# Regulatory non-coding small RNAs are diverse and abundant in an extremophilic microbial community

Diego R. Gelsinger*, Gherman Uritskiy*, Rahul Reddy, Adam Munn, Katie Farney, and Jocelyne DiRuggiero[#]

*Department of Biology, The Johns Hopkins University, Baltimore, Maryland, USA*

*Equal contributions

Running Head: Metatranscriptomic identification of sRNAs in extremohile microbiome

Key words: non-coding, sRNA, metatranscriptomics, metagenomics, gene regulation, extremophile, microbiome, microbial communities

# Corresponding author

Jocelyne DiRuggiero

Johns Hopkins University

Department of Biology

3400 N. Charles Street, Mudd Hall 235

Baltimore MD 21218, USA

jdiruggiero@jhu.edu

1

1  **ABSTRACT**

2  Regulatory small RNAs (sRNAs) represent a major class of regulatory molecules that

3  play large-scale and essential roles in many cellular processes across all domains of

4  life. Microbial sRNAs have been primarily investigated in a few model organisms and

5  little is known about the dynamics of sRNA synthesis in natural environments, and the

6  roles of these short transcripts at the community level. Analyzing the metatranscriptome

7  of a model extremophilic community inhabiting halite nodules (salt rocks) from the

8  Atacama Desert with SnapT – a new sRNA annotation pipeline – we discovered

9  hundreds of intergenic (itsRNAs) and antisense (asRNAs) sRNAs. The halite sRNAs

10  were taxonomically diverse with the majority expressed by members of the

11  *Halobacteria*. We found asRNAs with expression levels negatively correlated with that

12  of their putative overlapping target, suggesting a potential gene regulatory mechanism.

13  A number of itsRNAs were conserved and significantly differentially expressed (FDR

14  <5%) between 2 sampling time points allowing for stable secondary structure modeling

15  and target prediction. This work demonstrates that metatranscriptomic field experiments

16  link environmental variation with changes in RNA pools and have the potential to

17  provide new insights into environmental sensing and responses in natural microbial

18  communities through non-coding RNA mediated gene regulation.

19      **INTRODUCTION**

20      Non-coding RNAs (ncRNAs) are untranslated short transcripts that are found in the

21      three domains of life where they play essential roles in many cellular processes

22      (Gelsinger and DiRuggiero 2018b, Cech and Steitz 2014). In prokaryotes, a subset of

23      these ncRNAs, thereby called small RNAs (sRNAs), are specifically involved in gene

24      regulation through RNA-RNA mediated interactions, modulating core metabolic

25      functions and stress related responses (Gottesman and Storz 2011). These sRNAs

26      range from 50 to 500 nucleotides in size and can be of two types: *trans*-encoded

27      sRNAs, also called intergenic sRNAs (itsRNAs), which bind their mRNA targets via

28      imperfect base-pairing and can target multiple genes, including key transcription factors

29      and regulators (Wagner and Romby 2015). itsRNAs can activate or inhibit translation

30      initiation by interacting with the ribosome binding site (RBS) and/or modulating mRNA

31      stability (Wagner and Romby 2015). In contrast, *cis*-encoded antisense RNAs (asRNAs)

32      are transcribed on the DNA strand opposite their target gene and thus can act via

33      extensive base pairing; they have been found to repress transposons and toxic protein

34      synthesis (Wagner and Romby 2015).

35      The functional roles of microbial sRNAs have been extensively studied in a few model

36      organisms and very little is known about the dynamics of sRNA synthesis in natural

37      environments and the roles of these short transcripts at the community level (Carrier,

38      Lalaouna, and Massé 2018, Gelsinger and DiRuggiero 2018b). To our knowledge, only

39      two studies have reported the discovery of sRNAs in natural microbial communities

40      (Shi, Tyson, and DeLong 2009, Bao et al. 2015). This paucity of knowledge suggests

41      that an abundance of sRNAs remain to be discovered, in particular in extreme

42      environments where they likely play essential roles in stress response (Clouet-d'Orval et

43      al. 2018), inter-species communication, and/or cross-species RNA interference

44      (Toyofuku, Nomura, and Eberl 2019, Cai et al. 2018, Tsatsaronis et al. 2018).

45      In hyper-arid deserts, microbial communities find refuge inside rocks as a survival

46      strategy against the extreme conditions of their environment (Pointing and Belnap

47      2012). Such community inhabits halite (salt) nodules in Salars of the Atacama Desert,

48      Chile, which is one of the oldest and driest deserts on Earth (Crits-Christoph et al. 2016,

49      Finstad et al. 2017). The halite endolithic (within rock) community harbors mostly

50  members of the Archaea (*Halobacteria*), unique *Cyanobacteria*, diverse heterotrophic

51  bacteria, and a novel type of algae (Crits-Christoph et al. 2016, Finstad et al. 2017). The

52  main source of liquid water for this community is from salt deliquescence (Davila et al.

53  2008) and it is sustained by $CO_2$ fixed via photosynthesis (Crits-Christoph et al. 2016,

54  Davila et al. 2015). While previous studies have demonstrated the role of sRNAs in the

55  stress response of one of the members of this community, the halophilic archaeon

56  *Haloferax volcanii* (Gelsinger and DiRuggiero 2018a, Kliemt, Jaschinski, and Soppa

57  2019), there is no information on any of the other members.

58  Here we used a combination of genome-resolved metagenomics and

59  metatranscriptomics to investigate the role of sRNAs in the adaptive response of

60  microorganisms inhabiting halite nodules. We developed an analytical pipeline, SnapT,

61  built on our previous work on sRNAs with model organisms (Gelsinger and DiRuggiero

62  2018a), to enable the discovery of sRNAs at the community level. We found hundreds

63  of sRNAs (both itsRNAs and asRNAs) in the halite community, including conserved

64  sRNAs, validating our experimental approach. A number of itsRNAs were significantly

65  differentially regulated between 2 sampling time points and, for a subset of these, we

66  were able to perform structure and target prediction, deciphering their potential

67  regulatory roles. Coupling metagenomics and metatranscriptomics with SnapT allows

68  for the potential to uncover the complex regulatory networks that govern the state of a

69  microbial community.

70

71  **MATERIAL AND METHODS**

72  *Sample and weather data collection and nucleic acid extraction*

73  Halite nodules were harvested in Salar Grande, an ancient evaporated lake in the

74  Northern part of the Atacama Desert (Robinson et al. 2015) in February 2016 and 2017,

75  3 and 15 months after a major rain event (Uritskiy et al. 2019). All nodules were

76  harvested within a 50m$^2$ area as previously described (Robinson et al. 2015). The

77  colonization zone of each nodule was grounded into a powder, pooling from 1-3 nodules

78  until sufficient material was collected, and stored in the dark in dry conditions until DNA

79  extraction in the lab. Samples used for RNA were stored in *RNAlater* at 4°C until RNA

80  extraction in the lab. Genomic DNA was extracted as previously described (Robinson et

81   al. 2015, Crits-Christoph et al. 2016) with the DNAeasy PowerSoil DNA extraction kit

82   (QIAGEN). Total RNA was extracted from the fixed samples by first isolating the cells

83   through gradual dissolving of the salt particles as previously described (Robinson et al.

84   2015, Crits-Christoph et al. 2016) and lysing them through mechanical bead beating

85   with the RNAeasy PowerSoil RNA extraction kit (QIAGEN). Total RNA was then

86   extracted from the lysate with a Quick-RNA miniprep kit (Zymo Research). RT-PCR was

87   used to validate the absence of contaminating DNA in the total RNA used for RNA-seq

88   libraries (Fig. S11).

89

90   *Library preparation*

91   Whole genome sequencing libraries were prepared using the KAPA HyperPlus kit

92   (Roche) as previously described (Uritskiy et al. 2019) and sequenced with paired 150bp

93   reads on the HiSeq 2000 platform at the Johns Hopkins Genetic Resources Core

94   Facility (GRCF). Total RNA-seq libraries were prepared with the SMARTer Stranded

95   RNA-seq kit (Takara and Bell), using 25ng of RNA input and 12 cycles for library

96   amplification. We sequenced 22 libraries from replicate samples from 2016 and 24

97   libraries from replicate samples from 2017.

98

99   *WMG sequence processing*

100  The de-multiplexed WMG sequencing reads were processed with the complete

101  metaWRAP v0.8.2 pipeline (Uritskiy, DiRuggiero, and Taylor 2018) with recommended

102  databases on a UNIX cluster with 48 cores and 1024GB of RAM available. Detailed

103  scripts   for   the   entire   analysis   pipeline   can   be   found   at

104  https://github.com/ursky/timeline_paper.

105

106  *SnapT for sRNA community identification*

107  An   analytic   pipeline,   SnapT   for   Small   ncRNA   Annotation   Pipeline   for

108  (meta)Transcriptomic data, was adapted and developed from our previous work

109  (Gelsinger and DiRuggiero 2018a) to find, annotate, and quantify intergenic and

110  antisense sRNA transcripts from transcriptomic or metatranscriptomic data. Detailed

111  scripts for the pipeline can be found at https://github.com/ursky/SnapT and search

4

112 criteria were as follows: intergenic transcripts were at least 30 nt away from any gene or
113 ORF on both strands; antisense transcripts were 30 nt away from any gene on their
114 strand, but overlapped with a gene on the opposite strand by at least 10 nt; small
115 peptides (<100 nt) were not counted as genes if they were encoded in a transcript that
116 was more than 3 times their length; non-coding transcripts could not contain any
117 reading frame greater than 1/3 of their lengths; predicted non-coding transcripts near
118 contig edges were discarded and the minimum distance to the edge of a contig was
119 dynamically computed such that the tips of contigs were not statistically enriched in
120 annotated ncRNAs; small ncRNAs were between 50 nt and 500 nt in length; sRNA
121 transcripts could not have significant homology with any protein in the NCBI_nr
122 database (query cover>30%, Bitscore>50, evalue<0.0001, and identity>30%) and with
123 any tRNA, RNase P, or signal recognition particle (SRP) model in the Rfam non-coding
124 RNA database.

125

126 *Taxonomic assignment and distribution of sRNAs*
127 The taxonomic origin of each annotated sRNA was taken to be as that of the contig on
128 which it lay. The taxonomy of each contig was estimated by taking the weighted
129 average of the taxonomic assignment of the genes encoded on it, as determined
130 through the JGI IMG functional and taxonomic annotation service.

131

132 *Metatranscriptomic Correlation and Differential Expression Analysis*
133 We used a read count-based differential expression analysis to identify differentially
134 expressed sRNA and mRNA transcripts. The program featureCounts (Liao, Smyth, and
135 Shi 2014) was used to rapidly count reads that map to the assembled RNA transcripts
136 (described above) as previously described (Gelsinger and DiRuggiero 2018a). In order
137 to account for organism abundance changes (as opposed to true transcript changes),
138 we normalized the transcript read counts to the total read counts from the contig on
139 which the transcript lies on. The read counts were then used in the R differential
140 expression software package DESeq2 (Love, Huber, and Anders 2014) to calculate
141 differential expression by determining the difference in read counts between 2016
142 normalized read counts from 2017 normalized read counts. The differentially expressed

143    RNAs were filtered based on the statistical parameter of False Discovery Rate (FDR)

144    and those that were equal to or under a FDR of 5% were classified as true differentially

145    expressed transcripts. We carried out differential expression analysis using a pairwise

146    Wald test to find any possible differences between years (Love, Huber, and Anders

147    2014). In parallel, normalized expression values were calculated using stringtie in

148    transcripts per million (TPM). TPM of transcripts were normalized in the same way as

149    read counts, except using contig TPM. TPM of transcripts was used for ranking of

150    expression within samples as opposed to differential expression analysis.

151

152    *Regulatory element motif identification of sRNAs, structure and target prediction*

153    50 nucleotides upstream from the sRNA transcript start coordinates were searched for

154    transcription motifs (BRE and TATA-box for archaea and -35 and -10 consensus

155    sequences for bacteria) using both multiple sequence alignments and visualization with

156    WebLogo and motif searching with MEME (Gelsinger and DiRuggiero 2018a).

157    Conserved sRNAs were identified using blastn against the NCBI nt database.

158    Secondary structures of conserved sRNAs were predicted using sRNAs that had an e-

159    value maximum of 1E-3, a sequence similarity of 70% or more, and 50% or more

160    coverage with a NCBI nt database blastn hit; a minimum of 14 alignments were used in

161    the program LocARNA using global alignment settings (Will et al. 2012). Lastly, putative

162    targets were predicted for itsRNAs by searching for optimal sRNA-mRNA hybridization

163    using the IntaRNA program with the no seed parameter (Mann, Wright, and Backofen

164    2017) and the reference genes for each respective MAG. Targets were ranked by

165    lowest p-value. Expression levels for putative targets of antisense sRNAs were obtained

166    from co-expression analysis of transcripts (Gelsinger and DiRuggiero 2018a). The

167    sRNA and putative target mRNA TPM expression values were tracked across the

168    replicates, and the Pearson correlation was computed.

169

170    *Enrichment cultures*

171    Three types of culture medium were inoculated in triplicate with ~2 g of grounded halite

172    colonization zones and incubated at 42°C with shaking at 220 rpm (Amerex Gyromax

173    737) for 1 to 2 weeks. Cells were harvested by centrifugation and nucleic acids

6

174 extracted as described above. Media were: GN101 medium (Kish et al. 2009)
175 containing 250 g of salt per L and 10 g of peptone as carbon source; Hv-YPC medium
176 (Dyall-Smith 2009) containing 250 g of salt per L and 8.5 g of yeast extract, 1.7 g of
177 peptone, and 1.7 of casamino acids as carbon sources; and IO containing 250 g of salt
178 and the same carbon sources as the Hv-YPC medium. The taxonomic distribution of the
179 cultures was obtained with 16S rRNA gene sequencing as previously described
180 (Uritskiy et al. 2019).

181

182 <u>sRNA validation</u>
183 Total RNA extracted from environmental samples and enrichment cultures was
184 converted into cDNA using the SuperScript III First-Strand Synthesis System
185 (ThermoFisher). The cDNA was then amplified using primers designed for sRNAs
186 identified in the halite community **(Table S1)**, as previously described (Meslier et al.
187 2018). Amplicons were sequenced using Sanger sequencing (GENEWIZ, South
188 Plainfield, NJ).

189

190 <u>Data availability</u>
191 Raw sequencing data are available from the National Centre for Biotechnology
192 Information under NCBI project ID PRJNA484015. The metagenome co-assembly and
193 functional annotation are available from the JGI Genome Portal under IMG taxon OID
194 3300027982. Metatranscriptome data GEO # in process. Scripts for functional
195 annotation, statistical analyses, differential expression, and figures are available at
196 https://github.com/ursky/srna_metatranscriptome_paper.

197

198 **RESULTS**
199 ***Landscape of predicted sRNAs in the halite community and validation***
200 We discovered hundreds of ncRNAs in an extremophilic community inhabiting halite
201 nodules (salt rocks) in the Atacama Desert by using SnapT
202 (https://github.com/ursky/SnapT), a pipeline adapted from our previous work on a model
203 haloarchaeon present in the halite community (**Table 1; data S1**) (Gelsinger and
204 DiRuggiero 2018a). We used metatranscriptomics data from multiple replicate samples

7

205    collected in the field in 2016 and 2017  (21 and 24 replicates for 2016 and 2017,

206    respectively; **Fig. S1**). Using SnapT, we aligned reads from stranded RNA-seq libraries

207    to our reference co-assembled metagenome from a previous study (Uritskiy et al. 2019)

208    and assembled the reads into transcripts (**Fig S2**). The transcripts were then intersected

209    with the metagenome annotation as well as open reading frames to select for either

210    novel transcripts on the opposite strand of coding transcripts (asRNAs) or for novel

211    transcripts that fell into intergenic regions (itsRNAs). Putative ncRNA transcripts were

212    then further enriched (**Fig. S2**) using a threshold at 5x and 10x assembly coverage in

213    order to identify intergenic and antisense ncRNAs, respectively. (**Fig. S3; Table 1).** The

214    size of these ncRNAs was then filtered from 50 to 500 nucleotides to produce a final set

215    of non-coding sRNAs. The size distribution of these sRNAs was primarily between 50

216    and 200 nt for itsRNAs and above 200 nt for asRNAs. (**Fig. S4**).

217    The halite ncRNAs were taxonomically assigned to diverse members of the community;

218    their distribution between Archaea (54%) and Bacteria (46%) (**Table 1**) was similar to

219    that of the total metatranscriptomic reads for the community (**Fig. 1B and C**). In

220    contrast, the taxonomic profile of the metagenome showed a larger contribution of

221    bacterial reads and in particular of reads assigned to *Cyanobacteria* and *Bacteroidetes*

222    (**Fig. 1A**). Because of the use of strand specific RNA-seq libraries, we could confidently

223    identify both intergenic (it)sRNA, located between coding regions, and antisense

224    (a)sRNA, overlapping with their putative target (**Table 1**). We found 3 times more

225    itsRNAs in the Archaea than in the Bacteria, whereas asRNAs were more abundant in

226    the Bacteria and more often associated with members of the *Cyanobacteria* (38%) and

227    *Bacteriodetes* (15%) (**Table 1**; **Fig. 1D and E**). We also found 79 ncRNAs, that belong

228    to 6 known families of RNAs present in the Rfam database (**Fig. S5**; **data S2**) (Kalvari

229    et al. 2017), validating our experimental and computational approach. This database is

230    a collection of RNA families, each represented by multiple sequence alignments,

231    consensus secondary structures, and covariance models. Of the Rfam-conserved

232    ncRNAs, 70% were assigned to archaea and included RNaseP RNAs, signal

233    recognition particle RNAs (SRP RNAs), and tRNAs. Of the Rfam-conserved bacterial

234    ncRNAs, most were from SRP RNAs and tRNA conserved families. In addition, a

235    cobalamin riboswitch and the regulatory sRNA, CyVA-1, were detected in low

236   abundance in the halite *Cyanobacteria*. We also found 3 ncRNAs (4%) from Eukarya, a

237   tRNA, a U4 spliceosomal RNA, and a RNase for mitochondrial RNA processing (MRP).

238   Using blastn analysis (max e-value of 1E-3, sequence similarity of 70% or more,

239   coverage of 50% or more), we discovered another 155 ncRNAs that were conserved in

240   the NCBI nt datasbase, with 60% from archaea and 40% from bacteria (**Table 1**). The

241   majority were asRNAs (109), with only 44 itsRNAs. The conserved asRNAs most highly

242   expressed (standardized tpm> 100) were all SPR RNAs in haloarchaea that were not

243   found in the Rfam database. Of the conserved itsRNAs, we identified 3 tRNAs, 13 SRP

244   RNAs, and 22 ncRNAs that were found in the genome of multiple species, all

245   *Halobacteria*, but with no function assigned. The most highly expressed and conserved

246   itsRNAs (standardized tpm> 100; 13 ncRNAs) were SRP RNAs not included in the

247   Rfam database.

248   Another validation of our findings was the presence of canonical promoter elements

249   upstream of archaeal itsRNAs, suggesting that they were indeed *bona fide* transcripts

250   that could recruit basal transcription factors (**Fig. S6).** We did not find significant

251   promoter elements upstream of the bacterial itsRNAs, which might reflect the diversity

252   of promoter elements across the various bacterial taxa we identified in the halite

253   community. In contrast, no promoter elements were identified in the upstream regions of

254   asRNA from both domains of life.

255   When looking at the expression levels of all itsRNAs normalized to contig abundances,

256   we found that they were similar for both the 2016 and 2017 samples and slightly higher

257   than that of the asRNAs, whereas the expression profile of the asRNAs was more

258   variable across samples for both years (**Fig. S7**). Remarkably, the expression levels of

259   itsRNAs and asRNAs for both years was 2-fold higher than that of protein encoding

260   genes. Whereas there is an inherent bias in our approach to identify sRNAs at the

261   community level (coverage threshold in SnapT) compared to protein encoding genes,

262   this finding strongly indicates potential functional relevance for a number of these

263   sRNAs.

264   We experimentally validated a number of sRNAs using RT-PCR with environmental and

265   enrichment cultures (**Table S1**). Enrichments were performed with several media

9

266   containing high (25%) and relatively low (18%) salt, and various combinations of carbon

267   sources. Amplicon sequencing of the enrichments revealed that high salt and diverse

268   carbon sources resulted in higher diversity of taxa, although haloarchaea dominated in

269   all enrichments (**Fig. S8**). All validated sRNAs belong to haloarchaea with the exception

270   of one from *Cyanobacteria*. Sequences of the PCR products confirmed that they were

271   sRNAs and validated our computational approach.

272   ***Relationship with target genes and putative function of community asRNAs***

273   Using our strand-specific RNA-seq data, we were able to identify the overlap position of

274   asRNAs to their antisense transcripts. We found that, in both Archaea and Bacteria, the

275   majority of asRNAs start within the span of their cognate gene and end near the 5' end

276   of its mRNA. In both domains there is also an enrichment for asRNA-mRNA overlaps

277   near the 5' end of the mRNA. A similar trend has previously been reported in two

278   species of archaea (Gelsinger and DiRuggiero 2018a, de Almeida et al. 2019).

279   We compared the expression level of asRNAs with that of their putative target genes

280   and found that highly expressed asRNAs were associated with lowly express genes

281   (**Fig. 2A**). Of gene pairs with asRNA expression >100 tpm and gene expression <0.1

282   tpm, most where from haloarchaea (77%), with 12% of *Cyanobacteria*, and 11% of

283   other bacteria (*Bacteriodetes* and A*cinetobacter*) (**data S3**). Gene functions were

284   enriched for transport (16%) and cell membrane/wall metabolism (5%), while most were

285   hypothetical proteins (44%). Of the genes potentially negatively regulated by their

286   cognate asRNAs, we found an archaeal regulator of the IclR family and potassium

287   uptake protein TrkA. Only 2 asRNAs with high expression levels (>100 standardized

288   tpm) were associated with genes with relatively high expression levels (>1 standardized

289   tpm), while still being negatively correlated (Fig. 2A). The corresponding genes encoded

290   for an iron complex outermembrane receptor protein from *Salinibacter* and a ABC-type

291   sodium efflux pump permease subunit from a *Halobacteria*. When applying a stringent

292   cut-off, we found 9 statistically significant and negatively correlated asRNA:gene pairs

293   **(Fig. 2B).** Four were from *Bacteroidetes*, 4 from *Halobacteria*, and 1 from an

294   unidentified bacterium. At the functional level, transport systems, and in particular iron

295   transport systems, were particularly enriched (**data S3**). In contrast, we did not find any

296   significant positive regulation between asRNAs and their cognate genes. When

297   adjusted for the carrying organism's abundance, expressed as the average RNA read

298   coverage of the contigs, we found that overall itsRNAs were more highly expressed

299   than asRNAs (**Fig. 2C and 2D**). Highly expressed sRNAs, for both types, were mostly

300   carried by haloarchaea.

301   ***Differential expression of itsRNAs at the community level and target prediction***

302   Analysis of itsRNAs expression levels showed a clear separation between the 2016 and

303   2017 samples (**Fig. 3a**). We carried out a differential expression analysis and found that

304   109 (18%) of the regulatory itsRNAs were significantly differentially expressed (FDR

305   <5%) between samples collected in 2016 and 2017 (**Fig. 3 and data S4**), 3 and 15

306   months after a major rain event in the desert (Uritskiy et al. 2019). Of these, 72% were

307   annotated as archaea and 28 % as bacteria and 16 were conserved in multiple

308   genomes (14 from *Halobacteria* and 2 from *Cyanobacteria*). Conservation of

309   differentially expressed itsRNAs allowed for structure modeling and, when high quality

310   MAGs (>70% completion and <5% contamination) were available from the

311   metagenome, target prediction (**Fig. 4 and Fig. S10**). A number of non-differentially

312   expressed itsRNAs were also conserved, providing additional opportunity for structure

313   prediction; these included itsRNAs from *Halococcus* (STRG. 48671.1; 69 nt), *Halobellus*

314   *limi* (STRG136887.1; 209 nt), and from a member of the *Nanohaloarchaea*

315   *(*STRG.4577.1; 266 nt) (**Fig. S10A).**

316   All predicted structures displayed loop and stems regions that had high sequence

317   conservation (light purple regions on sequence–structure-based alignment reliability

318   [STAR] profile plots) and high structure conservation (dark purple), and line plots

319   representing the reliability of the predictions as calculated by LocaRNA (**Fig. 4 and Fig.**

320   **S10B**). Density plots combined with dumbbell plots were used for visualizing predicted

321   interactions between itsRNAs and their putative targets, using IntaRNA data from the

322   top 100 most reliable interaction predictions with the lowest free energy of hybridization

323   (Mann, Wright, and Backofen 2017) **(Fig. 4)**. High confidence assignments were

324   obtained for 4 differentially expressed itsRNAs from *Cyanobacteria*, *Halapricum salinun*,

325   and a member of the *Halobacteria* **(data S5**) More than one interaction peak were

326   derived from density plots; peak 1 (green) corresponded to the highest interaction

327  density, which mapped to loop regions in the itsRNA secondary structure with high

328  sequence and structure conservation, respectively, and was thus a confident

329  assignment as an interaction region, whereas Peak 2 (yellow) was a less confident

330  assignment structurally despite high interaction density (**Fig. 4 and Fig. S10B**).

331  Using this information, we identified the most probable targets for *Cyanobacteria*

332  STRG.5354.4 candidate itsRNA (229 nt). This itsRNA was conserved as a 6S

333  regulatory RNA in the rfam database, which in bacteria is found to inhibit transcription

334  by binding directly to the housekeeping holoenzyme form of RNA polymerase

335  (Wassarman 2018). Of the top 50 most probable targets for STRG.5354.4, which were

336  those with the lowest free energy of hybridization between itsRNA and targets, were

337  cation:H+ antiporters [shown to be involved in osmoregulation (Krulwich, Hicks, and Ito

338  2009)], a PleD family two-component response regulator, the photosystem I PsaB

339  protein, chemotaxis transducers, and proteins involved in energy metabolism. Most

340  probable targets for differentially expressed itsRNA, STRG. 86294.1 (281 nt) from

341  *Halapricum salinum* included various transporters and putative membrane and cell wall

342  associated proteins; notably an ammonium transporter (Amt family), an alkanesulfonate

343  monooxygenase SsuD from a gene cluster expressed under sulfate or cysteine

344  starvation (Eichhorn, van der Ploeg, and Leisinger 1999)**,** and several proteins involved

345  in cofactors and vitamin metabolism. Predicted targets with the lowest free energy of

346  hybridization for STRG.49508.3 candidate itsRNA (99 nt) from *Halobacteria* were

347  elongation factor 1-alpha, which promotes the GTP-dependent binding of aminoacyl-

348  tRNA to the A-site of ribosomes during protein biosynthesis, several ribosomal proteins,

349  and a number of hypothetical proteins. Target prediction for *Cyanobacteria*

350  STRG.5356.1 candidate itsRNA (242 nt) included molecular chaperones (DnaK and

351  DnaJ classes), a cell division protease FtsH, and a number of uncharacterized proteins.

352

353  **DISCUSSION**

354  The roles of regulatory sRNAs have been extensively studied in bacterial, and to a

355  lesser extent, in archaeal model systems (Carrier, Lalaouna, and Massé 2018,

356  Gelsinger and DiRuggiero 2018b) but, to date, only two studies have reported the

357  discovery of sRNAs in microbial communities. In one study, Shi *et al.* (Shi, Tyson, and

358    DeLong 2009) used metatranscriptomic data to identify unique microbial sRNAs in the

359    ocean's water column while the study by Bao et al. (Bao et al. 2015) revealed extensive

360    antisense transcription in the human gut microbiota, also using metatranscriptomic

361    datasets. Efforts have also been made to mine publically available databases for sRNA

362    discovery (Weinberg et al. 2017) but this was still addressing the role of sRNAs in single

363    microorganisms.

364    One major difficulty in obtaining metatranscriptomic data from natural microbial

365    communities, in particular from extreme environments, is the low amount of biomass

366    that can be collected, resulting in low RNA yields (Uritskiy and DiRuggiero 2019). This,

367    in turn, prevents attempts at ribo-depletion, resulting in a decreased number of non-

368    ribosomal RNA reads available for analysis. Nevertheless, using SnapT, a flexible

369    pipeline to process metagenomics and metatranscriptomic data, we report the discovery

370    of hundreds of diverse sRNAs from an extremophilic community inhabiting halite

371    nodules in the Atacama Desert. In the process, we applied extensive quality control with

372    coverage thresholding, correction for contig edge mis-annotation, and the removal of

373    potential non-ncRNAs through sequence and homology searches. While this approach

374    might potentially result in false negatives, and may bias our findings toward the most

375    highly expressed sRNAs in the community, it also insured the robustness of our sRNA

376    predictions by minimizing the number false positives. The identification of ncRNAs in the

377    halite community that belong to the Rfam database **(Kalvari et al. 2017)**, together with

378    experimental validation of a number of sRNAs with environmental and enrichment

379    cultures, substantiated our analytical approach. Additionally, expression levels of

380    sRNAs 2-fold higher than that of protein encoding genes, strongly indicates potential

381    functional relevance for a number of these sRNAs.

382    The taxonomic composition of the halite sRNAs matched that of the community's

383    metatranscriptomic profile, reflecting the contribution of the most active members,

384    including *Cyanobacteria*, *Bacteriodetes*, and a number of *Halobacteria*. We found

385    significantly more itsRNAs in the archaea than in the bacteria and the trend was reverse

386    for the asRNAs. This novel finding is representative of published work in model

387    organisms where a wide range of sRNAs has been found so far in prokaryotes, from

13

388    less than a dozen to more than a thousand per genome (Carrier, Lalaouna, and Massé

389    2018, Gelsinger and DiRuggiero 2018b).

390    Antisense sRNAs overlap their putative targets providing insights into their functional

391    role (Wagner and Romby 2015). In the halite community, we found that asRNAs

392    expression levels were negatively correlated with that of their putative targets, with

393    highly expressed asRNAs overlapping lowly expressed protein encoding genes. A

394    similar trends was reported in the haloarchaeon *Haloferax volcanii*, when investigating

395    oxidative stress responsive sRNAs, and most of the putative targets were transposase

396    genes (Gelsinger and DiRuggiero 2018a). Putative target gene functions in our study

397    were mostly from haloarchaea and enriched for transport systems, cell membrane and

398    cell wall metabolism, with a large number of hypotheticals. Of particular interest, was an

399    archaeal IclR transcription regulator; these regulators are known to be involved in

400    diverse physiological functions, including multidrug resistance, degradation of

401    aromatics, and secondary metabolites production (Molina-Henares et al. 2006) and are

402    distributed in a wide range of prokaryotes, including Archaea (Perez-Rueda et al. 2018).

403    Also of interest, was a Trk potassium uptake system, also found in both bacteria and

404    archaea and essential for the maintenance of high intracellular potassium in salt-in

405    strategists (Oren 2013). In contrast, we did not find any significant positive regulation

406    between asRNAs and their cognate genes, which might be due to the inherent quality of

407    our data set, i.e. no ribo-depletion and heterogeneity across replicates (Uritskiy and

408    DiRuggiero 2019). Alternatively, it might also reflect promiscuous transcription

409    processes as argued when considering the functionality of asRNAs (Lloréns-Rico et al.

410    2016). Other arguments in favor of spurious transcription was the size distribution for

411    asRNAs found in the halite community, which was significantly larger than that of

412    itsRNAs, low expression level when adjusted for organism abundance when compared

413    to itsRNAs, and the absence of canonical regulatory elements in the upstream regions

414    of asRNAs. However, we found also putative target functions that reflected the

415    environmental challenges faced by members of this extremophile community, such as

416    osmoregulation and nutrient uptake, indicating that these asRNAs might indeed regulate

417    fundamental biological functions at the community level.

418    We previously showed that the halite community dramatically shifted its taxonomic and

419    functional composition after a major rain event in 2015, and while it recovered at the

420    functional level in 2017, 15 months after the rain, members of the communities were

421    permanently replaced (Uritskiy et al. 2019). Here we found that 18% of the halite

422    community itsRNAs were significantly differentially expressed (FDR <5%) between

423    samples collected in 2016 and 2017 (3 and 15 months after the rain, respectively),

424    potentially indicating a transcriptional response to changes in environmental conditions.

425    Intergenic sRNAs are of particular interest because they can target multiple genes,

426    including key transcription factors and regulators (Gottesman and Storz 2011). As a

427    consequence, a single sRNA can modulate the expression of large regulons and thus

428    have a significant effect on metabolic processes (Carrier, Lalaouna, and Massé 2018).

429    However, they do not overlap their target genes or bind their targets mRNAs with

430    perfect complementary, which make finding targets for these sRNAs very challenging

431    without genetic tools (Gelsinger and DiRuggiero 2018b).

432    To solve this problem at the community level, we focused on itsRNAs that were

433    conserved and for which we could perform structural prediction. The intersection of this

434    small subset of sRNAs with high quality MAGs that could be used as reference

435    genomes, yielded confident target predictions for 4 differentially expressed itsRNAs,

436    giving insights into metabolic functions potentially regulated by sRNAs at the community

437    level. These included transporters, particularly related to osmotic stress, nutrient uptake

438    and starvation, and pathways for chemotaxis and energy production and conversion.

439    These pathways reflect the environmental challenges members of the halite

440    communities are subjected to, including osmotic adjustments to climate perturbation

441    (Uritskiy et al. 2019) and competition for nutrients in a near-close system with primary

442    production as the major source of organic carbon (Crits-Christoph et al. 2016). Using

443    the genomic context of sRNAs from the ocean's water column microbial communities,

444    Shi *et al*. (Shi, Tyson, and DeLong 2009) reported similar metabolic functions,

445    underlying specific regulatory needs for natural communities. In contrast, genes with

446    antisense transcription to asRNAs identified in the human gut microbiome were mostly

447    transposase genes with a small component of bacterial house-keeping genes (Bao et

15

448     al. 2015). It important to note that no computational target prediction, using sRNA
449     conserved predicted structure, was reported in either study.

450     Regulation of transcription by 6S sRNA has been shown to increase competitiveness
451     and long-term survival in bacteria (Wassarman 2018), suggesting an important role for
452     *Cyanobacteria* candidate sRNA STRG.5354.4, identified as a 6S sRNA. Because of
453     high RNA-seq coverage of the *Cyanobacteria* MAGs, we could show that 40% of the top
454     50 targets for sRNA STRG.5354.4 were differentially regulated and more highly
455     expressed in 2016, suggesting positive regulation by this sRNAs onto its putative
456     targets. Transcriptional factors and regulators were also found as putative targets of
457     differentially regulated itsRNAs in the halite community, underlying the capacity of
458     microbial sRNAs to modulate the expression of large regulons (Gottesman and Storz
459     2011, Gelsinger and DiRuggiero 2018b, Nitzan, Rehani, and Margalit 2017). Finally, a
460     candidate itsRNAs from the *Halobacteria* had a number of predicted targets associated
461     with ribosomal proteins and proteins involved in translation processes. This finding,
462     together with a recent study in *H. volcanii* **(Wyss et al. 2018)**, support the idea of sRNA
463     modulation of protein biosynthesis in the Archaea. A potential framework for
464     mechanisms for sRNA regulation of translation might be provided by a report, in the
465     haloarchaeon *Halobacterium salinarum*, of modular translation subsystems that might
466     selectively translate a subset of the transcriptome under specific growth conditions
467     **(Raman et al. 2018)**.

468     ***Conclusion***
469     In this study, we characterized the taxonomic and functional landscape of sRNAs
470     across two domains of life in an extremophilic microbial community, demonstrating that
471     asRNAs and itsRNAs can be reliably identified from natural environmental communities.
472     To     facilitate     this     work,     we     built     a     flexible     pipeline,     SnapT
473     (https://github.com/ursky/SnapT), leveraged by our expertise of sRNA biology in a
474     model halophilic archaeon, and which is available to use with metatranscriptomic data
475     from any community. We demonstrated that we could perform target prediction and
476     correlate expression levels between itsRNAs and predicted target mRNAs, paving the
477     way for novel discoveries that have never been done at the community level. While
478     additional work with enrichment cultures remain to be done to fully characterize the

479  functional roles of sRNAs from the halite community, and their mechanism of action,

480  these differentially expressed sRNAs for which we found putative targets show the

481  power of community-level, culture-independent approach analysis for gene regulation

482  processes.

483

484  **Author Contributions**: JDR design the study, collected field samples, and wrote the

485  manuscript. DRG and GU collected field samples, conducted experiments, analyzed the

486  data, and edited the manuscript. RR and AM performed experimental target validation.

487  KF grew and characterized enrichment cultures. All authors approved the final version

488  for submission.

489

490  **Acknowledgments**: This work was supported by NASA grant 18-EXO18-0091.

491

492  **Conflict of Interest**: The authors declare that they have no conflict of interest.

493    **Tables and Figures**

494    **Table 1:** Summary of ncRNAs discovered in halite community

|  | Number (%)* | % in Archaea | % in Bacteria |
|---|---|---|---|
| **Total ncRNAs** | 1538 (100) | 54 | 46 |
| **Rfam ncRNAs** | 79 (5) | 73 | 27 |
| **Conserved sRNAs**\*\* | 155 (10) | 60 | 40 |
| **Antisense sRNAs** | 925 (60) | 40 | 60 |
| **Intergenic sRNAs** | 613 (40) | 75 | 25 |

495    *Percent from total ncRNAs
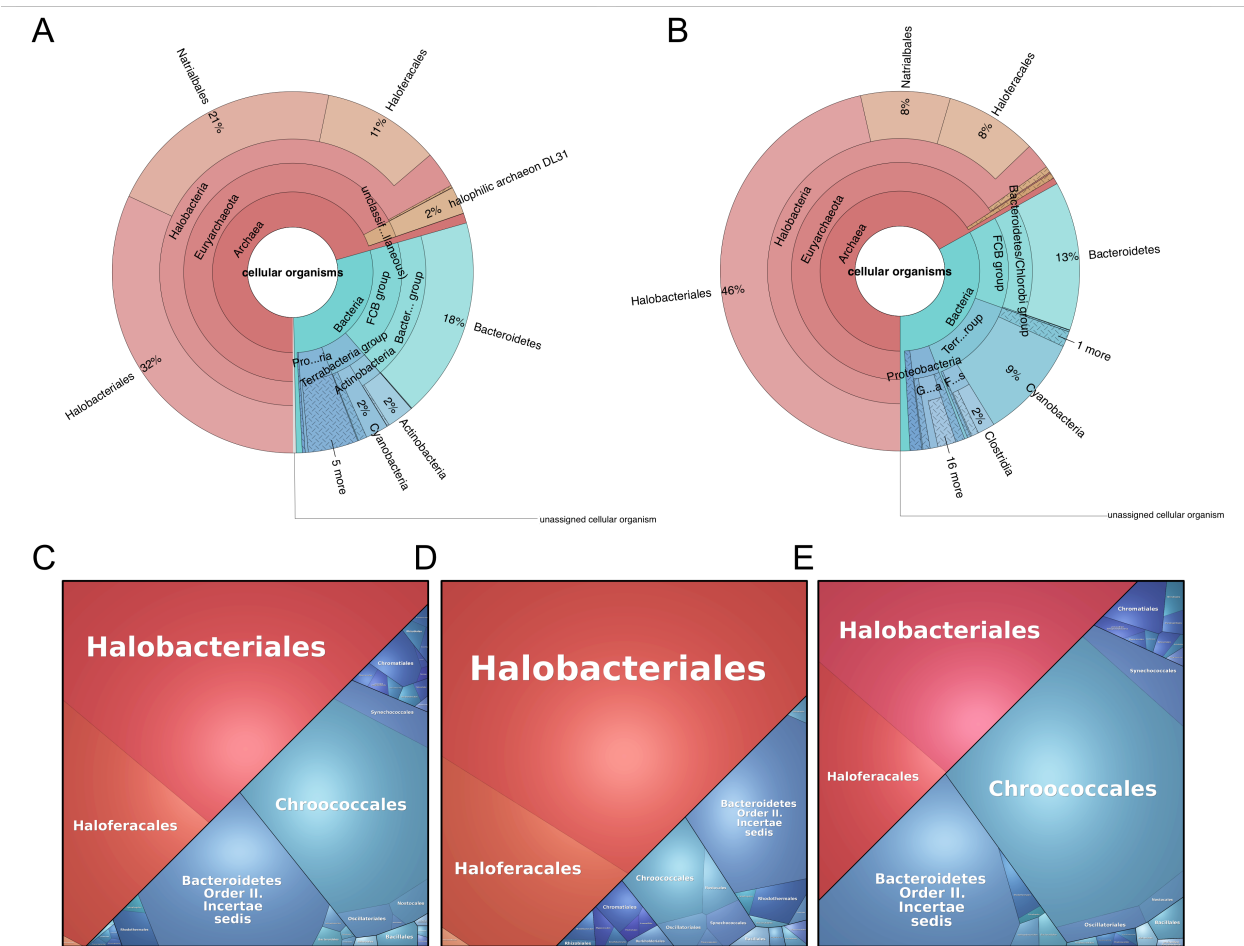496    ** Conserved other than Rfam ncRNAs
497

498

**Fig. 1** Taxonomic distribution. Krona graphs of (A) the halite metagenome based of DNA sequence reads and (B) the halite metatranscriptome based on RNA sequence reads; and Voronoi plots of (C) total sRNAs; (D) itsRNAs and (E) asRNAs discovered in the halite community.
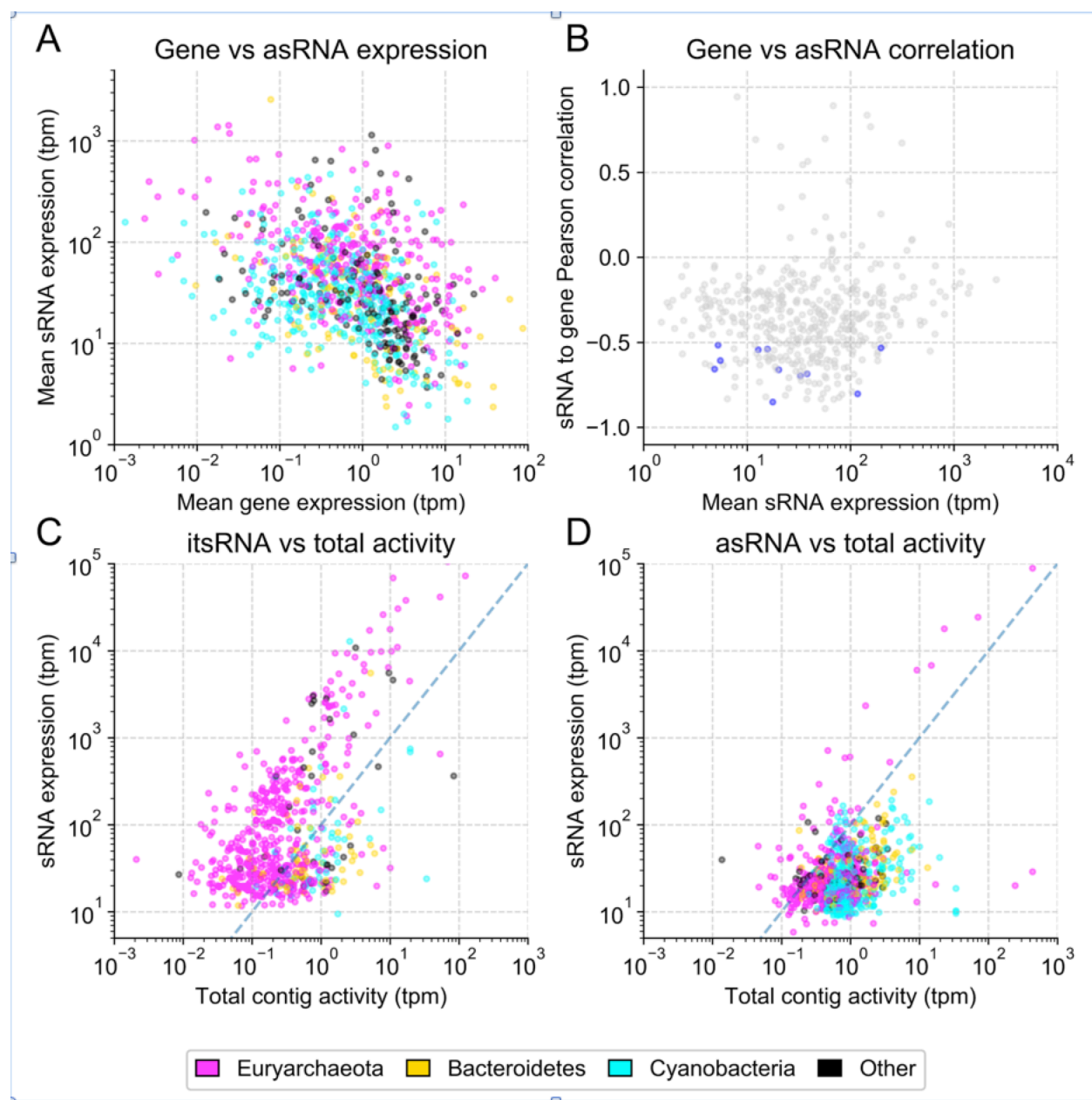
503

**Fig. 2** sRNA expression levels. (A) asRNAs and their putative targets mean expression
levels (TPM); (B) Pearson correlations in expression level between asRNAs and their
putative mRNA targets across the replicates, with significant correlations (pval<0.01)
highlighted in blue; (C) average expression of itsRNA and average expression of (D)
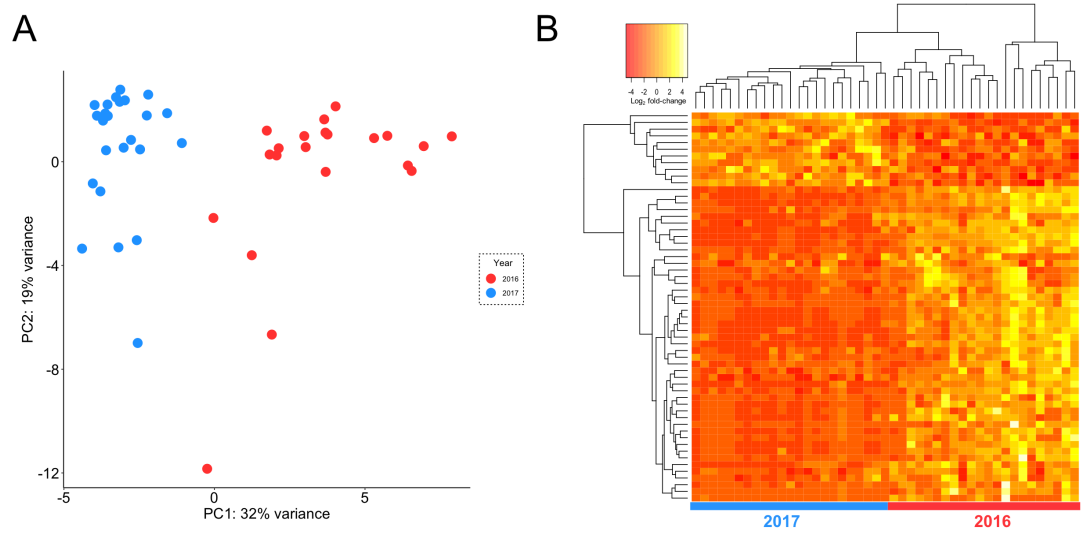asRNAs over the average expression of the contigs on which they are found.

509

510 **Fig. 3** itsRNA differential expression. (A) PCA plot showing itsRNA expression levels

511 clustered by year and (B) heat map of $\log_2$-transformed fold change for the top 50

512 significantly differentially expressed itsRNAs; each row is an itsRNA and each column a

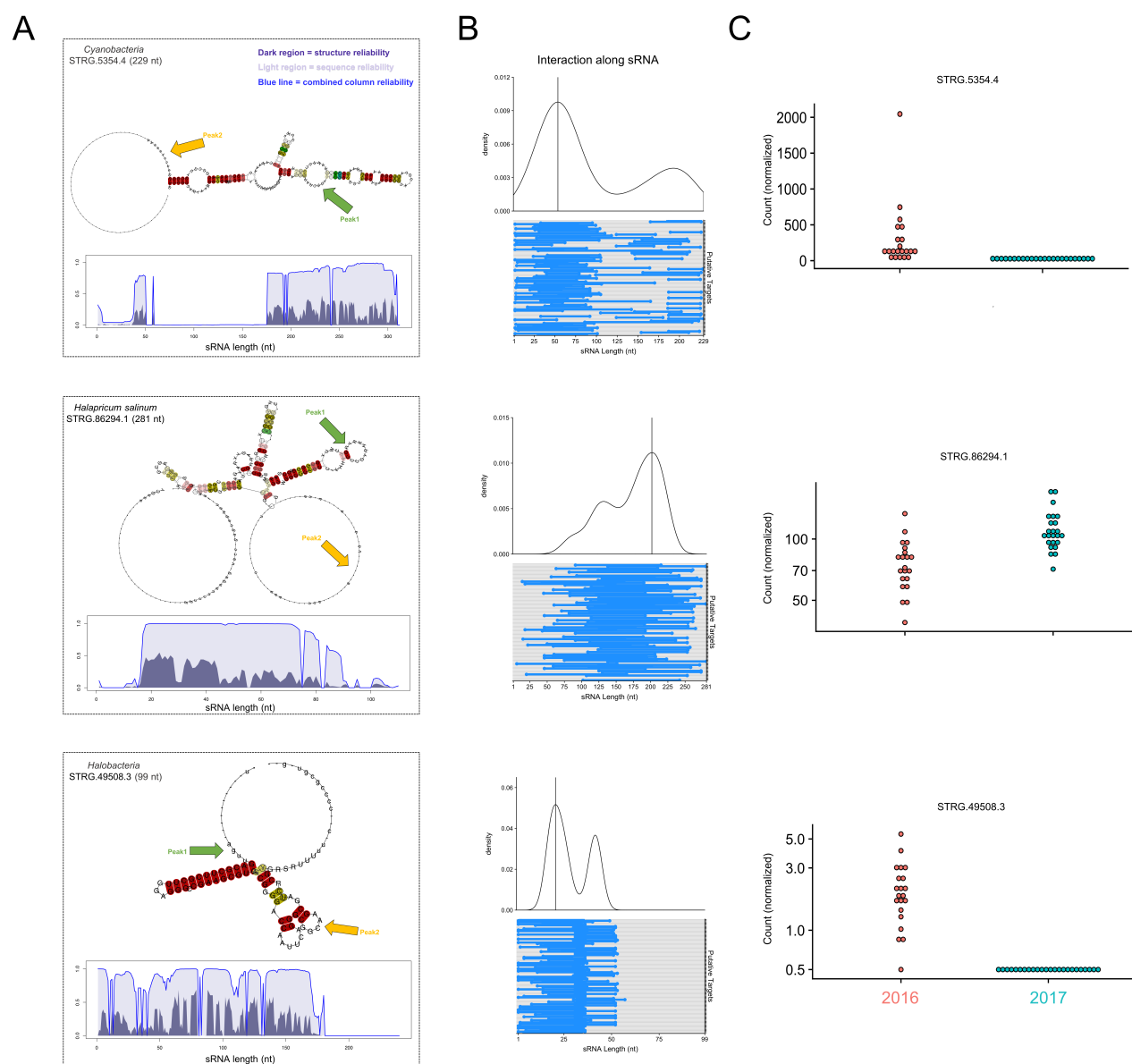513 sample collected in 2016 or 2017.

514

Fig. 4 Predicted structure, target identification, and expression levels for selected differentially expressed itsRNAs. (A) 2D-layout of consensus structures with base pairs coloring showing sequence and structure conservation and interactions peaks (green and yellow arrows); STAR profile plots with dark regions indicating structure reliability, light regions representing sequence reliability, and thin lines showing the combined column-reliability as computed by LocARNA-P. (B) Interaction plots of itsRNAs and their predicted targets. The top graphs are density plots calculated from the top 100 putative targets, and on the bottom are dumbbell plots of interactions (blue dumbbells) along the length of the itsRNA for the top 100 predicted mRNA targets; interaction peaks are

525    shown in green and yellow in the predicted structures; (C) Expression levels

526    represented as normalized count for each itsRNA in 2016 and in 2017 across all

527    samples.

528

529    **References**

530    Bao, G. H., M. J. Wang, T. G. Doak, and Y. Z. Ye. 2015. "Strand-specific community

531        RNA-seq reveals prevalent and dynamic antisense transcription in human gut

532        microbiota." *Front Microbiol* 6. doi: ARTN 89610.3389/fmicb.2015.00896.

533    Cai, Qiang, Lulu Qiao, Ming Wang, Baoye He, Feng-Mao Lin, Jared Palmquist, Sienna-

534        Da Huang, and Hailing Jin. 2018. "Plants send small RNAs in extracellular

535        vesicles to fungal pathogen to silence virulence genes." *Science* 360

536        (6393):1126. doi: 10.1126/science.aar4142.

537    Carrier, Marie-Claude, David Lalaouna, and Eric Massé. 2018. "Broadening the

538        Definition of Bacterial Small RNAs: Characteristics and Mechanisms of Action."

539        *Ann Rev Microbiol* 72 (1):141-161. doi: 10.1146/annurev-micro-090817-062607.

540    Cech, Thomas R, and Joan A Steitz. 2014. "The Noncoding RNA Revolution -Trashing

541        Old Rules to Forge New Ones." *Cell* 157 (1):77-94. doi:

542        10.1016/j.cell.2014.03.008.

543    Clouet-d'Orval, Béatrice, Manon Batista, Marie Bouvier, Yves Quentin, Gwennaele

544        Fichant, Anita Marchfelder, and Lisa-Katharina Maier. 2018. "Insights into RNA-

545        processing pathways and associated RNA-degrading enzymes in Archaea."

546        *FEMS Microbiol Rev* 42 (5):579-613. doi: 10.1093/femsre/fuy016.

547    Crits-Christoph, A., D.R. Gelsinger, B. Ma, J. Wierzchos, J. Ravel, C. Ascaso, O.

548        Artieda, A. Davila, and J. DiRuggiero. 2016. "Functional analysis of the archaea,

549        bacteria, and viruses from a halite endolithic microbial community." *Env.*

550        *Microbiol.* 18:2064-2077. doi: 10.1111/1462-2920.13259.

551    Davila, A.F., B. Gomez-Silva, A. de los Rios, C. Ascaso, H. Olivares, C.P. McKay, and

552        J. Wierzchos. 2008. "Facilitation of endolithic microbial survival in the hyperarid

553        core of the Atacama Desert by mineral deliquescence." *J. Geophys. Res.* 113

554        (G01028):G01028, doi:10.1029/2007JG000561.

555    Davila, A.F., I. Hawes, J. Garcia, D.R. Gelsinger, J. DiRuggiero, C. Ascaso, A. Osano,

556        and J. Wierzchos. 2015. " In situ metabolism in halite endolithic microbial

557        communities of the hyperarid Atacama Desert." *Front Microbiol*

558        http://dx.doi.org/10.3389/fmicb.2015.01035.

559    de Almeida, João Paulo Pereira, Ricardo Z. N. Vêncio, Alan P. R. Lorenzetti, Felipe Ten

560        Caten, José Vicente Gomes-Filho, and Tie Koide. 2019. "The Primary Antisense

561        Transcriptome of *Halobacterium salinarum* NRC-1." *Genes* 10 (4):280. doi:

562        10.3390/genes10040280.

563    Dyall-Smith, M. 2009. "The Halohandbook – Protocols for haloarchaeal genetics."

564        *Available at*

565        *https://www.researchgate.net/publication/278741334_The_Halohandbook_v73*

566    Eichhorn, Eric, Jan R. van der Ploeg, and Thomas Leisinger. 1999. "Characterization of

567        a Two-component Alkanesulfonate Monooxygenase from *Escherichia coli*." *J*

568        *Biol l Chem* 274 (38):26639-26646. DOI:10.1074/jbc.274.38.26639

569    Finstad, K. M., A. J. Probst, B. C. Thomas, G. L. Andersen, C. Demergasso, A.

570        Echeverria, R. G. Amundson, and J. F. Banfield. 2017. "Microbial Community

571        Structure and the Persistence of Cyanobacterial Populations in Salt Crusts of the

572        Hyperarid Atacama Desert from Genome-Resolved Metagenomics." *Front*

573        *Microbiol* 8. doi: ARTN 143510.3389/fmicb.2017.01435.

574    Gelsinger, D. R., and J. DiRuggiero. 2018a. "Transcriptional Landscape and Regulatory

575        Roles of Small Noncoding RNAs in the Oxidative Stress Response of the

576        Haloarchaeon Haloferax volcanii." *J Bacteriol* 200 (9). doi: ARTN e00779-

577        1710.1128/JB.00779-17.

578    Gelsinger, Diego, and Jocelyne DiRuggiero. 2018b. "The Non-Coding Regulatory RNA

579        Revolution in Archaea." *Genes* 9 (3):141. doi: 10.3390/genes9030141.

580    Gottesman, S., and G. Storz. 2011. "Bacterial small RNA regulators: versatile roles and

581        rapidly evolving variations." *Cold Spring Harb Perspect Biol* 3 (12). doi:

582        10.1101/cshperspect.a003798.

583    Kalvari, Ioanna, Joanna Argasinska, Natalia Quinones-Olvera, Eric P. Nawrocki, Elena

584        Rivas, Sean R. Eddy, Alex Bateman, Robert D. Finn, and Anton I. Petrov. 2017.

585          "Rfam 13.0: shifting to a genome-centric resource for non-coding RNA families."

586          *NAR* 46 (D1):D335-D342. doi: 10.1093/nar/gkx1038.

587 Kish, A., G. Kirkali, C. Robinson, R. Rosenblatt, P. Jaruga, M. Dizdaroglu, and J.

588          DiRuggiero. 2009. "Salt shield: intracellular salts provide cellular protection

589          against ionizing radiation in the halophilic archaeon, Halobacterium salinarum

590          NRC-1." *Environ Microbiol* 11 (5):1066. DOI: 10.1111/j.1462-2920.2008.01828.x

591 Kliemt, Jana, Katharina Jaschinski, and Jörg Soppa. 2019. "A Haloarchaeal Small

592          Regulatory RNA (sRNA) Is Essential for Rapid Adaptation to Phosphate

593          Starvation Conditions." *Front Microbiol* 10:1219-1219. doi:

594          10.3389/fmicb.2019.01219.

595 Krulwich, Terry A., David B. Hicks, and Masahiro Ito. 2009. "Cation/proton antiporter

596          complements of bacteria: why so large and diverse?" *Mol Microbiol* 74 (2):257-

597          260. doi: 10.1111/j.1365-2958.2009.06842.x.

598 Liao, Yang, Gordon K. Smyth, and Wei Shi. 2014. "featureCounts: an efficient general

599          purpose program for assigning sequence reads to genomic features."

600          *Bioinformatics* 30 (7):923-930. doi: 10.1093/bioinformatics/btt656.

601 Lloréns-Rico, Verónica, Jaime Cano, Tjerko Kamminga, Rosario Gil, Amparo Latorre,

602          Wei-Hua Chen, Peer Bork, John I. Glass, Luis Serrano, and Maria Lluch-Senar.

603          2016. "Bacterial antisense RNAs are mainly the product of transcriptional noise."

604          *Science adv* 2 (3):e1501363-e1501363. doi: 10.1126/sciadv.1501363.

605 Love, M. I., W. Huber, and S. Anders. 2014. "Moderated estimation of fold change and

606          dispersion for RNA-seq data with DESeq2." *Genome Biol* 15 (12):550. doi:

607          10.1186/s13059-014-0550-8.

608 Mann, Martin, Patrick R. Wright, and Rolf Backofen. 2017. "IntaRNA 2.0: enhanced and

609          customizable prediction of RNA–RNA interactions." *NAR* 45 (W1):W435-W439.

610          doi: 10.1093/nar/gkx279.

611 Meslier, V., M.C. Casero, M. Daily, J. Wierchos, C. Ascaso, O. Artieda, P.R.

612          McCullough, and J. DiRuggiero. 2018. "Fundamental drivers for endolithic

613          microbial community assemblies in the hyperarid Atacama Desert." *Env.*

614          *Microbiol.* 20:1765-1781. https://doi.org/10.1111/1462-2920.14106.

615  Molina-Henares, Antonio J., Tino Krell, Maria Eugenia Guazzaroni, Ana Segura, and
616      Juan L. Ramos. 2006. "Members of the IclR family of bacterial transcriptional
617      regulators function as activators and/or repressors." *FEMS Microbiol Rev* 30
618      (2):157-186. doi: 10.1111/j.1574-6976.2005.00008.x.
619  Nitzan, Mor, Rotem Rehani, and Hanah Margalit. 2017. "Integration of Bacterial Small
620      RNAs in Regulatory Networks." *Ann Rev Biophys* 46 (1):131-148. doi:
621      10.1146/annurev-biophys-070816-034058.
622  Oren, Aharon. 2013. "Life at high salt concentrations, intracellular KCl concentrations,
623      and acidic proteomes." *Front Microbiol* 4:315-315. doi:
624      10.3389/fmicb.2013.00315.
625  Perez-Rueda, Ernesto, Rafael Hernandez-Guerrero, Mario Alberto Martinez-Nuñez,
626      Dagoberto Armenta-Medina, Israel Sanchez, and J. Antonio Ibarra. 2018.
627      "Abundance, diversity and domain architecture variability in prokaryotic DNA-
628      binding transcription factors." *PloS one* 13 (4):e0195332-e0195332. doi:
629      10.1371/journal.pone.0195332.
630  Pointing, Stephen B., and Jayne Belnap. 2012. "Microbial colonization and controls in
631      dryland systems." *Nature Rev Microbiol* 10:551. doi: 10.1038/nrmicro2831.
632  Raman, A.V. , A.n López García de Lomana, U. Kusebauch, M. Pan, S. Turkarslan, R.
633      L. Moritz, and N. S. Baliga. 2018. "Context-Specific Regulation of Coupled
634      Transcription-Translation Modules Predicts Pervasive Ribosome Specialization."
635      *Available at SSRN:* http://dx.doi.org/10.2139/ssrn.3155765
636  Robinson, C. K., J. Wierzchos, C. Black, A. Crits-Christoph, B. Ma, J. Ravel, C. Ascaso,
637      O. Artieda, S. Valea, M. Roldan, B. Gomez-Silva, and J. DiRuggiero. 2015.
638      "Microbial diversity and the presence of algae in halite endolithic communities are
639      correlated to atmospheric moisture in the hyper-arid zone of the Atacama
640      Desert." *Environ Microbiol* 17:299-315. doi: 10.1111/1462-2920.12364.
641  Shi, Y. M., G. W. Tyson, and E. F. DeLong. 2009. "Metatranscriptomics reveals unique
642      microbial small RNAs in the ocean's water column." *Nature* 459 (7244):266-
643      U154. doi: 10.1038/nature08055.
644  Takara, T. J., and S. P. Bell. 2009. "Putting two heads together to unwind DNA." *Cell*
645      139 (4):652-654. doi: 10.1016/j.cell.2009.10.037.

646    Toyofuku, Masanori, Nobuhiko Nomura, and Leo Eberl. 2019. "Types and origins of
647            bacterial membrane vesicles." *Nature Rev Microbiol* 17 (1):13-24. doi:
648            10.1038/s41579-018-0112-2.

649    Tsatsaronis, James A., Sandra Franch-Arroyo, Ulrike Resch, and Emmanuelle
650            Charpentier. 2018. "Extracellular Vesicle RNA: A Universal Mediator of Microbial
651            Communication?" *Trends Microbiol* 26 (5):401-410. doi:
652            10.1016/j.tim.2018.02.009.

653    Uritskiy, G.V., and J. DiRuggiero. 2019. "Applying Genome-Resolved Metagenomics to
654            Deconvolute the Halophilic Microbiome." *Gene* 10:220. doi:
655            10.3390/genes10030220.

656    Uritskiy, Gherman, Samantha Getsin, Adam Munn, Benito Gomez-Silva, Alfonso Davila,
657            Brian Glass, James Taylor, and Jocelyne DiRuggiero. 2019. "Halophilic microbial
658            community compositional shift after a rare rainfall in the Atacama Desert." *ISMEJ*
659            doi: 10.1038/s41396-019-0468-y.

660    Uritskiy, Gherman V., Jocelyne DiRuggiero, and James Taylor. 2018. "MetaWRAP—a
661            flexible pipeline for genome-resolved metagenomic data analysis." *Microbiome* 6
662            (1):158. doi: 10.1186/s40168-018-0541-1.

663    Wagner, E. Gerhart H., and Pascale Romby. 2015. "Chapter Three - Small RNAs in
664            Bacteria and Archaea: Who They Are, What They Do, and How They Do It." In
665            *Advances in Genetics*, edited by Theodore Friedmann, Jay C. Dunlap and
666            Stephen F. Goodwin, 133-208. Academic Press.

667    Wassarman, Karen M. 2018. "6S RNA, a Global Regulator of Transcription." *Microbiol*
668            *Spectr* 6 (3). doi: 10.1128/microbiolspec.RWR-0019-2018.

669    Weinberg, Zasha, Christina E. Lünse, Keith A. Corbino, Tyler D. Ames, James W.
670            Nelson, Adam Roth, Kevin R. Perkins, Madeline E. Sherlock, and Ronald R.
671            Breaker. 2017. "Detection of 224 candidate structured RNAs by comparative
672            analysis of specific subsets of intergenic regions." *NAR* 45 (18):10811-10823.
673            doi: 10.1093/nar/gkx699.

674    Will, Sebastian, Tejal Joshi, Ivo L. Hofacker, Peter F. Stadler, and Rolf Backofen. 2012.
675            "LocARNA-P: Accurate boundary prediction and improved detection of structural
676            RNAs." *RNA* 18 (5):900-914. doi: 10.1261/rna.029041.111.

677    Wyss, Leander, Melanie Waser, Jennifer Gebetsberger, Marek Zywicki, and Norbert
678         Polacek. 2018. "mRNA-specific translation regulation by a ribosome-associated
679         ncRNA in Haloferax volcanii." *Scientific Rep* 8 (1):12502. doi: 10.1038/s41598-
680         018-30332-w.
681