

# 1 **The Spot 42 RNA: A regulatory small RNA with roles in the central** 2 **metabolism**

3 Cecilie Bækkedal and Peik Haugen\*

4

5 Department of Chemistry, The Norwegian Structural Biology Centre (NorStruct) and Centre for  
6 Bioinformatics (SfB), UiT – The Arctic University of Norway, 9037 Tromsø, Norway

7

8 Key words: sRNA, small RNA, Spot 42, *spf*, non-coding RNA, gamma proteobacteria, pirin.

9 \*Correspondence to: Peik Haugen; E-mail: peik.haugen@uit.no

10 Disclosure statement: the authors have no conflict of interest and nothing to disclose.

11

12 The Spot 42 RNA is a 109 nucleotide long (in *Escherichia coli*) noncoding small regulatory RNA  
13 (sRNA) encoded by the *spf* (spot fourty-two) gene. *spf* is found in gamma-proteobacteria and the  
14 majority of experimental work on Spot 42 RNA has been performed using *E. coli*, and recently  
15 *Aliivibrio salmonicida*. In the cell Spot 42 RNA plays essential roles as a regulator in carbohydrate  
16 metabolism and uptake, and its expression is activated by glucose, and inhibited by the cAMP-CRP  
17 complex. Here we summarize the current knowledge on Spot 42, and present the natural distribution  
18 of *spf*, show family-specific secondary structural features of Spot 42, and link highly conserved  
19 structural regions to mRNA target binding.

## 20 **Introduction**

21 The *spf* gene is highly conserved in *Escherichia*, *Shigella*, *Klebsiella*, *Salmonella* and *Yersinia*  
22 (genera) within the *Enterobacteriaceae* family.<sup>1</sup> In *E. coli* the *spf* gene is flanked by *polA*  
23 (upstream) and *yihA* (downstream),<sup>2,3</sup> and a CRP binding sequence and -10 and -35 promoter  
24 sequences are found upstream of *spf*. *spf* is also highly conserved within the *Vibrionaceae* family,  
25 and was recently identified in 76 *Vibrionaceae* genomes that were available at that time (e.g.,  
26 *Vibrio*, *Aliivibrio*, *Photobacterium* and *Grimontia* genera).<sup>4</sup> In e.g., *Vibrio cholerae*, *Vibrio vulnificus*,  
27 *Aliivibrio fischeri* and *A. salmonicida* the *spf* gene is flanked by *polA* (upstream) and a sRNA gene  
28 encoding the novel VSsrna24 RNA (downstream).

29 Spot 42 was first described in 1973 as an unstable RNA species of 109 nucleotides in *E.*  
30 *coli*.<sup>5,6</sup> It was discovered by polyacrylamide gel electrophoresis and 2-D fingerprinting in an  
31 attempt to study the accumulation of small RNAs in *E. coli* during amino acid starvation. In these  
32 experiments the electrophoretic mobility of Spot 42 was similar to that of 5S rRNA. In 1979, Spot  
33 42 was reported to accumulate under growth in the presence of glucose (i.e., when adenosine  
34 3',5'-cyclic monophosphate (cAMP) is low).<sup>7,8</sup> During growth with a non-glucose carbon source  
35 (i.e., when cAMP concentrations are high) Spot 42 concentrations were significantly lower. Later  
36 experiments showed that overexpression of Spot 42 (tenfold increase) resulted in impaired  
37 growth and lowered ability to adapt to shifts to richer media.<sup>9</sup> Further, shift from glucose to  
38 succinate as the carbon source resulted in a long lag period and slow growth rate, the reason for  
39 the abnormal responses was caused by an elevated number of excessive Spot 42 RNA gene  
40 products rather than excess of the gene itself. A deletion study of *spf* in *E. coli* cells resulted in  
41 viable *spf* null mutants, which indicated that Spot 42 was non-essential, at least under controlled  
42 lab conditions.<sup>10</sup>

43 It was for some years unclear if the function of Spot 42 was mediated through the 109  
44 nucleotide RNA itself or if the function was mediated through the 14 amino acids long peptide  
45 which is hypothetically encoded from within the sRNA sequence. This confusion was based on  
46 the observation that Spot 42 contains structural features similar to other non-coding RNAs  
47 found in *E. coli* (such as 6S RNA and lambda bacteriophage), as well as features that are typically  
48 found in mRNAs (i.e., polypurine sequence followed by AUG, 14 amino acids codons and an UGA  
49 termination codon).<sup>7</sup> Using a filter binding assay and other methods Rice et al. showed that Spot  
50 42 is not an mRNA.<sup>11</sup> In this approach the affinity between Spot 42 and the 70S ribosome was  
51 tested. Here, Spot 42 showed very inefficient binding to purified 70S ribosomes, which lead to  
52 the conclusion that the function of Spot 42 is mediated by the RNA itself.

53 The direct responsiveness of Spot 42 levels to glucose and cAMP is due to repression of  
54 *spf* expression by a cAMP-CRP (cAMP-receptor protein) complex.<sup>2</sup> The reduction of Spot 42 in  
55 cells grown in secondary carbon sources is a result of binding of the cAMP-CRP complex to the  
56 *spf* promoter, which negatively regulates transcription of Spot 42. Later, the proximity of *spf* to  
57 *polA* (gene encoding DNA polymerase I) led Dahlberg and co-workers to test whether the  
58 products of these genes could influence each other.<sup>12</sup> They found that by reducing levels of Spot  
59 42, either by deletion of *spf* or by manipulating the growth conditions, the DNA pol A activity  
60 was reduced. The underlying mechanism for this observation remains however unknown.

61 Spot 42 can interact directly with mRNA targets through base pairing. The first Spot 42  
62 target was discovered by Møller et al., who showed that Spot 42 specifically binds to a short

63 complementary region at the translation initiation region of *galK* (encodes a galactokinase)  
64 mediated through binding of the posttranscriptional regulator Hfq.<sup>1</sup> *galK* is the third gene in the  
65 galactose operon, which contains four genes (*galETKM*) and produces a polycistronic mRNA.  
66 Spot 42 mediates discoordinate expression of the *gal* operon (i.e., the individual genes in the  
67 operon are not similarly expressed) by binding to the *galK* Shine-Dalgarno region, thereby  
68 blocking ribosome binding and translation of the *galK* gene. The physiological significance of the  
69 discoordinate expression is unclear, but suggests that Spot 42 plays a role in fine-tuning gene  
70 expression to optimize the utilization of carbon sources. Recently, Wang et al. showed that Spot  
71 42 represses expression of *galK* through direct binding to the 5' end of the *galK* mRNA, and also  
72 mediates transcription termination of *galT* in the *galT-galK* junction.<sup>13</sup>

73 Beisel and Storz demonstrated with microarray analysis and reporter fusions that Spot  
74 42 plays a broader role in metabolism by regulating at least 14 operons.<sup>14</sup> These operons contain  
75 a number of genes involved in uptake and catabolism of non-favored carbon sources. During  
76 overexpression of Spot 42 sixteen different genes showed consistently twofold reduced or  
77 elevated levels of mRNA. The identified reduced genes are mostly involved in central and  
78 secondary metabolism, as well as uptake and catabolism of non-preferred carbon sources and  
79 oxidation of NADH. In 2012 Beisel et al. performed computational target analysis using the  
80 three conserved regions of Spot 42 as input. Compared to when using full-length Spot 42  
81 sequence as input the target identification was improved and additional targets were revealed.<sup>15</sup>  
82 The target analysis combined with assays of reporter fusions identified seven novel Spot 42  
83 mRNA targets, all involved in catabolite repression. Mutational analysis showed that the  
84 interactions of the three conserved regions of Spot 42 are critical in target regulation and that  
85 regulation through multiple conserved regions of Spot 42 as well as increased base-pairing in  
86 these regions strengthen the target regulation.

87 The evolution of sRNAs in *E. coli* and their regulatory interactions with mRNAs was  
88 recently studied using computational methods.<sup>16</sup> Compared to cis-acting sRNA and other non-  
89 coding RNA (housekeeping RNA), trans-acting sRNA was the latest to appear in evolution.  
90 Furthermore, after Enterobacteriales diverged into a separate lineage within gamma-  
91 proteobacteria, the trans-acting sRNAs likely appeared in relatively high numbers compared to  
92 the cis-acting sRNAs that evolved more evenly among all orders within gamma-proteobacteria.  
93 The evolutionary age of 15 sRNAs and 49 corresponding sRNA-mRNA interactions were  
94 examined. Here, Spot 42 was found to be the most ancient sRNA. Of the six Spot 42 mRNA  
95 targets considered, only two (*xytF* and *galK*) evolved before Spot 42, albeit all the  
96 complementary mRNA binding sites appeared after Spot 42.

97           The observation that *A. salmonicida* contains the *spf* gene (which encodes the Spot 42  
98 RNA), but lacks the *galK* operon (the natural Spot 42 target in *E. coli*), have inspired scientists to  
99 study the role of Spot 42 in this fish pathogen.<sup>4</sup> *A. salmonicida* is unable to utilize galactose (lacks  
100 *gal* operon) in minimal medium and addition of galactose has little effect on the growth rate.  
101 When cells are grown in glucose the level of Spot 42 is increased 16-40 fold, but is in contrast  
102 decreased threefold when cAMP is added, indicating that Spot 42 have similar roles as in *E. coli*,  
103 i.e., in carbohydrate metabolism. It has been hypothesized that Spot 42 works in concert with a  
104 novel sRNA gene, called *VSsrna24*, located 262 nt downstream of *spf*. The *VSsrna42* RNA is  
105 approximately 60 nt in length and has an expression pattern opposite to that of Spot 42.  
106 Furthermore, in a *spf* deletion mutant a gene encoding a pirin-like protein was upregulated 16  
107 fold. Pirin has key roles in the central metabolism by regulating the activity of pyruvate  
108 dehydrogenase E1 and therefore select whether pyruvate will be fermented, or subjected to  
109 respiration through the TCA cycle and electron transport.

110           Although the Spot 42 RNA was discovered more than 40 years ago there are still a  
111 number of unanswered question related to this highly interesting RNA, e.g.: What is the natural  
112 distribution of the Spot 42 gene (*spf*) in Bacteria? What is the complete set of biological roles of  
113 Spot 42, and does Spot 42 play the proposed key role in the central metabolism? How does Spot  
114 42 interact with its apparently many mRNA targets? In this work we have summarized the  
115 current literature on Spot 42, and extended this knowledge by surveying the known natural  
116 distribution of *spf*, we have identified family-specific structural features of Spot 42, and  
117 evaluated if highly conserved structural regions can be linked to mRNA binding.

## 118 **Results**

### 119 ***spf* is restricted to 5 orders of gamma-proteobacteria**

120           The distribution of *spf* in nature is shown in **Fig. 1**. The basis for the figure was available  
121 nucleotide sequences of *spf* included in the Rfam database (677 sequences), and *spf* sequences  
122 identified in this study by using the Blastn server and *spf* sequences from selected taxa as  
123 queries. All previously known cases of *spf* originate from gamma-proteobacteria, and after  
124 fruitless searches in all other domains of Bacteria we therefore concentrated our efforts on  
125 specific searches within gamma-proteobacteria, both by using *spf* sequences from the closest  
126 neighbors, and by manual inspection of the known genic location of *spf*, i.e., in the intergenic  
127 region between *polA* and *engB*. The result of our search was finally mapped onto a phylogenetic  
128 tree generated using the iTOL web service.

129

130 The result show that *spf* is exclusively found in five orders of gamma-bacteria, i.e., in  
131 Enterobacteriales, Aeromonadales, Alteromonadales, Vibrionales and Chromatiales. These  
132 orders, except Chromatiales, share the same closest common ancestor (arrow in **Fig. 1**), and  
133 constitutes a clade. *spf* has still not been found in Pasteurellales, which is likely due to that  
134 Pasteurellales genomes are underrepresented in the European Nucleotide Archive (ENA)  
135 compared to e.g., the sister Enterobacteriales. We suspect that *spf* will be discovered in  
136 Pasteurellales as more genomes are being sequenced. In addition to known cases of *spf* our  
137 Blastn search revealed previously unreported cases within genera of Enterobacteriales and  
138 Alteromonadales. In Enterobacteriales *spf* was identified in the genera *Morganella* and  
139 *Raoultella*, as well as in draft genomes of *Budvicia*, *Cedecea*, *Hafnia*, *Leminorella*, *Plesimonas* and  
140 *Yokenella*. And, in genera where *spf* was already known to occur, *spf* was in this work identified  
141 in *Enterobacter radicincitans* and *Escherichia blattae*. Similarly, in Alteromonadales *spf* is found  
142 in the five families *Ferrimonadaceae*, *Shewanellaceae*, *Moritellaceae*, *Pseudoalteromonadaceae*  
143 and *Alteromonadaceae*, and *spf* was in this study identified in the three genera *Glaccola*,  
144 *Alteromonas* and *Pseudoalteromonas* by our blast searches, whereas *spf* was found in *Moritella*  
145 *viscosa* by manual inspection of the intergenic region *polA/engB*. Interestingly, in Chromatiales,  
146 *spf* is exclusively found in the genera *Rheinheimera* and *Arsukibacterium*, which is represented in  
147 ENA by six and two available draft genomes, all containing *spf*. Given that the phylogeny as  
148 shown in **Fig. 1** is correct then it is tempting to speculate that *spf* was acquired by lateral  
149 transfer, perhaps from a donor within the clade marked by an arrow in **Fig. 1**.

150  
151 We also wanted to answer the following question: Is *spf* optional or ubiquitous within  
152 the individual orders and families? Spot 42 appears to play central roles in the carbohydrate  
153 metabolism, and we therefore hypothesized that it might be present in all representatives of the  
154 same order, family or genus once it has been identified in one genome. To answer this question  
155 we used the list of complete bacterial genomes found at the NCBI Genomes resource  
156 (<http://www.ncbi.nlm.nih.gov/genome/>), and searched for presence of *spf* in all representatives  
157 of the current orders, families and genera. Our result show that *spf* is found in 699 of 741  
158 complete genomes distributed among 34 genera (a detailed list is provided in **Table S1**). *spf* is  
159 missing in representatives of the two genera *Glaccola* and *Pseudoalteromonas* of  
160 Alteromonadales. In both of these genera *spf* is found in one of three complete genomes. All  
161 three genomes of *Glaccola* have the same genic organization with *polA* and *engB* as neighbors  
162 (*spf* is usually located between these two genes). In *Pseudoalteromonas*, *spf* is only found in one  
163 genome, i.e., in *Pseudoalteromonas atlantica*, where *polA* and *engB* are located next to each other.  
164 The two other genomes with no *spf* have a different genic organization (syteny) at this region.  
165 Finally, *spf* has not been found in any of the complete genomes within the following genera:

166 *Buchnera*, *Candidatus Moranella*, *Candidatus Riesia* and *Wigglesworthia* (from Enterobacteriales),  
167 *Oceanomoas* and *Tolomonas* (from Aeromonadales), *Marinobacter*, *Sacchrophagus*, *Colwellia*,  
168 *Idiomarina* and *Psychromonas* (from Alteromonadales), and all genera of Chromatiales (i.e., *spf*  
169 found in 6 draft genomes of the genus *Rheinheimera* and 2 draft genomes of *Arsukibacterium*). In  
170 summary, of a total of 741 genomes from the 5 orders *Enterobacteriales*, *Aeromonadales*,  
171 *Alteromonadales*, *Chromatiales* and *Vibrionales*, 699 complete genomes contain *spf*, whereas 42  
172 lack *spf*. The result is in agreement with conserved, but not necessarily indispensable roles of *spf*.

173

#### 174 **The Spot 42 RNA consensus secondary structure**

175 We next mapped the level of identity among all known *spf* sequences (120 in total when redundant  
176 sequences have been removed) onto a consensus secondary structure model of Spot 42 (based on  
177 structure probing by Møller et al.<sup>1</sup>) to find clues to possible structural regions that might be  
178 important for target identification and interaction, in general (**Fig. 2**). The Spot 42 RNA consists of  
179 one long hairpin structure located at the 5' end (from now on referred to as the 5' hairpin; 45–59 nt  
180 in length), and a second smaller hairpin separated from the 5' hairpin by a 9 - 20 nt long single-  
181 stranded region. In addition, a rho-independent terminator is located immediately downstream of  
182 the second hairpin. Structural regions of Spot 42 from the families *Vibrionaceae*, *Aeromonadaceae*  
183 and *Shewanellaceae* differ from the general “consensus” and are shown in separate boxes in **Fig. 2**.  
184 The sRNA gene is, in general, highly conserved with 76 of 108 positions (when using the “consensus”  
185 sequence as the reference) being 80–100% identical across all orders (shown as uppercase bold  
186 letters in **Fig. 2**). Notably, the 5' hairpin is highly conserved, i.e., 80–100% identity from positions  
187 1–41, which indicate that these positions are interesting candidates for having general roles in target  
188 binding, perhaps with the terminal loop functioning as the seed sequence. The single-stranded region  
189 separating the 5' hairpin and the second hairpin is less conserved, with 80–100% identity in three  
190 positions and 60–79% identity in six positions, and is therefore perhaps less likely to have general  
191 roles in target recognition. *spf* is as expected most conserved within families. The *Shewanellaceae spf*  
192 differs most from the “consensus”. Here, the 5' hairpin contains two bulges with eight additional nt  
193 (inserted between pos. 39 and pos. 47). The *Vibrionaceae* and *Aeromonadaceae* sequences also  
194 differ to some extent from the “consensus”. In summary, Spot 42 is a highly conserved sRNA across  
195 five orders. The 5' hairpin represents the most conserved region and is therefore expected to have  
196 general roles in target recognition and interaction.

197

## 198 **Spot 42 structure conservation and potential base pairing with targets**

199 We next wanted to investigate if the highly conserved nucleotide positions of Spot 42 (as  
200 described above) are implicated in target binding (i.e., base-pairing between Spot 42 and mRNA  
201 target). Interactions between Spot 42 and *galK* mRNA has been determined using structure  
202 probing,<sup>1</sup> whereas potential base-pairing to other targets is based on bioinformatics predictions  
203 followed by experimental work.<sup>4,14,15</sup>

204 **Fig. 3** shows schematically potential base-pairing between Spot 42 and experimentally verified  
205 mRNA targets for the following genes: *galK*, *pirin*, *fucl*, *xylF*, *sthA*, *gltA*, *srlA*, *nanC*, *paaK*, *ascF*,  
206 *caiA*, *fucP*, *atoD*, *puuE* and *nanT*. Interestingly, for all except two genes (i.e., *sthA* and *fucP*) the  
207 most conserved region of the 5' hairpin (i.e., pos. 1-41) can potentially participate in extensive  
208 base-pairing with the corresponding mRNAs. This suggests that the 5' hairpin, is essential for  
209 target recognition and binding. Moreover, the first six positions of Spot 42 (5' single stranded  
210 region) can potentially base-pair with ten of fifteen targets (*galK*, *pirin*, *fucl*, *xylF*, *gltA*, *nanC*,  
211 *paaK*, *ascF*, *atoD* and *nanT*), and the terminal loop of the 5' hairpin can base-pair with eight of  
212 fifteen targets (*galK*, *pirin*, *fucl*, *xylF*, *srlA*, *caiA*, *puuE* and *nanT*). The second hairpin is only partly  
213 conserved. In agreement with this observation base-pairing with targets are rarer and only  
214 observed for two targets (*galK* and *pirin*). This is in agreement with results from Beisel et al.<sup>15</sup>  
215 Using three unstructured regions (the 5' single stranded region, the 5' hairpin and the single-  
216 stranded region separating the hairpins) as input during computational target identification, they  
217 improved identification of direct targets, compared to when using the full-length sequence of  
218 Spot 42. In summary, highly conserved nucleotide positions of Spot 42 have the potential to  
219 participate in extensive base-pairing with known mRNA targets.

220

## 221 **sRNA genes in the intergenic region downstream of *polA***

222 Interestingly, *spf* is not the only sRNA gene located in the intergenic region downstream of *polA*  
223 (see **Fig. 4**). In *Vibrionaceae* a gene encoding the sRNA VSsrna24 is located approximately 600 nt  
224 downstream of *spf*. Expression of VSsrna24 is repressed by glucose, and is hypothesized to have  
225 roles in the central carbohydrate metabolism.<sup>4</sup> The sRNAs sX13,<sup>17</sup> ErsA<sup>18</sup> and Smr7C,<sup>19,20</sup> are  
226 found in *Xanthomonadaceae*, *Pseudomonas* and Rhizobiales, respectively, but neither has the  
227 same function or structure as Spot 42. sX13 and Smr7C share secondary structure features  
228 comprising three stem-loops with C-rich motifs and are Hfq-independent.<sup>17,21</sup> ErsA is Hfq-  
229 mediated and regulated by sigma factor 22, in contrast to Spot 42 that is dependent on sigma

230 factor 70. If any of these four sRNA genes originates from a common ancestral gene or not is  
231 currently unknown.

## 232 **Concluding Remarks**

233 We have conducted a survey on Spot 42 RNA in order to learn about its natural distribution,  
234 conservation patterns, and mRNA target recognition. We demonstrated that Spot 42, which was  
235 first identified in *E. coli* (Enterobacteriales), is also common in four other orders, i.e.,  
236 Aeromonadales, Alteromonadales, Chromatiales and Vibrionales. Using blastn analysis we  
237 discovered novel *spf* sequences. Of a total of 741 complete genomes from the 5 orders  
238 Enterobacteriales, Aeromonadales, Alteromonadales, Chromatiales and Vibrionales, 699  
239 genomes contain *spf*. Furthermore, a total of 30 draft genomes distributed among 11 genera  
240 (from all orders except Aeromonadales) contain *spf*. As shown in **Fig. 1**, within gamma-  
241 proteobacteria, Aeromonadales, Alteromonadales, Enterobacteriales and Vibrionales share the  
242 same last common ancestor, whereas Chromatiales does not, which suggest that *spf* was  
243 introduced into Chromatiales by lateral transfer by a donor from the clade marked by an arrow.  
244 We made a consensus secondary structure model of Spot 42 based on all known *spf* sequences and  
245 compared this to a schematically figure showing potential base-pairing between Spot 42 and known  
246 mRNA targets. Our results show that highly conserved nucleotide positions, in general, have  
247 potential to participate in extensive base-pairing with target mRNAs. This is in agreement with  
248 an earlier study by Beisel et al. which suggested that the strength of Spot 42 regulation is  
249 directly dependent on the number of nucleotides and the number of highly conserved structural  
250 regions which are involved in base-pairing between Spot 42 and its target.<sup>15</sup>

251 It is intriguing to us that although Spot 42 was discovered more than 40 years, there are  
252 still many unanswered questions. As more sequence data are being produced from high-  
253 throughput sequencing techniques and better tools and search algorithms are being developed,  
254 the known natural distribution of *spf* will certainly expand to new orders, families and genera  
255 (and perhaps phyla). And detailed knowledge on target recognition (other than *galk*) and roles  
256 in cellular processes will come from functional and bioinformatics studies. One particularly  
257 interesting aspect of Spot 42 is its apparent central role (via pirin) in the central metabolism by  
258 directing pyruvate towards fermentation or respiration through the tricarboxylic acid (TCA)  
259 cycle and electron transport.

## 260 **Materials and Methods**

### 261 **Homology search**



262 All previously known *spf* sequences were retrieved from Rfam  
263 (<http://rfam.sanger.ac.uk/family/RF00021>).<sup>22</sup> Blastn searches in all domains of Bacteria were  
264 performed using *spf* sequences from 43 selected taxa as query sequences. All complete bacterial  
265 genomes found at the NCBI Genomes resource (<http://www.ncbi.nlm.nih.gov/genome/>) were  
266 checked for the presence of *spf*. More thorough blastn searches were performed in gamma-  
267 proteobacteria, as *spf* were exclusively found in this bacterial class. This was done as follows:  
268 Representative *spf* sequences from all *spf*-containing genera were used as queries in blast  
269 searches. All blast “hits” had a low E-value (i.e., high statistical support; typically below 1e-11).  
270 In other words, *spf* was identified with a high degree of confidence, or, *spf* was not found. In one  
271 case a hit with a poor E-value was found (0.65). Here, we did a manual inspection to decide the  
272 presence/absence of *spf*. First, the NCBI Sequence Viewer  
273 (<http://www.ncbi.nlm.nih.gov/projects/sviewer/>) was used to locate the intergenic region  
274 between *polA* and *engB* (genes that are known to flank *spf*). Next, a manual text search revealed  
275 the presence of highly conserved 5' hairpin, and thereafter the entire *spf*. The  
276 presence/absence of *spf* in all complete genomes from gamma-proteobacteria is provided in  
277 **Table S1**. The presence of *spf* was next mapped on the tree of life, which was produced using the  
278 iTol web tool.<sup>23</sup>

## 279 **Alignments and nucleotide diversity**

280 The sequences from the Rfam list and the newly discovered sequences of *spf* were automatically  
281 aligned and manually examined using Jalview.<sup>24</sup> An alignment containing only one version of  
282 each nucleotide variation of *spf* (no redundant *spf* sequences) was used to examine the  
283 variations on nucleotide level between families, genera and species. A consensus *spf* sequence  
284 was made based on the alignment and was mapped onto an *E. coli* secondary structure (**Fig. 2**).<sup>1</sup>  
285 The *spf* alignment in Rfam includes the first 10 nucleotide upstream of the 5' end of *spf*. However,  
286 the promoter region of *spf* was not considered in this work, and was not included in the  
287 alignment. Existing literature on experimentally verified mRNA targets of Spot 42 were used to  
288 map mRNA targets onto the secondary structure of Spot 42 (**Fig. 3**).<sup>4, 14, 15</sup>

## 289 **Funding**

290 This work was supported by internal grants from UiT- The Arctic University of Norway.

291

292

## 293 Supplemental Material

294 Supplemental data for this article can be accessed on the publisher's website.

295

## 296 References

- 297 1. Møller T, Franch T, Udesen C, Gerdes K, Valentin-Hansen P. Spot 42 RNA mediates  
298 discoordinate expression of the *E. coli* galactose operon. *Genes Dev* 2002; 16: 1696-1706.
- 299 2. Polayes DA, Rice PW, Garner MM, Dahlberg JE. Cyclic AMP-cyclic AMP receptor protein as  
300 a repressor of transcription of the *spf* gene of *Escherichia coli*. *J Bacteriol* 1988; 170:  
301 3110-3114.
- 302 3. Joyce CM, Grindley ND. Identification of two genes immediately downstream from the  
303 *polA* gene of *Escherichia coli*. *J Bacteriol* 1982; 152: 1211-1219.
- 304 4. Hansen GA, Ahmad R, Hjerde E, Fenton CG, Willassen NP, Haugen P. Expression profiling  
305 reveals Spot 42 small RNA as a key regulator in the central metabolism of *Aliivibrio*  
306 *salmonicida*. *BMC Genomics* 2012; 13: 37.
- 307 5. Ikemura T, Dahlberg JE. Small ribonucleic acids of *Escherichia coli*. I. Characterization by  
308 polyacrylamide gel electrophoresis and fingerprint analysis. *J Biol Chem* 1973; 248:  
309 5024-5032.
- 310 6. Ikemura T, Dahlberg JE. Small ribonucleic acids of *Escherichia coli*. II. Noncoordinate  
311 accumulation during stringent control. *J Biol Chem* 1973; 248: 5033-5041.
- 312 7. Sahagan BG, Dahlberg JE. A small, unstable RNA molecule of *Escherichia coli*: spot 42  
313 RNA. I. Nucleotide sequence analysis. *J Mol Biol* 1979; 131: 573-592.
- 314 8. Sahagan BG, Dahlberg JE. A small, unstable RNA molecule of *Escherichia coli*: spot 42  
315 RNA. II. Accumulation and distribution. *J Mol Biol* 1979; 131: 593-605.
- 316 9. Rice PW, Dahlberg JE. A gene between *polA* and *glnA* retards growth of *Escherichia coli*  
317 when present in multiple copies: physiological effects of the gene for spot 42 RNA. *J*  
318 *Bacteriol* 1982; 152: 1196-1210.
- 319 10. Hatfull GF, Joyce CM. Deletion of the *spf* (spot 42 RNA) gene of *Escherichia coli*. *J Bacteriol*  
320 1986; 166: 746-750.
- 321 11. Rice PW, Polayes DA, Dahlberg JE. Spot 42 RNA of *Escherichia coli* is not an mRNA. *J*  
322 *Bacteriol* 1987; 169: 3850-3852.
- 323 12. Polayes DA, Rice PW, Dahlberg JE. DNA polymerase I activity in *Escherichia coli* is  
324 influenced by spot 42 RNA. *J Bacteriol* 1988; 170: 2083-2088.
- 325 13. Wang X, Ji SC, Jeon HJ, Lee Y, Lim HM. Two-level inhibition of galK expression by Spot 42:

- 326 Degradation of mRNA mK2 and enhanced transcription termination before the galK  
327 gene. Proc Natl Acad Sci U S A 2015; 112(24): 7581-6.
- 328 14. Beisel CL, Storz G. The base-pairing RNA spot 42 participates in a multioutput  
329 feedforward loop to help enact catabolite repression in *Escherichia coli*. Mol Cell 2011;  
330 41: 286-297.
- 331 15. Beisel CL, Updegrave TB, Janson BJ, Storz G. Multiple factors dictate target selection by  
332 Hfq-binding small RNAs. EMBO J 2012; 31: 1961-74.
- 333 16. Peer A, Margalit H. Evolutionary patterns of *Escherichia coli* small RNAs and their  
334 regulatory interactions. RNA 2014; 20: 994-1003.
- 335 17. Schmidtke C, Abendroth U, Brock J, Serrania J, Becker A, Bonas U. Small RNA sX13: a  
336 multifaceted regulator of virulence in the plant pathogen *Xanthomonas*. PLoS Pathog  
337 2013; 9(9): e1003626.
- 338 18. Ferrara S1, Carloni S, Fulco R, Falcone M, Macchi R, Bertoni G. Post-transcriptional  
339 regulation of the virulence-associated enzyme AlgC by the  $\sigma(22)$  -dependent small RNA  
340 ErsA of *Pseudomonas aeruginosa*. Environ Microbiol 2015; 17(1): 199-214.
- 341 19. del Val C, Rivas E, Torres-Quesada O, Toro N, Jiménez-Zurdo JI. Identification of  
342 differentially expressed small non-coding RNAs in the legume endosymbiont  
343 *Sinorhizobium meliloti* by comparative genomics. Mol Microbiol 2007; 66(5): 1080-91.
- 344 20. Valverde C, Livny J, Schlüter JP, Reinkensmeier J, Becker A, Parisi G. Prediction of  
345 *Sinorhizobium meliloti* sRNA genes and experimental detection in strain 2011. BMC  
346 Genomics 2008; 9: 416.
- 347 21. Torres-Quesada O1, Oruezabal RI, Peregrina A, Jofré E, Lloret J, Rivilla R, Toro N,  
348 Jiménez-Zurdo JI. The *Sinorhizobium meliloti* RNA chaperone Hfq influences central  
349 carbon metabolism and the symbiotic interaction with alfalfa. BMC Microbiol 2010;  
350 10:71.
- 351 22. Burge SW, Daub J, Eberhardt R, Tate J, Barquist L, Nawrocki EP, Eddy SR, Gardner PP,  
352 Bateman A. Rfam 11.0: 10 years of RNA families. Nucleic Acids Res 2013; 41: D226-232.
- 353 23. Letunic I, Bork P. Interactive Tree Of Life v2: online annotation and display of  
354 phylogenetic trees made easy. Nucleic Acids Res 2011; 39: W475-478.
- 355 24. Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. Jalview Version 2-a  
356 multiple sequence alignment editor and analysis workbench. Bioinformatics 2009; 25:  
357 1189-1191
- 358 25. Gao B, Mohan R, Gupta RS. Phylogenomics and protein signatures elucidating the  
359 evolutionary relationships among the Gammaproteobacteria. Int J Syst Evol Microbiol  
360 2009; 59: 234-47
- 361

362 **Figure legends**

363 **Figure 1**

364 The natural distribution of *spf*. *spf* is restricted to five orders of gamma-proteobacteria (shown in  
365 bold letters), four of which share the same closest common ancestor (indicated by an arrow).  
366 The circular phylogenetic tree (made using the iTol web tool) shows all major branches of  
367 Bacteria. The gamma-proteobacteria phylogeny in the right panel is based on Gao et al.<sup>25</sup> Here,  
368 numbers in parentheses indicate the number of complete genomes that contain *spf* (first  
369 number) and the total number of available complete genomes (second number) in each order.  
370 In addition, *spf* is found in 8 Chromatiales draft genomes (asterisk).

371

372 **Figure 2**

373 Secondary structure consensus model of the Spot 42 RNA. The structure model was made by  
374 aligning all known *spf* sequences, and by mapping the consensus sequence onto a secondary  
375 structure model of the *E. coli* Spot 42 (based on Møller et al.<sup>1</sup>). The structure consists of a relatively  
376 long 5' hairpin, a 9 nt long single-stranded region followed by a second hairpin and a rho-  
377 independent terminator. Level of identity is shown using different type of letters in the structure.  
378 Uppercase bold letters indicate 80–100 % identity, uppercase regular letters indicate 60–79%  
379 identity, and lowercase letters indicate <60% identity. Structural segments with family-specific (i.e.,  
380 *Vibrionaceae*, *Aeromonadaceae* and *Shewanellaceae*) variations are shown in separate colored  
381 boxes. Here, circles indicate U or A insertions (compared to the “consensus”). Grey square around a  
382 letter symbolizes aberration from the consensus structure.

383 **Figure 3**

384 Potential base-pairing between the Spot 42 RNA and experimentally verified mRNA targets from  
385 the following genes: (A) *galK*, (B) *pirin*, (C) *fucI*, (D) *xylF* and *sthA*, (E) *gltA* and *srlA* and (F) *nanC*,  
386 (G) *paaK*, *ascF*, *caiA* and *fucP*, (H) *atoD* and *puuE* and (I) *nanT*. **Fig. 3** is based on data from  
387 Møller et al.,<sup>1</sup> Hansen et al.,<sup>4</sup> Beisel and Storz,<sup>14</sup> and Beisel et al.<sup>15</sup>

388 **Figure 4**

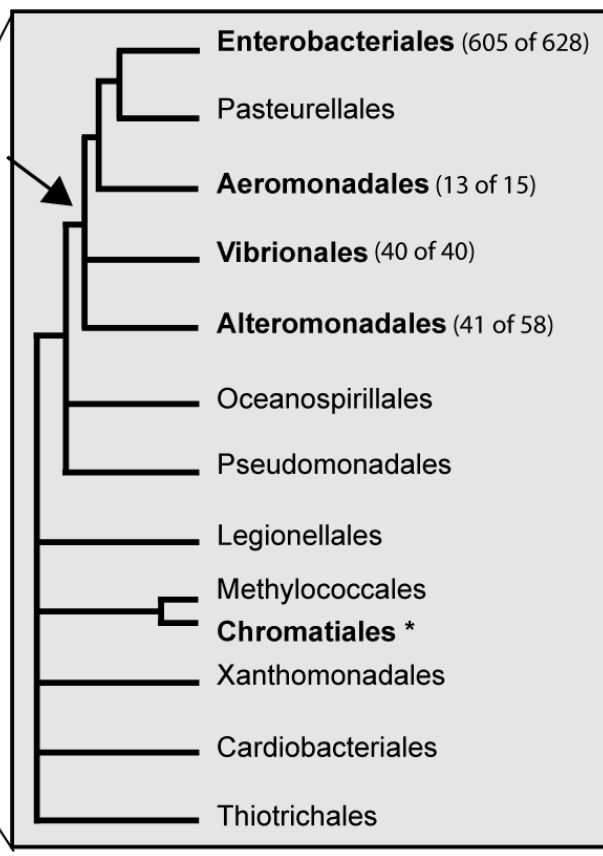
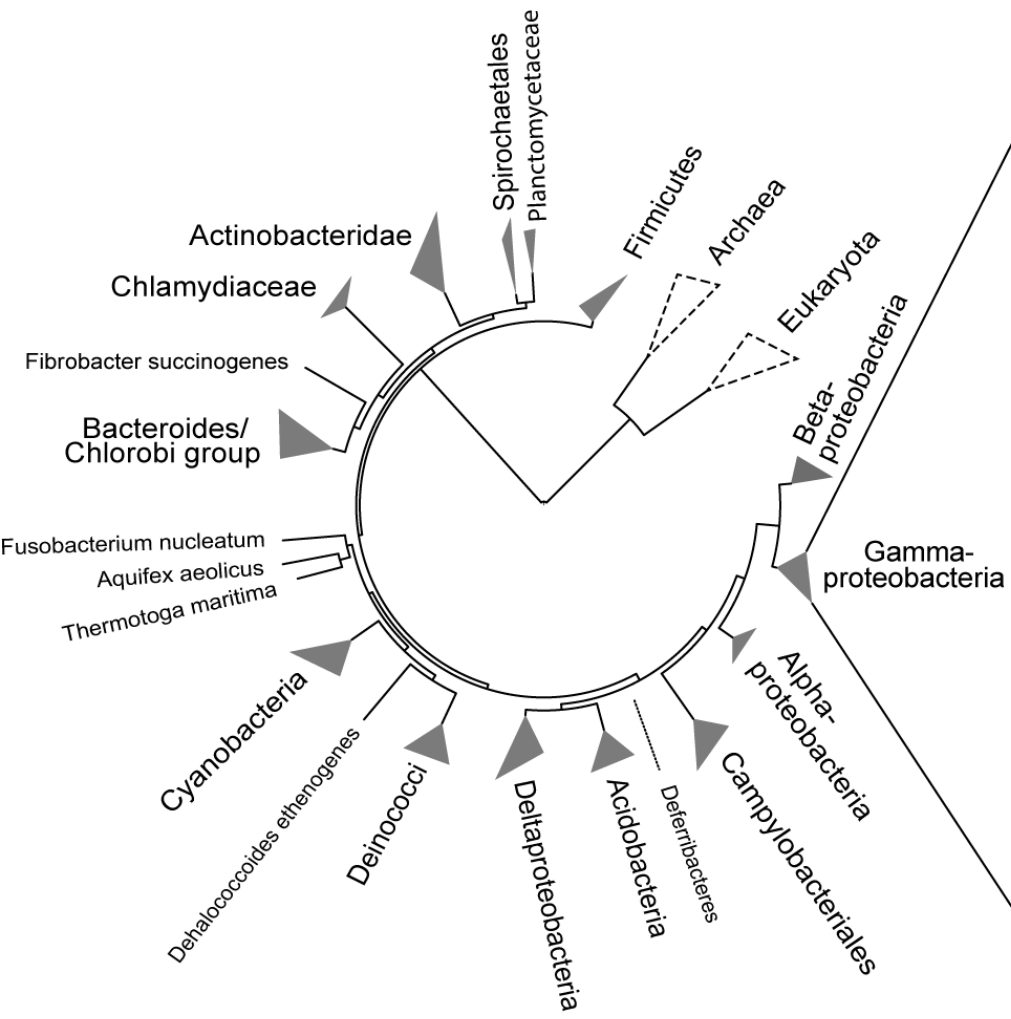
389 sRNA genes in the intergenic region downstream of *polA*. The figure shows currently known  
390 sRNA genes which have been found in the same intergenic region as *spf*. The scale bar shows  
391 distance in nucleotides. (A) Representative species containing *spf* are shown. The *VSsrna24*  
392 sRNA gene is located downstream of *spf* in *V. cholerae* and *A. salmonicida*. Question mark

393 denotes hypothetical protein. (B) Genomic location of the sRNA genes *ersA* in *Pseudomonas*  
394 *aeruginosa*, *sX13* in *Xanthomonas campestris* and *SMc02857* in *Sinorhizobium meliloti*.

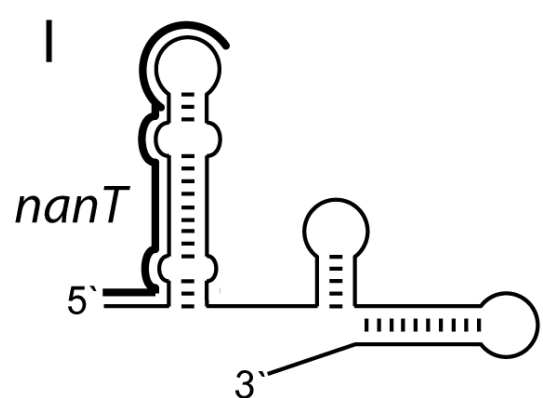
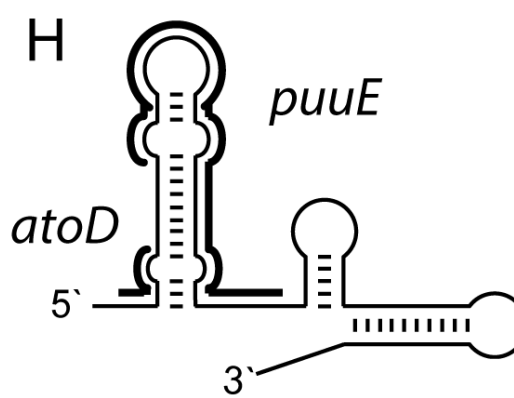
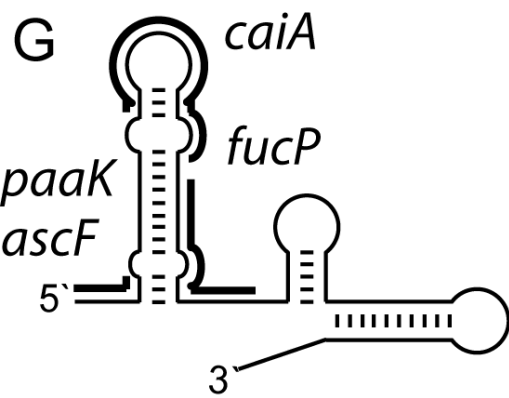
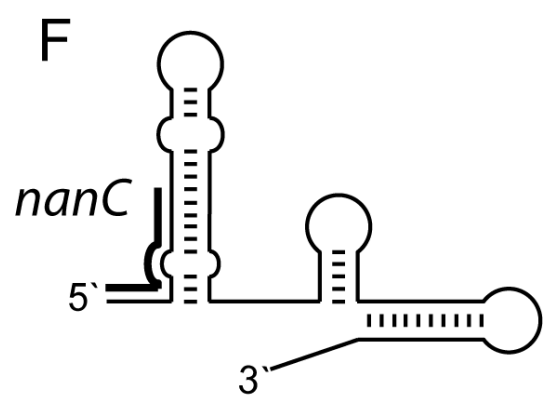
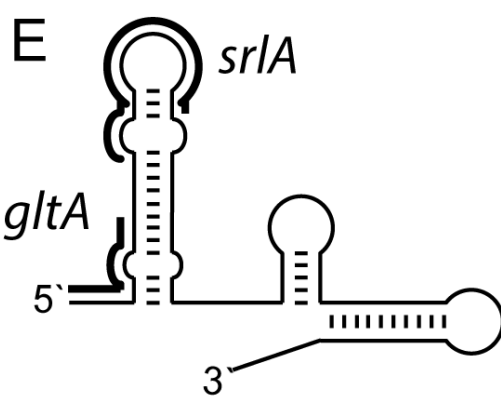
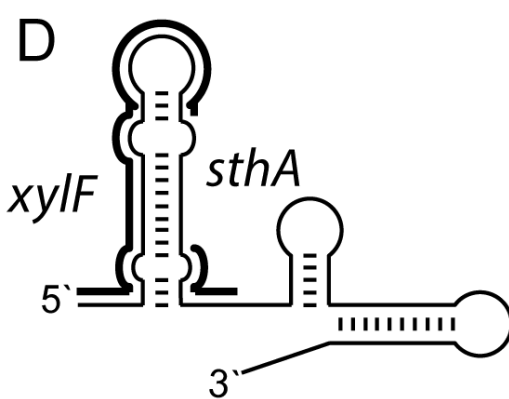
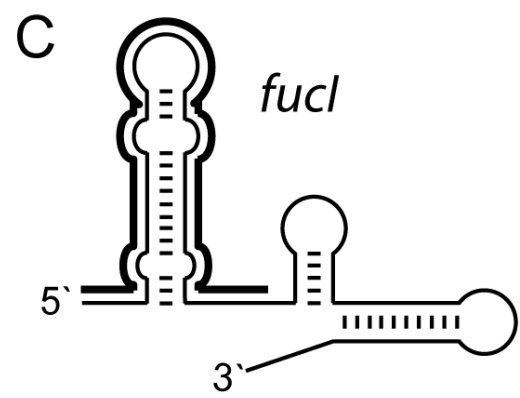
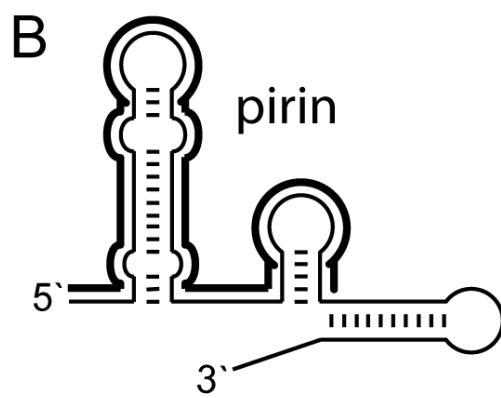
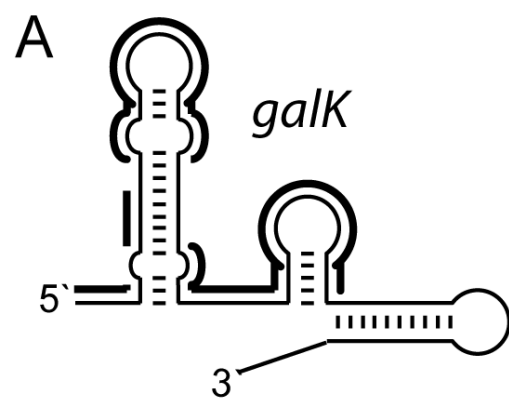
395

396

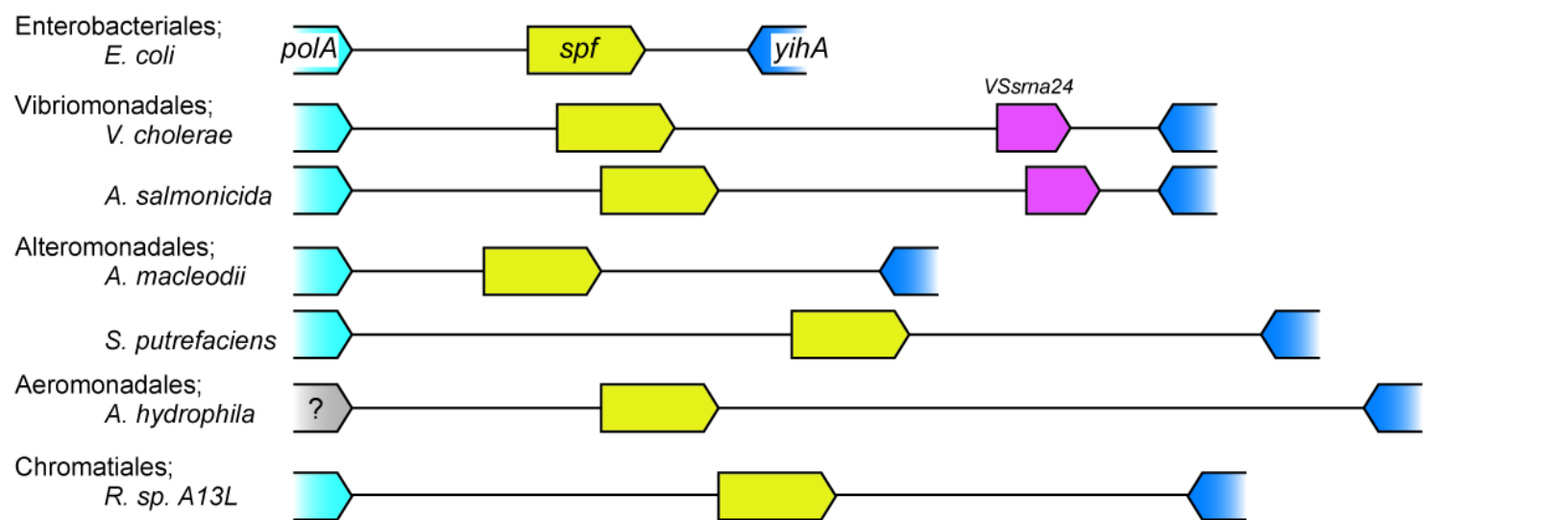
397









**A****Order; species****B**