

Basis of genetic adaptation to heavy metal stress in the acidophilic green alga

Chlamydomonas acidophila

Fernando Puente-Sánchez^a, Silvia Díaz^b, Vanessa Penacho^c, Angeles Aguilera^d and Sanna Olsson^{e,f*}

Running title: Genetic adaptation of extremophilic *C. acidophila*

^aSystems Biology Program, Centro Nacional de Biotecnología (CNB-CSIC), Calle Darwin 3, 28049, Madrid, Spain

^bDept. Physiology, genetics and microbiology, Complutense University of Madrid (UCM), Calle José Antonio Novais 12, 28040 Madrid, Spain

^cBioarray, S.L. Parque Científico y Empresarial de la UMH. Edificio Quorum III. Avenida de la Universidad s/n, 03202 Elche, Alicante, Spain

^dCentro de Astrobiología (CSIC-INTA), Carretera de Ajalvir Km 4, 28850 Torrejón de Ardoz, Madrid, Spain

^eINIA Forest Research Centre (INIA-CIFOR), Dept. Forest Ecology and Genetics, Carretera de la Coruña km 7.5, 28040 Madrid, Spain

^fDept. Agricultural Sciences, P.O. Box 27, 00014 University of Helsinki, Finland

*corresponding author : Sanna Olsson, INIA Forest Research Centre (INIA-CIFOR), Dept. Forest Ecology and Genetics, Carretera de A Coruña km 7.5, E-28040 Madrid, Spain, tel. +34 913476773, fax: +34 913476767, e-mail: sanna.olsson@helsinki.fi

1 ABSTRACT

2 To better understand heavy metal tolerance in *Chlamydomonas acidophila*, an extremophilic green alga, we
3 assembled its transcriptome and measured transcriptomic expression before and after Cd exposure in this and
4 the neutrophilic model microalga *Chlamydomonas reinhardtii*. Genes possibly related to heavy metal tolerance
5 and detoxification were identified and analyzed as potential key innovations that enable this species to live in
6 an extremely acid habitat with high levels of heavy metals. In addition we provide a data set of single
7 orthologous genes from eight green algal species as a valuable resource for comparative studies including
8 eukaryotic extremophiles.

9

10 Our results based on differential gene expression, detection of unique genes and analyses of codon usage all
11 indicate that there are important genetic differences in *C. acidophila* compared to *C. reinhardtii*. Several efflux
12 family proteins were identified as candidate key genes for adaptation to acid environments. This study suggests
13 for the first time that exposure to cadmium strongly increases transposon expression in green algae, and that
14 oil biosynthesis genes are induced in *Chlamydomonas* under heavy metal stress. Finally, the comparison of the
15 transcriptomes of several acidophilic and non-acidophilic algae showed that the *Chlamydomonas* genus is
16 polyphyletic and that acidophilic algae have distinctive aminoacid usage patterns.

17

18 Key words: heavy metals, transcriptomics, green algae, Río Tinto, extremophiles, transposons

19

20 1 INTRODUCTION

21 Cd is a widespread environmental pollutant which is even at low concentrations extremely toxic to aquatic
22 microorganisms, in particular microalgae (Brayner et al., 2011; Wang et al., 2013). In spite of its harmfulness
23 there exist very few studies on transcriptomic alterations caused by increased levels of this or other heavy
24 metals in green algae, (Hutchins et al., 2010; Jamers et al., 2013; Zhang et al., 2014). Cd binds to organic

25 molecules by forming bonds with sulfur and nitrogen, thereby inactivating proteins causing a broad range of
26 adverse effects. It is easily absorbed and bio-accumulated by lower organisms and transferred to higher trophic
27 levels in food chain. It has been shown to inhibit growth (Okamoto et al., 1996), chlorophyll and chloroplast
28 synthesis (Lamai et al., 2005), cause disintegration of the cell wall as well as induce a large increase in
29 superoxide dismutase (SOD) activity, indicative of oxidative stress (Okamoto et al., 1996). Additionally, Cd
30 replaces zinc and selenium at the active sites of enzymes, competes with other ions in membrane transport,
31 and decreases RNA and DNA synthesis as well as photosynthetic pigments and proteins (Prasad et al., 1999;
32 Wang et al., 2013).

33

34 The extremophilic green alga *Chlamydomonas acidophila* grows in very acidic environments (pH 2.3-3.4). Metal
35 sequestration in vacuoles seems to be an important mechanism in cadmium tolerance and detoxification in *C.*
36 *acidophila* (Aguilera and Amils, 2005) but there is evidence that also unique genetic features in *C. acidophila*
37 contribute to its high heavy metal tolerance (Olsson *et al.*, 2015; Olsson et al., 2017). The strain analyzed in this
38 work was isolated from Río Tinto (SW Spain), one of the most extensive examples of natural extreme acidic
39 environments (Fernández-Remolar et al., 2003). The river has a low pH (ranging from 0.8 to 2.5) buffered by
40 ferric iron and with high concentrations of heavy metals (Aguilera et al., 2006). These extreme conditions are
41 produced by the metabolic activity of chemolithotrophic prokaryotes that are found in high numbers in its
42 waters (González-Toril et al., 2003). Despite these extreme environmental conditions, Río Tinto shows an
43 unexpected degree of eukaryotic diversity (Amaral-Zettler et al., 2011). Cd was chosen for this study due to its
44 toxicity and also because it is found in very high concentrations in Río Tinto, with local average amounts that
45 can reach ca. 40 mg/L (Aguilera et al., 2007).

46

47 Research on extremophilic organisms significantly contribute to our understanding of the evolution of heavy
48 metal tolerance in plants and algae. The results enable detection of novel genes potentially useful for

49 biotechnology and phytoremediation of contaminated water resources. In spite of this, there is very limited
50 genetic data available for *C. acidophila* while the genome of *C. reinhardtii* has been sequenced and annotated
51 (Merchant et al., 2007; Manichaikul et al., 2009). For *C. reinhardtii*, there also exist several physiological,
52 molecular, and genetic studies including experimental verification of the functionality of the predicted ORFs
53 (Ghamsari et al., 2011). To increase genomic resources in *C. acidophila* we assembled an improved
54 transcriptome for this non-model species. We compared it to the transcriptomes of the model microalga
55 *Chlamydomonas reinhardtii* from the same genus and other publicly available green algal transcriptomes. To
56 explain how *C. acidophila* is able to survive extreme environments we used transcriptomic sequencing and qRT-
57 PCR to detect transcriptional changes caused by high Cd concentrations in *C. reinhardtii* and *C. acidophila* and
58 identified possible adaptive key genes. The high level of genes with unknown function as well as lack of an
59 annotated genome assembly makes the identification of important genes involved in heavy metal detoxification
60 in *C. acidophila* challenging. In spite of these difficulties we provide new information on heavy metal tolerance
61 in this organism, extremophiles and green algae in general.

62

63

64 2 MATERIAL AND METHODS

65 2.1 Sample collection, cultivation and exposure to Cadmium

66 *Chlamydomonas acidophila* strain RT46 was collected from water samples taken in 2010 at the CEM sampling
67 station of Río Tinto (SW Spain) (Aguilera et al., 2006), and isolated to grow in the presence of antibiotics,
68 vancomycin 50 µg/mL, cefotaxime 100 µg/mL and chloramphenicol 15 µg/mL (Sigma), on agar plates made with
69 0.22 µm-filtered river water. Individual colonies were transferred into K medium (Keller et al., 1987), pH 2. A strain
70 of *Chlamydomonas reinhardtii* (CC-1374, SAG 77.81) was purchased from the Chlamydomonas Resource Center
71 (University of Minnesota) and grown in K medium, pH 7. The K medium was supplied with the same antibiotics
72 as the ones used for cell isolation (vancomycin 50 µg/mL, cefotaxime 100 µg/mL and chloramphenicol 15 µg/mL).

73 The algae were grown under ca. $70 \mu\text{E s}^{-1} \text{m}^{-2}$ irradiance provided by day-light fluorescent tubes, 16:8 h LD cycle
74 and 22 °C of temperature. The cultivations were refreshed every two weeks in corresponding growth media and
75 cells undergoing exponential growth were grown to be treated with metalloids solutions. To reach exponential
76 growth 5 ml of *Chlamydomonas* cultivate was transferred into an 1 L Erlenmeyer bottle with 500 ml medium. After 10
77 days of growth 15 ml of cultivate was transferred into three 2 L Erlenmeyer bottles with 980 ml medium in each.

78

79 For the transcriptomic sequencing a Cd solution ($\text{CdCl}_2 \times 2 \frac{1}{2} \text{H}_2\text{O}$) with a final concentration of 245 μM was used.
80 Earlier studies on *Chlamydomonas* showed a peak of gene expression between three and four hours in genes
81 involved in cadmium tolerance (Hanikenne et al., 2005; Olsson et al. 2017). Therefore time points for cell
82 collection were set before the treatment, at 3h and 6h after Cd exposure. The cells were collected in 50 ml Falcon
83 tubes, centrifuged for 10 min in 5000 rpm, the supernatant was discarded and the pellets frozen at -80 °C until
84 RNA extraction. For qRT-PCR cultures were treated with following solutions: 1 μM Cd solution ($\text{CdCl}_2 \times 2 \frac{1}{2} \text{H}_2\text{O}$),
85 1mM Cu ($\text{CuSO}_4 \times 5\text{H}_2\text{O}$), 10 mM Fe ($\text{FeSO}_4 \times 7\text{H}_2\text{O}$), 1mM As (III) (AsNaO_2) or 5 mM As(V) (Na_2HAsO_4) and cells
86 were collected at 1, 3 or 24 h after exposure.

87

88 **2.2 RNA extraction and sequencing**

89 Total RNA was extracted with TRI Reagent® Solution (Ambion, Life Technologies, CA, USA) following
90 manufacturer's protocol. RNA quality and quantity were estimated using an Agilent 2100 bioanalyzer (Agilent
91 Technologies). RNA library preparation and high-throughput sequencing were carried out in the NGS sequencing
92 Unit (Scientific Park Foundation, Madrid, Spain) using Illumina GAiix sequencing platform. One full lane of 75
93 basepair long reads for each sample was sequenced to provide sufficient coverage for a representative overview
94 of the expression profile. The generated transcriptome library was non-normalized to allow detection of
95 differences on the gene expression level between the different treatments and untreated cultures.

96

97 **2.3 Data preprocessing *de novo* hybrid assembly**

98 All raw transcriptomic reads were filtered and trimmed with PRINSEQ lite (version 0.18.3 (Schmieder and
99 Edwards, 2011) in order to remove duplicates and low quality reads (using default parameters except for the
100 following: -min qual mean 25 -derep 12, -ns max p 1 -derep 12 -lc method dust -lc threshold 7 -trim tail left 6 -
101 trim tail right 6 -trim ns left 2 -trim ns right 2 -trim qual left 25 -trim qual right 25).

102

103 The single-end Illumina reads from *C. acidophila* obtained in this study were combined with the 454 reads
104 obtained in Olsson et al., (2015). Paired-end Illumina reads were simulated from 454 reads by using the
105 *run_simulate_illuminaPE_from_454ds.sh* script included in the Trinity suite. The resulting reads were
106 subsequently normalized in silico with the *normalize_by_kmer_coverage.pl* included in Trinity. The paired-end
107 normalized reads coming from the 454 dataset were pooled together with the single-end Illumina reads obtained
108 in this study, and assembled with Trinity (release 2013_08_14) (Grabherr et al., 2011) using Jellyfish (Marcais and
109 Kingsford, 2011) for k-mer counting with a maximum memory of 40G, minimum contig length of 200, paired
110 fragment length of 350 and a maximum butterfly heap space of 20G. Contigs with a BLASTn identity of more than
111 90% to the *E. coli*, *C. reinhardtii* and human transcriptomes were discarded.

112

113 **2.4 Abundance estimation and transcriptome coverage analysis**

114 The RSEM software package (version 1.1.18.modified) (Li and Dewey, 2011) was used to estimate the gene and
115 isoform expression values. For *C. acidophila*, a reference transcriptome was generated from the Trinity assembly
116 by using the RSEM commands *extract-transcript-to-gene-map-from-trinity* and *rsem-prepare-reference* with
117 default parameters. For *C. reinhardtii*, the reference transcriptome v4.0 (Merchant et al., 2007) available from
118 Phytozome (<http://www.phytozome.net/>) was used as a reference for estimating transcript expression. Reads
119 from the six samples were aligned separately to the reference transcriptomes by using Bowtie (version 0.12.7)
120 (Langmead et al., 2009) and expression values for each sample were obtained with RSEM. The resulting

121 expression counts were normalized with the trimmed mean of M-values method, as implemented in the edgeR
122 package (version 2.15.0) (Robinson et al., 2010). The transcripts with a log₂ fold change higher than 6 and FPKM
123 (Fragments Per Kilobase Million) of more than 20 in at least one sample were selected for further analysis. For *C.*
124 *acidophila*, only the longest transcript per Trinity subcomponent was reported.

125

126 In order to assess the coverage of each sequence in our *C. acidophila* assembly, reads from the three *C. acidophila*
127 samples were pooled and aligned against the reference transcriptome. We used the *align-reads.pl* script included
128 in the Trinity package (release 2013_08_14), resorting to Bowtie (version 0.12.7) to perform the alignment. The
129 script also utilized Samtools (version 0.1.18) (Li et al., 2009) for SAM-format alignment manipulations. The output
130 file bowtie out.coordSorted.bam, which contains both properly mapped pairs and single unpaired fragment
131 reads, was used as input for *Qualimap* (version 0.6) (García-Alcalde et al., 2012) in order to estimate transcript
132 coverage.

133

134 **2.5 Taxonomic and functional annotation**

135 All transcripts were annotated via BLASTx searches (Altschul et al., 1997). For taxonomic annotation GenBank's
136 non-redundant protein database (nr) was used. For functional annotation two other major databases, Uniprot's
137 Swiss-Prot and TrEMBL protein databases were used in addition to the nr database to get more accurate
138 information on genetic functions. Taxonomic and functional information from the multiple databases for each
139 differentially expressed contig was collected into a table preferring the most accurate functional annotation
140 from Swiss-Prot when available using the methods and scripts modified from de Wit et al., (2012).

141

142 **2.6 Protein prediction and orthology search with OrthoFinder across green algal transcriptomes**

143 To identify orthologous gene groups among green algae, representative transcriptome files were downloaded
144 from Phytozome v.11 (<http://www.phytozome.net/>) for six available species: *Chlamydomonas reinhardtii*,

145 *Coccomyxa subellipsoidea*, *Micromonas pusilla*, *Micromonas* sp. RCC299, *Osterococcus lucimarinus* and *Volvox*
146 *carteri*. The *de novo* assemblies of *C. acidophila* and *D. acidophila* (Puente-Sánchez et al., 2016) were translated
147 to amino acids with TransDecoder (v. 3.0.0, The Broad Institute). Orthologous sequences from these eight
148 species were grouped with the clustering software OrthoFinder (Emms and Kelly, 2015). The resulting
149 alignments were filtered to contain only the longest isoform of *C. acidophila* and *D. acidophila* when several
150 isoforms of the same gene (belonging to same component and subcomponent in the *de novo* assembly built
151 with Trinity) were present in the same orthologous group. Orthologous groups related to heavy-metal tolerance
152 were subject to further analyses while orthologous groups representing putative single-copy nuclear genes (an
153 orthologous group with exactly one gene / species) present in all species were used to build a phylogeny.

154

155 **2.7 Genes present in *C. acidophila* but not in *C. reinhardtii***

156 To find an explanation for the different responses to heavy metals in extremophiles and neutrophiles, two
157 approaches to identify genes that are present in *C. acidophila* but not *C. reinhardtii* were employed. First,
158 screening for genes related to heavy metal tolerance and detoxification was done based on keywords in the
159 annotation of the *C. acidophila* transcriptome. Only transcripts that had other organisms than *C. reinhardtii* as
160 first BLAST match were included. To verify that the identified candidate genes are not present in *C. reinhardtii*, a
161 local BLASTn search against *C. reinhardtii* transcripts was performed. Reciprocal BLAST was performed to
162 confirm the matches and confirmed isoforms were used in downstream analyses.

163

164 Secondly, to identify genes specific to acidophiles, orthologous groups containing both of the extremophiles (*C.*
165 *acidophila* and *D. acidophila*) but not *C. reinhardtii* were extracted. As a precaution to exclude contaminant
166 sequences, in the absence of reference genomes for the extremophiles, only the orthologous groups containing
167 at least one additional green algal species were kept. In addition transcripts with organellar annotations
168 (mitochondrial or chloroplast) were excluded. Phylogenetic analyses were made for the transcripts with an

169 annotation related to heavy metal tolerance and detoxification after confirming their absence in *C. reinhardtii*
170 by a BLASTx against nr database with a cut-off E-value of $\leq 10^{-3}$.

171

172 **2.8 Phylogenetic analyses**

173 Sequences were aligned with Mafft (Katoh et al., 2002). For individual genes the alignments were manually
174 edited in PhyDE® v1.0 (Müller et al., 2005) by excluding ends of the alignments which could not be confidently
175 aligned due to length differences and ambiguities in homology assessment. The concatenated data matrix of
176 488 single orthologous groups was trimmed with Trimal (Capella-Gutierrez et al., 2009) using the option -
177 gappyout. Bayesian analyses were performed with MrBayes v3.2.1 (Ronquist et al., 2012), applying the
178 suggested search strategies for amino acids (Huelsenbeck et al., 2001; Ronquist and Huelsenbeck, 2003). For
179 the individual genes four runs with four chains (1×10^6 iterations each) were run simultaneously while for the
180 concatenated matrix of 488 single orthologous groups four runs with two chains (1×10^6 iterations each) were
181 run. Chains were sampled every 1000 generations and the respective trees written to a tree file. Calculations of
182 the consensus tree and of the posterior probability of clades were performed based upon the trees sampled
183 after the chains converged. The concatenated matrix was also analyzed using RAxML (Stamatakis, 2006;
184 Stamatakis et al., 2008) defining the used model automatically with the option -m PROTGAMMAAUTO.
185 Consensus topologies and support values from the different methodological approaches were compiled and
186 drawn using TreeGraph2 (Stöver and Müller, 2010).

187

188 **2.9 Quantitative reverse transcription PCR (qRT-PCR)**

189 For qRT-PCR protocols established by Díaz et al., (2007) were followed, applying the modifications detailed in
190 Olsson et al., (2017). Actin (ACT1) and 18S were used as endogenous control genes. All qRT-PCR reactions were
191 carried out in an iQTM5 multicolor Real-Time PCR detection System (Bio-Rad) apparatus with the following
192 cycling conditions: (i) 5 min at 95°C to denature reverse transcriptase, (ii) 40 cycles of 95°C for 30 s, 55°C for 30

193 s and 72°C after 20 s. Both NTC (no template control) and RT minus control were negative. The real-time
194 dissociation curve was used to check primer specificity and to confirm the presence of a unique PCR product.
195 Standard curves were obtained using 10-fold serial cDNA dilutions and determining the Ct (cycle threshold)
196 values. The standard line parameters (amplification efficiency, slope and correlation coefficient) are reported in
197 Table 1. Analysis of relative gene expression was carried out according to the Standard-curve quantification
198 method (Larinov et al., 2005) from, at least, four independent experiments (each performed in duplicates).
199 Primers for qRT-PCR were designed using the program Primer3 ([http://frodo.wi.mit.edu/cgi-](http://frodo.wi.mit.edu/cgi-bin/primer3/primer3_www.cgi)
200 [bin/primer3/primer3_www.cgi](http://frodo.wi.mit.edu/cgi-bin/primer3/primer3_www.cgi)) with default settings. All primers used in this study are listed in Table 2.

201

202 **2.10 Codon usage bias and GC content analyses**

203 Complete CDSs (coding DNA sequences) were extracted from the eight algal transcriptomes (*Chlamydomonas*
204 *acidophila*, *Dunaliella acidophila*, *Chlamydomonas reinhardtii*, *Coccomyxa subellipsoidea*, *Micromonas pusilla*,
205 *Micromonas sp. RCC299*, *Osterococcus lucimarinus* and *Volvox carteri*) by using the Transdecoder software
206 included in the Trinity suite. GC content and codon and aminoacid usage for each CDS were calculated with
207 GCUA (General Codon Usage Analysis; McInerney, 1998). For each gene, only the longest transcript was
208 included in the analysis.

209

210 **2.11 GO-terms enriched/depleted in particular aminoacids in acidophiles versus non-acidophiles**

211 In order to detect Gene Ontology (GO) terms with a significant enrichment/depletion of particular aminoacids
212 in acidophiles versus non-acidophiles, the following procedure was followed. Firstly, we selected the proteins
213 with i) less than 2% of glutamate, ii) less than 2% of aspartate, iii) more than 4% of cysteine, iv) more than 15%
214 of serine. These proteins will henceforth be referred to as “extreme” proteins. The particular aminoacids and
215 the percentage cutoff values were selected after inspection of the aminoacid utilization profiles shown in the
216 figure obtained in the previous section and shown in Additional File 6. For each of the four aminoacids, we then

217 counted the number of appearances of each GO-term in the “extreme” proteomes and in the “non-extreme”
218 proteomes of the two acidophilic species and the six non-acidophilic species, respectively. This was information
219 used to build the following contingency table for each GO-term, which was subjected to the Fisher’s Exact test
220 in order to assess whether that particular GO-term was significantly enriched ($p < 0.05$) in the extreme fraction
221 of the proteome in acidophiles versus non-acidophiles, that is, was significantly enriched/depleted in that
222 particular aminoacid in acidophiles versus non-acidophiles.

223

# GO-term appearances in the extreme proteome of acidophiles	# GO-term appearances in the non-extreme proteome of acidophiles
# GO-term appearances in the extreme proteome of non-acidophiles	# GO-term appearances in the non-extreme proteome of non-acidophiles

224

225

226 3 RESULTS AND DISCUSSION

227 **3.1 High-throughput sequencing, assembly and taxonomic annotation of *C. acidophila* transcripts**

228 Six single-end Illumina Hi-Seq libraries were sequenced in order to monitor the transcriptomic response of
229 *Chlamydomonas reinhardtii* and *Chlamydomonas acidophila* to cadmium stress right after cadmium exposure,
230 three hours after exposure and six hours after exposure. A total of 131,128,472 raw reads were generated, of
231 which 66,677,308 passed quality filtering, with the duplication level being consistent with that found in other
232 studies (Gómez-Álvarez et al., 2009).

233

234 In order to obtain a high quality draft transcriptome for *C. acidophila*, the reads obtained in this study were
235 pooled together and co-assembled with the reads obtained in Olsson *et al.*, (2015). This yielded 151449
236 transcripts of unique isoforms corresponding to 47411 unique Trinity subcomponents (which can be interpreted
237 as distinct genes), with a N50 of 3212 nucleotides, and average isoform coverage of 54.62X. The pre-processing

238 and assembly statistics are summarized in Tables 3A and 3B. The hybrid assembly significantly improved the
239 assembly results and genome fraction coverage over the existing assembly from the earlier study (GenBank
240 accession GBAH000000000) for which only 454 reads were used.

241

242 **3.2 Differential expression analysis of transcripts**

243 For both species, the gene expression after 3h and 6h of cadmium exposure was compared to the gene
244 expression right before cadmium exposure. H43 (Rubinelli et al., 2002) and *Cds1* (Hanikenne et al., 2005) are
245 among the few genes that have been identified to be induced by cadmium in *C. reinhardtii*. In addition a novel
246 phytochelatin synthase CaPCS2 was recently showed to be strongly induced by Cd in *C. acidophila* (Olsson et al.,
247 2017). Transcripts homologous to these genes were not found to be differentially expressed in this study,
248 possibly due to the strict cutoff values applied. The time and concentration of the exposure might also greatly
249 affect the transcriptomic response of green algae to Cd. Hanikenne et al., (2005) observed a peak of expression
250 in the half-size ABC transporter gene *Cds1* at 4 hours after 200-400 μM Cd exposure and argued that the
251 transcript levels of this gene were too low to be detected under the experimental conditions (2 h exposure to
252 25 μM cadmium) used earlier by Rubinelli et al. (2002). On the other hand, Olsson et al., (2017) reported a very
253 strong induction of the gene CaPCS2 in as low concentration as 1 μM . Furthermore, different isoforms might
254 result in different expression values.

255

256 In this study we focused on genes showing differential gene expression when exposed to very strong Cd
257 exposure. To complement the gene expression profile of selected candidate genes qRT-PCR was performed
258 using different concentrations and time points.

259

260 **3.3 Differentially expressed genes in *C. reinhardtii***

261 The low number of transcripts detected to be differentially expressed in *C. reinhardtii* (Additional File 1) is likely
262 due to the high amount of Cd used in the experiment, which was chosen to give a visible effect on the
263 transcriptomic expression in *C. acidophila*.

264

265 The transcripts with highest increase in expression after Cd exposure between control and one of the cadmium
266 treated samples include transcripts coding for an apoptosis-inducing factor, NSG6 protein, NifU-like protein 5,
267 and a vacuolar protein sorting-associated protein, in addition to transcripts with unknown function. Induction
268 of stress related genes, as well as genes operating in metal uptake and export as a response to cadmium has
269 been observed previously in *C. reinhardtii* (Jamers *et al.*, 2013) as well as in cyanobacteria (Houot *et al.*, 2007).
270 Other differentially expressed genes could not be directly linked to heavy metal detoxification but can
271 nonetheless be related to stress responses. For example, NSG6 is involved in gametogenesis and induced under
272 nitrogen starvation (Abe *et al.*, 2004). Interestingly, gametogenesis in *C. reinhardtii* leads to an increased
273 production of lipids with use as biofuels (Miller *et al.*, 2010). Here we show that, apart from nitrogen starvation,
274 this process can also be induced by heavy metal stress, opening the way for novel engineering strategies in the
275 search for high oil yields

276

277 **3.4 Differentially expressed genes in *C. acidophila***

278 The top fifty up regulated higher transcripts are summarized in Additional File 2, and included several
279 transposable elements. Significant upregulation was observed in a transcript annotated as retrotransposon
280 copia (FPKM in 0h 0, 3h 240.83, 6h 840.66 in comp17295_c0_seq16), a transcript annotated as retrovirus-
281 related Pol polyprotein from transposon TNT 1-94 (the annotation varies according to isoform, highest FPKM in
282 comp17071_c1_seq27: 0h 1.84, 3h 156.11, 6h 676.71), retrovirus-related Pol polyprotein from transposon 297
283 (comp18064_c0_seq15) and Transposon Ty3-I Gag-Pol polyprotein (comp16440_c0_seq16). Retrotransposons
284 are assumed to be a major driving force for genome evolution through genome organization and gene

285 regulation in plants (Flavell *et al.* 1992), some being transferred horizontally (Cheng *et al.* 2009 and references
286 therein). There are indications that retrovirus and retrotransposons are involved in gene regulation and
287 detoxification of heavy metals. Retrovirus-related Pol polyprotein from transposon TNT 1-94 has been shown to
288 alter its methylation status in *Populus alba* when grown on heavy metal contaminated soil (Cicatelli *et al.*,
289 2014). Castrillo *et al.*, (2013) showed that heavy-metal stress induced transposon activity in plants. Exposure to
290 Cd strongly increased transposon expression in *C. acidophila*, which suggests for the first time that heavy-metal
291 stress induces transposon activity also in green algae.

292

293 There are several transcripts with the annotation arsenite resistance protein ArsB among the differentially
294 expressed genes (comp14907_c0 or comp15936_c0). The annotations are partly incongruent, comp14907_c0
295 getting annotated as arsenite resistance protein ArsB or ubiquitin-like modifier-activating enzyme ATG7, while
296 the annotations for comp15936_c0 are arsenite resistance protein ArsB or arsenate reductase. The automated
297 annotation is complicated by the fact that the nomenclature of the ACR3 family ArsB protein overlaps with ArsB
298 of *E.coli* belonging to the ArsB family (Wu *et al.*, 1992). It was verified from the alignments including all isoforms
299 (data not shown) that all isoforms of one component belong together, and the different annotations are due to
300 lack of highly similar sequences in GenBank of some sequence parts. To avoid confusions in this manuscript the
301 isoforms of comp14907_c0 are referred to arsenical-resistance protein ACR3 and isoforms of comp15936_c0 as
302 ACR3 family arsenite transporter based on the annotation of the consensus sequences of these isoforms. Most
303 differentially expressed transcripts of both comp15936_c0 and comp14907_c0 are strongly induced by Cd, with
304 the exception of comp15936_c0_seq52. However, according to qRT-PCR analyses the ACR3 family arsenite
305 transporter comp_15936_c0 is down-regulated by Cd (Table 4). The incongruent results between the measures
306 based on the gene expression data and the qRT-PCR are likely due to the differences in the used Cd
307 concentrations.

308

309 An oil globule associated protein (comp13235_c0_seq1) was detected to be induced by copper by Olsson *et al.*,
310 (2015), and is now shown to be also induced by Cd (FPKM in 0h 0, 3h 10.058, 6 h 30.682). This again highlights
311 the role of heavy metals as inductors of oil production in *Chlamydomonas* (see previous section). Furthermore,
312 the ability of *C. acidophila* to tolerate extreme acidity and heavy metal concentrations might help it avoid the
313 contamination issues that commonly hamper microalgal biodiesel production (Siaut *et al.*, 2011; Wang *et al.*,
314 2016). While a detailed study of the oil production potential of *C. acidophila* is beyond the scope of this
315 manuscript, our findings warrant further investigation on its biotechnological applications.

316

317 **3.5 Species phylogeny based on orthologous sequences**

318 We identified 488 single orthologues present in all eight species (Additional File 3), which were used to build a
319 species phylogeny. According to the phylogenomic analyses the genus *Chlamydomonas* is not monophyletic
320 (Fig. 1). This is not so surprising since *Chlamydomonas* is known to be polyphyletic and in need for revision, first
321 shown by Buchheim *et al.*, (1990) and confirmed by several later studies (e.g. Leliaert *et al.*, 2012; Nakada *et al.*,
322 2016). However, earlier phylogenies have been based on few molecular markers and now the polyphyly of
323 *Chlamydomonas* is shown for the first time on a phylogenomic level. *Micromonas pusilla* was resolved as best
324 root in the species tree. All clades got full support both with MrBayes and RAxML.

325

326 **3.6 Identification of genes unique to *C. acidophila***

327 Some genes can be important in heavy metal tolerance and metal homeostasis even if their expression is not
328 altered in the presence of the metal. Most phytochelatin synthases, for example, are known to be constitutively
329 expressed but post-translationally activated by heavy metals in plants (Cobbett and Goldsborough, 2002; Rea *et al.*,
330 2004). To better understand the mechanisms that enable *C. acidophila* to live in its extremely acid
331 environment we therefore identified genes involved in heavy metal tolerance and detoxification that are

332 present in *C. acidophila* but do not have an orthologue in *C. reinhardtii*, irrespective of their expression. Two
333 approaches were employed.

334

335 First we identified thirteen candidate genes based on annotations of the transcripts and verified by reciprocal
336 Blast searches as explained in material and methods (Table 5). Of these, in addition to the ACR3 family
337 members discussed above, transcripts with following annotations were up-regulated: several isoforms of
338 mitochondrial carrier protein MTM1 (comp10226_c0), which carries manganese for the mitochondrial
339 superoxide dismutase, and of the MATE efflux family protein DETOXIFICATION 44 (comp12911_c0). To test for
340 changes in gene expression caused by a low Cd exposure and metal specificity, qRT-PCR was performed on a
341 selection of these candidate genes unique to *C. acidophila* (Table 4). Cells were collected for qRT-PCR at 1, 3 or
342 24 h after exposure to Cd (1 μ M), Cu (1mM), Fe (10 mM), As (III) (1mM) or As(V) (5 mM). Due to degraded
343 cDNA Cd and Fe are represented by only two time points each.

344

345 Cd was noted to somewhat affect the expression of 18S and therefore the relative mRNA expression levels of
346 target genes were normalized against the levels of actin. Surprisingly, none of the tested transcripts were
347 detected to be significantly induced by any of the added metals but significant down-regulation can be
348 observed in most of them (Pair Wise Fixed Reallocation Randomisation test, $p < 0.01$). These incongruences
349 could be due to known methodological caveats in RNA-seq including different expression of the different
350 isoforms, gene duplications or artifacts in assembly and annotation (Conesa et al. 2016). They demonstrate the
351 importance of detailed functional studies of individual genes, although automated studies with massive input
352 can offer useful information about general trends and serve as first step for further studies.

353

354 Secondly, we employed OrthoFinder to cluster genes and detect those unique to acidophiles (*C. acidophila* and
355 *Dunaliella acidophila*). Eighteen genes present in both extremophilic species and at least one further algal

356 species were extracted from the resulting orthologous groups (Table 6). Three of them (phytochelatin synthase
357 CaPCS2, Arsenical-resistance protein ACR3 and multidrug efflux transporter AcrB) were detected both by the
358 first method based on key word search for metal tolerance from the annotations and the second method based
359 on filtering of orthologous groups.

360

361 Some of the key candidate genes highlighted in this study have been shown to enhance heavy metal tolerance
362 in *C. acidophila* or other organisms. The phytochelatin synthase CaPCS2 was shown to be strongly induced by
363 Cd in *C. acidophila* and cloning and expression of the gene in *Escherichia coli* clearly improved its cadmium
364 resistance (Olsson et al. 2017). Cobalamin has been shown to protect against oxidative stress in the acidophilic
365 iron-oxidizing bacterium *Leptospirillum* (Ferrer et al., 2016). Arsenical-resistance protein ACR3 is suggested to
366 be a key trait to its arsenic tolerance in the arsenic hyperaccumulator *Pteris vittata* (Indriolo et al., 2010) and it
367 might similarly enhance the tolerance to heavy metals in *C. acidophila*. Arsenite resistance efflux pump ArsB,
368 which pumps arsenite and antimonite, but not arsenate or cadmium, was first described in *E. coli* (Wu et al.,
369 1992. Our results suggest that these genes could be key traits for heavy-metal hypertolerance in *C. acidophila*. It
370 has been proposed in other extremophiles as well that just a few key genes would be responsible for their
371 hypertolerance to heavy-metals, for example Fer1 in the acidophilic archaeon *Ferroplasma acidarmanus* (Baker-
372 Austin et al., 2007).

373

374 **3.7 Phylogenetic distribution of candidate key genes involved in heavy-metal hyper-tolerance in *C. acidophila***

375 Most of the transcripts not present in *C. reinhardtii* with an annotation related to heavy metal tolerance are
376 most closely related to genes in other green algae or vascular plants (Fig. 2A and Additional File 4). However,
377 some transcripts get a first Blast hit in other algae, fungi, prokaryotes and amoebzoa. The phylogenetic
378 distribution patterns in these genes can be explained by ancient gene duplications, loss in some lineages or

379 horizontal gene transfer, and according to our results there are more than one explanation for the origin of
380 these genes.

381

382 The mitochondrial carrier MTM1 (comp_10226, Fig. 2B), DETOXIFICATION 44 protein (comp_12911, Additional
383 File 4) and arsenite-antimonite efflux family (comp_15332, Additional File 4) include both green algae and
384 Chromalveolata among the most closely related genes. The phytochelatin synthase CaPCS2 (comp 11852,
385 Additional File 4), which is located within a clade of prokaryotic genes, has been functionally characterized and
386 shown to likely originate from horizontal gene transfer from bacteria (Olsson *et al.*, 2017). The cobalamin
387 biosynthesis protein CobW contains two gene copies in *C. acidophila* (comp15241_c0_seq3 and
388 comp15241_c0_seq6), of which one is similar to *C. reinhardtii* but the other is more similar to bacterial
389 homologues (Fig. 2C).

390

391 Similarly, the genes not present in *C. reinhardtii* extracted with OrthoFinder have variable phylogenetic
392 distribution patterns (Fig. 3). Some of the genes are nested in a clade containing mainly bacteria (e.g.
393 dioxygenase) and could be horizontally transferred. But for others, like transmembrane protein
394 comp9629_c0_seq1 are closely related only to green algae and gene loss is a more likely explanation for their
395 presence in *C. acidophila* but absence in *C. reinhardtii*.

396

397 **3.8 Codon code and aminoacid usage analysis**

398 The transcripts belonging to each of the analyzed species (*Chlamydomonas acidophila*, *Chlamydomonas*
399 *reinhardtii*, *Coccomyxa subellipsoidea*, *Dunaliella acidophila*, *Micromonas pusilla*, *Micromonas sp. RC299*,
400 *Ostreococcus lucimarinus* and *Volvox carteri*) clearly clustered together with regards to their Relative
401 Synonymous Codon Usage (Fig. 4a), showing the presence of distinct codon usage biases, even within
402 phylogenetically close species. For the most part, those differences did not result in different aminoacid usage

403 (Fig. 4b). The majority of the transcripts clustered together regardless of their source organism, except for a
404 large set of transcripts from *C. reinhardtii* and the two *Micromonas* species, which clustered independently.
405 Both *C. acidophila* and *D. acidophila* showed similar utilization profiles for several aminoacids, particularly an
406 enrichment in serine and cysteine, and a depletion in glutamic and aspartic acids when compared to the non-
407 acidophilic species (Fig. 5a, Additional File 5). This depletion in Glu and Asp in acidophiles was also observed by
408 Goodarzi et al. (2008), but their study only included bacterial and archaeal genomes. To the best of our
409 knowledge, this is the first study which proposes that the same can also be true in eukaryotes. We further
410 calculated which GO-terms were significantly depleted in Glu and Asp, or enriched in Cys and Ser, in the
411 acidophilic species when compared to the non-acidophilic species (Fig. 5b, Additional File 6). In the four cases,
412 the significant GO-terms chiefly belonged to the binding, catalytic activity, and transporter activity base
413 categories. These four modifications (lower Glu, lower Asp, higher Cys and higher Ser contents) are likely
414 related to optimizations for acidic environments. For example, Glu and Asp are negatively charged in neutral
415 conditions, but become neutral at lower pHs, which cancels their ability to stabilize proteins via salt bridges
416 (Anderson et al., 1990). On the other hand, the higher content of cystein could contribute to metal
417 detoxification and provide extra stability via disulfide bonds.

418

419 While the differences in total amino acid usage between organisms in acidophilic environments are usually
420 caused mostly by a limited number of amino acids (Goodarzi et al. 2008), the observed differences in codon
421 usage might be due to stronger selection pressure for codon optimization in extreme environments. Natural
422 selection acting through external environmental factors can shape the genomic pattern of synonymous codon
423 usage in extremophilic prokaryotes (Lynn et al. 2002; Zeldovic 2007), and our study is the first to suggest this to
424 be true also in eukaryotes.

425

426 The GC contents of a large amount of transcripts in all codon positions are different (Fig. 6) and can't be
427 explained by, for example, GC-content differences in a few horizontally transferred genes. Differences in the
428 first codon position affect the amino acid usage while differences in the third position are likely due to
429 preferences for different synonymous codons. Our results confirm that the overall genome-wide GC content is
430 the most significant parameter in explaining codon bias differences between organisms, suggested by
431 Hershberg and Petrov (2008).

432

433 **3.9 Conclusions**

434 The results of this study, including the most complete published transcriptome of *C. acidophila* and a set of
435 identified orthologous genes between eight green algae, increase the genomic information available on green
436 algae and extremophilic eukaryotes, highlight the adaptations mechanisms used by algae to thrive in acidic
437 environments, and provide a valuable resource for comparative studies on green algae from different habitats.
438 Further work should focus on detailed analyses of individual genes and applied exploitation of the results,
439 including engineering heavy metal tolerance in green algae for environmental and economic interests.

440

441

442 **5 ACKNOWLEDGEMENTS**

443 This work was supported by the Spanish Ministry of Economy and Competitiveness (MINECO) [CGL2011-22540,
444 AYA2011-24803]; the European Research Council (ERC) Advanced Grant [250350]. F. Puente-Sánchez was
445 supported by the Spanish MINECO/FEDER [CTM2013-48292-C3-2-R]. We acknowledge the Data Intensive
446 Academic Grid (DIAG) computing infrastructure (funded by National Science Foundation [0959894]) as well as
447 CSC – Finnish IT Center for Science and the Finnish grid infrastructure (FGI) for the allocation of computational
448 resources. Kimmo Mattila is acknowledged for help with setting up the OrthoFinder analysis pipeline. None of
449 the co-authors declare a conflict of interest.

450

451

452 6 REFERENCES

453 Abe, J., Kubo, T., Takagi, Y., Saito, T., Miura, K., Fukuzawa, H., and Matsuda, Y. (2004) The transcriptional program
454 of synchronous gametogenesis in *Chlamydomonas reinhardtii*. *Curr Genet* 46: 304–315.

455 Aguilera, A., and Amils, R., 2005. Tolerance to cadmium in *Chlamydomonas sp.* (Chlorophyta) strains isolated
456 from an extreme acidic environment, the Tinto River (SW, Spain). *Aquatic Toxicol* 75: 316–329.

457 Aguilera, A., Manrubia, S.C., Gómez, F., Rodríguez, N., and Amils, R. (2006) Eukaryotic community distribution
458 and its relationship to water physicochemical parameters in an extreme acidic environment, Rio Tinto (SW,
459 Spain). *Appl Environ Microbiol* 72: 5325–5330.

460 Aguilera, A., Zettler, E., Gomez, F., Amaral-Zettler, L., Rodríguez, N., and Amils, R. (2007). Distribution and
461 seasonal variability in the benthic eukaryotic community of Río Tinto (SW, Spain), an acidic, high metal extreme
462 environment. *Syst Appl Microbiol* 30: 531–546.

463 Anderson, D E., Becktel, W.J., and Dahlquist, F.W. (1990). pH-induced denaturation of proteins: a single salt
464 bridge contributes 3-5 kcal/mol to the free energy of folding of T4 lysozyme. *Biochem* 29: 2403–2408.

465 Altschul, S., Madden, T., Schaffer, A., Zhang, J., Zhang, Z., Miller, W., and Lipman D. (1997) Gapped BLAST and
466 PSI-BLAST: a new generation of protein database search programs. *Nucl Acid Res* 25: 3389–3402.

467 Amaral-Zettler, L.A., Zettler, E.R., Theroux, S.M., Palacios, C., Aguilera, A., and Amils, R. (2011) Microbial
468 community structure across the tree of life in the extreme Río Tinto. *ISME J* 5: 42–50.

469 Baker-Austin, C., Dopson, M., Wexler, M., Sawers, G.R., Stemmler, A., Rosen, B.R., and Bond, P.L. (2007) Extreme
470 arsenic resistance by the acidophilic archaeon '*Ferroplasma acidarmanus*' Fer1. *Extremophiles* 11: 425–434.

471 Brayner, R., Dahoumane, S.A., Nguyen, J.N., Yéprémian, C., Djediat, C., Couté, A., and Fiévet, F. (2011)

472 Ecotoxicological studies of CdS nanoparticles on photosynthetic microorganisms. *J Nanosci Nanotechnol* 11:

473 1852–1858.

474 Buchheim, M.A., Turnel, M., Zimmer, E.A., and Chapman, R. (1990) Phylogeny of *Chlamydomonas* (Chlorophyta)
475 based on cladistic analysis of nuclear 18S rRNA sequence data. *J Phycol* 26: 689–699.

476 Capella-Gutiérrez, S., Silla-Martínez, J.M., and Gabaldón, T. (2009) trimAl: a tool for automated alignment
477 trimming in large-scale phylogenetic analyses. *Bioinformatics* 25: 1972–1973.

478 Castrillo, G., Sánchez-Bermejo, E., De Lorenzo, L., Crevillén, P., Fraile-Escaciano, A., T.C., M., *et al.* (2013) WRKY6
479 transcription factor restricts arsenate uptake and transposon activation in *Arabidopsis*. *Plant Cell* 25: 2944–
480 2957.

481 Cheng, X., Zhang, D., Cheng, Z., Keller, B., Ling, and H.-Q. (2009) A New Family of Ty1-copia-Like
482 retrotransposons originated in the tomato genome by a recent horizontal transfer event. *Genetics* 181: 1183–
483 1193.

484 Cicatelli, A., Todeschini, V., Lingua, G., Biondi, S., Torrigiani, P., and Castiglione, S. (2014) Epigenetic control of
485 heavy metal stress response in mycorrhizal versus non-mycorrhizal poplar plants. *Environ Sci Pollut Res Int* 21:
486 1723–37.

487 Cobbett, C.S., and Goldsbrough, P. (2002) Phytochelatins and metallothioneins: roles in heavy metal
488 detoxification and homeostasis. *Annu Rev Plant Physiol Plant Mol Biol* 53: 159–182.

489 Conesa, A., Madrigal, P., Tarazona, S., Gomez-Cabrero, D., Cervera, A., McPherson, A., Wojciech Szcześniak, M.,
490 Gaffney, D.J., Elo, L.L., Zhang, X., and Mortazavi, A. (2016) A survey of best practices for RNA-seq data analysis.
491 *Genome Biol* 17:13

492 De Wit, P., Pespeni, M.H., Ladner, J.T., Barshis, D.J., Seneca, F., Jaris, H., *et al.* (2012) The simple fool’s guide to
493 population genomics via RNA-seq: an introduction to high-throughput sequencing data analysis. *Mol Ecol*
494 *Resour* 12: 1058–1067.

495 Díaz, S., Amaro, F., Rico, D., Campos, V., Benítez, L., Martín-González, A., *et al.* (2007) *Tetrahymena*
496 *metallothioneins* fall into two discrete subfamilies. *PLoS ONE* 2: e291

497 Emms, D.M., and Kelly, S. (2015) OrthoFinder: solving fundamental biases in whole genome comparisons
498 dramatically improves orthologous group inference accuracy. *Genome Biol* 16: 157.

499 Fernández-Remolar, D.C., Rodríguez, N., Gómez, F., and Amils, R. (2003) Geological record of an acidic
500 environment driven by the iron hydrochemistry: the Tinto River system. *J Geophys Res* 108: 5080–5095.

501 Ferrer, A., Rivera, J., Zapata, C., Norambuena, J., Sandoval, Á., Chávez, R., Orellana, O., and Levicán, G. (2016)
502 Cobalamin protection against oxidative stress in the acidophilic iron-oxidizing bacterium *Leptospirillum* group II
503 CF-1. *Front Microbiol* 7: 748.

504 Flavell, A.J., Dunbar, E., Anderson, R., Pearce, S.R., Hartley, R., and Kumar, A. (1992) Ty1-copia group
505 retrotransposons are ubiquitous and heterogeneous in higher plants. *Nucl Acids Res* 20:3639–44.

506 García-Alcalde, F., Okonechnikov, K., Carbonell, J., Cruz, L.M., Götz, S., Tarazona, S., *et al.* (2012) Qualimap:
507 evaluating next generation sequencing alignment data. *Bioinformatics* 28: 2678–2679.

508 Ghamsari, L., Balaji, S., Shen, Y., Yang, X., Balcha, D., Fan, C., *et al.* (2011) Genome-wide functional annotation
509 and structural verification of metabolic ORFeome of *Chlamydomonas reinhardtii*. *BMC Genomics*, 12: S4.

510 Gómez-Álvarez V, Teal T.K., and Schmidt, T.M. (2009) Systematic artifacts in metagenomes from complex
511 microbial communities. *ISME J* 3:1314–1317.

512 González-Toril, E., Llobet-Brossa, E., Casamayor, E.O., Amann, R., and Amils, R. (2003) Microbial ecology of an
513 extreme acidic environment, the Tinto River. *Appl Environ Microbiol* 69: 4853–4865.

514 Goodarzi, H., Torabi, N., Najafabadi, H.S., and Archetti, M. (2008) Amino acid and codon usage profiles: adaptive
515 changes in the frequency of amino acids and codons. *Gene* 407: 30–41.

516 Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., *et al.* (2011) Full-length
517 transcriptome assembly from RNA-seq data without a reference genome. *Nat Biotechnol* 29:644–52.

518 Hanikenne, M., Motte, P., Wu, M.C.S., Wang, T., Loppes, R., and Matagne, R.F. (2005) A mitochondrial half-size
519 ABC transporter is involved in cadmium tolerance in *Chlamydomonas reinhardtii*. *Plant Cell Environ* 28: 863–
520 873.

521 Huelsenbeck, J.P., Ronquist, F., Nielsen, R., and Bollback, J.P. (2001) Bayesian inference of phylogeny and its
522 impact on evolutionary biology. *Science* 294: 2310–2314.

523 Hershberg, R., and Petrov, D.A. (2008) Selection on codon bias. *Annu Rev Genet* 42: 287–99.

524 Houot, L., Floutier, M., Marteyn, B., Michaut, M., Picciocchi, A., Legrain, P., *et al.* (2007) Cadmium triggers an
525 integrated reprogramming of the metabolism of *Synechocystis* PCC6803, under the control of the Slr1738
526 regulator. *BMC Genomics* 8: 350.

527 Hutchins, C.M., Simon, D.F., Zerges, W., Wilkinson, K.J. (2010) Transcriptomic signatures in *Chlamydomonas*
528 *reinhardtii* as Cd biomarkers in metal mixtures. *Aquat Toxicol* 100: 120–127.

529 Indriolo, E., Na, G., Ellis, D., Salt, D.E., and Banks, J.O. (2010) A Vacuolar arsenite transporter necessary for
530 arsenic tolerance in the arsenic hyperaccumulating fern *Pteris vittata* is missing in flowering plants. *The Plant*
531 *Cell* 6: 2045–2057.

532 Jamers, A., Blust, R., Coen, W.D., Griffin, J.L., and Jones O.A.H. (2013) An omics based assessment of cadmium
533 toxicity in the green alga *Chlamydomonas reinhardtii*. *Aquat Toxicol* 126: 355–364.

534 Katoh, K., Misawa, K., Kuma, K., and Miyata, T. (2002) MAFFT: a novel method for rapid multiple sequence
535 alignment based on fast Fourier transform. *Nucleic Acid Res* 30: 3059–3066.

536 Keller, M.D., Selvin, R.C., Claus, W., and Guillard, R.R.L. (1987) Media for the culture of oceanic
537 ultraphytoplankton. *J Phycol* 23: 633–638.

538 Lamai, C., Kruatrachue, M., Pokethitiyook, P., Upatham, E.S., and Soonthornsarathool, V. (2005) Toxicity and
539 accumulation of lead and cadmium in the filamentous green alga *Cladophora fracta* Kützing: a laboratory study.
540 *Science Asia* 31: 121–127.

541 Langmead, B., Trapnell, C., Pop, M., and Salzberg, S.L. (2009) Ultrafast and memory-efficient alignment of short
542 DNA sequences to the human genome. *Genome Biol* 10: R25.

543 Larinov, A., Krause, A., Miller, W. (2005) Standard curve based method for relative real time PCR data
544 processing. *BMC Bioinformatics* 6: 62

545 Leliaert, F., Smith, D.R., Moreau, H., Herron, M.D., Verbruggen, H., Delwiche, C.F., and De Clerck, O. (2012)
546 Phylogeny and molecular evolution of the green algae. *Crit Rev Plant Sci* 31: 1–46.

547 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., *et al.*, (2009) The Sequence alignment/map
548 (SAM) format and SAMtools. *Bioinformatics* 25: 2078–2079.

549 Li, B., and Dewey, C.N. (2011) RSEM: Accurate transcript quantification from RNA-Seq data with or without a
550 reference genome. *BMC Bioinformatics* 12: 323.

551 Lynn, D.J., Singer, G.A.C., and Hickey, D.A. (2002) Synonymous codon usage is subject to selection in
552 thermophilic bacteria. *Nucl Acid Res* 30: 4272–4277.

553 Manichaikul, A., Ghamsari, L., Hom, E.F., Lin, C., Murray, R.R., Chang, R.L., *et al.* (2009) Metabolic network
554 analysis integrated with transcript verification for sequenced genomes. *Nat Methods* 6: 589–592.

555 Marcais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting of occurrences of
556 k-mers. *Bioinformatics* 27: 764–770.

557 McInerney, J.O. (1998) GCUA: general codon usage analysis. *Bioinformatics* 14: 372–373.

558 Merchant, S.S., Prochnik, S.E., Vallon, O., Harris, E.H., Karpowicz, S.J., Witman, G.B., *et al.* (2007) The
559 *Chlamydomonas* genome reveals the evolution of key animal and plant functions. *Science* 318: 245–250.

560 Miller, R., Wu, G., Deshpande, R. R., Vieler, A., Gärtner, K., Li, X., *et al.* (2010). Changes in transcript abundance
561 in *Chlamydomonas reinhardtii* following nitrogen deprivation predict diversion of metabolism. *Plant Physiol*
562 154: 1737–1752.

563 Müller, K., Quandt, D., Müller, J., and Neinhuis, C. (2005) PhyDE®: Phylogenetic Data Editor, version 0.995.
564 www.phyde.de

565 Nakada, T., Tomita, M., Wu, J.-T., and Nozaki, H. (2016) Taxonomic revision of *Chlamydomonas* subg.
566 *Amphichloris* (Volvocales, Chlorophyceae), with resurrection of the genus *Dangeardinia* and descriptions of
567 *Ixipapillifera* gen. nov. and *Rhysamphichloris* gen. nov. *J Phycol* 52: 283–304.

568 Okamoto, O.K., Asano, C.S., Aidar, E., and Colepicolo, P. (1996) Effects of cadmium on growth and superoxide
569 dismutase activity of the marine microalga *Tetraselmis gracilis* (Prasinophyceae). *J Phycol* 32: 74–79.

570 Olsson, S., Puente-Sánchez, F., Gómez-Rodríguez, and M., Aguilera, A. (2015) Transcriptional response to copper
571 excess and identification of genes involved in heavy metal tolerance in the extremophilic microalga
572 *Chlamydomonas acidophila*. *Extremophiles* 19: 657–672.

573 Olsson, S., Penacho, V., Puente-Sánchez, F., Díaz, S., and Aguilera, A. (2017) Horizontal gene transfer of
574 phytochelatin synthases from bacteria to extremophilic green algae. *Microbial Ecol* 73: 50–60.

575 Puente-Sánchez, F., Olsson, S., and Aguilera, A. (2016). Comparative transcriptomic analysis of the response of
576 *Dunaliella acidophila* (Chlorophyta) to short-term cadmium and chronic natural metal-rich water exposures.
577 *Microbial Ecol* 72: 595–607.

578 Prasad, M.N.V., and Strzalka, K. (1999) Impact of heavy metals on photosynthesis. In: Heavy metal stress in
579 plants: from molecules to ecosystems. Prasad MNV, and Hagemeyer J, (eds). Berlin, Germany: Springer, pp.
580 117–128.

581 Rea, P.A., Vatamaniuk, O.K., and Rigden DJ. (2004) Weeds, Worms, and More. Papain’s Long-Lost Cousin,
582 Phytochelatin Synthase. *Plant Physiol* 136: 2463–2474.

583 Robinson, M.D., McCarthy, D.J., and Smyth, G.K. (2010) EdgeR: a Bioconductor package for differential
584 expression analysis of digital gene expression data. *Bioinformatics* 26: 139–140.

585 Ronquist, F., and Huelsenbeck, J.P. (2003) MRBAYES 3: Bayesian phylogenetic inference under mixed models.
586 *Bioinformatics* 19: 1572–1574.

587 Ronquist, F., Teslenko, M., van der Mark, P., Ayres, D., Darling, A., Höhna, S., *et al.* (2012) MrBayes 3.2: efficient
588 Bayesian phylogenetic inference and model choice across a large model space. *Syst Biol* 61: 539–542.

589 Rubinelli, P., Siripornadulsil, S., Gao-Rubinelli, F., and Sayre, R.T. (2002) Cadmium- and iron-stress-inducible gene
590 expression in the green alga *Chlamydomonas reinhardtii*: evidence for H43 protein function in iron assimilation.
591 *Planta* 215: 1–13.

592 Schmieder, R., and Edwards, R. (2011) Quality control and preprocessing of metagenomic datasets.
593 Bioinformatics 27: 863–864.

594 Siaux, M., Cui n , S., Cagnon, C., Fessler, B., Nguyen, M., Carrier, P., *et al.* (2011) Oil accumulation in the model
595 green alga *Chlamydomonas reinhardtii*: characterization, variability between common laboratory strains and
596 relationship with starch reserves. BMC Biotechnol 11: 7.

597 Stamatakis, A. (2006) RAxML-VI-HPC: Maximum likelihood-based phylogenetic analyses with thousands of taxa
598 and mixed models. Bioinformatics 22: 2688–2690.

599 Stamatakis, A., Hoover, P., and Rougemont, J. (2008) A rapid bootstrap algorithm for the RAxML Web servers.
600 Syst Biol 57: 758–771.

601 St ver BC, and M ller KF. (2010) TreeGraph 2: Combining and visualizing evidence from different phylogenetic
602 analyses. BMC Bioinformatics 11: 7.

603 Wang, L., Yang, F., Chen, H., Fan, Z., Zhou, Y., Lu, J., and Zheng, Y. (2016) Antimicrobial cocktails to control
604 bacterial and fungal contamination in *Chlamydomonas reinhardtii* cultures. Biotechniques 60: 145–149.

605 Wang, S., Zhang, D., and Pan, X. (2013) Effects of cadmium on the activities of photosystems of *Chlorella*
606 *pyrenoidosa* and the protective role of cyclic electron flow. Chemosphere 93: 230–237.

607 Wit Wu, J., Tisa, L.S., and Rosen, B.P. (1992) Membrane topology of the ArsB protein, the membrane subunit of
608 an anion-translocating ATPase. J Biol Chem 267: 12570–12576.

609 Zeldovich, K.B., Berezovsky, I.N., and Shakhnovich, E.I. (2007) Protein and DNA Sequence Determinants of
610 Thermophilic Adaptation. PLoS Comput Biol 3: e5.

611 Zhang, W., Tana, N., and Li, S.F. (2014) NMR-based metabolomics and LC-MS/MS quantification reveal metal-
612 specific tolerance and redox homeostasis in *Chlorella vulgaris*. Mol Biosyst 10: 149–160.

FIGURE CAPTIONS

Figure 1. Phylogenetic relationships based on 488 nuclear single orthologous genes clustered with OrthoFinder and present in all eight species (*Chlamydomonas acidophila*, *Dunaliella acidophila*, *Chlamydomonas reinhardtii*, *Coccomyxa subellipsoidea*, *Micromonas pusilla*, *Micromonas sp. RCC299*, *Osterococcus lucimarinus* and *Volvox carteri*). The trees represent the majority consensus of trees sampled after stationarity in the Bayesian analysis. Posterior probability values from the Bayesian inference are indicated above, the corresponding bootstrap values of the maximum likelihood analysis below the branches.

Figure 2. Simplified phylogenetic analyses of transcripts coding for genes with an annotation related to heavy metal tolerance and present in *C. acidophila* but not in *C. reinhardtii*. The phylograms represent the majority consensus of trees sampled after stationarity in the Bayesian analysis. PP values equal or greater than 0.50 are shown above branches. The scale bar indicates relative distance between different sequences based on mutation rate. A) peroxisome isogenesis protein comp_10128 B) mitochondrial carrier comp 10226 C) cobalamin biosynthesis protein CobW comp_15241.

Figure 3. Simplified phylogenetic analyses of transcripts coding for genes that are involved in heavy metal tolerance and are present in *C. acidophila* but not in *C. reinhardtii* extracted from the results from the search for orthologous genes with OrthoFinder. The phylograms represent the majority consensus of trees sampled after stationarity in the Bayesian analysis. PP values equal or greater than 0.50 are shown above branches. The scale bar indicates relative distance between different sequences based on mutation rate. A) 2-hydroxyacyl-CoA lyase comp18202_c0_seq7 B) Dioxygenase comp13804_c0_seq4 C) Transmembrane protein 230 comp9629_c0_seq1 D) Cocaine esterase comp3348_c0_seq1 E) SDR-family protein with acetoacetyl-CoA reductase activity comp_14433_c0_seq1.

Figure 4. a) Correspondence analysis showing the distribution of transcripts (points) according to their Relative Synonymous Codon Usage distribution and b) Correspondence analysis showing the distribution of transcripts (points) according to the aminoacid usage bias of their predicted ORFs. The percentage of inertia explained by each axis is indicated in the axis caption. Transcripts are coloured by their source genome.

Figure 5. a) Distribution of Glu, Asp, Cys and Ser contents in the predicted ORFs from the eight species analyzed in this study. The ORFs included in the green area (low Glu, low Asp, high Cys and high Ser) were subjected to a GO-term enrichment analysis between the two acidophilic species and the rest. b) Summary of the Molecular Function GO-terms found to be significantly enriched ($p < 0.05$) in acidophiles versus non-acidophiles, when focusing in the low Glu, low Asp, high Cys, and high Ser fractions of the proteomes. Full results are provided in Additional File 6.

Figure 6. Scatterplots showing the GC content on each transcript (points) in the three different codon positions. Transcripts are coloured by their source genome.

TABLES WITH CAPTIONS

Table 1. Quantitative real-time RT-PCR standard-curve parameters for selected transcripts present in *C. acidophila* but not in *C. reinhardtii* and the expression control (housekeeping) genes 18S rRNA and actin. S= slope, R²= correlation coefficient, E= amplification efficiency.

Gene	S	R2	E
ACR3	-2.858	0.99	2.24
Arsenite transporter	-2.885	0.97	2.22
AcrB	-3.04	0.99	2.13
Glutathione-regulated potassium-efflux family	-2.823	1	2.26
MATE efflux protein	-3.079	0.99	2.11
Arsenite-antimonite efflux family	-2.96	0.96	2.18

Table 2. Primers used for quantitative real-time RT-PCR used in this study. For each region, forward (F) and reverse (R) primers are indicated, as well as product size.

Gene	Primer name	5' Sequence 3'	F/R	Product size (bp)
Multidrug efflux transporter AcrB comp_16471	comp16471_AcrB-F	GTAGGCATTCCCTTGCTGTC	F	89
	comp16471_AcrB-R	CCAAGGACCAAACAAGCAT	R	
ACR3 comp_14907	comp14907_ACR3-F	ACTTTGGCTTCTGGGAGGT	F	106
	comp14907_ACR3-R	TTTCACCATAAGCCCAGACC	R	
Arsenite transporter comp_15936	comp15936_ArsB-F	AATGTTACGGCAAAGCGAAC	F	100
	comp15936_ArsB-R	CAGTCACTGGCGAGCTCATA	R	
MATE efflux protein comp_12911	comp12911_MATE-F	ACTTTGGGTTTCATGGCTTTG	F	98
	comp12911_MATE-R	CACTCCTGCCAGTCCTAACC	R	
Arsenite-antimonite efflux family comp_15332	comp15332_MATE-F	CTAACACTCCTGTGGCAGCA	F	125
	comp15332_MATE-R	CAGCCTGTAAAGCCCTTTTG	R	
Glutathione-regulated potassium-efflux family comp_16013	comp16013_K-efflux-F	CGCTAGAAATCCCAACCAG	F	87
	comp16013_K-efflux-R	GCATTCTTTGCACCTCCAT	R	

Table 3. Sequence statistics on A) Illumina sequencing and for comparison, statistics on 454 reads from Olsson et al., (2015) are also shown. B) transcriptome assembly with Trinity using a hybrid assembly strategy combining 454 reads with Illumina reads.

Illumina sequencing					
Library	Condition	Raw reads	Input bases (Gb)	Trimmed reads	Discarded sequences (including duplicates)
J1	Reinhardtii-0h	21243002	1.61	451530	9504620
J2	Reinhardtii-3h	21819884	1.66	459756	10645037
J3	Reinhardtii-6h	19474228	1.48	438115	9709139
J4	Acidophila-0h	23656624	1.80	477750	11356978
J5	Acidophila-3h	22928585	1.74	503621	11287054
J6	Acidophila-6h	22006149	1.67	466332	11948336
454 data from Olsson <i>et al.</i> , (2015)					
Input 454 reads		Simulated Illumina reads		Normalized simulated pairs	
1021062		458306001		7717263	

B)

Hybrid assembly	
Input pairs	7717263
Input SE reads	33998990
Bases in assembly (Mb)	293
Trinity genes	47411
Trinity isoforms	151449
Isoform median length	1398
Isoform mean length	1936.66
Range of isoform lengths	201-19360
Isoform N50	3212
Isoform mean coverage	54.62X
Isoform std coverage	107.15X
Isoforms after filtering	129188
Isoforms after filtering with nr BLAST matches	87676

Table 4. Results of gene expression analysis by qRT-PCR of selected genes present in *C. acidophila* but not in *C. reinhardtii*. The cells were collected at 1, 3 or 24 hours after exposure to 1 μ M Cd solution ($\text{CdCl}_2 \times 2 \frac{1}{2} \text{H}_2\text{O}$), 1mM Cu ($\text{CuSO}_4 \times 5\text{H}_2\text{O}$), 10 mM Fe ($\text{FeSO}_4 \times 7\text{H}_2\text{O}$), 1mM As (III) (AsNaO_2) or 5 mM As(V) (Na_2HAsO_4). The relative

mRNA expression levels of target genes were normalized against the levels of actin and 18S. The fold induction and SD for each target gene is shown. ACR3 comp14907 = arsenical-resistance protein ACR3, ACR3 comp15936 = ACR3 family arsenite transporter, MATE comp12911 = MATE efflux protein, Arsenite-antimonite comp15332 = Arsenite-antimonite efflux family, AcrB comp16471 = Multidrug efflux transporter AcrB, MATE comp12911 = MATE efflux protein. Nd = Not defined.

	ACR3 comp14907	ACR3 comp15936	MATE comp12911	Arsenite- antimonite comp15332	AcrB comp1647 1	MATE comp12911
Cd 1h	-12,98 ± 1,97	-513 ± 126,1	-675,6 ± 79,1	-2327 ± 574	-761 ± 112	-8,1 ± 2,5
Cd 24h	-6,36 ± 0,8	-14,3 ± 3,6	-20,62 ± 3,0	-2862 ± 599	-609 ± 113	-12,5 ± 3,3
As(III) 1h	0,52 ± 0,06	-1,19 ± 0,09	-1,04 ± 0,3	-10,2 ± 1,4	-13,5 ± 2	-0,7 ± 0,1
As(III) 3h	-293,79 ± 36,59	-1824 ± 138	-731,7 ± 135,4	-15158 ± 1135	-3142 ± 378	-4,2 ± 1,2
As(III) 24h	-11,7 ± 1,45	-3445 ± 261	-314,9 ± 46,1	-3531 ± 864	-476 ± 70	-18 ± 5,4
Cu 1h	4,15 ± 0,52	-3687 ± 1012	3,38 ± 1,1	-3,8 ± 1,2	-6639 ± 786	-55 ± 11,7
Cu 3h	4,65 ± 0,5	-6358 ± 2424	5,6 ± 0,4	-11,48 ± 3,3	-690 ± 81	-1,5 ± 0,4
Cu 24h	1,04 ± 0,12	-170,2 ± 35,7	-26,02 ± 3,9	-554,9 ± 89	-207 ± 30	-47 ± 12
Fe 3h	Nd	Nd	1,31 ± 0,2	-819,8 ± 94,7	Nd	Nd
Fe 24h	0,36 ± 0,05	1,18 ± 0,2	0,8 ± 0,1	-2,2 ± 0,3	Nd	Nd

Table 5. Transcripts coding for genes that are involved in heavy metal tolerance present in *C. acidophila* but not in *C. reinhardtii* based on transcript annotations.

Contig name	Putative function	BLAST top match organism	BLAST match accession	E-value
comp10128_c0_seq1	Peroxisome isogenesis	<i>Coccomyxa subellipsoidea</i>	XP_005647114	3.01E-39
comp10226_c0_seq5	Mitochondrial carrier	<i>Coccomyxa subellipsoidea</i>	XP_005652123	3.95E-27
comp11852_c0_seq1	Phytochelatin synthase	<i>Calothrix sp.</i>	YP_007140091	6.44E-30
comp12911_c0_seq40	Protein DETOXIFICATION 44	<i>Chlorella variabilis</i>	EFN56963	6.62E-22
comp13602_c0_seq11	NRAMP family protein	<i>Volvox carteri</i>	XP_002947173	5.89E-153
comp14042_c0_seq2	ABC-ATPase	<i>Coccomyxa subellipsoidea</i>	XP_005643834	1.62E-92
comp14907_c0_seq53	Arsenical-resistance protein ACR3	<i>Coccomyxa subellipsoidea</i>	XP_005649016	1.99E-29
comp15241_c0_seq3	Cobalamin biosynthesis CobW	<i>Chlamydomonas reinhardtii</i>	XP_001699037	7.26E-60
comp15241_c0_seq6	Cobalamin biosynthesis CobW	<i>Burkholderia vietnamiensis</i>	YP_001117931	1.28E-80
comp15332_c0_seq3	Arsenite-antimonite efflux family	<i>Guillardia theta</i>	EKX52062	3.20E-66
comp15936_c0_seq52	ACR3 family arsenite transporter ArsB	<i>Coccomyxa subellipsoidea</i>	XP_005649501	6.10E-54
comp16013_c0_seq1	Glutathione-regulated potassium-efflux system	<i>Volvox carteri</i>	XP_002953483	8.63E-37
comp16471_c0_seq3	Multidrug efflux transporter AcrB	<i>Zea mays</i>	AFW59203	4.15E-58
comp17557_c1_seq9	Multidrug resistance-associated protein	<i>Coccomyxa subellipsoidea</i>	XP_005651467	1.59E-144

Table 6. Transcripts coding for genes that are involved in heavy metal tolerance present in *C. acidophila* but not in *C. reinhardtii* filtered from orthologous groups defined with OrthoFinder. Orthologous groups with annotations related to heavy metal tolerance and detoxification are marked with *. In the case of groups including several transcripts, the Blast hit organism and accession refers to the first one.

Orthologous group	Contig name	Putative function	BLAST top match organism	BLAST match accession	E-value
OG0001276	comp13064_c0_seq2, comp15004_c0_seq1	Tripeptidyl-peptidase 1	<i>Polysphondylium pallidum</i>	EFA84081	1.38e-17
OG0001752	comp12567_c0_seq1	Amino acid permease 2	<i>Capsella rubella</i>	EOA20485	3.33e-60
OG0001782	comp15790_c0_seq1	Alpha-1,3-glucosyltransferase	<i>Coccomyxa subellipsoidea</i>	XP_005651392	9.52e-92
OG0003420	comp16380_c0_seq1	Metal-nicotianamine transporter	<i>Amborella trichopoda</i>	ERN09450	5.88e-32
OG0003495	comp18202_c0_seq7	2-hydroxyacyl-CoA lyase	<i>Galdieria sulphuraria</i>	XP_005708092	0.0
OG0004374	comp18062_c0_seq1	Abhydrolase domain-containing protein	<i>Dictyostelium purpureum</i>	XP_002957250	2.33e-05
OG0004475*	comp11852_c0_seq1	Phytochelatin synthase	<i>Calothrix sp.</i>	YP_007140091	6.44e-30
OG0005070	comp16077_c0_seq9	G-box-binding factor 1	<i>Brassica napus</i>	CAA58774	1.33e-10
OG0005487*	comp14907_c0_seq15	Arsenical-resistance protein ACR3	<i>Coccomyxa subellipsoidea</i>	XP_005649016	1.45e-85
OG0005928*	comp13804_c0_seq4	dioxygenase	<i>Volvox carteri</i>	XP_002957190	2.32e-53
OG0006489	comp13735_c0_seq3	Snurportin-1	<i>Physcomitrella patens</i>	XP_001763666	5.34e-49
OG0006590*	comp16471_c0_seq1	multidrug efflux transporter AcrB	<i>Arabidopsis thaliana</i>	OAP00250	5.35e-63
OG0007003	comp14473_c0_seq1	Ankyrin-1	<i>Aegilops tauschii</i>	EMT31987	3.47e-56
OG0007890*	comp9629_c0_seq1	Transmembrane protein 230	<i>Physcomitrella patens</i>	XP_001772694	2.89e-17
OG0008459*	comp3348_c0_seq1	Cocaine esterase	<i>Achromobacter xylosoxidans</i>	WP_006387564	1.35e-80
OG0009876*	comp14433_c0_seq1	SDR-family protein with acetoacetyl-CoA reductase activity	<i>Sphingobium japonicum</i>	YP_003545425	9.92e-41
OG0010052	comp17871_c0_seq9	Histidine kinase	<i>Synechocystis sp.</i>	WP_009631601	1.298e-39

SUPPLEMENTARY MATERIAL

Additional File 1. *C. reinhardtii* transcripts with a log₂ fold change higher than 6 and FPKM of more than 20 in at least one sample. Normalized expression of the transcripts in each condition (0h, 3h and 6h after exposure to cadmium), their putative annotation and sequence length are shown. Transcripts are ordered by FPKM values at 6h.

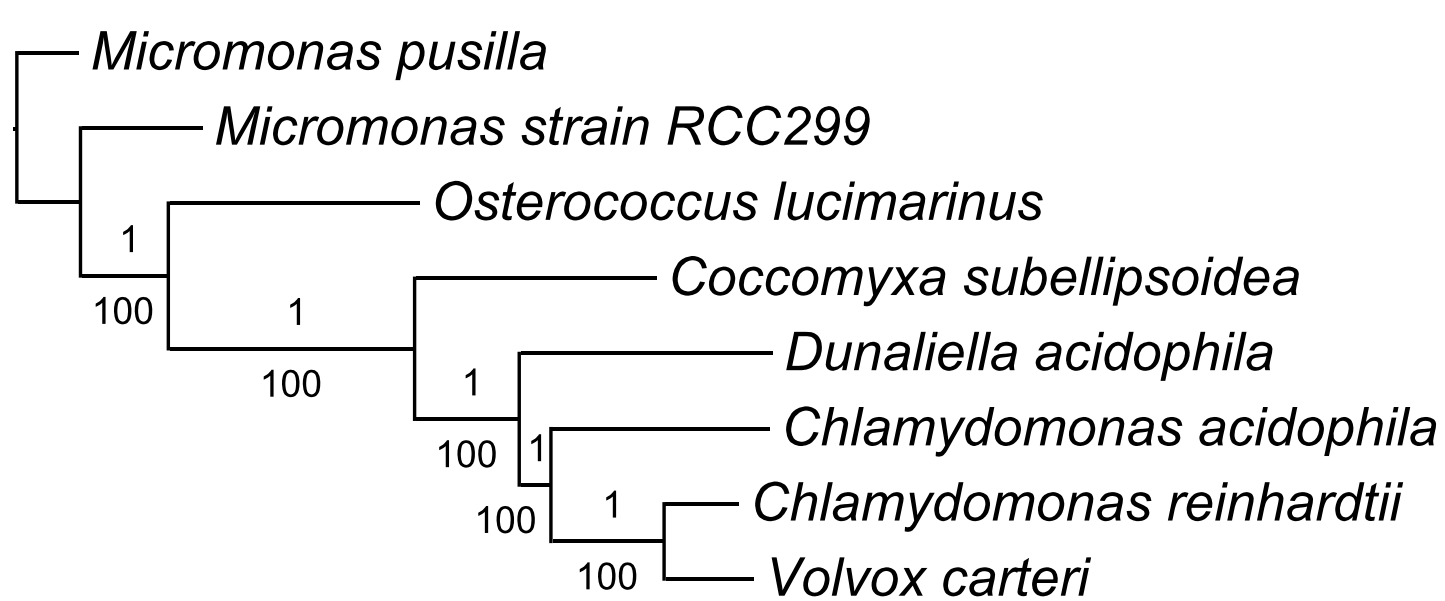
Additional File 2. *C. acidophila* RT46 transcripts with a log₂ fold change higher than 6 and FPKM of more than 20 in at least one sample. Normalized expression of the transcripts in each condition (0h, 3h and 6h after exposure to cadmium) and their putative annotation are shown. Transcripts are ordered by FPKM values at 6h.

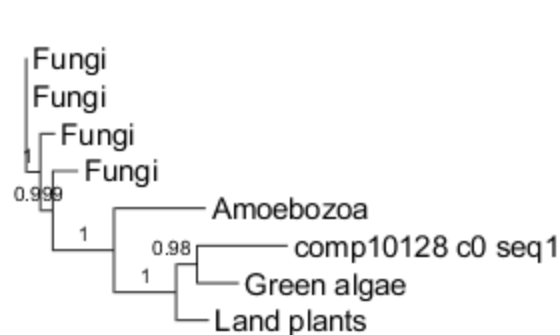
Additional File 3. Data matrix used for the phylogenomic analysis, containing 488 single orthologous groups of proteins clustered with OrthoFinder present in all species (*Chlamydomonas acidophila*, *Chlamydomonas reinhardtii*, *Coccomyxa subellipsoidea*, *Dunaliella acidophila*, *Micromonas pusilla*, *Micromonas sp. RC299*, *Ostrosoccus lucimarinus* and *Volvox carteri*).

Additional File 4. Simplified phylogenetic analyses of transcripts coding for genes with an annotation related to heavy metal tolerance and present in *C. acidophila* but not in *C. reinhardtii*. The trees represent the majority consensus of trees sampled after stationarity in the Bayesian analysis. PP values equal or greater than 0.50 are shown above branches. A) ACR3 comp_14907 and comp_15936 B) DETOXIFICATION 44 protein comp_12911 C) NRAMP family comp_13602 D) ABC-ATPase comp_14042 E) Arsenite-antimonite efflux family comp_15332 F) Glutathione-regulated potassium-efflux comp_16013 G) Multidrug efflux transporter AcrB comp_16471 H) Multidrug resistance-associated protein comp_17557.

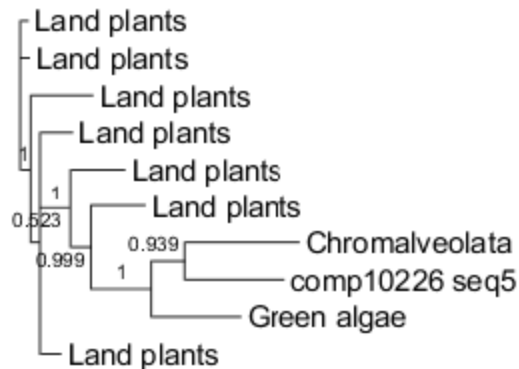
Additional File 5. Aminoacid utilization profiles for the different species. There is one cumulative frequency plot for each aminoacid.

Additional File 6. GO terms enriched in acidophiles versus non-acidophiles in the low-Glu, low-Asp, high-Cys, and high-Ser fractions of the proteome.

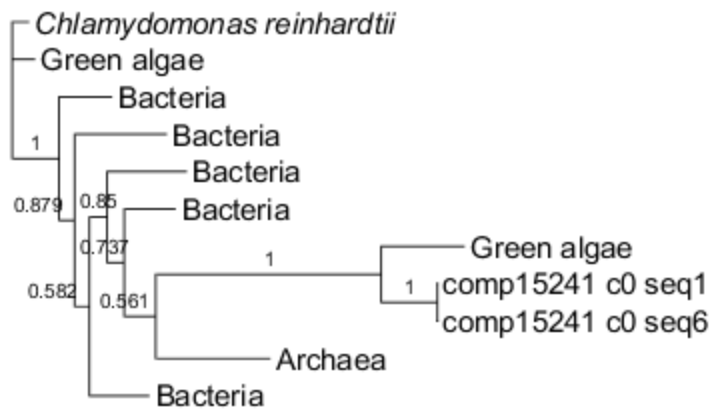




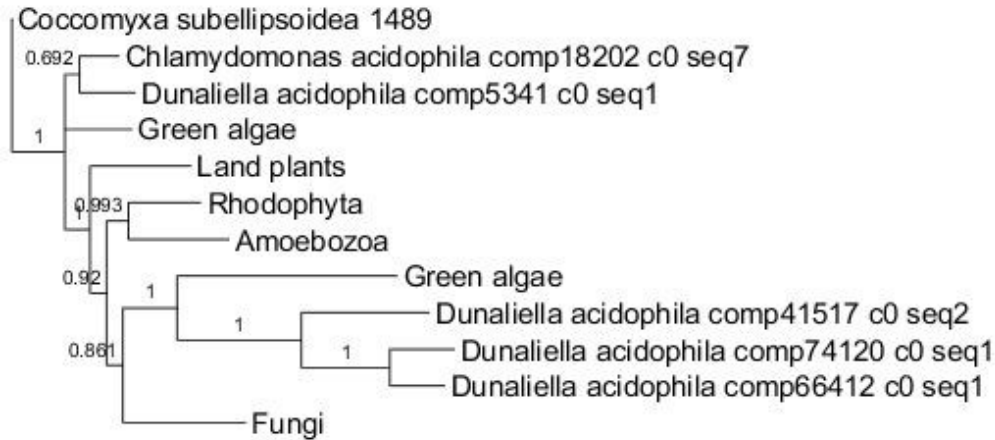
A



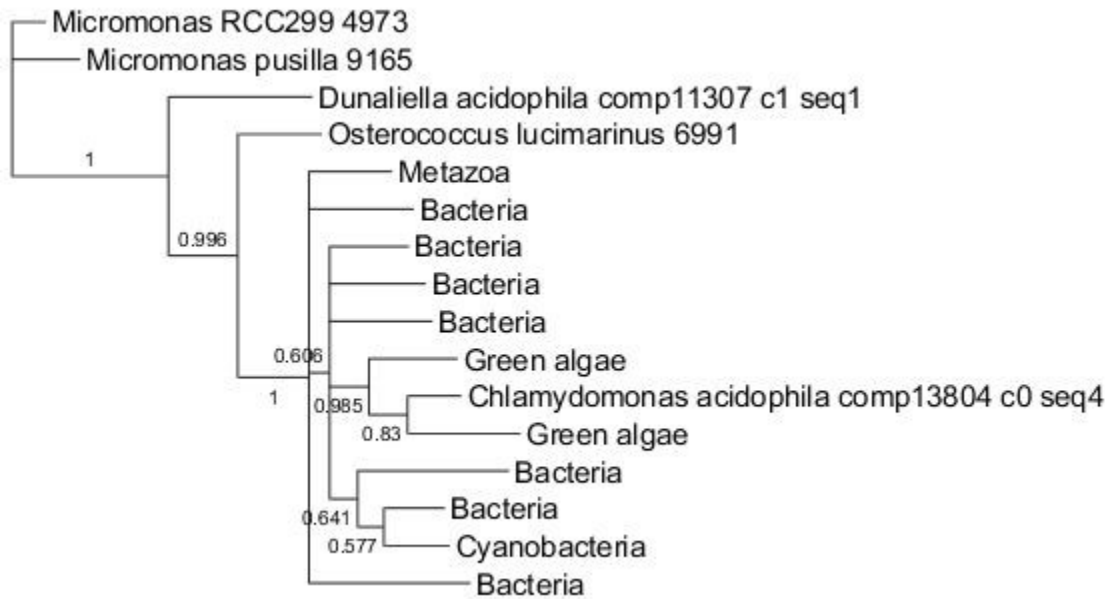
B



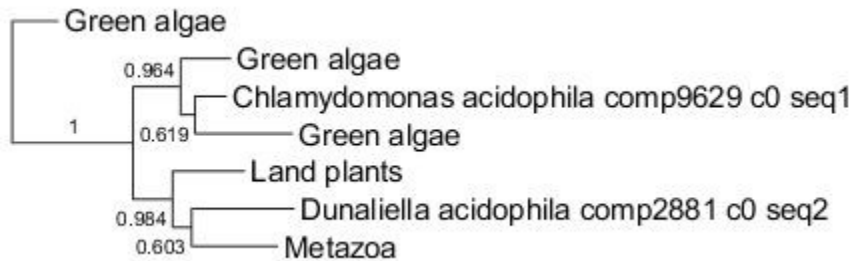
C



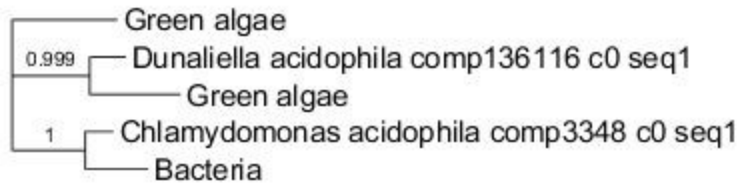
A) OG0003495: 2-hydroxyacyl-CoA lyase comp18202_c0_seq7



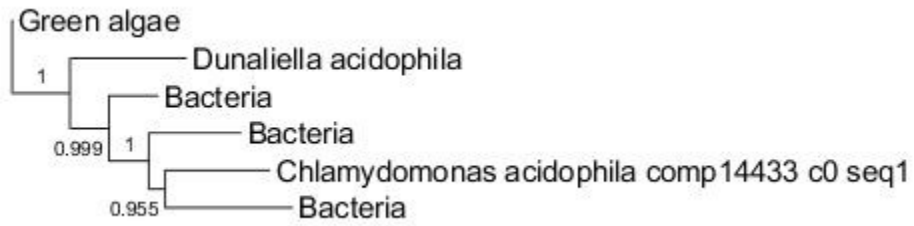
B) OG0005928: Dioxygenase comp13804_c0_seq4



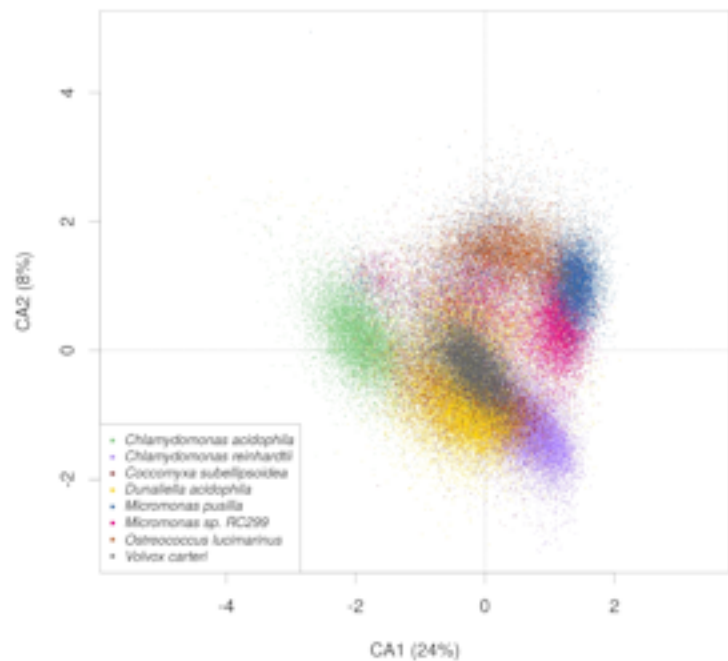
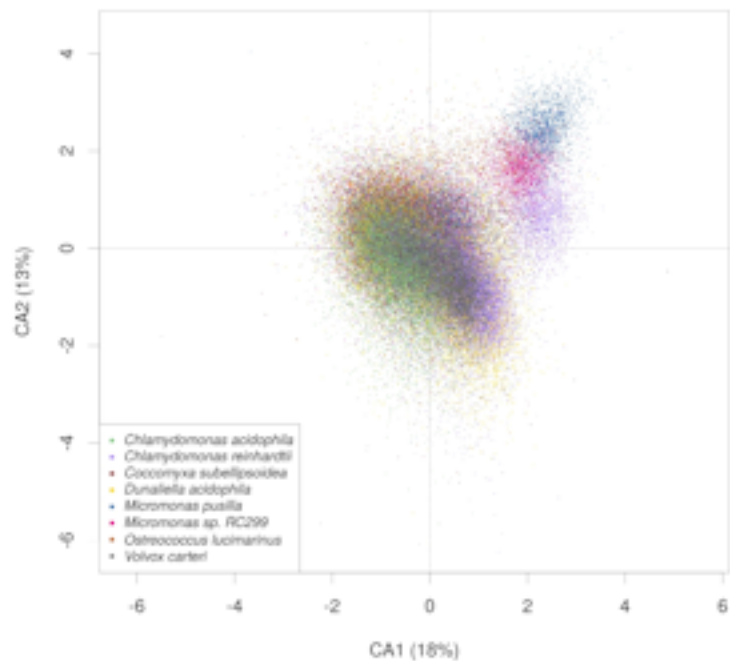
C) OG0007890: Transmembrane protein 230 comp9629_c0_seq1

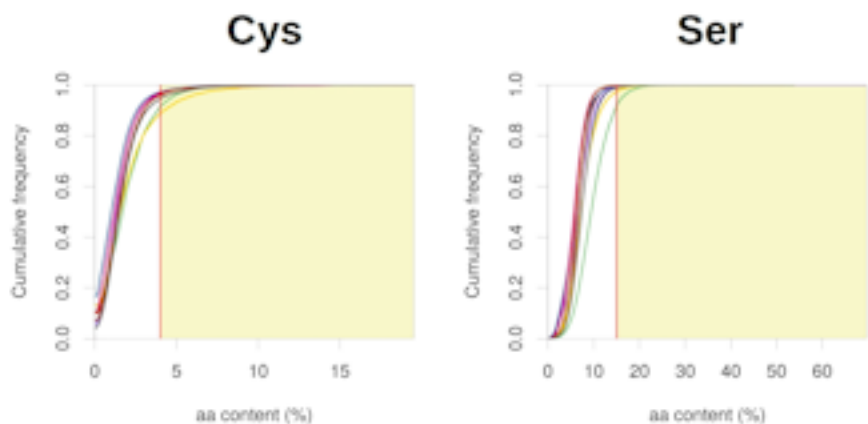
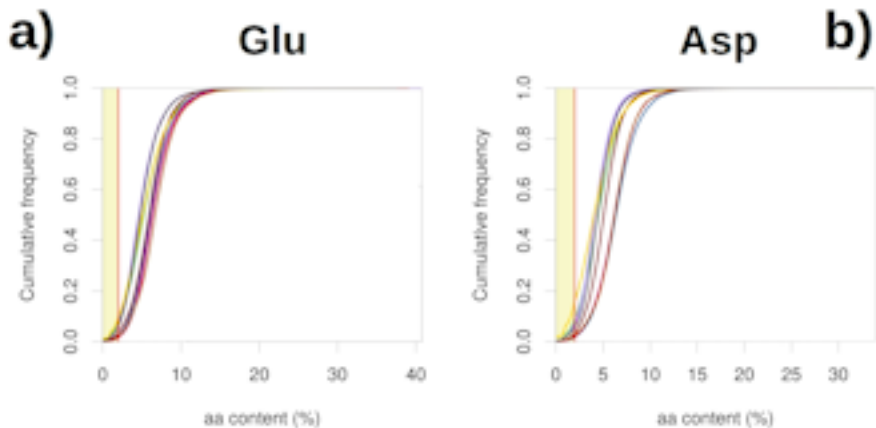


D) OG0008459: Cocaine esterase comp3348_c0_seq1



E) OG0009876: SDR-family protein with acetoacetyl-CoA reductase activity comp_14433_c0_seq1

a)**Codon Usage Bias****b)****Aminoacid Usage Bias**



- *Chlamydomonas acidophila*
- *Chlamydomonas reinhardtii*
- *Coccomyxa subellipsoidea*
- *Dunaliella acidophila*
- *Micromonas pusilla*
- *Micromonas sp. RC299*
- *Ostreococcus lucimarinus*
- *Volvox carteri*

b) Glu-depleted GO-terms in acidophiles Asp-depleted GO-terms in acidophiles



Cys-enriched GO-terms in acidophiles

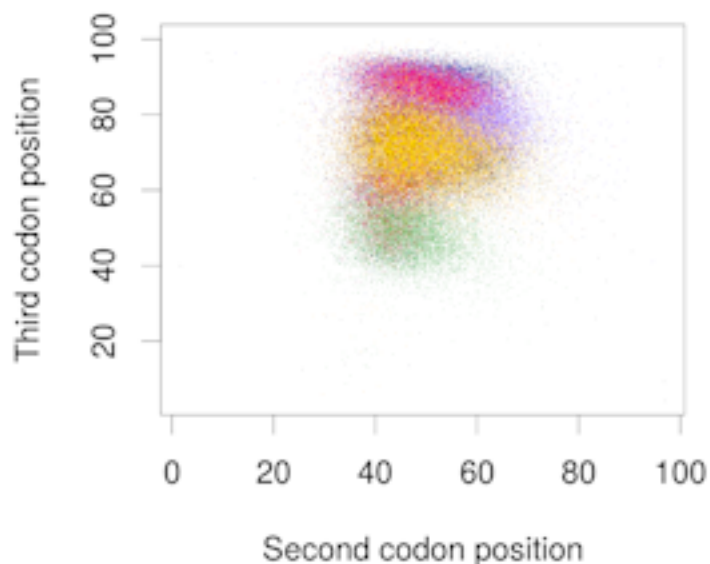
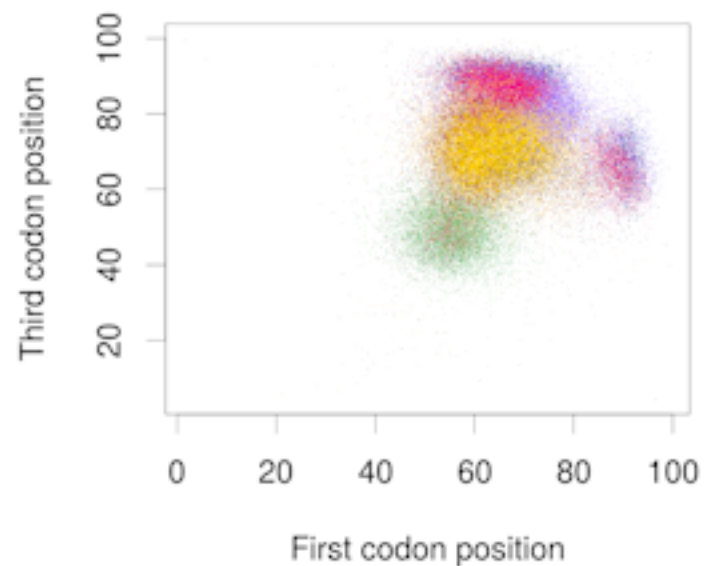
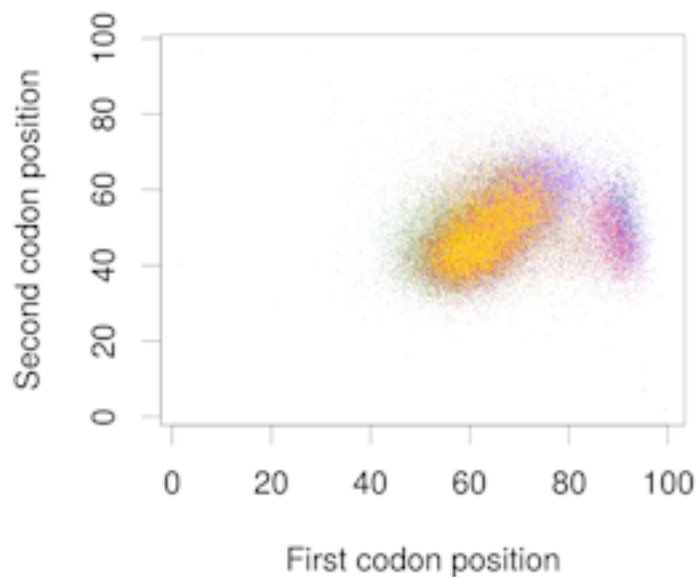


Ser-enriched GO-terms in acidophiles



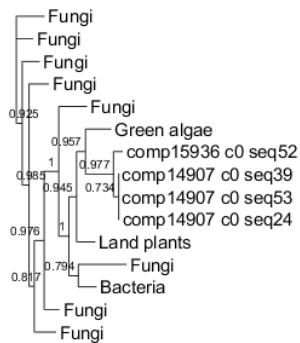
- binding
- catalytic activity
- transporter activity
- other

Average GC content in the three codon positions

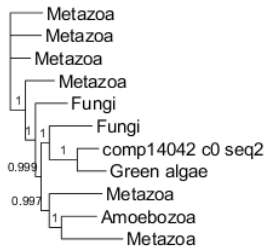


- *Chlamydomonas acidophila*
- *Chlamydomonas reinhardtii*
- *Coccomyxa subellipsoidea*
- *Dunaliella acidophila*
- *Micromonas pusilla*
- *Micromonas sp. RC299*
- *Ostreococcus lucimarinus*
- *Volvox carteri*

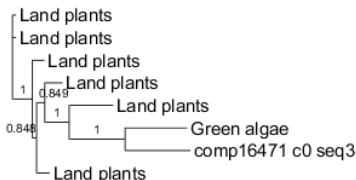
Additional File 4. Simplified phylogenetic analyses of additional transcripts coding for genes with an annotation related to heavy metal tolerance and present in *C. acidophila* but not in *C. reinhardtii*. The phylograms represent the majority consensus of trees sampled after stationarity in the Bayesian analysis. PP values equal or greater than 0.50 are shown above branches. The scale bar indicates relative distance between different sequences based on mutation rate. A) ACR3 comp_14907 and comp_15936 B) MATE efflux protein comp_12911 C) NRAMP family comp_13602 D) ABC-ATPase comp_14042 E) arsenite-antimonite efflux family comp_15332 F) glutathione-regulated potassium-efflux comp_16013 G) multidrug efflux transporter AcrB comp_16471 H) multidrug resistance-associated protein comp_17557.



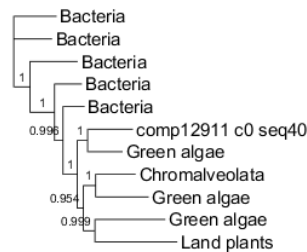
A) ACR3 comp_14907 and comp_15936



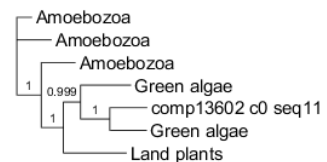
D) ABC-ATPase comp_14042



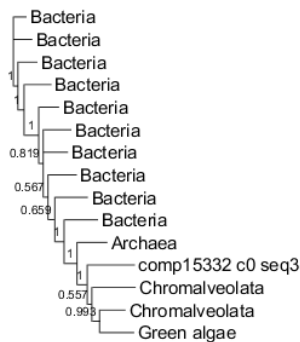
G) multidrug efflux transporter AcrB comp_16471



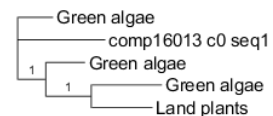
B) MATE efflux protein comp_12911



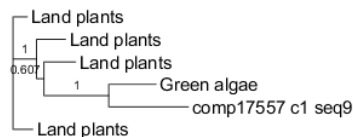
C) NRAMP family comp_13602



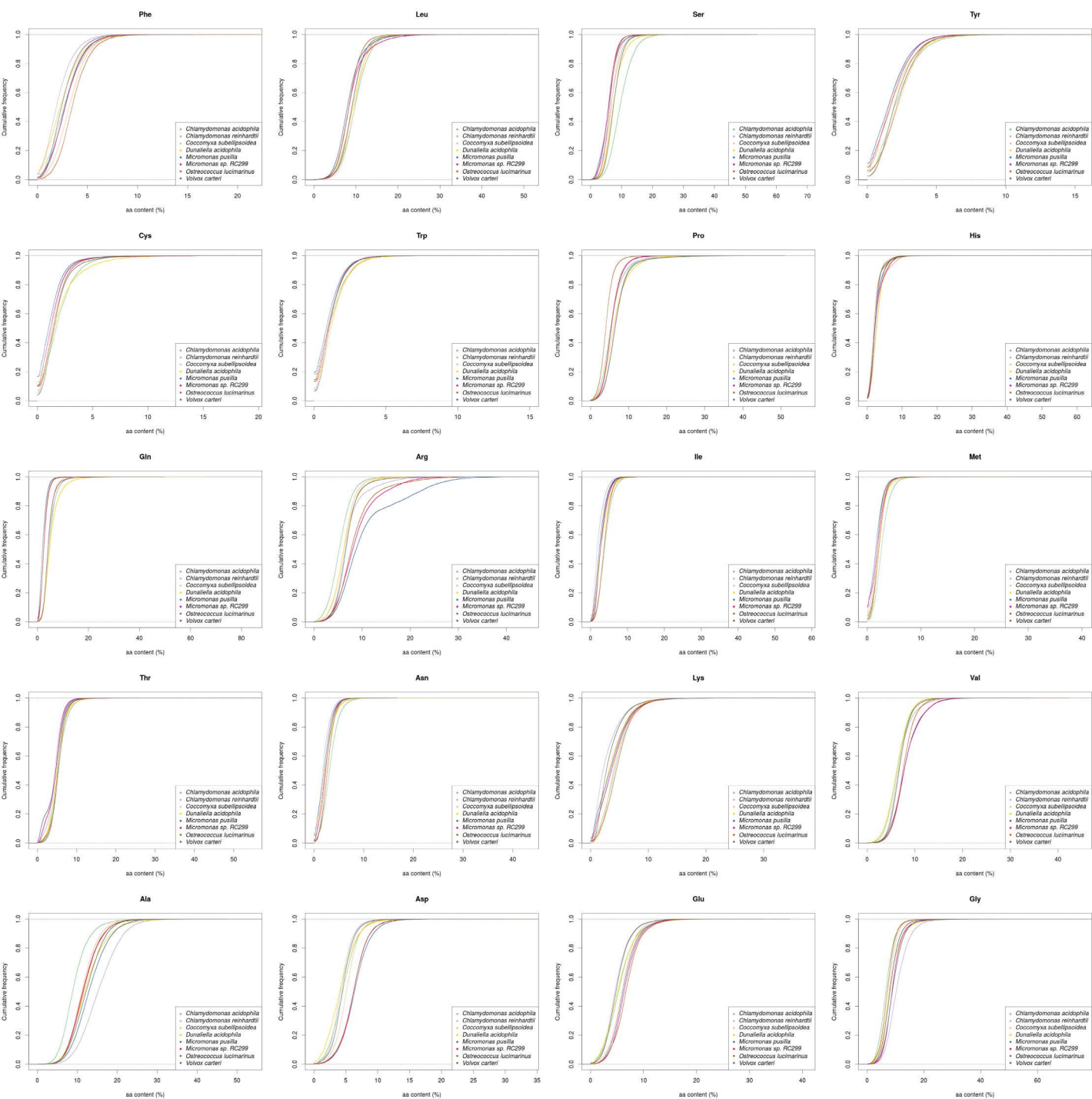
E) arsenite-antimonite efflux family comp_15332



F) glutathione-regulated potassium-efflux comp_16013



H) multidrug resistance-associated protein comp_17557



Additional File 5. Aminoacid utilization profiles for the different species. There is one cumulative frequency plot for each aminoacid.