

# 1 **Detection of orthologous genes with expression shifts linked to** 2 **nickel hyperaccumulation across Eudicots**

3 Mélina Gallopin<sup>1\*</sup>, Christine Drevet<sup>1</sup>, Vanesa S. Garcia de la Torre<sup>1</sup>, Sarah Jelassi<sup>1</sup>, Marie  
4 Michel<sup>1,4</sup>, Claire Ducos<sup>1,5</sup>, Cédric Saule<sup>1</sup>, Clarisse Majorel<sup>2</sup>, Valérie Burtet-Sarramegna<sup>2</sup>, Yohan  
5 Pillon<sup>3</sup>, Paul Bastide<sup>4</sup>, Olivier Lespinet<sup>1</sup>, Sylvain Merlot<sup>1\*</sup>

6 <sup>1</sup>*Université Paris-Saclay, CEA, CNRS, Institute for Integrative Biology of the Cell (I2BC), 91198,*  
7 *Gif-sur-Yvette, France*

8 <sup>2</sup>*Institute of Exact and Applied Sciences (ISEA), Université de la Nouvelle-Calédonie, BP R4,*  
9 *Nouméa Cedex, 98851, New Caledonia*

10 <sup>3</sup>*Laboratoire des Symbioses Tropicales et Méditerranéennes (LSTM), IRD, INRAE, CIRAD,*  
11 *Institut Agro, Univ. Montpellier, Montpellier, France*

12 <sup>4</sup>*IMAG, Université de Montpellier, CNRS, 34000, Montpellier, France*

13 <sup>5</sup> *Present address: Theories and Approaches of Genomic Complexity (TAGC), INSERM U1090,*  
14 *Parc scientifique de Luminy, case 928, 13009 Marseille, France*

15 <sup>6</sup> *Present address: Hôpital Gustave Roussy, Bat. B2M, 114 rue Edouard Vaillant, 94805, Villejuif,*  
16 *France*

17 \* For correspondence: [melina.gallopin@i2bc.paris-saclay.fr](mailto:melina.gallopin@i2bc.paris-saclay.fr); [sylvain.merlot@i2bc.paris-saclay.fr](mailto:sylvain.merlot@i2bc.paris-saclay.fr)

18 **Keywords:** differential expression analysis, RNA-seq, phylogenetic comparative methods,  
19 ortholog groups, nickel hyperaccumulation

20

## 21 **Abstract**

22 The remarkable capacity of plants to tolerate and accumulate tremendous amount of nickel is a  
23 complex adaptative trait that appeared independently in more than 700 species distributed in about  
24 fifty families. Nickel hyperaccumulation is thus proposed as a model to investigate the evolution  
25 of complex traits in plants. However, the mechanisms involved in nickel hyperaccumulation are  
26 still poorly understood in part because comparative transcriptomic analyses struggle to identify  
27 genes linked to this trait from a wide diversity of species. In this work, we have implemented a  
28 methodology based on the quantification of the expression of orthologous groups and phylogenetic  
29 comparative methods to identify genes which expression is correlated to the nickel  
30 hyperaccumulation trait. More precisely, we performed *de novo* transcriptome assembly and reads  
31 quantification for each species on its own transcriptome using available RNA-Seq datasets from  
32 15 nickel hyperaccumulator and non-accumulator species. Assembled contigs were associated to  
33 orthologous groups built using proteomes predicted from completed plant genome sequences. We  
34 then analyzed the transcription profiles of 5953 orthologous groups from distant species using a  
35 phylogenetic ANOVA. We identified 31 orthologous groups with an expression shift associated  
36 with nickel hyperaccumulation. These orthologous groups correspond to genes that have been  
37 previously implicated in nickel accumulation, and to new candidates involved in this trait. We thus  
38 believe that this method can be successfully applied to identify genes linked to other complex traits  
39 from a wide diversity of species.

40

## 41 **Introduction**

42 Comparative biology is a fundamental strategy to study the evolution of complex developmental  
43 and physiological traits. The development of the RNA-Seq technology has opened the possibility  
44 to perform comparative studies at the molecular level in a wide diversity of plant species (Leebens-  
45 Mack *et al.*, 2019). In plants, RNA-Seq data have been used to compare distant species for  
46 phylogenetic tree inference and to study the evolution of complex traits such as the type of  
47 photosynthesis or the capacity to establish symbioses (Jiao *et al.*, 2011; Wickett *et al.*, 2014; Yang  
48 *et al.*, 2015; Heyduk *et al.*, 2019; Radhakrishnan *et al.*, 2020; Rich *et al.*, 2021). In most of these  
49 studies, RNA-Seq data are used to reveal gene loss, gene duplication and genetic variations linked  
50 to a specific trait in different phylae. However, because of the difficulty to use a unique sequence  
51 as reference, comparative studies rarely take full advantage of the quantitative information  
52 enclosed in RNA-Seq data to compare gene expression between distant species and identify genes  
53 which expression is linked to a particular trait (Roux *et al.*, 2015; Voelckel *et al.*, 2017; García de  
54 la Torre *et al.*, 2021; Rich *et al.*, 2021).

55 Metal hyperaccumulation represents an interesting case study to identify genes linked to a complex  
56 adaptative trait over a wide diversity of plant species (Manara *et al.*, 2020). Metal  
57 hyperaccumulation is defined as the capacity of plant species to accumulate in their leaves a high  
58 concentration of metal, such as nickel, manganese, zinc or cadmium, that is normally toxic for the  
59 vast majority of plants (van der Ent *et al.*, 2013). Today, about 700 plant species are known to  
60 hyperaccumulate metals but the large majority (*ie* 500 species) hyperaccumulates nickel (Reeves  
61 *et al.*, 2018). Nickel hyperaccumulators are distributed in about 50 families among more than 300  
62 families of dicotyledon plants suggesting that the nickel hyperaccumulation trait appeared  
63 independently in several clades along plant evolution (Krämer, 2010; Cappa and Pilon-Smits,

64 2014). Comparative transcriptomic analysis of zinc and cadmium hyperaccumulators and non-  
65 accumulator species of the Brassicaceae family has first revealed that metal hyperaccumulation is  
66 linked to the high and constitutive expression of several genes involved in metal transport and  
67 homeostasis (Hammond *et al.*, 2006; Weber *et al.*, 2006; Hanikenne *et al.*, 2008; Halimaa *et al.*,  
68 2014). Our knowledge of the molecular mechanisms involved in nickel hyperaccumulation is still  
69 limited but, as for the hyperaccumulation of zinc and cadmium, the hyperaccumulation of nickel  
70 likely evolved from the high and constitutive expression of genes involved in metal homeostasis.  
71 Halimaa *et al.* (2014) used SOLiD-based RNA-Seq approach to compare the expression of genes  
72 from three accessions of *Noccaea caerulescens* (Brassicaceae) with various abilities to tolerate and  
73 accumulate metals in order to identify genes linked to metal hyperaccumulation including nickel.  
74 In this study, the authors used the genome of the related model species *Arabidopsis thaliana* as a  
75 common reference to align RNA-Seq reads. Using an Illumina-based RNA-Seq approach, Meier  
76 *et al.* (2018) compared the expression of genes from various populations of *Senecio coronatus*  
77 (Asteraceae) hyperaccumulating (NiH) or not (NA) nickel. The authors generated a *S. coronatus*  
78 reference transcriptome by *de novo* assembly to quantify gene expression and identify differentially  
79 expressed genes in both type of populations. More recently, we used the same RNA-Seq technology  
80 to identify genes differentially expressed in pairs of NiH and closely related NA species from five  
81 distant plant families (García de la Torre *et al.*, 2021). Then, to identify convergent mechanisms  
82 involved in nickel accumulation, we used orthologous relationship between genes from these  
83 distant families and a multiple testing correction to identify orthologous groups (OG) containing  
84 genes differentially expressed between NiH and NA species in at least 3 plant families.

85 These comparative approaches rely on the possibility to have access to pairs of closely related  
86 species or populations with contrasting capacity to accumulate metals in order to use a common

87 reference sequence. However, the identification of such pairs of species is not possible in all plant  
88 clades. In addition, the output of these analyses strongly depends on the specific pair of species  
89 that have been selected for the study. Finally, none of these studies take into account the  
90 phylogenetic tree and the genetic drift associated with the selected species.

91 Phylogenetic relationships between species are known to induce correlations between trait  
92 measurements that can affect the analyses when ignored (Felsenstein, 1985). Phylogenetic  
93 Comparative Methods (PCMs) precisely aim at taking these relationships into account, and have  
94 been extensively studied over the last few decades (Harmon, 2019). In the context of gene  
95 expression, Bedford & Hartl (2009) studied several stochastic models, including the Ornstein-  
96 Uhlenbeck (OU) process. This process can be seen as modeling the evolution of a quantitative trait  
97 under stabilizing selection towards an optimal value (Hansen, 1997). Building on this process,  
98 Rohlf and Nielsen (2015) proposed a phylogenetic ANOVA framework that takes into account  
99 both phylogenetic and individual variations. Individual variations represent both intra-specific  
100 variations and measurement errors, and ignoring them can lead to severe bias in PCMs (Silvestro  
101 *et al.*, 2015; Cooper *et al.*, 2016). This framework has been used to detect OG with significant  
102 mean expression shifts across groups of species from various clades of animal kingdom (Rohlf  
103 and Nielsen, 2015; Stern and Crandall, 2018; Chen *et al.*, 2019; Catalán *et al.*, 2019).

104 In this work, we have implemented a methodology to identify orthologous groups (OG) with  
105 expression shifts linked to the nickel hyperaccumulation trait in a wide diversity of plant species.

106 This method uses RNA-Seq datasets to produce reference transcriptomes by *de novo* assembly and  
107 quantify gene expression in each species. Genes from the different species are then associated to  
108 OG and the expression of OG identified in all species is then analyzed by PCM. Using this  
109 methodology, we have identified OGs previously associated with nickel hyperaccumulation as well

110 as new candidate genes involved in this trait. We believe that this methodology can be used more  
111 generally to identify genes associated to complex traits in a wide diversity of species.

112

## 113 **Materials and Methods**

### 114 *Collection of RNA-Seq datasets*

115 RNA-Seq datasets used in this study were collected from NCBI bioprojects PRJNA476917 (García  
116 de la Torre *et al.*, 2021), PRJNA312157 (Meier *et al.*, 2018) and PRJNA657163. The selected  
117 samples correspond to RNA extracted from leaves of nickel hyperaccumulator (NiH) and non-  
118 accumulator (NA) species or accessions from the families Asteraceae, Brassicaceae, Cunoniaceae,  
119 Phyllanthaceae, Rubiaceae and Salicaceae (representing five orders of Eudicots) and sequenced  
120 with the Illumina HiSeq2000 paired-end sequencing technology. Information on the different  
121 samples is summarized in **Table 1**.

### 122 *Transcriptome assembly and expression quantification*

123 Several *de novo* assembled transcriptomes used in this study were previously published and  
124 available from the bioproject PRJNA476917 (García de la Torre *et al.*, 2021). The transcriptomes  
125 of *Senecio coronatus* and *Microthlaspi perfoliatum* were assembled *de novo* with QIAGEN CLC  
126 Genomics Workbench v9 using the same assembly parameters (similarity  $\geq 0.95$ ; length fraction  $\geq$   
127 0.75) as used in García de la Torre *et al.* (2021). For these assemblies, we used the paired-end  
128 Illumina reads SRX1901479 (*S. coronatus*), SRX8947157 and SRX8947158 (*M. perfoliatum*).  
129 For each species, the reads of each sample were mapped to the corresponding *de novo*  
130 transcriptome using CLC Genomics Workbench v9 (similarity  $\geq 0.875$ ; length fraction  $\geq 0.75$ ).

131

132 ***Construction of ortholog group seeds and annotation***

133 The sequence of predicted proteomes encoded by 12 plant genomes were downloaded from the

134 PLAZA 4.0 Dicots database (Van Bel *et al.*, 2018): *Arabidopsis thaliana* and *Brassica rapa*

135 (Brassicaceae), *Gossypium raimondii* and *Theobroma cacao* (Malvaceae), *Carica papaya*

136 (Caricaceae), *Prunus persica* (Rosaceae), *Cucumis melo* (Cucurbitaceae), *Glycine max* (Fabaceae),

137 *Ricinus communis* (Euphorbiaceae), *Populus trichocarpa* (Salicaceae), *Solanum lycopersicum*

138 (Solanaceae) and *Coffea canephora* (Rubiaceae). We used the meta-approach MARIO to build

139 ortholog group (OG) seeds using the 12 proteome sequences (Pereira *et al.*, 2014).

140 We performed the annotation of these OGs using the HMMER package (Eddy, 1998). For each

141 OG resulting from the MARIO output, we first performed a multiple alignment with MUSCLE

142 (Edgar, 2004) and created a HMM profile using hmmbuild. A profile database of the OG seeds

143 was created using hmmcompress. To annotate the OGs, we searched for the closest homologs in the

144 SwissProt database with hmmsearch using the HMM profiles as queries. We extracted the function,

145 EC numbers, and GO terms from the hmmsearch hits with the lowest e-value (e-value  $\leq 10e-45$ )

146 and transferred these annotations to the corresponding OGs.

147

148 ***Assignment of contigs to Orthologous Groups***

149 For each assembled transcriptome, we searched in each contig for the longest ORF on the forward

150 and reverse strands and translated this ORF to obtain the putative encoded protein. The assignment

151 of contigs to OGs was performed with the HMMER package (Eddy, 1998). We performed

152 hmmscan using the translation of the longest ORF of each contig as a query against the OG seeds

153 profile database. We assigned contigs to the OG profile having the lowest e-value (e-value  $\leq 1e-$   
154 10, coverage  $\geq 20\%$ ).

155  
156 ***OG expression matrix construction***  
157 We built an OG expression matrix associating each OG to its level of expression in each sample.  
158 The expression level of an OG was calculated as the sum of the read counts corresponding to the  
159 contigs assigned to this OG in each sample. We also summed the lengths of all contigs assigned  
160 to a single OG to obtain an OG length matrix.

161  
162 ***Data normalization and transformation***  
163 We computed the normalization factors  $n_i$  for each sample  $i$  using the TMM method, implemented  
164 in the function `calcNormFactor` of the edgeR package (Robinson and Oshlack, 2010). To take the  
165 length of the OGs (sums of the lengths of each contigs within each OG ) and the size of libraries  
166 into account, we then computed log2 RPKM [reads per kilobase per million reads, (Mortazavi *et*  
167 *al.*, 2008)] using the following formula:

168 
$$y_{gi} = \log_2 \left( \frac{\frac{c_{gi} + 0.5}{l_{gi}}}{C_i n_i + 1} \times 10^9 \right)$$

169 where  $C_i = \sum_{g=1}^G c_{gi}$  is the library size of sample  $i$ ,  $c_{gi}$  is the read count for OG  $g$  in sample  $i$ , and  
170  $l_{gi}$  is the length of the OG  $g$  in sample  $i$ . Note that the length can vary between samples for the  
171 same OG, as samples are taken from different species. To ensure that the ratio inside the log is



172 strictly less than 1 and greater than zero,  $c_{gi}$  and  $C_i n_i$  were offset away from zero by adding 0.5  
173 and 1 respectively (Law *et al.*, 2014).

174 To perform PCA on the OG expression matrix, we used the `plot.PCA` function of the R package  
175 DESeq2. We used the `rlog` function to log-transform the data prior to the PCA.

176  
177 ***Phylogenetic tree***

178 A custom dated phylogenetic tree was built from the plant tree backbone at the family level  
179 proposed by Magallón *et al.* (2015). The tree topology and time divergence used in our study are  
180 based on (Barrabé *et al.*, 2014; Igea *et al.*, 2015; Razafimandimbison *et al.*, 2017) for Rubiaceae,  
181 (Pillon *et al.*, 2014) for Cunoniaceae and (Huang *et al.*, 2016) for Brassicaceae.

182  
183 ***Phylogenetic ANOVA***

184 We used the R package `phylolm` (Ho and Ané, 2014) to perform phylogenetic ANOVA, using the  
185 "OU with fixed root" model, with measurement error. In a phylogenetic regression using an OU  
186 with fixed root model, the residuals are assumed to be correlated with a correlation between two  
187 species that depends on their shared evolutionary time (*ie* the time between the root of the tree and  
188 the most recent common ancestor of the two species). In addition to the phylogenetic residuals, we  
189 included in the model independent identically distributed residuals to capture additional non-  
190 phylogenetic variance, that we fitted to the data. The factor of interest was the nickel  
191 hyperaccumulation capacity of the species or accession. The design matrix also included an  
192 intercept, and a factor representing the country of origin of the sample. For each OG, a p-value was  
193 computed, corresponding to a t-test (with correlated observations) on the coefficient associated to  
194 the hyperaccumulator factor. A Benjamini-Hochberg multiple testing correction was applied to the

195 vector of p-values. We used a threshold of 0.01, and selected OGs with a log<sub>2</sub> fold change  $\geq 1.5$  or  
196  $\leq -1.5$ .

197 **Results**

198  
199 ***A methodology for detection of orthologous genes with mean expression shifts in nickel***  
200 ***hyperaccumulators from distant plant families***

201 In this work we wanted to develop a methodology to identify genes whose expression is linked to  
202 nickel hyperaccumulation across a wide diversity of plant species (**Figure 1**). We took advantage  
203 of RNA-Seq datasets previously generated from nickel hyperaccumulator (NiH) and related non-  
204 accumulator (NA) species or populations to generate *de novo* transcriptome assemblies for each  
205 species and then use these transcriptomes as references to quantify gene expression for each sample  
206 corresponding to these species. This methodology also uses the concept of orthologous groups  
207 (OG) to annotate genes putatively playing conserved functions in distant plant families (Altenhoff  
208 *et al.*, 2012). Finally, we quantified the expression of OG in NiH and NA species groups and  
209 analyzed the data with a Phylogenetic Comparative Method (PCM) to identify OG with an  
210 expression shift linked to the nickel hyperaccumulation trait.

211  
212 ***De novo transcriptome assembly and quantification of contig expression from nickel***  
213 ***hyperaccumulator and non-accumulator species***

214 We used available RNA-Seq from nickel hyperaccumulator (NiH) and non-accumulator (NA)  
215 species or accessions that were generated using Illumina paired-end technology (**Table 1**). For  
216 most of the NiH and NA species used in this study, the transcriptomes were previously assembled  
217 *de novo* using the CLC Genomic Workbench software (García de la Torre *et al.*, 2021). We  
218 assembled the transcriptomes of *Senecio coronatus* and *Microthlaspi perfoliatum* using the same  
219 parameters. The number of contigs and the median size of the contigs for each species is given in

220 **Table 2.** The number of contigs obtained for *M. perfoliatum* is significantly higher than for the  
221 other species probably due to the tetraploid nature of this species and because we assembled reads  
222 from both roots and shoots samples.

223 The RNA-Seq reads from each replicate corresponding to the leaf samples of all NiH and NA  
224 species or populations (**Table 1**) were mapped to the corresponding *de novo* transcriptome. This  
225 generated a read count table for each contig in all sample replicate from each species.

226

### 227 *Construction of orthologous group database and assignation of contigs*

228 In the second step of our methodology, we wanted to establish orthologous relationships between  
229 genes expressed in the NiH and NA species in order to annotate genes potentially playing similar  
230 biological functions across these species (Altenhoff *et al.*, 2012). However, because the translation  
231 of *de novo* assembled transcriptomes generates a large number of truncated peptides that could  
232 affect the analysis of the orthologous relationships, we decided to first build up an orthologous  
233 group (OG) library using proteomes predicted from sequenced plant genomes. We selected 12  
234 sequenced plant genomes chosen along the dicotyledon phylogenetic tree and including species  
235 belonging to the same families as the nickel hyperaccumulators (see Methods). These proteomes,  
236 containing from 27000 to 56000 peptides, were used to create Orthologous Group (OG) seeds using  
237 the MARIO meta-approach (Pereira *et al.*, 2014). MARIO combines the results of four methods:  
238 Best Reciprocal Hits, Inparanoid (O'Brien, 2004), OrthoFinder (Emms and Kelly, 2015) and  
239 Phylogeny (Lemoine *et al.*, 2007) to establish orthologous relationships between proteins and  
240 compute a consensus OG annotation. Using this method, we obtained 17830 OGs containing from  
241 2 to 1466 proteins. We could attribute a function to 11301 OG (63 %) using the Swissprot database  
242 and 4486 OGs (25%) were associated to at least one EC number. 7779 OGs (43%) are represented

243 by at least one peptide in all model species. An HMM profile database was created from these OG  
244 seeds.

245 We then assigned each contig of the assembled transcriptomes from NiH and NA species to the  
246 OG seeds constructed with MARIO. For each contig, the longest Open Reading Frame (ORF) was  
247 identified and translated into a peptide. Each peptide was then associated to the closest OG using  
248 the HMM profile database (see Methods). Depending on the species, 44 to 69% of the contigs were  
249 assigned to an OG (**Table 2**). In total, among the 17830 OGs generated with the proteome of model  
250 plants, 15941 OGs are represented by at least one contig from at least one plant species of interest.  
251 It is important to note here that the RNA-Seq data used in this study only represent gene expressed  
252 in leaves. More importantly, for cross-species comparison, 5953 OGs are represented by at least  
253 one contig in each studied species.

254  
255 ***Identification of differentially expressed Ortholog Groups between nickel hyperaccumulator and***  
256 ***non-accumulator species***

257 We focused our differential expression analysis on the 5953 OGs represented in all plant species.  
258 We built an OG expression matrix (see Methods), each row representing one of the 5953 OGs and  
259 each column representing an RNA-Seq sample from 9 nickel hyperaccumulator (NiH) and 9 non-  
260 accumulator (NA) species or populations (2 to 3 biological replicates per species). We first  
261 performed principal component analysis (PCA) of the dataset after normalization and a log  
262 transformation of the expression data to explore the correlation between samples (**Figure 2**). The  
263 representation of the two principal components having the most important effect on total variation  
264 indicates that the RNA-Seq samples do not cluster with respect to the NiH trait but rather with  
265 respect to the plant family to which they belong. This result further illustrates that it is important

266 to take into consideration the phylogenetic relationships between species in the transcriptomic  
267 comparison between the NiH and NA groups.

268 To identify OGs differentially expressed between NiH and NA species from distant plant families,  
269 we used a phylogenetic mixed model (Lynch, 1991; Housworth et al., 2004) implemented in the R  
270 package phylolm (Ho and Ané, 2014).

271 Using this approach, we identified 31 OGs differentially expressed between NiH and NA species  
272 with a log<sub>2</sub>FC threshold  $\geq 1.5$  or  $\leq -1.5$ . and an FDR  $\leq 0.01$  (**Table 3**). We used the length  
273 normalized unit RPKM to compare the expression of OG between samples. We also performed the  
274 analysis with the TPM unit (transcripts per million) (Wagner *et al.*, 2012), but the results did not  
275 differ significantly (data not shown). 17 OGs are more expressed in the group of NiH species. The  
276 expression in NiH and NA species of OG 4147 and OG 7137, coding for a histidinol dehydrogenase  
277 and a membrane transporter of the SLC40A family respectively, is presented in **Figure 3**. These  
278 data illustrate the higher expression of these OG in several NiH species or populations compared  
279 to NA species. In contrast, 14 OG are less expressed in NiH species including OGs coding for a  
280 calcium-transporting ATPase (OG 10076) and a probable receptor-like serine/threonine-protein  
281 kinase related to Lr10 (OG 13722).

## 282 **Discussion**

283  
284 The comparison of the expression of genes playing conserved function over a wide diversity of  
285 plant species is still a challenging task in comparative and evolution biology. The goal of this work  
286 was to implement a method using the annotation of orthologous groups (OG) and a Phylogenetic  
287 Comparative Method (PCM) to compare the transcriptomes of evolutionary distant species to  
288 identify genes whose expression level is linked with the nickel hyperaccumulation trait.

### 289 290 ***Identification of orthologous genes with expression shifts linked to nickel hyperaccumulation***

291 Genes linked to the nickel hyperaccumulation trait have been previously searched by comparing  
292 gene expression in pairs of closely related species or populations of the same species with  
293 contrasted capacity to accumulate nickel (Halimaa *et al.*, 2014; Meier *et al.*, 2018; García de la  
294 Torre *et al.*, 2021; Enomoto *et al.*, 2021). In this study, we have used available and comprehensive  
295 RNA-Seq datasets to identify OGs with expression shifts between nickel hyperaccumulator (NiH)  
296 and non-accumulator (NA) groups of species. Interestingly, we identified OG 7137 (**Table 3**,  
297 **Figure 3**) corresponding to the SLC40A membrane transporter family, also known as IREG or  
298 Ferroportin (FPN) transporters, as an OG more expressed in NiH compared to NA species.  
299 IREG/FPN transporters are able to transport divalent metal ions including nickel across membranes  
300 (Schaaf *et al.*, 2006; Morrissey *et al.*, 2009; Billesbølle *et al.*, 2020). In *Arabidopsis thaliana*, two  
301 transporters belonging to OG 7137, AtIREG1 and AtIREG2, have been shown to localize on the  
302 plasma membrane and the vacuolar membrane respectively (Schaaf *et al.*, 2006; Morrissey *et al.*,  
303 2009). High expression of vacuolar localized IREG transporters in transgenic *A. thaliana* increases  
304 nickel tolerance (Schaaf *et al.*, 2006; Merlot *et al.*, 2014). The high expression of IREG/FPN genes

305 was previously shown to be linked to the nickel hyperaccumulation trait in species from several  
306 plant families (Meier *et al.*, 2018; García de la Torre *et al.*, 2021). The identification of OG 7137  
307 thus validates the capacity of our methodology to identify OG linked to nickel hyperaccumulation.  
308 It also indicates that the high expression of IREG transporters in leaves is a robust characteristic of  
309 nickel hyperaccumulators that can be observed independently of the method used for cross species  
310 transcriptomic comparison. Our analysis also revealed a higher expression of OG 8462 and OG  
311 6278, corresponding to the Ferric reductase FRO2 and the Cationic amino acid transporters CAT4  
312 and CAT2 respectively, in the NiH species. These results confirm previous observations made in  
313 the Asteraceae species *S. coronatus* (Meier *et al.*, 2018) and further suggest that the high expression  
314 of these genes corresponds to convergent mechanisms implicated in metal hyperaccumulation.  
315 Recently, the cationic amino acid transporter CAT4, primarily localizing on the vacuolar  
316 membrane, was associated with the histidine level trait in *A. thaliana* (Angelovici *et al.*, 2017).  
317 Indeed, previous studies have highlighted the role of histidine as an important metal ligand  
318 involved in nickel hyperaccumulation in Brassicaceae species (Krämer *et al.*, 1996; Kozhevnikova  
319 *et al.*, 2014). Interestingly, our analysis also indicated that a high expression of Histidinol  
320 dehydrogenase (OG4147), the last step in histidine biosynthesis, is also associated with the nickel  
321 hyperaccumulation trait (**Table 3, Figure 3**). These results support a role for histidine synthesis  
322 and transport in nickel hyperaccumulation. The high expression of MATE and ABC transporters  
323 related to DTX27 (OG 6266) and ALS3 (OG 411), potentially involved in metal transport and  
324 tolerance (Larsen *et al.*, 2004; Liu *et al.*, 2009), was not previously associated with nickel  
325 hyperaccumulation. The validation of the contribution of these transporters to nickel  
326 hyperaccumulation will require further support.



327 On the contrary, our results suggest a lower expression of calcium-transporting ATPase (OG  
328 10076) linked to the nickel hyperaccumulation trait. The rationale behind this result is not clear but  
329 it could be an indirect consequence of the edaphic conditions on which nickel hyperaccumulators  
330 are evolving. Indeed, NiH species are naturally growing on ultramafic soils that are rich in metals,  
331 including nickel, and with a strong calcium/magnesium imbalance affecting the development of  
332 most plant species (*aka* serpentine syndrome; Konečná *et al.*, 2020). Therefore, the lower  
333 expression of calcium-transporting ATPase in NiH species might be linked to the adaptation of  
334 these species to the serpentine syndrome.

335

### 336 *Advantages and limitations of this methodology to compare distant species*

337 One motivation of this work was to valorize the quantitative information contained in numerous  
338 RNA-Seq datasets already available from various plant species to identify gene functions  
339 associated with nickel hyperaccumulation. Most of these RNA-Seq datasets were produced with  
340 Illumina HiSeq paired-end technology. While the high number and quality of reads generated with  
341 this technology allow a good estimation of gene expression, their assembly frequently generates  
342 truncated ORF, thus affecting the annotation of OGs. To circumvent this limitation, we generated  
343 OG seeds using proteomes predicted from sequenced plant genomes. In a near future, the  
344 development of long-read sequencing technologies will undoubtedly ease the assembly of full-  
345 length transcripts (Amarasinghe *et al.*, 2020). In this work, we produced the OG database using the  
346 Mario method (Pereira *et al.*, 2014). While this strategy allowed us to compare transcriptomes from  
347 distant species, the database is specific to this work and the results are thus difficult to compare  
348 with other studies. The use of more general OG databases for comparative and functional genomics

349 such as Bgee for animals (Bastian *et al.*, 2021) or PLAZA for plants (Van Bel *et al.*, 2022), would  
350 thus favor comparisons between studies.

351 One main advantage of our method is that gene expression is quantified in each species using its  
352 own transcriptome. Therefore, this method does not absolutely require the identification of pairs of  
353 closely related species or populations with contrasted traits to quantify and compare gene  
354 expression using a common reference sequence. The identification of such pairs of closely related  
355 species with contrasted traits may represent a limiting condition in some genera. For example, the  
356 nickel hyperaccumulator *Blepharidium guatemalense* represents a monotypic genus (Navarrete  
357 Gutiérrez *et al.*, 2021), and all species of the Cuban genus *Leucocroton* are able to hyperaccumulate  
358 nickel (Reeves *et al.*, 1996). In addition, the output of pairwise comparisons to identify genes linked  
359 to a specific trait such as metal hyperaccumulation, strongly depends on the particular pair of  
360 species chosen for the comparison (Halimaa *et al.*, 2014; Meier *et al.*, 2018; García de la Torre *et*  
361 *al.*, 2021). Our methodology based on the use of several species belonging to two contrasted  
362 phenotypic groups (*eg* NiH and NA) is less sensitive to the choice of species to identify conserved  
363 or convergent mechanisms linked in a complex trait.

364 The orthologous conjecture proposes that orthologous genes are functionally more similar than  
365 paralogous genes. While this conjecture is still disputed (Stambouliau *et al.*, 2020), the annotation  
366 of groups of orthologous genes or orthologous groups is a widely used strategy for comparative  
367 analysis over a wide diversity of species (Van Bel *et al.*, 2022). To compute the level of expression  
368 of OGs in each species, we decided to sum the counts of all contigs associated to the same OG as  
369 previously used (Lallemand *et al.*, 2019). This choice is consistent with the hypothesis that contigs  
370 belonging to the same OG more likely encode for proteins playing the same function. However,  
371 OG may also contain inparalog genes, resulting from gene duplication in some species lineages,  
372 that might have acquired specific function. This is the case for OG 7137 containing for example

373 *AtIREG1* and *AtIREG2* playing different roles in *A. thaliana* (see above). It is important to notice  
374 that this OG does not appear to be more expressed in NiH species from several families. This might  
375 be the consequence of the presence of inparalog genes in this OG with contrasted expression levels  
376 in species from some families, thus affecting the calculation of the OG expression level.  
377 Alternatively, the high expression of IREG in NiH species might not be a mechanism observed in  
378 NiH species from all plant families. A more detailed analysis of the expression of contigs  
379 composing this OG would be necessary to validate these hypotheses. However, the identification  
380 of OG 7137 suggested that the method using the sum of read counts does not prevent the  
381 identification of OG containing inparalog genes.

382 Altogether, our results suggest that this methodology based on the quantification of the expression  
383 of orthologous groups allows the identification of genes with expression shifts linked to nickel  
384 hyperaccumulation from distant plant species. Even though RNA-Seq sequencing technologies and  
385 comparative genomics resources are evolving rapidly, we believe that the methodology presented  
386 in this work could be used as a framework to identify genes linked to a specific complex trait in a  
387 wide diversity of plant species or other organisms.

388

### 389 **Data and code availability**

390 The *Microthlaspi perfoliatum* Transcriptome Shotgun Assembly project has been deposited at  
391 DDBJ/ENA/GenBank under the accession GITW00000000. The version described in this paper is  
392 the first version, GITW01000000. The scripts and processed datasets are available at  
393 <https://github.com/i2bc/plant-nickel-accumulation>

394

## 395 **Contributions**

396  
397 M.G., O.L., S.M. conceived the study, V.S.GdlT., S.J., M.M., C.D., C.S., C.M, V.B. Y.P. collected  
398 the data, M.G., C.D., P.B., S.M performed the analyses and interpreted the results, M.G., O.L.,  
399 S.M. wrote the manuscript. All authors have read and approved the final manuscript.

400

## 401 **Acknowledgement / Funding**

402 The authors would like to thank Claire Toffano-Nioche (I2BC) for valuable discussions and  
403 support, and Bruno Fogliani (IAC/UNC) for his expertise on New Caledonian species. This work  
404 was supported by the X-TreM grant from the CNRS MITI to SM and the special funding  
405 MODELCOG from I2BC to MG. RNA-Seq sequencing for *Microthlaspi* and *Phyllanthus* species  
406 was financed by the ANR grant ANR-13-ADAP-0004 to SM, VB and BF.

407

## 408 **Bibliography**

409 **Altenhoff AM, Studer RA, Robinson-Rechavi M, Dessimoz C.** 2012. Resolving the Ortholog  
410 Conjecture: Orthologs Tend to Be Weakly, but Significantly, More Similar in Function than  
411 Paralogs (JA Eisen, Ed.). PLoS Computational Biology **8**, e1002514.

412 **Amarasinghe SL, Su S, Dong X, Zappia L, Ritchie ME, Gouil Q.** 2020. Opportunities and  
413 challenges in long-read sequencing data analysis. Genome Biology **21**, 30.

414 **Angelovici R, Batushansky A, Deason N, Gonzalez-Jorge S, Gore MA, Fait A, DellaPenna**  
415 **D.** 2017. Network-Guided GWAS Improves Identification of Genes Affecting Free Amino Acids.

- 416 Plant Physiology **173**, 872–886.
- 417 **Barrabé L, Maggia L, Pillon Y, Rigault F, Mouly A, Davis AP, Buerki S.** 2014. New  
418 Caledonian lineages of Psychotria (Rubiaceae) reveal different evolutionary histories and the  
419 largest documented plant radiation for the archipelago. Molecular Phylogenetics and Evolution  
420 **71**, 15–35.
- 421 **Bastian FB, Roux J, Niknejad A, et al.** 2021. The Bgee suite: integrated curated expression  
422 atlas and comparative transcriptomics in animals. Nucleic Acids Research **49**, D831–D847.
- 423 **Bedford T, Hartl DL.** 2009. Optimization of gene expression by natural selection. Proceedings  
424 of the National Academy of Sciences **106**, 1133–1138.
- 425 **Van Bel M, Diels T, Vancaester E, Kreft L, Botzki A, Van de Peer Y, Coppens F,**  
426 **Vandepoele K.** 2018. PLAZA 4.0: an integrative resource for functional, evolutionary and  
427 comparative plant genomics. Nucleic Acids Research **46**, D1190–D1196.
- 428 **Van Bel M, Silvestri F, Weitz EM, Kreft L, Botzki A, Coppens F, Vandepoele K.** 2022.  
429 PLAZA 5.0: extending the scope and power of comparative and functional genomics in plants.  
430 Nucleic Acids Research **50**, D1468–D1474.
- 431 **Billesbølle CB, Azumaya CM, Kretsch RC, Powers AS, Gonen S, Schneider S, Arvedson T,**  
432 **Dror RO, Cheng Y, Manglik A.** 2020. Structure of hepcidin-bound ferroportin reveals iron  
433 homeostatic mechanisms. Nature **586**, 807–811.
- 434 **Cappa JJ, Pilon-Smits EAH.** 2014. Evolutionary aspects of elemental hyperaccumulation.  
435 Planta **239**, 267–275.
- 436 **Catalán A, Briscoe AD, Höhna S.** 2019. Drift and Directional Selection Are the Evolutionary  
437 Forces Driving Gene Expression Divergence in Eye and Brain Tissue of Heliconius Butterflies.

- 438 Genetics **213**, 581–594.
- 439 **Chen J, Swofford R, Johnson J, Cummings BB, Rogel N, Lindblad-Toh K, Haerty W, di**  
440 **Palma F, Regev A.** 2019. A quantitative framework for characterizing the evolutionary history  
441 of mammalian gene expression. *Genome Research* **29**, 53–63.
- 442 **Cooper N, Thomas GH, Venditti C, Meade A, Freckleton RP.** 2016. A cautionary note on the  
443 use of Ornstein Uhlenbeck models in macroevolutionary studies. *Biological Journal of the*  
444 *Linnean Society* **118**, 64–77.
- 445 **Eddy SR.** 1998. Profile hidden Markov models. *Bioinformatics* **14**, 755–763.
- 446 **Edgar RC.** 2004. MUSCLE: multiple sequence alignment with high accuracy and high  
447 throughput. *Nucleic Acids Research* **32**, 1792–1797.
- 448 **Emms DM, Kelly S.** 2015. OrthoFinder: solving fundamental biases in whole genome  
449 comparisons dramatically improves orthogroup inference accuracy. *Genome Biology* **16**, 157.
- 450 **Enomoto T, Yoshida J, Mizuno T, Watanabe T, Nishida S.** 2021. Differences in mineral  
451 accumulation and gene expression profiles between two metal hyperaccumulators, *Noccaea*  
452 *japonica* and *Noccaea caerulescens* ecotype Ganges, under excess nickel condition. *Plant*  
453 *Signaling & Behavior* **16**, 1945212.
- 454 **van der Ent A, Baker AJM, Reeves RD, Pollard AJ, Schat H.** 2013. Hyperaccumulators of  
455 metal and metalloid trace elements: Facts and fiction. *Plant and Soil* **362**, 319–334.
- 456 **Felsenstein J.** 1985. Phylogenies and the Comparative Method. *The American Naturalist* **125**, 1–  
457 15.
- 458 **García de la Torre VS, Majorel-Loulergue C, Rigai GJ, et al.** 2021. Wide cross-species  
459 RNA-Seq comparison reveals convergent molecular mechanisms involved in nickel

- 460 hyperaccumulation across dicotyledons. *New Phytologist* **229**, 994–1006.
- 461 **Halimaa P, Lin Y-F, Ahonen VH, et al.** 2014. Gene expression differences between *Noccaea*  
462 *caerulescens* ecotypes help to identify candidate genes for metal phytoremediation.  
463 *Environmental science & technology* **48**, 3344–53.
- 464 **Hammond JP, Bowen HC, White PJ, Mills V, Pyke KA, Baker AJM, Whiting SN, May ST,**  
465 **Broadley MR.** 2006. A comparison of the *Thlaspi caerulescens* and *Thlaspi arvense* shoot  
466 transcriptomes. *New Phytologist* **170**, 239–260.
- 467 **Hanikenne M, Talke IN, Haydon MJ, Lanz C, Nolte A, Motte P, Kroymann J, Weigel D,**  
468 **Krämer U.** 2008. Evolution of metal hyperaccumulation required cis-regulatory changes and  
469 triplication of HMA4. *Nature* **453**, 391–395.
- 470 **Hansen TF.** 1997. Stabilizing Selection and the Comparative Analysis of Adaptation. *Evolution*  
471 **51**, 1341.
- 472 **Harmon LJ.** 2019. *Phylogenetic Comparative Methods: Learning From Trees* (LJ Harmon, Ed.).  
473 Center for Open Science.
- 474 **Heyduk K, Ray JN, Ayyampalayam S, Moledina N, Borland A, Harding SA, Tsai C-J,**  
475 **Leebens-Mack J.** 2019. Shared expression of crassulacean acid metabolism (CAM) genes pre-  
476 dates the origin of CAM in the genus *Yucca* (J Cushman, Ed.). *Journal of Experimental Botany*  
477 **70**, 6597–6609.
- 478 **Ho LST, Ané C.** 2014. A Linear-Time Algorithm for Gaussian and Non-Gaussian Trait  
479 Evolution Models. *Systematic Biology* **63**, 397–408.
- 480 **Housworth EA, Martins EP, Lynch M.** 2004. The Phylogenetic Mixed Model. *The American*  
481 *Naturalist* **163**, 84–96.

- 482 **Huang C-H, Sun R, Hu Y, *et al.*** 2016. Resolution of Brassicaceae Phylogeny Using Nuclear  
483 Genes Uncovers Nested Radiations and Supports Convergent Morphological Evolution.  
484 *Molecular Biology and Evolution* **33**, 394–412.
- 485 **Igea J, Bogarín D, Papadopulos AST, Savolainen V.** 2015. A comparative analysis of island  
486 floras challenges taxonomy-based biogeographical models of speciation. *Evolution* **69**, 482–491.
- 487 **Jiao Y, Wickett NJ, Ayyampalayam S, *et al.*** 2011. Ancestral polyploidy in seed plants and  
488 angiosperms. *Nature* **473**, 97–100.
- 489 **Konečná V, Yant L, Kolář F.** 2020. The Evolutionary Genomics of Serpentine Adaptation.  
490 *Frontiers in Plant Science* **11**.
- 491 **Kozhevnikova AD, Seregin I V, Erlikh NT, Shevyreva TA, Andreev IM, Verweij R, Schat  
492 H.** 2014. Histidine-mediated xylem loading of zinc is a species-wide character in *Noccaea*  
493 *caerulescens*. *New Phytol* **203**, 508–519.
- 494 **Krämer U.** 2010. Metal Hyperaccumulation in Plants. *Annual Review of Plant Biology* **61**, 517–  
495 534.
- 496 **Krämer U, Cotter-Howells JD, Charnock JM, Baker AJM, Smith JAC.** 1996. Free histidine  
497 as a metal chelator in plants that accumulate nickel. *Nature* **379**, 635–638.
- 498 **Lallemand F, Martin-Magniette M, Gilard F, Gakière B, Launay-Avon A, Delannoy É,  
499 Selosse M.** 2019. In situ transcriptomic and metabolomic study of the loss of photosynthesis in  
500 the leaves of mixotrophic plants exploiting fungi. *The Plant Journal* **98**, 826–841.
- 501 **Larsen PB, Geisler MJB, Jones CA, Williams KM, Cancel JD.** 2004. ALS3 encodes a  
502 phloem-localized ABC transporter-like protein that is required for aluminum tolerance in  
503 *Arabidopsis*. *The Plant Journal* **41**, 353–363.



- 504 **Law CW, Chen Y, Shi W, Smyth GK.** 2014. voom: precision weights unlock linear model  
505 analysis tools for RNA-seq read counts. *Genome Biology* **15**, R29.
- 506 **Leebens-Mack JH, Barker MS, Carpenter EJ, et al.** 2019. One thousand plant transcriptomes  
507 and the phylogenomics of green plants. *Nature* **574**, 679–685.
- 508 **Lemoine F, Lespinet O, Labedan B.** 2007. Assessing the evolutionary rate of positional  
509 orthologous genes in prokaryotes using synteny data. *BMC Evolutionary Biology* **7**, 237.
- 510 **Liu J, Magalhaes J V, Shaff J, Kochian L V.** 2009. Aluminum-activated citrate and malate  
511 transporters from the MATE and ALMT families function independently to confer Arabidopsis  
512 aluminum tolerance. *The Plant Journal* **57**, 389–399.
- 513 **Lynch M.** 1991. Methods for the analysis of comparative data in evolutionary biology. *Evolution*  
514 **45**, 1065–1080.
- 515 **Magallón S, Gómez-Acevedo S, Sánchez-Reyes LL, Hernández-Hernández T.** 2015. A  
516 metacalibrated time-tree documents the early rise of flowering plant phylogenetic diversity. *New*  
517 *Phytologist* **207**, 437–453.
- 518 **Manara A, Fasani E, Furini A, DalCorso G.** 2020. Evolution of the metal hyperaccumulation  
519 and hypertolerance traits. *Plant, Cell & Environment* **43**, 2969–2986.
- 520 **Meier SK, Adams N, Wolf M, Balkwill K, Muasya AM, Gehring CA, Bishop JM, Ingle RA.**  
521 2018. Comparative RNA-seq analysis of nickel hyperaccumulating and non-accumulating  
522 populations of *Senecio coronatus* (Asteraceae). *The Plant Journal* **95**, 1023–1038.
- 523 **Merlot S, Hannibal L, Martins S, Martinelli L, Amir H, Lebrun M, Thomine S.** 2014. The  
524 metal transporter PgIREG1 from the hyperaccumulator *Psychotria gabriellae* is a candidate gene  
525 for nickel tolerance and accumulation. *J Exp Bot* **65**, 1551–1564.

- 526 **Morrissey J, Baxter IR, Lee J, Li L, Lahner B, Grotz N, Kaplan J, Salt DE, Guerinot M**  
527 **Lou.** 2009. The Ferroportin Metal Efflux Proteins Function in Iron and Cobalt Homeostasis in  
528 Arabidopsis. *The Plant Cell* **21**, 3326–3338.
- 529 **Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B.** 2008. Mapping and quantifying  
530 mammalian transcriptomes by RNA-Seq. *Nature Methods* **5**, 621–628.
- 531 **Navarrete Gutiérrez DM, Pollard AJ, van der Ent A, Cathelineau M, Pons M-N, Cuevas**  
532 **Sánchez JA, Echevarria G.** 2021. *Blepharidium guatemalense*, an obligate nickel  
533 hyperaccumulator plant from non-ultramafic soils in Mexico. *Chemoecology* **31**, 169–187.
- 534 **O’Brien KP.** 2004. Inparanoid: a comprehensive database of eukaryotic orthologs. *Nucleic*  
535 *Acids Research* **33**, D476–D480.
- 536 **Pereira C, Denise A, Lespinet O.** 2014. A meta-approach for improving the prediction and the  
537 functional annotation of ortholog groups. *BMC Genomics* **15**, S16.
- 538 **Pillon Y, Hopkins HCF, Rigault F, Jaffré T, Stacy EA.** 2014. Cryptic adaptive radiation in  
539 tropical forest trees in New Caledonia. *New Phytologist* **202**, 521–530.
- 540 **Radhakrishnan G V, Keller J, Rich MK, et al.** 2020. An ancestral signalling pathway is  
541 conserved in intracellular symbioses-forming plant lineages. *Nature Plants* **6**, 280–289.
- 542 **Razafimandimbison SG, Kainulainen K, Wikström N, Bremer B.** 2017. Historical  
543 biogeography and phylogeny of the pantropical Psychotrieae alliance (Rubiaceae), with particular  
544 emphasis on the Western Indian Ocean Region. *American Journal of Botany* **104**, 1407–1423.
- 545 **Reeves RD, Baker AJM, Borhidi A, Berazaín R.** 1996. Nickel-accumulating plants from the  
546 ancient serpentine soils of Cuba. *New Phytologist* **133**, 217–224.
- 547 **Reeves RD, Baker AJM, Jaffré T, Erskine PD, Echevarria G, Ent A.** 2018. A global database

- 548 for plants that hyperaccumulate metal and metalloid trace elements. *New Phytologist* **218**, 407–  
549 411.
- 550 **Rich MK, Vigneron N, Libourel C, et al.** 2021. Lipid exchanges drove the evolution of  
551 mutualism during plant terrestrialization. *Science* **372**, 864–868.
- 552 **Robinson MD, Oshlack A.** 2010. A scaling normalization method for differential expression  
553 analysis of RNA-seq data. *Genome Biology* **11**.
- 554 **Rohlf R V, Nielsen R.** 2015. Phylogenetic ANOVA: The Expression Variance and Evolution  
555 Model for Quantitative Trait Evolution. *Systematic Biology* **64**, 695–708.
- 556 **Roux J, Rosikiewicz M, Robinson-Rechavi M.** 2015. What to compare and how: Comparative  
557 transcriptomics for Evo-Devo (M Robinson-Rechavi, Ed.). *Journal of Experimental Zoology Part*  
558 *B: Molecular and Developmental Evolution* **324**, 372–382.
- 559 **Schaaf G, Honsbein A, Meda AR, Kirchner S, Wipf D, von Wirén N.** 2006. AtIREG2  
560 Encodes a Tonoplast Transport Protein Involved in Iron-dependent Nickel Detoxification in  
561 *Arabidopsis thaliana* Roots. *Journal of Biological Chemistry* **281**, 25532–25540.
- 562 **Silvestro D, Kostikova A, Litsios G, Pearman PB, Salamin N.** 2015. Measurement errors  
563 should always be incorporated in phylogenetic comparative analysis (T Münkemüller, Ed.).  
564 *Methods in Ecology and Evolution* **6**, 340–346.
- 565 **Stamboulian M, Guerrero RF, Hahn MW, Radivojac P.** 2020. The ortholog conjecture  
566 revisited: the value of orthologs and paralogs in function prediction. *Bioinformatics* **36**, i219–  
567 i226.
- 568 **Stern DB, Crandall KA.** 2018. The Evolution of Gene Expression Underlying Vision Loss in  
569 Cave Animals (C Wilke, Ed.). *Molecular Biology and Evolution* **35**, 2005–2014.

- 570 **Voelckel C, Gruenheit N, Lockhart P.** 2017. Evolutionary Transcriptomics and Proteomics:  
571 Insight into Plant Adaptation. *Trends in Plant Science* **22**, 462–471.
- 572 **Wagner GP, Kin K, Lynch VJ.** 2012. Measurement of mRNA abundance using RNA-seq data:  
573 RPKM measure is inconsistent among samples. *Theory in Biosciences* **131**, 281–285.
- 574 **Weber M, Tramczynska A, Clemens S.** 2006. Comparative transcriptome analysis of toxic  
575 metal responses in *Arabidopsis thaliana* and the Cd<sup>2+</sup>-hypertolerant facultative metallophyte  
576 *Arabidopsis halleri*. *Plant, Cell and Environment* **29**, 950–963.
- 577 **Wickett NJ, Mirarab S, Nguyen N, *et al.*** 2014. Phylotranscriptomic analysis of the origin and  
578 early diversification of land plants. *Proceedings of the National Academy of Sciences* **111**,  
579 E4859–E4868.
- 580 **Yang Y, Moore MJ, Brockington SF, *et al.*** 2015. Dissecting Molecular Evolution in the Highly  
581 Diverse Plant Clade Caryophyllales Using Transcriptome Sequencing. *Molecular Biology and*  
582 *Evolution* **32**, 2001–2014.
- 583
- 584

585 **Table 1.** Description of RNA-Seq samples

Species	Families	Location	NiH	Short name	Bioprojects	SRA samples
<i>Senecio coronatus</i> Agnes mine	Asteraceae	South Africa	Yes	ScorA	PRJNA312157	SRX1901460 SRX1901471
<i>Senecio coronatus</i> Kaapsehoop	Asteraceae	South Africa	Yes	ScorC	PRJNA312157	SRX1901463 SRX1901464 SRX1901465
<i>Senecio coronatus</i> Galaxy mine	Asteraceae	South Africa	No	ScorB	PRJNA312157	SRX1901479 SRX1901480 SRX1901481
<i>Senecio coronatus</i> Pullen Farm	Asteraceae	South Africa	No	ScorD	PRJNA312157	SRX1901469 SRX1901470 SRX1901472
<i>Noccaea caerulescens</i> Firmiensis	Brassicaceae	France	Yes	Ncfi	PRJNA474900	SRX4174673 SRX4174674
<i>Noccaea montana</i>	Brassicaceae	France	No	Nmon	PRJNA474900	SRX4174677 SRX4174678
<i>Microthlaspi</i> <i>perfoliatum</i>	Brassicaceae	France	No	Mper	PRJNA657163	SRX8947159 SRX8947160 SRX8947161
<i>Geissois pruinosa</i>	Cunoniaceae	New Caledonia	Yes	Gpru	PRJNA476928	SRX4261243 SRX4261244 SRX4261245
<i>Geissois racemosa</i>	Cunoniaceae	New Caledonia	No	Grac	PRJNA476928	SRX4261246 SRX4261247 SRX4261248
<i>Phyllanthus luciliae</i>	Phyllanthaceae	New Caledonia	Yes	Phlu	PRJNA645979	SRX8723728 SRX8723729 SRX8723730
<i>Phyllanthus conjugatus</i>	Phyllanthaceae	New Caledonia	No	Phco	PRJNA645979	SRX8723731 SRX8723732 SRX8723733
<i>Psychotria grandis</i>	Rubiaceae	Cuba	Yes	Pgra	PRJNA476927	SRX4261234 SRX4261235
<i>Psychotria costivenia</i>	Rubiaceae	Cuba	Yes	Pcos	PRJNA476927	SRX4261236 SRX4261237
<i>Psychotria revoluta</i>	Rubiaceae	Cuba	No	Prev	PRJNA476927	SRX4261238 SRX4261239
<i>Psychotria gabriellae</i>	Rubiaceae	New Caledonia	Yes	Pgab	PRJNA476924	SRX4261225 SRX4261226 SRX4261227
<i>Psychotria</i> <i>semperflorens</i>	Rubiaceae	New Caledonia	No	Psem	PRJNA476924	SRX4261228 SRX4261229 SRX4261230
<i>Homalium kanaliense</i>	Salicaceae	New Caledonia	Yes	Hkan	PRJNA476925	SRX4261218 SRX4261219 SRX4261220
<i>Homalium betulifolium</i>	Salicaceae	New Caledonia	No	Hbet	PRJNA476925	SRX4261221 SRX4261222 SRX4261223

586

**Table 2.** *De novo* assemblies of contigs and assignment to OGs

Species	Number of contigs	Contig median length	ORF median length	Number of classified contigs (%)
Scor	46726	475	330	27423 (59)
Ncfi	41843	620	390	23608 (56)
Nmon	64367*	436	345	44460 (69)
Mper	144815	420	312	89067 (62)
Gpru	41188*	486	288	20998 (51)
Grac	37663*	649	366	20410 (54)
Phlu	42028	482	303	22136 (53)
Phco	45524	437	288	22414 (49)
Pgra	31362*	728	456	19506 (62)
Pcos	37309*	542	321	19158 (51)
Prev	45787*	534	285	20254 (44)
Pgab	46860*	558	273	20493 (44)
Psem	41854*	600	285	19379 (46)
Hkan	42325*	541	300	22500 (53)
Hbet	40774*	537	306	21890 (54)

587

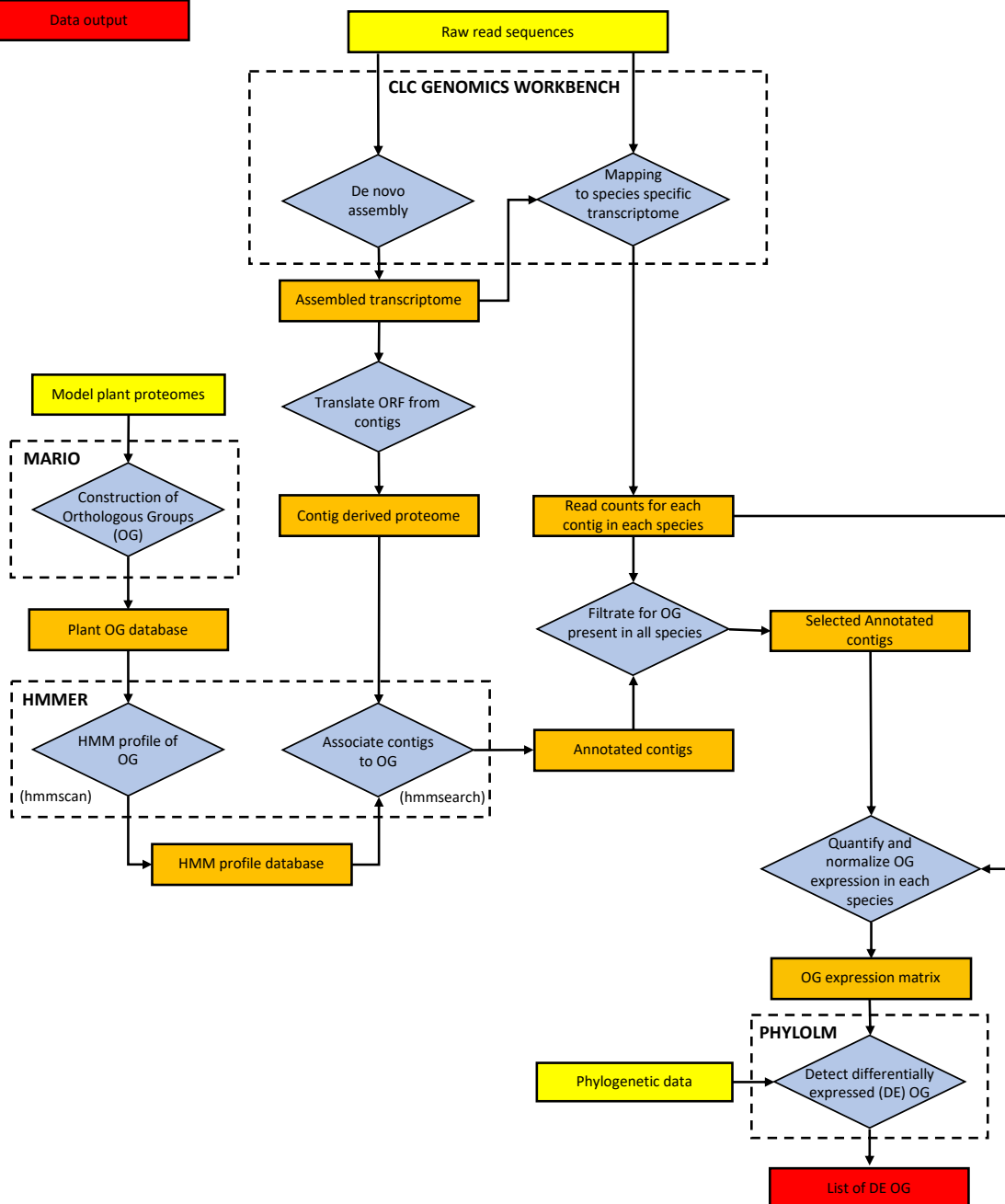
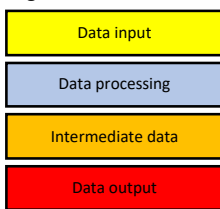
588 \*contigs with low expression (TPM<1) were filtered out.

589

**Table 3.** List of differentially expressed OG

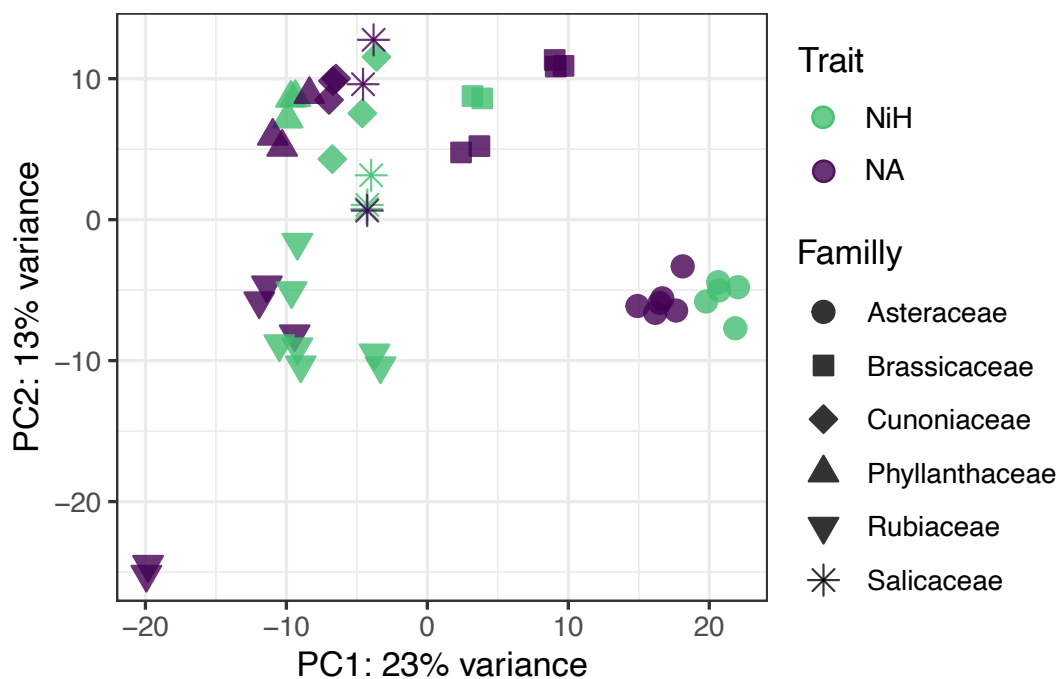
OG id	log2FC	p-value	OG predicted function* (EC number)	Best Hit*
6278	3.33	9.3e-20	Cationic amino acid transporter 4, vacuolar	CAAT4_ARATH
6266	3.01	2.2e-08	Protein Detoxification 27	DTX27_ARATH
4147	2.79	6.6e-18	Histidinol dehydrogenase, chloroplastic (1.1.1.23)	HISX_BRAOC
411	2.17	6.0e-06	Protein Aluminum Sensitive 3	ALS3_ARATH
12179	2.09	1.3e-04	Purple acid phosphatase 3 (3.1.3.2)	PPA3_ARATH
1871	2.00	9.7e-08	-	-
1781	1.97	1.5e-14	Ferredoxin-NADP reductase, chloroplastic (1.18.1.2)	FENR2_PEA
8462	1.92	1.7e-06	Ferric reduction oxidase 2 (1.16.1.7)	FRO2_ARATH
1441	1.80	8.8e-14	Vegetative incompatibility protein HET-E-1 (2.7.11.1)	HETE1_PODAS
7137	1.79	1.2e-04	Solute carrier family 40 member 2	S40A2_ARATH
10886	1.77	9.5e-07	-	-
11472	1.74	1.4e-04	Phenolic glucoside malonyltransferase 2 (2.3.1.-)	PMAT2_ARATH
2692	1.65	4.8e-10	-	-
1778	1.60	1.3e-05	Probable glutamate carboxypeptidase 2 (3.4.17.21)	GCP2_ARATH
3868	1.54	9.6e-08	Probable folate-biopterin transporter 7	FBT7_ARATH
3738	1.51	4.3e-05	Shikimate O-hydroxycinnamoyltransferase (2.3.1.133)	HST_TOBAC
559	1.51	1.5e-05	Polynucleotide 5'-hydroxyl-kinase NOL9	NOL9_ARATH
11324	-1.63	3.7e-05	Protein Defective in Meristem Silencing 3	DMS3_ARATH
4494	-1.76	6.2e-05	DNA ligase 4 (6.5.1.1)	DNLI4_ARATH
10076	-1.76	2.6e-04	Calcium-transporting ATPase 12, plasma membrane-type (3.6.3.8)	ACA12_ARATH
3681	-1.77	1.4e-08	Dicer-like protein 4 (3.1.26.3)	DCL4_ARATH
2967	-1.93	3.4e-08	Uncharacterized protein HI_0077	Y077_HAEIN
693	-1.94	3.2e-10	-	-
8515	-1.96	2.8e-04	Putative axial regulator YABBY 2	YAB2_ARATH
3961	-2.07	2.4e-06	Probable GTP-binding protein OBG, mitochondrial	OBGM_ARATH
3821	-2.09	1.8e-11	Probable protein phosphatase 2C 11 (3.1.3.16)	P2C11_ARATH
1916	-2.12	7.2e-05	Linoleate 13S-lipoxygenase 2-1, chloroplastic (1.13.11.12)	LOX2_ARATH
8093	-2.33	3.4e-05	Glycolipid transfer protein 2	GLTP2_ARATH
11720	-2.56	3.0e-05	Beta-fructofuranosidase, insoluble isoenzyme 1	INV1_DAUCA
5933	-3.17	4.1e-08	ACT domain-containing protein ACR4	ACR4_ARATH
13722	-3.19	1.1e-09	Rust resistance kinase Lr10 (2.7.11.1)	LRK10_WHEAT

\*From UniProtKB/Swiss-Prot database

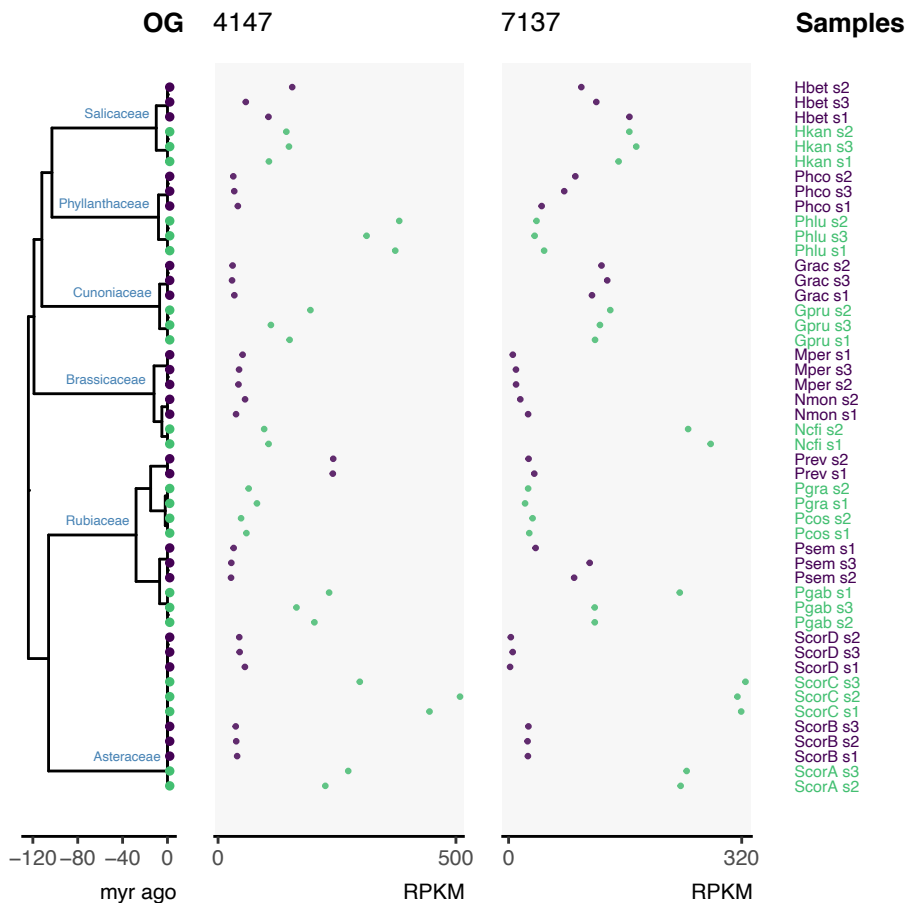
**Legend**

**Figure 1:** Workflow used to compare the expression of Orthologous groups in distant species (a) RNA-Seq datasets from different species were used to assemble transcriptomes *de novo* and to quantify gene (contig) expression in these transcriptomes using QIAGEN CLC Genomics Workbench. (b) Proteomes from sequenced plants species was used to generate a plant Orthologous group (OG) seed database using MARIO. (c) HMM profiles from the OG seed were produced by HMMER and used to annotate contigs from *de novo* transcriptomes. (d) The PHYLOLM package was used to detect differentially expressed OG using the normalized OG expression matrix.





**Figure 2.** Principal Component Analysis performed on the OG expression table after normalization and log transformation. The shapes of the symbols correspond to the family of the species and the color to the nickel hyperaccumulation (NiH, green) or non-accumulating (NA, purple) trait.



**Figure 3.** Representation of the expression of OG 4147 and OG 7137 along a plant phylogenetic tree. The expression of OG is expressed as Reads Per Kilobase of transcript, per Million mapped reads (RPKM). The phylogenetic tree is scaled in million years (myr), the six families considered are labeled on the tree. The name of the samples is given on the right. The color of the symbols correspond to NiH (green) or NA (purple) species or populations.