

# Lawrence Berkeley National Laboratory

## LBL Publications

### Title

A predicted CRISPR-mediated symbiosis between uncultivated archaea

### Permalink

<https://escholarship.org/uc/item/3q2101xs>

### Journal

Nature Microbiology, 8(9)

### ISSN

2058-5276

### Authors

Esser, Sarah P

Rahlff, Janina

Zhao, Weishu

et al.

### Publication Date

2023-09-01

### DOI

10.1038/s41564-023-01439-2

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution-NonCommercial License, available at <https://creativecommons.org/licenses/by-nc/4.0/>

Peer reviewed

## A predicted CRISPR-mediated symbiosis between uncultivated archaea

Sarah P. Esser<sup>1,2,\*</sup>, Janina Rahlff<sup>2,\*,+</sup>, Weishu Zhao<sup>3,++</sup>, Michael Predl<sup>4,5</sup>, Julia Plewka<sup>1,2</sup>, Katharina Sures<sup>1,2</sup>, Franziska Wimmer<sup>6</sup>, Janey Lee<sup>7</sup>, Panagiotis S. Adam<sup>2</sup>, Julia McGonigle<sup>8</sup>, Victoria Turzynski<sup>1,2</sup>, Indra Banas<sup>1,2</sup>, Katrin Schwank<sup>2,+++</sup>, Mart Krupovic<sup>9</sup>, Till L. V. Bornemann<sup>1,2</sup>, Perla Abigail Figueroa-Gonzalez<sup>1,2</sup>, Jessica Jarett<sup>7</sup>, Thomas Rattei<sup>4,5</sup>, Yuki Amano<sup>10</sup>, Ian K. Blaby<sup>7</sup>, Janfang Cheng<sup>7</sup>, William J. Brazelton<sup>8</sup>, Chase L. Beisel<sup>6,11</sup>, Tanja Woyke<sup>7</sup>, Ying Zhang<sup>3</sup>, and Alexander J. Probst<sup>1,2,12,13,#</sup>

<sup>1</sup>Environmental Metagenomics, Research Center One Health of the University Alliance Ruhr, Faculty of Chemistry, University of Duisburg-Essen, 45151 Essen

<sup>2</sup>Group for Aquatic Microbial Ecology, Environmental Microbiology and Biotechnology, University of Duisburg-Essen, 45141 Essen

<sup>3</sup>Department of Cell and Molecular Biology, College of the Environment and Life Sciences, University of Rhode Island, Kingston, RI 02881, USA

<sup>4</sup>Computational Systems Biology, Centre for Microbiology and Environmental Systems Science, University of Vienna, 1030 Vienna, Austria

<sup>5</sup>Doctoral School in Microbiology and Environmental Science, University of Vienna, 1030 Vienna, Austria

<sup>6</sup>Helmholtz Institute for RNA-based Infection Research (HIRI), Helmholtz-Centre for Infection Research (HZI), 97080 Würzburg, Germany

<sup>7</sup>DOE Joint Genome Institute, Lawrence Berkeley National Laboratory, One Cyclotron Rd, Berkeley, CA, 94720, USA

<sup>8</sup>School of Biological Sciences, University of Utah, Salt Lake City, Utah, USA

<sup>9</sup>Institut Pasteur, Université Paris Cité, CNRS UMR6047, Archaeal Virology Unit, 75015 Paris, France

<sup>10</sup>Nuclear Fuel Cycle Engineering Laboratories, Japan Atomic Energy Agency, Tokai, Ibaraki, 319-1194 Japan

<sup>11</sup>Medical faculty, University of Würzburg, 97080 Würzburg Germany

30 <sup>12</sup>Centre of Water and Environmental Research (ZWU), University of Duisburg-Essen, 45141  
31 Essen, Germany

32 <sup>13</sup>Centre of Medical Biotechnology (ZMB), University of Duisburg-Essen, 45141 Essen, Germany

33

34 \*authors contributed equally

35 †present address: Centre for Ecology and Evolution in Microbial Model Systems (EEMiS),  
36 Department of Biology and Environmental Science, Linnaeus University, Kalmar, Sweden

37 ††present address: Shanghai Jiao Tong University, School of Life Sciences and Biotechnology,  
38 International Center for Deep Life Investigation (IC-DLI), Shanghai Jiao Tong University, Shanghai,  
39 200240, China

40 †††present address: University of Regensburg, Biochemistry III, 93053 Regensburg, Germany

41

42 #corresponding author: alexander.probst@uni-due.de

43

44 **ABSTRACT**

45 CRISPR-Cas systems defend prokaryotic cells from invasive DNA of viruses, plasmids, and other  
46 mobile genetic elements. Here we show using metagenomics, metatranscriptomics and single-  
47 cell genomics that CRISPR systems of widespread, uncultivated archaea can also target  
48 chromosomal DNA of archaeal episymbionts of the DPANN superphylum. Using meta-omics  
49 datasets from Crystal Geyser (USA) and Horonobe Underground Research Laboratory (Japan), we  
50 find that CRISPR spacers of the hosts *Candidatus* Altiarchaeum crystalense and *Ca.* A.  
51 horonobense, respectively, match putative essential genes in their episymbionts' genomes of the  
52 genus *Ca.* Huberiarchoaeum, and that some of these spacers are expressed *in situ*. Metabolic  
53 interaction modeling also reveals complementation between host-episymbiont systems, based  
54 on which we propose that episymbionts are either parasitic or mutualistic depending on the  
55 genotype of the host. By expanding our analysis to 7,012 archaeal genomes, we suggest that  
56 CRISPR-Cas targeting of genomes associated with symbiotic archaea evolved independently in  
57 various archaeal lineages.

58

59

60 **INTRODUCTION**

61 Clustered regularly interspaced short palindromic repeats associated systems (CRISPR-Cas)  
62 facilitate adaptive prokaryotic immunity via cleavage of mobile genetic elements (MGEs), e.g.,  
63 viruses and plasmids<sup>1</sup>. CRISPR *loci* consist of a series of direct repeat (DR) sequences interspaced  
64 by short variable fragments, *i.e.*, spacers, flanked by *cas* genes. Upon exposure to novel MGEs,  
65 short DNA fragments from these invaders are incorporated into the CRISPR array as spacers. The  
66 spacers are then used as templates to form CRISPR RNAs (crRNAs) that guide effector Cas  
67 nucleases to complementary nucleic acid sequences. Spacer sequences can also be used to study  
68 infection histories *in silico* based on matches to protospacers, corresponding nucleic acid regions  
69 in the MGE<sup>2</sup>.

70 CRISPR systems exhibit remarkable diversity and functional plasticity including roles in  
71 non-defensive functions (reviewed by Koonin & Makarova 2022<sup>3</sup>). Six main types of CRISPR-Cas  
72 systems have been described, including different subtypes, *e.g.*, type I-A to I-F, depending on

73 signature genes and their arrangements<sup>4,5</sup>. Target identification in type I and II systems is  
74 dependent on the recognition of a short protospacer-adjacent motif (PAM) in the target DNA  
75 sequence, which elicits cleavage and clearance of the MGE's protospacer. Rather than relying on  
76 a defined PAM for target recognition, other CRISPR systems (*e.g.* type III) generally evaluate the  
77 extent of hybridization between the flanking portions of the crRNA (called protospacer-flanking  
78 sequence) and the target<sup>6,7</sup>. While CRISPR-Cas systems are widely distributed, they are more  
79 common in archaea (in ~85% of genomes) than in bacteria (in ~40% of genomes; reviewed by  
80 Makarova et al. 2020)<sup>5</sup>.

81         Branching from the archaeal tree of life, the DPANN superphylum including *i.e.*,  
82 Diapherotrites, Parvarchaeota, Aenigmarchaeota, Nanoarchaeota and Nanohaloarchaeota, and  
83 several other recently proposed phyla<sup>8,9</sup>, comprises a vast collection of microorganisms  
84 remarkably small in size and enigmatic due to the scarcity of cultivated representatives<sup>10,11</sup>.  
85 Insights into the physiological characteristics of DPANN archaea arise primarily from detailed  
86 analyses of co-cultivation with amenable microorganisms<sup>12,13</sup> and/or imaging of environmental  
87 samples<sup>14</sup>. These inferences, along with the limited metabolic potential contained in their  
88 comparatively small genomes, suggest that most DPANN archaea exist as (epi-)symbionts of  
89 other archaea<sup>13,15–18</sup> or even as intracellular symbionts<sup>19</sup>. The independent and autotrophic  
90 *Candidatus* Altiarchaeum sp. is host organism to another uncultivated DPANN archaeon,  
91 *Candidatus* Huberiarchaeum crystalense<sup>14,20</sup>.

92         Previous evidence suggested that certain DPANN archaea can fuse their cytoplasm with  
93 that of their hosts<sup>14,15,21–23</sup> and even exchange enzymes<sup>21</sup>. Cytoplasm fusion could in theory  
94 facilitate transfer of metabolites from the host to the symbiont consequently rendering such a  
95 symbiosis potentially parasitic. Hence, we investigated the symbiotic nature of the uncultivated  
96 DPANN *Ca.* Altiarchaeum and its uncultivated DPANN episymbiont *Ca.* Huberiarchaeum using  
97 meta-omics and metabolic modeling in two independent subsurface ecosystems (Fig. 1A). Based  
98 on the complex interaction of *Ca.* Altiarchaea and viruses in deep subsurface ecosystems, we  
99 examined the encoded CRISPR systems and analyzed the targets of their respective spacer  
100 populations in two independent subsurface ecosystems (accessible through Crystal Geyser (CG),  
101 Utah, USA and Horonobe Underground Research Laboratory (HURL), Hokkaido, Japan), where we

102 identified both the host, and the symbiont being associated based on fluorescence *in situ*  
103 hybridization (FISH). Our findings demonstrate that a substantial portion of the host spacer  
104 population is targeted towards the genomes of the episymbionts having the same PAM sequence  
105 as the respectively targeted viruses in the ecosystems. In addition, the host CRISPR systems also  
106 target the host chromosome, based on which genome-centric metabolic modeling predicted a  
107 mutualistic relationship between host and episymbiont as a function of metabolic  
108 complementation. Based on our results, we suggest that CRISPR-Cas systems play an integral role  
109 in mediating archaeal host-DPANN interactions.

110

## 111 **RESULTS**

### 112 ***A CRISPR-Cas system targets archaeal episymbiont genomes***

113 Two subsurface ecosystems separated by 8,255 km (Fig. 1A) and derived from different  
114 geological formations<sup>20,24</sup>, *i.e.*, a Wingate Sandstone-hosted aquifer of the Colorado Plateau at  
115 ~ 350 m depth (the CG ecosystem)<sup>20,25–27</sup> and a diatomaceous/siliceous mudstone-hosted aquifer  
116 of the HURL<sup>24</sup> at ~ 250 m depth, were dominated by two species of *Ca. Altiarchaea* (up to 24.5%  
117 and 51.6% of the community, respectively). We show their association with cells of *Ca.*  
118 *Huberiarchaea*, their DPANN episymbiont, using species-specific FISH (Fig. 1A). The *Ca.*  
119 *Altiarchaea* genomes retrieved from CG and HURL were shown to encode a I-B CRISPR system  
120 and an abundant CRISPR array, which could not be assigned to a specific *cas* gene cassette, as has  
121 been reported for other *Ca. Altiarchaea* species<sup>25,28</sup>. Confidence in assigning the CRISPR-Cas  
122 system to its correct metagenome-assembled genome (MAG; *Altiarchaea* genomes n=1; Table  
123 S1-S2) derives from the exceedingly high abundance of *Ca. Altiarchaea* genome fragments in the  
124 CG samples (Fig. S1). In addition, within 219 single-cell amplified genomes (SAGs; *Altiarchaea*  
125 SAGs n=7; Table S3) from CG, only *Ca. Altiarchaea* bore the corresponding consensus DR  
126 sequence (see Extended Data Fig. 1 for additional correlation-based evidence), which were  
127 remarkably well-conserved across ecosystems<sup>28</sup> (Fig. 1B).

128 Analyses of spacers from *Ca. Altiarchaeum crystalense* detected in 66 CG metagenomes  
129 over six years of surveillance (1.07 Tbps of sequencing data, Table S1) revealed 297,531 distinct  
130 spacer clusters (Fig. 1B), indicative of a complex CRISPR spacer repertoire system for this

131 organism (Figs. S2 and S3). Within these metagenomes, CRISPR type I-B spacers matched the  
132 protospacers of 64 viral DNA sequences corresponding to 14 distinct viral genus clusters (Fig. 1C,  
133 Table S6, Extended Data Fig. 2, Fig. S4-S6, details in Supplementary Results). The PAM sequence  
134 5'-TTN-3' was identified on viral targets matched by type I-B spacers (Fig. S7). However, we were  
135 unable to experimentally confirm this PAM using an established PAM assay<sup>29</sup> or to assess GFP  
136 repression<sup>30</sup> using the 5'-TTN-3' PAM in a cell-free transcription-translation (TXTL) system<sup>29</sup>, likely  
137 due to the divergent settings (including temperature) of the host environment compared to  
138 those used in the established assay<sup>29</sup>.

139         The finding that all virus-matched spacers detected in the exhaustive CG survey derived  
140 from the CRISPR I-B system and the ubiquitous nature of this system in *Ca.* Altiarchaea  
141 worldwide<sup>28</sup> suggests that the I-B system serves as a primary line of defense against viruses  
142 infecting these archaea. A substantial fraction of the spacers matched microbial genomes,  
143 including those of *Ca. A. crystalense*, *i.e.*, its own genome (self-targeting, up to 2.9% in sample  
144 CG16) and of its episymbiont *Ca. Huberiarchaeum crystalense* (up to 2.8% in sample CG08; Fig.  
145 1C-1D, Fig. 2A-D). The relative proportion of spacers matching the episymbiont was greater than  
146 that matching the host genome (Fig. 1C-D), indicative of biased acquisition, negative selection of  
147 self-targeting spacers or a positive selection for spacers from the episymbiont genome. The  
148 positions of these spacers in the CRISPR I-B array encoded in an altiarchaeal SAG suggest that  
149 these spacers prevailed in the system for extended periods (Fig. 2D). While 17% of the  
150 protospacers self-targeted through the I-B system showed a significant decrease in metagenomic  
151 coverage compared to untargeted scaffold regions (bootstrapped Wilcoxon paired signed rank  
152 test, target sites = 196, FDR-corrected p-value < 0.05), 30% of the I-B protospacers in *Ca. H.*  
153 *crystalense* genomes showed a significant drop in coverage, suggesting *in situ* targeting of the  
154 episymbiont in CG (bootstrapped Wilcoxon signed rank test, target sites = 73, FDR-corrected p-  
155 value < 0.05; Supplementary Results, Extended Data Fig. 3). The coverage of the significantly  
156 different targeted regions, compared to the average coverage of the scaffold decreases in *Ca. A.*  
157 *crystalense* and *Ca. H. crystalense* by 10.74% (median) and 36.99% (median), respectively  
158 (bootstrapped Wilcoxon signed rank test, n=990, FDR-corrected p-value < 0.05; details in  
159 Supplementary Results, Table S5). Supporting this difference, the PAM sequence detected next

160 to the protospacers in *Ca. H. crystalense* was identical to that of the virus-targeting PAM  
161 sequence (Fig. S7). Coverage drops as observed herein could also arise from misassemblies,  
162 regions excised in subpopulations, or elevated SNPs resulting in low recruitment of reads.

163 In contrast to the conserved PAM in the episymbiont and the viruses, the self-targeted  
164 protospacer regions were not associated with the 5'-TTN-3' PAM (Fig. S7). As shown for other  
165 microbial communities, self-targeting can result in cell suicide (reviewed in Heussler and O'Toole,  
166 2016)<sup>31</sup> or transcriptional regulation<sup>32</sup> of genes influencing the fitness of the microbial population  
167 and can thus reduce the strain variation within an ecosystem<sup>33</sup>. However, the lack of the PAM,  
168 the essential motif for successful targeting of DNA by CRISPR system type I (reviewed in Bhaya *et*  
169 *al.* 2011)<sup>34</sup> in the population genomes of *Ca. Altiarchaea*, might on the one hand prevent  
170 subpopulations of *Ca. Altiarchaea* from cell death by autoimmunity. On the other hand, the  
171 correct PAM could still lead to cell death in subpopulations, given that the PAM has not been  
172 silenced by mutations. Based on the overall results from metagenomics and metatranscriptomics  
173 we suggest that CRISPR-Cas systems may function similarly against viral DNA and chromosomal  
174 DNA of episymbionts.

175 Previous investigations, which were based on either species-specific FISH or electron  
176 microscopy, indicate that many DPANN archaea (including *Ca. A. crystalense* and its episymbiont)  
177 fuse their cytoplasms<sup>14,15,21-23</sup>. This direct interaction of the host's and the symbiont's  
178 cytoplasms, and a potentially predatory nature of the symbiont<sup>14,20</sup>, likely underlie the evolution  
179 of a direct assault on the episymbiont's genome by the *Altiarchaea* CRISPR system (Fig. 1A). To  
180 this end, we annotated genes of *Ca. H. crystalense* targeted by *Ca. A. crystalense*'s CRISPR type  
181 I-B system and identified several hypothetical proteins, proteins lacking annotation, and non-  
182 coding genomic regions (these categories sum up to 98.25%). Targeted genes included a CTP  
183 synthase and a DNA methyltransferase N-4/N-6 domain protein (Fig. 2C, Table S5).  
184 Methyltransferases protect DNA against cleavage by restriction enzymes<sup>35</sup>. Inactivation of such  
185 a methyltransferase might increase vulnerability of the episymbiont towards enzymatic cleavage  
186 by the host.

187

188



## 189 **CRISPR targeting in two independent host-episymbiont systems**

190 Targeting of episymbiont's genomes by altiarchaean CRISPR spacers was also observed in the  
191 HURL ecosystem. In contrast to *Ca. A. crystalense*'s CRISPR-Cas I-B dependent targeting of *Ca. H.*  
192 *crystalense* genomes in the CG environment, *Ca. Altiarchaeum horonobense* found within the  
193 HURL ecosystem appeared to employ CRISPR spacers of an unassigned array (*i.e.*, no *cas* genes  
194 in direct vicinity could be detected due to genome fragmentation but the DR sequence is identical  
195 to type III CRISPR-Cas systems of other Altiarchaea<sup>28</sup>) to potentially ward off *Ca. Huberiarchaeum*  
196 *juliae* episymbionts and viral invaders (Fig. 1D). While spacers of this unassigned array targeting  
197 the *Ca. H. juliae*'s genome exhibited greater diversity compared to the self-targeting  
198 counterparts of *Ca. A. horonobense*'s (Fig. 1D), their relative abundance in the metagenome was  
199 nearly two-fold lower (Fig. 1E). The ecosystem-specific involvement of CRISPR-Cas I-B along with  
200 the unassigned array targeting of the episymbiont genomes in two distinct subsurface  
201 ecosystems seemingly indicates an independent evolution of defense against intruding DNA,  
202 which aligns with previous investigations that demonstrated a strict biogeography of *Ca.*  
203 Altiarchaea core genomes and site-specific evolution<sup>36</sup>. Given the site-specific evolution of *Ca.*  
204 Altiarchaea an alternative explanation for the acquisition of spacers against foreign chromosomal  
205 DNA might be avoidance of spoilage of the host chromosome by intruding genes (horizontal gene  
206 transfer). In Haloarchaea, such a mechanism has been shown to indirectly control for unwanted  
207 horizontal gene transfer between strains of the same genus<sup>37</sup>.

208         Spacers of the unassigned CRISPR array were detected in much greater diversity than  
209 those of CRISPR I-B systems at the HURL site (Fig. 1D-E). While spacers of the unassigned array of  
210 CG-derived *Ca. A. crystalense* self-target chromosomal gene sequences, the spacers of the  
211 unassigned array of *Ca. A. horonobense*'s self-target intergenic regions (Fig. 1C, Table S7).  
212 Notably, it has been demonstrated in haloarchaea that self-targeting does not necessarily lead  
213 to cell suicide<sup>38</sup>. Assuming that the CRISPR-Cas interference is associated with a defense against  
214 the symbiont, a plethora of spacers present at CG might effectively repress the symbiont  
215 (host:symbiont = 11:1 based on metagenomic read mapping), while a lower abundance of  
216 spacers targeting *Ca. H. juliae* at HURL was associated with a higher presence of episymbionts  
217 (host:symbiont = 6:1).

## 218 **Episymbionts metabolically complement self-targeted hosts**

219 We applied genome-scale metabolic modeling to examine the different symbiotic interactions of  
220 *Ca. Altiarchaea* and *Ca. Huberiarchaea* implicated by variations in the CRISPR systems and host-  
221 symbiont ratios of the two ecosystems analyzed. MAGs (ten genomes of *Ca. A. crystalense*, ten  
222 of *Ca. H. crystalense*, one of *Ca. A. horonobense*, one of *Ca. H. julieae*), SAGs (seven of *Ca. A.*  
223 *crystalense* and one of *Ca. H. crystalense*) and transcriptomic data (extracted spacers of samples  
224 CG05, CG08 and CG16 from 2015) from CG and HURL ecosystems were used to render genome-  
225 scale metabolic reconstructions. Although we applied thorough manual data curation<sup>20,27</sup>, the  
226 genomes were fairly fragmented (average  $N50_{\text{host/CG}} = 8067.24$ , average  $N50_{\text{symbtont/CG}} =$   
227  $14983.73$ , average  $N50_{\text{host/HURL}} = 3604$ , average  $N50_{\text{symbtont/HURL}} = 4115$ ), and missing information  
228 due to fragmentation or binning errors cannot be excluded.

229 A consensus model was created for each ecosystem to cogently summarize and compare  
230 the metabolic capacities of *Ca. Altiarchaeum* and *Ca. Huberiarchaeum*, and constraint-based  
231 modeling of these metabolic networks facilitated an assessment of host-symbiont metabolic  
232 complementarity (Fig. 2E-F, Fig. S8). Models of the CG and HURL environment both revealed a  
233 significant reliance of *Ca. Huberiarchaea* upon its host's metabolism yet little to no dependency  
234 of the host upon the metabolism of *Ca. Huberiarchaea*. For example, glucose, amino acids,  
235 vitamins, and energy carrying compounds like adenosine triphosphate (ATP) were transferred  
236 from *Ca. Altiarchaeum* to *Ca. Huberiarchaeum* in both models (Table S8-S11; details in  
237 Supplementary Results), supporting the notion that *Ca. Altiarchaeum* is a primary producer, while  
238 *Ca. Huberiarchaeum* relies on its host for carbon and energy sources<sup>14</sup>.

239 Analyses of CG and HURL host-symbiont relationships also revealed highly variable  
240 metabolic collaborations between episymbionts and their hosts. In the CG ecosystem, a  
241 deoxycytidylate monophosphate (dCMP) deaminase was absent in *Ca. A. crystalense* but present  
242 in *Ca. H. crystalense*. This gene is essential to reach a non-zero biomass for *Ca. A. crystalense* in  
243 the model (see Supplementary Results), suggesting a collaborative effort of synthesizing  
244 pyrimidine (Fig. S8C). Similarly, HURL-borne *Ca. A. horonobense* genomes lacked deoxythymidine  
245 monophosphate (dTMP) synthase genes, while these genes were present in the genomes of *Ca.*  
246 *H. julieae* – once again implicating collaboration, namely in folate biosynthesis (Fig. S8C; Fig. 3).

247 At HURL, the self-targeting of genes in *Ca. Altiarchaea* did not impact the host's dependency on  
248 the symbiont's metabolism within both CRISPR systems (Fig. S8, Supplementary Results). At CG,  
249 however, eliminating the functions of genes self-targeted by the I-B system in metabolic models  
250 exposed additional modes of complementing *Ca. A. crystalense*'s metabolic demands by *Ca.*  
251 *Huberarchaea* via lysyl-tRNA synthetases and phenylalanyl-tRNA synthetases (Fig. 2D-E, Fig. 3,  
252 and Fig. S8A, and C-F). The respective protein sequences were not horizontally transferred  
253 between *Ca. Altiarchaea* and *Ca. Huberarchaea* based on phylogenetic analyses; instead, the  
254 phenylalanyl-tRNA synthetase of *Ca. Huberarchaeum* can be traced back with strong confidence  
255 to *Ca. Woesearchaeota* and *Ca. Pacearchaeota* (Supplementary Data).

256 While the protospacers of *Ca. Altiarchaea* viruses and *Ca. H. crystalense* harbored a  
257 definitive PAM (5'-TTN-3' associated with other I-B systems<sup>39</sup>), no such clear motif was detected  
258 in the host protospacers. Here, the second base of the putative PAM region, exhibited a four-fold  
259 greater than the average single nucleotide polymorphism (SNP) rate of genes (Fig. S9; details in  
260 Supplementary Results). Mutations in the PAM region diverging from the 5'-TTN-3' motif would  
261 prevent self-targeting at least for parts of the *altiarchaeal* population<sup>40</sup> and thus protect the host  
262 chromosome from CRISPR-Cas-mediated cleavage. In our model, removal of self-targeting would  
263 lessen the metabolic dependence on the symbiont and enable subpopulations of *Ca. Altiarchaea*  
264 to flourish more independently. The missing PAM sequence for self-targeting spacers and the  
265 increased SNP-rate in such regions compared to those targeting the episymbiont suggest that the  
266 population of *Ca. Altiarchaea* is adapting to escape the dependency of the symbiont. Considering  
267 that acquisition of self-targeting spacers is a stochastic process<sup>41</sup>, escape mutations or deletions  
268 within the essential targeted genes could have detrimental effects on the cell viability due to the  
269 deficits in the corresponding metabolic activities resulting in cell suicide (reviewed in Heussler  
270 and O'Toole, 2016)<sup>31</sup>. Episymbionts could provide a temporary relief to the host cell by  
271 complementing the metabolic deficiency, becoming a bona fide symbiont, at least until the  
272 metabolic autonomy of the host is reestablished. We thus hypothesize that interactions between  
273 hosts and episymbiont depend on the genotype of the host and can consequently be either  
274 mutualistic or parasitic. However, cultivation of the host-symbiont system along with establishing  
275 a genetic system to modify the host genome are necessary to test this hypothesis.

276

## 277 **Inter-phylum interactions in other symbiotic archaea**

278 To facilitate the overlay of our findings on other potential archaeal host-DPANN episymbiont  
279 relationships, we analyzed CRISPR spacer matches between all archaeal genomes publicly  
280 available in NCBI's GenBank (7,012 genomes: state May, 2021, Table S4). After having extracted  
281 106,641 spacer sequences, 39,875 distinct spacer-to-protospacer matches across all genomes  
282 were detected. Few contigs carrying CRISPR arrays (*e.g.*, for *Ca. Micrarchaeum*) also contained  
283 taxonomic hallmark genes, such as those coding for DNA-directed RNA polymerase subunit or  
284 ribosomal proteins, which provided additional confidence for the correct assignment of spacers  
285 to the fragmented public MAGs. The spacer hits accounted for both self-targeting and  
286 interspecies spacer interactions (Extended Data Fig. 4). Network analyses showed the genomes  
287 of the DPANN *Ca. Aenigmarchaeota* and *Ca. Altiarchaeota* (Fig. 4), as well as *Sulfolobus*,  
288 *Methanosarcina*, *Haloferax*, and *Halobacterium* spp. forming large clusters resulting from a  
289 wealth of interspecies hits and/or self-targeting (Extended Data Fig. 4), which was also previously  
290 shown for other archaea<sup>37,42</sup>. Well-established DPANN-host co-cultures, *e.g. Ignicoccus hospitalis*  
291 and *Nanoarchaeum equitans*<sup>43</sup>, did not exhibit CRISPR-Cas-derived targeting to either of the  
292 symbionts in our archaeal genome dataset.

293 Particularly for the hydrothermal system of Guaymas Basin, Gulf of California<sup>44</sup>, our  
294 approach enabled the *a priori* prediction of DPANN-host interactions based on CRISPR-Cas  
295 genome targeting (Fig. 4). Analyses of the spacer-protospacer matches from the read data of  
296 Guaymas Basin revealed frequent targeting (160 spacer-protospacer matches) of *Ca.*  
297 *Aenigmarchaeota* by *Ca. Bathyarchaeota*. Genes targeted by these spacer-protospacer matches,  
298 *e.g.*, encode for the LamGL domain-containing protein, which is *inter alia* responsible for the  
299 binding of sulfated glycolipids<sup>45,46</sup>, and the ribonucleoside triphosphate reductase, amenable for  
300 catalysis of the conversion of ribonucleotides into deoxyribonucleotides<sup>47</sup>.

301 Another host-DPANN interaction unveiled by these analyses involves *Ca. Micrarchaeota*  
302 spacers matching a Thermoprotei archaeon, with both of these genomes arising from the same  
303 ecosystem but a few centimeters apart in depth<sup>44</sup>. Comparing those targeted gene-encoding  
304 regions to the targeted genetic regions in *Ca. Huberiarchaeum* by spacers of *Ca. Altiarchaeum*  
305 (CTP synthase and DNA methylase, see above), no acquisition pattern of spacers directed against

306 genomic regions that encode specific functions could be detected. Overall, these findings suggest  
307 spacer-protospacer matches are a useful tool for identifying *in-silico* host-symbiont interactions  
308 of uncultivated archaea based on metagenomic analyses.

309

## 310 **DISCUSSION**

311 The findings discussed here demonstrate that archaeal CRISPR-Cas systems acquire  
312 resistance not only to genomes of foreign MGEs<sup>1</sup> and closely related species<sup>32</sup> but also to archaea  
313 of other phyla, particularly episymbionts belonging to the DPANN superphylum. Our results  
314 suggest that CRISPR-Cas-mediated adaptive immunity might lead to complex interactions  
315 between the host and symbiont at the population level, possibly drawing the host into  
316 maintaining a collaborative relationship with the symbiont due to balancing the self-targeting  
317 nature of the host's CRISPR system and the potential defense against the episymbiont. Based on  
318 our results from single-cell genomic data, metagenomes, and metatranscriptomes, we suggest  
319 that a double-edged sword drives the evolution of microbial populations, *i.e.*, CRISPR-Cas-  
320 mediated defenses likely render a major fraction of the DPANN episymbiont population truly  
321 parasitic, while the remainder seem to support the host in a mutualistic fashion.

322 Future studies should be set out with the aim of cultivating of the host-symbiont system  
323 to validate the herein proposed CRISPR-Cas interference. However, cultivation attempts might  
324 selectively enrich for systems with mutualistic relationship, and *in silico* screening of currently  
325 existing host-DPANN co-cultures for spacer targeting of the episymbiont's genome were devoid  
326 of such an interaction, including the well-known archaeal system *Ignicoccus hospitalis* and  
327 *Nanoarchaeum equitans*<sup>13</sup>. Consequently, genetic engineering of the host and the symbiont will  
328 be necessary to eventually clarify the relationship between hosts and DPANN symbionts – may it  
329 be mutualistic, parasitic, or a mixed population model as suggested by our findings.

330

## 331 **METHODS**

332 **RNA extraction and metatranscriptomic sequencing.** Samples for transcriptomics were collected  
333 along with DNA samples as previously published<sup>20</sup>. For the samples CG05, CG08 and CG16 we  
334 filtered approx. 189, 151, and 151 L of geyser-erupted water, respectively. The MoBio PowerMax

335 Soil DNA kit, now re-branded as the Qiagen DNeasy PowerMax Soil kit (Qiagen, Germantown,  
336 MD), was used to perform all metagenomic RNA extractions. Filters were aseptically cut into  
337 pieces, and 20 mL of lysis buffer from the kit was added for removal of cells from the filters. The  
338 manufacturer's alternative protocol, entitled "Alternative PowerMax Protocol for  
339 Isolation of RNA and DNA from Low Biomass Soil with Low Humics" was adjusted as follows:  
340 briefly, 10 mL of Bead Solution was added to the thawed filter and vortexed at maximum speed  
341 for 5 minutes to remove cells. The cell solution was transferred to a bead tube, and 5 ml of  
342 phenol:chloroform:isoamyl alcohol (25:24:1), pH 6.6, was added and homogenized by vortexing  
343 for 10 minutes. The manufacturer's protocol was followed thereafter. The metagenomic RNA  
344 extracts underwent DNase treatment using the Qiagen DNase Max kit (Qiagen), following  
345 manufacturer's standard protocol. Quality control and quantification of all RNA extracts were  
346 performed using the Agilent Bioanalyzer RNA 6000 Nano kit (Agilent, Santa Clara, CA).  
347 Sequencing libraries were created using the Illumina TruSeq Stranded mRNA Library Prep Kit,  
348 following manufacturer's protocol (Illumina, San Diego, CA). Libraries were sequenced on the  
349 Illumina HiSeq 2500 platform (Illumina).

350 **Sample preparation for FISH.** Groundwater for FISH analysis was sampled to visualize the  
351 Altiarchaea-Huberarchaea relationship within the CG and HURL ecosystem. Water from CG was  
352 sampled onto a 0.2  $\mu\text{m}$  filter with a syringe filter holder until the filter started clogging and  
353 afterwards fixed by slowly pressing 3% formaldehyde (Thermo Fisher Scientific, MA, USA)  
354 through the filter to exchange the sample water with fixative. Fixation was performed for one  
355 hour in the dark. Within the filter holder, a washing step with 3x 20 mL Phosphate Buffered Saline  
356 (PBS) (conc. 1 v/v%) was done, followed by alternating washing and incubation with ethanol with  
357 50, 70 and 100% (v/v)% for 10 minutes at room temperature. The filter holder was opened in a  
358 sterile environment, and the filter was stored in a petri dish with the biofilm facing upwards and  
359 then air dried for 10 minutes. Filter samples for FISH from CG of the sampling campaign in 2021  
360 were covered and stored in RNA $later^{\text{TM}}$  (Invitrogen by Thermo Fisher Scientific; Ref: AM7021).

361 **Imaging of FISH samples.** FISH was performed according to Schwank et al.<sup>14</sup> with the following  
362 modifications. DAPI (4',6-Diamidino-2-Phenylindole) was used at concentrations of 4  $\mu\text{g}$  per mL  
363 without dilution in the washing buffer. Visualization was performed with an Axio Imager M2m

364 epifluorescence microscope (X-Cite XYLIS Broad Spectrum LED Illumination System, Excelitas)  
365 equipped with an Axio Cam MRm and a Zen 3.4 Pro software (version 3.4.91.00000) (Carl Zeiss  
366 Microscopy GmbH, Jena, Germany). Imaging was carried out by using the 110x/1.3 oil objective  
367 EC-Plan NEOFLUAR (Carl Zeiss Microscopy GmbH) and three different filter sets (Carl Zeiss): 49  
368 DAPI for imaging *Ca. A. crystalense/horonobense* cells and *Ca. H. crystalense/julieae* cells, 43 Cy3  
369 for the detection of *Ca. Huberiarchaea* signals, and 09 for achieving 16S rRNA signals of *Ca.*  
370 *Altiarchaea*. The FISH images are shared within a Figshare folder  
371 (10.6084/m9.figshare.22739849).

372 **Metagenome assembly and genome reconstruction.** Omic datasets generated from sampling  
373 campaigns for CG<sup>20,27</sup> (Utah, USA) and HURL<sup>24</sup> (Japan) were downloaded from the NCBI Sequence  
374 Read Archive (SRA) in April 2019. Please refer to Table S1 (metagenomic and metatranscriptomic  
375 datasets), Table S2 (genome accessions of *Ca. Aliarchaeum* and *Ca. Huberiarchaeum*), and Table  
376 S3 (single cell genomic dataset) for all accession numbers of publicly available datasets and  
377 generated genomes used in this study. For all metagenomic datasets of CG and HURL, quality  
378 filtering and trimming of reads was done using BBduk  
379 (<https://github.com/BioInfoTools/BBMap/blob/master/sh/bbduk.sh>, version 37.09) and Sickel<sup>48</sup>  
380 (version 1.33). The MetaSPAdes<sup>49</sup> (version 3.10) and Bowtie2<sup>50</sup> utilities (--sensitive, version  
381 2.3.5.1) were applied to assemble reads and estimate coverage, respectively. Scaffolds < 1 kbp  
382 were excluded from further analysis. The interactive uBin<sup>51</sup> software (version 0.9.14) was used  
383 to segregate the genomes of *Ca. Altiarchaeum horonobense* and *Ca. Huberiarchaeum julieae*  
384 based on %GC content, taxonomy, and coverage information. To determine genome  
385 completeness and contamination we used checkM<sup>52</sup> (v1.2.2) (Table S2). Previously published *Ca.*  
386 *Huberiarchaeum* genomes generated from each of the CG and HURL ecosystems were used as  
387 probes to identify respective scaffolds at the protein level (≥80% similarity).

388 **Phylogeny of *Altiarchaeum* and *Huberiarchaeum*.** A reference dataset spanning the diversity of  
389 176 archaeal genomes was used to place *Ca. Huberiarchaea* and *Ca. Altiarchaea* phylogenetically.  
390 The accession numbers of all genomes within the reference datasets can be found in the  
391 Supplementary Data within the phylogenetic tree (with the suffix "GCA\_"). To avoid redundancy,  
392 all genomes annotated as *Ca. Altiarchaea* on NCBI (June 2019), previously published *Ca.*

393 Altiarchaea genomes<sup>36</sup>, and one representative genome from *Ca. Altiarchaeum* and *Ca.*  
394 *Huberiarchaeum* were consolidated for this work. Individual homology searches were executed  
395 across these datasets, using HMMER 3.2.1<sup>53</sup> with the Phylosift<sup>54</sup> marker HMM profiles and an e-  
396 value cutoff of  $1 \times 10^{-5}$ . All DNA sequences were aligned with MUSCLE v3.8.31<sup>55</sup> (default  
397 parameters) and manually curated to fuse fragmented genes and remove distant homologs and  
398 paralogous copies. One *Ca. Altiarchaea* genome (GCA\_003663105) was likely contaminated and  
399 thus removed from the final alignments. Sequence sets resulting from each of the four datasets  
400 were fused together (36 single-gene datasets; one of the 37 Phylosift marker genes  
401 (DNGNGWU00035) was omitted due to many missing taxa), realigned as before, trimmed with  
402 BMGE (BLOSUM30)<sup>56</sup>, and concatenated into one supermatrix (200 taxa; 5,974 amino acid  
403 positions). Phylogenies were reconstructed with IQ-TREE 2<sup>57</sup> (v2.0.5), first using ModelFinder<sup>58</sup>,  
404 then using that phylogeny as a guide, with the PMSF model<sup>59</sup> (LG+C60+F+G). Branch supports  
405 were calculated using 1,000 ultrafast bootstrap<sup>60</sup> and 1,000 SH-aLRT<sup>61</sup> replicates and the aBayes<sup>62</sup>  
406 test and trees were visualized in iTOL<sup>63</sup> (version 5).

407 **Naming of archaeal genomes.** Except for *Ca. Huberiarchaeum crystalense*, all host and  
408 episymbiont species were previously only classified at the genus level or – in case of the  
409 episymbiont from the HURL ecosystem – not classified at all. Using established average  
410 nucleotide identity (ANI) and Average Amino Acid Identity (AAI) cutoffs along with phylogenetic  
411 analyses (Fig. 1 and Supplementary Data), we established the host-symbiont pairs as *Ca.*  
412 *Altiarchaeum crystalense* and *Ca. Huberiarchaeum crystalense* from the CG ecosystem (named  
413 after the ecosystem Crystal Geysir) and *Ca. Altiarchaeum horonobense* (named after the  
414 sampling site Horonobe) and *Ca. Huberiarchaeum julieae* (named after subsurface microbiologist  
415 Julie Huber).

416 **Phylogenetic reconstruction of individual metabolic genes.** For the phylogenies of lysine and  
417 phenylalanine (subunit B) tRNA synthetases, the protein sequences inferred from both genes  
418 from *Ca. Altiarchaeum hamiconexum* and *Ca. Huberiarchaeum crystalense* were used for  
419 homology searches against local databases of 1808 archaeal and 25118 bacterial genomes (all  
420 genomes of the respective domain on NCBI as of 2019.06.01 dereplicated at species level) with  
421 DIAMOND v2.0.15.153<sup>64</sup>. The maximum number of target sequences (-k 400) was determined by



422 trying different numbers (100, 200, 400, 800, 1000, 0/all), aligning with MAFFT FFT-NS-2 (v7.505)  
423 and running a preliminary phylogeny (BioNJ or PhyML without tree topology optimization) in  
424 Seaview version 5.0.4<sup>65</sup>. We picked the number that we deemed to give a reasonable view of the  
425 origin of each sequence, without including too many divergent homologs or increasing the  
426 downstream computational load too much. The original query sequences were added to the set  
427 of hits and aligned with MAFFT E-INS-I. The datasets were curated semi-manually  
428 ([https://github.com/ProbstLab/Adam\\_Kolyfetis\\_2021\\_methanogenesis/blob/master/fuse\\_sequences.py](https://github.com/ProbstLab/Adam_Kolyfetis_2021_methanogenesis/blob/master/fuse_sequences.py)) to fuse fragmented sequences, realigned as before, and trimmed with BMGE<sup>56</sup>  
429 (BLOSUM30). Phylogenies were reconstructed with IQ-TREE 2<sup>57</sup> using ModelFinder<sup>58</sup> for the  
430 model selection and branch supports calculated using 1000 ultrafast bootstrap<sup>60</sup> and 1000 SH-  
431 aLRT replicates.

433 **CRISPR system extraction and viral sequence determination.** The CRISPR systems of 18 distinct  
434 *Ca. Altiarchaeum crystalense* genomes<sup>20,26</sup> (Table S2) and one *Ca. Altiarchaeum horonobense*  
435 genome<sup>24</sup> (Table S2) were extracted with CRISPRCasFinder<sup>66</sup> (version 1.2), and annotated *cas*  
436 genes were used to identify CRISPR-Cas cassettes. Two resulting consensus DR sequences were  
437 used as input for MetaCRIST<sup>67</sup> (-d 1 -c 1 -a 1 -h -r), analysis of metagenomic reads,  
438 metatranscriptomic reads, and single cell genome reads. Only spacers having adjacent repeat  
439 sequences bearing 100% similarity with the respective read were considered. All spacers shorter  
440 than 24 bps, longer than 57 bps, or harboring homopolymers of six or more identical bases in a  
441 row were excluded. Spacers were clustered to 97% nucleotide identity using CD-hit<sup>68</sup> (version  
442 4.8.1) and respective centroid sequences were used in downstream analyses.

443 To check if spacers were biased towards matching genome transcripts, the orientation of the  
444 CRISPR array was confirmed on all available *Ca. Altiarchaeum crystalense* genomes to identify  
445 the forward strand that corresponds to CRISPR-RNA by using CRISPRDirection2.0 with default  
446 settings<sup>69</sup>. To avoid false positive predictions of self-targeting and episymbiont targeting, we  
447 masked prophage region, predicted by VirSorter<sup>70</sup> (category 1-3, 4-6) and transposon regions,  
448 predicted by ISEScan<sup>71</sup>. The spacers from this analysis were blasted (nucleotide blast,  
449 bidirectional [default setting] and unidirectional [-strand plus] on the forward strand) against the  
450 CDS data of 18 *Ca. Altiarchaeum crystalense* genomes (including seven SAGs), eleven *Ca.*

451 Huberiarchaeum crystalense (including 2 SAGs), one genome of *Ca. Altiarchaeum horonobense*  
452 and *Ca. Huberiarchaeum julieae*, respectively. All unpublished viral genomes used in this study  
453 are deposited in the Figshare folder (10.6084/m9.figshare.22738568).

454 **Detection, dereplication and analysis of DNA viral scaffolds.** Assembled metagenomes were  
455 used to extract and predict viral and putatively viral sequences as previously performed<sup>28</sup>. In  
456 brief, predicted viral operational taxonomic units (vOTUs) >3kb were dereplicated via usearch<sup>72</sup>  
457 at 95% nucleotide identity resulting in centroid sequences for downstream analysis. VOTUs were  
458 identified via blastn<sup>73</sup> (--short, filtering for 80% similarity, version 2.9.0+) of Clustered regularly  
459 interspaced short palindromic repeat (CRISPR)-derived spacers against centroid vOTUs.  
460 Completeness and origin (host, viral, unclassified) of vOTUs was assessed using CheckV v.0.4.0<sup>74</sup>.  
461 Clustering of viral sequences with a recent viral Refseq database<sup>75</sup> (release July 2022). and  
462 previously detected Altiarchaea-targeting viruses<sup>28</sup> was performed using vConTACT2<sup>76,77</sup>  
463 v.0.11.3, VICTOR<sup>78</sup> (using nucleic acid sequences) and VIRIDIC<sup>79</sup> under default settings and for  
464 calculating intergenomic similarities. In VICTOR, all pairwise comparisons of the nucleotide  
465 sequences were conducted using the Genome-BLAST Distance Phylogeny<sup>80</sup> (GBDP) method  
466 under settings recommended for prokaryotic viruses<sup>78</sup>. The resulting intergenomic distances  
467 from VICTOR were used to infer a balanced minimum evolution tree with branch support via  
468 FastME including Subtree Pruning and Regrafting post-processing<sup>81</sup> for the distance formula D0.  
469 Branch support was inferred from 100 pseudo-bootstrap replicates each. Trees were rooted at  
470 the midpoint<sup>82</sup>. Visualization of viral clusters identified with vConTACT2 in conjunction with the  
471 viral RefSeq database was performed using Cytoscape v.3.9.02<sup>83</sup>. In addition, a circular proteomic  
472 tree with viral genomes using the Virus-Host DB: RefSeq release 217 was build using ViPTree  
473 version 3.5.<sup>84</sup>. Within ViPTree, dsDNA was selected as nucleic acid type and “any host” chosen  
474 as host category.

475 **Sliding window for coverage analysis of regions targeted by CRISPR spacers.** Variations in  
476 coverage over the genomes were investigated to deduce possible negative selection at targeted  
477 sites. Targeted scaffolds from individual genomes were mapped back to the raw reads (from  
478 sample CG05, CG08 and CG16) with Bowtie2<sup>50</sup> with default settings. Mappings were filtered to  
479 remove hits with more than three mismatches using SAMtools<sup>85</sup> (version 1.10). Genomecov from

480 BEDtools (version 2.27.1) was used to calculate coverage per position<sup>86</sup>. The first and last 150  
481 base pairs of each scaffold, and possible transposons and viruses were masked by setting the  
482 breadth to zero. Mean breadth from sliding windows of 35 base pairs were calculated. In  
483 addition, all position with a coverage lower than 10 were excluded. The median coverage of each  
484 scaffold ( $\delta$ ) serves to differentiate high and low breadth. Wilcoxon signed rank tests (standard  
485 function R<sup>87</sup>) were performed between targeted regions of a scaffold and the same amount of  
486 randomly drawn non-targeted windows from the same scaffold. Random sampling and the test  
487 were repeated 1000 times for each scaffold.

488 **Models for *Ca. Altiarchaea* and *Ca. Huberarchaea* host-symbiont interaction based on genomic**  
489 **information.** To infer metabolic interactions, genome-scale metabolic reconstructions of *Ca.*  
490 *Altiarchaeum crystalense/horonobense* and *Ca. Huberarchaeum crystalense/julieae* (see  
491 accession numbers Table S2) were based on MAGs and SAGs identified from CG (AltiCG-HuberCG  
492 model) and HURL (AltiHURL-HuberHURL model). The genome-scale metabolic models of AltiCG-  
493 HuberCG and AltiHURL-HuberHURL were represented in a YAML format following conventions  
494 defined by the PSAMM software package<sup>89,90</sup>. The AltiCG-HuberCG model included 515 genes of  
495 *Ca. A. crystalense* and 88 genes of *Ca. H. crystalense*, associated with 477 and 125 reactions,  
496 respectively (Table S8). The AltiHURL-HuberHURL model included 388 *Ca. A. horonobense* and  
497 78 *Ca. H. julieae* genes, associated with 495 and 128 reactions, respectively (Table S9). Each  
498 model contained two compartments (one for *Ca. Altiarchaeum* and one for *Ca. Huberarchaeum*),  
499 with either restricted or unlimited metabolite exchanges between the two compartments to  
500 model the metabolite availability upon cytoplasmic fusion of the two organisms.

501 Details of the model are represented in Tables S8 - S11. The CG model was based on the  
502 prediction of metabolic pathways using combined annotation of all MAGs and SAGs identified  
503 from this and a prior study<sup>14</sup>. Protein sequences annotated from the individual MAGs and SAGs  
504 were clustered at 100% amino acid identity using CD-HIT<sup>68,91</sup>, followed by a pangenome analysis  
505 to capture metabolic capacities represented by the entire population. Automated metabolic  
506 reconstruction was performed based on ortholog mapping to (i) existing models of other archaeal  
507 strains, i.e. *Pyrococcus furiosus*, *Thermococcus eurythermalis*, *Methanosarcina barkeri* and  
508 *Methanococcus maripaludis*<sup>92,93</sup>, and (ii) public databases, such as the Kyoto Encyclopedia of

509 Genes and Genomes<sup>94</sup> (KEGG), EggNOG<sup>95</sup> and Transporter Classification Database<sup>96</sup> (TCDB).  
510 Extensive manual curations were carried out following the automated reconstruction to integrate  
511 prior annotations of *Ca. Altiarchaeum*'s and *Ca. Huberiarchaeum*'s metabolism<sup>14,20</sup>, as well as  
512 latest biochemical evidence of enzymatic functions in archaeal organisms (Tables S8 and S9).  
513 Overall, literature evidence was assigned to 137 reactions in the model for AltiCG-HuberCG and  
514 144 reactions in the model for AltiHURL-HuberHURL through homologous mapping to  
515 experimentally verified enzymes. The biomass equations of *Ca. Altiarchaeum* and *Ca.*  
516 *Huberiarchaeum* were individually formulated in both models following a standard procedure  
517 (Table S8 and S9). The biosynthesis of macromolecules (e.g., DNA, RNA, protein, and lipids) were  
518 defined to account for the mM composition of each building block in assembling 1 g of a given  
519 component and the associated energy cost. The stoichiometry of DNA and RNA biosynthesis was  
520 derived based on the average composition of nucleotides in the genomes and coding genes,  
521 respectively. The energy cost for DNA and RNA synthesis was estimated as 2 mM of ATP per  
522 millimole of nucleotides according to the mechanism of polynucleotide biosynthesis<sup>97</sup>. The  
523 stoichiometry of protein biosynthesis was calculated based on the average composition of amino  
524 acids in the corresponding proteome, and the associated energy cost was estimated based on  
525 the mechanism of protein synthesis<sup>98</sup>, where one ATP was consumed for each tRNA charging,  
526 and two GTPs were consumed for extending one amino acid to a growing peptide chain. The  
527 tRNA charging equations were represented separately for each amino acid. The stoichiometry of  
528 lipid biosynthesis was formulated based on experimental measurements of the weight  
529 compositions of core lipids and header groups of *Ca. Altiarchaeum* or *Ca. Huberiarchaeum* of the  
530 respective system<sup>20</sup>. Following the definition of macromolecular synthesis functions, the biomass  
531 equations of *Ca. A. crystalense*, *Ca. H. crystalense*, *Ca. A. horonobense*, and *Ca. H. julieae* were  
532 formulated to represent the gram composition of DNA, RNA, proteins, lipids, and vitamins in 1 g  
533 of cell dry weight (gDW). Relative abundance (based on coverage) of the respective genomes was  
534 calculated via metagenomic read mapping with Bowtie2<sup>50</sup> (--sensitive mode). The CG- and HURL-  
535 specific *Ca. Altiarchaeum* and *Ca. Huberiarchaeum* biomass were then combined based on an  
536 estimation of their relative abundance in the respective ecosystems using the metagenomic data.  
537 Specifically, the combined *Altiarchaeum-Huberiarchaeum* biomass has a relative

538 Huberiarchaeum:Altiarchaeum ratio between 0.06 and 0.12 in the CG system, and a ratio of 0.205  
539 in the HURL system (as estimated via stringent read mapping, see Supplementary Results).

540 **Metabolic modeling and reconstruction.** Consensus models of *Ca. Altiarchaeum* to *Ca.*  
541 Huberiarchaeum were constructed based on collections of MAGs and SAGs from CG (20 MAGs,  
542 8 SAGs) and HURL (2 MAGs) (Table S2) to capture the metabolic potential of each population.  
543 Candidate genes were first identified based on a pangenome analysis, which was performed  
544 following ortholog identification using a bidirectional best hit approach<sup>88</sup>. All representative  
545 genes from the MAGs or SAGs of a given ecosystem served as candidates for that ecosystem's  
546 metabolic reconstruction. Complementary metabolic characteristics were identified between *Ca.*  
547 Altiarchaeum and *Ca. Huberiarchaeum* via a *fastgapfill* implementation in the PSAMM software  
548 (version v1.0) package<sup>89,90</sup> using the cplex solver (v12.7.1.0). Simulations targeted the growth  
549 optimization of *Ca. Altiarchaeum* while applying the metabolic reactions of *Ca. Huberiarchaeum*  
550 as a reference, which facilitated the identification of *Ca. Huberiarchaeum*-encoded  
551 complementary functions essential for *Ca. Altiarchaeum* - and vice versa. Combined *Ca.*  
552 Altiarchaeum and *Ca. Huberiarchaeum* metabolic models were formulated with exchange  
553 constraints representative of environmental *in situ* geochemical measurements corresponding  
554 to either CG or HURL (Table S9 and S10). Comparative analyses based on computational  
555 simulations were carried out in the presence or absence of CRISPR-targeted genes. This enabled  
556 the identification of changes in metabolite transfer and/or metabolic collaboration between *Ca.*  
557 Altiarchaeum and *Ca. Huberiarchaeum* (Fig. S8) upon targeting specific genes with spacers.

558 Metabolic gaps in the production of biomass components by *Ca. Altiarchaea* were identified using  
559 the PSAMM *fluxcheck* and *gapcheck* functions<sup>89,90</sup> Candidate gap-filling reactions for unblocking  
560 each biomass component were identified using the PSAMM *fastgapfill* implementation with the  
561 KEGG reaction database<sup>94</sup> as a reference, and subsequently curated before being incorporated  
562 into the models. A total of 17 gap-filling reactions were included in the *Ca. Altiarchaea*  
563 compartment of both CG and HURL models, including functions in the citrate cycle, amino acids-  
564 , lipids-, and cofactor-biosynthesis. The overall stoichiometric consistency, formula and charge  
565 balance of the model were validated using the PSAMM *masscheck*, *formulacheck*, and  
566 *chargecheck* functions<sup>89,90</sup>. The exchange reactions, compound sources or sinks, biomass

567 equations and reactions involving compounds with undefined group R or X were excluded from  
568 formula and charge checks but instead manually inspected to ensure proper formulation.

569 Metabolic simulations were performed with PSAMM version 1.0 using the IBM ILOG  
570 CPLEX Optimizer version 12.7.1.0 ([https://www.ibm.com/products/ilog-cplex-optimization-](https://www.ibm.com/products/ilog-cplex-optimization-studio)  
571 [studio](https://www.ibm.com/products/ilog-cplex-optimization-studio)). Simulation of the *Ca. Altiarchaeum* – *Ca. Huberiarchaeum* metabolism was formulated  
572 with exchange constraints that represent the corresponding *in situ* geochemical measurements  
573 in the CG<sup>20</sup> and HURL<sup>24</sup>. These geochemical measurements included the ion concentrations in  
574 porewater and the compositions of headspace gas (Table S8 - S11). Some measurements, e.g.,  
575 CO<sub>2</sub> and H<sub>2</sub> at the CG site, were not available, but the compounds were required for biomass  
576 production in the *Ca. Altiarchaeum* – *Ca. Huberiarchaeum* system, and thus they were added to  
577 the exchange without implicit constraints. To simulate the fusion of the cytoplasm between *Ca.*  
578 *Altiarchaeum* and *Ca. Huberiarchaeum*, unlimited metabolite exchange was introduced to allow  
579 the free transfer of all small-molecular metabolites (excluding macromolecules, such as DNA,  
580 RNA, protein, lipids, and biomass) between the *Ca. Altiarchaeum* and *Ca. Huberiarchaeum* cell  
581 compartments.

582 To identify complementary metabolic processes between *Ca. Altiarchaea* and *Ca.*  
583 *Huberiarchaea*, the PSAMM *fastgapfill* implementation<sup>89,90</sup> was applied to optimize the *Ca.*  
584 *Altiarchaea* biomass while using all metabolic reactions in the *Ca. Huberiarchaea* compartment  
585 as a reference database, and vice versa, using corresponding models for CG or HURL. A list of  
586 metabolic reactions, including metabolite exchange functions between *Ca. Altiarchaea* and *Ca.*  
587 *Huberiarchaea*, was identified from this automated gap filling procedure to reveal the potential  
588 metabolic interactions between the two archaea at each site. The predicted complementary  
589 metabolites were subsequently confirmed by showing that the removal of any metabolite  
590 exchange would lead to a non-viable *Ca. Altiarchaea* or *Ca. Huberiarchaea* (biomass production  
591 is zero), suggesting that these metabolite exchanges reflect minimal essential interactions  
592 between *Ca. Altiarchaea* and *Ca. Huberiarchaea* of a given site (Table S8 and S9). Genes  
593 corresponding to the CRISPR type I-B and the unassigned CRISPR array spacer targeting in both  
594 CG and HURL systems were mapped to the metabolic reconstructions of AltiCG-HuberCG and  
595 AltiHURL-HuberHURL, respectively, for the identification of putative targets for simulating the

596 metabolic influences of attacks targeted by the CRISPR system (Table S10 and S11). To identify  
597 changes in the *Ca. Altiarchaea* – *Ca. Huberarchaea* metabolic collaboration when considering  
598 attacks of respective genes by CRISPR-Cas systems, comparisons were made between the  
599 exchange unlimited model (where all metabolites (with the exception of macromolecules) were  
600 allowed to transfer freely between *Ca. Altiarchaeum* and *Ca. Huberarchaeum*) and the exchange  
601 limited model (where only the complementary metabolites were allowed to transfer between  
602 *Ca. Altiarchaeum* and *Ca. Huberarchaeum*). Flux variability analysis (FVA) was applied to the  
603 optimization of the combined *Altiarchaeum*-*Huberarchaeum* biomass in the limited or unlimited  
604 models. Pathways that are required for complementing the effect of CRISPR spacer attacks were  
605 identified by comparing the FVA results of the limited and unlimited models. If the deletion of a  
606 spacer attacked gene would result in a zero-biomass flux in the limited model while a non-zero  
607 biomass flux in the unlimited model, a complementary pathway to the corresponding gene  
608 deletion was explored by identifying the enabling functions in the unlimited model. Note that the  
609 FVA was performed in PSAMM using the CPLEX Optimizer version 12.7.1.0, a zero range is defined  
610 as any fluxes within  $1E^{-6}$  from zero.

611 **PAM analysis of *Ca. Altiarchaea*, *Ca. Huberarchaea*, and viruses.** Applying CRISPRTarget<sup>99</sup>  
612 (accessed in June 2020) with default settings, protospacer adjacent motifs (PAMs) were identified  
613 within the genomes of *Ca. Altiarchaea*, *Ca. Huberarchaea*, and viruses using spacers bearing 80%  
614 sequence similarity. CRISPRTarget results were screened with WebLogo<sup>100,101</sup> (v2.8.2) in batches  
615 of 10,000 8-bps sequences.

616 **SNP analysis.** To identify *Ca. Altiarchaeum* *crystalense* SNPs, reads from samples CG05, CG08,  
617 and CG16 (samples for which also transcriptomic datasets were available, and which were used  
618 in the metabolic modeling) were aligned to nine different MAGs (Table S2) and analyzed  
619 individually by using BMAP (<https://sourceforge.net/projects/bbmap/>, version 38.92) (default  
620 parameters). SNPs were predicted using the VarScan<sup>102</sup> pileup2snp command (v2.4.3; default  
621 settings) with observations and coverage thresholds set to a minimum of two and eight,  
622 respectively. SNPs bearing the reference allele 'N' were excluded if all base called reads showed  
623 this 'N'.

624 **Synthesis of *cas* genes derived from *Ca. Altiarchaea* MAGs.** The CRISPR-Cas gene cassette (Cas1,  
625 Cas2, Cas3, Cas4, Cas5, Cas8b) of one single-cell amplified genomes (SAG) of *Ca. Altiarchaeum*  
626 was used in gene synthesis. The Cas6 gene was annotated in two other SAGs of *Ca. Altiarchaea*,  
627 once with 438 and 468 amino acids respectively. To synthesize these genes, the sequences were  
628 first codon optimized using the BOOST design software v.1.3.9<sup>103</sup> and an *E. coli* codon frequency  
629 table. The synthetic DNA fragments were obtained from Twist Bioscience, CA, USA, which were  
630 later PCR amplified and cloned into the NcoI and XhoI sites of the pET28b vector using the  
631 NEBuilder HiFi Assembly kit (E2621X, New England BioLabs). The PCR was performed using the  
632 KAPA HiFi HotStart ReadyMix (Roche Sequencing, Pleasanton, CA, USA) according to the  
633 manufacturer recommended cycling protocol. The sequences of the refactored *cas* genes were  
634 verified by Pacific Bioscience sequencing. The synthetic building blocks and PCR primer  
635 sequences are listed in Table S13.

636 **CRISPR-Cas activity assay in TXTL.** The activity of the *Ca. Altiarchaea* type I-B CRISPR-Cas system  
637 was tested in a cell-free transcription-translation (TXTL) system. Circular or linear DNA constructs  
638 that were added to a TXTL reaction were transcribed and translated, and RNAs and proteins were  
639 produced<sup>104</sup>. The reaction conditions of the TXTL reactions performed here were adapted from  
640 Wimmer et al.<sup>29</sup>. A deGFP reporter plasmid was generated with Site Directed Mutagenesis (SDM)  
641 using p70a\_deGFP\_PacI<sup>29</sup> as backbone and introducing a TTTTC motif 12 nucleotides upstream  
642 of the p70a promoter driving the deGFP expression. The TTTTC motif was used as a putative PAM  
643 sequence because this motif was found next to a sequence matching a type I-B spacer (see main  
644 text). Constructs encoding single spacer arrays driven by the constitutive promoter J23119  
645 contained either a spacer targeting the p70a promoter region of the reporter plasmid or a non-  
646 targeting spacer. These constructs were generated by Golden Gate adding spacer sequences in a  
647 plasmid which contained two repeat sequences interspaced by two BbsI restriction sites. The  
648 construct p70a-T7RNAP<sup>104</sup> encoding the T7 RNA polymerase and Isopropyl  $\beta$ -D-1-  
649 thiogalactopyranoside (IPTG; Carl Roth, Karlsruhe, Germany) was added to the TXTL reaction to  
650 ensure expression of the *cas* genes. Two Master mixes containing plasmids encoding for Cascade-  
651 forming *cas* proteins were prepared using the stoichiometry Cas8b1-Cas77-Cas51-Cas61, namely  
652 one for the 245 and the 268 amino acids long Cas6. A volume of 3  $\mu$ L TXTL reaction were prepared



653 in Costar 3357 96-well V-bottom plates (Corning, NY, USA) with Costar 2080 cover mats (Corning)  
654 using the liquid handling machine Echo525 (Beckman Coulter, Brea, CA, USA) including the  
655 following components: 2.25  $\mu$ L myTXTL Sigma 70 MasterMix (Arbor Biosciences, MI, USA), 0.2 nM  
656 p70a-T7RNAP, 0.5 mM IPTG, 3 nM Cascade Master mix, 1 nM Cas3 plasmid, and 1 nM targeting  
657 or non-targeting spacer plasmid. After a 4 h pre-incubation period at 29°C to allow the  
658 ribonucleoprotein complex of Cascade and crRNA to form, 1 nM deGFP reporter plasmid  
659 containing the TTTTC motif was added to the TXTL reactions. The reactions were incubated at  
660 29°C for additional 16 h while measuring deGFP expression with BioTek Synergy H1 plate reader  
661 (BioTek, Winooski, VT, USA) at 485/528 nm excitation/emission<sup>105</sup>. Targeting spacer-mediated  
662 binding of the Cascade complex to the target region in the deGFP driving promoter or target  
663 plasmid degradation by Cas3 would lead to inhibition of deGFP production. The non-targeting  
664 spacer does not affect deGFP production and was used as a control. The fluorescence background  
665 values were measured with reactions containing solely myTXTL Sigma 70 MasterMix and  
666 nuclease-free H<sub>2</sub>O and were subtracted from the endpoint deGFP values of the TXTL reactions.  
667 Significance between deGFP values derived from the non-targeting and targeting samples was  
668 calculated with Welch's t-test. All results showed a p-value > 0.05 and were therefore seen as  
669 non-significant. Hence, we concluded that the type I-B systems do not exhibit binding or  
670 degradation activity under the tested conditions. This could be due to the conditions used here  
671 not reflecting the conditions at the sampling site of *Ca. Altiarchaea*, or the motif TTTTC being a  
672 non-recognized PAM. All reactions were performed in triplicates.

673 **PAM assay in TXTL (PAM-DETECT).** To reveal the PAM diversity recognized by the type I-B system  
674 of *Ca. Altiarchaea*, PAM-DETECT (PAM DETermination with Enrichment-based Cell-free TXTL) was  
675 performed. A detailed protocol can be found in Wimmer et al.<sup>29</sup>. A plasmid containing a PAM  
676 library of five randomized nucleotides was used as a target plasmid. A single spacer array plasmid  
677 is constructed as mentioned above harboring a spacer targeting the target plasmid adjacent to  
678 the randomized nucleotides. Upon recognition of a PAM sequence, the Cascade complex binds  
679 to its target and thereby covers a PstI recognition site included in the target region. Cascade-  
680 bound target plasmids are protected from PstI digestion leading to an enrichment of recognized  
681 PAMs, detected by next-generation-sequencing (NGS; specified below). Separate 6  $\mu$ L TXTL

682 reactions were prepared containing one or the other Cascade Master mix mentioned above. TXTL  
683 reactions contained: 4.5  $\mu\text{L}$  myTXTL Sigma 70 Master mix, 0.2 nM pET28a\_T7RNAP<sup>29</sup>, 0.5 mM  
684 IPTG, 3 nM Cascade Master mix, 1 nM targeting spacer plasmid (targeting PAM library plasmid)  
685 and 1 nM PAM library plasmid (pPAM\_library)<sup>29</sup>. After incubation at 29°C for 16 h, the TXTL  
686 samples were diluted 1:400 in nuclease-free H<sub>2</sub>O. A volume of 500  $\mu\text{L}$  of the dilution was digested  
687 with 0.09 units  $\mu\text{L}^{-1}$  PaeI (NEB) in 1x CutSmart Buffer (NEB) at 37°C for 1 h. A “non-digested”  
688 control was prepared using 500  $\mu\text{L}$  of the dilution and adding nuclease-free H<sub>2</sub>O instead of PaeI.  
689 PaeI was inactivated at 65°C for 20 min and proteins were digested with 0.05 mg mL<sup>-1</sup> Proteinase  
690 K (Cytiva, Marlborough, MA, USA) at 45°C for 1 h. Proteinase K was inactivated at 95°C for 5 min  
691 and remaining plasmids were extracted with standard ethanol precipitation. To prepare  
692 sequencing libraries, Illumina adapters with unique dual indices were added in two amplifications  
693 steps using KAPA HiFi HotStart Library Amplification Kit (KAPA Biosystems, Wilmington, MA, USA)  
694 and purification by AMPure XP (Beckman Coulter) after every amplification step. A volume of 15  
695  $\mu\text{L}$  of the ethanol-purified samples was used in a 50  $\mu\text{L}$  PCR reaction with 19 cycles to add Illumina  
696 sequencing primer sites. The flow cell binding sequence was added in the second PCR reaction  
697 using 1 ng purified amplicons generated with the first PCR in a 50  $\mu\text{L}$  reaction and 18 cycles. NGS  
698 was performed on an Illumina NovaSeq 6000 sequencer with 50 bp paired-end reads and 2.0  
699 million reads per sample. PAM wheels were generated according to Leenay et al.<sup>106</sup> and Ondov  
700 et al.<sup>107</sup> and are not depicted here as no PAM enrichment was observed. Absence of PAM  
701 enrichment might be due to the reaction conditions of PAM-DETECT deviating from the  
702 conditions at the sampling site of *Ca. Altiarchaea*. PAM-DETECT assays were performed in  
703 duplicates.

704 **CRISPR-Cas interactions across archaeal diversity.** All archaeal genomes housed in the publicly  
705 accessible NCBI database (May 2021; Table S4) were screened for viral sequence contaminants  
706 using VirSorter<sup>70</sup> (default settings), and all respective hits, as well as annotated plasmids, were  
707 excluded from consideration. The CRISPRCasFinder<sup>66</sup> (version 2.0.3) utility was used to extract  
708 spacers, DR, and *cas* genes from each genome individually with the help of the *cas* gene database  
709 (-ArchaCas). All CRISPR arrays detected were masked in their respective genomes to avoid false  
710 positives, and spacers were filtered for homopolymers and sequence length as described above.

711 All spacer sequences were queried<sup>73</sup> against all archaeal genomes to an 80% nucleotide similarity  
712 threshold, and interactions between genomes based on CRISPR spacer matches were visualized  
713 in Cytoscape<sup>83</sup> (version v.3.9.02). The taxonomy of each genome was pulled from the NCBI  
714 taxonomy database and in single cases validated using Genome Taxonomy Database<sup>108–110</sup>(GTDB-  
715 Tk classify, version v0.3.3, database r89). To avoid false positive predictions of self-targeting and  
716 episymbiont targeting, we masked prophage region, predicted by VirSorter<sup>70</sup> (category 1-3, 4-6).

717

718

## 719 **ACKNOWLEDGEMENTS**

720 This effort was funded by the Ministerium für Kultur und Wissenschaft des Landes Nordrhein-  
721 Westfalen (“Nachwuchsgruppe Dr. Alexander Probst”) and the German Science Foundation  
722 under project NOVAC (grant number DFG PR1603/2-1) and through SPP 2141 (grant number DFG  
723 BE6703/1-1). Genome-scale metabolic modeling was supported by the National Science  
724 Foundation under Grant No. 1553211. The Ministry of Economy, Trade and Industry of Japan  
725 funded a part of the work as “The project for validating assessment methodology in geological  
726 disposal system” (2019 FY, Grant Number: JPJ007597). The work  
727 (proposal: 10.46936/10.25585/60000800) conducted by the U.S. Department of Energy Joint  
728 Genome Institute (<https://ror.org/04xm1d337>), a DOE Office of Science User Facility, is  
729 supported by the Office of Science of the U.S. Department of Energy operated under Contract  
730 No. DE-AC02-05CH11231. M.P. is supported by the Austrian Science Funds (Project MAINTAIN,  
731 DOC 69 doc.funds). P.S.A was supported by a postdoctoral fellowship from the Alexander von  
732 Humboldt Foundation. J.P. was supported by Lundin Energy Norway AS within the framework of  
733 the GeneOil Project. J.R. was supported by the German Science Foundation (grant number  
734 RA3432/1-1, project number 446702140). Support by the German Federal Ministry of Education  
735 and Research within the project “MultiKulti” (BMBF funding code: 161L0285E) is acknowledged.  
736 We thank Ken Dreger for exemplary server administration, and Bettina Siebers, Ivan Berg, Jillian  
737 F. Banfield and Benjamin Meyer for insightful discussion.

738

739

740 **CONFLICT OF INTEREST**

741 The authors declare no conflict of interest.

742

743 **AUTHOR CONTRIBUTIONS**

744 SPE and AJP performed genome-resolved metagenomics, while SPE and JR performed viromics.  
745 JR analyzed viral genomes with input from MK. CRISPR-Cas analyses were done by SPE, JR and  
746 AJP. SNP analysis was performed by MP and TR. Genome-scale modeling was conducted by WZ  
747 and YZ with input from SPE, PAFG, and AJP. Phylogenomic analyses were carried out by PSA. TLVB  
748 provided bioinformatic assistance, and KSch and VT performed microscopy and initial metabolic  
749 analyses. JM and WB re-sampled Crystal Geyser and, JL, TW, and AJP conducted RNA extraction  
750 and sequencing, and SPE analyzed transcriptomes. FW and CB performed binding, cleavage and  
751 PAM assays and JL, JJ, YA, TW, and AJP generated/provided raw data. KS and SPE analyzed the  
752 archaeal CRISPRCas interactions from published NCBI archaeal genomes. AJP conceptualized the  
753 work. SPE, JR, WZ, YZ, and AJP wrote the manuscript with input from all authors.

754

755 **DATA AVAILABILITY STATEMENT**

756 Metagenomic datasets generated from the Crystal Geyser (CG)<sup>20,27</sup> ecosystem (Utah, USA) in  
757 2009, 2014, and 2015 (n = 66), and the Horonobe Underground Research Laboratory (HURL)<sup>24</sup>  
758 (Hokkaido, Japan) environment (n = 2) were downloaded from the NCBI' Sequence Read Archive  
759 (SRA) in April 2019 (Table S1). SAGs generated in a previous study<sup>20</sup> (n = 219) were retrieved from  
760 the JGI's Integrated Microbial Genomes and Microbiomes database<sup>111</sup> (Table S3). The  
761 metagenome-derived genomes of *Ca. A. crystalense* and *Ca. H. crystalense* from CG are publicly  
762 accessible from NCBI (accession numbers in Table S2). The genomes of *Ca. A. horonobense* and  
763 *Ca. H. julieae* from HURL were newly reconstructed in this investigation (Table 2). All previously  
764 unpublished genomes used in this study are available in a Figshare folder  
765 10.6084/m9.figshare.22339555 and all viral genomes are available here:  
766 10.6084/m9.figshare.22738568.

767

768

769 **CODE AVAILABILITY STATEMENT**

770 The code used in this publication is based on previously published code. Please refer to the  
771 method section for information regarding the software and versions used.

772

773 **REFERENCES**

- 774 1. Garneau, J. E. *et al.* The CRISPR/Cas bacterial immune system cleaves bacteriophage and  
775 plasmid DNA. *Nature* **468**, 67 (2010).
- 776 2. Andersson, A. F. & Banfield, J. F. Virus population dynamics and acquired virus resistance in  
777 natural microbial communities. *Science* **320**, 1047–50 (2008).
- 778 3. Koonin, E. V. & Makarova, K. S. Evolutionary plasticity and functional versatility of CRISPR  
779 systems. *PLOS Biol.* **20**, e3001481 (2022).
- 780 4. Makarova, K. S. *et al.* An updated evolutionary classification of CRISPR–Cas systems. *Nat.*  
781 *Rev. Microbiol.* **13**, 722 (2015).
- 782 5. Makarova, K. S. *et al.* Evolutionary classification of CRISPR–Cas systems: a burst of class 2  
783 and derived variants. *Nat. Rev. Microbiol.* **18**, 67–83 (2020).
- 784 6. Maniv, I., Jiang, W., Bikard, D. & Marraffini, L. A. Impact of different target sequences on  
785 type III CRISPR-Cas immunity. *J. Bacteriol.* **198**, 941 (2016).
- 786 7. Marraffini, L. A. & Sontheimer, E. J. Self versus non-self discrimination during CRISPR RNA-  
787 directed immunity. *Nature* **463**, 568–571 (2010).
- 788 8. Dombrowski, N., Lee, J.-H., Williams, T. A., Offre, P. & Spang, A. Genomic diversity, lifestyles  
789 and evolutionary origins of DPANN archaea. *FEMS Microbiol. Lett.* **366**, (2019).
- 790 9. Rinke, C. *et al.* Insights into the phylogeny and coding potential of microbial dark matter.  
791 *Nature* **499**, 431–437 (2013).
- 792 10. Castelle, C. J. *et al.* Biosynthetic capacity, metabolic variety and unusual biology in the CPR  
793 and DPANN radiations. *Nat. Rev. Microbiol.* **16**, 629–645 (2018).

- 794 11. Sakai, H. D. *et al.* Insight into the symbiotic lifestyle of DPANN archaea revealed by  
795 cultivation and genome analyses. *Proc. Natl. Acad. Sci.* **119**, e2115449119 (2022).
- 796 12. Jahn, U. *et al.* *Nanoarchaeum equitans* and *Ignicoccus hospitalis*: new insights into a unique,  
797 intimate association of two archaea. *J. Bacteriol.* **190**, 1743–1750 (2008).
- 798 13. Huber, H. *et al.* A new phylum of Archaea represented by a nanosized hyperthermophilic  
799 symbiont. *Nature* **417**, 63–67 (2002).
- 800 14. Schwank, K. *et al.* An archaeal symbiont-host association from the deep terrestrial  
801 subsurface. *ISME J.* **13**, 2135–2139 (2019).
- 802 15. Hamm, J. N. *et al.* Unexpected host dependency of Antarctic *Nanohaloarchaeota*. *Proc.*  
803 *Natl. Acad. Sci.* **116**, 14661 (2019).
- 804 16. Munson-McGee, J. H. *et al.* *Nanoarchaeota*, their *Sulfolobales* host, and *Nanoarchaeota*  
805 virus distribution across Yellowstone National Park hot springs. *Appl. Environ. Microbiol.* **81**,  
806 7860–7868 (2015).
- 807 17. Jarett, J. K. *et al.* Single-cell genomics of co-sorted *Nanoarchaeota* suggests novel putative  
808 host associations and diversification of proteins involved in symbiosis. *Microbiome* **6**, 161  
809 (2018).
- 810 18. Wurch, L. *et al.* Genomics-informed isolation and characterization of a symbiotic  
811 *Nanoarchaeota* system from a terrestrial geothermal environment. *Nat. Commun.* **7**, 12115  
812 (2016).
- 813 19. Hamm, J. N. *et al.* The intracellular lifestyle of an archaeal symbiont. *bioRxiv*  
814 2023.02.24.529834 (2023).

- 815 20. Probst, A. J. *et al.* Differential depth distribution of microbial function and putative  
816 symbionts through sediment-hosted aquifers in the deep terrestrial subsurface. *Nat.*  
817 *Microbiol.* **3**, 328–336 (2018).
- 818 21. Heimerl, T. *et al.* A complex endomembrane system in the archaeon *Ignicoccus hospitalis*  
819 tapped by *Nanoarchaeum equitans*. *Front. Microbiol.* **8**, 1072 (2017).
- 820 22. Comolli, L. R. & Banfield, J. F. Inter-species interconnections in acid mine drainage microbial  
821 communities. *Front. Microbiol.* **5**, 367 (2014).
- 822 23. Baker, B. J. *et al.* Enigmatic, ultrasmall, uncultivated Archaea. *Proc. Natl. Acad. Sci.* **107**,  
823 8806–8811 (2010).
- 824 24. Hensdorf, A. W. *et al.* Potential for microbial H<sub>2</sub> and metal transformations associated with  
825 novel bacteria and archaea in deep terrestrial subsurface sediments. *ISME J.* **11**, 1915–1929  
826 (2017).
- 827 25. Probst, A. J. *et al.* Biology of a widespread uncultivated archaeon that contributes to carbon  
828 fixation in the subsurface. *Nat. Commun.* **5**, 5497 (2014).
- 829 26. Probst, A. J. *et al.* Genomic resolution of a cold subsurface aquifer community provides  
830 metabolic insights for novel microbes adapted to high CO<sub>2</sub> concentrations. **19**, 459–474  
831 (2017).
- 832 27. Emerson, J. B., Thomas, B. C., Alvarez, W. & Banfield, J. F. Metagenomic analysis of a high  
833 carbon dioxide subsurface microbial community populated by chemolithoautotrophs and  
834 bacteria and archaea from candidate phyla. *Environ. Microbiol.* **18**, 1686–1703 (2016).
- 835 28. Rahlff, J. *et al.* Lytic archaeal viruses infect abundant primary producers in Earth's crust.  
836 *Nat. Commun.* **12**, 4642 (2021).



- 837 29. Wimmer, F., Mougiakos, I., Englert, F. & Beisel, C. L. Rapid cell-free characterization of  
838 multi-subunit CRISPR effectors and transposons. *Mol. Cell* **82**, 1210-1224.e6 (2022).
- 839 30. Marshall, R. *et al.* Rapid and scalable characterization of CRISPR technologies using an *E. coli*  
840 cell-free transcription-translation system. *Mol. Cell* **69**, 146-157.e3 (2018).
- 841 31. Heussler, G. E. & O'Toole, G. A. Friendly Fire: Biological functions and consequences of  
842 chromosomal targeting by CRISPR-Cas systems. *J. Bacteriol.* **198**, 1481–1486 (2016).
- 843 32. Stern, A., Keren, L., Wurtzel, O., Amitai, G. & Sorek, R. Self-targeting by CRISPR: gene  
844 regulation or autoimmunity? *Trends Genet.* **26**, 335–340 (2010).
- 845 33. Akujkar, M. & Lovley, D. R. Interference with histidyl-tRNA synthetase by a CRISPR spacer  
846 sequence as a factor in the evolution of *Pelobacter carbinolicus*. *BMC Evol. Biol.* **10**, 230  
847 (2010).
- 848 34. Bhaya, D., Davison, M. & Barrangou, R. CRISPR-Cas systems in Bacteria and Archaea:  
849 Versatile small RNAs for adaptive defense and regulation. *Annu. Rev. Genet.* **45**, 273–297  
850 (2011).
- 851 35. Wilson, G. G. Organization of restriction-modification systems. *Nucleic Acids Res.* **19**, 2539–  
852 2566 (1991).
- 853 36. Bornemann, T. L. V. *et al.* Genetic diversity in terrestrial subsurface ecosystems impacted by  
854 geological degassing. *Nat. Commun.* **13**, 284 (2022).
- 855 37. Turgeman-Grott, I. *et al.* Pervasive acquisition of CRISPR memory driven by inter-species  
856 mating of archaea can limit gene transfer and influence speciation. *Nat. Microbiol.* **4**, 177–  
857 186 (2019).

- 858 38. Stachler, A.-E. *et al.* High tolerance to self-targeting of the genome by the endogenous  
859 CRISPR-Cas system in an archaeon. *Nucleic Acids Res.* **45**, 5208–5216 (2017).
- 860 39. Vink, J. N. A., Baijens, J. H. L. & Brouns, S. J. J. PAM-repeat associations and spacer selection  
861 preferences in single and co-occurring CRISPR-Cas systems. *Genome Biol.* **22**, 281 (2021).
- 862 40. Pyenson, N. C., Gayvert, K., Varble, A., Elemento, O. & Marraffini, L. A. Broad targeting  
863 specificity during bacterial type III CRISPR-Cas immunity constrains viral escape. *Cell Host*  
864 *Microbe* **22**, 343-353.e3 (2017).
- 865 41. Chabas, H., Müller, V., Bonhoeffer, S. & Regoes, R. R. Epidemiological and evolutionary  
866 consequences of different types of CRISPR-Cas systems. *PLOS Comput. Biol.* **18**, e1010329  
867 (2022).
- 868 42. Brodt, A., Lurie-Weinberger, M. N. & Gophna, U. CRISPR loci reveal networks of gene  
869 exchange in archaea. *Biol. Direct* **6**, 65 (2011).
- 870 43. Paper, W. *et al.* *Ignicoccus hospitalis* sp. nov., the host of ‘*Nanoarchaeum equitans*’. *Int. J.*  
871 *Syst. Evol. Microbiol.* **57**, 803–808 (2007).
- 872 44. Dombrowski, N., Teske, A. P. & Baker, B. J. Expansive microbial metabolic versatility and  
873 biodiversity in dynamic Guaymas Basin hydrothermal sediments. *Nat. Commun.* **9**, 4999  
874 (2018).
- 875 45. Hohenester, E. & Yurchenco, P. D. Laminins in basement membrane assembly. *Cell Adhes.*  
876 *Migr.* **7**, 56–63 (2013).
- 877 46. Hohenester, E. Laminin G-like domains: dystroglycan-specific lectins. *Seq. Topol.* ●  
878 *Carbohydr.* **56**, 56–63 (2019).

- 879 47. Benner, S. A., Ellington, A. D. & Tauer, A. Modern metabolism as a palimpsest of the RNA  
880 world. *Proc. Natl. Acad. Sci.* **86**, 7054–7058 (1989).
- 881 48. Joshi, N. A. & Fass, J. N. Sickle: A sliding-window, adaptive, quality-based trimming tool for  
882 FastQ files (Version 1.33) [Software]. (2011).
- 883 49. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new versatile  
884 metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).
- 885 50. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**,  
886 357–359 (2012).
- 887 51. Bornemann, T. L. V., Esser, S. P., Stach, T. L., Burg, T. & Probst, A. J. uBin – a manual refining  
888 tool for genomes from metagenomes. *Environ. Microbiol.* (2023).
- 889 52. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:  
890 assessing the quality of microbial genomes recovered from isolates, single cells, and  
891 metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
- 892 53. Eddy, S. R. Accelerated Profile HMM Searches. *PLOS Comput. Biol.* **7**, e1002195 (2011).
- 893 54. Darling, A. E. *et al.* PhyloSift: phylogenetic analysis of genomes and metagenomes. *PeerJ* **2**,  
894 e243 (2014).
- 895 55. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
896 *Nucleic Acids Res.* **32**, 1792–1797 (2004).
- 897 56. Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): A new  
898 software for selection of phylogenetic informative regions from multiple sequence  
899 alignments. *BMC Evol. Biol.* **10**, 210 (2010).

- 900 57. Minh, B. Q. *et al.* IQ-TREE 2: New models and efficient methods for phylogenetic inference  
901 in the genomic era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 902 58. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermini, L. S.  
903 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods* **14**,  
904 587–589 (2017).
- 905 59. Wang, H.-C., Minh, B. Q., Susko, E. & Roger, A. J. Modeling site heterogeneity with posterior  
906 mean site frequency profiles accelerates accurate phylogenomic estimation. *Syst. Biol.* **67**,  
907 216–235 (2017).
- 908 60. Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2: Improving  
909 the ultrafast bootstrap approximation. *Mol. Biol. Evol.* **35**, 518–522 (2017).
- 910 61. Guindon, S. *et al.* New algorithms and methods to estimate maximum-likelihood  
911 phylogenies: Assessing the performance of PhyML 3.0. *Syst. Biol.* **59**, 307–321 (2010).
- 912 62. Anisimova, M., Gil, M., Dufayard, J.-F., Dessimoz, C. & Gascuel, O. Survey of branch support  
913 methods demonstrates accuracy, power, and robustness of fast likelihood-based  
914 approximation schemes. *Syst. Biol.* **60**, 685–699 (2011).
- 915 63. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new  
916 developments. *Nucleic Acids Res.* **47**, W256–W259 (2019).
- 917 64. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND.  
918 *Nat. Methods* **12**, 59 (2014).
- 919 65. Gouy, M., Tannier, E., Comte, N. & Parsons, D. P. Seaview Version 5: A multiplatform  
920 software for multiple sequence alignment, molecular phylogenetic analyses, and tree

- 921 reconciliation. in *Multiple Sequence Alignment: Methods and Protocols* (ed. Katoh, K.) 241–  
922 260 (Springer US, 2021).
- 923 66. Couvin, D. *et al.* CRISPRCasFinder, an update of CRISPRFinder, includes a portable version,  
924 enhanced performance and integrates search for Cas proteins. *Nucleic Acids Res* **46**, W246–  
925 W251 (2018).
- 926 67. Moller, A. G. & Liang, C. MetaCRASST: reference-guided extraction of CRISPR spacers from  
927 unassembled metagenomes. *PeerJ* **5**:e3788, e3788 (2017).
- 928 68. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation  
929 sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).
- 930 69. Biswas, A., Fineran, P. C. & Brown, C. M. Accurate computational prediction of the  
931 transcribed strand of CRISPR non-coding RNAs. *Bioinformatics* **30**, 1805–13 (2014).
- 932 70. Roux, S., Enault, F., Hurwitz, B. L. & Sullivan, M. B. VirSorter: mining viral signal from  
933 microbial genomic data. *PeerJ* **3**:e985, e985 (2015).
- 934 71. Xie, Z. & Tang, H. ISEScan: automated identification of insertion sequence elements in  
935 prokaryotic genomes. *Bioinformatics* **33**, 3340–3347 (2017).
- 936 72. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**,  
937 2460–2461 (2010).
- 938 73. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search  
939 tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 940 74. Nayfach, S. *et al.* CheckV assesses the quality and completeness of metagenome-assembled  
941 viral genomes. *Nat. Biotechnol.* **39**, 578–585 (2021).

- 942 75. Cook, R. *et al.* INfrastructure for a PHAge REference Database: Identification of large-scale  
943 biases in the current collection of cultured phage genomes. *PHAGE* **2**, 214–223 (2021).
- 944 76. Bolduc, B. *et al.* vConTACT: an iVirus tool to classify double-stranded DNA viruses that infect  
945 Archaea and Bacteria. *PeerJ* **5**, e3243 (2017).
- 946 77. Bin Jang, H. *et al.* Taxonomic assignment of uncultivated prokaryotic virus genomes is  
947 enabled by gene-sharing networks. *Nat. Biotechnol.* **37**, 632–639 (2019).
- 948 78. Meier-Kolthoff, J. P. & Göker, M. VICTOR: genome-based phylogeny and classification of  
949 prokaryotic viruses. *Bioinformatics* **33**, 3396–3404 (2017).
- 950 79. Moraru, C., Varsani, A. & Kropinski, A. M. VIRIDIC—A novel tool to calculate the  
951 intergenomic similarities of prokaryote-infecting Viruses. *Viruses* **12**, (2020).
- 952 80. Meier-Kolthoff, J. P., Auch, A. F., Klenk, H.-P. & Göker, M. Genome sequence-based species  
953 delimitation with confidence intervals and improved distance functions. *BMC*  
954 *Bioinformatics* **14**, 60 (2013).
- 955 81. Lefort, V., Desper, R. & Gascuel, O. FastME 2.0: A comprehensive, accurate, and fast  
956 distance-based phylogeny inference program. *Mol. Biol. Evol.* **32**, 2798–2800 (2015).
- 957 82. Farris, J. S. Estimating phylogenetic trees from distance matrices. *Am. Nat.* **106**, 645–668  
958 (1972).
- 959 83. Shannon, P. *et al.* Cytoscape: A software environment for integrated models of  
960 biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
- 961 84. Nishimura, Y. *et al.* ViPTree: the viral proteomic tree server. *Bioinformatics* **33**, 2379–2380  
962 (2017).

- 963 85. Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078–  
964 2079 (2009).
- 965 86. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic  
966 features. *Bioinformatics* **26**, 841–842 (2010).
- 967 87. Team, R. C. *R: A language and environment for statistical computing.* (Vienna, Austria,  
968 2013).
- 969 88. Zhang, Y. & Sievert, S. Pan-genome analyses identify lineage- and niche-specific markers of  
970 evolution and adaptation in *Epsilonproteobacteria*. *Front. Microbiol.* **5**, 110 (2014).
- 971 89. Dufault-Thompson, K., Steffensen, J. L. & Zhang, Y. Using PSAMM for the curation and  
972 analysis of genome-scale metabolic models. in *Metabolic Network Reconstruction and*  
973 *Modeling: Methods and Protocols* (ed. Fondi, M.) 131–150 (Springer New York, 2018).
- 974 90. Steffensen, J. L., Dufault-Thompson, K. & Zhang, Y. PSAMM: A portable system for the  
975 analysis of metabolic models. *PLoS Comput. Biol.* **12**, e1004732–e1004732 (2016).
- 976 91. Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein  
977 or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).
- 978 92. Gonnerman, M. C., Benedict, M. N., Feist, A. M., Metcalf, W. W. & Price, N. D. Genomically  
979 and biochemically accurate metabolic reconstruction of *Methanosarcina barkeri* Fusaro,  
980 iIMG746. *Biotechnol. J.* **8**, 1070–1079 (2013).
- 981 93. Goyal, N., Widiastuti, H., Karimi, I. A. & Zhou, Z. A genome-scale metabolic model of  
982 *Methanococcus maripaludis* S<sub>2</sub> for CO<sub>2</sub> capture and conversion to methane. *Mol. Biosyst.*  
983 **10**, 1043–1054 (2014).

- 984 94. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new perspectives  
985 on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–D361 (2017).
- 986 95. Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically  
987 annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids*  
988 *Res.* **47**, D309–D314 (2019).
- 989 96. Saier, M. H., Jr *et al.* The transporter classification database (TCDB): recent advances.  
990 *Nucleic Acids Res.* **44**, D372–D379 (2016).
- 991 97. Neidhardt, F. C., Neidhardt, F. C. N., Ingraham, J. L. & Schaechter, M. Physiology of the  
992 bacterial cell: A molecular approach. (Sinauer Associates, 1990).
- 993 98. Nelson, D. L., Nelson, R. D. & Cox, M. M. *Lehninger Principles of Biochemistry, Fourth Edition*  
994 *+ Lecture Notebook.* (W.H. Freeman, 2004).
- 995 99. Biswas, A., Gagnon, J. N., Brouns, S. J. J., Fineran, P. C. & Brown, C. M. CRISPRTarget:  
996 bioinformatic prediction and analysis of crRNA targets. *RNA Biol.* **10**, 817–827 (2013).
- 997 100. Crooks, G. E., Hon, G., Chandonia, J.-M. & Brenner, S. E. WebLogo: A sequence logo  
998 generator. *Genome Res.* **14**, 1188–1190 (2004).
- 999 101. Schneider, T. D. & Stephens, R. M. Sequence logos: a new way to display consensus  
1000 sequences. *Nucleic Acids Res.* **18**, 6097–6100 (1990).
- 1001 102. Koboldt, D. C. *et al.* VarScan 2: Somatic mutation and copy number alteration discovery  
1002 in cancer by exome sequencing. *Genome Res.* **22**, 568–576 (2012).
- 1003 103. Oberortner, E., Cheng, J.-F., Hillson, N. J. & Deutsch, S. Streamlining the design-to-build  
1004 transition with build-optimization software tools. *ACS Synth. Biol.* **6**, 485–496 (2017).



- 1005 104. Garamella, J., Marshall, R., Rustad, M. & Noireaux, V. The All E. coli TX-TL Toolbox 2.0: A  
1006 Platform for Cell-Free Synthetic Biology. *ACS Synth. Biol.* **5**, 344–355 (2016).
- 1007 105. Shin, J. & Noireaux, V. An *E. coli* cell-free expression toolbox: Application to synthetic  
1008 gene circuits and artificial cells. *ACS Synth. Biol.* **1**, 29–41 (2012).
- 1009 106. Leenay, R. T. *et al.* Identifying and visualizing functional PAM diversity across CRISPR-Cas  
1010 systems. *Mol. Cell* **62**, 137–147 (2016).
- 1011 107. Ondov, B. D., Bergman, N. H. & Phillippy, A. M. Interactive metagenomic visualization in  
1012 a Web browser. *BMC Bioinformatics* **12**, 385 (2011).
- 1013 108. Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to classify  
1014 genomes with the Genome Taxonomy Database. *Bioinformatics* **36**, 1925–1927 (2020).
- 1015 109. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny  
1016 substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996–1004 (2018).
- 1017 110. Parks, D. H. *et al.* A complete domain-to-species taxonomy for Bacteria and Archaea.  
1018 *Nat. Biotechnol.* **38**, 1079–1086 (2020).
- 1019 111. Chen, I.-M. A. *et al.* IMG/M v.5.0: an integrated data management and comparative  
1020 analysis system for microbial genomes and microbiomes. *Nucleic Acids Res.* **47**, D666–D677  
1021 (2019).
- 1022 112. Sharrar, A. M. *et al.* Novel large sulfur bacteria in the metagenomes of groundwater-fed  
1023 chemosynthetic microbial mats in the Lake Huron Basin. *Front. Microbiol.* **8**, 791 (2017).
- 1024 113. Bird, J. T., Baker, B. J., Probst, A. J., Podar, M. & Lloyd, K. G. Culture independent  
1025 genomic comparisons reveal environmental adaptations for Altiarchaeales. *Front.*  
1026 *Microbiol.* **7**, (2016).

1027 114. Posit team. Rstudio: Integrated development environment for R. (2022).

1028 115. Esser, S. P. & Probst, A. J. Genomes of *Ca. Altiarchaeum* and *Ca. Huberiarchoaeum* from  
1029 Crystal Geyser and Horonobe Underground Research Laboratory. (2023)  
1030 doi:[10.6084/m9.figshare.22339555](https://doi.org/10.6084/m9.figshare.22339555).

1031 116. Esser, S. P., Rahlff, J. & Probst, A. J. Viral operational taxonomic units (vOTUs) derived from  
1032 metagenomes from Crystal Geyser. (2023) doi:[10.6084/m9.figshare.22738568](https://doi.org/10.6084/m9.figshare.22738568).

1033 117. Turzynski, V., Esser, S. P. & Probst, A. J. FISH images of *Ca. Altiarchaeum* and *Ca.*  
1034 *Huberiarchoaeum* in Crystal Geyser and Horonobe Underground Research Laboratory. (2023)  
1035 doi:[10.6084/m9.figshare.22739849](https://doi.org/10.6084/m9.figshare.22739849).

1036

1037

1038

1039 **Figure 1. Phylogenetic positioning of *Ca. Altiarchaea* and *Ca. Huberarchaea*, sampling**  
1040 **locations, FISH analysis, and CRISPR-Cas targets | (A)** Phylogenetic tree of archaea highlighting  
1041 *Candidatus Altiarchaeum* and *Candidatus Huberiarchaeum* of the sampling locations Crystal  
1042 Geyser (CG, Utah, USA, orange) and Horonobe Underground Research Laboratory (HURL,  
1043 Hokkaido, Japan, green). Fluorescence pictures show *Ca. Altiarchaeum* (blue; H - host) as host and  
1044 its episymbiont *Ca. Huberiarchaeum* (orange; S - symbiont) in the respective ecosystems. Scale  
1045 bar 1  $\mu\text{m}$ . **(B)** *Ca. Altiarchaea* CRISPR systems, their associated conserved direct repeat (DR)  
1046 sequences (with exception of a point mutation marked in red), and the number of spacer clusters  
1047 (97% nucleotide identity) arising from the two sampling sites. **(C)** Logarithmic number of centroid  
1048 spacers derived from spacer clusters matching 64 extracted viral sequences (total number of  
1049 spacer matches: 0 of unassigned CRISPR system and 16561 of CRISPR system IB), 17 binned  
1050 genomes of *Ca. Altiarchaeum crystalense* (total number of spacer matches: 115 of unassigned  
1051 CRISPR system and 1,311 of CRISPR system IB) and 11 binned genomes of *Ca. Huberiarchaeum*  
1052 *crystalense* (total number of spacer matches: 0 of unassigned CRISPR system and 1,445 of CRISPR  
1053 system IB) originating from the CG site (Table S2). Spacers were derived from the complete 66-  
1054 sample metagenomic dataset. **(D)** Percentage of CRISPR system I-B spacer cluster abundances  
1055 matching to organisms that were previously detected in this ecosystem at the CG site. Listed are  
1056 the logarithmic genome abundances of the respective organisms. Error bars denote the standard  
1057 deviation of the abundance of matching spacer clusters for samples CG05, CG08, and CG16 of the  
1058 year 2015. These were displayed because also transcriptomic data was available. The dataset of  
1059 HURL is referring to one metagenome, as no other data was available. (Means and standard  
1060 deviation: CG *Altiarchaeum crystalense*:  $2.69 \pm 0.21$ ,  $1.93 \pm 0.24$ ; *Huberiarchaeum crystalense*:  
1061  $2.99 \pm 0.23$ ,  $2.19 \pm 0.23$ ; HURL *Altiarchaeum horonobense*: 0.029; *Huberiarchaeum julieae*: 0.019)  
1062 **(E)** Logarithmic number of centroid spacers derived from spacer clusters matching extracted viral  
1063 sequences (total number of spacer matches: 64 of unassigned CRISPR system and 22 of CRISPR  
1064 system IB), two binned genomes of *Ca. Altiarchaeum horonobense* (total number of spacer  
1065 matches: 19 of unassigned CRISPR system and 2 of CRISPR system IB) and one binned genome of  
1066 *Ca. Huberiarchaeum julieae* (total number of spacer matches: 7 of unassigned CRISPR system and

1067 0 of CRISPR system IB) originating from the HURL site. Spacers were derived from one  
1068 metagenomic dataset.

1069

1070 **Figure 2. Example of *Ca. Altiarchaea* CRISPR-Cas type IB loci, gene targets on host and**  
1071 ***episymbiont* genomes, and metabolic interactions between *Ca. Altiarchaea* and *Ca.***  
1072 ***Huberiarchaea* as inferred from genome-scaled metabolic modeling | (A)** Example of CRISPR  
1073 system I-B locus of *Ca. A. crystalense* with assembled CRISPR array from a single amplified genome  
1074 (accession no. 1088571). Red box highlights the analysed CRISPR array bearing the repeat  
1075 sequence GTTTAAATCGTACTATGTAGTATGGAAAC and its respective spacers within the array. **(B)**  
1076 Example of a *Ca. A. crystalense* DNA polymerase II large subunit locus self-targeted by  
1077 *altiarchaeal* spacers extracted from metagenomes (accession no. 2786546692). Red boxes on the  
1078 genomic region highlight spacer matching regions. Yellow genes are annotated as  
1079 uncharacterized proteins. **(C)** Example of a *Ca. H. crystalense* genome (accession no. 2785510793)  
1080 partially matched by *Ca. A. crystalense* spacers at the genetic loci of the 30S ribosomal protein  
1081 S11, CTP synthase, and an uncharacterized protein. **(D)** Example of a *Ca. H. crystalense* SAG  
1082 (accession no. 1088571) partially matched by *Ca. A. crystalense* spacers at the genetic loci of  
1083 uncharacterized proteins. **(E)** Metabolic interactions between *Ca. Altiarchaeum* and *Ca.*  
1084 *Huberiarchaeum* in Crystal Geysers, CG (17 genomes of *Ca. A. crystalense* and eleven of *Ca. H.*  
1085 *crystalense*, and spacers extracted from transcriptomes) and in **(F)** Horonobe Underground  
1086 Research Laboratory, HURL (one genome of *A. horonobense* and one *Ca. H. julieae*). Solid arrows  
1087 denote exchanges of putative essential metabolites between *Ca. Altiarchaeum* and *Ca.*  
1088 *Huberiarchaeum*. Dashed arrows indicate exchange of metabolites that are only required when  
1089 CRISPR spacers attack certain target genes (type I-B labeled with red diamonds and the  
1090 unassigned type labeled with green diamonds). **(E-F)** While most compounds were produced by  
1091 *Ca. Altiarchaea*, the production of dUMP requires an essential gene, ⑤-dCMP deaminase (EC  
1092 3.5.4.12), in *Huberiarchaea*. Circled numbers indicate key enzymes involved in symbiotic  
1093 metabolic interactions at CG: ①-Phenylalanyl-tRNA synthetase (EC 6.1.1.20), Lysyl-tRNA  
1094 synthetase (EC 6.1.1.6); ②, ③-(d)NDP kinase (EC 2.7.4.6); ④-dCMP kinase (EC 2.7.4.25); ⑤-  
1095 dCMP deaminase (EC 3.5.4.12); ⑥, ⑦-dTMP synthase (EC 2.1.1.45); ⑧-FAD-dependent dTMP

1096 synthase (EC 2.1.1.148). The production of tetrahydrofolate (THF) requires an essential gene  
1097 encoded by *Ca. Huberiarchoaeum julieae*, ⑬-dTMP synthase (EC 2.1.1.45). Circled numbers  
1098 denote key enzymes involved in the symbiotic metabolic interactions at HURL: ⑬-dTMP synthase  
1099 (EC 2.1.1.45); ⑫-FAD-dependent dTMP synthase (EC 2.1.1.148); ⑨, ⑩-dCTP deaminase (EC  
1100 3.5.4.13); ⑪-dUTPase (EC 3.6.1.23); and ⑭-Dihydrofolate synthase (EC 6.3.2.12).

1101

1102

1103 **Figure 3. Illustration of the newly discovered functionality of CRISPR-Cas systems within *Ca.***  
1104 ***Altiarchaea* | a. Viral targeting:** CRISPR-Cas system targets the genomes of MGEs that infect the  
1105 cell (current state of knowledge). **b. Targeting of episymbiont:** CRISPR-Cas system targets the  
1106 genome of the episymbiont *Ca. Huberiarchaeum* to defend against the parasite. **c. Self-targeting**  
1107 **and respective metabolic complementation:** Self-targeting of CRISPR-Cas in *Altiarchaea*  
1108 mediates metabolic patchiness, which is complemented by the episymbiont metabolism, leading  
1109 to mutualism. Please note, that this mutualism might be limited to a subset of organisms in the  
1110 host population. Arrows symbolize spacer-protospacer interactions. The Figure was created with  
1111 Biorender.com.

1112  
1113 **Figure 4. Directed spacer interaction of DPANN archaea derived from the analysis of 7,012**  
1114 **publicly available archaeal genomes |** Nodes correspond to archaeal genomes. Boomerang and  
1115 linear grey arrows indicate self-targeting and non-self (including interspecies) targeting spacers,  
1116 respectively. With the exception of *Thermoprotei* and *Bathyarchaeota*, all of the archaea pictured  
1117 belong to the DPANN superphylum. Colors represent the phylogenetic affiliation of genomes.  
1118 Genomes of *Ca. Altiarchaeum* and *Ca. Huberiarchaeum* derives primarily from CG. Genomes  
1119 coded according to their corresponding ecosystem: CG - Crystal Geyser<sup>20,25,27</sup>; LHB - Lake Huron  
1120 Basin<sup>112</sup>; WOR - White Oak River<sup>113</sup>; GUAY - Guaymas Basin<sup>44</sup>; HURL - Horonobe Underground  
1121 Research Laboratory<sup>24</sup>.

1122  
1123 **Extended Data Fig. 1 | Correlation of repeat abundance and abundance of *Ca. Altiarchaea***  
1124 **genomes.** Spearman rank correlation (two-tailed) of logarithmic abundances of *Ca. A. crystalense*  
1125 and logarithmic abundances of repeat sequences of the unassigned CRISPR array ( $p$ -value  $3.4 e^{-}$   
1126  $16$ ) and the I-B CRISPR system ( $p < 2.2 e^{-16}$ ) in metagenomes from CG ( $n=66$ ). The grey area depicts  
1127 the confidence interval of 0.95. The line indicates that the correlation of the genome abundance  
1128 and repeat abundance is linear. Visualization was performed with R<sup>87,114</sup> (version 3.6.1).

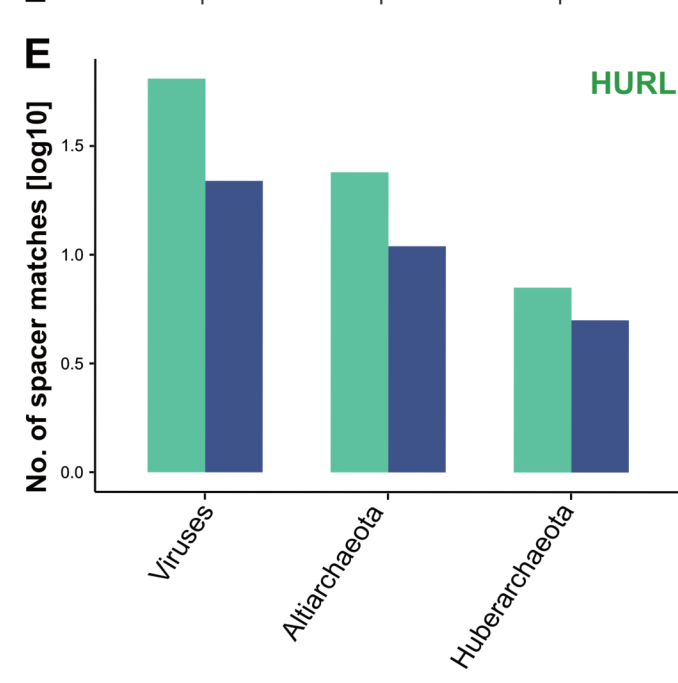
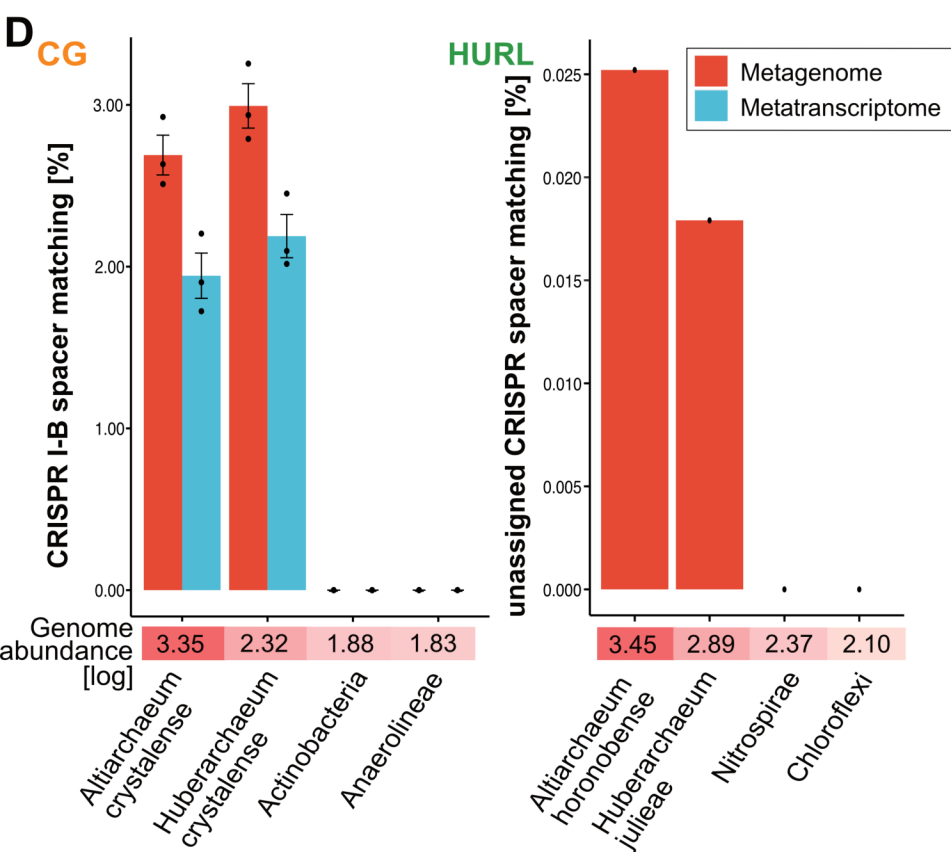
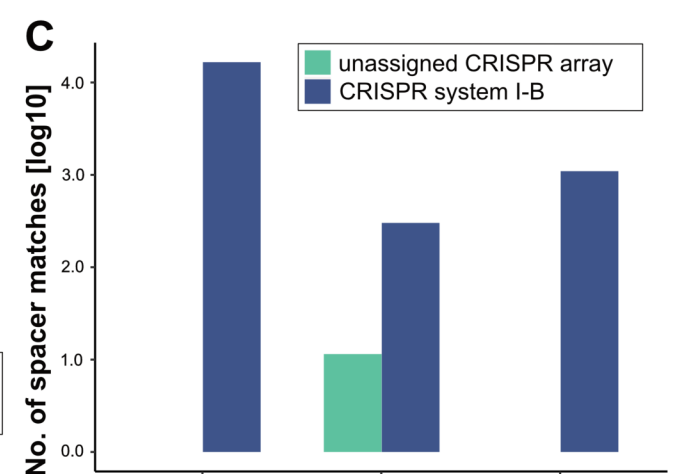
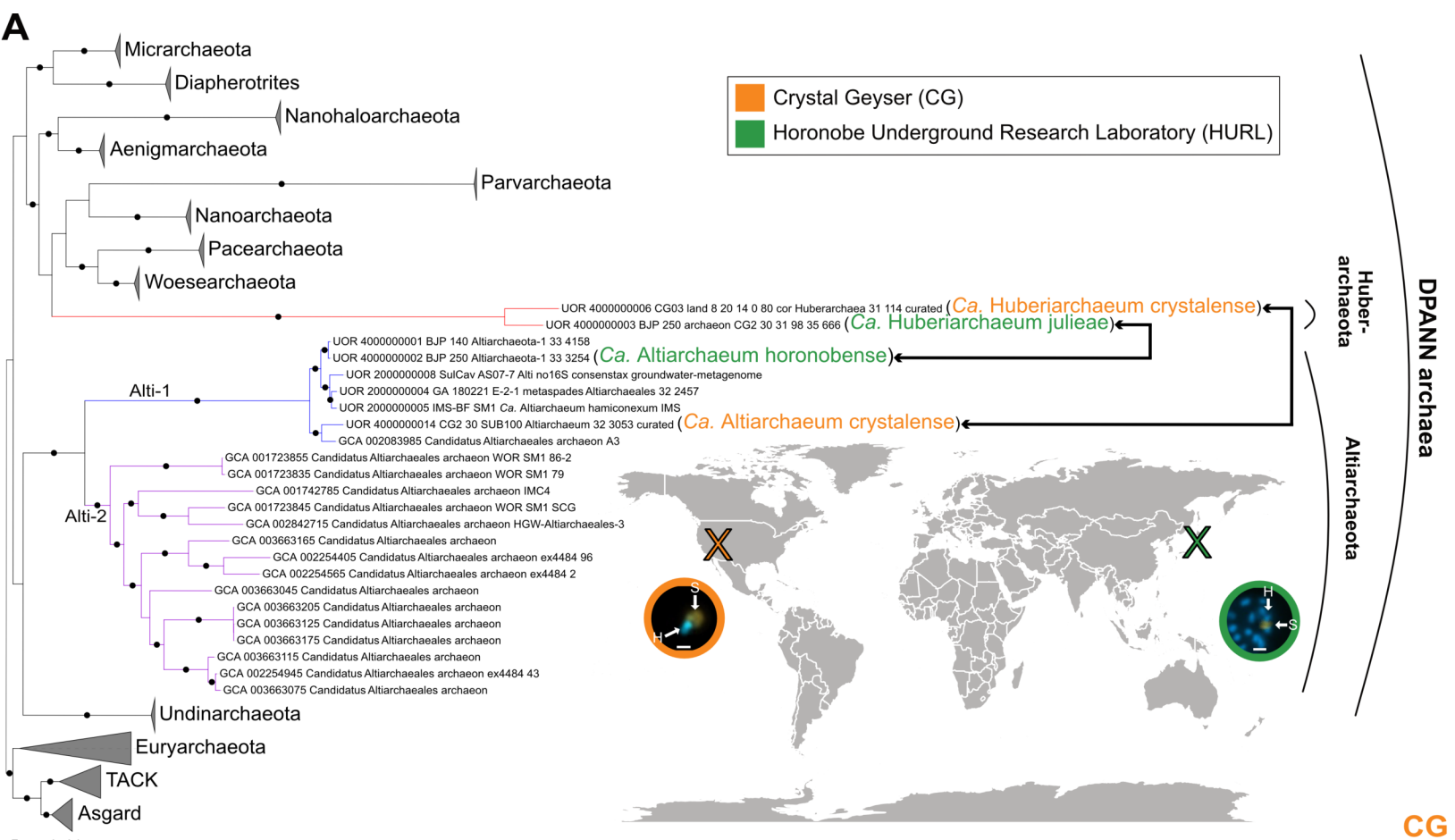
1129 **Extended Data Fig. 2 | Viral clusters predicted by VIRIDIC<sup>79</sup>.** Heatmap showing intergenomic  
1130 similarity for viral scaffolds of viral clusters (VC\_XY) and some singletons (black). Coloring of viral  
1131 OTUs (vOTUs) according to Table S6. VC\_09, \_12, \_13 determined by the other tools were not

1132 *found by VIRIDIC. Only scaffolds with intergenomic similarity of >10 between two viral scaffolds*  
1133 *are shown.*

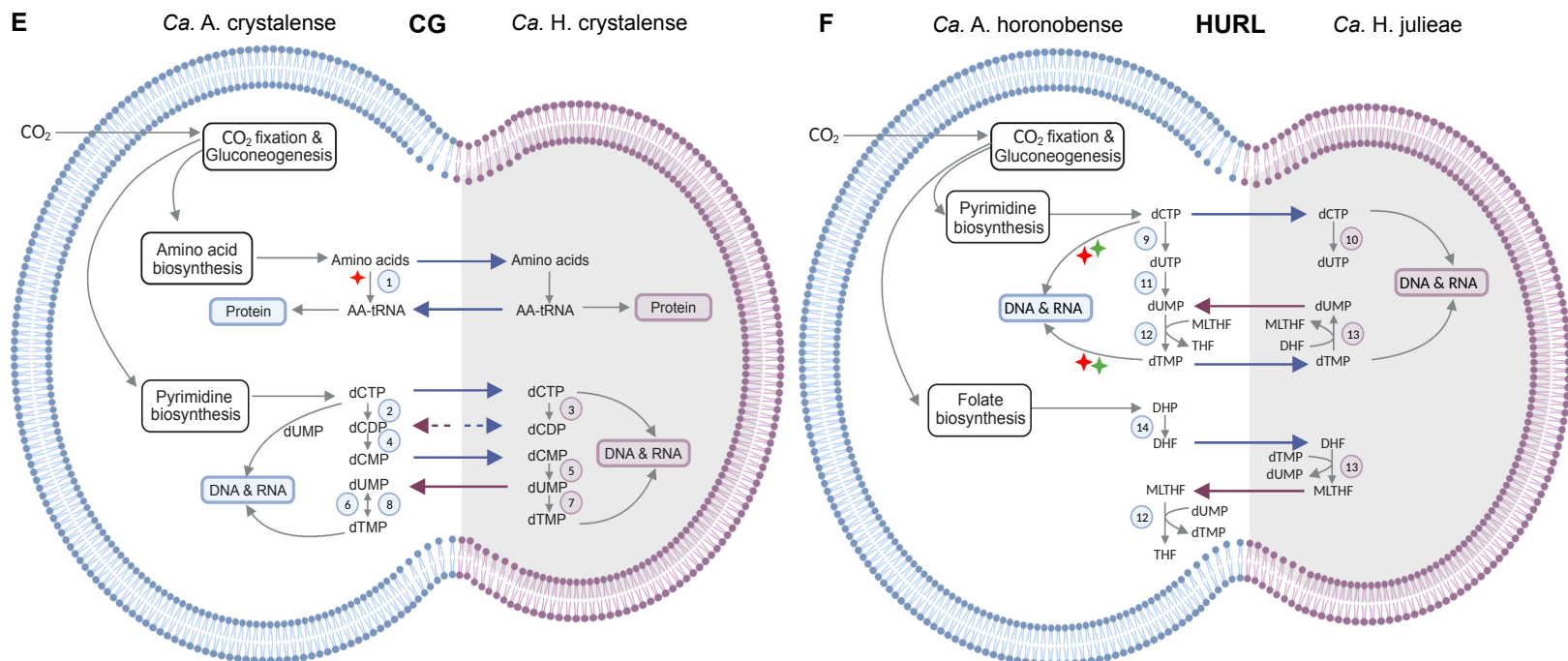
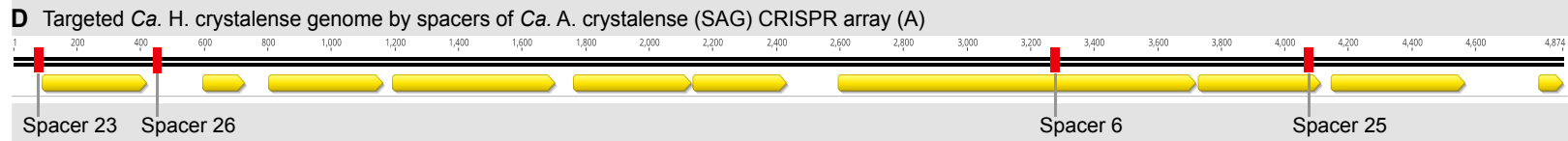
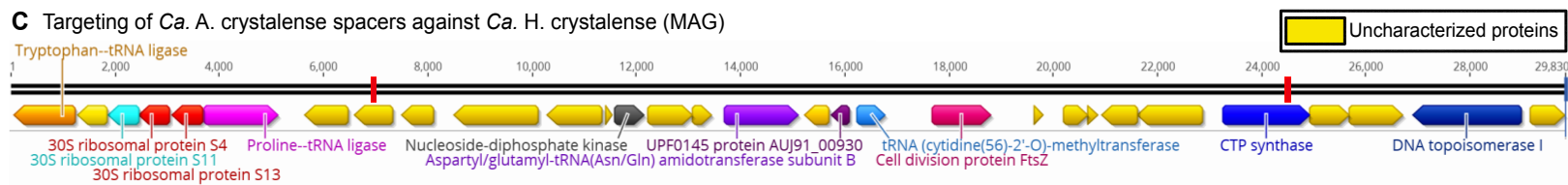
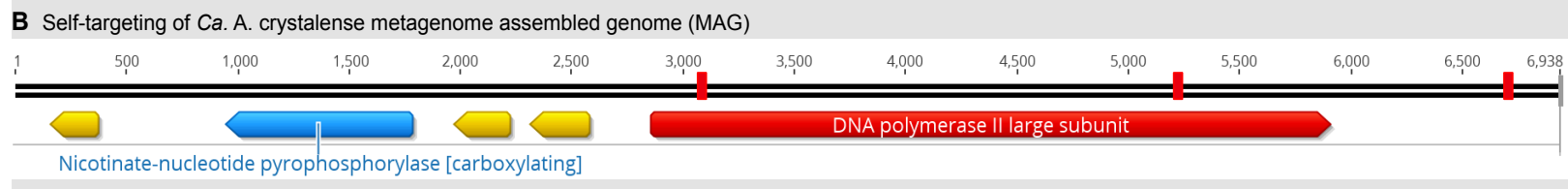
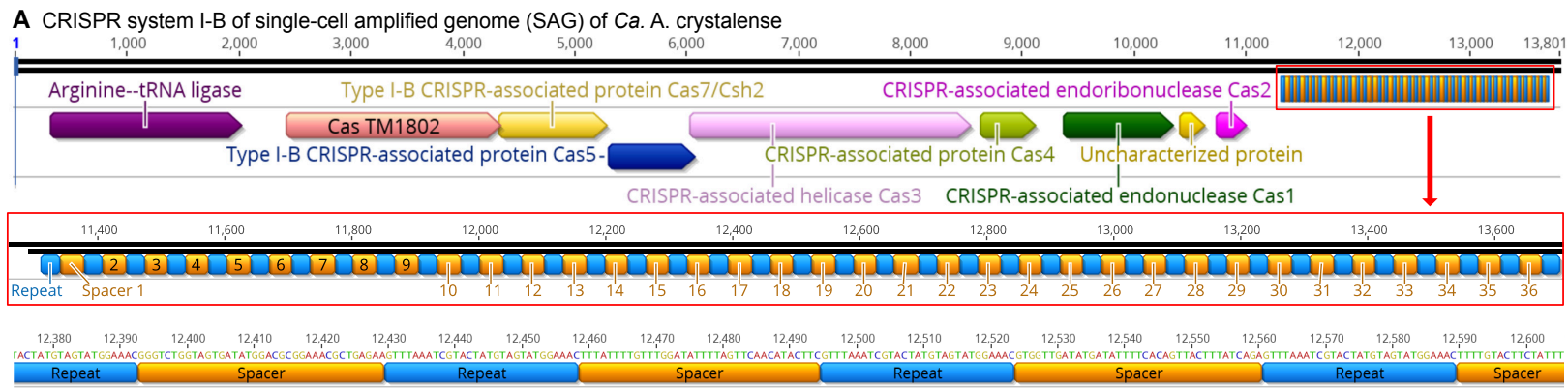
1134 ***Extended Data Fig. 3 | Coverage analyses of scaffolds targeted by spacers from Ca. Altiarchaea.***  
1135 *Coverage changes within targeted regions by CRISPR system IB of Ca. Altiarchaeum and Ca.*  
1136 *Huberiarchaeum based on metagenomic read mapping. The vertically grey marked regions are*  
1137 *spacer targeted regions of either Ca. Altiarchaeum or Ca. Huberiarchaeum, whereby the*  
1138 *horizontally dark grey lines are showing the average coverage of the scaffold. The colored graphs*  
1139 *show the coverage across the spacer targeted region of three samples from the minor eruption*  
1140 *phase, where Ca. Altiarchaeum is the most abundant organism (Fig. S1).*

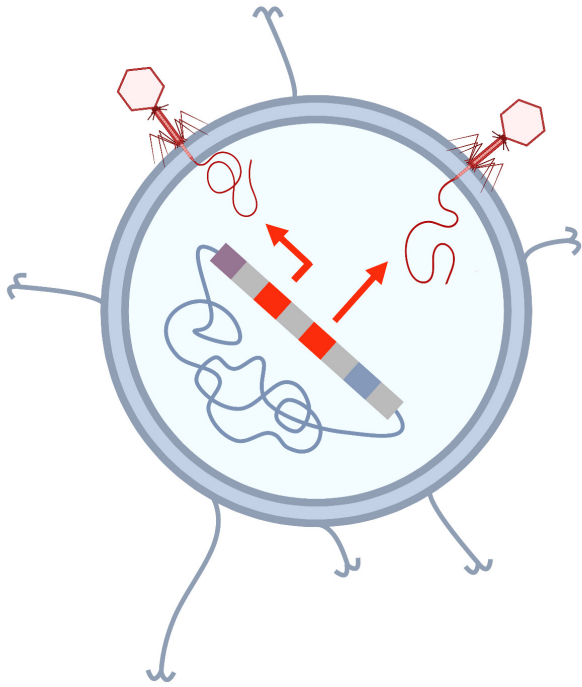
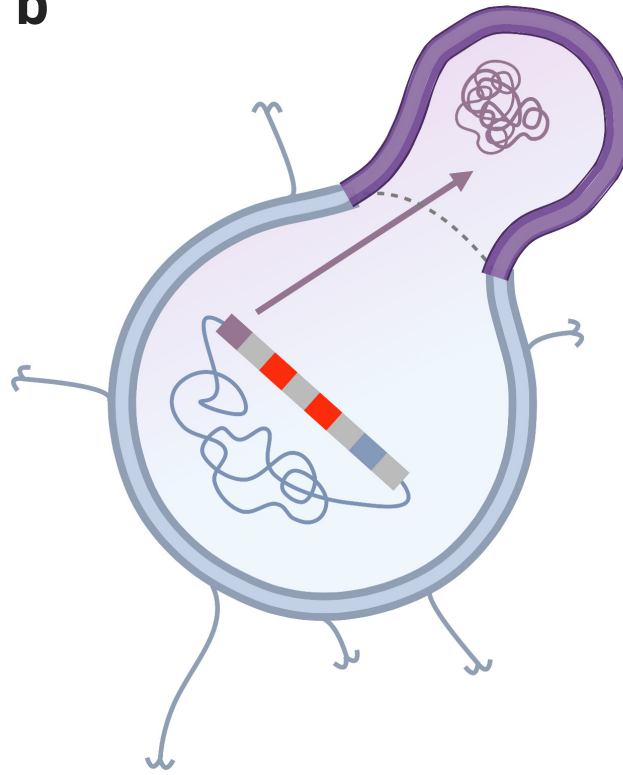
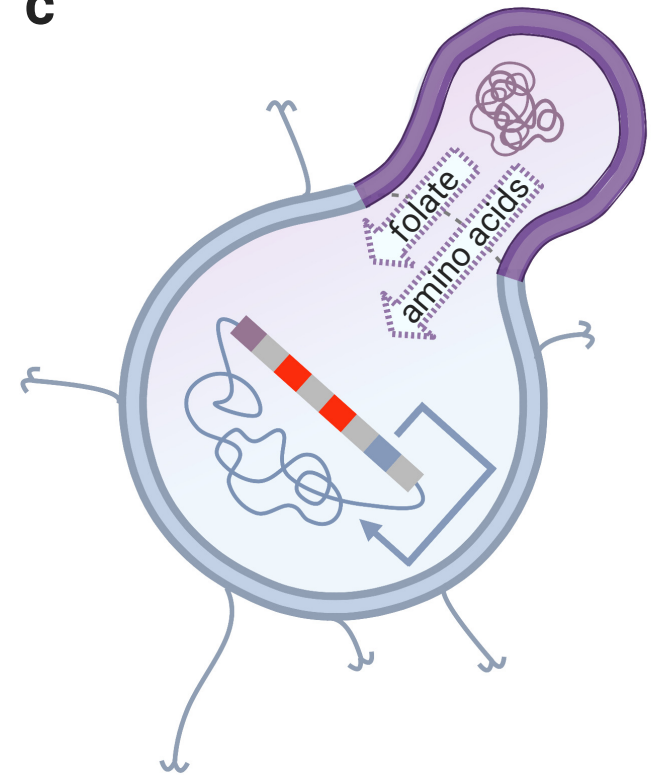
1141 ***Extended Data Fig. 4 | Spacer targeting analyses of publicly available archaeal genomes.***  
1142 *Directed spacer analysis of 7,012 publicly available archaeal genomes (Table S4) shows large*  
1143 *clusters of spacers targeting at species level. The targeting spacers (edges) of the genomes*  
1144 *Sulfolobus, Methanomicrobia and Halobacterium (nodes) form large clusters performing self-*  
1145 *targeting or targeting other genomes of the same family. The clustering was illustrated with*  
1146 *Cytoscape<sup>83</sup> (version 3.9.1). Please note that targeting within the same genus might limit the*  
1147 *interspecies recombination, as demonstrated in haloarchaea<sup>37</sup>, or reflect the presence of multiple*  
1148 *conserved genomic regions between the genomes.*

1149



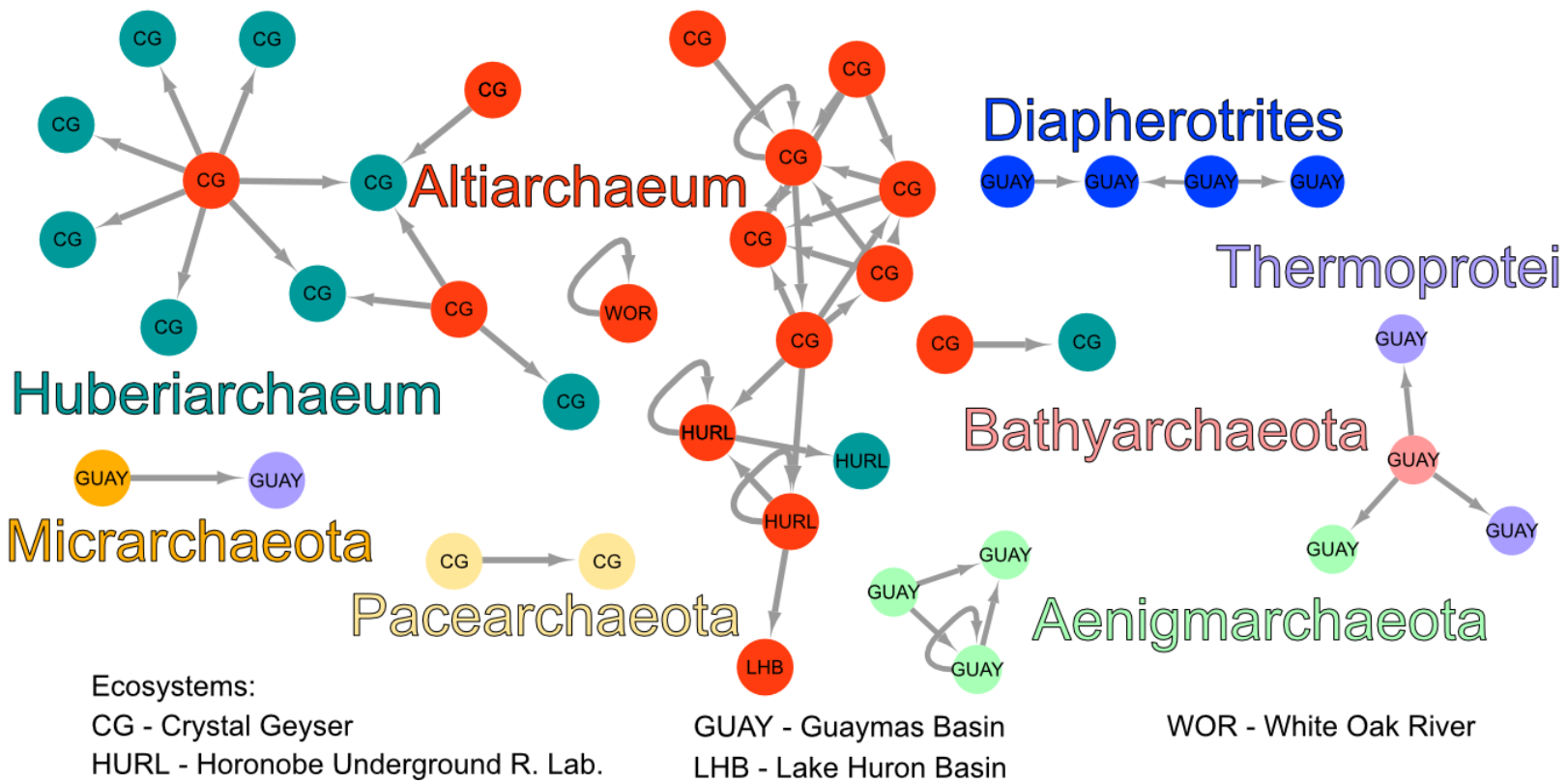


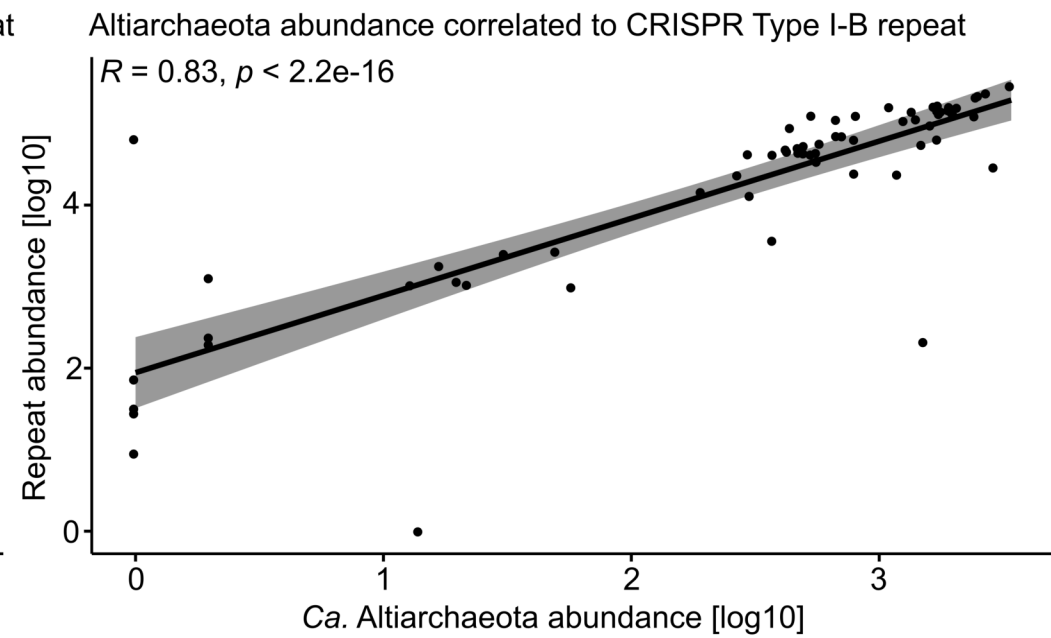
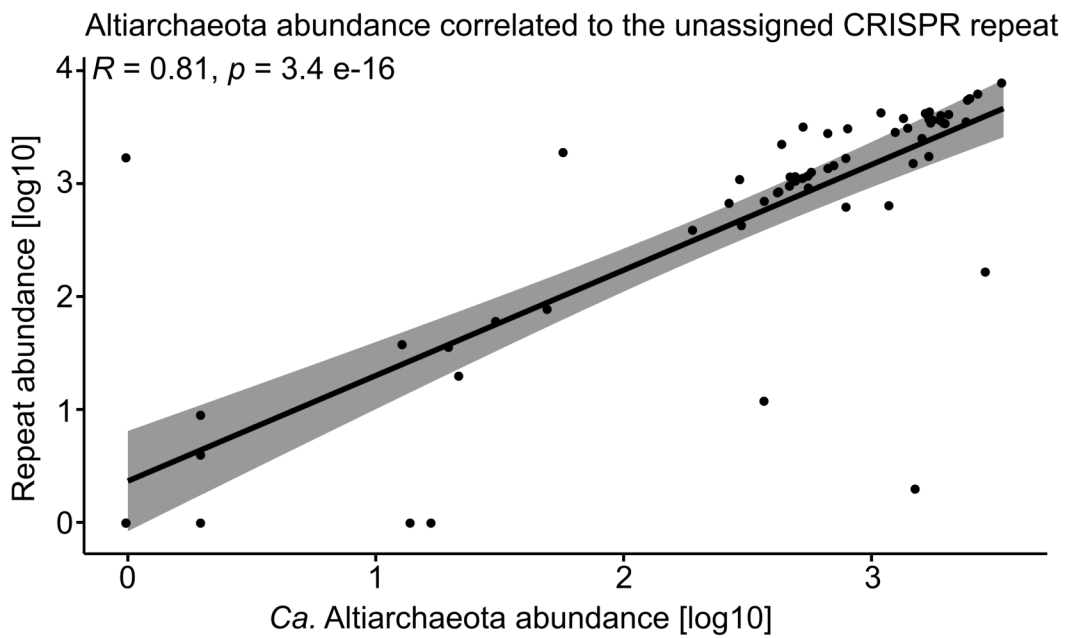


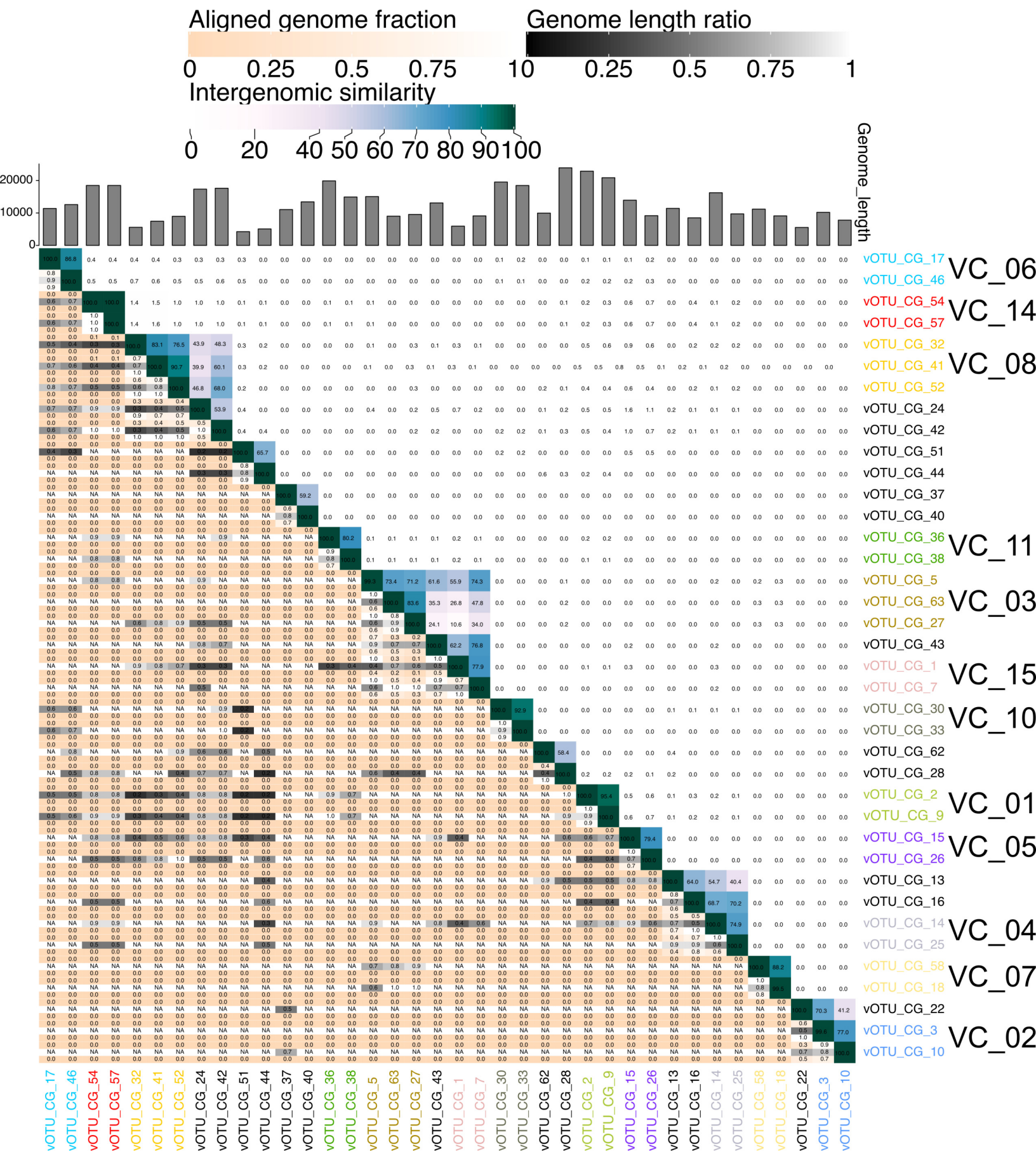
**a****b****c**

○ host   ○ episymbiont   ■ repeat   ■ spacer

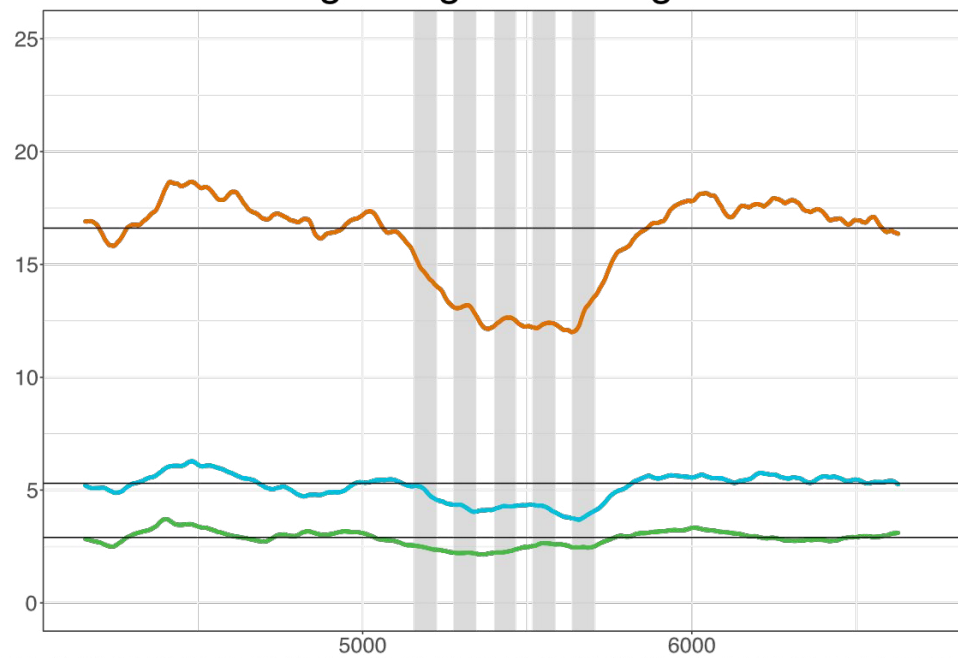
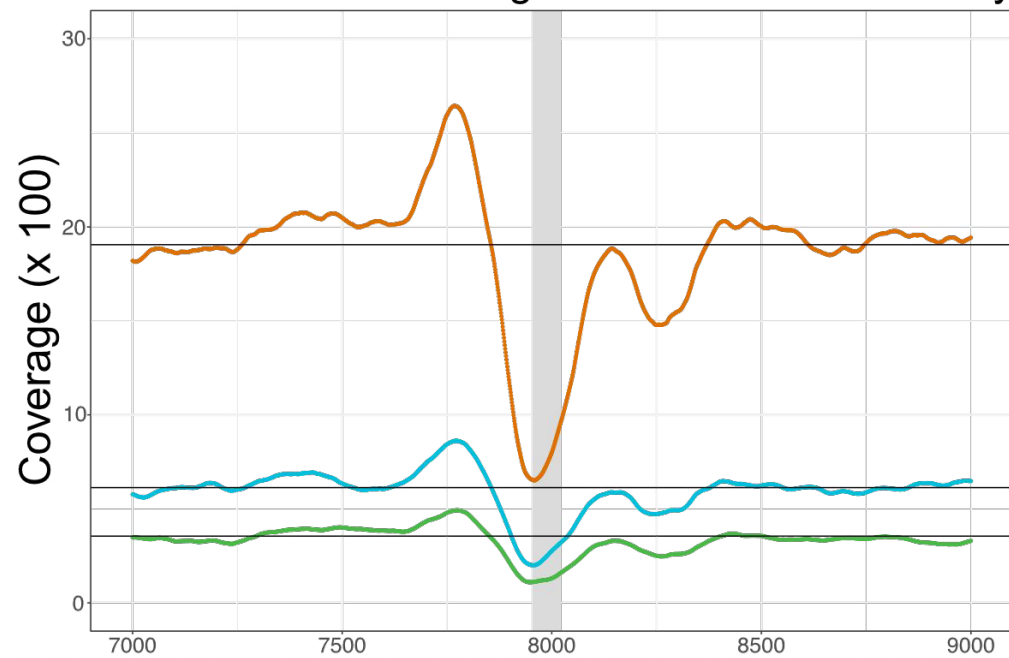
**a) Viral targeting. b) Targeting of episymbiont. c) Self-targeting and respective metabolic complementation.**



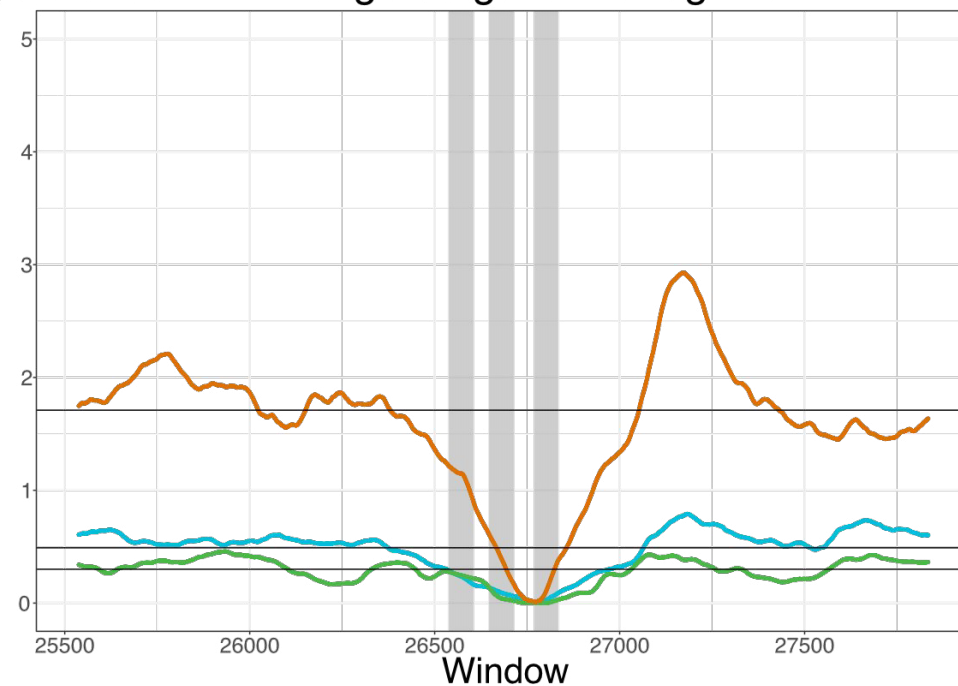
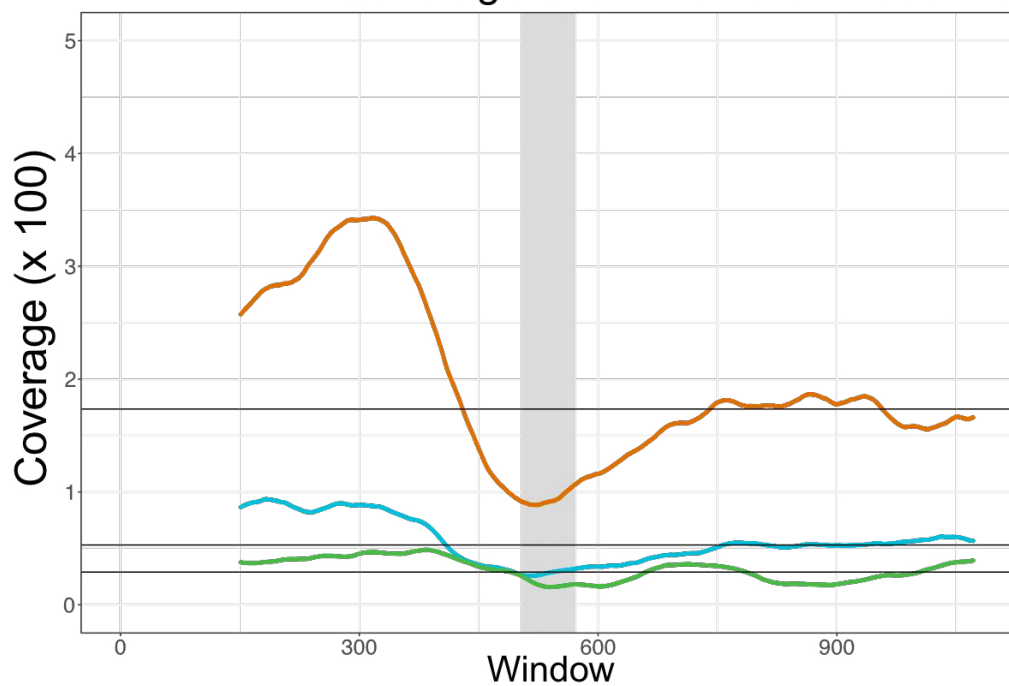




Coverage of *Ca. Altiarchaeum* crystalense within targeted genomic regions



Coverage of *Ca. Huberiarchaeum crystalense* within targeted genomic regions



CG05 CG08 CG16

