# Journal Pre-proof

Exploring glucosinolates diversity in Brassicaceae: A genomic and chemical assessment for deciphering abiotic stress tolerance

Anyse Pereira Essoh, Filipa Monteiro, Ana Rita Pena, M. Salomé Pais, Mónica Moura, Maria Manuel Romeiras

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

**Brassicaceae diversity to decipher abiotic stress tolerance**

Separation between genes in GLS core structure and *CYP450/MYB* gene families.

Recent diversification of aliphatic genes and an earliest for indolic genes

Distinct GLS chemo-profile between *Brassica* crops and *Diplotaxis* species (wild relatives)

*Glucosinolates* (GLS)

Are GLS related with abiotic stress tolerance?

Which genes are involved in GLS biosynthesis?

Can GLS pathways diversification be observed?

Can chemo-profile differences be depicted between Brassicaceae taxa?

UPGMA / PCA

Phylogenetic analysis CYP79 F1-F2 / CYP81 F1-F4

Chemodiversity profile

# Exploring glucosinolates diversity in Brassicaceae: a genomic and chemical assessment for deciphering abiotic stress tolerance

Anyse Pereira Essoh[1,2,3*], Filipa Monteiro[1,4* †], Ana Rita Pena[4], , M. Salomé Pais[5], Mónica Moura[2], Maria Manuel Romeiras[1,4,5 †]

[1]Linking Landscape, Environment, Agriculture and Food (LEAF), Instituto Superior de Agronomia, Universidade de Lisboa, Lisboa, Portugal.
[2]Research Centre in Biodiversity and Genetic Resources (CIBIO), InBIO Associate Laboratory, Faculdade de Ciências e Tecnologia, Universidade dos Açores, Ponta Delgada, Portugal.
[3]Nova School of Business and Economics, 2775-405 Campus de Carcavelos, Portugal.
[4]Centre for Ecology, Evolution and Environmental Changes (cE3c), Faculdade de Ciências,Universidade de Lisboa, Lisboa, Portugal.
[5]Academia das Ciências de Lisboa, Rua Academia das Ciências 19, 1200-168 Lisboa, Portugal.
* authors contributed equally to this work.

[†] **Correspondence:**
fimonteiro@fc.ul.pt; mmromeiras@isa.ulisboa.pt

## Running title: Brassicaceae diversity to decipher abiotic stress tolerance

## ABSTRACT

*Brassica* is one of the most economically important genus of the Brassicaceae family, encompassing several key crops like *Brassica napus* (cabbage) and broccoli (*Brassica oleraceae* var. *italic*a). This family is well known for their high content of characteristic secondary metabolites such as glucosinolates (GLS) compounds, recognize for their beneficial health properties and role in plants defense. In this work, we have looked through gene clusters involved in the biosynthesis of GLS, by combining genomic analysis with biochemical pathways and chemical diversity assessment. A total of 101 Brassicaceae genes involved in GLS biosynthesis were identified, using a multi-database approach. Through a UPGMA and PCA analysis on the 101 GLS genes recorded, revealed a separation between the genes mainly involved in GLS core structure synthesis and genes belonging to the *CYP450*s and *MYB*s gene families. After, a detailed phylogenetic analysis was conducted to better understand the disjunction of the aliphatic and indolic genes, by focusing on *CYP79F1-F2* and *CYP81F1-F4*, respectively. Our results point to a recent diversification of the aliphatic *CYP79F1* and *F2* genes in *Brassica* crops, while for indolic genes an earliest diversification is observed for *CYP81F1-F4* genes. Chemical diversity revealed that *Brassica* crops have distinct GLS chemo-profiles from other Brassicaceae genera; being highlighted the high contents of GLS found among the *Diplotaxis* species. Also, we have explored GLS-rich species as a new source of taxa with great agronomic potential, particularly in abiotic stress tolerance, namely *Diplotaxis*, the closest wild relatives of *Brassica* crops.

**Keywords:** chemical diversity; genomic diversity; GLS; abiotic stress; *Brassica* crops; *Diplotaxis*

## 1. Introduction

The Brassicaceae is one of the world's most economically important plant family (Ishida et al., 2014). It includes important crop species such as *Brassica oleracea* (e.g., cauliflower, Brussels sprouts, cabbage, broccoli, and Kai Lan), *Brassica rapa* (e.g., pakchoi, choy sum, and Chinese cabbage), *Nasturtium officinale* (e.g., watercress), and *Raphanus sativus* (e.g., daikon radish and red cherry radish). Other species such as *Diplotaxis tenuifolia* and *Eruca vesicaria*, commonly referred as 'rocket salads', have also attracted a considerable interest as culinary vegetables because of their strong flavor and content of putative health-promoting compounds (Verkerk et al., 2010). These species and their crop wild relatives (CWR – taxa closely related to crops) grown primarily in the Euro-Mediterranean region, which contains the highest proportion of agronomically important plants representing an important reservoir of genetic resources for crop improvement (Kell et al., 2008). CWR are likely to contain a great genetic diversity necessary to combat climate change because of the diversity of habitats in which they grow and the wide range of conditions they are adapted to (Ford-Lloyd et al., 2011).

Among the most important chemical compounds produced by Brassicaceae species are the Glucosinolates (GLS), which proved to have health promoting effects and importance in abiotic stress tolerance (Cartea and Velasco, 2007). They are constituted by a common structure comprising a β-D-thioglucose group, a sulfonated oxime moiety and a variable side-chain derived either from methionine, tryptophan, phenylalanine, or from other branched chain amino acids. GLS are found in 16 dicotyledonous plant families where, at least, 130 different structures have been identified so far (Fahey et al., 2001; Collett et al., 2014).

GLS are present at different concentrations throughout the plant organs. They can reach 1% of the dry weight in some tissues of *Brassica* (Fahey et al., 2001). Within a single species, up to 4 different GLSs dominate the GLS occurrence in the plant (Verkerk et al., 2008). The type, concentration and distribution of the GLS in the plants of Brassicaceae family vary according to a high number of factors, namely species (Bellostas et al., 2004), variety (Choi et al., 2014), plant organ (Brown et al., 2003; Bellostas et al., 2004) or plant age (Fahey et al., 1997; Brown et al., 2003) and developmental cycle. Moreover, environmental conditions such as season (Cartea et al., 2007), biotic (Verkerk et al., 2008) or abiotic stress factors such as salinity or drought, are also known to play a role on the production and content of these compounds (Khan et al., 2011; Martínez-Ballesta et al., 2015).

Recent studies have revealed that GLS and their derivatives have beneficial effects on humans. They can help in suppressing tumor growth of various types of cancers namely: breast, brain, blood, bone,

79   colon, gastric, liver, lung, oral, pancreatic and prostate (Zhang et al., 2003; Soundararajan and Kim,

80   2018). Significant reduction in plasma LDL-C levels has also been reported as being directly linked

81   to consumption of GLS-rich broccoli (Armah et al., 2015). Some GLS derived products are reported

82   to have antimicrobial effects and well documented health benefits (Cavaiuolo and Ferrante, 2014;

83   Bischoff, 2016). Exclusive or excessive feeding of vegetables and/or seeds from the *Brassica* plants

84   have been associated with toxic effects in livestock (VanEtten and Tookey, 1983; Tripathi and

85   Mishra, 2007) and strategies have been explored to reduce GLS content in *Brassica* vegetables to

86   increase their palatability for animal consumption (Verker et al., 2008).

87   The GLS biosynthetic pathway has been partially elucidated by studies on *Arabidopsis* (e.g. reviewed

88   in Grubb and Abel 2006; Halkier and Gershenzon 2006). The GLS, synthesized from amino acids,

89   are grouped in three subtypes according to their corresponding precursors: i) aliphatic GLS, derived

90   from alanine, leucine, isoleucine, valine, and methionine; ii) indole GLS, derived from tryptophan;

91   and iii) aromatic GLS, derived from phenylalanine and tyrosine (Fahey et al., 2001; Halkier and

92   Gershenzon, 2006). Different authors have reported on aliphatic GLS accounting for 70–97% of the

93   total GLS content in leaves of *Brassica oleracea* (Cartea et al., 2007), leaves and stems of *Brassica*

94   *napus* (Cleemput and Becker, 2012), leaves and seeds of *Brassica juncea* (Gupta et al., 2012;

95   Othmane, 2015), and sprouts and mature leaves of *Brassica rapa* (Wiesner et al., 2013). The

96   formation of the GLS core structure involves the action of enzymes from different families, namely

97   the CYP79 (Hansen et al., 2001; Chen et al., 2003), CYP83 (Bak and Feyereisen, 2001), UGT74

98   (Grubb et al., 2014), C-S-lyases (Mikkelsen et al., 2004) and of sulfotransferases (SOTs or STs)

99   (Piotrowski et al., 2004). These enzymes are involved in the biosynthesis of basic GLS structures

100   from elongated and non-elongated amino acids. The basic GLS structures are subjected to a range of

101   secondary side chain modification and transformation pathways catalyzed by enzymes such as flavin

102   monooxygenase (FMOOXs) (Hansen et al., 2007), GLS-AOPs (Mithen et al., 1995), GLS-OH

103   (Hansen et al., 2008) and CYP81Fs (Pfalz et al., 2009; 2011) to generate different types of GLS

104   structures, that are the last finalizing gene family involved in the indolic biosynthetic pathway

105   (Clarke, 2010; Fahey et al., 2001).

106   The most important mechanism for the wide production of secondary metabolites as glucosinolates

107   relies on whole-genome events, which occurred in Brassicaceae evolution history (Kliebenstein et al.,

108   2001a,b; Kroymann, 2011). The availability of the whole-genome sequences gives an opportunity for

109   using comparative genomics, which, in turn, can lead to a better understanding of the genome

110   evolution in this family. Whole-genome sequences are available for more than 100 plant species

111   (Tohge et al., 2014). The massive contribution, resulting from next-generation technologies, cannot

112 be currently matched by metabolomics, especially if high-quality and species-optimized approaches

113 are adopted (Fukushima et al., 2014). With the increasing number of whole-genome sequences and

114 the freely available genomic resources, the opportunities for conducting an analysis based on

115 comparative genomics is foreseen.

116 In this paper, we investigated gene clusters involved on the biosynthesis of GLS, by combining

117 genome analysis with biochemical pathways and compound structure assessment. Considering the

118 high diversity in GLS content in Brassicaceae species, we aim to: i) contribute to the global GLS

119 gene inventory in Brassicaceae; ii) compare gene diversity within the three GLS sub-pathways; iii)

120 assess a potential genetic basis for GLS divergence using 6 CYP genes (*CYP79F1-F2* and *CYP81F1-*

121 *F4*), known to be key genes of indolic and aliphatic GLS biosynthetic pathways, respectively; and iv)

122 increase the knowledge on the chemical diversity of GLS compounds in major *Brassica* crops

123 compared to the CWR of the genus *Diplotaxis*. By combining chemical data with genomic

124 sequences, we expect to provide information of interest for promoting the use of the neglected

125 *Diplotaxis* genus as a potential viable CWR of economically important *Brassica* crops.

126

## 2. Materials and Methods

128

### 2.1. GLS biosynthetic genes: compilation and gene ontology annotation

130 QuickGO (https://www.ebi.ac.uk/QuickGO/, Binns et al., 2009), AmiGO

131 (http://amigo.geneontology.org/amigo, Carbon et al., 2008) and MetaCyc (https://metacyc.org/, Caspi

132 et al., 2017) databases were used to filter genes involved in GLS biosynthetic process (GBP) by

133 searching the specific GO term (GO:0019761). Sequences representing the complete set of GLS

134 biosynthetic genes in *Arabidopsis thaliana* were acquired from The Arabidopsis Information

135 Resource (TAIR, www.arabidopsis.org, accessed on July 2019, Berardini et al., 2015), and further

136 complemented with a set of genes listed as GLS genes in the Brassica database (BRAD,

137 http://brassicadb.org, Wang et al., 2015), which is a web-based database of genetic data at the whole

138 genome scale for important *Brassica* crops. After, a complete assessment of GLS biosynthetic genes

139 in Brassicaceae species was retrieved, through searching of several public databases namely:

140 Arabidopsis Information Resource (TAIR), BrassicaDB, and nucleotide blast (Blastn) at NCBI,

141 restricting the search to orthologs within the Brassicaceae family. The genes sequences listed as GLS

142 genes in BRAD, were subjected to nucleotide Blast (Blastn on TAIR), to identify *Arabidopsis*

143 *thaliana* homologous genes with a threshold of E-value $\leq 10^{-10}$. The following step was to perform a

144   complete assessment of GLS biosynthetic genes in Brassicaceae species by using the BLASTN

145   algorithm in National Center for Biotechnology Information (NCBI) public database, restricting the

146   search for orthologs within the Brassicaceae family, with a threshold of E-value $\leq 10^{-10}$ and 50% of

147   query cover. Blast2GO v.5.2 (Götz et al., 2008) was used to assign GO terms to the sequences

148   dataset, to allow unigene annotation according to three main Gene Ontology categories, i.e. Cellular

149   Compartment, Molecular Function and Biological Process. A BlastX algorithm was used with the

150   following parameters: a constant expectation value threshold of $1.0E^{-10}$, 20 Blast Hits, HSP length

151   cutoff set at 33 and HSP Hit Coverage at 60. The different genomic information gathered from a

152   multi-databasing approach was represented by an Euler diagram using the online generator tool

153   available at https://www.meta-chart.com/. The resulting figure (Figure 2) was scaled, so that the area

154   of the shape was proportional to the number of genes it contained, and the overlapping shapes

155   represented the genes that were present in more than one database.

156

157   **2.2. Gene clustering analysis**

158   The collected GLS biosynthetic genes were used to perform a gene cluster analysis under two

159   different approaches: unsupervised Principal Component Analysis (PCA) and a UPGMA. The PCA

160   analysis was carried out using *factoextra* package in R version 3.6.1 through RStudio version

161   1.2.5001. To carry out the UPGMA analysis, a dataset containing the 78 GLS gene sequences

162   assigned to each of the sub-pathways was analyzed using MEGA X version 10.0.5 (Kumar et al.,

163   2018). A model assessment was performed to calculate the most adequate model to the dataset, and

164   subsequently, a UPGMA analysis was constructed using 10000 bootstraps. Phenograms were edited

165   using FigTree version 1.4.4 (Rambaut, 2009).

166

167   **2.3. Phylogenetic analysis of CYP79F and CYP81F genes**

168   Sequences from *Arabidopsis thaliana CYP79F1-F2* and *CYP81F1-F4* were retrieved from the TAIR

169   database. Brassicaceae orthologs were assessed by blasting genes from *Arabidopsis thaliana* against

170   the NCBI database using Blastn, with an E-value of $\leq 10^{-10}$ and 50% of query cover, restricted to the

171   Brassicaceae family. A total of 101 sequences were retrieved and analyzed, where only 69 were

172   marked as unique (i.e. not shared across genes). The final dataset comprised 25 sequences from

173   *CYP79F1-F2* [*CYP79F1*: n=8, *CYP79F2*: n=7 and shared: n=10] and 44 from *CYP81F1-F4*

174   [*CYP81F1*: n=9, *CYP81F2*: n= 9, *CYP81F3*: n= 10, *CYP81F4*: n= 9 and shared: n=7]. Sequences

175   were aligned using MAFFT version 7 auto strategy (Katoh et al., 2017) and then trimmed using

176   trimAl version 1.3 available at the Phylemon 2 suite (http://phylemon.bioinfo.cipf.es/) under the

177    automated1 algorithm. Model calculations were carried out using PartitionFinder2 (Lanfear, 2017)

178    and then phylogenetic estimations were made using RAxML version 8.2.10 through raxmlGUI

179    version 1.5b2 using a ML+ rapid bootstrap, autoMRE, using *Arabidopsis thaliana* CYP79 genes as

180    outgroups. Lastly, visualization and manipulation of the trees was done using FigTree version 1.4.4

181    (Rambaut, 2009).

182

183    **2.4. GLS compounds assessment**

184    Major agricultural brassica crops (i.e. *Brassica* sp., *Eruca vesicaria*) were selected and compiled for

185    GLS compounds analysis through an exhaustive literature review (*Brassica rapa* – Cartea et al.,

186    2012; *Brassica napus* – Velasco et al., 2008; *Brassica olearaceae* – Bhandari et al., 2015; *Eruca*

187    *vesicaria*- D'Antuono et al., 2008). *Diplotaxis* species were also included in the GLS chemodiversity

188    analysis as being probable precursors and wild relatives of *Brassica* crops (D'Antuono et al., 2008).

189    A matrix of presence/absence was built and then projected as a heat map using the Heatmap tool

190    freely available (http://www.hiv.lanl.gov/) using the Euclidean distance method with an average

191    linkage clustering, and 10000 bootstraps.

192

193    # 3. Results

194

195    **3.1. Genomic information on GLS genes**

196    The species diversity assessment carried out on GLS genes available at public databases enabled the

197    identification of 101 *Arabidopsis* genes that were blasted using Blastn (NCBI) restricted to

198    *Brassicaceae*. From the results obtained, 36 species contain information on orthologous genes

199    belonging to the GLS metabolic pathway. As expected, the most represented species was *Arabidopsis*

200    *thaliana*, which accounted for 32% of the total GLS available genes. Other species, in particular the

201    major crop species *Brassica napus, Brassica oleracea and Brassica rapa*, display 37% of the

202    genomic information available at public databases. *Raphanus sativus* (radish) comprised 8% of the

203    data, with other Brassicaceae model species, namely *Camelina sativa*, *Capsella rubella*, *Arabis*

204    *alpina* and *Eutrema salsugineum* complementing the remaining genomic information available on

205    GLS genes (Figure 1).
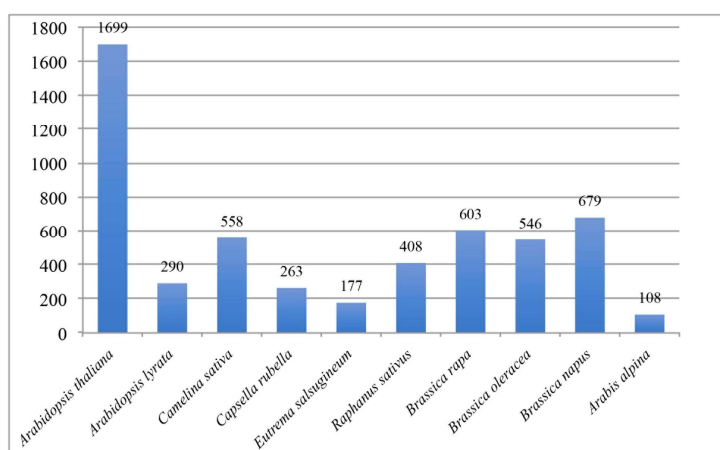
206
207
208
209
210
211

**Figure 1**- Genomic information of the GLS biosynthetic genes of the top ten Brassicaceae species registered at NCBI database. Number of sequences available per species is represented above each bar.

### 3.2. Global assessment of GLS biosynthetic genes identification

A global overview of the GLS biosynthetic pathway in the Brassicaceae family was developed using

a multi-database approach. Although several studies have already been performed to achieve a

similar pathway reconstruction analysis, we provide in our paper a global assessment of the GLS

biosynthetic pathway using not only genes described for *Arabidopsis* but also for *Brassica* species.

To do so, we retrieved all the genes belonging to the GLS biosynthetic pathway using its specific GO

term (GO:0019761) (Supplementary Table 1). From this thorough inventory, a total of 101 genes

were identified in *Arabidopsis thaliana* as being GLS biosynthetic genes: 52 from AmiGO, described

as being involved in the GLS biosynthesis (GO:0019761); 52 from Brassica database

(Brassicadb.org) classified as GLS genes and 67 from MetaCyc that were present in the GLS
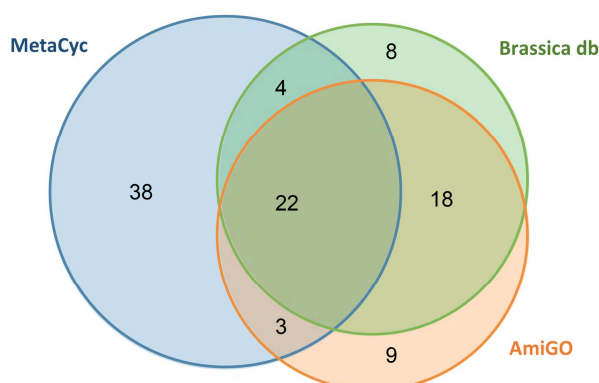
synthesis reaction cascade (Figure 2).



**Figure 2**- Euler diagram displaying GLS gene annotation gathered from a multi-database approach.

7

251  Using MetaCyc, it was possible to assign the genes according to each of the GLS sub pathways:

252  aromatic, indolic and aliphatic. From the total 101 genes, 78 were assigned to each of the sub-

253  pathways, while the remaining 23 were not reported as biosynthetic specific genes, and thus their

254  assignment remains unclear. This is probably due to putative functions related to substrate diversity

255  and regulation of GLS synthesis. Using this database, it was possible to identify 31 genes specific

256  from aliphatic GLS, 26 genes specific from indolic GLS synthesis, and 6 genes specific from

257  aromatic GLS biosynthesis (Table 1, for specific genes see Supplementary Table 2).

258

259  **Table 1-** GLS genes information according to sub-pathways of indolic, aliphatic and aromatic.
260  Number of genes - total of genes annotated in each sub-pathway; Number of specific genes - genes
261  exclusive to a given sub-pathway; Number of shared genes - genes shared in at least two sub-
262  pathways.

263

|  | Aliphatic | Aromatic | Indolic | Combined unigenes of the 3 pathways |
|---|---|---|---|---|
| Nº. Genes | 40 | 20 | 41 | 78 |
| Nº. Specific Genes | 31 | 6 | 26 | - |
| Nº. Shared Genes | 9 | 14 | 14 | 8 |

264

265  Gene Ontology (GO) assignment revealed a high diversity regarding the three multi-level

266  categorizations: Biological Process (BP), Molecular Function (MF) and Cellular Component (CC)

267  (Figure 3). On the first categorization level, the top-hits of biological process (with more than 45

268  sequences) were related to metabolic and cellular processes, followed by response to stimulus. In the

269  molecular function level, binding and transcription regulation activities were the most represented

270  after the catalytic activity; while in the cellular component level, genes were mainly grouped by

271  membrane and/or organelle. These GO terms tie in with GLS biosynthetic functions, like the

272  transcription regulation activities attributed to the MYB gene family, known to act as transcription

273  factors/regulators of GLS unique to the GLS-synthesizing Brassicales (*MYB34*, *MYB51* and MYB122

274  in indolic pathway; *MYB28*, *MYB29* and *MYB76* in aliphatic GLS). These hints at possible unknown

275  GLS functions need to be further explored to fully assign and determine the complete gene functions

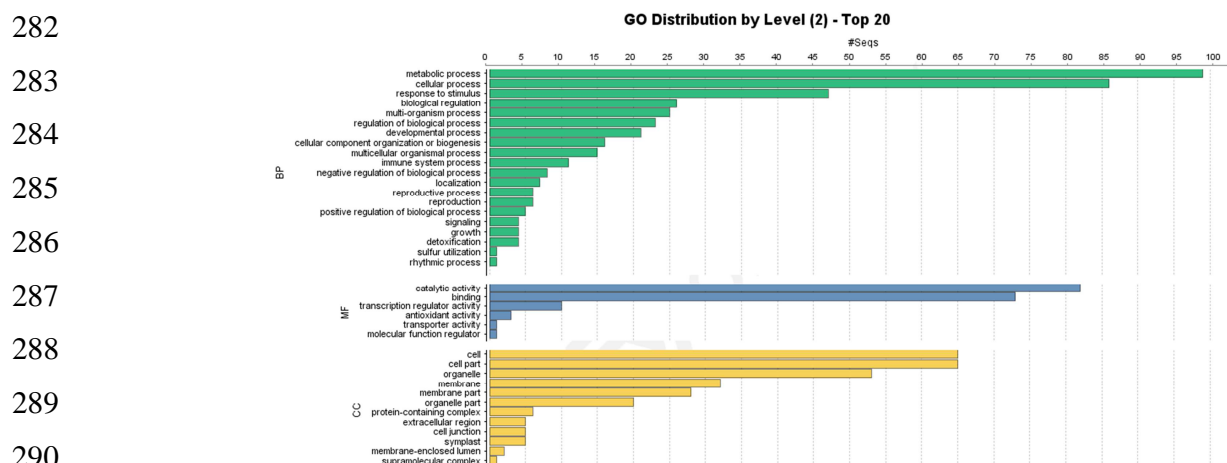276  in the Brassicaceae GLS biosynthetic pathway.

277

278

279

280

281

282

283

284

285

286

287

288

289

290



**Figure 3**- Gene Ontology (GO) terms assignment for the GLS biosynthetic genes. The graph displays the term enrichment levels of the annotated sequences along with the GO term hierarchy: Biological Process (BP, in green), Molecular Function (MF, in blue) and Cellular Component (CC, in yellow).
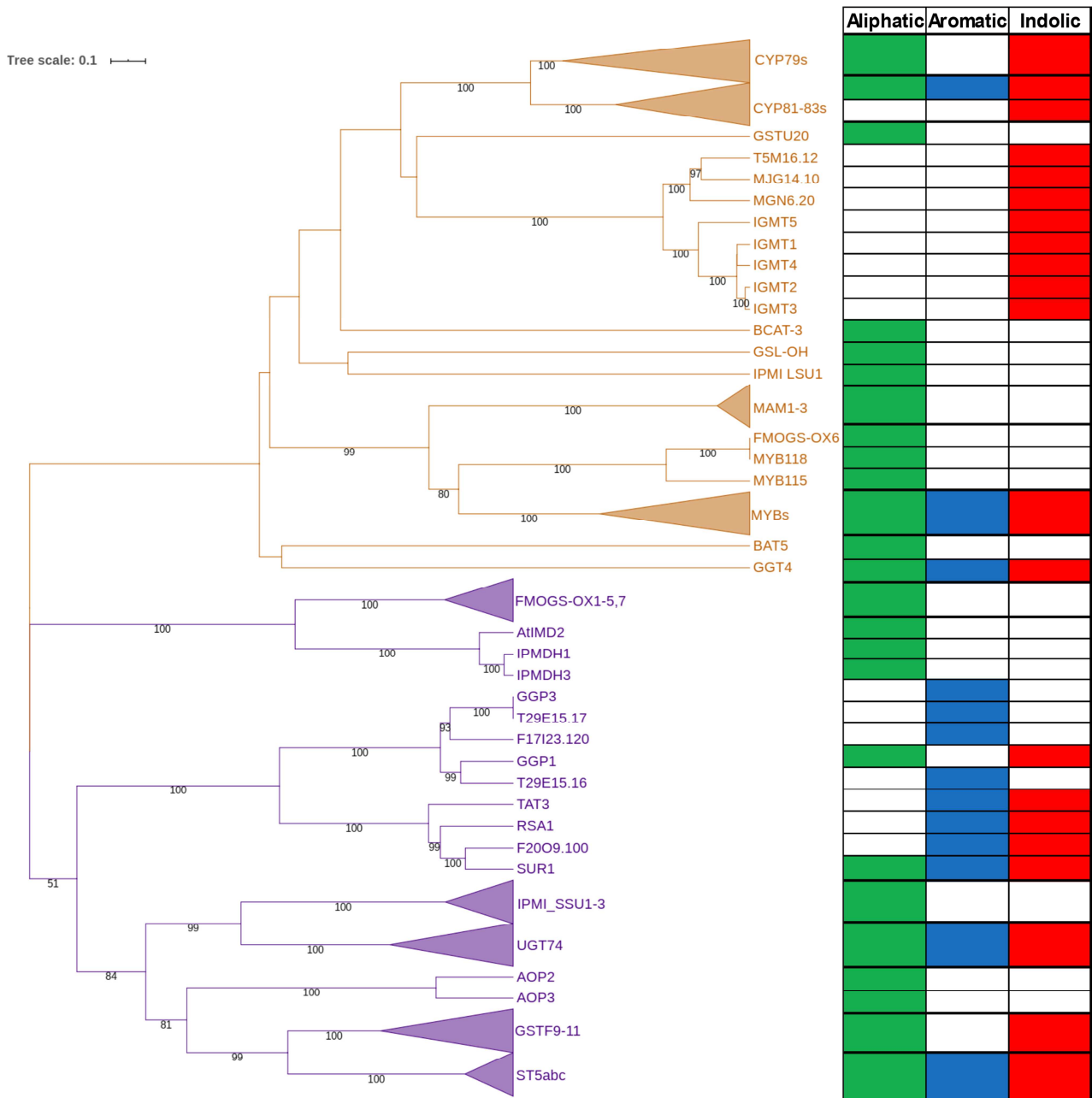
**3.3. Uncovering GLS sub-pathways gene specificity**

In order to detect sub-pathway specificity of GLS genes, two different methodologies have been

applied: 1) UPGMA as a bottom-up hierarchical clustering method to evaluate gene clustering based

on sequence alignment, regardless evolutionary features; 2) Principal Component Analysis (PCA) to

assess gene grouping discrimination associated to each sub-pathway.

An UPGMA phenogram (Figure 4) allowed the discrimination of two main clusters: one highlighted

in purple, which include many genes shared by the 3 subtypes pathways (aliphatic / indolic /

aromatic) and mostly related the synthesis of GLS core structure; and a second cluster in orange

with many genes belonging to the *CYP450*s and *MYB*s gene families, which are essentially genes

related to side chain elongation of GLS, regardless of being indolic or aliphatic, and which are known

to be responsible for the great diversity of existing compounds.

PCA analysis of the GLS genes (Supplementary Figure 2) showed a shared membership with no

discrimination between the three sub-pathways (indolic, aliphatic and aromatic). These results are

corroborated by the UPGMA phenogram where no pathway-specific clustering was identified. The

analysis of PCAs loading plots (Supplementary Figure 2), PCA1 reveals 26% of the total variation,

while PCA2 accounts for 13.5%. The PCA1 variation appears to be connected with a group

composed of CYP79 and CYP81 genes. Interestingly, these genes belong to different GLS pathways:

*CYP81F1-F4* is indolic-specific while *CYP79F1-F2* is exclusive to the aliphatic pathway. Only

*CYP83A1* and B1, and *CYP79A2* are shared within the three pathways.

315

316



**Figure 4**- UPGMA phenogram of the 101 GLS biosynthetic genes. Bootstrap values above 50 are represented on the branches. Detailed UPGMA tree is available at Supplementary Figure 1.

319

320

**3.4. Testing gene divergence as a baseline to GLS diversification**

Considering the two approaches by UPGMA and PCA, which disclosed a potential clustering of

*CYP79* and *CYP81* genes (*CYP79F1-F2* and *CYP81F1-F4*), a more detailed phylogenetic analysis

324   was conducted to understand the disjunction of the aliphatic and indolic genes. Considering that side

325   chain modifications of indolic GLS are controlled by four CYP81F enzymes (*CYP81F1- F4*) (Barco

326   and Clay, 2019), while *CYP79F1* and *CYP79F2* are involved in the biosynthesis of aliphatic GLS in

327   *Arabidopsis thaliana*, by exploring sequence diversification on Brassicaceae orthologs, a possible

328   differentiation of those genes could be uncovered on *Brassica* crops, which display a higher diversity

329   of aliphatic compounds than *Arabidopsis*. Overall, the obtained ML phylogeny is well supported,

330   resulting in two main clades, which separate the genes associated with the aliphatic pathway

331   (*CYP79F1* and *CYP79F2*) and those associated to the indolic pathway (*CYP81F1*, *CYP81F2*,

332   *CYP81F3* and *CYP81F4*) (Figure 5). The phylogenetic analysis clearly splits the two types of CYPs

333   analyzed, *CYP79* (in red) and *CYP81* (in blue), into two well-supported clusters. As such, it appears

334   that the gene divergence between these CYPs underlies the basis of indolic and aliphatic GLS

335   biosynthesis.



**Figure 5-** Phylogenetic tree from the Maximum Likelihood analysis of *CYP79F1-F2* and *CYP81F1-F4* genes in Brassicaceae with *A. thaliana CYP79F* genes as outgroups. Acronyms are present as the first letter of the genus and the second to species, e.g. At for *Arabidopsis thaliana*, and gene identification when possible. Upon lack of complete CYP annotation, accession numbers were used. Different copies of the same gene are identified by an "X" following sequential numbering, e.g. *A. thaliana* X1, *A. thaliana* X2. Only bootstrap values above 50 are presented. Accession numbers of the sequences analyzed are provided in Supplementary Table 3.

11

359

360

361 From the analysis of the *CYP79* genes, a cluster including all major *Brassica* crops (*Brassica rapa,*

362 *Brassica juncea, Brassica olearacea*) is evident (Figure 5, highlighted as *Brassica* cluster), where no

363 disjunction is observed from being *CYP79F1* or *CYP79F2*. It can be easily recognize that *Brassica*

364 crops are usually grouped in the same cluster, which reveals a common diversification of indolic

365 GLS that portrays *Brassica* chemotypes. Two apparent copies of *Brassica napa CYP79F2* are

366 grouped while other *Brassica* sp. sequences were assembled in different tree branches disclosing a

367 wide divergence on the *CYP79F1* and *F2* gene sequences which could be associated with the

368 diversity of aliphatic GLS in *Brassica* s.l. Regarding *CYP81F1-F4*, four clusters were obtained

369 matching essentially each of the *CYP81F* genes covered (Figure 5).

370

371 **3.5. Snapshot on GLS chemodiversity: *Brassica* crops and *Diplotaxis***

372 By performing a snapshot of the GLS chemodiversity using an average linkage clustering method

373 (Figure 6), a cluster including the *Brassica* crops (e.g. *Brassica olearacea*, *Brassica juncea*, *Brassica*

374 *rapa*, and *Brassica napus*) can be depicted, while *Diplotaxis* species appear to have a more complex

375 and diversified GLS chemical profile. Phylogenetic relationships indicate that *Diplotaxis* maintains

376 most of the primitive morphological characters while *Brassica* presents the most evolved ones with

377 *Erucastrum* occupying an intermediate position (Gómez-Campo and Tortosa, 1974; Gómez-Campo,

378 1980; Sánchez-Yélamo, 2009). By comparing the GLS chemotype diversification between *Brassica*

379 and *Diplotaxis* species, the latter shows a distinct GLS profile. In what concerns rocket crops,

380 collectively attributed to *Diplotaxis* and *Eruca,* the wild (*Diplotaxis tenuifolia*) and cultivated (*Eruca*

381 *sativa* and *Eruca vesicaria*) rockets are clustered together sharing a common GLS profile.

382 The results obtained revealed that *Brassica* and *Diplotaxis* have distinct GLS chemo-profiles. Within

383 *Brassica* species, a shared GLS profile is displayed, namely in what concerns aliphatic GLS such as

384 progoitrin, gluconapin, glucobrassicanapin that are specific to *Brassica* chemo-lineage. In *Diplotaxis*

385 and *E. vesicaria,* glucolepidin appears as the main distinctive GLS, followed by glucoerucin.

386 Moreover, such GLS are more diverse among *Diplotaxis* species than in *Brassica* species, possibly as

387 the result of crop selection events that have narrowed *Brassica* chemodiversity when compared to

388 *Diplotaxis* species, in which few domestication events occurred and several species are in the wild

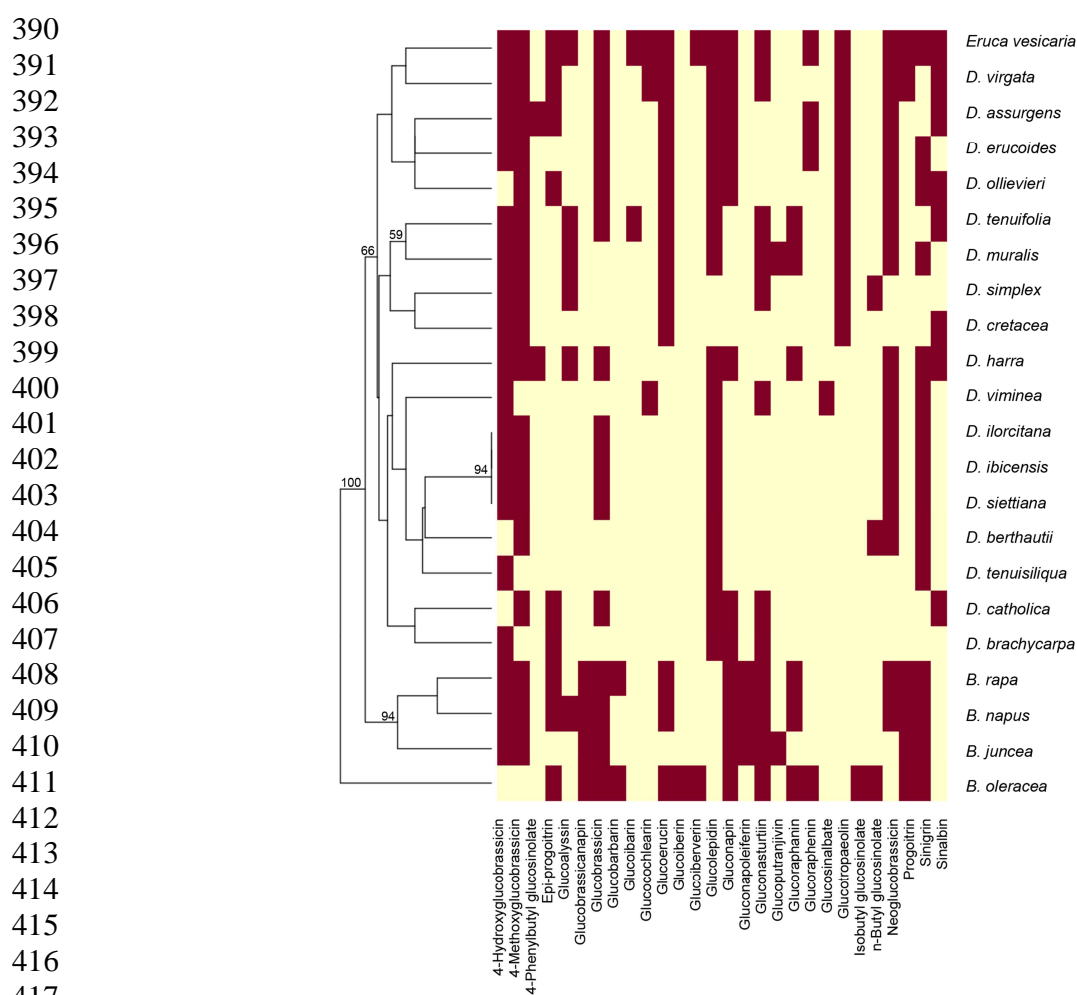389 exposed to habitat conditions and constraints.

**Figure 6**– Chemodiversity profiling of GLS in *Brassica* and rocket species (*Diplotaxis* and *Eruca*). Data matrix of GLS chemodiversity is provided in Supplementary Table 4. Colors indicate presence (red) and absence (light yellow) of a glucosinolate compound. Bootstraps values above 50 are presented in the clustering phenogram resulted from the Eucledian distances method.

## 4. Discussion

### 4.1. GLS biosynthetic pathway: gene signature of aliphatic and indolic vias

In this study, we have performed a comprehensive assessment of the GLS biosynthetic pathways in Brassicaceae family. A reconstruction analysis of GLS pathway and a global assessment using genes described for *Arabidopsis* and for *Brassica* species were established, using a multi-database approach (i.e. TAIR, NCBI, Brassicadb, MetaCyc). From a total of 101 genes identified, about 78 previously identified genes in *Arabidopsis* were classified into the three sub-pathways of the GLS biosynthetic pathway (Supplementary Table 1), while the remaining 23 were not possible to assign to any of the three-specific pathway of GLS, and need further study to uncover their functional role in specific pathways. With the upcoming availability of more and more genomic resources from Brassicaceae

13

434    species, a complete set of functional disclosure within each GLS is foreseen in a near future. Our

435    clustering approach of the 101 GLS genes recorded, showed a clear separation between the genes

436    involved mainly in GLS core structure synthesis shared between the 3 sub-pathways (aliphatic /

437    indolic / aromatic), and other genes belonging to the CYP450s and MYBs gene families (Figure 4).

438    The last cluster was therefore essentially composed of genes related to the side chain elongation

439    process of GLS, which is responsible for the great diversity of GLS compounds, namely the

440    biosynthesis genes (*CYP450*s) and the regulators of transcription (*MYB*s). In the PCA analysis,

441    *CYP450*s genes are the highest contributors for explaining the variation within GLS biosynthetic

442    pathway, namely: *CYP81 F1-F4*, which is indolic-specific, and *CYP79F1-F2*, unique to the aliphatic

443    pathway. These overall results allowed us to conclude that the GLS biosynthetic pathway depends on

444    upstream genes essentially involved in the core structure synthesis, while genes involved in the

445    synthesis of aliphatic and indolic GLS, apparently specific of Brassicaceae, may depend on two set of

446    gene clusters, known to be important for aliphatic - and indolic-specific pathways (*CYP79F1-F2* and

447    *CYP81F1-F4*, respectively). Likewise, aliphatic and indolic GLS are the two most important types of

448    GLS present in Brassicaceae.

449

450    **4.2. Gene diversification in GLS: the case of *CYP79F1-F2* and *CYP81F1-F4***

451    Aliphatic and indolic GLS are derived from aliphatic (methionine, alanine, valine, leucine, and

452    isoleucine) and indolic amino acids (tryptophan), respectively (Wittstock and Halkier, 2002). In

453    *Arabidopsis thaliana* and *Arabidopsis lyrata*, aliphatic GLS are formed exclusively from methionine

454    (Windsor et al., 2005). Species of the Brassicaceae have been useful models to understand the

455    dynamics and impacts of ancient polyploidy (genome doubling), with the entire family having

456    undergone a whole genome duplication (named At-α) and the *Brassica* crops suffered an additional

457    genome triplication (Br-α) (Schranz and Mitchell-Olds, 2006; Thomas et al., 2006). Several authors

458    have suggested that the genetic diversification of GLS in Brassicaceae is correlated with the

459    polyploid occurrence in this family, with the At-β WGD event at 77.5 Mya where indolic

460    glucosinolates appeared, and the At-α event at approximately 56 Mya (Kagale et al., 2014), where

461    chain elongation of Met-derived aliphatic GLS is present (Schranz et al., 2011). Moreover, it has

462    been pointed out that the diversity based on GLS composition in *Brassica* species could be related to

463    A, B and C genomes (Ishida et al., 2014). The three ancestral *Brassica* species with diploid genome

464    chromosomes: *Brassica nigra* (BB, 2n = 16) contain GLS with three carbon (C) side chains, derived

465    from a single elongation reaction: *Brassica oleracea* (CC, 2n = 18) contains GLS with either 3C or

466    4C side chains; and *Brassica rapa* (AA, 2n = 20) contains GLS with either 4C or 5C side chains

467   (Ishida et al., 2014). Recently, studies have focused on genome evolution underlying the basis of

468   GLS diversification. Bergh et al. (2016) reported that the genes that have undergone a high

469   diversification process encode the MAM (Methylthioalkylmalate) enzymes and also the CYP81 side-

470   chain modification enzymes responsible for a large part of the GLS chemotypes observed. MAM

471   synthase enzymes are central for the diversification of aliphatic GLS structures in *Arabidopsis*

472   *thaliana* and related species (Heidel et al, 2006); while *CYP81F* acts in the final step of the indolic

473   GLS pathway and have been reported as responsible for a wide array of natural variation among

474   *Arabidopsis thaliana* ecotypes (Pfalz et al., 2009). In this species, the biosynthesis of indolic GLS,

475   hydroxylations are catalyzed by cytochromes P450 of the *CYP81F* subfamily (Pfalz et al., 2009;

476   Barco and Clay, 2019), followed by methylation of the methyltransferases, IGMT1 and IGMT2 in

477   *Arabidopsis thaliana* (Pfalz et al., 2011). In indolic GLS, the CYP81Fs family (CYP81F1-F3) has

478   been identified as the encoder of the oxidizing enzyme that converts indolyl- 3-methyl GLS (I3M) to

479   4OH-I3M, while CYP81F4 acts in the hydroxylation at C1-position (Pfalz et al., 2011). Such

480   secondary modifications can present high variability among species in nature and they are the main

481   responsible for the diversity observed across more than 120 types of GLS that have been described to

482   date (Kliebenstein et al., 2001a). CYP81F2 has been suggested to have neofunctionalized in plant

483   innate immunity that subsequently was maintained in *Arabidopsis thaliana*, but lost in the ancestral

484   Brassicaceae species. The phylogenetic analysis performed revealed four clusters, each of them

485   associated to *CYP81F1* to *F4*. Since *Brassica* crops were grouped in the same cluster, it suggests a

486   common diversification of indolic GLS that portrays *Brassica* chemotypes, which are present in less

487   extent in aliphatic GLS.

488   In the aliphatic GLS pathway in Brassicaceae, *CYP79* is a key variable gene that has been considered

489   as a driving force in GLS diversification. Several steps catalyzed by *CYP79F1* and *CYP79F2* result

490   from gene duplication (Olson-Manning et al., 2013). CYP79F1 and CYP79F2 present slightly

491   different substrate specificities: CYP79F1 uses both short- and long-chain substrates, whereas

492   CYP79F2 tends to use only long-chain substrates. It has been considered that in *Brassica rapa*, like

493   in *Arabidopsis thaliana* all the gene counterparts participate in the formation of the GLS core

494   structure, except for *CYP79F2* (Wang et al., 2011). The absence of *CYP79F2* agrees well with the

495   fact that all profiles of aliphatic GLS in *Brassica rapa* are composed of short-chain GLS. From the

496   phylogenetic analysis we performed, it can be concluded that *Arabidopsis thaliana CYP79F1* and *F2*

497   genes are in the basis of the diversification of the remaining Brassicaceae species (Figure 6).

498   *Brassica* crops are grouped in a single cluster (highlighted in shaded red in Figure 6), which

499   represents a common genetic basis of the *CYP79F1-F2* responsible for the GLS diversification and

500    possibly links to the additional genome triplication (Br-α) event that these crops suffered throughout

501    their evolution. The annotation of *CYP79F1* and *F2* genes in Brassicaceae is limited as only recently

502    genome sequences are being released, pushed by the continuous lower costs of whole-genome

503    sequencing technologies. With our study, we were able to determine the disjunction of *Arabidopsis*

504    *thaliana CYP79F1-F2* with the remaining Brassicaceae species and in particular with *Brassica* crops,

505    which were grouped together in a lineage associated with aliphatic GLS.

506    Our results revealed the most recent diversification of *CYP79F1* and *F2* genes in *Brassica* crops,

507    where a single cluster including *Brassica* species is difficult to depict (Figure 6). This lack of clear

508    clustering from *CYP79F1* and *CYP79F2*, in opposite with what is observed in *CYP81F1-F4*, may

509    suggest the absence of a *CYP79F2* gene as reported for *Brassica rapa* (Wang et al., 2011), which

510    may not be the case for other *Brassica* species. This may suggest a shared genetic basis underlying

511    short-chain aliphatic GLS, since in *Arabidopsis thaliana* a *CYP79F2* knockout mutant presents a

512    considerable reduction of long-chain aliphatic GLS (Chen et al., 2003). Moreover, future annotation

513    efforts of Brassicaceae genes has to be performed as a way to clarify *CYP79F1* diversification within

514    *Brassica* crops that should be linked to a higher production of short-chain aliphatic GLS.

515

### 4.3. Chemical diversity of GLS in Brassicaceae

517    GLS production by Brassicaceae plants is considered as being influenced by environmental factors

518    such as soil, climate and cultivation conditions including fertilization, harvest time, and plant organ

519    (Martínez-Ballesta et al., 2013). In general, the diversity of GLS profiles is higher in *Brassica*

520    *oleracea* as opposed to *Brassica rapa* (Figure 7). The Brassicaceae plant tissues include one or more

521    major GLS mostly composed of aliphatic GLS. In general, Brassicaceae vegetables GLS contain an

522    alkyl side chain with 3–5 carbons (Ishida et al., 2014). From these ones, glucoiberin is present mostly

523    in *Brassica oleracea* vegetables (cabbage, broccoli, and cauliflower) while, gluconapin and

524    progoitrin are ubiquitous in many *Brassica* vegetables such as *Brassica rapa* (Chinese cabbage,

525    mustard spinach, and turnip), *Brassica oleracea*, *Brassica juncea* (mustard green), and *Brassica*

526    *napus* (rapeseed vegetable) (Ishida et al, 2014). Glucoerucin is mainly found in cultivated *Eruca*

527    *sativa* and wild rockets (*Diplotaxis tenuifolia*, *Diplotaxis* sp.) rockets.

528    In general, *Diplotaxis* spp. emerges as an extremely GLS-rich species, revealing likely taxonomic

529    affinities with taxa previously examined by other criteria suggesting a high potential for further

530    exploitation. The disclosure of a distinct GLS chemo-profile between *Brassica* crops and *Diplotaxis*

531    species (i.e. in *Brassica*, progoitrin, gluconapin, glucobrassicanapin are the most abundant GLSs,

532    while in *Diplotaxis* glucolepidin and glucoerucin are the most distinctive), opens a new perspective

533  for addressing more studies towards not only the characterization of new taxa from the later genus

534  but also the quantification of such GLS, since many of them, in high amounts, are considered to be

535  anti-nutritional even in vegetables (e.g. Augustine et al., 2013). GLS production and contents in

536  Brassicaceae plants are influenced by environmental factors such as soil, climate and cultivation

537  conditions including fertilization, harvest time, and plant position, besides its straight relation to both

538  biotic and abiotic stresses (Martínez-Ballesta et al., 2013; Ishida et al., 2014). Despite several reports

539  on a positive relationship between GLS production and abiotic stress, it is still unknown which are

540  the mechanisms of resistance to drought and salinity conditions. Determining a chemodiversity

541  profile associated with phenotypes adapted to extreme environmental conditions, such as drought and

542  salinity, could be a good strategy for prospecting GLS compounds and contents and quantity for

543  coping with abiotic stresses.

544

545  **4.4. Abiotic stress and GLS crosstalk in Brassicaceae: wild rockets as emergent taxa**

546  Variation in the amount and profile of GLS compounds has been correlated with abiotic stresses

547  (Variyar et al., 2014). Among the most important, salinity and drought stresses are known to

548  significantly affect crops productivity. Overall, GLS content increases markedly under salinity,

549  drought, high temperature and nitrogen (N) deficiency (Martínez-Ballesta et al., 2015).

550  Extensive studies in Brassicaceae family showed a positive correlation between salt stress and GLS

551  content, [e.g. in broccoli (López-Berenguer et al., 2009), canola (Khalifa, 2012), radish sprouts (Yuan

552  et al., 2010), pakchoi (Keling and Zhujun, 2010)]. An increase in the signature of GLS content has

553  also been reported for Brassicaceae taxa under drought stress, namey in *Brassica napus*

554  (Champolivier and Merrien, 1997), *Brassica oleracea* (Radovich et al., 2005), *Brassica rapa* (Zhang

555  et al., 2008), *Brassica juncea* (Tong et al., 2014), and *Brassica carinata* (Ngwene et al., 2017).

556  However, recent studies in wild rocket (*D. tenuifolia*), demonstrated that salinity conditions did not

557  affect the total amount of GLS profile (Bonasia et al., 2017; Cocetta et al., 2018). Bonasia et al.,

558  (2017) showed that the aliphatic-GLSs proidrin, epiproidrin, and glucoerucin contents were

559  unaffected by salt stress (Bonasia et al., 2017), with glucoerucin emerging as a GLS compound

560  specific of *Diplotaxis*, of *Eruca vesicaria* and of *E. sativa* (Barillari et al., 2005). Furthermore,

561  glucoerucin could be linked to a distinctive chemical signature of the *Diplotaxis-Eruca* lineage

562  involved in salt tolerance, setting it apart from the *Brassica* crops chemo-lineage (Figure 6).

563  Under drought stress, indole glucosinolate biosynthetic genes revealed to be up-regulated in wild

564  rocket (Cavaiuolo et al., 2017), which accounts for a possible tolerance mechanism as described for

565  other brassicas under stress (Martínez-Ballesta et al., 2015). In this tolerance mechanism, *MYB* genes

566     (particularly MYB28 an MYB29) may play a role in variations of GLS contents. Salehin et al. (2019)

567     confirmed that MYB28 and MYB29 are important transcription factors regulating the synthesis of

568     indole GLS, where a *cyp79f1f2* double mutant revealed to be less tolerant to drought, probably due to

569     the loss of aliphatic GLS compounds, corroborating former studies (Martínez-Ballesta et al., 2015).

570     Moreover, Martínez-Ballesta et al. (2015) highlighted that pathways involved in the physiological

571     responses to salt stress are connected to GLS metabolism. Under salt stress, an increase in short-

572     chain aliphatic GLS was observed which has been further associated to a higher expression of

573     aquaporins, involved on osmoregulation pathways (Martínez-Ballesta et al., 2014), and thus could

574     contribute to water saving process (Martínez-Ballesta et al., 2015). Overall, short-chain aliphatic

575     GLS may contribute to water saving under salt stress (Martínez-Ballesta et al., 2015), while under

576     drought indolic GLS seems to be the most affected (Salehin et al., 2019).

577     When compared to Brassica crops, wild rockets seem to display a different GLS profile that could be

578     associated to an abiotic stresses tolerance. Indeed, the neglected and underutilized rocket species, i.e.

579     *Eruca sativa* (rocket), *Diplotaxis tenuifolia* and *Diplotaxis muralis* (wild rocket), as well as other

580     wild taxa distributed and adapted to extreme ecological conditions (i.e. severe salinity and drought

581     conditions), may be considered as potential targets to understand abiotic stress tolerance mechanisms.

582     *Diplotaxis* is considered an unexplored Brassicaceae crop wild relative (CWR), with *Brassica* crops

583     having evolved from the *Diplotaxis–Erucastrum* complex (Arias and Pires, 2012), which makes

584     *Diplotaxis* species an important reservoir of genetic resources for crop improvement.

585

## 5. Conclusions

587     Overall, we have <u>analysed</u> gene clusters involved in the biosynthesis of GLS, by combining genome

588     analysis with biochemical pathways and chemical diversity assessment. An integrated approach was

589     performed by assessing a global GLS gene inventory in Brassicaceae and its diversity, analysing a

590     potential genetic basis for GLS divergence using 6 CYP genes (CYP79F1-F2 and CYP81F1-F4),

591     known to be key genes of indolic and aliphatic GLS biosynthetic pathways, linked to a chemical

592     diversity evaluation of GLS compounds in major *Brassica* crops compared to the wild relative genus

593     *Diplotaxis*. Our results point to a recent diversification of the aliphatic CYP79F1 and CYP79F2

594     genes in *Brassica* crops, while for indolic genes a clear separation is observed for CYP81F1-F4

595     genes, revealing an earliest divergence on this GLS sub-pathway. Chemical diversity assessment

596     allowed recognizing that *Brassica* and *Diplotaxis* have distinct GLS chemo-profiles, highlighting that

597     the latter genus includes extremely GLS-rich species. Considering the enormous potential of

598     biodiversity for finding new traits useful in breeding programs, screening of GLS-enriched

599    Brassicaceae species is of particular interest. Despite that GLS profiles may vary among species and

600    according to plant development and/or environmental factors, a highly diverse and unexplored

601    chemodiversity has been recognized within *Diplotaxis*. The discovery of the genomic information

602    behind such GLS diversity could constitute a potential for discovering new phytochemical and

603    nutraceutical sources potentially transferable to *Brassica* crops. Also, understanding the relationship

604    between Brassicaceae GLS genes and abiotic stress tolerance will be useful to contribute as source of

605    genes for improving new Brassicaceae vegetable varieties to cope with effects of global climate

606    changes.

607
608    **Conflict of Interest**
609    The authors declare that the research was conducted in the absence of any commercial or financial
610    relationships that could be construed as a potential conflict of interest.

611
612    **Authors Contributions**

613    Conceptualization, F.M. and M.M.R.; methodology, F.M., A.P.E., A.R.P., M.M.R.; Bioinformatic
614    analysis, A.P.E., F.M and A.R.P.; Results analysis, A.P.E., F.M., A.R.P., M.M., M.M.R.; writing—
615    original draft preparation, A.P.E., F.M., A.R.P., M.M.R.; writing—review and editing, A.P.E., F.M.,
616    A.R.P., M.S.P., M.M.R. and M.M. All authors have approved the submitted version of this
617    manuscript.

627

628    **References**
629    Armah, C. N., Derdemezis, C. Traka, M. H., Dainty, J. R., Doleman, J. F., Saha, S., et al. 2015. Diet
630    rich in high glucoraphanin broccoli reduces plasma LDL cholesterol: evidence from randomised
631    controlled trials. *Molecular Nutrition & Food Research* 59, 918–926.
632    https://doi.org/10.1002/mnfr.201400863.

633    Augustine, R., Mukhopadhyay, A., Bisht, N. C. 2013. Targeted silencing myb28transcription factor
634    gene directs development of low glucosinolate lines in oilseed (*Brassica juncea*). *Plant
635    Biotechnology Journal* 11, 855–866. https://doi.org/10.1111/pbi.12078.

636    Bak, S., and Feyereisen, R. 2001. The involvement of two p450 enzymes, CYP83B1 and CYP83A1,
637    in auxin homeostasis and glucosinolate biosynthesis. *Plant Physiology* 127, 108–118.
638    https://doi.org/10.1104/pp.127.1.108.

639    Barco, B. and Clay, N.C. 2019. Evolution of glucosinolate diversity via whole-genome duplications,
640    gene rearrangements, and substrate promiscuity. *Annual Review of Plant Biology* 70 (1): 585–604.
641    https://doi.org/10.1146/annurev-arplant-050718-100152.

642    Barillari, J., Canistro, D., Paolini, M., Ferroni, F., Pedulli, G. F., and Iori, R. (2005).Direct
643    antioxidant activity of purified glucoerucin, the dietary secondary metabolite contained in rocket
644    (*Eruca sativa* Mill.) seeds and sprouts. *Food Chemistry* 53, 2475–2482. doi: 10.1021/jf047945a

645    Bellostas, N., Sorensen, J., Sorensen, H. 2004. Qualitative and quantitative evaluation of
646    glucosinolates in cruciferous plants during their life cycles. *Agroindustria* 3, 267–272.

647    Berardini, T. Z., Reiser, L., Li, D., Mezheritsky, Y., Muller, R., Strait, E., Huala, E. 2015. The
648    *Arabidopsis* information resource: making and mining the 'gold standard' annotated reference plant
649    genome. *Genesis* 53, 474–485. https://doi.org/10.1002/dvg.22877.

650    Bergh, E. van den, Hofberger, J. A., Schranz, M. E. 2016. Flower power and the mustard bomb:
651    comparative analysis of gene and genome duplications in glucosinolate biosynthetic pathway
652    evolution in Cleomaceae and Brassicaceae. *American Journal of Botany* 103, 1212–1222.
653    https://doi.org/10.3732/ajb.1500445.

654    Bhandari, S., Jo, J., Lee, J. 2015. Comparison of glucosinolate profiles in different tissues of nine
655    Brassica crops. *Molecules* 20, 15827–15841. https://doi.org/10.3390/molecules200915827.

656    Binns, D., Dimmer, E., Huntley, R., Barrell, D., O'Donovan, C., Apweiler, R. 2009. QuickGO: a
657    web-based tool for gene ontology searching. *Bioinformatics* 25, 3045–3'46.
658    https://doi.org/10.1093/bioinformatics/btp536.

659    Bischoff, K. L. 2016. Glucosinolates. *Nutraceuticals*, 551–554. https://doi.org/10.1016/b978-0-12-
660    802147-7.00040-1.

661    Bonasia, A., Lazzizera, C., Elia, A., Conversa, G.  2017. Nutritional, biophysical and physiological
662    characteristics of wild rocket genotypes as affected by soilless cultivation system, salinity level of
663    nutrient solution and growing period. *Frontiers in Plant Science* 8:300. doi:10.3389/fpls.2017.00300.

664    Brown, P. D, Tokuhisa, J. G., Reichelt, M., Gershenzon, J. 2003. Variation of glucosinolate
665    accumulation among different organs and developmental stages of *Arabidopsis thaliana*.
666    *Phytochemistry* 62, 471–481. https://doi.org/10.1016/s0031-9422(02)00549-6.

667    Carbon, S., Ireland, A., Mungall, C. J., Shu, S. Q., Marshall, B., Lewis, S. 2008. AmiGO: online
668    access to ontology and annotation data. *Bioinformatics* 25, 288–289.
669    https://doi.org/10.1093/bioinformatics/btn615.

670    Cartea, M. E., and Velasco, P. 2007. Glucosinolates in Brassica foods: bioavailability in food and
671    significance for human health. *Phytochemistry Reviews* 7, 213–229. https://doi.org/10.1007/s11101-
672    007-9072-2.

673 Cartea, M. E., de Haro, A., Obregón, S., Soengas, P., Velasco, P. 2012. Glucosinolate variation in
674 leaves of *Brassica rapa* crops. *Plant Foods for Human Nutrition* 67, 283–288.
675 https://doi.org/10.1007/s11130-012-0300-6.

676 Caspi, R., Billington, R., Fulcher, C., Keseler, I. M., Kothari, A., Krummenacker, M., Latendresse,
677 M. et al. 2017. The MetaCyc database of metabolic pathways and enzymes. *Nucleic Acids*
678 *Research* 46 (D1): D633–39. https://doi.org/10.1093/nar/gkx935.

679 Cavaiuolo, M., and Ferrante A. 2014. Nitrates and glucosinolates as strong determinants of the
680 nutritional quality in rocket leafy salads. *Nutrients* 6, 1519–1538. https://doi.org/10.3390/nu6041519.

681 Cavaiuolo, M., Cocetta, G., Spadafora, N. D., Müller, C. T., Rogers, H. J., Ferrante, A. 2017.Gene
682 expression analysis of rocket salad under preharvest and postharvest stresses: a transcriptomic
683 resource for *Diplotaxis tenuifolia*. *PLoS ONE* 12(5): e0178119.
684 https://doi.org/10.1371/journal.pone.0178119

685 Champolivier, L., and Merrien, A. 1997. Corrigendum to 'Effects of water stress applied at different
686 growth stages to *Brassica napus L.* var. *oleifera* on yield, yield components and seed quality'.
687 *European Journal of Agronomy* 6 (3–4): 309–10. https://doi.org/10.1016/s1161-0301(97)00003-8.

688 Charron, C. S., and Carl E. S. 2004. Glucosinolate content and myrosinase activity in rapid-cycling
689 *Brassica oleracea* grown in a controlled environment. *Journal of the American Society for*
690 *Horticultural Science* 129 (3): 321–30. https://doi.org/10.21273/jashs.129.3.0321.

691 Chen, S., Erich, G., Kirsten, J., Peter, N., Bodil, J., Olsen, C. et al. 2003. CYP79F1 and CYP79F2
692 have distinct functions in the biosynthesis of aliphatic glucosinolates in *Arabidopsis*. *The Plant*
693 *Journal* 33(5): 923–937. https://doi.org/10.1046/j.1365-313x.2003.01679.x.

694 Choi, S., Park, S., Lim Y., Kim, S.,Park, J., and An, G. 2014. Metabolite profiles of glucosinolates in
695 cabbage varieties (*Brassica oleracea* var. *capitata*) by season, color, and tissue
696 position. *Horticulture, Environment, and Biotechnology* 55; 237–247.
697 https://doi.org/10.1007/s13580-014-0009-6.

698 Clarke, D. 2010. Glucosinolates, structures and analysis in food. *Analytical Methods* 2; 310.
699 https://doi.org/10.1039/b9ay00280d.

700 Cleemput, S., and Becker, H. 2011. Genetic variation in leaf and stem glucosinolates in resynthesized
701 lines of winter rapeseed (*Brassica napus L.*). *Genetic Resources and Crop Evolution* 59; 539–46.
702 https://doi.org/10.1007/s10722-011-9701-x.

703 Cocetta, G., Mishra, S., Raffaelli, A., Ferrante, A. 2018. Effect of heat root stress and high salinity on
704 glucosinolates metabolism in wild rocket. *Journal of Plant Physiology* 231: 261-270,
705 https://doi.org/10.1016/j.jplph.2018.10.003.

706 Collett, M.,Stegelmeier, B., and Tapper, B. 2014. Could nitrile derivatives of turnip (*Brassica rapa*)
707 glucosinolates be hepato- or cholangiotoxic in cattle? *Journal of Agricultural and Food*
708 *Chemistry* 62; 7370–7375. https://doi.org/10.1021/jf500526u.

709 D'Antuono, L. F., Elementi, S., and Neri, R. 2008. Glucosinolates in *Diplotaxis* and *Eruca* leaves:
710 diversity, taxonomic relations and applied aspects. *Phytochemistry* 69 (1): 187–199.
711 https://doi.org/10.1016/j.phytochem.2007.06.019.

712 Fahey, J. W., Zhang, Y., and Talalay, P. 1997. Broccoli sprouts: an exceptionally rich source of
713 inducers of enzymes that protect against chemical carcinogens. *Proceedings of the National Academy*
714 *of Sciences* 94 (19): 10367–72. https://doi.org/10.1073/pnas.94.19.10367.

715 Fahey, J.W., Zalcmann, A.T., and Talalay, P. 2001. The chemical diversity and distribution of
716 glucosinolates and isothiocyanates among plants. *Phytochemistry* 56(1):5-51.
717 https://doi.org/10.1016/S0031-9422(00)00316-2.

718 Ford-Lloyd, B. V., Schmidt, M., Armstrong, S. J., Barazani, O. Z., Engels, J., Hadas, R. et al. 2011.
719 Crop Wild Relatives—Undervalued, Underutilized and under Threat? *BioScience* 61(7); 559–565.
720 https://doi.org/10.1525/bio.2011.61.7.10.

721 Fukushima, A., Kusano, M., Mejia, R., Iwasa, M., Kobayashi, M.,Hayashi, N.,Watanabe-Takahashi,
722 A., et al. 2014. Metabolomic characterization of knockout mutants in Arabidopsis: development of a
723 metabolite profiling database for knockout mutants in *Arabidopsis*. *Plant Physiology* 165; 948–961.
724 https://doi.org/10.1104/pp.114.240986.

725 Gómez-Campo, C., and Tortosa M. 1974. The taxonomic and evolutionary significance of some
726 juvenile characters in the Brassiceae. *Botanical Journal of the Linnean Society* 69 (2): 105–124.
727 https://doi.org/10.1111/j.1095-8339.1974.tb01619.x.

728 Gómez-Campo C.1980 Morphology and morphotaxonomy of the tribe Brassiceae. In: Tsunoda S,
729 Hinata K, Gómez-Campo C. (ed) Brassica crops and wild allies. *Japan Scientific Societies Press*,
730 Tokyo, Pp.3-31.

731 Götz, S., Garcia-Gomez, J. M.,Terol, J., Williams, T. D., Nagaraj, S. H.,Nueda, M. J.,Robles, M.,
732 Talon, M.,Dopazo, J., and Conesa, A. 2008. High-throughput functional annotation and data mining
733 with the Blast2GO suite. *Nucleic Acids Research* 36 (10): 3420–35.
734 https://doi.org/10.1093/nar/gkn176.

735 Grubb, C. D., and Abel, S. 2006. Glucosinolate metabolism and its control. *Trends in Plant*
736 *Science* 11 (2): 89–100. https://doi.org/10.1016/j.tplants.2005.12.006.

737 Grubb, C. D., Zipp, B. J., Kopycki, J., Schubert, M., Quint, M., Lim, E., Bowles, D.J., Pedras, M. S.
738 C., and Abel, S. 2014. Comparative analysis of Arabidopsis UGT74 glucosyltransferases reveals a
739 special role of UGT74C1 in glucosinolate biosynthesis. *The Plant Journal* 79 (1): 92–105.
740 https://doi.org/10.1111/tpj.12541.

741 Gupta, S., Sangha, M. K ., Kaur, G., Atwa, A. K. A., Banga, S., and Banga, S. S.. 2012. Variability
742 for leaf and seed glucosinolate contents and profiles in a germplasm collection of the *Brassica*
743 *juncea*. *Biochemistry & Analytical Biochemistry* 01 (07). https://doi.org/10.4172/2161-
744 1009.1000120.

745 Halkier, B. A., and Gershenzon, J. 2006. Biology and biochemistry of glucosinolates. *Annual Review*
746 *of Plant Biology* 57 (1): 303–333. https://doi.org/10.1146/annurev.arplant.57.032905.105228.

747    Hansen, C. H., Wittstock, U., Olsen, C.E., Hick, A. J., Pickett, J. A., and Halkier, B.A.. 2001.
748    Cytochrome P450 CYP79F1 from *Arabidopsis* catalyzes the conversion of dihomomethionine and
749    trihomomethionine to the corresponding aldoximes in the biosynthesis of aliphatic
750    glucosinolates. *Journal of Biological Chemistry* 276 (14): 11078–85.
751    https://doi.org/10.1074/jbc.m010123200.

752    Hansen, B. G., Kliebenstein, D. J., and Halkier, B. A. 2007. Identification of a flavin-monooxygenase
753    as the S-oxygenating enzyme in aliphatic glucosinolate biosynthesis in *Arabidopsis*. *The Plant*
754    *Journal* 50 (5): 902–10. https://doi.org/10.1111/j.1365-313x.2007.03101.x.

755    Hansen, B. G., Kerwin, R. E., Ober, J. A., Lambrix, V. M., Mitchell-Olds, T., Gershenzon, J.,
756    Halkier, B. A., and Kliebenstein, D. J. 2008. A novel 2-oxoacid-dependent dioxygenase involved in
757    the formation of the goiterogenic 2-hydroxybut-3-enyl glucosinolate and generalist insect resistance
758    in *Arabidopsis*. *Plant Physiology* 148 (4): 2096–2108. https://doi.org/10.1104/pp.108.129981.

759    Heidel, A.J., Clauss, M. J., Kroymann, J., Savolainen O., and Mitchell-Olds, T. 2006. Natural
760    Variation in MAM Within and Between Populations of *Arabidopsis lyrata* Determines Glucosinolate
761    Phenotype. *Genetics* 173 (3): 1629–36. https://doi.org/10.1534/genetics.106.056986.

762    Ishida, M., Hara, M., Fukino, N., Kakizaki, T., and Morimitsu, Y. 2014. Glucosinolate metabolism,
763    functionality and breeding for the improvement of Brassicaceae vegetables. *Breeding Science* 64 (1):
764    48–59. https://doi.org/10.1270/jsbbs.64.48.

765    Kell, S.P., Knüpffer, H., Jury, S.L., Ford-Lloyd, B.V. and Maxted, N. 2008. Crops and wild relatives
766    of the Euro-Mediterranean region: making and using a conservation catalogue. In: Maxted, N., Ford-
767    Lloyd, B.V., Kell, S.P., Iriondo, J., Dulloo, E. and Turok, J. (eds) Crop Wild Relative Conservation
768    and Use. Pp. 69–109. CAB International, Wallingford, UK.

769    Khalifa, N.S. 2012. Protein expression after NaCl treatment in two tomato cultivars differing in salt
770    tolerance. *Acta Biologica Cracoviensia Series Botanica* 54 (2). https://doi.org/10.2478/v10182-012-
771    0020-0.

772    Khan, M.A.M., Ulrichs, C., and Mewis, I. 2011. Drought stress - impact on glucosinolate profile and
773    performance of phloem feeding cruciferous insects. *Acta Horticulturae*, 917; 111–17.
774    https://doi.org/10.17660/actahortic.2011.917.13.

775    Kliebenstein, D. J., Kroymann, J., Brown, P., Figuth, A., Pedersen, D.,Gershenzon, J., and Mitchell-
776    Olds, T. 2001a. Genetic control of natural variation in *Arabidopsis* glucosinolate accumulation. *Plant*
777    *Physiology* 126 (2): 811–25. https://doi.org/10.1104/pp.126.2.811.

778    Kliebenstein, D. J., Kroymann, J., Brown, P., Figuth, A., Pedersen, D.,Gershenzon, J., and Mitchell-
779    Olds, T. 2001b. Gene duplication in the diversification of secondary metabolism: tandem 2-
780    oxoglutarate-dependent dioxygenases control glucosinolate biosynthesis in *Arabidopsis*. *The Plant*
781    *Cell* 13 (3): 681. https://doi.org/10.2307/3871415.

782    Kumar, S., Stecher, G., Li,M., Knyaz, C., and Tamura K. 2018. MEGA X: molecular evolutionary
783    genetics analysis across computing platforms. Eds Battistuzzi, F. U. *Molecular Biology and*
784    *Evolution* 35 (6): 1547–49. https://doi.org/10.1093/molbev/msy096.

785    Lanfear, R., Frandsen, P. B., Wright, A. M., Senfeld, T., and Calcott, B. 2016. PartitionFinder 2: new
786    methods for selecting partitioned models of evolution for molecular and morphological phylogenetic
787    analyses. *Molecular Biology and Evolution* 34 (3): 772–73. https://doi.org/10.1093/molbev/msw260.

788    López-Berenguer, C., Martínez-Ballesta, M.C., Moreno, D. A., Carvajal, M., and García-Viguera, C.
789    2009. Growing hardier crops for better health: salinity tolerance and the nutritional value of broccoli.
790    *Journal of Agricultural and Food Chemistry* 57 (2): 572–78. https://doi.org/10.1021/jf802994p.

791    Martínez-Ballesta, M. C., Moreno, D. A., Carvajal, M. 2013. The physiological importance of
792    glucosinolates on plant response to abiotic stress in *Brassica*. *International Journal of Molecular*
793    *Sciences* 14(6): 11607-11625. https://doi.org/10.3390/ijms140611607.

794    Martínez-Ballesta, M. C., Muries, B., Moreno, D. A., Domínguez-Perles, R., García-Viguera, C., and
795    Carvajal, M. 2014. Involvement of a glucosinolate (sinigrin) in the regulation of water transport in
796    *Brassica oleracea* grown under salt stress. *Physiologia Plantarum* 150, 145–160. doi:
797    10.1111/ppl.12082

798    Martínez-Ballesta, M. C., Moreno-Fernández, D. A., Castejón, D., Ochando, C., Morandini, P. A.,
799    and Carvajal M. 2015. The impact of the absence of aliphatic glucosinolates on water transport under
800    salt stress in *Arabidopsis thaliana*. *Frontiers in Plant Science* 6:524. doi: 10.3389/fpls.2015.00524.

801    Merah, O. 2015. Genetic variability in glucosinolates in seed of *Brassica juncea*: interest in mustard
802    condiment. *Journal of Chemistry*. 2015: 606142. http://dx.doi.org/10.1155/2015/606142.

803    Mikkelsen, M. D., Naur, P., and Halkier, B. A. 2004. *Arabidopsis* mutants in the c-s lyase of
804    glucosinolate biosynthesis establish a critical role for indole-3-acetaldoxime in auxin homeostasis.
805    *The Plant Journal* 37 (5): 770–77. https://doi.org/10.1111/j.1365-313x.2004.02002.x.

806    Mithen, R., Clarke, J., Lister, C., and Dean, C. 1995. Genetics of aliphatic glucosinolates. III. Side
807    chain structure of aliphatic glucosinolates in *Arabidopsis thaliana*. *Heredity* 74 (2): 210–15.
808    https://doi.org/10.1038/hdy.1995.29.

809    Moreira-Rodríguez, M., Nair ,V., Benavides, J., Cisneros-Zevallos. L., and Jacobo-Velázquez, D.
810    2017a. UVA, UVB light doses and harvesting time differentially tailor glucosinolate and phenolic
811    profiles in broccoli sprouts. *Molecules* 22 (7): 1065. https://doi.org/10.3390/molecules22071065.

812    Moreira-Rodríguez, M., Nair ,V., Benavides, J., Cisneros-Zevallos. L., and Jacobo-Velázquez, D.
813    2017b. UVA, UVB light, and methyl jasmonate, alone or combined, redirect the biosynthesis of
814    glucosinolates, phenolics, carotenoids, and chlorophylls in broccoli sprouts. *International Journal of*
815    *Molecular Sciences* 18 (11): 2330. https://doi.org/10.3390/ijms18112330.

816    Ngwene, B., Neugart. S., Baldermann, S., Ravi, B. and Schreiner, M. 2017. Intercropping induces
817    changes in specific secondary metabolite concentration in ethiopian kale (*Brassica carinata*) and
818    african nightshade (*Solanum scabrum*) under controlled conditions. *Frontiers in Plant Science* 8
819    https://doi.org/10.3389/fpls.2017.01700.

820    Pfalz, M., Vogel, H., and Kroymann. J. 2009. The gene controlling the indole glucosinolate
821    modifier1 quantitative trait locus alters indole glucosinolate structures and aphid resistance in
822    *Arabidopsis*. *The Plant Cell* 21 (3): 985–99. https://doi.org/10.1105/tpc.108.063115.

823  Pfalz, M., Mikkelsen, M. D., Bednarek, P., Olsen, C. E., Halkier, B. A. and Kroymann., J. 2011.
824  Metabolic engineering in *Nicotiana benthamiana* reveals key enzyme functions in *Arabidopsis* indole
825  glucosinolate modification. *The Plant Cell* 23 (2): 716–29. https://doi.org/10.1105/tpc.110.081711.

826  Piotrowski, M., Schemenewitz, A., Lopukhina, A., Müller, A., Janowitz, T., Weiler, E. W. and
827  Oecking, C. 2004. Desulfoglucosinolate sulfotransferases from *Arabidopsis thaliana* catalyze the
828  final step in the biosynthesis of the glucosinolate core structure. *Journal of Biological Chemistry* 279:
829  50717–25. https://doi.org/10.1074/jbc.m407681200.

830  Radovich, T. J. K., Kleinhenz, M. D. and Streeter, J. G. 2005. Irrigation timing relative to head
831  development influences yield components, sugar levels, and glucosinolate concentrations in cabbage.
832  *Journal of the American Society for Horticultural Science* 130 (6): 943–49.
833  https://doi.org/10.21273/jashs.130.6.943.

834  Rambaut A. 2009. FigTree version 1.3.1 [computer program] http://tree.bio.ed.ac.uk .

835  Salehin, M., Li, B., Tang, M., Katz, E., Song, L., Ecker, J. R., Kliebenstein, D. J., and Estelle, M.
836  2019. Auxin-sensitive Aux/IAA proteins mediate drought tolerance in *Arabidopsis* by regulating
837  glucosinolate levels. *Nature Communications* 10 (1). https://doi.org/10.1038/s41467-019-12002-1.

838  Sánchez-Yélamo, M. D. 2009. Relationships in the *Diplotaxis–Erucastrum–Brassica* complex
839  (Brassicaceae) evaluated from isoenzymatic profiles of the accessions as a whole. Applications for
840  characterisation of phytogenetic resources preserved *ex situ*. *Genetic Resources and Crop Evolution*
841  56 (7): 1023–1036. https://doi.org/10.1007/s10722-009-9423-5.

842  Soundararajan, P., and Kim., J. 2018. Anti-carcinogenic glucosinolates in cruciferous vegetables and
843  their antagonistic effects on prevention of cancers. *Molecules* 23: 2983.
844  https://doi.org/10.3390/molecules23112983.

845  Tohge, T., Perez de Souza, L., and Fernie, A. R. 2014. Genome-enabled plant metabolomics. *Journal*
846  *of Chromatography B* 966: 7–20. https://doi.org/10.1016/j.jchromb.2014.04.003.

847  Tong, Y., Gabriel-Neumann, E., Ngwene, B., Krumbein, A., George, E., Platz, S., Rohn, S., and
848  Schreiner, M. 2014. Topsoil drying combined with increased sulfur supply leads to enhanced
849  aliphatic glucosinolates in *Brassica juncea* leaves and roots. *Food Chemistry* 152 (June): 190–96.
850  https://doi.org/10.1016/j.foodchem.2013.11.099.

851  Tripathi, M.K., and Mishra, A.S. 2007. Glucosinolates in animal nutrition: a review. *Animal Feed*
852  *Science and Technology* 132 (1–2): 1–27. https://doi.org/10.1016/j.anifeedsci.2006.03.003.

853  VanEtten, C. H., and H. L. Tookey. 1983. Glucosinolates. In *Handbook of Naturally Occurring Food*
854  *Toxicants*, 15–30. Boca Raton, Florida: CRC Press.

855  Variyar, P.S., Banerjee, A., Akkarakaran, J. J., and Suprasanna., P. 2014. Role of glucosinolates in
856  plant stress tolerance. *Emerging Technologies and Management of Crop Stress Tolerance*, 271–91.
857  https://doi.org/10.1016/b978-0-12-800876-8.00012-6.

858  Velasco, P., Soengas, P., Vilar, M., Cartea, M. E. and del Rio., M. 2008. Comparison of
859  glucosinolate profiles in leaf and seed tissues of different *Brassica napus* crops. *Journal of the*

860 *American Society for Horticultural Science* 133 (4): 551–58.

861 https://doi.org/10.21273/jashs.133.4.551.

862 Verkerk, R., S. Tebbenhoff, and M. Dekker. 2010. Variation and distribution of glucosinolates in 42

863 cultivars of *Brassica oleracea* vegetable crops. *Acta Horticulturae*, no. 856: 63–70.

864 https://doi.org/10.17660/actahortic.2010.856.7.

865 Verkerk, R., Schreiner, M., Krumbein, A., Ciska, E., Holst, B., Rowland, I., De Schrijver, R., et al.

866 2008. Glucosinolates in Brassica vegetables: the influence of the food supply chain on intake,

867 bioavailability and human health. *Molecular Nutrition & Food Research* 53: S219–S219.

868 https://doi.org/10.1002/mnfr.200800065.

869 Wang, H., Wu, J., Silong Sun, Liu, B., Cheng, F., Sun, R. and Wang., X. 2011. Glucosinolate

870 biosynthetic genes in *Brassica rapa*. *Gene* 487: 135–142. https://doi.org/10.1016/j.gene.2011.07.021.

871 Wang, X., Wu, J., Liang, J., Cheng, F. and Wang., X. 2015. Brassica database (BRAD) version 2.0:

872 integrating and mining Brassicaceae species genomic resources. *Database* 2015: bav093.

873 https://doi.org/10.1093/database/bav093.

874 Wiesner, M., Zrenner, R., Krumbein, A., Glatt, H. and Schreiner, M. 2013. Genotypic variation of

875 the glucosinolate profile in pak choi (*Brassica rapa* ssp. *chinensis*). *Journal of Agricultural and Food*

876 *Chemistry* 61: 1943–1953. https://doi.org/10.1021/jf303970k.

877 Windsor, A. J., Reichelt, M., Figuth, A., Svatoš, A., Kroymann, J., Kliebenstein, D. J., Gershenzon,

878 J. and Mitchell-Olds., T. 2005. Geographic and evolutionary diversification of glucosinolates among

879 near relatives of *Arabidopsis thaliana* (Brassicaceae). *Phytochemistry* 66: 1321–33.

880 https://doi.org/10.1016/j.phytochem.2005.04.016.

881 Wittstock, U., and B. A. Halkier. 2002. Glucosinolate research in the *Arabidopsis* Era. *Trends in*

882 *Plant Science* 7: 263–70. https://doi.org/10.1016/s1360-1385(02)02273-2.

883 Yuan, G., Wang, X., Guo, R. and Wang., Q. 2010. Effect of salt stress on phenolic compounds,

884 glucosinolates, myrosinase and antioxidant activity in radish sprouts. *Food Chemistry* 121: 1014–

885 1019. https://doi.org/10.1016/j.foodchem.2010.01.040.

886 Zhang, Y., Tang, L. and Gonzalez., V. 2003. Selected isothiocyanates rapidly induce growth

887 inhibition of cancer cells. *Molecular Cancer Therapy* 2 (10): 1045–1052.

888 Zhang, H., Schonhof, I., Krumbein, A., Gutezeit, B., Li, L., Stützel, H. and Schreiner. M. 2008.

889 Water supply and growing season influence glucosinolate concentration and composition in turnip

890 root (*Brassica rapa* ssp. *rapifera L.*). *Journal of Plant Nutrition and Soil Science* 171: 255–265.

891 https://doi.org/10.1002/jpln.200700079.

892

893 **Figures**

894

895  **Figure 1**- Genomic information of the GLS biosynthetic genes of the top ten Brassicaceae species
896  registered at NCBI database. Number of sequences available per species is represented above each
897  bar.
898
899  **Figure 2**- Euler diagram displaying GLS gene annotation gathered from a multi-database approach.
900
901  **Figure 3**- Gene Ontology (GO) terms assignment for the GLS biosynthetic genes. The graph displays
902  the term enrichment levels of the annotated sequences along with the GO term hierarchy: Biological
903  Process (BP, in green), Molecular Function (MF, in blue) and Cellular Component (CC, in yellow).
904
905  **Figure 4**- UPGMA phenogram of the 101 GLS biosynthetic genes. Bootstrap values above 50 are
906  represented on the branches. Detailed UPGMA tree is available at Supplementary Figure 1.
907
908  **Figure 5-** Phylogenetic tree from the Maximum Likelihood analysis of CYP79F1-F2 and CYP81F1-
909  F4 genes in Brassicaceae with *A. thaliana* CYP79F genes as outgroups. Acronyms are present as the
910  first letter of the genus and the second to species, e.g. At for *Arabidopsis thaliana*, and gene
911  identification when possible. Upon lack of complete CYP annotation, accession numbers were used.
912  Different copies of the same gene are identified by an "X" following sequential numbering, e.g. *A.*
913  *thaliana* X1, *A. thaliana* X2. Only bootstrap values above 50 are presented. Accession numbers of
914  the sequences analyzed are provided in Supplementary Table 3.
915
916  **Figure 6**– Chemodiversity profiling of GLS in *Brassica* and rocket species (*Diplotaxis* and *Eruca*).
917  Data matrix of GLS chemodiversity is provided in Supplementary Table 4. Colors indicate presence
918  (red) and absence (yellow) of a glucosinolate compound. Bootstraps values above 50 are presented in
919  the clustering phenogram resulted from the Eucledian distances method.
920

## Tables

922  **Table 1-** GLS genes information according to sub-pathways of indolic, aliphatic and aromatic.
923  Number of genes - total of genes annotated in each sub-pathway; Number of specific genes - genes
924  exclusive to a given sub-pathway; Number of shared genes - genes shared in at least two sub-
925  pathways.
926

## Supplementary Material

928
929  **Supplementary Figure 1**- Detailed UPGMA phenogram of the 101 GLS biosynthetic genes.
930  Bootstrap values above 70 are represented on the branches.
931
932  **Supplementary Figure 2**- PCA analysis using the 101 GLS genes (**A**) and PCAs loading plots (**B**) of
933  PCA 1 (above) and PCA 2 (below).
934
935  **Supplementary Table 1-** Gene compilation of GLS biosynthetic pathway using a multi-databasing
936  approach. For each gene, accession numbers annotated for *A. thaliana* are provided, alongside the
937  number of sequences available at NCBI database, restricted to Brassicaceae.
938

939 **Supplementary Table 2-** GLS gene classification according to each sub-pathway (indolic, aliphatic
940 and aromatic).

941

942 **Supplementary Table 3**- *CYP79F1-F2* and *CYP81F1-F4* sequences retrieved from NCBI database
943 to perform phylogenetic analysis in Brassicaceae available species. Accession numbers, species and
944 number of sequences are provided, together with code identification used in the Maximum
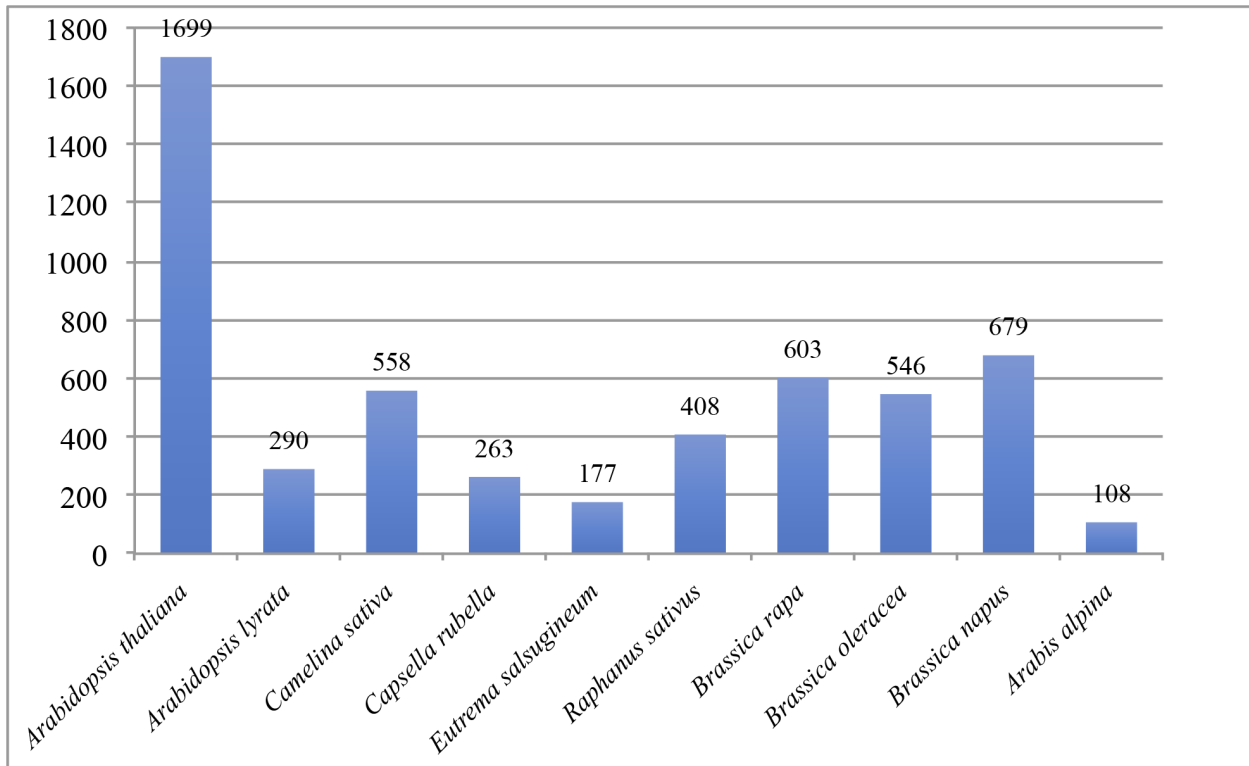945 Likelihood tree.

946

947 **Supplementary Table 4-** Data matrix of GLS used for a chemodiversity snapshot on *Brassica*
948 species (*B. napus, B. olearacea, B. rapa, B. juncea*), *Eruca vesicaria* and several wild rocket
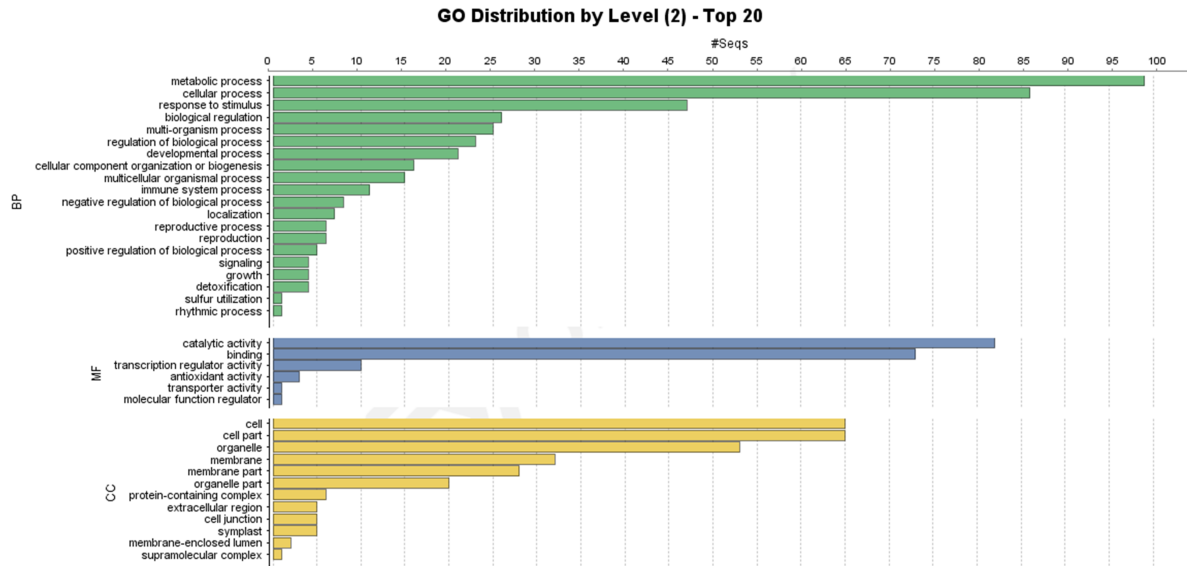949 *Diplotaxis* species.

950

**Table 1-** GLS genes information according to sub-pathways of indolic, aliphatic and aromatic. Number of genes - total of genes annotated in each sub-pathway; Number of specific genes - genes exclusive to a given sub-pathway; Number of shared genes - genes shared in at least two sub-pathways.

| | Aliphatic | Aromatic | Indolic | Combined unigenes of the 3 pathways |
|---|---|---|---|---|
| **Nº. Genes** | 40 | 20 | 41 | 78 |
| **Nº. Specific Genes** | 31 | 6 | 26 | - |
| **Nº. Shared Genes** | 9 | 14 | 14 | 8 |

**GO Distribution by Level (2) - Top 20**

Tree scale: 0.1

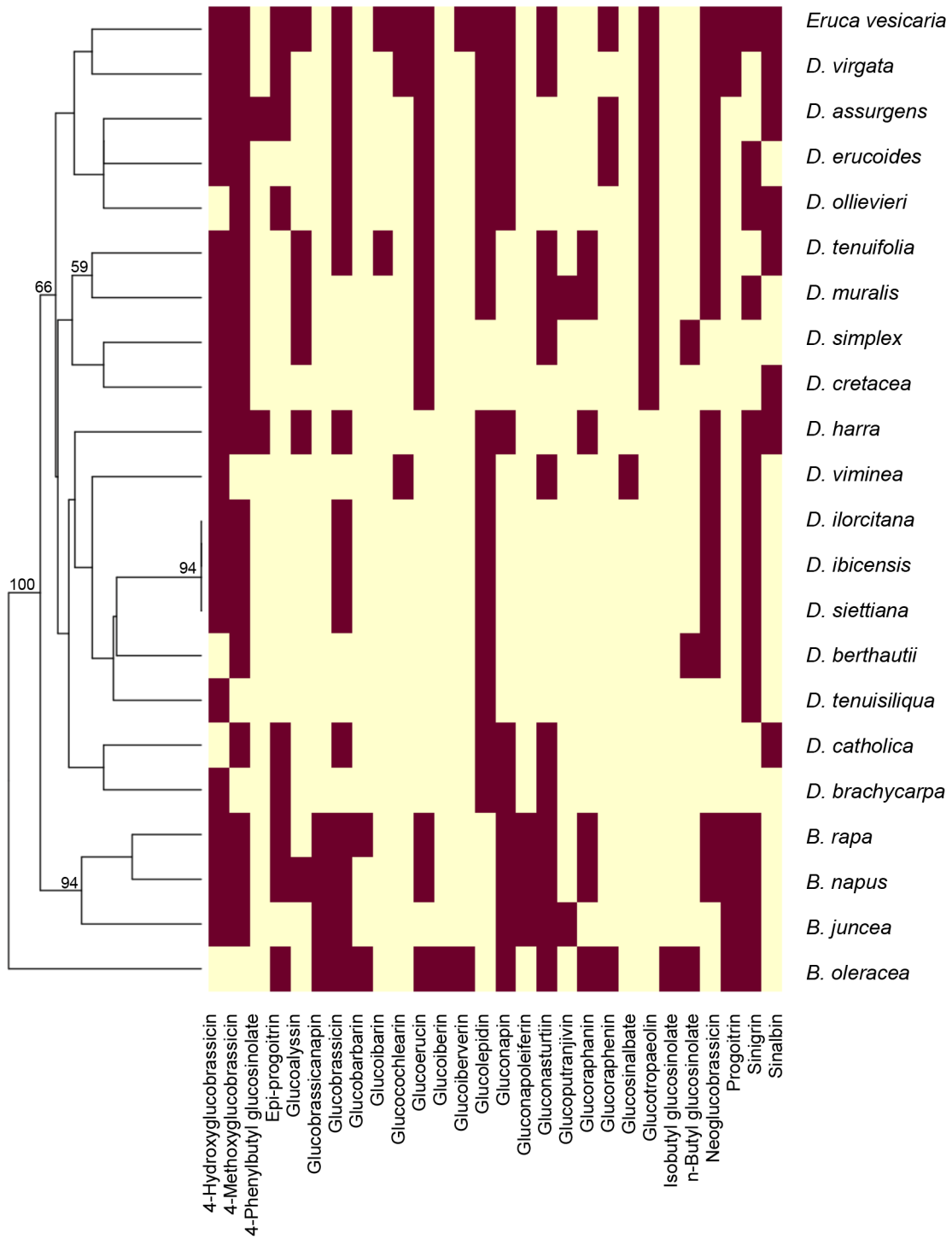| | Aliphatic | Aromatic | Indolic |
|---|---|---|---|
| CYP79s | ■ | | ■ |
| CYP81-83s | ■ | ■ | ■ |
| GSTU20 | ■ | | |
| T5M16.12 | | | ■ |
| MJG14.10 | | | ■ |
| MGN6.20 | | | ■ |
| IGMT5 | | | ■ |
| IGMT1 | | | ■ |
| IGMT4 | | | ■ |
| IGMT2 | | | ■ |
| IGMT3 | | | ■ |
| BCAT-3 | ■ | | |
| GSL-OH | ■ | | |
| IPMI LSU1 | ■ | | |
| MAM1-3 | ■ | | |
| FMOGS-OX6 | ■ | | |
| MYB118 | ■ | | |
| MYB115 | ■ | | |
| MYBs | ■ | ■ | ■ |
| BAT5 | ■ | | |
| GGT4 | ■ | ■ | ■ |
| FMOGS-OX1-5,7 | ■ | | |
| AtIMD2 | ■ | | |
| IPMDH1 | ■ | | |
| IPMDH3 | ■ | | |
| GGP3 | | ■ | |
| T29E15.17 | | ■ | |
| F17I23.120 | | ■ | |
| GGP1 | ■ | | ■ |
| T29E15.16 | | ■ | |
| TAT3 | | ■ | ■ |
| RSA1 | | ■ | ■ |
| F20O9.100 | | ■ | ■ |
| SUR1 | ■ | ■ | ■ |
| IPMI_SSU1-3 | ■ | | |
| UGT74 | ■ | ■ | ■ |
| AOP2 | ■ | | |
| AOP3 | ■ | | |
| GSTF9-11 | ■ | | ■ |
| ST5abc | ■ | ■ | ■ |

**Highlights**

- Brassicaceae genes involved in GLS biosynthesis were identified using a multi-database approach
- UPGMA and PCA separation between genes in GLS core structure and CYP450/MYB gene families.
- Phylogenetics revealed a recent diversification of aliphatic genes and an earliest for indolic.
- Distinct GLS chemo-profiles between *Brassica* crops and *Diplotaxis* species, wild relatives.
- GLS-rich species as a new source of taxa with great agronomic potential for abiotic stress tolerance.

## CONTRIBUTION

The Brassicaceae family is one of the world's most economically important plant groups. They include important crop species (e.g., *Brassica* spp.), weeds (e.g., *Capsella*, *Lepidium*, *Sisymbrium*, and *Thlaspi*), ornamentals (e.g., *Hesperis*, *Lobularia*, and *Matthiola*), and the model organism for flowering plants *Arabidopsis thaliana*. Among the most important chemical compounds produced by Brassicaceae species, are glucosinolates (GLS) with proven and widely documented health promoting effects. Glucosinolates have been the subject of several studies in Brassicaceae as important chemical compounds, particularly in chemical assessment in commercial crops, and also on the characterization of its biochemical pathway reconstruction. However, an integrated approach covering genomic, phylogenetic and chemical analysis in GLS pathway in Brassicaceae remains limited. There are several novel and important aspects to our paper, namely it is the first time where a taxa approach is performed on GLS pathway genes in Brassicaceae species, while in *A. thaliana* its assessment has been extensively studied.

In our paper, we looked through gene clusters involved in the biosynthesis of GLS, by combining genome analysis with biochemical pathways and chemical diversity assessment. Considering the high diversity in GLS content in Brassicaceae species, an integrated approach was performed by assessing a global GLS gene inventory in Brassicaceae and its diversity, analysing a potential genetic basis for GLS divergence using 6 CYP genes (CYP79F1-F2 and CYP81F1-F4), known to be key genes of indolic and aliphatic GLS biosynthetic pathways, linked to a chemical diversity evaluation of GLS compounds in major *Brassica* crops compared to the wild relative genus (*Diplotaxis*).

Our results point to a recent diversification of the aliphatic CYP79F1 and F2 genes in *Brassica* crops, while for indolic genes a clear separation is observed for CYP81F1-F4 genes, revealing an earliest divergence on this GLS sub-pathway. Chemical diversity snapshot allowed recognizing that *Brassica* and *Diplotaxis* have distinct GLS chemo-profiles, highlighting that the latter genus appears as an extremely GLS-rich species. Given the importance of GLS in abiotic stress tolerance, we have explored *Diplotaxis* species, the closest wild relatives of *Brassica* crops, as a new source of taxa with great agronomic potential. Understanding the genomic diversity responsible for the corresponding GLS biosynthetic pathways linked to the chemical diversity could bring insights for exploring new opportunities for using GLS-rich species, yet unexplored.

In summary, this work provides an integrated framework to analyse the chemical diversity of GLS in Brassicaceae, and provides data that complement current state of the art studies performed in GLS within Brassicaceae to answer a wide range of scientific questions in the fields of the genomic basis of chemical diversity and on species diversity assessment using an integrative approach.

**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.