and
860 biotic factors. Aliphatic synthesis pathway shown in teal, and the indolic pathway shown in pink,
861 is regulated by R2R3-MYB transcription factors. Known interactions between MYBs and specific
862 genes within each respective pathway are highlighted as follows: ● = MYB28, ★ = MYB29, ◎
863 = MYB76, ■ = MYB34, ◆ = MYB51, ❖ = MYB122. Genes with identified orthologs in the *E.*
864 *sativa* genome annotation are written in black; those with no identified homologous sequence are
865 written in grey. Abbreviations: *SULTR*, sulfate transporter; ATP, adenosine triphosphate; *ATPS*,
866 *ATP sulfurylase*; *APR*, *APS reductase*; *APK*, *APS kinase*; *SiR*, *sulfite reductase*; *OASTL*, *O-*
867 *acetylserine lyase*; *GCL*, *glutamate cysteine ligase*; γ -EC, γ -glutamyl-cysteine; *GSHS*, *GSH*
868 *synthetase*; *GR*, *glutathione reductase*; *GPX*, *glutathione peroxidase*; ASC, ascorbate; DHA,
869 dehydroascorbate; *DHAR*, *DHA reductase*; MDHA, monodehydroascorbate; *MDHAR*, *MDHA*
870 *reductase*; *APX2*, *ascorbate peroxidase 2*; H₂O₂, hydrogen peroxide; PAPS, 5'-phosphoadenosine-
871 3'-phosphosulfate; SAM, S-adenosyl methionine; dcSAM, decarboxylated SAM; MTA,
872 methylthioadenosine; *SPDS*, *spermidine synthase*; *SPMS*, *spermine synthase*; *SLIM1*, *sulfur*
873 *limitation 1*; *IQD1*, *protein IQ domain 1*; *BCAT*, *methionine aminotransferase*; *MAM*,
874 *methylthioalkylmalate synthase*; *CYP*, *cytochrome P450*; *GST*, *glutathione-S-transferase*; *SUR1*,
875 *C-S lyase 1*; *UGT*, *UDP-glycosyltransferase*; *FMO_{GS-OX}*, *flavin-containing monooxygenase*;
876 *ASA1*, *anthranilate synthase alpha subunit 1*; *TSB1*, *tryptophan synthase beta chain 1*; IAA,
877 indole-3-acetic acid; *NIT*, *nitrilase*; *ESM1*, *epithiospecifier modifier protein 1*; *NSP*, *nitrile*
878 *specifier protein*.
879



881 **Figure 2. *Eruca sativa* genome annotation and differentially expressed gene number Venn**
882 **diagrams for each inbred line**

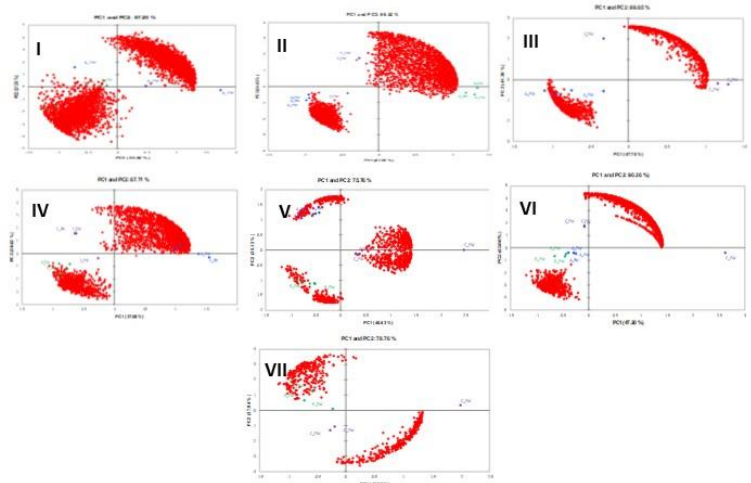
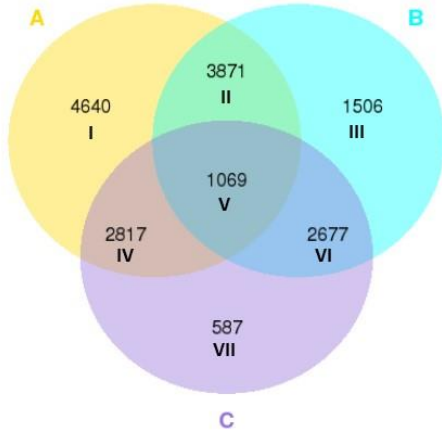
883 Venn diagrams of the *Eruca sativa* reference genome annotation gene identification sources (**a**)
884 and functional annotation databases used to assign putative gene identities (**b**). Also shown are
885 Venn diagrams of global differentially expressed genes (DEGs) at an early harvest (**EH**), second
886 harvest (**SC**), pre-wash (**PW**), post-wash (**D0**), and seven-day shelf life (**D7**) time points relative
887 to a first harvest (**H**) time point of three elite breeding lines: **A** (**c**), **B** (**d**), and **C** (**e**). The numbers
888 of DEGs identified under each condition are contained within the ellipses and their overlaps.
889

a

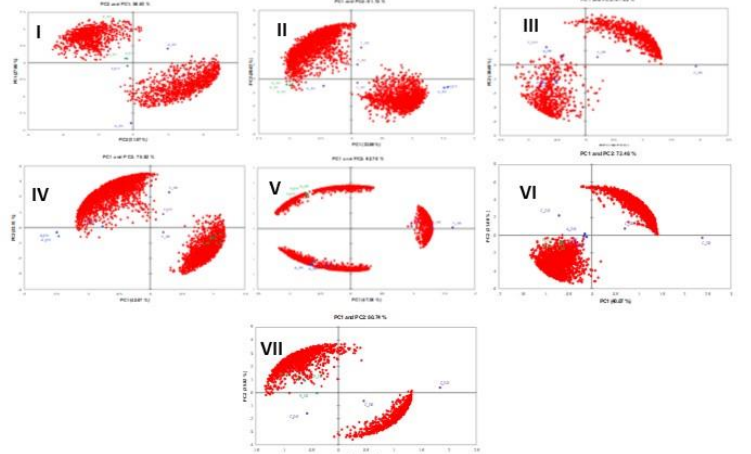
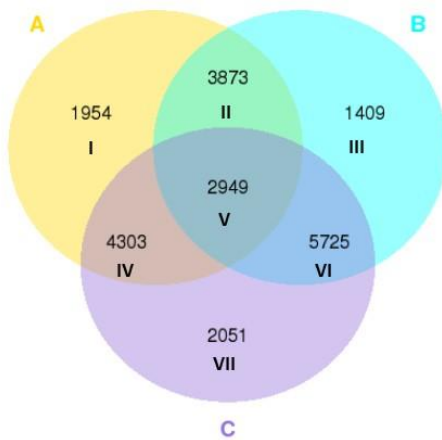


b

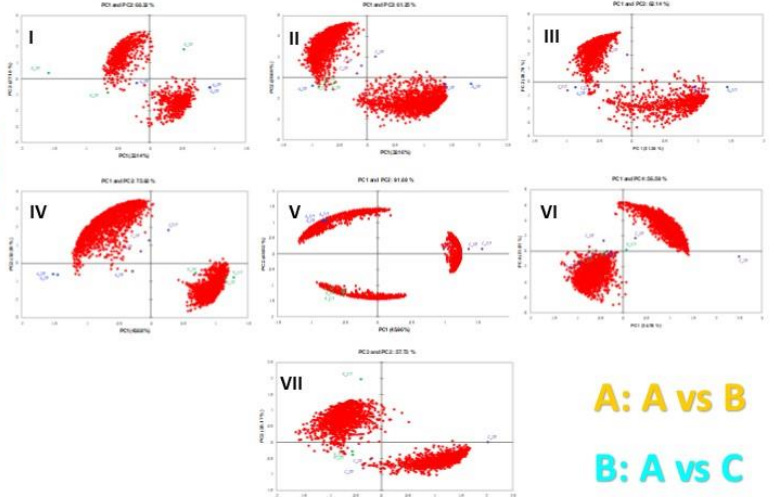
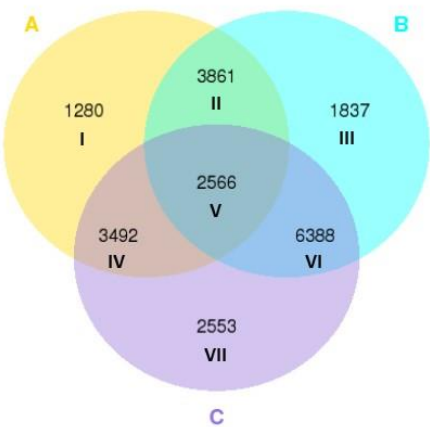
PW



D0



D7



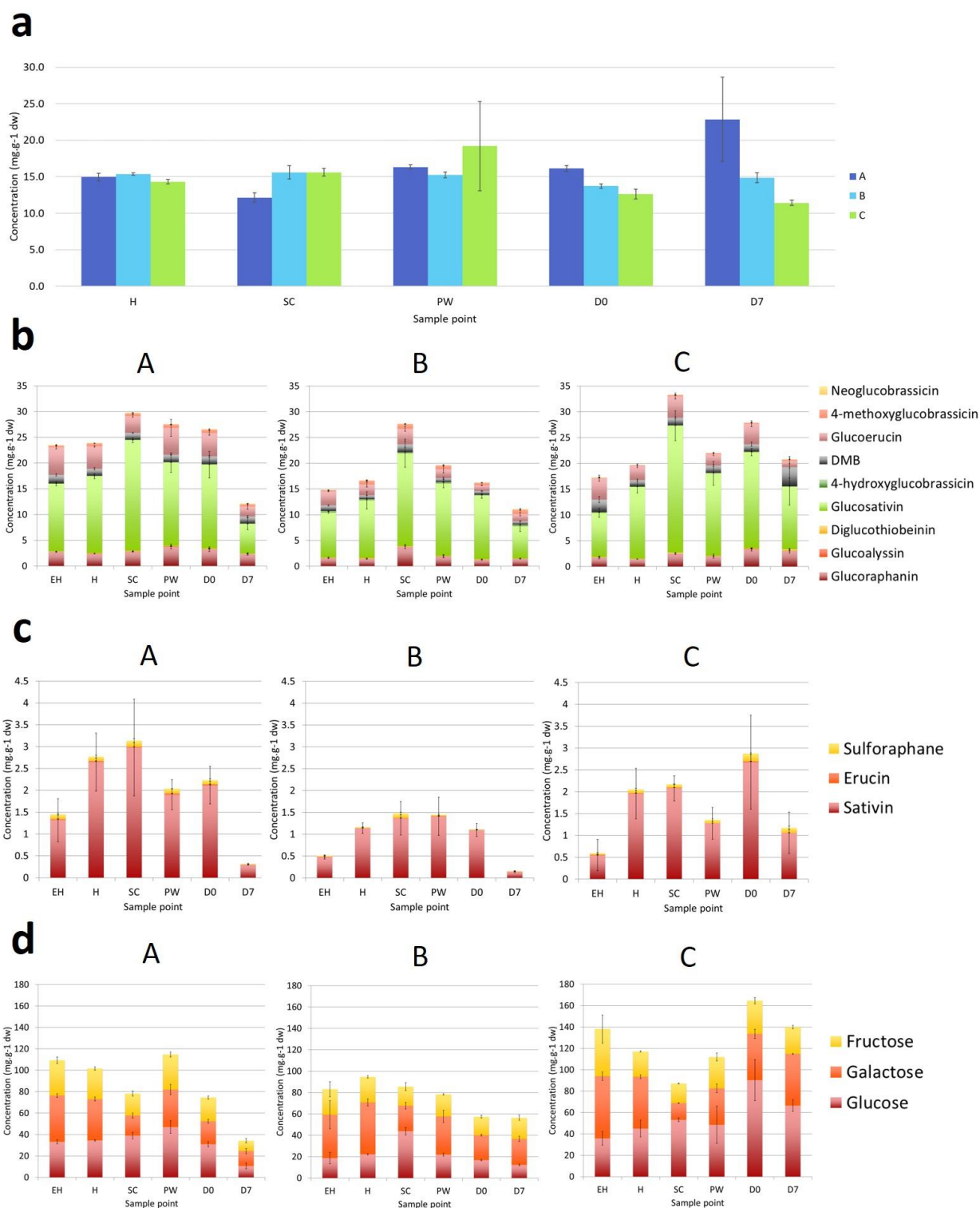
A: A vs B

B: A vs C

C: B vs C

892 **Figure 3. Venn diagrams of global differentially expressed genes between inbred lines with**
893 **Principal Component Analyses**

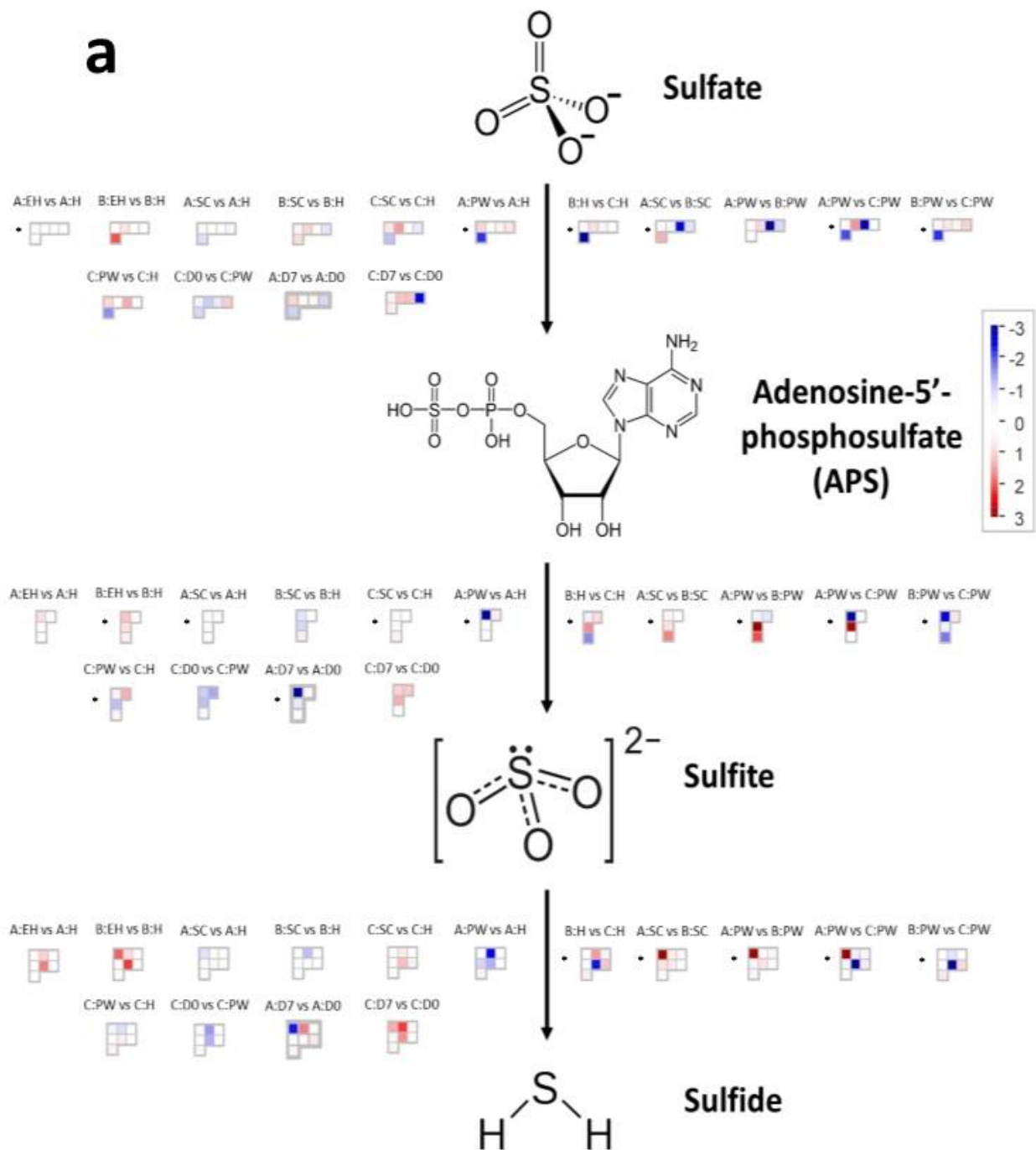
894 Venn diagrams of the numbers of differentially expressed genes (DEGs) found between three elite
895 inbred lines of *Eruca sativa*. Different harvest ontogeny (**a**) and postharvest (**b**) time points are
896 presented with accompanying Principal Component Analysis (PCA) plots that highlight the
897 multidimensional separations of each DEG cluster between lines **A**, **B**, and **C**. Roman numerals
898 are indicative of corresponding PCA analysis of each Venn diagram segment. Red data points
899 within each loadings-scores biplot represent gene expression data (FPKM). Blue circles = **A**; green
900 circles = **B**; purple circles = **C**. See insets for Venn diagram color coding. Abbreviations: early
901 harvest (**EH**), harvest (**H**), second harvest (**SC**), pre-wash (**PW**), post-wash (**D0**), and seven-day
902 shelf life (**D7**).
903



904
905
906
907
908

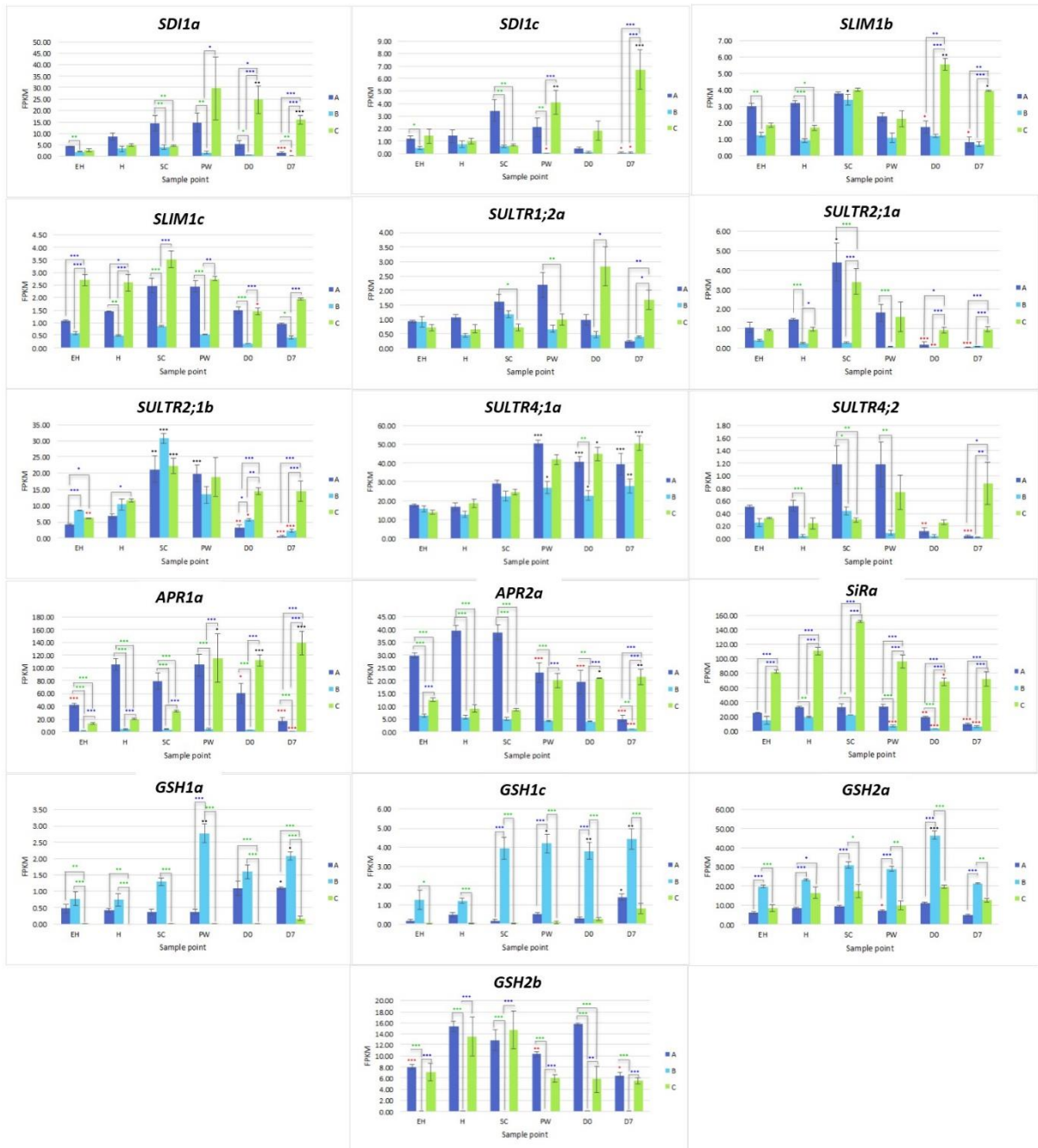
Figure 4. Elemental sulfur and phytochemical compositions of rocket inbred lines
Elemental sulfur (a), glucosinolate (b), glucosinolate hydrolysis product (c), and monosaccharide (d) concentrations observed in elite inbred lines of *Eruca sativa* (A, B, and C). Concentrations are

909 expressed as $\text{mg}\cdot\text{g}^{-1}$ of dry weight. Error bars represent standard error of the mean of each analyte
 910 detected. See insets for compound color coding. For ANOVA and Tukey's HSD pairwise
 911 significance values see Supplementary Data File S1. Abbreviations: early harvest (EH), harvest
 912 (H), second harvest (SC), pre-wash (PW), post-wash (D0), and seven-day shelf life (D7).
 913
 914



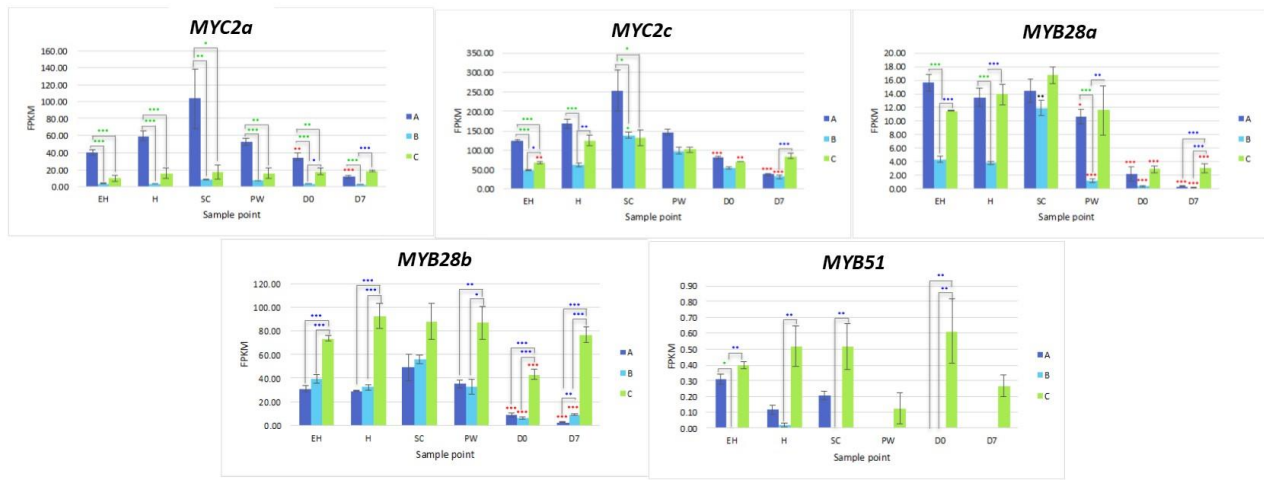
915

b



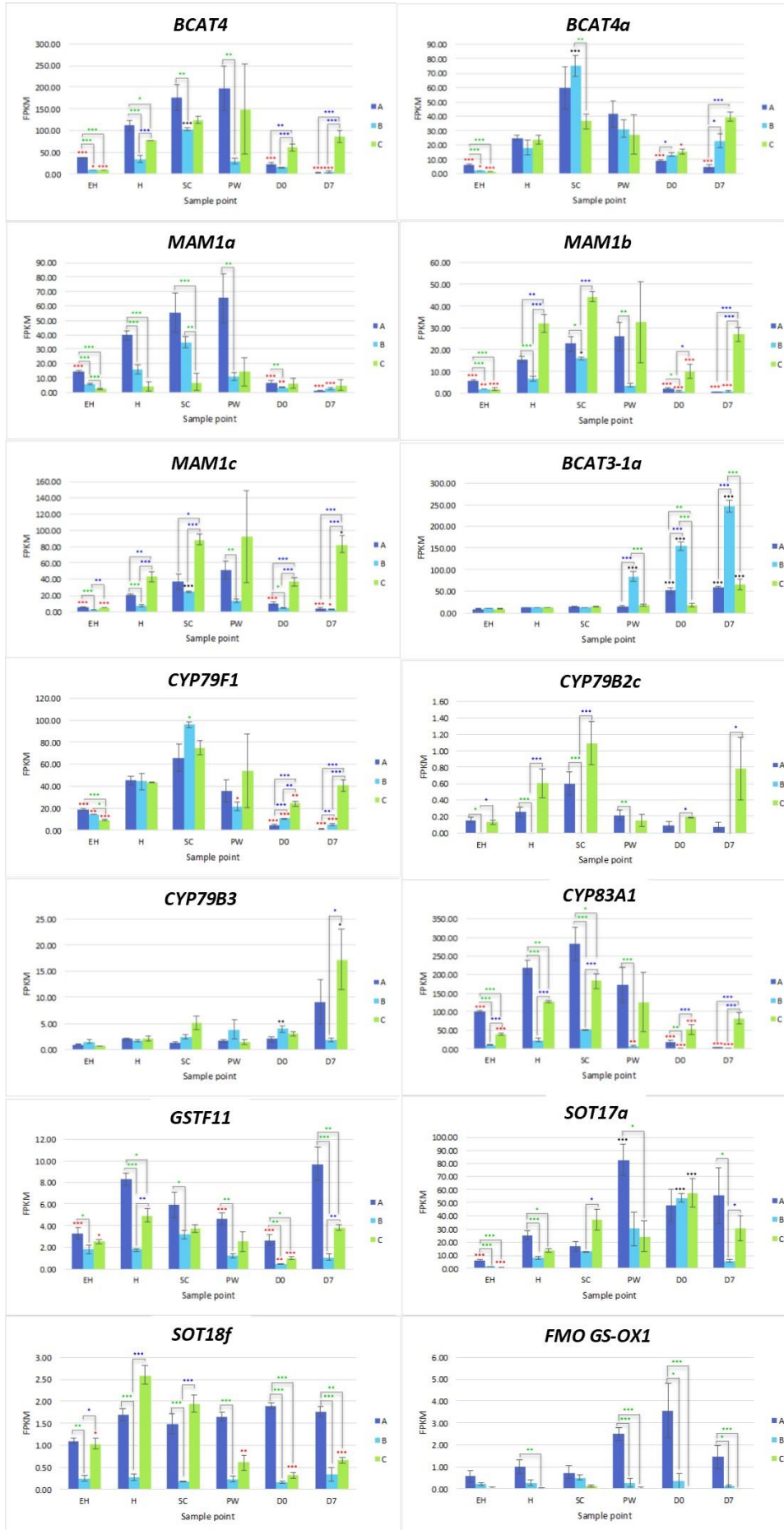
916

C

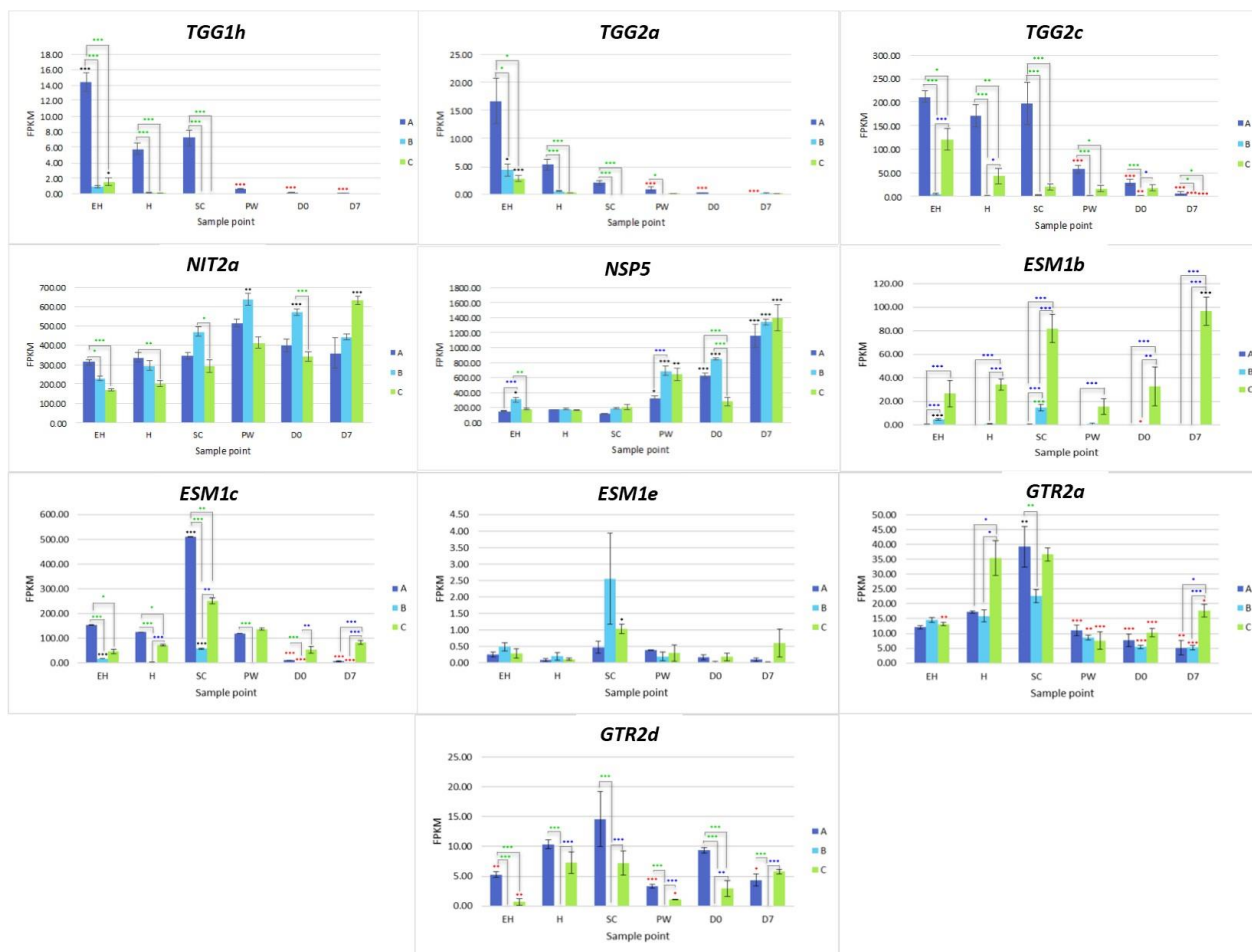


917

d



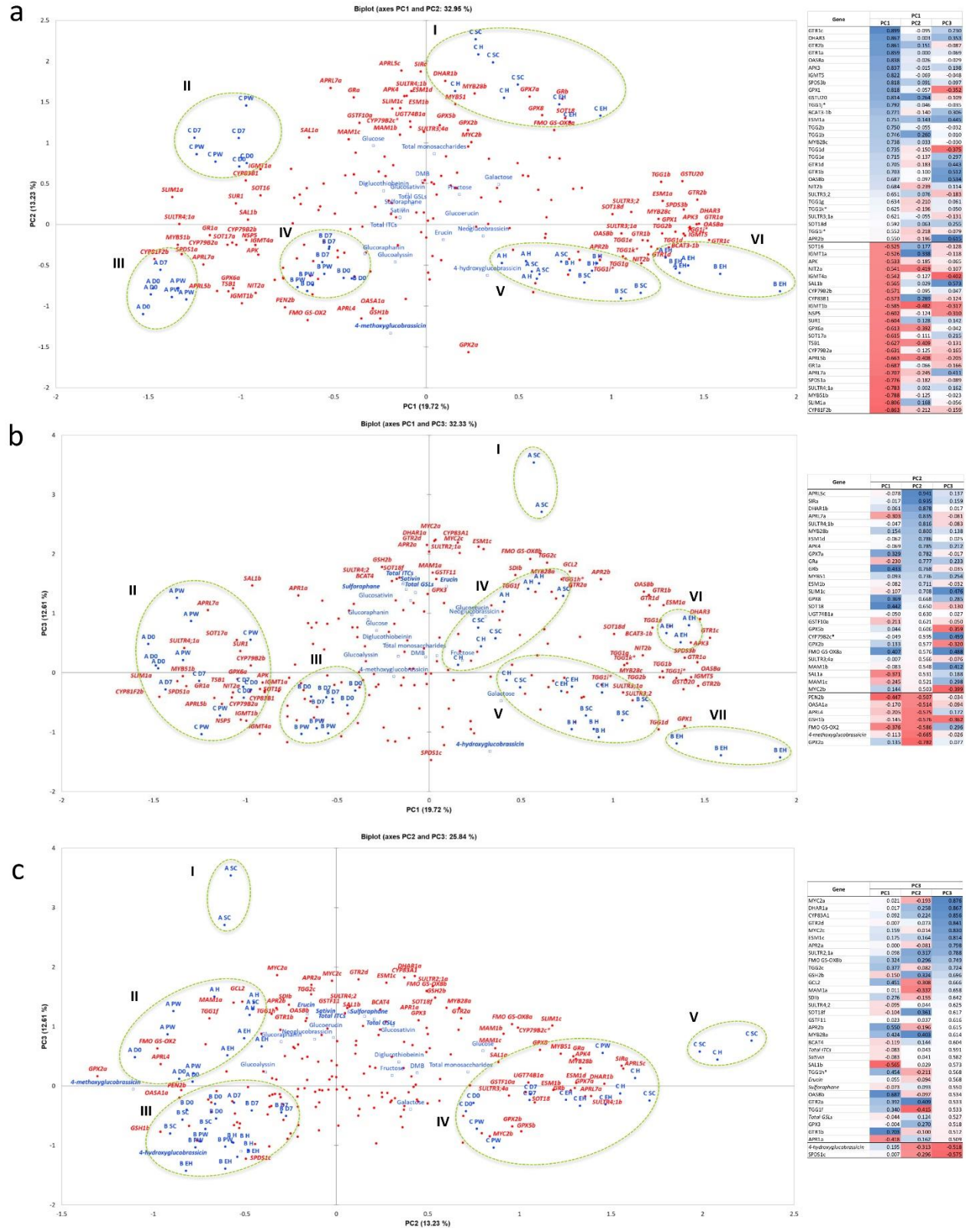
e



919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937

Figure 5. RNaseq expression data for sulfur metabolism and glucosinolate biosynthesis-related genes

RNaseq expression data (FPKM) for genes involved with sulfate assimilation (a), sulfate transport and redox response (b), glucosinolate transcription factors (c), glucosinolate biosynthesis (d), and glucosinolate hydrolysis and transport (e) in three elite inbred lines of *Eruca sativa* (A = dark blue; B = light blue; C = green). For (a) a custom MapMan (version 3.6.0RC1) annotation file of the *E. sativa* reference genome was created using Mercator4 (version 1.0, plaBi dataBase, Institute of Biology, Aachen, Germany), and used to visualize differential expression of genes within the sulfate assimilation pathway. In (b) to (e), standard errors of the mean expression values are represented by error bars. Asterisks denote levels of significance of up and down regulation within sample points (between each inbred line) and relative to the point of harvest for each respective sample point: * = $P \leq 0.05$; ** = $P \leq 0.01$; *** = $P \leq 0.001$; green = significant up regulation between lines A, B, and C; blue = significant down regulation between lines A, B, and C; black = significant up regulation relative to H; red = significant down regulation relative to H. Abbreviations: early harvest (EH), harvest (H), second harvest (SC), pre-wash (PW), post-wash (D0), and seven-day shelf life (D7).



939 **Figure 6. Principal Component Analysis of sulfur assimilation pathway, glucosinolate**
 940 **biosynthesis, and glucosinolate hydrolysis gene expression data**

941 Principal Component Analysis of sulfur assimilation pathway, glucosinolate biosynthesis, and
 942 glucosinolate hydrolysis gene expression data (FPKM) for three *Eruca sativa* elite inbred lines (**A**,
 943 **B**, and **C**) across ontogenic and postharvest sample points. Biplot (**a**) displays Principal
 944 Components (PCs) 1 and 2, which represent 33% of variation within the data. Biplot (**b**) displays
 945 PC1 and PC3, explaining 32.3% of variation within the data; and (**c**) displays PC2 and PC3,
 946 explaining 25.8% of the variability. The PCA plots presented are the results of Varimax rotation.
 947 Each biplot is accompanied by a factor loadings table sorted according to PC1 (**a**), PC2 (**b**), and
 948 PC3 (**c**); italics denotes a supplementary variable, and * = a putative novel gene within the
 949 reference annotation. Blue coloration denotes high factor loading scores, red denotes low. Only
 950 genes with loading values >0.5 were included, and each is represented within the biplots in red
 951 (bold italics). Red circles represent individual genes included in the analysis ($n = 177$). Blue circles
 952 represent sample point variables and have accompanying labels (blue bold). Blue squares denote
 953 phytochemical data regressed onto the PCA as supplementary variables. Bold data labels indicate
 954 phytochemical components with >0.5 factor loadings scores. Green dotted ellipses denote clusters
 955 of variables and are numbered using Roman numerals, which are quoted within the text.
 956 Abbreviations: early harvest (**EH**), harvest (**H**), second harvest (**SC**), pre-wash (**PW**), post-wash
 957 (**D0**), and seven-day shelf life (**D7**).

958

959 **Supplementary Materials**

960

961 **Supplementary Table S1**

962

Table S1. Summary of genome assembly and annotation of *Eruca sativa*

Genome assembly	>= 0 bp	>= 1,000 bp	Largest contig	Total (>= 500 bp)				
Contig number	1,041,818	12,352	1,477,633	49,933				
Total length	850,956,505	562,271,846		586,731,295				
Assembly-related statistics								
GC%	N50	NG50	N75	NG75	L50	LG50	L75	LG75
36.25	196,831	136,378	87,576	2,634	789	1,256	1,889	7,243

963

964

965

966 **Supplementary Table S2**

967

Table S2. Transposable elements content in the reference genome

Type	Denovo + Repbase*		TE proteins [§]		Combined TEs [£]	
	Length (bp)	% in genome	Length (bp)	% in genome	Length (bp)	% in genome
DNA	69,251,054	8.14	21,607,510	2.54	76,517,426	8.99
LINE	20,290,781	2.38	17,153,783	2.02	28,200,567	3.31
SINE	2,134,305	0.25	0	0	2,134,305	0.25
LTR	311,377,915	36.59	91,001,347	10.69	317,124,290	37.27
Other [^]	106,176	0.01	0	0	106,176	0.01
Unknown ^{^^}	155,033,031	18.22	0	0	155,033,031	18.22
Total	547,675,259	64.36	129,571,414	15.23	563,873,839	66.26

* = RepeatMasker based on the uclust algorithm combined with the known Repbase and *de novo* repeat library created by RepeatModeler / Repeat Scout / LTR_finder

§ = RepeatProteinMask based on Repbase

£ = The non-redundant set of results combining Denovo+Repbase TEs and TE proteins

^ = Repeats that can be classified by RepeatMasker, but not included by classes above

^^ = Repeats that could not be classified by RepeatMasker

968

969

970

971

972

973

974

975

976

977

978

Table S3. Predicted protein-coding genes within the *E. sativa* reference genome

Gene set	Number	Average gene length (bp)	Average CDS length (bp)	Average exons per gene	Average exon length (bp)	Average intron length (bp)	
Augustus	50,179	1,701.56	1,024.89	4.48	228.88	194.57	
Glimmer HMM	73,989	1,335.36	725.53	3.03	239.16	299.86	
<i>De novo</i> *	SNAP	80,264	1,231.13	728.3	4.02	180.98	166.27
Geneid	100,165	2,127.25	585.82	3.07	191.03	745.87	
Genscan	71,813	3,942.04	785.32	3.92	200.11	1,079.4	
Homolog [^]	<i>Arabidopsis lyrata</i>	32,667	1,867.18	1,084.09	4.86	223.12	202.93
	<i>Arabidopsis thaliana</i>	27,416	1,870.34	1,218.4	5.13	237.58	157.91
	<i>Brassica napus</i>	101,040	1,764.75	1,001.16	4.91	204.06	195.48
	<i>Boechera stricta</i>	27,416	2,006.68	1,181.2	5.09	231.86	201.61
	<i>Capsella rubella</i>	26,521	1,958.82	1,248.6	5.19	240.53	169.46
	<i>Raphanus sativus</i>	49,733	2,064.75	1,194.41	4.94	241.57	220.66
RNAseq ^{^^}	Cufflinks	43,200	2,848.16	1,723.58	6.03	285.65	223.41
	PASA	37,870	1,744.08	1,034.72	4.77	216.9	188.13
EVM	59,643	1,665.12	926.2	4.17	222.2	233.23	
PASA-update	59,491	1,656.25	929.33	4.17	222.9	229.36	
Final set	45,438	1,889.6	1,069.44	4.76	224.81	218.3	

* = The combined results by EVM of *5-ab initio* gene predictions

[^] = The combined results by EVM of homology-based gene prediction

^{^^} = The combined results by EVM of transcriptome data sets

979 **Supplementary Table S4**

980

981

982

Table S4. Numbers of genes with homology or functional assignment

Database	Number of annotated genes	Gene annotation %
NR	44,630	98.2
Swiss-Prot	35,309	77.7
KEGG	32,589	71.7
InterPro	All	81.6
	Pfam	76
	GO	56.8
Annotated	44,655	98.3
Total	45,438	-

983

984

985

986

987

988

989

990

991

992

993

994

995

996

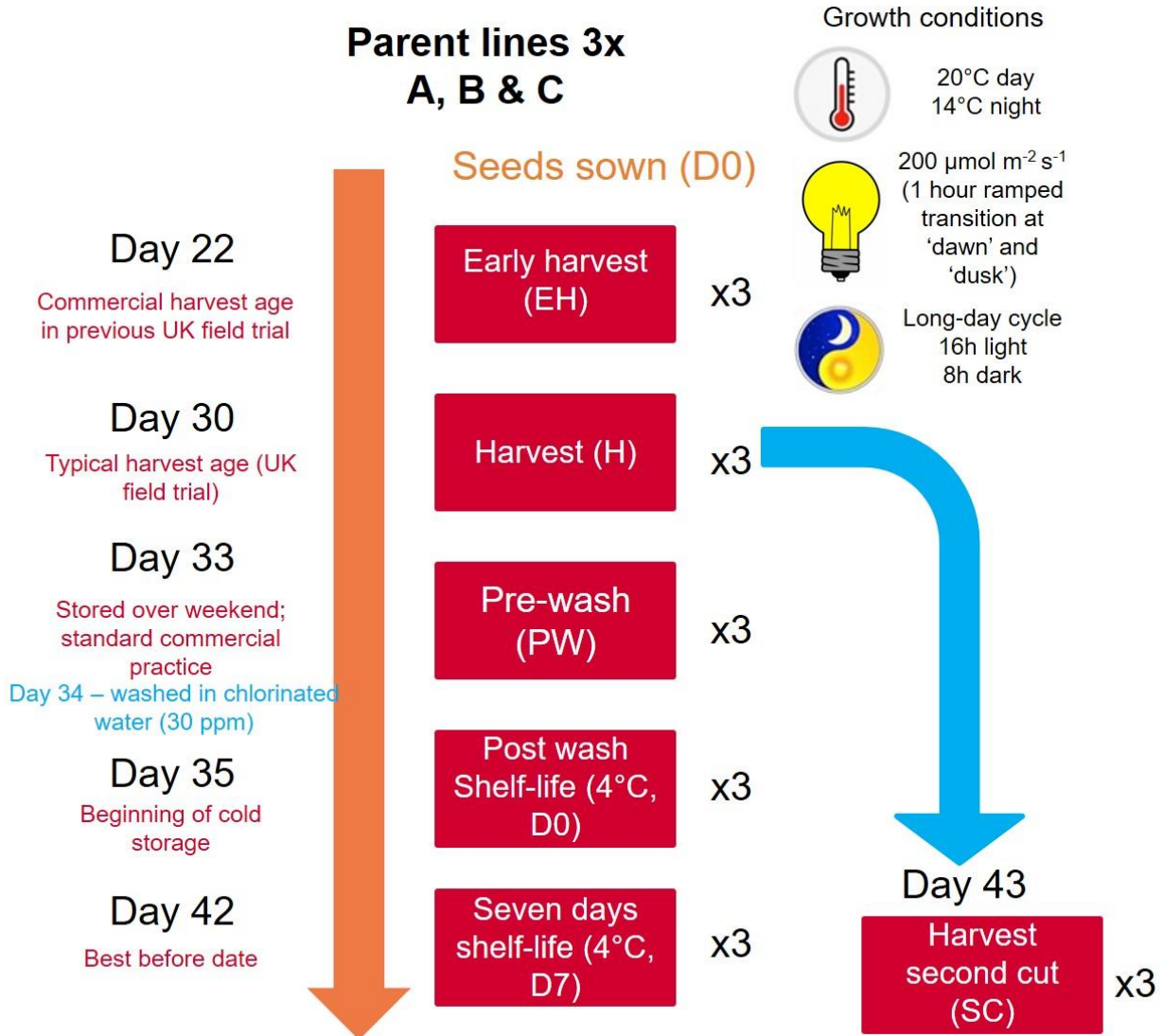
997

998 **Supplementary Table S5**

Table S5. qRT-PCR primers and efficiencies

Gene	Oligo	Length	tm	GC%	Sequence	Product size	Efficiency
<i>MYB122-1a</i>	LEFT	20	62.45	60.00	GGCGAACCTACCCGACAAA	163	2.13
	RIGHT	23	60.31	43.48	GTGGACCATTTGTTGCCATGAAT		
<i>MYB51a</i>	LEFT	19	60.00	57.89	GGCGAACTCTCCCGAAAA	164	1.93
	RIGHT	21	62.29	52.38	TGCAGCCCATTTGTTTCCGTG		
<i>BCAT4</i>	LEFT	20	62.35	55.00	TCGTCTCCGCCGTCAAACAA	198	1.91
	RIGHT	20	61.55	55.00	ACCCCGCGTTATCCTTGTA		
<i>SOT16</i>	LEFT	21	63.06	57.14	CCCAACACAACGGACTGTT	187	1.77
	RIGHT	20	62.74	60.00	AGAGGGTTCGTTGCGTCGTC		
<i>CYP83B1</i>	LEFT	23	59.73	43.48	TTTGATATTGTTGTACCCGGGA	150	2.03
	RIGHT	20	59.12	55.00	TCTCTTCCGAGACGTGTCC		
<i>SUR1</i>	LEFT	22	61.94	50.00	TCATTCAGGCTGCACTTCTCA	105	1.9
	RIGHT	23	61.17	47.83	GCCTATCACACCAAATCGACA		
<i>TGG1b</i>	LEFT	20	61.91	55.00	AGCTTCTCATGGCCTCGCT	190	1.93
	RIGHT	22	63.20	59.09	CCCCTCCCTCCTCAATCTGGT		
<i>TGG1j</i>	LEFT	21	61.51	52.38	TGGTACATGGAGCCGTTACA	159	2.08
	RIGHT	21	61.76	57.14	TGACTGGGCGTACTGAGTGAC		
<i>UGT74B1a</i>	LEFT	20	62.99	65.00	CAGCATCGACGCTACTCCG	247	1.71
	RIGHT	20	61.26	60.00	GGGAACAGCGGGGAGAGTAA		
<i>TGG1d</i>	LEFT	21	62.26	52.38	TCAGAAGACCGTTGCCAAGCT	104	1.84
	RIGHT	22	61.55	50.00	ACGTTGACACCTTCTCCTTGA		
<i>ACT11</i>	LEFT	20	62.99	60.00	TTCACCACCACAGCAGAGCG	165	2.19
	RIGHT	20	62.59	65.00	CCTCTCCCCTCCGATGGTGA		

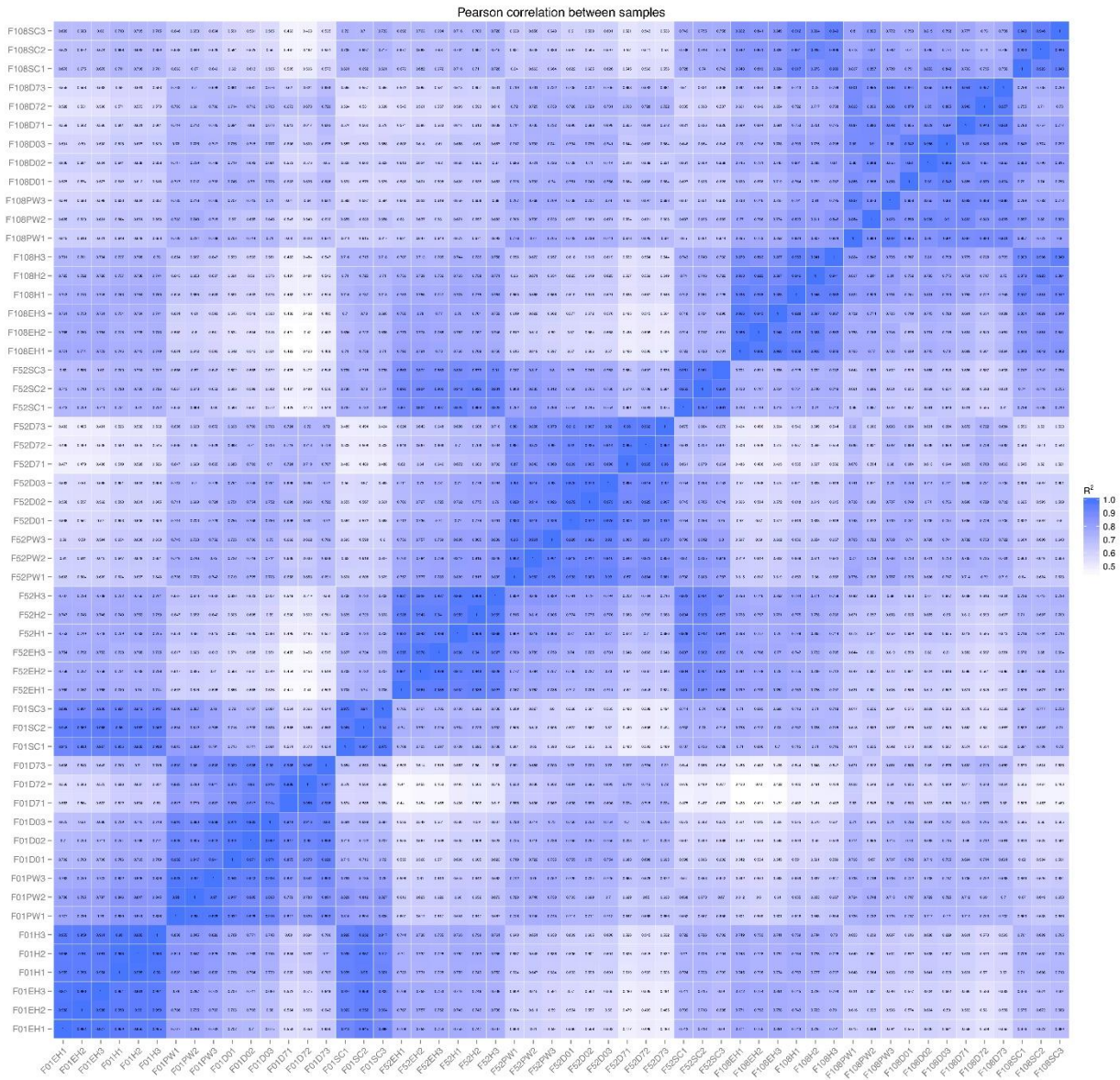
999
1000



1001
1002
1003
1004
1005
1006
1007
1008

Supplementary Figure S1

RNAseq experimental design and sampling diagram. Three elite inbred lines of *Eruca sativa* were grown under controlled environment conditions and sampled at each of the six time points indicated (in triplicate). Abbreviations: early harvest (**EH**), harvest (**H**), second harvest (**SC**), pre-wash (**PW**), post-wash (**D0**), and seven-day shelf life (**D7**).



1009
1010
1011
1012
1013
1014
1015
1016

Supplementary Figure S2

Pearson correlation matrix of RNAseq biological sample replicate gene expression values. Replicates of each sample showed a high degree of correlation ($r^2 = >0.884$) indicating robust reproducibility of gene expression between the individual plants tested at each respective sample point.



1017
1018

1019 **Supplementary Figure S3**

1020 qRT-PCR (green) vs. RNAseq (orange) gene expression of ten randomly selected glucosinolate
1021 biosynthesis and hydrolysis-related genes. Data are expressed as the normalized log₂ fold-change
1022 in expression relative to the reference gene *ACT11*. ANOVA revealed no significant difference
1023 between the two data sets. Abbreviations: early harvest (**EH**), harvest (**H**), second harvest (**SC**),
1024 pre-wash (**PW**), post-wash (**D0**), and seven-day shelf life (**D7**).

1025
1026 **Supplementary Data File S1**

1027 Analysis of Variance outputs with Tukey's HSD pairwise comparisons between sample points and
1028 each respective rocket breeding line: Tab 1 – glucosinolate analysis; Tab 2 – glucosinolate
1029 hydrolysis product analysis; Tab 3 – sugar analysis. Highlighted values are significant at the
1030 following levels: $P < 0.05$ (yellow), $P = 0.01$ (orange), and $P = 0.001$ (green). Tab 4 contains a
1031 Pearson's correlation analysis matrix for sulfur and glucosinolate-related gene expression values
1032 and phytochemical observations. Values in bold are significant correlations at the $P = 0.001$
1033 threshold.

1034
1035 **Supplementary Data File S2**

1036 RNAseq read counts, log₂-fold changes, P -values, and adjusted P -values (padj) for sulfur
1037 metabolism, glucosinolate biosynthesis, hydrolysis, and transport genes for each of the three rocket
1038 lines and the respective sample points. Significant up/down regulation is denoted by green/red
1039 highlighting, respectively. KEGG annotation numbers and UniProt gene descriptions for
1040 orthologous genes in *Arabidopsis thaliana* are provided.