1  **Developing robust protein analysis profiles to identify bacterial acid phosphatases in**

2  **genomes and metagenomic libraries**

3

4  Zulema Udaondo[1,2*], Estrella Duque[1*], Abdelali Daddaoua[3], Carlos Caselles[1], Amalia

5  Roca[4], Paloma Pizarro-Tobias[4], and Juan L. Ramos[1]

6

7  [1]Estación Experimental del Zaidín, CSIC, E-18008 Granada, Spain

8  [2]Department of Biomedical Informatics, University of Arkansas for Medical Sciences,

9  Little Rock, AR 72205, USA

10  [3]Department of Biochemistry and Molecular Biology II, Faculty of Pharmacy, University

11  of Granada, Granada, Spain

12  [4]Bio-Iliberis R&D, Peligros, Granada, Spain

13

14  *These two co-authors contributed equally to this work

15

16  Corresponding author: Juan L. Ramos

17  Contact information juanluis.ramos@eez.csic.es

18  Phone: +34 958181600 extension 289

19

20

21

22

23

24

25

26

27

28

29

30

31 ABSTRACT

32 Phylogenetic analysis of more than 4000 annotated bacterial acid phosphatases was

33 carried out. Our analysis enabled us to sort these enzymes into the following three

34 types: 1) class B acid phosphatases, which were distantly related to the other types, 2)

35 class C acid phosphatases, and 3) generic acid phosphatases (GAP). While class B

36 phosphatases are found in a limited number of bacterial families, which include known

37 pathogens, class C acid phosphatases and GAP proteins are found in a variety of

38 microbes that inhabit soil, fresh water and marine environments. As part of our analysis

39 we developed three profiles, named Pfr-B-Phos, Pfr-C-Phos and Pfr-GAP, to describe the

40 three groups of acid phosphatases. These sequence-based profiles were then used to

41 scan genomes and metagenomes to identify a large number of formerly unknown acid

42 phosphatases. A number of proteins in databases annotated as hypothetical proteins

43 were also identified by these profiles as putative acid phosphatases.  To validate these

44 *in silico* results, we cloned genes encoding candidate acid phosphatases from genomic

45 DNA, or recovered from metagenomic libraries or genes synthesized in vitro based on

46 protein sequences recovered from metagenomic data. Expression of a number of these

47 genes, followed by enzymatic analysis of the proteins, further confirmed that sequence

48 similarity searches using our profiles could successfully identify previously unknown acid

49 phosphatases.

50

51

52 **INTRODUCTION**

53 Phosphorous is a major component of cells in all living organisms and all prokaryotic and

54 eukaryotic cells have developed mechanisms for the uptake of inorganic phosphate,

55 which is used in the biosynthesis of phospholipids, sugar phosphates, nucleotides and

56 other molecules (Barea and Richardson, 2015). Despite phosphorous being one of the

57 most abundant non-metallic elements in the earth's crust, it is frequently found in forms

58 that are not bioavailable - a reality that often leads to phosphorous nutrient limitation

59 (Ågren *et al.*, 2012; Sosa *et al.*, 2019). Inorganic phosphorous forms are often solubilized

60 by plants and microorganisms (bacteria and fungi) through the production of weak acids

61  (Barea and Richardson, 2015). However, a number of common organic phosphorous

62  compounds (i.e., phytic acid, sugar phosphates, nucleotides, phospholipids and others)

63  must be first hydrolysed by phosphatases to yield inorganic phosphate, which can

64  subsequently be taken up by microorganisms and plants to be used as a phosphorous

65  source (Hayes *et al.*, 2000; Alori *et al.*, 2017; Thomashow *et al.*, 2018). Evidence suggests

66  that phosphatase activity in soils and aquatic environments is of ecological relevance

67  and is a driver of the productivity of terrestrial ecosystems (Turner *et al.*, 2013; Margalef

68  *et al.*, 2017) and influence primary and secondary production in fresh waters and marine

69  environments (Martiny *et al.*, 2019).

70  There are two types of phosphatases among the phosphoric ester hydrolases which are

71  defined based on their optimal pH. Alkaline phosphatases are a broad group of well

72  characterized enzymes that use different mechanisms and co-factors to carry out their

73  function (Mullaney and Ullah, 2003; Ragot *et al.*, 2015; Lidbury *et al.*, 2017; Neal *et al.*,

74  2018). Acid phosphatases are, in general, non-specific phosphatases with broad

75  substrate specificity and are often secreted across the outer membrane or are located

76  in the periplasmic space (Thaller *et al.*, 1997). At least three different types of

77  prokaryotic phosphatases that function at acidic pH have been distinguished mainly

78  based on their sequences; they are known as types A, B and C (Thaller *et al.*, 1997;

79  Lidbury *et al.*, 2017; Neal *et al.*, 2018). It was noted that the B class phosphatases are

80  generally associated with pathogenic microbes while the other types are widely

81  distributed in nature (Neal *et al.*, 2018). While the importance of acid phosphatases to

82  the acquisition of phosphorous in soils, fresh waters and marine environments (Neal *et*

83  *al.*, 2018); Margalef *et al* (2017) compiled phosphatase activities from a large number

84  studies of natural ecosystems and made 329 observations for acid phosphatases versus

85  72 for alkaline phosphatases, highlighting the environmental importance of  acid

86  phosphatases.

87  The work described here aims to contribute further to the understanding of organic

88  phosphorous mobilisation in the environment by acid phosphatases. To this end we

89  developed three robust profiles that can unequivocally identify the different types of

90  acid phosphatase types. We have empirically validated the profiles by cloning and

91    expression of putative acid phosphatases rescued from genomes, or recovered from

92    functional metagenomic libraries or genes synthesized *in vitro* based on protein

93    sequences recovered from metagenomic libraries (Fierer *et al.*, 2013; Berini *et al.*, 2017;

94    Duque *et al.*, 2018). This profiling methodology will serve as a valuable resource for the

95    identification of these important enzymes within the preponderance of already

96    sequenced genomes and widely available metagenomic data. Furthermore, this study

97    provides a proof-of-concept for the successful use of profiles to characterize enzymes

98    involved in biogenic cycles.

99

100   **RESULTS AND DISCUSSION**

101   As a first step towards the identification of bacterial acid phosphatases we retrieved

102   4644 sequences annotated as bacterial acid phosphatases (either due to protein name

103   or Pfam domain composition) from the Uniprot Database (UniProt: a worldwide hub of

104   protein knowledge, 2019). A phylogenetic tree was constructed with a refined set of

105   3741 protein sequences (see Experimental Procedures) and the results are shown in

106   Figure 1. The bacterial acid phosphatase tree has, as expected, three clear branches; one

107   represented by the outer blue circle which corresponds to class B (Figure 1), another

108   represented by the outer purple circle that corresponds to class C and the other

109   represented by the outer green circle that corresponds to Generic Acid Phosphatases

110   class A (GAP) (Figure 1). Supplementary Table 1 contains information collected from the

111   Uniprot database for each of the refined datasets of acid phosphatases. The

112   phylogenetic tree from Supplementary Figure 1 shows that acid phosphatases from class

113   GAP, B and C belong to three well defined monophyletic groups. The unrooted tree also

114   revealed that sequences from class A and C are closest relatives and therefore

115   sequences from class B are from an evolutionary point of view more distant from the

116   other two.

117   The blue branch of the tree grouped 512 sequences that corresponded with annotated

118   acid phosphatases of class B, the purple branch included 1701 sequences of annotated

119   class C acid phosphatases; while the other set, which we named GAP, comprised 1528

120   non-specific class A acid phosphatases. The analysis of the sequences at the family level

4

showed that class C and GAP proteins were found widely distributed among microbes that inhabit soils, fresh water and marine environments. In contrast, class B acid phosphatases are present in a limited number of microbial families which include Enterobacteriaceae, Pasteurellaceae, Morganellaceae, Aeromonadaceae and Vibrionaceae (Figure 1 and 2), of which some are pathogens (supplementary Table 2A, 2B, 2C). Conversely, it should be noted that class C and GAP acid phosphatases were also present in some Enterobacteriaceae. For example, in *Salmonella* and *Klebsiella* genomes GAP proteins were identified and in a number of *Enterobacter* species (mainly *cloacae*) class C proteins were found. In contrast in the *Escherichia coli* species, despite being a broad taxonomic group (Abram *et al.*, 2020), only class B acid phosphatases were identified (supplementary Table 2).

Bacterial acid phosphatases have previously been identified through a number of signatures; for example, the database of families and domain proteins PROSITE (Sigrist *et al.*, 2002) identified bacterial acid phosphatase sequences based on short sequence pattern motifs defined by the signature PS01157 (pattern G-S-Y-P-S-G-H-T). The compendium of protein fingerprints *PRINTS* database (Attwood *et al.*, 2000), contains the signature PR00483 which corresponds to a 5-element fingerprint from bacterial acid phosphatases derived from an initial alignment of a limited number of sequences. Four profiles were available from TIGRFAM database that were constructed using a limited set of acid phosphatase sequences (TIGR03397, 01675, 01672 and 01668); however, these profiles were found to have no discriminatory power. Other databases, such as Pfam domain protein database (El-Gebali *et al.*, 2019) and Simple Modular Architecture Research Tool (SMART) (Letunic and Bork, 2018) contain a number of entries related to identification and classification of bacterial acid phosphatases. Nonetheless, none of the above motifs and classifications distinguish unequivocally between the three classes of bacterial acid phosphatases.

To establish a new criterion defining the three kinds of acid phosphatases represented in the phylogenetic tree, we decided to explore the construction of PROSITE generalized profiles, which are not available in the PROSITE database (https://prosite.expasy.org). Profiles are weight matrices that are useful for grouping proteins into families (Gromiha,

151   2010) and use quantitative motif descriptors which are given as linear sequences that

152   comprise weighted match or mismatch residues and insert sequences in a profile

153   position (Sigrist *et al.*, 2002). Given that the phylogenetic tree defined three branches,

154   according to differences in their amino acid sequences, we expected that a Profile for

155   each of the branches would result in a net gain in specificity for identification and

156   assignation of the entire collection of bacterial acid phosphatases.

157   To construct the three new profiles, we proceeded as suggested by PROSITE

158   ([https://prosite.expasy.org/prosuser.html#meth_prf](https://prosite.expasy.org/prosuser.html#meth_prf)). To create these profiles, we used

159   the three sets of proteins identified in each of the branches of the tree, the profile for

160   class B phosphatases (Prf-B-Phos) was constructed using a set of 512 seed sequences,

161   for the profile for class C we used 1701 sequences while for the profile for GAP (Pfr-

162   GAP), due to the high variability in the sequence similarity and sequence length from

163   members of this class, we used a filtered set of 948 out of the 1528 sequences from the

164   previous analysis (supplementary Tables 3A, 3B and 3C). The three profiles obtained in

165   this study are publicly available in supplementary Table 4. The generation of a profile

166   requires a multiple-alignment of the seed sequences as input, which was performed

167   using Muscle (Edgar, 2004). The consensus sequences derived from the multiple-

168   alignments (supplementary Tables 3A, 3B and 3C) showed conserved regions with high-

169   sequence identity scattered throughout the full sequence of the proteins. This reflects

170   the existence of several functional constraint regions, with lower site-specific

171   substitution ratio distributed along the protein sequences belonging to each of the

172   classes. This is in contrast with most sequence patterns where high-sequence identity

173   regions are restricted to active sites, cofactor binding domains or specific DNA binding

174   regions (Fuglebakk *et al.*, 2012).

175   The multiple-alignment revealed that the short patterns used previously to define these

176   acid phosphatases were in a wider sequence identity context and this warranted the

177   construction of profiles to encompass the full gamut of acid phosphatase sequences

178   belonging to these families. We used the script *pfmake* to translate the multiple-

179   alignment into a matrix table of positions and convert frequency distributions into

180   positive specific amino acid weights and gaps according to the original algorithms of

181    Sibbald and Argos (Sibbald and Argos, 1990) and Gribskov *et al* (1987). Once the profiles

182    were constructed we proceeded to calibrate and validate the profiles as recommended

183    by PROSITE (described in Experimental Procedures), for this the profiles were run against

184    a database to produce a list of sorted scores. It has been previously empirically

185    determined that cut-off values of Z-scores equal or greater than 8.5 are biologically

186    significant and warrant the correct assignment of a protein to a family (Gallegos *et al.*,

187    1997; Sigrist *et al.*, 2002; Godoy *et al.*, 2010).

188    As a proof of concept, the three profiles were used as input for *pfsearch v2.3* from the

189    PTOOLS suite to scan the complete set of Uniref100 proteins (downloaded from the

190    UniProt database on May 24, 2019). As a result, 6000 proteins were matched with Pfr-

191    GAP (Figure 3 and Supplementary Figure 2), 2132 protein sequences were matched by

192    the Pfr-B-Phos (Figure 3 and Supplementary Figure 3) and 10494 with Pfr-C-Phos (Figure

193    3 and Supplementary Figure 4).

194    We found that Pfr-B-Phos identified acid phosphatases preferentially from

195    enterobacteria, vibrios and other microorganisms mainly from orders Pasteurellales and

196    Bacillales (see Supplementary Figure 3) whose life style indicated a close relationship

197    with eukaryotes, as mentioned above, and confirming previous studies (Gandhi and

198    Chandra, 2012; Neal *et al.*, 2018). Conversely, we found that Pfr-C-Phos and Pfr-GAP

199    identified acid phosphatases from a variety of different sources in a highly-specific and

200    sensitive manner, including Acidobacteria, Actinobacteria, alpha, beta, gamma and

201    epsilon proteobacteria, Firmicutes, Verrumicrobia, and Bacterioidetes among many

202    others (see supplementary Figure 2, and supplementary Figure 4). The results obtained

203    with the three profiles against Uniref100 database (Suzek *et al.*, 2015) demonstrated the

204    ability of Pfr-GAP, Pfr-C-Phos and Pfr-B-Phos to discriminate between all classes of acid

205    phosphatases displayed in the phylogenetic tree and within a wide taxonomic range. It

206    is worth noting that although the three profiles were developed using only bacterial

207    sequences, presumed eukaryotic acid phosphatases were also found in all cases. The

208    complete set of raw hits sorted by output score is shown in supplementary Table 5.

209    Remarkably, although these are non-filtered results, the high accuracy of the three

210    profiles allowed the identification of proteins belonging to each of the classes at very

7

211  low score numbers. The specificity of the profiles also identified a large number of

212  putative acid phosphatase sequences which were annotated in the Uniref100 database

213  as "uncharacterized protein".

214  To further validate the new profiles, we decided to test if Pfr-GAP, Pfr-C-Phos, and Pfr-

215  B-Phos could identify acid phosphatases within available annotated whole genomes, in

216  metagenomic libraries in which proteins are annotated as Hypothetical Proteins of

217  unknown function, as well as proteins recovered from functional metagenomic libraries

218  after screening for positive phosphatase activity. We found that the Pfr-GAP, Pfr-C-Phos

219  and Pfr-B-Phos profiles could indeed identify a number of potential acid phosphatases

220  in all of these screens. Specifically, we found that in the annotated reference genomes

221  collected from the NCBI database 4649 proteins were identified by the Prf-A GAP profile,

222  862 by the Pfr-C-Phos profile and 128 proteins by the Pfr-B-phos profile (Figure 3). For

223  most type strains the number of GAP acid phosphatases and class C was between 1 and

224  3, although we found 13 genomes with 6 GAP acid phosphatases and 2 genomes with

225  up to 5 class C acid phosphatase. In those genomes in which an acid phosphatase of class

226  B was present, a single gene was always found, except in one case in which a duplication

227  was identified, and another genome which bore 4 class B acid phosphatase genes. As

228  validation of the proof of concept, we rescued acid phosphatases from the genomes of

229  two microorganisms (i.e., *Pyrococcus*, and *Bacillus subtilis* strain 168). A search using the

230  three profiles with *pfsearch* against the genomes of *Pyrococcus furiosus* DSM 3638*,* and

231  *Bacillus subtilis* str. 168 identified the protein sequences PF0040 and BSU_36530 as

232  putative GAP acid phosphatases, encoded in each genome respectively. *Bacillus subtilis*

233  BSU_36530 was previously annotated as undecaprenyl diphosphatase, while PF0040

234  from *P. furiosus* was annotated as an acidic acid phosphatase. To confirm these 'hits'

235  empirically, we used whole chromosomal DNA from these microorganisms and cloned

236  the amplified DNA into pET28 as described in Experimental Procedures (Table 1).

237  As an initial step for confirmation of phosphatase activity, we spread the cells on LB

238  medium supplemented with BCIP and found that colonies turned deep blue, suggesting

239  that the cloned genes encoded, as expected, phosphatases. A single random clone

240  bearing the gene from each of the two microorganisms was kept. Then, cells were grown

8

241 in liquid LB and acid phosphatase activity determined over a wide pH range in
242 permeabilised cells as described in Experimental Procedures. The results revealed that
243 the optimal pH was in the range of 5 to 6 (Table 2).

244 Our laboratory previously screened a functional metagenomic library from
245 hydrocarbon-polluted soil after land farming and identified a clone, named FOS M2-62,
246 that had robust phosphatase activity (see Experimental Procedures). The fosmid of this
247 clone was sequenced and our profiles were used to identify it as a putative GAP acid
248 phosphatase. We subsequently cloned it into pET28 to generate pET28_FOS M2-62.
249 Phosphatase assays revealed that the AP-M2-62 protein had high activity between pH 5
250 and pH 7 (Table 2), but lower activity at pH greater than 7 or lower than 5. This suggests
251 that AP-M2-62 is indeed an acid phosphatase.

252 We then explored the ability of our constructed profiles to identify hypothetical proteins
253 as putative acid phosphatase from metagenomic libraries. To this end we screened
254 1,552,866 hypothetical proteins from soil metagenomes and 4,925,568 sequences from
255 marine metagenomes (downloaded in June 2019 from the NCBI database) and we found
256 that the search yielded a total of 539 hypothetical proteins from the soil metagenome
257 and 351 hypothetical proteins from marine metagenomes using Pfr-GAP profile
258 (Supplementary Table 5). The Pfr-C-Phos profile was able to find 242 proteins from
259 marine metagenomes and 23 from terrestrial metagenomes. The Pfr-B-Phos profile was
260 able to find only 11 proteins from marine metagenomes. These results are in line with
261 the initial phylogenetic tree results in the sense that class B proteins are poorly
262 represented in marine and terrestrial ecosystems.

263 This data confirmed that among non-characterized acid phosphatases, generic acid
264 phosphatases and class C phosphatases were more abundant than class B, and that class
265 C and GAP can be considered cosmopolitan proteins as they can be found in a wide range
266 of niches. We found that among the set of non-characterized proteins 1 acid
267 phosphatase could be rescued per 3,000 sequences in soil metagenomes while 1 acid
268 phosphatase protein was found every 14,000 sequences in marine metagenomes.
269 Because the quality of metagenomic sequences is non-homogeneous and because our

270    data are raw hit counts, at present we cannot make any conclusions regarding the

271    biogeographic distribution of acid phosphatases based on metagenomic data.

272    Considering the apparent abundance of these sequences, we explored whether the

273    identified sequences were indeed acid phosphatases. To this end we choose two

274    sequences with the highest Z-score from each acid phosphatase family (Supplementary

275    Table 6) and synthesised the corresponding genes. We then cloned and expressed them

276    in *Escherichia coli* and enzyme activity was determined in permeabilised whole cells

277    using the Britton-Robinson poly-buffer. We found that the six metagenomic acid

278    phosphatase had optimal activity at acidic pH (Table 2 and supplementary Table 7).

279    These results further validate the ability of the profiles to find acid phosphatase enzymes

280    from metagenomes. It is worth mentioning that although the MET_A1 enzyme exhibited

281    the highest activity at pH 5.5 to 6, it had significant activity pH in the pH range between

282    5 and 9 (Table 2).

283    In order to further characterize in more detail, the kinetics properties of the

284    metagenomic acid phosphatases, we purified three proteins (see Experimental

285    Procedures) and the kinetics parameters determined using isothermal titration

286    calorimetry (ITC) (Watt, 1990; Williams and Toone, 1993). The initial rate of reaction ($V_o$)

287    with different concentrations of pNPP was determined from the slope of the linear

288    portion of the curve of integrated heats versus time as described by Bianconi (2003). We

289    found that values for $V_o$ followed typical Michaelis-Menten kinetics and $K_{cat}$ and $K_M$ were

290    calculated by fitting the curve to the Michaelis-Menten kinetics equation using non-

291    linear regression (Ababou and Ladbury, 2006). For MET_A_1, M2-F62, and MET_C_1,

292    values of $K_M$ were 49.3 ± 2.6 µM, 29.7 ± 0.02 µM and 23.8 ± 6.9 µM, respectively; and

293    $k_{cat}$ were 0.63  $s^{-1}$, 0.55 $s^{-1}$ , and 0.26 $s^{-1}$, respectively. Our results revealed that the

294    substrate affinities were in the low micromolar range with up to 2-fold differences; $k_{cat}$

295    values differed by up to 2.5-fold. The $K_M$  values we determined are lower than those

296    measured    for    acid    phosphatases    from    different    sources    using    classical

297    spectrophotometric assays (Reilly *et al.*, 2009; Zhang *et al.*, 2013; Wang *et al.*, 2018).

298

299

300 **CONCLUSIONS**

301 In conclusion, we have constructed a phylogenetic tree for acid phosphatases that

302 grouped them into three branches. For each of the branches a Prosite profile was

303 constructed and validated; the three profiles were shown to be effective in the

304 differentiation of the three sets of acid phosphatases. These profiles were able to assign

305 a set of proteins annotated as hypothetical proteins in databases as being acid

306 phosphatases (Suppl. Table 4). We tested our 'hits' empirically and confirmed

307 phosphatase activity at acidic pH. Use of these profiles and the underlying strategy could

308 serve as a powerful approach to explore the role that acid phosphatases play in primary

309 productivity in edaphic and aquatic environments.

310

311

312 **EXPERIMENTAL PROCEDURES**

313

314 *Phylogenetic tree construction*

315 Sequences were downloaded from the Uniprot database by filtering proteins that

316 belong to the Domain = bacteria and the annotation = acid phosphatase and 5'

317 nucleotidase lipoprotein ep4 family; the later corresponds to class C acid phosphatases.

318 Using these filters (on April 26, 2019) we retrieved 4644 protein sequences. Muscle

319 v3.8.1551 (Edgar, 2004) alignment software with parameter - maxiters 1000 was used

320 to align the set of 4644 protein sequences and construct the phylogenetic tree. Very

321 divergent sequences were filtered and removed from the alignment until a final set of

322 3741 amino acid sequences were kept. The final set of sequences was aligned again

323 using Muscle v3.8.1551 with the same parameters. Aligned sequences were used as

324 input for the IQ-TREE software v1.6.10 (Nguyen *et al.*, 2015) with parameters -nt AUTO,

325 -bb 1000 -m TESTMERGE. The maximum likelihood tree was constructed following the

326 model of evolution WAG with parameters F+R10 (IQ-TREE uses ModelFinder).

327 Phylogenetic trees were plotted using the Interactive Tree of Life (iTOL) suite software

328 v4 (Letunic and Bork, 2016).

329

### *Profile construction*

To construct PROSITE "generalized" profiles, first we established the "seed protein sequences" that would determine the sensitivity and average quality of the profiles. Once visualized, the phylogenetic tree branches annotated as class B phosphatases, class C and generic acid phosphatases were aligned separately and filtered according to observed divergences in the alignment. Then *pfw* and *pfmake* scripts from PFTOOLS v2.3 (Gribskov *et al.*, 1987; Sigrist *et al.*, 2002; Bucher *et al.*, 2015) were used to compute new weights for each individual sequence from the multiple sequence alignment and to construct the profile respectively. The matrix BLOSUM 45 was selected for the construction of the profile.

*Pfsearch* and *pfscan* were used to calibrate each profile against a calibration database. The calibration database was made from the entire collection of Swiss-Prot protein sequences filtered by Taxonomy = bacteria. The database contained a total of 334,009 sequences that where shuffled randomly with a sliding window of 20 residues using the script fasta-shuffle-letters from MEME suite v5.0.2 (Bailey *et al.*, 2015).

Searches with the three profiles using Uniref100 database, a local database of representative bacterial and archaea sequences and hypothetical protein databases from metagenomic samples, were all done using *pfscan* script from PFTOOLS v2.3 with parameters -z -f (Bucher *et al.*, 2013).

### *Sequences in databases*

Uniref100 database was downloaded to be used locally in May, 2019. The set of protein FASTA sequences from representative strains was downloaded from the NCBI database in August, 2019. The set of representative strains was obtained via genome browse from NCBI https://www.ncbi.nlm.nih.gov/genome/browse#!/overview/ and then filtered by "archaea" AND "bacteria" AND "representative genome". The two sets of hypothetical proteins used in these analyses were obtained from NCBI protein database using filters: "soil metagenome" AND "hypothetical protein" and "marine metagenome" AND "hypothetical protein"

*Construction of a functional soil metagenomic library*

Soil samples were taken from hydrocarbon polluted soil after land farming. High-molecular-weight DNA extraction was performed from the soil using the commercial GNOME DNA kit (MP, Biomedicals) according to the manufacturer's instructions. DNA fragments of approximately 40 kb were recovered and ligated into the pCC1FOS vector (Epicentre®), and the product was transduced into *E. coli* EPI300 (Raleigh *et al.*, 2002) according to the manufacturer's protocol. Screening for phosphatase activity was performed by replicating the metagenomic library onto agar LB plates with 40 mg per mL of 5-bromo-4-chloro-3-indolyl phosphate (BCIP Applichem, Darmstadt, Germany) as substrate, supplemented with 12.5 µg per mL chloramphenicol and 0.01% *L*-arabinose. Following replication, the colonies were incubated for 24 h at 37°C. A total of 64 clones with phosphatase activity were identified and detected as pale to dark blue colonies. A single clone, named M2-62, that turned deep blue on these plates was used for further analysis in this study.


*Cloning of putative acid phosphatases in* **Escherichia coli.**

DNA from *Pyrococcus furiosus* DSM 3638 and *Bacillus subtilis* DSM 204 were obtained from the DSMZ culture collection. The *Bacillus subtilis* gene was PCR amplified with the following primers 5`-TTGAACTACGAAATTTTTAAAGCAATCC-3` and 5`-TTCTTAGAAATTTTGATCGGTTGG-3`, while the *Pyrococcus* gene was amplified using the following pair of primers 5'-ATGCTGGCAATACTTACGGCAA-3`and 5´-TCACTTATCCACTTTAAAAAAGATGCGC-3´; amplified DNA was subsequently cloned into pTOPO and further subcloned into pET28 after digestion with NdeI and EcoRI. Plasmids were transformed into *E. coli* BL21 (DE3) (Studier *et al.*, 2009). For amplification of the open reading frame encoding the AP-M2-62 protein, fosmid DNA was prepared and the following primers: 5'-CATATGAAAAAAATACCTGAACCCTTC-3' (forward) and 5'-GGATCCTCAGTGCTGGGTCAG-3' (reverse) were used. Following PCR amplification, under standard conditions, the fragment was cloned into the pMBL vector to yield pMBL_FOSM2-62. The plasmid was subsequently digested with NdeI/BamHI and the

389    806 bp fragment bearing the ORF AP-M2-62 was cloned into pET28b (+) digested with

390    the same enzymes (Table 1).

391

*Cloning of putative metagenomic acid phosphatases in* **Escherichia coli**

393    Protein sequences retrieved from metagenomic libraries with a high Z-score for GAP,

394    class B and class C were manually curated.  The protein sequences were then translated

395    into DNA sequences with optimized codon usage for *E. coli,* synthesized in vitro by

396    Genescript, cloned into pET28 and expressed from the P$_{lac}$ .

397

*Growth of* **Escherichia coli** *and in vivo acid phosphatase activity. Escherichia coli* BL21

399    (DE3) transformed with the corresponding plasmid was grown in 100 mL conical flasks

400    containing 25 mL of LB supplemented with 0.025 mg/mL kanamycin (pET28). Cultures

401    were incubated at 37 °C with shaking until they reached a turbidity at 660 nm (OD$_{660}$) of

402    0.6, at which point 0.1 mM isopropyl-$\beta$-D-thiogalactopyranoside (IPTG) was added, to

403    induce expression, incubation was continued overnight. After growth of *E. coli* the

404    turbidity of the cultures was adjusted to 1 in 600 µL of lysis buffer (100 mM acetate, pH

405    5.5,  CaCl$_2$, 1 mM, and Tween 80, 0.01%  (or a drop of toluene) (Lassen *et al.*, 2001). The

406    assay was performed by combining 100 µL of permeabilized cells with 10 µL of a solution

407    of 100 mM *p*-nitrophenyl phosphate (pNNP) dissolved in 0.1 M Na-acetate buffer, pH

408    5.5. The reaction mixture was incubated for 30 min at 25°C. Subsequently, 100 µL of 0.5

409    M sodium hydroxide in water was added to stop the reaction. The samples were then

410    centrifuged in a bench centrifuge (5 min at 10000 rpm) and the absorbance at 405 nm

411    was measured in a spectrophotometer. To determine the optimal pH range the Britton-

412    Robinson poly-buffer (40 mM boric acid, 40 mM acetic acid and 40 mM phosphoric acid)

413    was adjusted with NaOH to a pH between 2 and 9 (Souri *et al.*, 2013). Other conditions

414    for the acid phosphatase assays are those mentioned above.

415

*Protein purification.*

417    For protein purification, cells were suspended in 25 mL of buffer A (50 mM  Hepes pH

418    6.9; 300 mM NaCl; 1 mM dithiothreitol) with EDTA-free protease inhibitor mixture. Cells

419 were lysed by two passes through a French Press at a p.s.i. of 1000. The cell suspension
420 was then centrifuged at 20,000 x g for 1 hour. The pellet was discarded and the
421 supernatant was filtered and loaded onto a 5 mL His-Trap chelating column (GE
422 Healthcare, St. Gibes, UK). The proteins were eluted with a 10 to 500 mM gradient of
423 imidazol in buffer A. The purity of the eluate was determined by running 12% SDS-PAGE
424 gels. Homogenous protein preparations were dialyzed overnight against buffer A but
425 supplemented with 10% [v/v] glycerol). Dialyzed protein was collected at a
426 concentration of about 1 mg/mL and stored in 1 mL aliquots at -80 °C.

427

428

435

436 CONFLICY OF INTEREST:
437 The authors declare no conflict of interest

438

439 **References**

440

Ababou, A. and Ladbury, J.E. (2006) Survey of the year 2004: literature on applications
of isothermal titration calorimetry. *Journal of Molecular Recognition* **19**: 79–89.

Abram, K., Udaondo, Z., Bleker, C., Wanchai, V., Wassenaar, T.M., Robeson, M.S., and
Ussery, D.W. (2020) What can we learn from over 100,000 Escherichia coli
genomes? *bioRxiv* 708131.

Ågren, G.I., Wetterstedt, J.Å.M., and Billberger, M.F.K. (2012) Nutrient limitation on terrestrial plant growth – modeling the interaction between nitrogen and phosphorus. *New Phytologist* 953–960.

Alori, E.T., Glick, B.R., and Babalola, O.O. (2017) Microbial Phosphorus Solubilization and Its Potential for Use in Sustainable Agriculture. *Front Microbiol* **8**:.

Attwood, T.K., Croning, M.D.R., Flower, D.R., Lewis, A.P., Mabey, J.E., Scordis, P., et al. (2000) PRINTS-S: the database formerly known as PRINTS. *Nucleic Acids Res* **28**: 225–227.

Bailey, T.L., Johnson, J., Grant, C.E., and Noble, W.S. (2015) The MEME Suite. *Nucleic Acids Res* **43**: W39–W49.

Barea, J.-M. and Richardson, A.E. (2015) Phosphate Mobilisation by Soil Microorganisms. In *Principles of Plant-Microbe Interactions: Microbes for Sustainable Agriculture*. Lugtenberg, B. (ed). Cham: Springer International Publishing, pp. 225–234.

Berini, F., Casciello, C., Marcone, G.L., and Marinelli, F. (2017) Metagenomics: novel enzymes from non-culturable microbes. *FEMS Microbiol Lett* **364**.

Bianconi, M.L. (2003) Calorimetric Determination of Thermodynamic Parameters of Reaction Reveals Different Enthalpic Compensations of the Yeast Hexokinase Isozymes. *J Biol Chem* **278**: 18709–18713.

Bucher, P., Cerutti, L., Pagni, M., and Schuepbach, T. (2013) PfTools Software Suite.

Bucher, P., Karplus, K., Moeri, N., and Hofmann, K. (2015) A Flexible Motif Search Technique Based on Generalized Pro les. *Computers and Chemistry* **20**: 3–24.

Duque, E., Daddaoua, A., Cordero, B.F., Udaondo, Z., Molina-Santiago, C., Roca, A., et al. (2018) Ruminal metagenomic libraries as a source of relevant hemicellulolytic enzymes for biofuel production. *Microbial Biotechnology* **11**: 781–787.

Edgar, R.C. (2004) MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* **32**: 1792–1797.

El-Gebali, S., Mistry, J., Bateman, A., Eddy, S.R., Luciani, A., Potter, S.C., et al. (2019) The Pfam protein families database in 2019. *Nucleic Acids Res* **47**: D427–D432.

Fierer, N., Ladau, J., Clemente, J.C., Leff, J.W., Owens, S.M., Pollard, K.S., et al. (2013) Reconstructing the Microbial Diversity and Function of Pre-Agricultural Tallgrass Prairie Soils in the United States. *Science* **342**: 621–624.

Fuglebakk, E., Echave, J., and Reuter, N. (2012) Measuring and comparing structural fluctuation patterns in large protein datasets. *Bioinformatics* **28**: 2431–2440.

Gallegos, M.T., Schleif, R., Bairoch, A., Hofmann, K., and Ramos, J.L. (1997) Arac/XylS family of transcriptional regulators. *Microbiol Mol Biol Rev* **61**: 393–410.

Godoy, P., Molina-Henares, A.J., Torre, J.D.L., Duque, E., and Ramos, J.L. (2010) Characterization of the RND family of multidrug efflux pumps: in silico to in vivo confirmation of four functionally distinct subgroups. *Microbial Biotechnology* **3**: 691–700.

Gribskov, M., McLachlan, A.D., and Eisenberg, D. (1987) Profile analysis: detection of distantly related proteins. *PNAS* **84**: 4355–4358.

Gromiha, M.M. (2010) Protein bioinformatics: from sequence to function, Academic Press.

Hayes, J.E., Richardson, A.E., and Simpson, R.J. (2000) Components of organic phosphorus in soil extracts that are hydrolysed by phytase and acid phosphatase. *Biol Fertil Soils* **32**: 279–286.

Lassen, S.F., Breinholt, J., Østergaard, P.R., Brugger, R., Bischoff, A., Wyss, M., and Fuglsang, C.C. (2001) Expression, Gene Cloning, and Characterization of Five Novel Phytases from Four Basidiomycete Fungi: *Peniophora lycii, Agrocybe pediades, a Ceriporia* sp., and *Trametes pubescens*. *Appl Environ Microbiol* **67**: 4701–4707.

Letunic, I. and Bork, P. (2018) 20 years of the SMART protein domain annotation resource. *Nucleic Acids Res* **46**: D493–D496.

Letunic, I. and Bork, P. (2016) Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res* **44**: W242–W245.

Lidbury, I.D.E.A., Fraser, T., Murphy, A.R.J., Scanlan, D.J., Bending, G.D., Jones, A.M.E., et al. (2017) The 'known' genetic potential for microbial communities to degrade organic phosphorus is reduced in low-pH soils. *MicrobiologyOpen* **6**: e00474.

Margalef, O., Sardans, J., Fernández-Martínez, M., Molowny-Horas, R., Janssens, I.A., Ciais, P., et al. (2017) Global patterns of phosphatase activity in natural soils. *Sci Rep* **7**: 1–13.

Martiny, A.C., Lomas, M.W., Fu, W., Boyd, P.W., Chen, Y.L., Cutter, G.A., et al. (2019) Biogeochemical controls of surface ocean phosphate. *Science Advances* **5**: eaax0341.

Mullaney, E.J. and Ullah, A.H.J. (2003) The term phytase comprises several different classes of enzymes. *Biochemical and Biophysical Research Communications* **312**: 179–184.

Neal, A.L., Blackwell, M., Akkari, E., Guyomar, C., Clark, I., and Hirsch, P.R. (2018) Phylogenetic distribution, biogeography and the effects of land management upon bacterial non-specific Acid phosphatase Gene diversity and abundance. *Plant Soil* **427**: 175–189.

Nguyen, L.-T., Schmidt, H.A., von Haeseler, A., and Minh, B.Q. (2015) IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol Biol Evol* **32**: 268–274.

Ragot, S.A., Kertesz, M.A., and Bünemann, E.K. (2015) *phoD* Alkaline Phosphatase Gene Diversity in Soil. *Appl Environ Microbiol* **81**: 7281–7289.

Raleigh, E.A., Elbing, K., and Brent, R. (2002) Selected Topics from Classical Bacterial Genetics. *Current Protocols in Molecular Biology* **59**: 1.4.1-1.4.14.

Reilly, T.J., Chance, D.L., Calcutt, M.J., Tanner, J.J., Felts, R.L., Waller, S.C., et al. (2009) Characterization of a Unique Class C Acid Phosphatase from Clostridium perfringens. *Appl Environ Microbiol* **75**: 3745–3754.

Sibbald, P.R. and Argos, P. (1990) Weighting aligned protein or nucleic acid sequences to correct for unequal representation. *J Mol Biol* **216**: 813–818.

Sigrist, C.J.A., Cerutti, L., Hulo, N., Gattiker, A., Falquet, L., Pagni, M., et al. (2002) PROSITE: A documented database using patterns and profiles as motif descriptors. *Brief Bioinform* **3**: 265–274.

Sosa, O.A., Repeta, D.J., DeLong, E.F., Ashkezari, M.D., and Karl, D.M. (2019) Phosphate-limited ocean regions select for bacterial populations enriched in the carbon–phosphorus lyase pathway for phosphonate degradation. *Environmental Microbiology* **21**: 2402–2414.

Souri, E., Kaboodari, A., Adib, N., and Amanlou, M. (2013) A New extractive spectrophotometric method for determination of rizatriptan dosage forms using bromocresol green. *DARU J Pharm Sci* **21**: 12.

Studier, F.W., Daegelen, P., Lenski, R.E., Maslov, S., and Kim, J.F. (2009) Understanding the Differences between Genome Sequences of Escherichia coli B Strains REL606 and BL21(DE3) and Comparison of the E. coli B and K-12 Genomes. *Journal of Molecular Biology* **394**: 653–680.

Suzek, B.E., Wang, Y., Huang, H., McGarvey, P.B., and Wu, C.H. (2015) UniRef clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics* **31**: 926–932.

Thaller, M.C., Schippa, S., Bonci, A., Cresti, S., and Rossolini, G.M. (1997) Identification of the gene (aphA) encoding the class B acid phosphatase/phosphotransferase of Escherichia coli MG1655 and characterization of its product. *FEMS Microbiol Lett* **146**: 191–198.

Thomashow, L.S., LeTourneau, M.K., Kwak, Y.-S., and Weller, D.M. (2018) The soil-borne legacy in the age of the holobiont. *Microbial Biotechnology* **12**: 51–54.

Turner, B.L., Lambers, H., Condron, L.M., Cramer, M.D., Leake, J.R., Richardson, A.E., and Smith, S.E. (2013) Soil microbial biomass and the fate of phosphorus during long-term ecosystem development. *Plant Soil* **367**: 225–234.

U. Gandhi, N. and B. Chandra, S. (2012) A COMPARATIVE ANALYSIS OF THREE CLASSES OF BACTERIAL NON-SPECIFIC ACID PHOSPHATASES AND ARCHAEAL PHOSPHOESTERASES: EVOLUTIONARY PERSPECTIVE. *Acta Inform Med* **20**: 167–173.

UniProt: a worldwide hub of protein knowledge (2019) *Nucleic Acids Res* **47**: D506–D515.

Wang, Z., Tan, X., Lu, G., Liu, Y., Naidu, R., and He, W. (2018) Soil properties influence kinetics of soil acid phosphatase in response to arsenic toxicity. *Ecotoxicology and Environmental Safety* **147**: 266–274.

Watt, G.D. (1990) A microcalorimetric procedure for evaluating the kinetic parameters of enzyme-catalyzed reactions: Kinetic measurements of the nitrogenase system. *Analytical Biochemistry* **187**: 141–146.

Wickham H (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. ISBN 978-3-319-24277-4, https://ggplot2.tidyverse.org.
Williams, B.A. and Toone, E.J. (1993) Calorimetric evaluation of enzyme kinetic parameters. *J Org Chem* **58**: 3507–3510.

Zhang, G.-Q., Chen, Q.-J., Sun, J., Wang, H.-X., and Han, C.-H. (2013) Purification and characterization of a novel acid phosphatase from the split gill mushroom *Schizophyllum commune*. *Journal of Basic Microbiology* **53**: 868–875.

441

442

443

444

445

446 **Table 1: Strains and plasmids used in this study**.

447

448

| Strains or plasmids | Genotype or relevant characteristics | Reference |
|---|---|---|
| *Escherichia coli* EPI 300 | *recA1, endA1, araD139, rpsL, nupG, trfA* | Epicenter (Studier *et al.*, 2009) |
| *Escherichia coli* BL21(DE3) | F'/ *ompI, hsdS, gal, dam, met* | |
| **Plasmids** | | |
| pMBL | Vector for cloning PCR amplicons, Ap | Dominion |
| pET28a | Expression vector, 6xHis, Km | Novagen |
| pET28::FOS M2-62 | pET28 containing the complete gene encoding acid phosphatase FOSM 2-62 | This study |
| pET28:BSU | pET28 containing the complete gene encoding acid phosphatase from *Bacillus subtilis* | This study |
| pET28:PYR | pET28 containing the complete gene encoding acid phosphatase from *Pyrococcus furiosus* | This study |
| pET28:MET_A1 | pET28 containing the complete gene encoding the MEAT_A1 GAP acid phosphatase deduced from environmental metagenomes | This study |
| pET28:MET_A2 | pET28 containing the complete gene encoding the MEAT_A2 GAP acid phosphatase deduced from environmental metagenomes | This study |
| pET28:MET_B1 | pET28 containing the complete gene encoding the MEAT_B1 class B acid phosphatase deduced from environmental metagenomes | This study |
| pET28:MET_B2 | pET28 containing the complete gene encoding the MEAT_B2 class B acid phosphatase deduced from environmental metagenomes | This study |
| pET28:MET_C1 | pET28 containing the complete gene encoding the MEAT_C1 class C acid phosphatase deduced from environmental metagenomes | This study |
| pET28:MET_C2 | pET28 containing the complete gene encoding the MEAT_C2 class C acid phosphatase deduced from environmental metagenomes | This study |

449

450 Ap and Km stand for resistance to ampicillin and kanamycin.

451

452

453

**Table 2. Relative acid phosphatase activity of genes amplified from genomic DNA and recovered from metagenomic libraries at different pHs.**

456
457
458

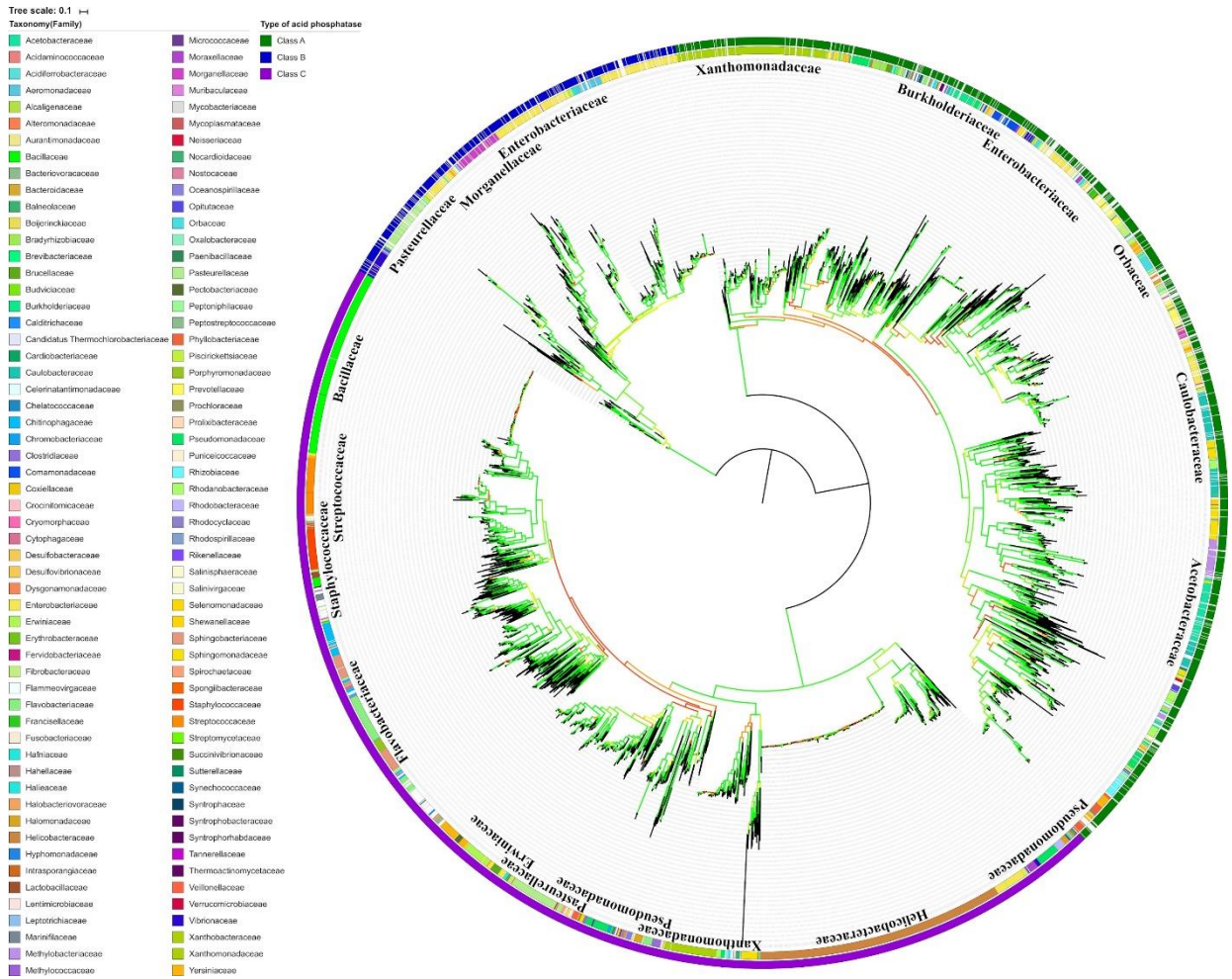| | Enzyme source | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| pH | MET_A1 | MET_A2 | MET_B1 | MET_B2 | MET_C1 | MET_2 | M2-62 | Bacillus |
| 2 | 1 | 5 | 30 | 5 | 3 | 2 | 2 | 2 |
| 3 | 9 | 15 | 30 | 59 | 23 | 8 | 8 | 15 |
| 4 | 41 | 23 | 41 | 56 | 77 | 21 | 21 | 22 |
| 5 | 79 | 90 | 50 | 55 | 106 | 30 | 85 | 90 |
| 5.5 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| 6 | 97 | 16 | 73 | 43 | 61 | 80 | 98 | 97 |
| 7 | 93 | 7 | 59 | 35 | 16 | 10 | 93 | 81 |
| 8 | 71 | 1 | 39 | 17 | 7 | 9 | 47 | 30 |
| 9 | 33 | 2 | 32 | 5 | 5 | 6 | 16 | 9 |

459

460

The set of acid phosphatases were expressed in *Escherichia coli* and the assays carried out as described in Materials and Methods at different pH in Britton-Robisson poly-buffer. Activities are expressed as relative activity, the maximum activity for all of the enzymes was at pH 5.5 and the corresponding value is considered 100% in each case. Results shown are the average of at least three replicates with standard deviations below 20% of the given values. Supplementary Table 5 shows the activity for each enzyme at pH 5.5 in nanomoles of *p*-nitrophenol produced per minute per milligram of cell dry weight at 25°C.

469

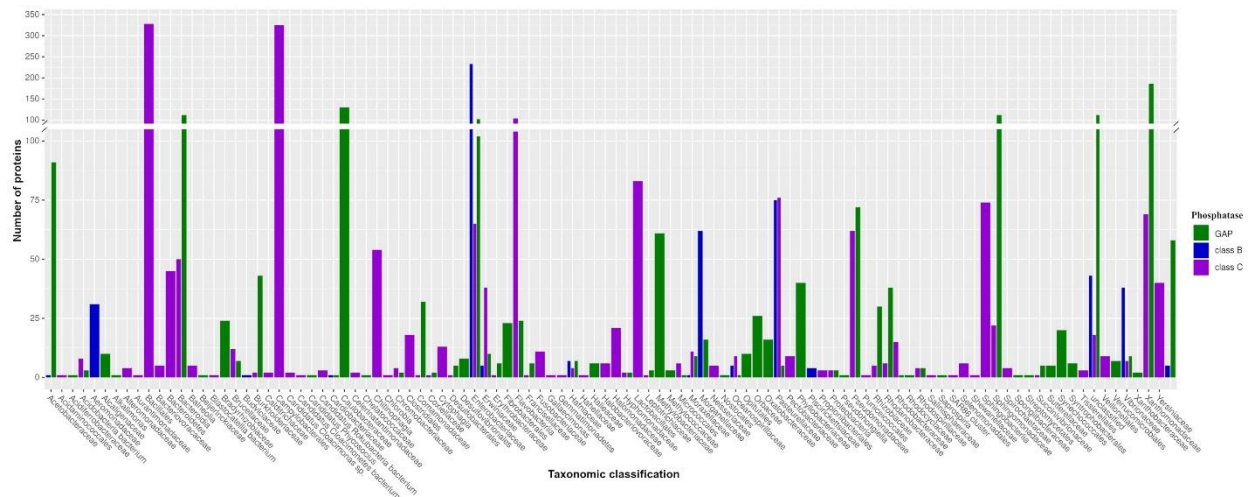**Figure 1. Maximum likelihood phylogenetic tree of bacterial acid phosphatases.** The maximum likelihood tree was inferred from a simultaneous comparison of 3741 protein sequences of bacterial acid phosphatases. Tree topology and branch lengths were calculated by maximum likelihood using the WAG+F+R10 model of evolution for amino acid sequences in lQ-TREE software Nguyen et al., 2015. The tree was rooted by using clade B as an outgroup that shows a clear separation between the three clades of acid phosphatase proteins. Colours of the branches represent levels of significance obtained in the bootstrapping analysis using 1000 bootstrap replications. Green indicates percentages close to 100% of confidence in the bootstraping analysis. The unrooted tree obtained using the same sample set it is shown in Supplementary Figure 1.

484

485

486 **Figure 2. Taxonomic distribution of the number of sequences used to construct the**

487 **three acid phosphatase profiles (Prf-GAP, Prf-B and Prf-C).** Sequences were

488 downloaded from the Uniprot database according to their functional annotation. The

489 number of proteins per taxonomic group were plotted using ggplot2 library in R
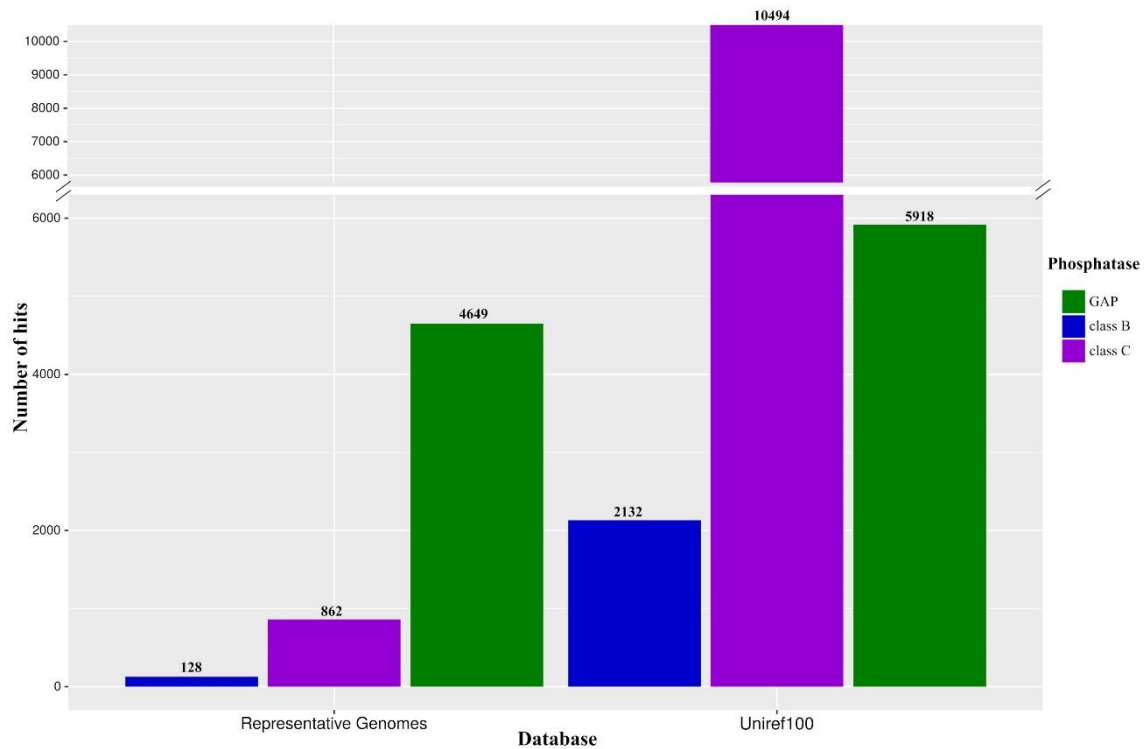
490 (Wickham et al., 2016).

491

492

493

494

495

24

**Figure 3. Number of hits found by the constructed profiles of three classes of acid phosphatases (Prf-GAP, Prf-B and Prf-C)**. Using *pfscan* tool on protein sequences from the Uniref100 database (the last three columns on the right) and a local database of proteomes of 5639 representative bacterial and archaea genomes (the three most left columns) downloaded from NCBI.

**SUPPLEMENTARY FIGURES AND TABLES. They will provide upon request to Juan L. Ramos due to the large size of the files**