

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

## Supporting Information (SI) Appendix

Form III RubisCO-mediated transaldolase variant of the Calvin cycle in a  
chemolithoautotrophic bacterium

Evgenii N. Frolov, Ilya V. Kublanov, Stepan V. Toshchakov, Evgenii A. Lunev, Nikolay V.  
Pimenov, Elizaveta A. Bonch-Osmolovskaya, Alexander V. Lebedinsky, Nikolay A. Chernyh

Corresponding author: Evgenii N. Frolov

E-mail: [evgenii\\_frolov\\_89@mail.ru](mailto:evgenii_frolov_89@mail.ru)

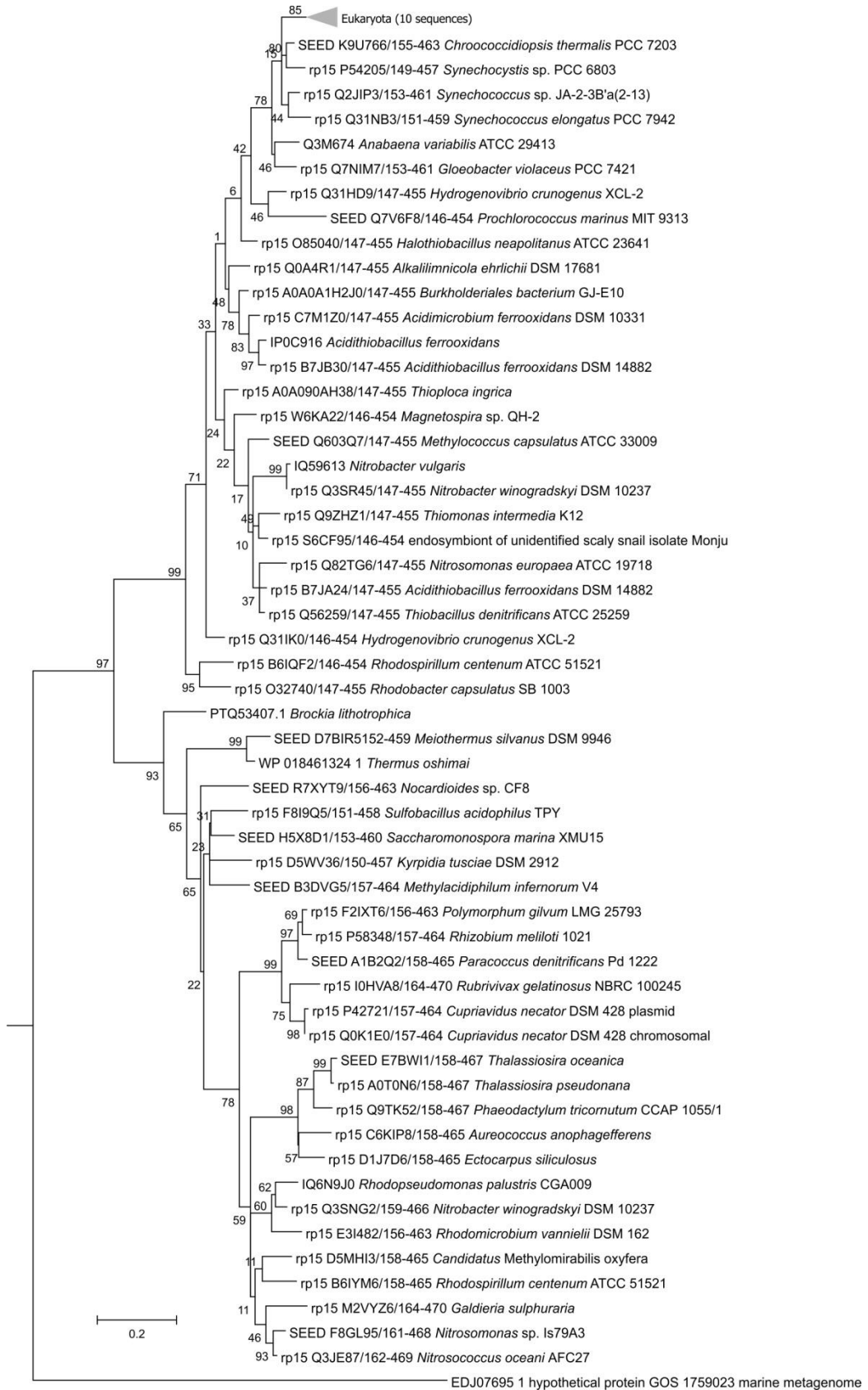
### **This PDF file includes:**

SI Text  
Figs. S1 to S10  
Table S1  
References for SI reference citations

26 **SI Text**

27 **General genomic properties of *T. acidiphilum*.** The genome of *T. acidiphilum* consists of one  
28 circular chromosome of 1,774,794 bp. The genome contains 1757 protein-coding genes, two  
29 identical rRNA operons and 46 tRNA genes. While the 16S rRNA genes of *T. acidiphilum* and  
30 those of *T. narugense* share 99.5% similarity, whole genome-genome ANI and AAI values were  
31 86% and 90%, respectively, which is far below the species border (1). Analysis of the core *in*  
32 *silico* proteome and species-specific proteins of *T. acidiphilum* and *T. narugense* by OthoVenn  
33 server (2) showed that the proteins formed 1665 clusters of homologous genes, while 79 and 127  
34 single copy protein-coding genes as well as 2 and 5 clusters of paralogous genes, were specific  
35 for *T. acidiphilum* and *T. narugense*, respectively. The increased number of species-specific  
36 proteins in *T. narugense* reflects higher genome mobility. Indeed, the *T. narugense* genome is  
37 124 kb larger than that of *T. acidiphilum* and possesses greater number of mobile genetic  
38 elements (9 complete and 4 partial transposase genes in *T. narugense* vs. only 4 partial  
39 transposase genes in *T. acidiphilum* (Fig. S9)). Either higher level of genome mobility or viral  
40 load of *T. narugense*'s environment or both reasons resulted in obviously beneficial evolutionary  
41 acquisition of at least three CRISPR-Cas protein operons of types III-A, III-D and I-B, as well as  
42 corresponding CRISPR repeat clusters. At the same time, neither CRISPR arrays nor functional  
43 CRISPR-Cas gene clusters were discovered in *T. acidiphilum* genome. In turn, the absence of  
44 genomic defense system in *T. acidiphilum* resulted in acquisition of several phage-related genes  
45 (Fig. S9). Search for genomic islands with three different tools (see Supplemental Materials and  
46 Methods), as well as analysis of tetranucleotide frequency biases, produced somewhat diverging  
47 results: none of the predicted HGT regions was univocally supported by all four methods (Figs.  
48 S9-S10). Despite their ambiguity, the HGT signatures observed in the region of the *cbb1* gene  
49 cluster and in 10-20 kb vicinity of the *cbb2* and *cbb3* gene clusters provide an opportunity to  
50 speculate that the *cbb* gene clusters, *cbb1* in particular, may have been acquired by distant lateral

51 transfer but the elapsed evolutionary time was sufficient to ameliorate the gene sequences, at  
52 least to a considerable extent.  
53



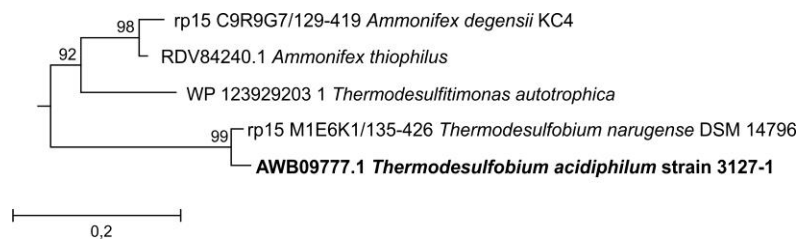
55

56

**Fig. S1.** RubisCO Form I phylogenetic subtree. For a collapsed view of the complete tree,

57

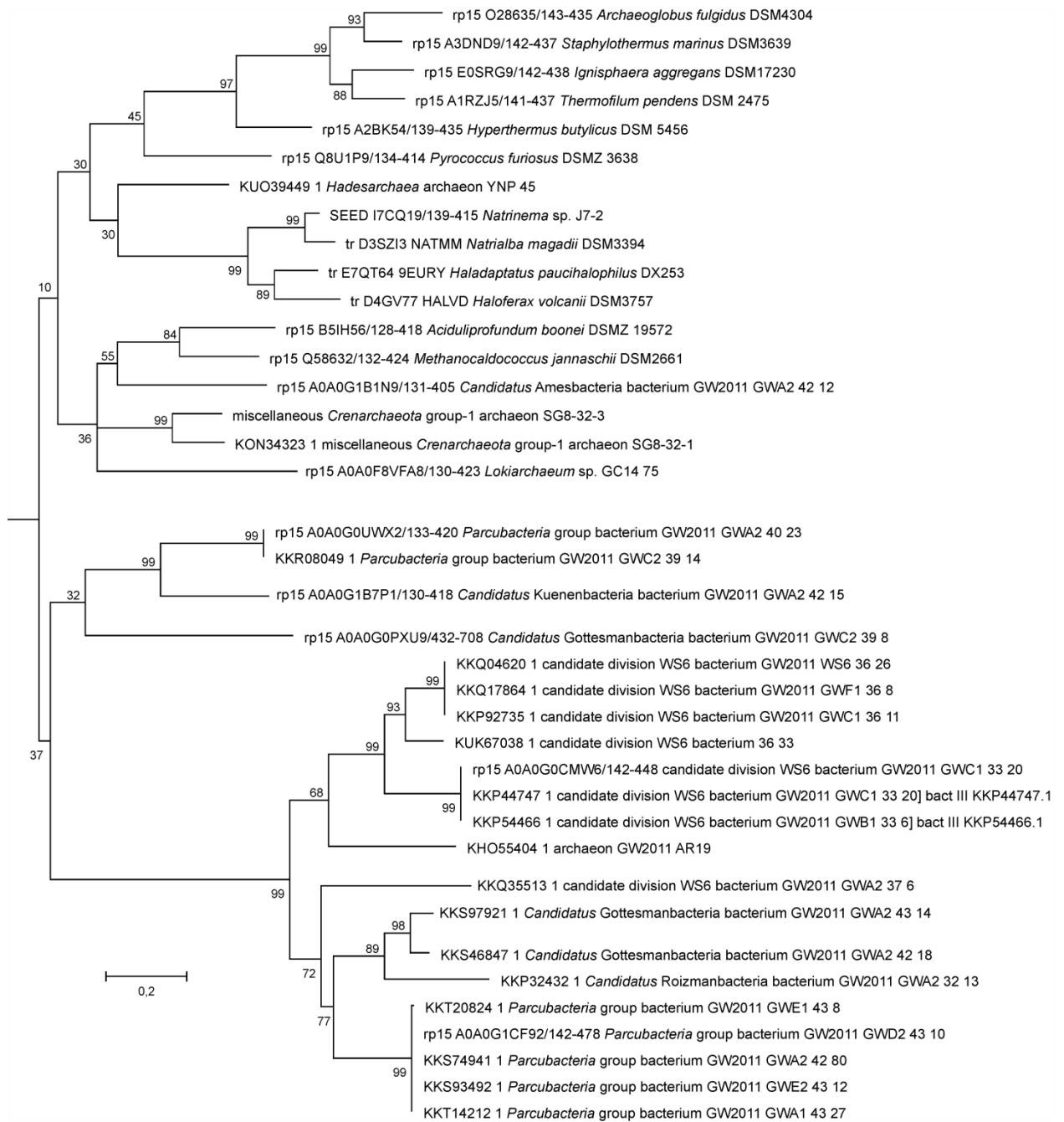
see Fig. 1.



58

59 **Fig. S2.** RubisCO Form III deep branch comprised by *Thermodesulfobium* and

60 *Ammonifex* sequences. For a collapsed view of the complete tree, see Fig. 1.



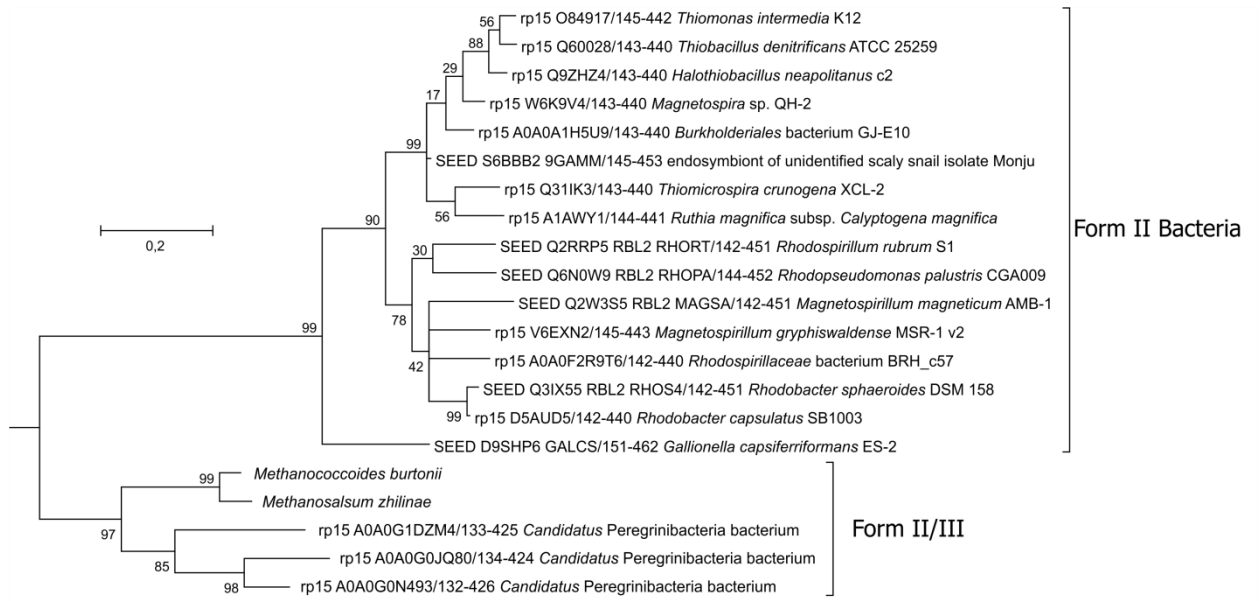
61

62

63 1.

**Fig. S3.** Branch of RubisCO Form III. For a collapsed view of the complete tree, see Fig.

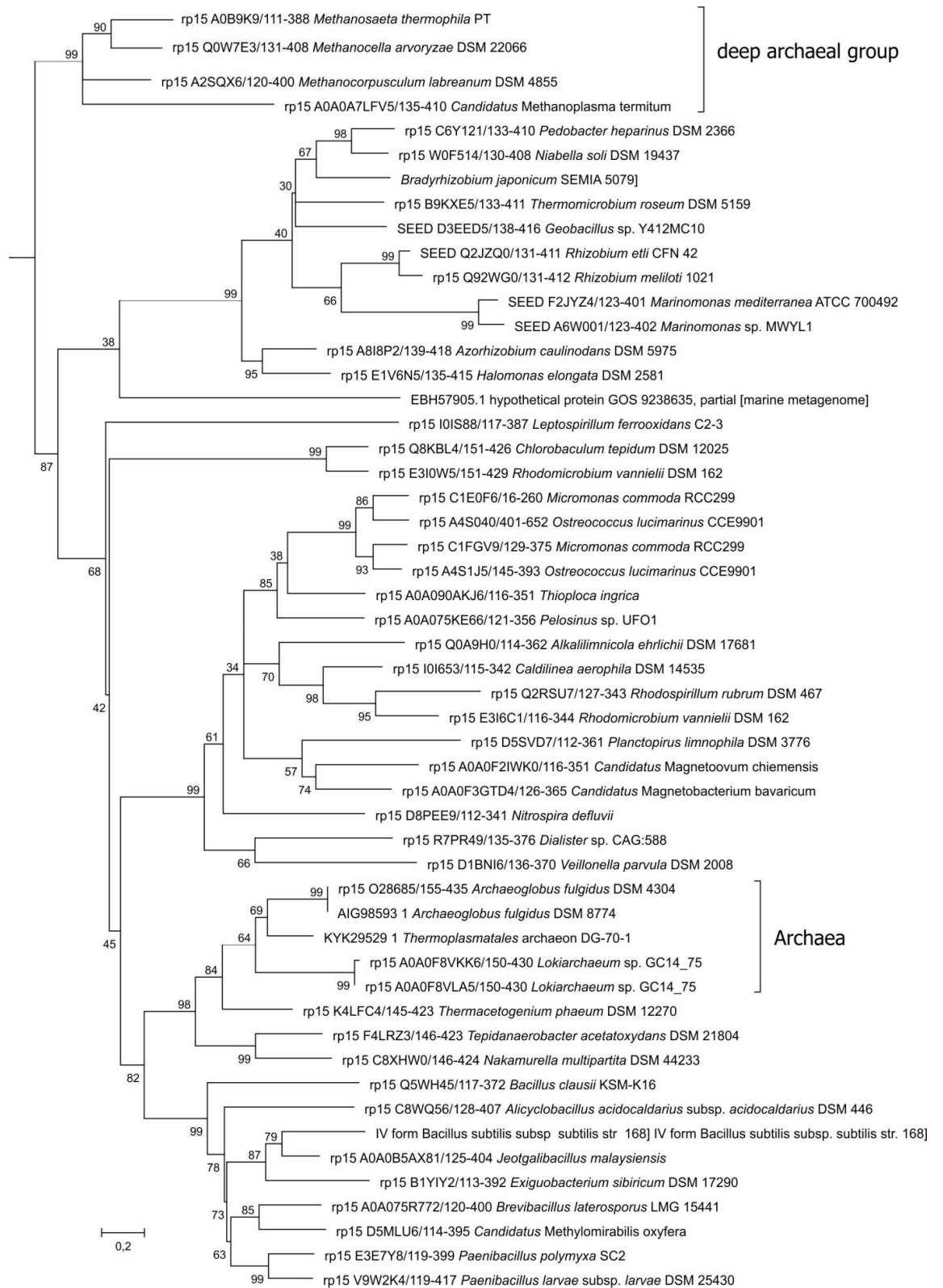
64



65

66

**Fig. S4.** RubisCO Form II. For a collapsed view of the complete tree, see Fig. 1.

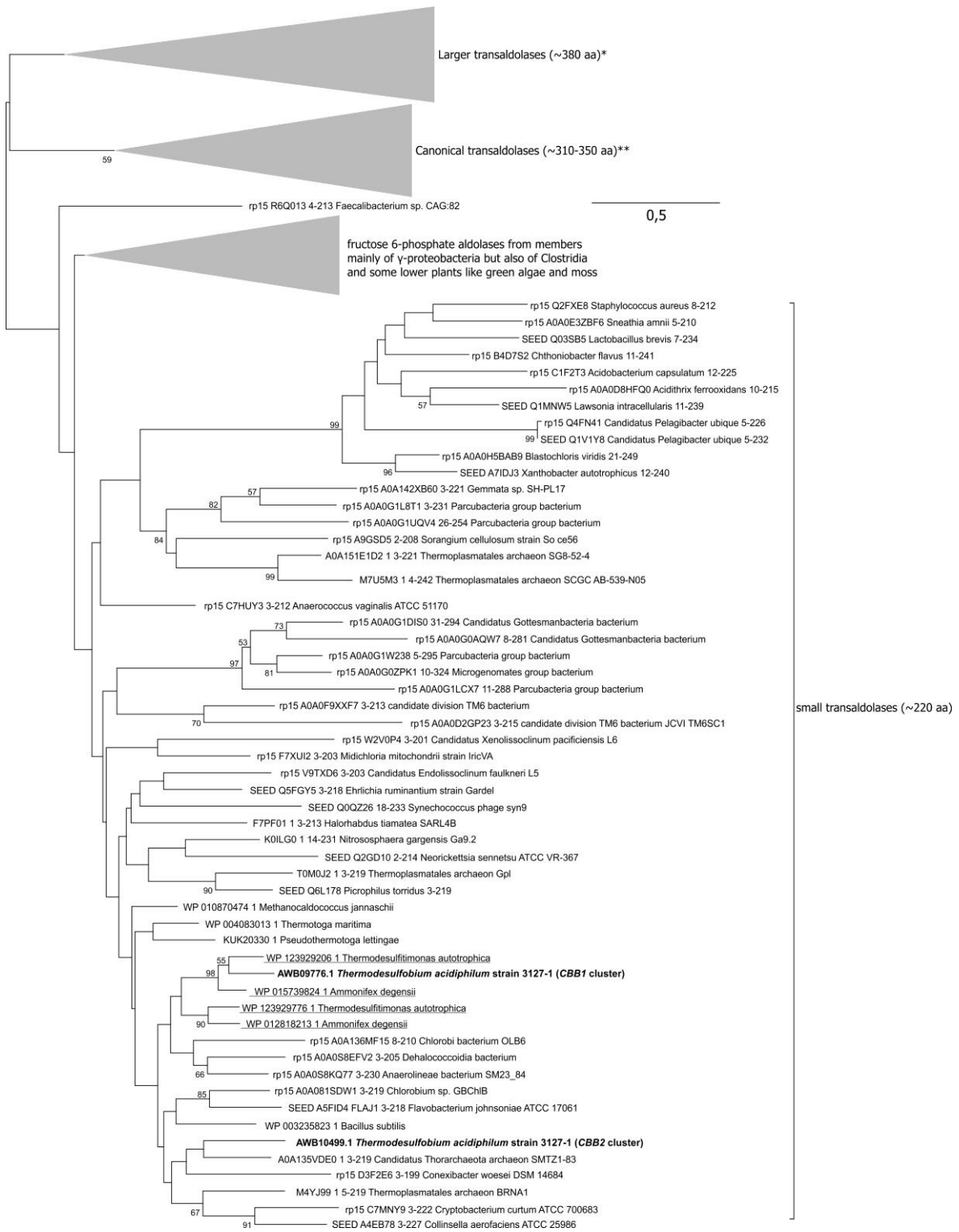


67

68 **Fig. S5.** RubisCO Form IV and deep archaeal branch. For a collapsed view of the

69 complete tree, see Fig. 1.





70

71 **Fig. S6.** Phylogenetic tree of transaldolase family proteins. *T. acidiphilum* proteins are in bold.





144 SSCH\_790022 is from *Syntrophaceticus schinkii*; fourth best blastp hit of Adeg\_1859 in NCBI nr among cultivated microorganisms (fourth  
145 best hit after Adeg\_0665, DXX99\_02705 and DXX99\_08975).

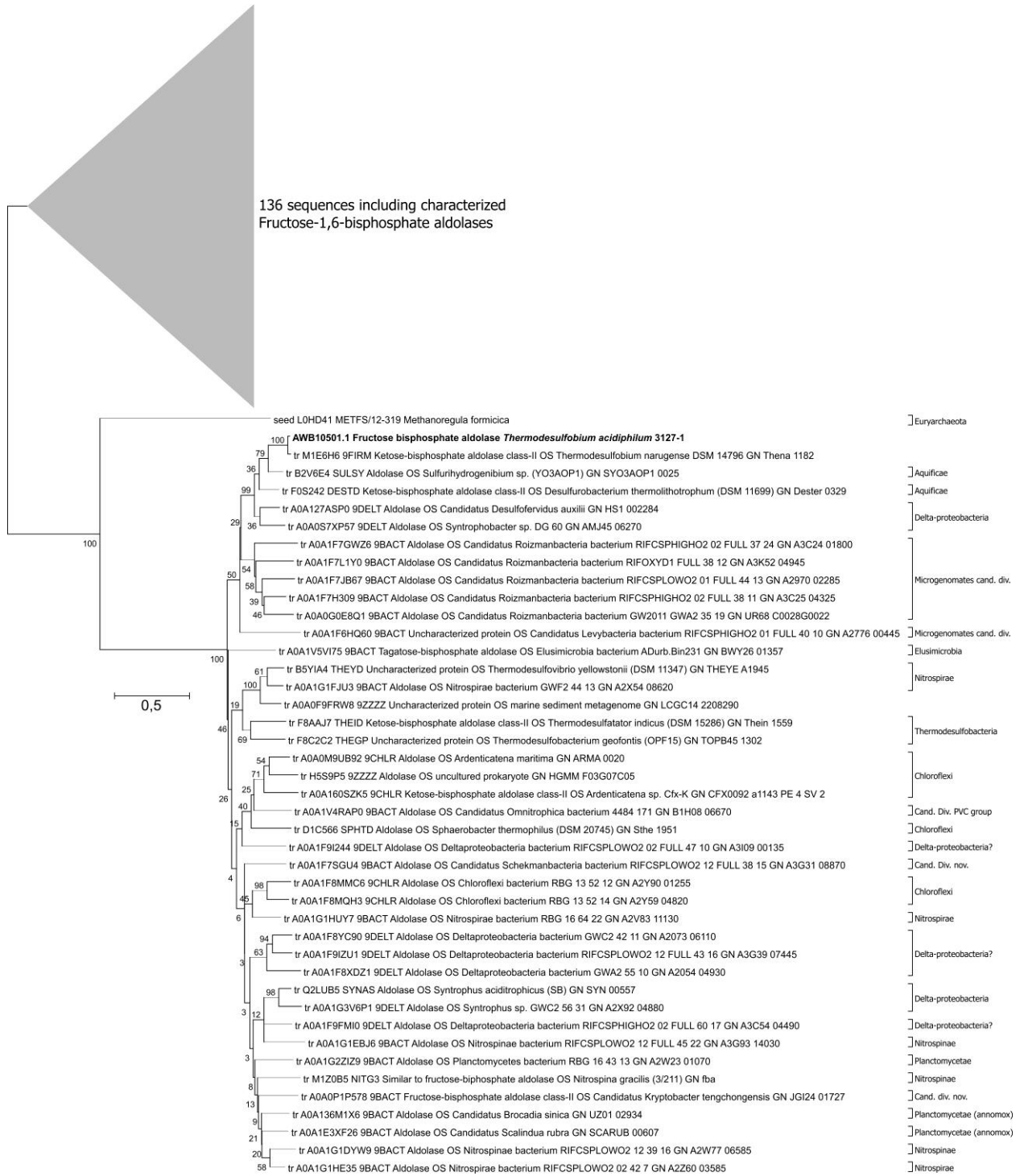
146 The other sequences belong to fructose-1,6-bisphosphate aldolases/phosphatases studied and/or discussed in the works on the FBPAP active site (3, 4):

147 Igni\_0363 is from *Ignicoccus hospitalis*;

148 CENSYA\_RS02585 is from *Cenarchaeum symbiosum*;

149 STK\_03180 is from *Sulfurisphaera tokodaii* (formerly, *Sulfolobus tokodaii*);

150 Tneu\_0133 is from *Pyrobaculum neutrophilum* (formerly, *Thermoproteus neutrophilus*).



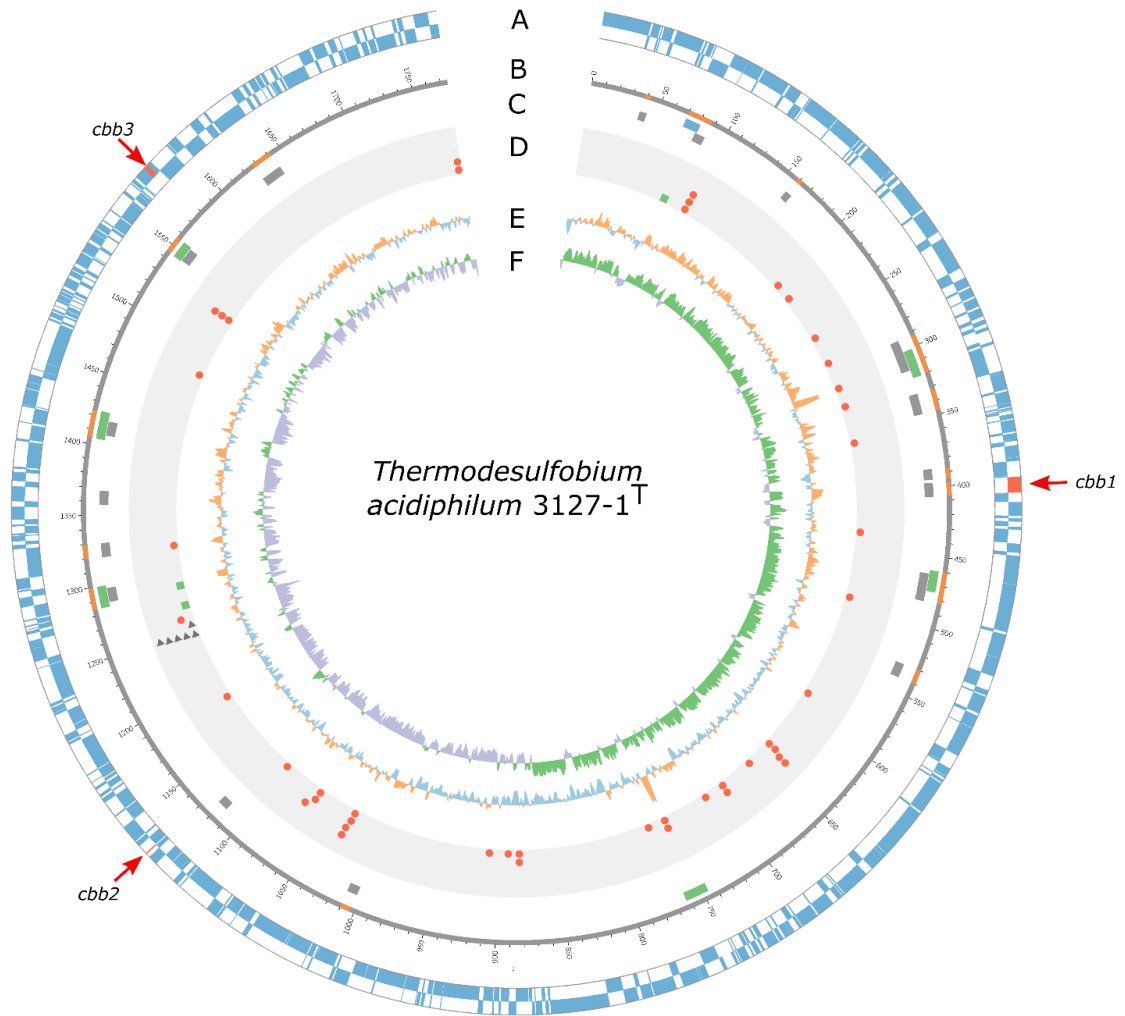
152

153

**Fig. S8.** Phylogenetic tree of class II aldolase proteins. The *T. acidiphilum* protein

154

(TDSAC\_1156) is in bold.



155

156 **Fig. S9.** Chromosome map of *T.acidophilum*.

157 (A) Predicted CDS of *T. acidophilum* genome. Positive strand CDS are shown on the  
 158 outer circle, negative strand CDS are on the inner circle. *cbb* gene clusters are highlighted with  
 159 red;

160 (B) Chromosome coordinates; regions corresponding to predicted genomic islands are  
 161 highlighted with orange;

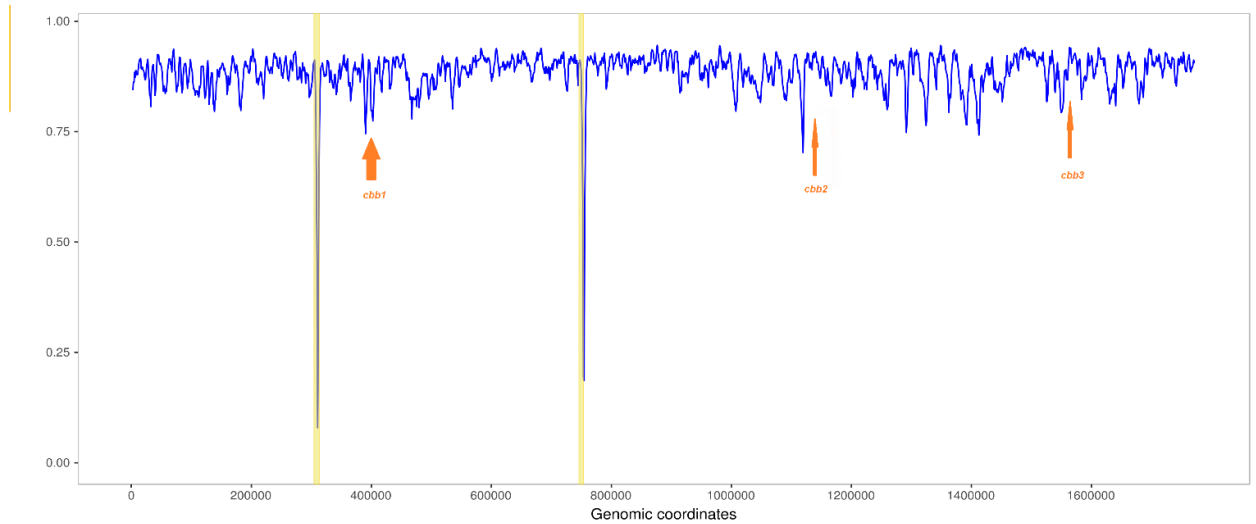
162 (C) Regions corresponding to genomic islands predicted by IslandViewer (blue),  
 163 SeqWordSniffer (green) and AlienHunter (gray);

164 (D) Genomic island-associated genomic features: tRNA (orange circles), transposase  
 165 genes and pseudogenes (green rectangles), phage-related proteins (gray triangles);

166 (E) GC-bias. Estimated in a window of 2000 nt with the step size of 200 nt. Values are  
167 shown relative to average GC-content of *T. acidiphilum* genome. Values inferior to the average  
168 GC are highlighted by light blue, values superior to average GC are in light orange.

169 (F) GC-skew. Estimated in a window of 2000 nt with step size of 200 nt. Negative values  
170 are highlighted by purple, positive values are highlighted by green.

171



172

173

174

175

176

**Fig. S10.** Tetranucleotide frequency bias of *T.acidiphilum* genome. Correlation of tetranucleotide frequencies against genome-wide tetranucleotide signature. The positions of key CBB cycle gene clusters are indicated by arrows. Regions with significant deviation of tetranucleotide patterns are highlighted by yellow rectangles.



**Table S1.** Predicted genes of *cbb* gene clusters in *T. acidiphilum* genome

Locus tags	Gene name	Predicted function	Abbreviations for enzymes	Best BlastP hit*	% identity	Score	E-value	2 <sup>nd</sup> BlastP hit*	% identity	Score	E-value	COGs and Pfam domains
<b><i>cbb1</i> gene cluster</b>												
TDSAC_0400	<i>tal1</i>	Transaldolase [EC 2.2.1.2]	TAL1	<i>T. narugense</i> (Thena_0422)	97	432	2e-153	<i>Ammonifex degensii</i> (Adeg_1864)	70	321	6e-109	COG0176, pfam00923
TDSAC_0401	<i>cbbL-III</i>	Form III ribulose bisphosphate carboxylase [EC 4.1.1.39]	RubisCO	<i>T. narugense</i> (Thena_0423)	95	845	0.0	<i>A. degensii</i> (Adeg_1863)	68	578	0.0	COG1850, pfam00016, pfam02788
TDSAC_0402	<i>cbbP</i>	Phosphoribulokinase [EC 2.7.1.19]	PRK	<i>T. narugense</i> (Thena_0424)	97	601	0.0	<i>Desulfotomaculum putei</i> (BUB67_RS15390)	50	295	2e-96	COG0572, pfam00485
TDSAC_0403	<i>cbbT-n</i>	Transketolase, N-terminal section [EC 2.2.1.1]	TK-N	<i>T. narugense</i> (Thena_0425)	97	588	0.0	<i>A. degensii</i> (Adeg_1861)	72	425	3e-148	COG3959, pfam00456
TDSAC_0404	<i>cbbT-c</i>	Transketolase, C-terminal section [EC 2.2.1.1]	TK-C	<i>T. narugense</i> (Thena_0426)	94	676	0.0	<i>A. degensii</i> (Adeg_1860)	62	418	1e-143	COG3958, pfam02779, pfam02780
TDSAC_0405	<i>cbbG1</i>	NAD-dependent glyceraldehyde-3- phosphate dehydrogenase [EC 1.2.1.12]	GAPDH1	<i>T. narugense</i> (Thena_0427)	91	629	0.0	<i>T. narugense</i> (Thena_1627)	62	420	2e-144	COG0057, pfam02800 pfam00044
TDSAC_0406	<i>cbbK1</i>	Phosphoglycerate kinase [EC 2.7.2.3]	PGK1	<i>T. narugense</i> (Thena_0428)	95	784	0.0	<i>T. narugense</i> (Thena_1628)	65	558	0.0	COG0126, pfam00162
TDSAC_0407	<i>cbbF1</i>	Fructose-1,6- bisphosphatase, type I [EC 3.1.3.11]	FBPase1	<i>T. narugense</i> (Thena_0429)	93	621	0.0	<i>Thermodesulfobac- terium hydrogeniphilum</i> (CC87_RS03320)	68	474	3e-166	COG0158, pfam00316
TDSAC_0408	<i>cbbI</i>	Ribose 5-phosphate isomerase B [EC 5.3.1.6]	RPI	<i>T. narugense</i> (Thena_0430)	97	291	1e-99	<i>Clostridiales bacterium</i> (AYC61_RS17895)	62	176	5e-54	COG0698, pfam02502
TDSAC_0409	<i>cbbE</i>	Ribulose-phosphate 3- epimerase [EC 5.1.3.1]	RuPE	<i>T. narugense</i> (Thena_0431)	97	413	6e-146	<i>Campylobacter peloridis</i> (CPEL_1132)	44	167	1e-48	COG0036, pfam00834
<b><i>cbb2</i> gene cluster</b>												
TDSAC_1154	<i>tal2</i>	Transaldolase [EC 2.2.1.2]	TAL2	<i>T. narugense</i> (Thena_1180)	91	400	2e-140	<i>Pseudothermotoga lettingae</i> (Tlet_1124)	55	251	2e-81	COG0176, pfam00923
TDSAC_1155	<i>cbbF2</i>	Fructose-1,6- bisphosphatase, type I [EC 3.1.3.11]	FBPase2	<i>T. narugense</i> (Thena_1181)	97	648	0.0	bacterium BMS3Bbin07 (BMS3Bbin07_00863)	69	474	3e-166	COG0158, pfam00316
TDSAC_1156	<i>cbbA</i>	Aldolase II class [EC 4.1.2.13]	FBPA	<i>T. narugense</i> (Thena_1182)	94	923	0.0	<i>Sulfurihydrogenibium</i> sp. YO3AOP1 (SYO3AOP1_0025)	69	650	0.0	COG0191, pfam01116, pfam00596

cbb3 gene cluster												
TDSAC_1598	<i>cbbG3</i>	NAD-dependent glyceraldehyde-3- phosphate dehydrogenase [EC 1.2.1.12]	GAPDH3	<i>T. narugense</i> (Thena_1627)	97	669	0.0	<i>T. narugense</i> (Thena_0427)	63	427	3e-147	COG0057, pfam02800 pfam00044
TDSAC_1599	<i>cbbK3</i>	Phosphoglycerate kinase [EC 2.7.2.3]	PGK3	<i>T. narugense</i> (Thena_1628)	99	796	0.0	<i>T. narugense</i> (Thena_0428)	98	559	0.0	COG0126, pfam00162
TDSAC_1600	<i>tpi</i>	Triosephosphate isomerase [EC 5.3.1.1]	TPI	<i>T. narugense</i> (Thena_1629)	93	461	1e-155	<i>Clostridium sp.</i> (HMPREF1092_01812)	48	189	2e-56	COG0149, pfam00121

178

179 \*When analyzing results of NCBI blasp searches, hits to sequences from the metagenomic assemblies PNIV01000000 (*Fervidicoccus fontis* ARK-12)

180 and PNIY01000000 (*Thermodesulfobium narugense* ARK-09) were neglected, because these are evidently misassemblies. The former one contains

181 962 (66.5%) genes 99-100% identical to genes of *Fervidicoccus fontis* Kam940<sup>T</sup> and 175 (12% genes) 99-100% identical to genes of *T. acidiphilum*

182 3127-1<sup>T</sup>. In the cases we checked, the *Fervidicoccus fontis* ARK-12 genes related to *Thermodesulfobium* genes were in short contigs of the

183 missassembly, and these contigs lacked any genes related to *Fervidicoccus fontis* Kam940<sup>T</sup> genes. On the other hand, the latter above-mentioned

184 assembly PNIY01000000 (*T. narugense* ARK-09), originating from the same metagenome, missed some of its native genes, attributed to *F. fontis*

185 ARK-12.

## References

1. E. N. Frolov *et al.*, *Thermodesulfobium acidiphilum* sp. nov., a new thermoacidophilic sulfate-reducing chemoautotrophic bacterium from a Kamchatkan thermal site. *Int J Syst Evol Microbiol* **67**, 1482-1485 (2017).
2. Y. Wang, D. Coleman-Derr, G. Chen, Y. Q. Gu, OrthoVenn: a web server for genome wide comparison and annotation of orthologous clusters across multiple species. *Nucleic Acids Res* **43(W1)**, W78-W84 (2015).
3. Say, R. F. & Fuchs, G. Fructose 1,6-bisphosphate aldolase/phosphatase may be an ancestral gluconeogenic enzyme. *Nature* **464**, 1077–1081 (2010).
4. J. Du, F. R. Say, W. Lü, G. Fuchs, O. Einsle, O. Active-site remodelling in the bifunctional fructose-1,6-bisphosphate aldolase/phosphatase. *Nature* **478**, 534-537 (2011)