1

2 **sRNA-target Prediction Organizing Tool (SPOT) integrates computational and**

3 **experimental data to facilitate functional characterization of bacterial small RNAs**

4

5 Alisa M. King[1], Carin K. Vanderpool*[1] and Patrick H. Degnan*[2]

6 [1]Department of Microbiology, University of Illinois, Urbana, IL 61801

7 [2]Department of Microbiology and Plant Pathology, University of California, Riverside, Riverside,

8 CA 92521

9

10 Running Title: Integrated pipeline for small RNA target prediction

11
12 *[1] Corresponding author
13 Carin K. Vanderpool, Ph.D.
14 Department of Microbiology
15 University of Illinois at Urbana-Champaign
16 C226 CLSL, MC-110
17 601 S. Goodwin Ave.
18 Urbana, IL 61801
19 (t) 217-333-7033
20 cvanderp@illinois.edu
21
22 *[2] Corresponding author
23 Patrick H. Degnan, Ph.D.
24 Department of Microbiology and Plant Pathology
25 University of California, Riverside
26 Webber Hall 2317A
27 Riverside, CA 92521
28 (t) 951-827-4415
29 Patrick.Degnan@ucr.edu
30

31

32

33

34

35    **ABSTRACT**

36    Small RNAs (sRNAs) post-transcriptionally regulate mRNA targets, typically under conditions of

37    environmental stress. Although hundreds of sRNAs have been discovered in diverse bacterial

38    genomes, most sRNAs remain uncharacterized, even in model organisms. Identification of

39    mRNA targets directly regulated by sRNAs is rate-limiting for sRNA functional characterization.

40    To address this, we developed a computational pipeline that we named SPOT for sRNA-target

41    Prediction Organizing Tool. SPOT incorporates existing computational tools to search for sRNA

42    binding sites, allows filtering based on experimental data, and organizes the results into a

43    standardized report. SPOT sensitivity (Correctly Predicted Targets/Total Known Targets) was

44    equal to or exceeded any individual method when used on 12 characterized sRNAs.  Using

45    SPOT, we generated a set of target predictions for the sRNA RydC, which was previously

46    shown to positively regulate *cfa* mRNA, encoding cyclopropane fatty acid synthase. SPOT

47    identified *cfa* along with additional putative mRNA targets, which we then tested experimentally.

48    Our results demonstrated that in addition to *cfa* mRNA, RydC also regulates *trpE* and *pheA*

49    mRNAs, which encode aromatic amino acid biosynthesis enzymes. Our results suggest that

50    SPOT can facilitate elucidation of sRNA target regulons to expand our understanding of the

51    many regulatory roles played by bacterial sRNAs.

52

53

54

55

56

57

58

59

60

61  **IMPORTANCE**

62  Small RNAs (sRNAs) regulate gene expression in diverse bacteria by interacting with mRNAs to

63  change their structure, stability or translation. Hundreds of sRNAs have been identified in

64  bacteria, but characterization of their regulatory functions is limited by difficulty with sensitive

65  and accurate identification of mRNA targets. Thus, new robust methods of bacterial sRNA target

66  identification are in demand. Here, we describe our Small RNA-target Prediction Organizing

67  Tool, which streamlines the process of sRNA target prediction by providing a single pipeline that

68  combines available computational prediction tools with customizable results filtering based on

69  experimental data.  SPOT allows the user to rapidly produce a prioritized list of predicted sRNA-

70  target mRNA interactions that serves as a basis for further experimental characterization. This

71  tool will facilitate elucidation of sRNA regulons in bacteria, allowing new discoveries regarding

72  the roles of sRNAs in bacterial stress responses and metabolic regulation.

73

74 **INTRODUCTION**

75 Bacterial small RNAs (sRNAs) range in size from 30 to 300 nucleotides (nts). Regulation of

76 mRNA targets by sRNAs via base-pairing dependent mechanisms alters translation or mRNA

77 stability (1, 2). Most of the time, base-pairing interactions involve the 5' or 3' untranslated region

78 (UTR) of the target mRNA but can also involve sites within the coding region of the target

79 mRNA. Small RNA-dependent translational repression often occurs via interactions that directly

80 interfere with ribosome binding to the mRNA. However, sRNAs have also been shown to

81 activate mRNA targets through various mechanisms, including interference with mRNA decay

82 (3, 4). In recent years it has become evident that sRNAs are ubiquitous and play an important

83 role in mediating and regulating many basic cellular processes and stress responses. Hundreds

84 of small RNAs have been identified in numerous bacterial species such as *Bacillus subtilis* (5),

85 *Listeria monocytogenes* (6), and *Salmonella enterica* (7, 8). With the advancement of current

86 technologies, the number of sRNAs identified in diverse organisms will surely increase.

87 Consequently, there is a pressing need to develop new and better tools for sRNA

88 characterization. In particular, there is a need for methods to address a major rate-limiting step

89 in novel sRNA functional characterization, which is high-fidelity identification of mRNA targets.

90       A variety of computational and experimental methods have been used to predict and

91 validate sRNA-mRNA target interactions. The computational tools currently available for sRNA

92 target prediction, such as TargetRNA (9), sTarPicker (10), IntaRNA (11, 12), and CopraRNA

93 (13), albeit powerful, have their limitations, the most problematic of which is the high rate of

94 false positives. TargetRNA, sTarPicker, and IntaRNA all scan the entire genome and search for

95 putative targets based on interaction hybridization energies. CopraRNA uses the same

96 methodology as IntaRNA for predicting targets based on thermodynamic favorability of the

97 interactions but goes a step further and also considers the conservation of those interactions

98 across species, giving more weight to predictions that are conserved (13). When CopraRNA,

99 IntaRNA, and TargetRNA were used in a side-by-side comparison, CopraRNA was found to

4

100    have the highest positive predictive value (PPV) of 44% and reported the lowest rate of false-

101    positives for known sRNAs across 18 enterobacterial species (14). Although CopraRNA

102    possesses the highest PPV out of all tools, there were still substantial false positives reported.

103    Moreover, CopraRNA is limited to identifying conserved sRNA-target RNA interactions and

104    cannot identify species-specific interactions. As a result, caution should be used with these

105    individual algorithms and they are frequently used in tandem with other target identification

106    methods (14).

107         Experimental methods, including transcriptomic studies, have often been used to identify

108    sRNA candidate targets. Transcriptomics methods uncover gene expression changes caused

109    by absence or overproduction of an sRNA. While microarrays and RNA-sequencing have been

110    successfully used to deduce sRNA targets, in many cases, separating direct effects from

111    indirect effects is laborious and time-consuming. Moreover, the data obtained from

112    transcriptomic studies can only reveal targets that are expressed under the specific growth

113    conditions examined. As such, bona fide target genes that are poorly expressed or that are

114    regulated by mechanisms that do not result in a substantial change in mRNA stability may be

115    missed as sRNA targets. To address these issues, affinity purification methods have been

116    developed to enhance identification of sRNA-mRNA interacting partners. For example, RIL-Seq

117    (RNA interaction by ligation and sequencing) (15) identifies sRNA-mRNA partners that bind to

118    the RNA chaperone Hfq (16) by co-immunoprecipitation, ligation, deep sequencing and analysis

119    of RNA chimeras, which often represent true interacting partners. MAPS (MS2-affinity

120    purification coupled with RNA-sequencing) (17) uses sRNA "bait" that is tagged with an MS2

121    aptamer and can be purified by interaction with the MS2 coat protein. RNA targets that are

122    copurified with the sRNA bait are then identified by deep sequencing. Even with the variety of

123    tools available for sRNA target identification, it is still not entirely clear which tools are the most

124    effective for sRNA target identification.

125    In order to streamline the use of multiple existing sRNA prediction algorithms, we

126    developed a software pipeline called SPOT (sRNA-target Prediction Organizing Tool) that uses

127    several algorithms in parallel to search for sRNA-mRNA interactions. The software collates

128    predictions and allows integration of experimental data using customizable results filters. First,

129    we used two well-characterized *E. coli* sRNAs, SgrS (18) and RyhB (19, 20), to assess the

130    effectiveness of SPOT as the targets of these sRNAs are well defined. Next, we extended the

131    application of the SPOT pipeline to UTRs of mRNAs to identify potential sRNAs involved in

132    regulation. We then applied the same parameters and analyses to a less characterized *E. coli*

133    sRNA, RydC. Employing a combinatorial approach through SPOT predictions and experimental

134    validation, we were able to identify two new RydC targets, *pheA* and *trpE*, which were

135    downregulated and upregulated, respectively, by RydC.

136    **MATERIALS AND METHODS**

137    **Software pipeline**

138    A software pipeline was constructed in PERL to provide a single interface for running four

139    sRNA-mRNA target prediction algorithms in parallel and collating their results (Fig. 1). Source

140    codes for TargetRNA2 v2.01 (9), sTarPicker (10), IntaRNA v1.0.4 (12), and CopraRNA v 1.2.9

141    (13) were downloaded and installed on a multicore local server. The pipeline is comprised of 4

142    steps described briefly here.

143    1. Reference genome files are retrieved from RefSeq or local customized genome files can

144    be used, provided they are in an appropriate RefSeq format (GBK file or PTT and FNA files).

145    2. Simultaneous searches are initiated for TargetRNA2, sTarPicker, and IntaRNA according

146    to user defined search parameters (e.g., window size, seed size, significance cutoffs).

147    Optionally, if RefSeq IDs and corresponding sRNA sequences from related genomes are

148    provided, a CopraRNA search is initiated.

149    3. The pipeline tracks the progress of each job and once each search is completed the raw

150    results files are read into memory.

151    4. User-defined results filtering parameters are applied (e.g., list with known binding

152    coordinates, differential expression, operon data) and the raw results in memory are collated

153    into a unified report.

154    The collated results report includes Excel-formatted data tables, functional enrichment

155    predictions for consensus mRNA targets as well as binding plots. Both the collated results and

156    individual search results can be downloaded once the job is complete. In addition, users can

157    elect to have an email notification sent when the job is complete. The pipeline also includes an

158    option to re-run the results collation steps using different results filters. This enables users to

159    make minor adjustments to the results reporting without waiting for the individual searches to be

160    re-run.

161    The SPOT program and installation instructions are available on GitHub

162    (https://github.com/phdegnan/SPOT). In addition, an Amazon Web Service (AWS) cloud

163    Amazon Machine Image (AMI) with all of the required software installed is available (search for

164    SPOTv1). The SPOT User Manual is also included in Supplementary Material.

165    **Generation of test data sets**

166    Known sRNA-mRNA interactions were collected from ecocyc.org (21), the literature, and

167    experiments herein for 12 sRNAs with ≥4 confirmed targets: RyhB (b4451, RF00057), Spot42

168    (spf, b3864, RF00021), SgrS (b4577, RF00534), RybB (b4417, RF00110), FnrS (b4699,

169    RF01796), GcvB (b4443, RF00022), OmrA (b4444, RF00079), CyaR (b4438, RF00112), MicA

170    (b4442, RF00078), MicF (b4439, RF00033), DicF (b1574, RF00039), and RydC (b4597,

171    RF00505) (Table S1). The confirmed sRNA-mRNA binding interactions were used as true

172    positives, to investigate the reliability and sensitivity of the pipeline.

173    In order to test CopraRNA, homologs for the 12 *E. coli* MG1655 sRNAs were identified in

174    related genomes using Infernal (22). For all sRNAs excluding DicF, the genomes of *E.*

175    *fergusonii* ATCC 35469 (NC_011740), *Citrobacter koseri* ATCC BAA-895 (NC_009792), and

176    *Salmonella enterica* sv. Typhimurium LT2 (NC_003197) were queried with the Infernal algorithm

177    and each covariance model. For the sRNA DicF, a phylogenetically restricted sRNA, *E. coli*

178    O157:H7 str. Sakai (NC_002695) and *E. coli* str. APEC O1 (NC_008563) were queried. In

179    cases where genomes encoded ≥1 prediction (e.g., OmrA), the prediction with the lowest E

180    value was used.

181         In addition, we compiled a list of 85 *E. coli* sRNAs to investigate the ability of the pipeline

182    to be used to predict mRNA-sRNA interactions using a putative mRNA target as the search

183    query (Table S2). This includes 65 RefSeq annotated sRNAs (NC_000913.3), an additional 19

184    sRNAs annotated in ecocyc.org (21), and the sRNA IepX (23). Note that 552 additional

185    predicted *E. coli* sRNAs, *cis* regulatory elements and other putative RNAs corresponding to

186    known RFAMs (n=172) or identified from expression studies (n=360) were not included (24, 25).

187         Finally, sRNA-mRNA interaction coordinates and the 5' UTRs of 11 mRNAs with ≥2

188    known interacting sRNAs were collected from ecocyc.org (21): *csgD* (b1040,n=5), *flhD*

189    (b1892,n=4), *ompA* (b0957,n=3), *ompC* (b2215,n=3), *ompF* (b0929,n=2), *ompX* (b0814,n=2),

190    *phoP* (b1130,n=2), *rpoS* (b2741,n=4), *sdhC* (b0721,n=3), *sodB* (b1656,n=2), and *tsx*

191    (b0411,n=2).

192    **Media and reagents**

193    *E. coli* strains were cultured in lysogeny broth (LB) medium or on LB agar plates at 37°C, unless

194    stated otherwise. For construction of reporter fusions by λ Red, recovery of recombinants was

195    carried out on M63 minimal medium containing 5% sucrose, 0.001% L-arabinose (Ara), 0.2%

196    glycerol, and 40 μg/ml 5-Bromo-4-Chloro-3-Indolyl β-D-Galactopyranoside (X-Gal). For β-

197    galactosidase assays, bacterial cells were grown in Tryptone Broth (TB) medium supplemented

198    with 100 μg/ml ampicillin (Amp) overnight at 37°C and then subcultured in TB broth containing

199    100 μg/ml ampicillin (Amp) with 0.002% L-arabinose. Where necessary, media were

200    supplemented with antibiotics at the following concentrations: 100 μg/ml ampicillin (Amp), 25

201    μg/ml chloramphenicol (Cm), and 25 μg/ml kanamycin (Kan). Expression of RydC was induced

202    with either 0.1 or 0.5 mM Isopropyl β-D-1-thiogalactopyranoside (IPTG) from the PLlacO-1

203    promoter.

**Strain construction**

205    Strains and plasmids used in this study are listed in Table S3. All strains used in this study are

206    derivatives of *E. coli* K12 strain MG1655. Oligonucleotide primers and 5'-biotinylated probes

207    used in this study are listed in Table S4 and were all acquired from Integrated DNA

208    Technologies (IDT). Chromosomal mutations were made by λ Red recombination (26, 27), and

209    marked alleles were moved between strains by P1 *vir* transduction (28). PCR products were

210    generated using the Expand™ High Fidelity PCR System (Sigma-Aldrich, St. Louis, MO)

211    according to the manufacturer's instructions. All mutations were verified by amplifying PCR

212    fragments using GoTaq polymerase (Promega, Madison, WI) and sequencing.

213         The translational *lacZ* reporter fusions under the control of the PBAD promoter were

214    constructed by PCR amplifying a fragment of interest using forward and reverse primers

215    containing 5' homologies to PBAD and *lacZ* (Table S3). PCR products were recombined into

216    PM1205 using λ Red homologous recombination and counter-selection against *sacB* as

217    described previously (29). The fusions used in this study were inserted into the *lac* locus of

218    PM1205. Some *lacZ* reporter fusions used in this study were constructed using the one-step

219    recombination method (30).

220         Plasmids harboring mutated *rydC* alleles under the control of the PLlacO-1 promoter

221    were constructed using the Quickchange II XL Site-Directed Mutagenesis Kit (Agilent

222    Technologies, Santa Clara, CA) with oligonucleotides AKP59 (P$_{LlacO-1}$-*rydC3*), AKP68 (P$_{LlacO-1}$-

223    *rydC5*), and AKP69 (P$_{LlacO-1}$-*rydC345*) that contained mismatched bases at the desired locations

224    and transformed into XL10-Gold Ultracompetent cells (Table S3).

**RNA-seq analysis**

226    *E. coli* K12 MG1655 strain AK250 (Δ*rydC, lacI*$^{q+}$) harboring vector (pBR322) or P$_{lac}$-*rydC*

227    plasmid was grown to OD$_{600}$~0.5 in LB broth media at 37°C and then induced with 0.1 mM IPTG

9

228 for 10 min. The hot phenol method (31) was used to extract total RNA after 2 and 10 minutes of

229 induction. Samples were then treated with TURBO™DNase (Ambion) kit according to the

230 manufacturer's protocol and resolved by gel electrophoresis on 1.2 % agarose gel to confirm

231 integrity of the 16S and 23S bands. Ribosomal RNA removal, library construction and

232 sequencing were performed at the W. M. Keck Center for Comparative and Functional

233 Genomics at the University of Illinois at Urbana-Champaign. Ribosomal RNA was removed from

234 1 µg of total RNA using Ribozero rRNA Removal Meta-Bacteria Kit (Illumina, Inc), and the

235 mRNA-enriched fraction was converted to indexed RNA-seq libraries (single reads) with the

236 TruSeq Stranded RNA Sample Preparation Kit (Illumina, Inc). The prepared libraries were then

237 pooled in equi-molar concentrations and were quantified by qPCR with the Library

238 Quantification kit Illumina compatible (Kapa Biosystems) and sequenced for 101 cycles plus

239 seven cycles for the index read on a HiSeq2000 using TruSeq SBS version 3 reagents. The

240 output fastq files were generated using Casava 1.8.2 (Illumina) and analyzed with Rockhopper

241 (32). Genes were considered differentially expressed in RydC pulse-expression strains if they

242 met a significance cutoff (q-value) of $\geq 0.005$ and a fold-change value of >1.5 or <0.5. Some

243 genes outside this range were studied because they met other criteria (e.g., prediction of a

244 RydC-mRNA interaction by multiple algorithms).

245 **β-galactosidase assays**

246 Bacterial strains were cultured overnight at 37ºC (shaking) in TB medium containing 100 µg/ml

247 Amp. Subsequent to overnight growth, cultures were diluted 1:100 into fresh TB media

248 containing 100 µg/ml Amp and 0.002% Ara and cultured at 37 ºC. After reaching an $OD_{600}$ of

249 0.3, 0.1 or 0.5 mM IPTG was added to induce expression of the plasmids and grown for an

250 additional hour until an $OD_{600}$ of 0.5 - 0.6 was reached. All β-galactosidase assays were

251 performed as described in previous protocols (33). In short, the samples were suspended in Z-

252 buffer, with reactions conducted at 28ºC with 4 mg/ml 2-nitrophenyl β-D-galactopyranoside

253 (ONPG) as a substrate and 1 M $Na_2CO_3$ to end the reactions.

254 **RESULTS**

255 **Integrated pipeline for sRNA target prediction algorithms**

256 A number of algorithms and tools for identifying putative sRNA-mRNA interactions have been

257 developed (9, 10, 12, 13). However, no single target prediction tool is 100% accurate, the tools

258 implement distinct user-defined parameters, each tool uses a different format for reporting

259 results, and tools are hosted on distinct web platforms. Our approach was to create a single

260 pipeline incorporating existing computational tools to search for sRNA binding sites, producing a

261 collated and standardized results report (Fig. 1). We incorporated the TargetRNA2 (9),

262 sTarPicker (10), IntaRNA (12), and CopraRNA (13) tools into this pipeline because they are

263 widely used and have open source code. Input for the pipeline minimally includes a fasta

264 sequence for the sRNA and the RefSeq number for the target genome. Additional RefSeq

265 genome IDs and homologous sRNA sequences can be provided if the user wishes to include

266 CopraRNA results in the analysis. The pipeline interface also allows the user to define a set of

267 parameters for the individual algorithms and results filters. In particular, the results can be

268 filtered for genes with known binding sites or sets of genes that were identified as putative

269 targets by experimental methods (e.g., RNA-seq, MAPS [MS2 affinity purification coupled with

270 RNA sequencing] (17), RIL-seq [RNA interaction by ligation and sequencing] (15)). For

271 instance, output from the RNA-seq analysis tool Rockhopper (32) can be used directly as a

272 results filter. The program then follows four basic steps (1) download/validate input files, (2)

273 simultaneously initiate computational tools, (3) track job progress and read individual raw

274 results, (4) filter and collate results into a single report (Fig. 1). Finally, an option is provided that

275 allows users to re-collate the results from an initial analysis using different results filter settings.

276 **Pipeline Optimization with SgrS and RyhB targets**

277 SgrS and RyhB are two well-characterized model sRNAs in *E. coli* critical for glucose-phosphate

278 (34) and iron limitation (35) stress responses, respectively. Numerous studies have confirmed 8

279 mRNA targets of SgrS (18, 36, 37, 38, 39) and 18 of RyhB (20, 40, 41, 42). We used these two

280    sRNAs to test the utility and sensitivity of the pipeline. For RyhB, the entire 90-nt sequence was

281    used as query for the bioinformatics search. For SgrS, only the 3' 80-nts of the 227-nt sRNA

282    was used as query, since this is the region involved in target RNA binding. Our initial

283    optimization of the pipeline focused primarily on three parameters: 'seed size', 'window size' and

284    'significance cutoffs.' Each application utilizes distinct defaults for these parameters. For

285    example, 'seed size,' defined as the number of contiguous base pairing interactions required to

286    define an sRNA-mRNA match is set to a default value of 7 in TargetRNA2 and IntaRNA and 5 in

287    sTarPicker. We varied the seed sizes for each algorithm and determined how different seed

288    sizes impact the sensitivity of detection of true targets for SgrS and RyhB. Sensitivity is defined

289    as Correctly Predicted Targets/Total Known Targets (i.e., true positive rate). For TargetRNA2, a

290    seed size of 7 gave the highest sensitivity for correct target predictions, with 38% and 56%

291    correct predictions, for SgrS and RyhB, respectively (Fig. 2A). For sTarPicker, the seed size

292    giving the optimal sensitivity was 6, with 63% and 72% of known binding interactions identified

293    for SgrS and RyhB, respectively. IntaRNA yielded the highest sensitivity of all three algorithms,

294    again at a seed size of 6. IntaRNA correctly identified 100% of known SgrS interactions and

295    94% of known RyhB interactions (Fig. 2A). Based on these results, we used seed size settings

296    of 7 for TargetRNA2 and 6 for IntaRNA and sTarPicker for all other analyses.

297         Next, we evaluated how altering the window size and significance cutoffs impacted the

298    accuracy of predictions (Fig. 2B, C). The window size refers to the size of the region upstream

299    and downstream of every start codon in the genome that is searched for potential base pairing

300    with the query sRNA. Default window sizes for each tool vary dramatically. The default

301    TargetRNA2 window size is 80 nt upstream and 20 nt downstream (80/20) of each start codon

302    (9). The default for sTarPicker is 150/100, IntaRNA suggests 75/75 and CopraRNA uses

303    200/100. Likewise, the tools have different metrics to determine the significance of a match

304    either providing a $P$ value (TargetRNA2, IntaRNA, CopraRNA) or a probability measure

305    (STarPicker). TargetRNA2 generates $P$ values for predicted interactions based on the sRNA-

306    mRNA hybridization energy scores of a randomized mRNA pool (32). IntaRNA utilizes *P* values

307    based on transformation of the energy scores calculated for all putative target binding sites with

308    energy score ≤0 (14). CopraRNA combines individual IntaRNA *P* value predictions among

309    clusters of genes to generate a weighted *P* value and false discovery rate (FDR)-corrected *Q*

310    value (14). In contrast, sTarPicker uses a machine learning approach to generate probabilities

311    as a proportion of base classifiers (n=1000) that support each proposed interaction (10). The

312    sTarPicker authors report that probabilities ≥ 0.5 correspond to likely sRNA-mRNA interactions.

313    SPOT provides the user with the ability to alter the search window and significance thresholds

314    used by all the algorithms included in the pipeline (Fig. 1). We chose two sets of parameters

315    that we define as "Stringent" and "Relaxed," and tested the performance of each set of

316    parameters in correctly identifying known RyhB and SgrS target binding sites (Fig. 2B, C).

317    Stringent parameters incorporated a window size of 80-nt upstream and 20-nt downstream

318    (80/20) of start codons as the search region, and significance thresholds of 0.05 for

319    TargetRNA2, 0.5 for sTarPicker, "top" (e.g., the top 100 predictions) for IntaRNA, and 0.01 for

320    CopraRNA (Fig. 2B, C). Relaxed parameters used a comparatively larger window size of

321    150/100 and thresholds of 0.5, 0.001, "un," (e.g., all predictions) and 0.01 for TargetRNA2,

322    sTarPicker, IntaRNA, and CopraRNA, respectively.

323        Using stringent search parameters, 10/18 known RyhB target binding sites and 2/8

324    known SgrS target binding sites were correctly predicted by ≥2 algorithms (Figure 2B, C,

325    indicated by 2 or more pink cells and absence of blue cells). Using relaxed parameters, the

326    correctly predicted interactions rose to 17/18 and 6/8 for RyhB and SgrS, respectively. Thus, for

327    both RyhB and SgrS, relaxed parameters substantially increased the number of correctly

328    identified binding sites (Fig. 2B, C). Notably, use of relaxed parameters was necessary to

329    capture true binding sites like the SgrS binding site on *yigL* mRNA, which is located further from

330    the start codon than is typical. The relaxed parameters improve the sensitivity of individual

331    methods but may result in the downside of identifying more false positives. IntaRNA has high

13

332 sensitivity for true positives (correct identification of known sRNA binding sites) under the

333 relaxed settings, but also gives a high rate of likely false positives, illustrated by the fact that

334 IntaRNA predicts >3400 binding interactions that are not predicted by any other algorithm.

335 Mitigating this downside of using relaxed parameters, we saw that in the majority of instances

336 the correct RyhB and SgrS binding sites were predicted by ≥2 methods and incorrect

337 predictions by ≥2 methods occurred rarely (RyhB= 1/18, SgrS= 0/8) (Fig. 2B, C).

338  For SgrS and RyhB, at least a dozen mRNAs have been experimentally defined as 'non-

339 targets' for each sRNA (18). In other words, predicted sRNA-mRNA interactions were tested

340 and shown not to mediate regulation of the mRNA in question. These examples served as

341 controls that allowed us to calculate False Positive Rates. Together with the Sensitivity

342 measures for each algorithm and the pipeline, we generated receiver operating characteristic

343 (ROC) curves to assess the accuracy of the methods alone and in combination (Fig. 2D, E).

344 Ideally tools should yield high true positive rates and low false positive rates, resulting in values

345 falling in the upper left quadrant of the ROC curve. Our results indicate that when 2 methods

346 converge on the same prediction, the pipeline achieves ≥75% sensitivity and ≤ 50% false

347 positive rate for both sRNAs. This is a marked improvement in most instances over the single

348 algorithms used here (Fig. 2D, E). In particular, using a 2-method threshold mitigates the very

349 high false positive rate from IntaRNA. We note that making the IntaRNA *P* value cutoff more

350 stringent (e.g., 0.05) decreases the false positive rate dramatically, but at a cost to sensitivity

351 (Fig. S1). Similarly, requiring 3 or 4 algorithms to identify the same predicted interaction

352 decreases the false positive rate of predictions for RyhB and SgrS, however, the sensitivity

353 decreases by more than 25% (Fig. 2D, E). Collectively, these analyses suggest that use of

354 relaxed search parameters and a combined evidence approach requiring a minimum of 2

355 algorithms to predict the same binding interaction is an effective means of improving sRNA

356 target prediction sensitivity.

357     The SPOT pipeline accepts several results filters to facilitate analysis of the predictions.

358     First, users can provide the program a list of binding site locations for known mRNA targets

359     (e.g., true positives). Second, users can include genes on the list that lack known binding sites

360     in order to limit the results reporting to select genes of interest, for example those that emerged

361     from experimental analyses (e.g., RNA-seq). Integration of experimental data with

362     computational predictions is another valuable way of reducing potential false positive

363     predictions.

364     Based on our results and observations during the optimization of SgrS and RyhB target

365     identification, we designed SPOT to prioritize the target binding site predictions (Fig. S2). First,

366     known binding sites correctly predicted by ≥2 algorithms (1) or 1 algorithm (2) are reported. Any

367     gene targets with predictions that are discordant with known binding sites (3) are reported next.

368     Then any additional targets with the same predicted target site found by ≥2 algorithms are

369     ranked next (4). This is followed by targets that were only predicted by a single algorithm, in the

370     following order: CopraRNA (5), TargetRNA, sTarPicker (6), and IntaRNA (7). Using the results

371     filters, a user can narrow or widen their searches, for example, by limiting the predictions made

372     by single algorithms or by applying secondary filters on binding site regions.

373     **Application of SPOT to additional sRNAs**

374     To evaluate the robustness of the defined pipeline parameters and our ranking methods, we ran

375     similar analyses on 9 additional sRNAs with ≥4 known targets. Overall, we found that the SPOT

376     pipeline sensitivity (e.g., the percent of correctly identified interactions) was equal to or

377     exceeded any individual method (average = 84% ± 8.5%, Fig. 3A, Fig. S3A). As before, we

378     found that correct identification by ≥2 methods occurred in the majority of instances (Fig. 3A, red

379     bars). The full list of target predictions generated by ≥2 methods for all 11 sRNAs (Fig. S3B) are

380     included as Supplemental Dataset 1. On average the primary analysis by the pipeline took 1hr

381     15min ± 35min, using as many as 6 processing cores simultaneously. Re-collation of the results

382     using different filters only took an average of 29s ± 6s.

**Extended application of the SPOT pipeline – mRNA as query sequence**

The four individual algorithms are intended to identify the interaction of an sRNA with mRNA targets. However, a user may be interested in determining which known sRNAs interact with a specific mRNA of interest. Normally this would require running an individual search for each of the 10s to 100s of sRNAs from that organism. As part of our pipeline we have designed a feature that allows a user to input a custom annotation file for their reference genome. Therefore, instead of providing the list of mRNA targets, sRNAs can be provided to the algorithm and the relevant mRNA sequence, e.g., a 5' untranslated region (UTR) of interest can be used as the query. We carried out this "reverse" analysis on 11 *E. coli* 5' UTRs that have already been demonstrated to interact with ≥2 different sRNAs. The results are comparable to the analysis using sRNAs as targets – known sRNA interactions were identified with an average sensitivity of 85% ± 24% (Fig. 3B). Moreover, using the 2-algorithm cutoff we were able to use this approach to predict 5 to 14 additional sRNAs that putatively bind the UTRs and could affect their regulation (Supplemental Dataset 1). We note that due to technical constraints the reverse search method can only be used with TargetRNA2, sTarPicker, and IntaRNA at this time. This approach is a novel feature that will facilitate ongoing sRNA research.

**Examination of novel RydC target predictions**

We next sought to use the SPOT pipeline to identify additional targets for the poorly characterized sRNA RydC. RydC was reported to repress *yejA* mRNA (encoding an uncharacterized ABC transporter (43)) and *csgD* mRNA (encoding the master regulator of curli biogenesis (44)), but the molecular mechanisms of RydC-mediated repression were not reported. Fröhlich, *et al.*, (2013) demonstrated that RydC activates *cfa* mRNA, encoding cyclopropane fatty acid synthase. This activation involves RydC-dependent protection of *cfa* mRNA from RNase E-mediated degradation (3). Despite identification of these targets, the physiological function of RydC remains unclear. We used SPOT to identify additional targets of RydC as a means to gain further insight into its physiological role in *E. coli*.

16

409     Our strategy for RydC target identification was to combine computational and

410     experimental data to generate an experimentally tractable list of putative targets for further

411     validation. Experimental identification of putative targets was accomplished by pulse expression

412     of RydC from an inducible promoter followed by identification of RydC-dependent changes in

413     gene expression by RNA-seq. Vector control and $P_{lac}$-*rydC* plasmids were maintained in a

414     Δ*rydC* host strain grown in rich medium (LB) at 37°C. Expression of *rydC* was induced by

415     addition of IPTG to cultures, and total RNA harvested at 10 minutes after induction. RNA-seq

416     data output fastq files were analyzed with Rockhopper and exported as .xls files (Supplemental

417     Dataset 2).

418     To identify putative RydC targets, the SPOT pipeline was applied to RydC using both

419     stringent and relaxed parameters, with the former being more restrictive for window size and

420     algorithm thresholds as described above. Similar to analyses for SgrS and RyhB, the relaxed

421     parameters yielded a greater number of predictions than the stringent parameters. Potential

422     targets that were predicted by ≥ 3 algorithms with the relaxed parameters are shown in Fig. 4

423     above the bold line. The RydC binding site for a validated target, *cfa* mRNA, was correctly

424     predicted by 3 algorithms in the relaxed run. TargetRNA2 predicted a binding site that was

425     inconsistent with the known binding site. The *cfa* prediction was absent in the stringent run,

426     since the base pairing interaction between RydC and *cfa* mRNA takes place outside the window

427     specified in the stringent run (Fig. 4). Some of the putative targets predicted by ≥ 3 algorithms

428     were also differentially expressed in RydC pulse expression RNA-seq experiments (indicated

429     under "Fold Change," Fig. 4). Another set of genes were predicted as targets by ≥ 2 algorithms,

430     and differentially expressed in RNA-seq experiments (Fig. 4, see targets below the bold line).

431     Genes chosen for further analysis are listed in Table 1, along with information about their

432     functions, differential expression in RNA-seq, predicted binding interactions, and algorithm

433     predictions. Several other genes that did not meet the criteria for inclusion in Fig. 4 were also

434     chosen for analysis because they had been described previously as RydC targets or because

17

435    they encode proteins belonging to functional categories related to known RydC targets (Table

436    1).

**Testing pipeline predictions for RydC**

438    To test the targets selected for further validation for regulation by RydC, we constructed

439    translational fusions to putative targets. These fusions were placed under the control of an

440    arabinose-inducible promoter ($P_{BAD}$) to eliminate possible indirect transcriptional effects. For

441    each target, the entire 5' UTR and part of the coding sequence (length variable, depending on

442    the location of the predicted RydC binding site) was fused to *'lacZ* (Fig. 5A). Strains containing

443    the reporter fusions were transformed with vector control or $P_{lac}$-*rydC* plasmids and reporter

444    activity measured after induction with IPTG.

445          In *Salmonella,* RydC was demonstrated to activate *cfa* translation by occluding an

446    RNase E cleavage site to stabilize the *cfa* mRNA (3). Conservation of RydC-*cfa* mRNA

447    interactions between *E. coli* and *Salmonella* as well as SPOT identification of *cfa* as a putative

448    RydC target (Fig. 4, Table 1) suggest that *E. coli* RydC regulates *cfa* in a similar manner. To

449    confirm this, we constructed two translational fusions: $P_{BAD}$-*cfa'-'lacZ*-Long, which contains the

450    RydC binding site, and $P_{BAD}$-*cfa'-'lacZ*-Short, which lacks the RydC binding site (Fig. 5B). RydC

451    production strongly activated the long fusion, increasing activity by >20-fold compared to the

452    vector control strain (Fig. 5B). As expected, activity of the short fusion lacking the RydC binding

453    site was unaffected upon RydC induction (Fig. 5B). These results support the model that *cfa*

454    mRNA is a directly regulated by RydC in both *S. enterica* and *E. coli*.

455          Strains harboring reporter fusions to 13 other putative targets (listed in Table 1) were

456    transformed with vector control and $P_{lac}$-*rydC* plasmids and β-galactosidase assays were

457    performed after a period of RydC induction (Fig. 5C). Only two of the target fusions were

458    differentially regulated by the criteria we selected (≥1.5-fold or ≤0.5-fold) in RydC-expressing

459    cells compared to the vector control (Fig. 5C). These two targets were *pheA* and *trpE*, which

460    both encode proteins involved in aromatic amino acid biosynthesis. Previous studies (43, 44)

461    reported RydC-dependent translational repression of the *yejA* and *csgD* mRNAs*,* though we

462    note that specific and direct base pairing interactions with RydC were not demonstrated. Our

463    translational fusions to these putative targets did not show any differential regulation in

464    response to RydC expression (Fig. 5C).

465     **RydC regulates genes in aromatic amino acid biosynthetic pathways**

466    In RNA-Seq experiments, levels of *pheA* mRNA were reduced to ~30% of control levels when

467    RydC was ectopically expressed (Supplemental Dataset 2). Likewise, in RydC-producing cells,

468    activity of the P$_{BAD}$-*pheA'-'lacZ* fusion was ~30% that of the vector control (Fig. 5C). The

469    predicted RydC-*pheA* mRNA base pairing interaction involves the 5' end of RydC and the

470    coding region of *pheA,* directly adjacent to the start codon (Fig. 6A). The P$_{BAD}$-*pheA'-'lacZ* fusion

471    encompasses all of the 5' UTR and 645-nt of the coding region. A reporter derived from this has

472    mutations that disrupt the predicted base pairing with RydC, resulting in the P$_{BAD}$-*pheA67'-'lacZ*

473    fusion with mutations G9C/G10C (Fig. 6A). A *rydC* allele with compensatory mutations

474    (C4G/C5G) was constructed and named RydC5. The mutations in RydC5 abrogated regulation

475    of the wild-type P$_{BAD}$-*pheA'-'lacZ* fusion. Likewise, the mutations in P$_{BAD}$-*pheA67'-'lacZ*

476    prevented regulation by wild-type RydC. The compensatory mutant pair: P$_{BAD}$-*pheA67'-'lacZ* and

477    RydC5 had restored regulation, albeit not to fully wild-type levels. Together, these data suggest

478    that RydC targets *pheA* mRNA for translational repression. Due to the location of the base

479    pairing interaction in the translation initiation region, mechanism is likely direct occlusion of

480    ribosome binding to *pheA* mRNA by RydC.

481           Another new putative RydC target is *trpE,* which encodes a component of the

482    anthranilate synthase involved in tryptophan biosynthesis. A P$_{BAD}$-*trpE'-'lacZ* fusion

483    encompassing the 30-nt *trpE* mRNA 5' UTR and 42-nt of *trpE* coding sequence was activated

484    upon RydC production by slightly less than 2-fold (Fig. 5C, 7B). The predicted RydC-*trpE* mRNA

485    base pairing interaction involves sequences near the 3' end of RydC and sequences within the

486    *trpE* coding sequence. Point mutations in the *trpE* reporter fusion (C20G/C22G) resulted in the

19

487 mutant reporter P$_{BAD}$-*trpE20'-'lacZ,* which was not substantially upregulated when wild-type

488 RydC was produced (Fig. 7B). Because of the unusual pseudoknot structure of RydC (3, 44)

489 mutations in the 3' end of RydC have a dramatic impact on RydC stability (45), thus we were not

490 able to test a RydC compensatory mutant that would restore pairing to the *trpE20* mutant fusion.

491 However, we did construct a second *trpE* fusion, P$_{BAD}$-*trunc-trpE'-'lacZ,* which was truncated to

492 remove the putative RydC binding site (Fig. 7B). This fusion was no longer activated by RydC at

493 all. These observations suggest that sequences early in the *trpE* coding sequence are important

494 for RydC-mediated increase in *trpE* translation.

495 **DISCUSSION**

496 Over the years, many sRNAs have been discovered and characterized using both

497 computational and experimental methods. Although target discovery of sRNAs still remains the

498 rate-limiting step in sRNA characterization, many new techniques have been developed to

499 overcome that obstacle. Some techniques take a purely computational approach to target

500 prediction, including the target-prediction algorithms we have included in SPOT, (9, 10, 11, 12,

501 13) and others we have not included (47-57). Experimental techniques to identify bacterial

502 sRNA targets have also expanded. Many of these use affinity purification or co-

503 immunoprecipitation approaches, with or without crosslinking (15, 17, 20, 58, 59, 60). To help

504 streamline the process of sRNA target identification, the SPOT pipeline was constructed to be

505 used in conjunction with other identification methods. In this study, we showed that the SPOT

506 pipeline achieved ≥ 75% sensitivity and ≤ 50% false positive rate when at least 2 methods

507 converged on a prediction for the well-characterized sRNAs SgrS and RyhB (Fig. 2D-E).

508 Expanding our analysis to other bacterial sRNAs, we found that the pipeline sensitivity was

509 equal to or exceeded that of any individual method (average = 84% ± 8.5%, Fig. 3A, Fig. S2).

510 As before, we found that correct identification by ≥2 methods occurred in the majority of

511 instances (Fig. 3A). Furthermore, SPOT can be applied to the reverse situation where a user

512 can search for potential sRNAs that regulate their UTR of interest. We found through these

513    analyses that for 11 *E. coli* 5' UTRs with ≥2 known interactions with sRNAs, the analysis gave

514    an average sensitivity of 85% ± 24% (Fig. 3A).

515    To test the utility of SPOT in identifying novel sRNA-mRNA target interactions, we used

516    it to predict targets of the poorly characterized sRNA RydC, which had been described to

517    regulate three genes: *yejA* (43), *cfa* (3), and *csgD* (44). Through SPOT analyses and filtering

518    based on experimental data, we generated a list of putative RydC targets (Table 1).

519    Reassuringly, SPOT identified the true RydC target, *cfa* mRNA, and correctly predicted the

520    known binding site on this target (Table 1, Supplemental Dataset 1). The other two reported

521    targets, *yejA* and *csgD,* were not identified by the SPOT computational pipeline, nor were these

522    genes differentially regulated in our RydC pulse-expression RNA-seq analyses (Supplemental

523    Dataset 2). Since no specific direct binding interactions were shown for RydC-*yejA* or RydC-

524    *csgD*, we postulate that the previously observed regulation of these targets by RydC may be

525    indirect. The SPOT pipeline also correctly identified 2 additional RydC targets, *pheA* and *trpE*

526    (Table 1, Figs. 5C, 6, 7). RydC represses *pheA* translation, likely by a mechanism common to

527    repressing sRNAs. Binding of RydC to sequences around the Shine-Dalgarno region would

528    prevent ribosome binding and inhibit translation initiation. The mechanism of RydC-dependent

529    activation of *trpE* appears to be more complex. The *trpE* gene is part of the *trpLEDCBA* operon

530    responsible for L-tryptophan biosynthesis, which is regulated by both the *trpR* repressor and an

531    attenuation mechanism (46). Depending on the availability of L-tryptophan, the ribosome can

532    either stalls at or moves quickly through Trp codons in the *trpL* ORF. When Trp is abundant, the

533    ribosome rapidly completes translation of *trpL,* which prevents co-transcriptional formation of an

534    antiterminator hairpin and allows formation of a transcription terminator just upstream of the *trpE*

535    coding sequence. When Trp is limiting, ribosome stalling at the Trp codons allows formation of

536    an antiterminator structure, which promotes transcription elongation into downstream Trp

537    biosynthesis structural genes. While sequences within the *trpE* coding sequence have not been

538    implicated in the Trp-dependent attenuation mechanism, it is possible that the sequences

21

539     including the RydC binding site are responsible for yet another layer of regulation of these

540     genes, perhaps at the level of translation. Alternatively, sequences in the *trpE* coding sequence

541     could have long-range interactions with the upstream terminator or antiterminator sequences

542     and RydC binding could modulate those interactions.

543         Our study and evaluation of a combinatorial approach to identify mRNA targets of

544     sRNAs of interest represents a step toward accelerating a rate-limiting step in sRNA

545     characterization. The SPOT pipeline is able to streamline the process of running individual

546     algorithms, which can take hours to days, by reducing the run times significantly for all 4

547     algorithms at once (under 2 hours). Since the pipeline runs all 4 algorithms simultaneously, a

548     more narrowed down, comprehensive list is generated, negating the need for manually selecting

549     targets from individual algorithm runs. However, every method has drawbacks and though

550     SPOT is a powerful tool, it has limitations as well. For instance, a 50% false positive rate (the

551     average for well-characterized sRNAs analyzed in this study) is still high even though it is

552     markedly better than the false positive rates of predictions made by any single algorithm. As

553     experimental approaches for sRNA-mRNA target identification continue to improve, the power

554     and accuracy of SPOT's combinatorial approach to sRNA-target binding site predictions will

555     likewise improve. Another factor impacting the accurate prediction of sRNA binding sites by

556     SPOT is the user-defined search window. The majority of early examples of sRNA-mediated

557     regulation involved sRNAs binding in translation initiation regions of target mRNAs. Thus, most

558     existing sRNA target prediction algorithms have default windows set to search around start

559     codons. As more sRNA-mRNA interactions are validated and mechanisms of regulation studied,

560     we and others have found increasing numbers of examples of sRNA-mRNA interactions that

561     occur outside this window. Some of these interactions are primary or only interactions

562     responsible for sRNA-mediated regulation of the mRNA, e.g., RydC-*cfa* mRNA (3), SgrS-*yigL*

563     mRNA (39), which both involve mRNA sequences far upstream of the start codons. Yet other

564     interactions involving mRNA sequences far from translation initiation regions represent

565     secondary or auxiliary binding interactions that nevertheless play important roles in regulation

566     (18, 38).

567          For the sRNA SgrS, there are two binding sites for its interaction with *asd* mRNA (18),

568     but SPOT was only able to predict the primary binding site. We expect that there are other

569     examples where the algorithms have failed to identify alternate or additional binding sites. This

570     is currently an area of development and once implemented, will serve as a valuable asset in

571     identifying putative targets for a sRNA of interest.

572          Taken together, the combinatorial approach revealed two new targets, *pheA* and *trpE*, in

573     the RydC regulon. Interestingly, both PheA and TrpE are involved in the chorismate metabolic

574     pathway, with PheA using chorismate as a substrate in L-Tyrosine/L-Phenylalanine biosynthesis

575     and TrpE for L-Tryptophan biosynthesis. Interestingly, RydC repressed *pheA* whereas it

576     activated *trpE,* an unusual case since both are involved in amino acid biosynthesis in divergent

577     pathways. In the case for *trpE,* the mechanism of positive regulation is unique in that the base

578     pairing interaction takes place 12-22 nt downstream of the start codon. RydC could possibly

579     serve as a sRNA modulator of the biosynthetic pools of amino acids by activating/repressing

580     *trpE*/*pheA* mRNA expression when necessary. As an aside, chorismate is also a substrate for

581     production of the *E. coli* siderophore enterobactin, which is synthesized under iron limiting

582     conditions. Mutations in *fur*, *tyrA*, *pheA*, or *pheU* resulted in increased enterobactin production

583     since the chorismate pools were used for enterobactin synthesis (62). These observations

584     suggest that there may be conditions where RydC impacts the iron starvation stress response,

585     perhaps forming a regulatory network that intersects with that of the well-characterized iron

586     starvation stress response sRNA, RyhB. To better understand these potential connections,

587     future work will be aimed at characterizing the regulators and conditions controlling synthesis of

588     RydC.

589          With the implementation of the SPOT pipeline, combined with RNA-Seq and MAPS data,

590     we were able to add to the RydC regulon and expand its network. Whether this regulatory

591    network is exhaustive remains to be determined. We note that there were other RydC-mRNA

592    binding interactions predicted by SPOT that were not analyzed further here. Moreover, there are

593    additional sRNA-mRNA interactions predicted by SPOT for the other sRNAs that were run

594    through the pipeline (Supplementary Dataset 1) and it is likely that more bona fide interactions

595    are among those predictions. All in all, we developed a streamlined method for sRNA-mRNA

596    binding site predictions that leverages the strengths of many pre-existing algorithms. We

597    showed the robustness of SPOT for identification of true sRNA-mRNA interactions using well-

598    characterized and poorly characterized sRNAs. We anticipate that SPOT will become a valuable

599    tool for many investigators who have found interesting sRNAs and wish to identify potential

600    mRNA targets for further characterization.

601    **Acknowledgements**

608    **Funding**

**REFERENCES**

1. Storz G, Vogel J, Wassarman KM. 2011. Regulation by small RNAs in bacteria: expanding frontiers. Mol Cell 43:880–891.

2. Papenfort K, Vanderpool CK. 2015. Target activation by regulatory RNAs in bacteria. FEMS Microbiol Rev 39:362–378.

3. Fröhlich KS, Papenfort K, Fekete A, Vogel J. 2013. A small RNA activates CFA synthase by isoform-specific mRNA stabilization. EMBO J 32:2963–2979.

4. Soper T, Mandin P, Majdalani N, Gottesman S, Woodson SA. 2010. Positive regulation by small RNAs and the role of Hfq. Proc Natl Acad Sci U S A 107:9602–9607.

5. Dambach M, Irnov I, Winkler WC. 2013. Association of RNAs with Bacillus subtilis Hfq. PLoS One 8:e55156.

6. Mollerup MS, Ross JA, Helfer A-C, Meistrup K, Romby P, Kallipolitis BH. 2016. Two novel members of the LhrC family of small RNAs in Listeria monocytogenes with overlapping regulatory functions but distinctive expression profiles. RNA Biol 13:895–915.

7. Ryan D, Mukherjee M, Suar M. 2017. The expanding targetome of small RNAs in Salmonella Typhimurium. Biochimie 137:69–77.

8. Vogel J. 2009. A rough guide to the non-coding RNA world of Salmonella. Mol Microbiol 71:1–11.

9. Kery MB, Feldman M, Livny J, Tjaden B. 2014. TargetRNA2: identifying targets of small regulatory RNAs in bacteria. Nucleic Acids Res 42:W124–9.

10. Ying X, Cao Y, Wu J, Liu Q, Cha L, Li W. 2011. sTarPicker: a method for efficient prediction of bacterial sRNA targets based on a two-step model for hybridization. PLoS One 6:e22705.

11. Busch A, Richter AS, Backofen R. 2008. IntaRNA: efficient prediction of bacterial sRNA targets incorporating target site accessibility and seed regions. Bioinformatics 24:2849–2856.

12. Mann M, Wright PR, Backofen R. 2017. IntaRNA 2.0: enhanced and customizable prediction of RNA-RNA interactions. Nucleic Acids Res 45:W435–W439.

13. Wright PR, Richter AS, Papenfort K, Mann M, Vogel J, Hess WR, Backofen R, Georg J. 2013. Comparative genomics boosts target prediction for bacterial small RNAs. Proc Natl Acad Sci U S A 110:E3487–96.

25

646   14.   Wright PR, Georg J, Mann M, Sorescu DA, Richter AS, Lott S, Kleinkauf R, Hess WR,
647         Backofen R. 2014. CopraRNA and IntaRNA: predicting small RNA targets, networks and
648         interaction domains. Nucleic Acids Res 42:W119–23.

649   15.   Melamed S, Peer A, Faigenbaum-Romm R, Gatt YE, Reiss N, Bar A, Altuvia Y,
650         Argaman L, Margalit H. 2016. Global Mapping of Small RNA-Target Interactions in
651         Bacteria. Mol Cell 63:884–897.

652   16.   August JT, Eoyang L, De Fernandez MT, Hasegawa S, Kuo CH, Rensing U, Shapiro L.
653         1970. Phage-specific and host proteins in the replication of bacteriophage RNA. Fed
654         Proc 29:1170–1175.

655   17.   Lalaouna D, Massé E. 2015. Identification of sRNA interacting with a transcript of
656         interest using MS2-affinity purification coupled with RNA sequencing (MAPS)
657         technology. Genom Data 5:136–138.

658   18.   Bobrovskyy M, Vanderpool CK. 2016. Diverse mechanisms of post-transcriptional
659         repression by the small RNA regulator of glucose-phosphate stress. Mol Microbiol
660         99:254–273.

661   19.   Massé E, Vanderpool CK, Gottesman S. 2005. Effect of RyhB small RNA on global iron
662         use in Escherichia coli. J Bacteriol 187:6962–6971.

663   20.   Lalaouna D, Carrier M-C, Semsey S, Brouard J-S, Wang J, Wade JT, Massé E. 2015. A
664         3' external transcribed spacer in a tRNA transcript acts as a sponge for small RNAs to
665         prevent transcriptional noise. Mol Cell 58:393–405.

666   21.   Keseler IM, Mackie A, Santos-Zavaleta A, Billington R, Bonavides-Martínez C, Caspi R,
667         Fulcher C, Gama-Castro S, Kothari A, Krummenacker M, Latendresse M, Muñiz-
668         Rascado L, Ong Q, Paley S, Peralta-Gil M, Subhraveti P, Velázquez-Ramírez DA,
669         Weaver D, Collado-Vides J, Paulsen I, Karp PD. 2017. The EcoCyc database: reflecting
670         new knowledge about Escherichia coli K-12. Nucleic Acids Res 45:D543–D550.

671   22.   Nawrocki EP, Eddy SR. 2013. Infernal 1.1: 100-fold faster RNA homology searches.
672         Bioinformatics 29:2933–2935.

673   23.   Castillo-Keller M, Vuong P, Misra R. 2006. Novel mechanism of Escherichia coli porin
674         regulation. J Bacteriol 188:576–586.

675   24.   Shinhara A, Matsui M, Hiraoka K, Nomura W, Hirano R, Nakahigashi K, Tomita M, Mori
676         H, Kanai A. 2011. Deep sequencing reveals as-yet-undiscovered small RNAs in
677         Escherichia coli. BMC Genomics 12:428.

678   25.   Rau MH, Bojanovič K, Nielsen AT, Long KS. 2015. Differential expression of small RNAs
679         under chemical stress and fed-batch fermentation in E. coli. BMC Genomics 16:1051.

680 26.  Datsenko KA, Wanner BL. 2000. One-step inactivation of chromosomal genes in
681      Escherichia coli K-12 using PCR products. Proc Natl Acad Sci U S A 97:6640–6645.

682 27.  Yu D, Ellis HM, Lee EC, Jenkins NA, Copeland NG, Court DL. 2000. An efficient
683      recombination system for chromosome engineering in Escherichia coli. Proc Natl Acad
684      Sci U S A 97:5978–5983.

685 28.  Zubay G, Morse DE, Schrenk WJ, Miller JH. 1972. Detection and isolation of the
686      repressor protein for the tryptophan operon of Escherichia coli. Proc Natl Acad Sci U S A
687      69:1100–1103.

688 29.  Mandin P, Gottesman S. 2009. A genetic approach for finding small RNAs regulators of
689      genes of interest identifies RybC as regulating the DpiA/DpiB two-component system.
690      Mol Microbiol 72:551–565.

691 30.  Thomason L, Court DL, Bubunenko M, Costantino N, Wilson H, Datta S, Oppenheim A.
692      2007. Recombineering: genetic engineering in bacteria using homologous
693      recombination. Curr Protoc Mol Biol Chapter 1:Unit 1.16.

694 31.  Aiba H, Adhya S, de Crombrugghe B. 1981. Evidence for two functional gal promoters in
695      intact Escherichia coli cells. J Biol Chem 256:11905–11910.

696 32.  McClure R, Balasubramanian D, Sun Y, Bobrovskyy M, Sumby P, Genco CA,
697      Vanderpool CK, Tjaden B. 2013. Computational analysis of bacterial RNA-Seq data.
698      Nucleic Acids Res 41:e140.

699 33.  Miller JH. 1972. Experiments in Molecular Genetics.

700 34.  Bobrovskyy M, Vanderpool CK. 2014. The small RNA SgrS: roles in metabolism and
701      pathogenesis of enteric bacteria. Front Cell Infect Microbiol 4:61.

702 35.  Salvail H, Massé E. 2012. Regulating iron storage and metabolism with RNA: an
703      overview of posttranscriptional controls of intracellular iron homeostasis. Wiley
704      Interdiscip Rev RNA 3:26–36.

705 36.  Vanderpool CK, Gottesman S. 2004. Involvement of a novel transcriptional activator and
706      small RNA in post-transcriptional regulation of the glucose phosphoenolpyruvate
707      phosphotransferase system. Mol Microbiol 54:1076–1089.

708 37.  Rice JB, Vanderpool CK. 2011. The small RNA SgrS controls sugar-phosphate
709      accumulation by regulating multiple PTS genes. Nucleic Acids Res 39:3806–3819.

710 38.  Rice JB, Balasubramanian D, Vanderpool CK. 2011. Small RNA binding-site multiplicity
711      involved in translational regulation of a polycistronic mRNA. Proc Natl Acad Sci U S A
712      109:E2691–8.

713   39.   Papenfort K, Sun Y, Miyakoshi M, Vanderpool CK, Vogel J. 2013. Small RNA-mediated
714          activation of sugar phosphatase mRNA regulates glucose homeostasis. Cell 153:426–
715          437.

716   40.   Massé E, Vanderpool CK, Gottesman S. 2005. Effect of RyhB small RNA on global iron
717          use in Escherichia coli. J Bacteriol 187:6962–6971.

718   41.   Kim JN, Kwon YM. 2013. Genetic and phenotypic characterization of the RyhB regulon
719          in Salmonella Typhimurium. Microbiol Res 168:41–49.

720   42.   Bos J, Duverger Y, Thouvenot B, Chiaruttini C, Branlant C, Springer M, Charpentier B,
721          Barras F. 2013. The sRNA RyhB regulates the synthesis of the Escherichia coli
722          methionine sulfoxide reductase MsrB but not MsrA. PLoS One 8:e63647.

723   43.   Antal M, Bordeau V, Douchin V, Felden B. 2005. A small bacterial RNA regulates a
724          putative ABC transporter. J Biol Chem 280:7901–7908.

725   44.   Bordeau V, Felden B. 2014. Curli synthesis and biofilm formation in enteric bacteria are
726          controlled by a dynamic small RNA module made up of a pseudoknot assisted by an
727          RNA chaperone. Nucleic Acids Res 42:4682–4696.

728   45.   Dimastrogiovanni D, Fröhlich KS, Bandyra KJ, Bruce HA, Hohensee S, Vogel J, Luisi
729          BF. 2014. Recognition of the small regulatory RNA RydC by the bacterial Hfq protein.
730          Elife 3.

731   46.   Yanofsky C, Platt T, Crawford IP, Nichols BP, Christie GE, Horowitz H, VanCleemput M,
732          Wu AM. 1981. The complete nucleotide sequence of the tryptophan operon of
733          Escherichia coli. Nucleic Acids Res 9:6647–6668.

734   47.   Eggenhofer F, Tafer H, Stadler PF, Hofacker IL. 2011. RNApredator: fast accessibility-
735          based prediction of sRNA targets. Nucleic Acids Res 39:W149–54.

736   48.   Mückstein U, Tafer H, Hackermüller J, Bernhart SH, Stadler PF, Hofacker IL. 2006.
737          Thermodynamics of RNA-RNA binding. Bioinformatics 22:1177–1182.

738   49.   Hofacker IL, Fontana W, Stadler PF, Bonhoeffer LS, Tacker M, Schuster P. 1994. Fast
739          folding and comparison of RNA secondary structures. Monatshefte für Chemie /
740          Chemical Monthly 125:167–188.

741   50.   Rehmsmeier M, Steffen P, Hochsmann M, Giegerich R. 2004. Fast and effective
742          prediction of microRNA/target duplexes. RNA 10:1507–1517.

743   51.   Krüger J, Rehmsmeier M. 2006. RNAhybrid: microRNA target prediction easy, fast and
744          flexible. Nucleic Acids Res 34:W451–4.

745   52.   Bernhart SH, Tafer H, Mückstein U, Flamm C, Stadler PF, Hofacker IL. 2006. Partition
746          function and base pairing probabilities of RNA heterodimers. Algorithms Mol Biol 1:3.

747  53.  Tulpan D, Andronescu M, Chang SB, Shortreed MR, Condon A, Hoos HH, Smith LM.
748      2005. Thermodynamically based DNA strand design. Nucleic Acids Res 33:4951–4964.
749  54.  Andronescu M, Zhang ZC, Condon A. 2005. Secondary structure prediction of
750      interacting RNA molecules. J Mol Biol 345:987–1001.
751  55.  Wenzel A, Akbasli E, Gorodkin J. 2012. RIsearch: fast RNA-RNA interaction search
752      using a simplified nearest-neighbor energy model. Bioinformatics 28:2738–2746.
753  56.  Gerlach W, Giegerich R. 2006. GUUGle: a utility for fast exact matching under RNA
754      complementary rules including G-U base pairing. Bioinformatics 22:762–764.
755  57.  Tafer H, Hofacker IL. 2008. RNAplex: a fast tool for RNA-RNA interaction search.
756      Bioinformatics 24:2657–2663.
757  58.  Han K, Tjaden B, Lory S. 2016. GRIL-seq provides a method for identifying direct targets
758      of bacterial small regulatory RNA by in vivo proximity ligation. Nat Microbiol 2:16239.
759  59.  Waters SA, McAteer SP, Kudla G, Pang I, Deshpande NP, Amos TG, Leong KW,
760      Wilkins MR, Strugnell R, Gally DL, Tollervey D, Tree JJ. 2017. Small RNA interactome of
761      pathogenic E. coli revealed through crosslinking of RNase E. EMBO J 36:374–387.
762  60.  Lioliou E, Sharma CM, Altuvia Y, Caldelari I, Romilly C, Helfer A-C, Margalit H, Romby
763      P. 2013. In vivo mapping of RNA-RNA interactions in Staphylococcus aureus using the
764      endoribonuclease III. Methods 63:135–143.
765  61.  Jagodnik J, Chiaruttini C, Guillier M. 2017. Stem-Loop Structures within mRNA Coding
766      Sequences Activate Translation Initiation and Mediate Control by Small Regulatory
767      RNAs. Mol Cell 68:158–170.e3.
768  62.  Foster MS, Carroll JN, Niederhoffer EC. 1994. Phenylalanine- and tyrosine-dependent
769      production of enterobactin in Escherichia coli. FEMS Microbiol Lett 117:79–83.
770
771
772
773
774
775
776
777
778
779
780
781

782 **Table 1. List of putative RydC targets chosen for further testing[a, b, c]**

783

| Gene | Putative function | Fold change P$_{lac}$-rydC/vector | Predicted interactions | Algorithm predictions | Reference |
|---|---|---|---|---|---|
| cfa | cyclopropane fatty acyl phospholipid synthase | 31.00 | cfa -110 to -98 RydC +14 to +2 | T, S, I, C | (3) |
| grpE | nucleotide exchange factor | 0.32 | grpE +16 to +29 RydC +64 to +51 | T, S, I, C | This study |
| moaB | part of moaABCDE operon | 0.63 | moaB +27 to +42 RydC +44 to +29 | T, S, I, C | This study |
| araH | arabinose ABC transporter membrane subunit | 0.67 | araH +9 to +25 RydC +17 to +1 | T, S, I, C | This study |
| yhjD | putative transporter | 0.30 | yhjD -62 to -11 RydC +53 to +2 | S, I, C | This study |
| ygaU | potassium binding protein (kbp) | 0.67 | ygaU +78 to +94 RydC +31 to +15 | T, S, I | This study |
| yibT | protein YibT | 0.30 | yibT -24 to -10 RydC +28 to +14 | T, S, I | This study |
| trpE | anthranilate synthase subunit | 1.00 | trpE +12 to +22 RydC +47 to +37 | T, S, I | This study |
| pheA | fused chorismate mutase/prephenate dehydratase | 0.30 | pheA +4 to +11 RydC +10 to +3 | S, I, C | This study |
| cysQ | 3′ (2′), 5′-bisphosphate nucleotidase | 0.59 | cysQ +48 to +67 RydC +21 to +2 | T, S, I | This study |
| purK | 5-(carboxyamino) imidazole ribonucleotide synthase | 3.49 | purK -38 to -19 RydC +59 to +5 | S, I, C | This study |
| csgD | DNA binding transcriptional dual regulator | 0.75 | csgD -19 to +3 RydC +26 to +5 | I | (44) |
| yejA | putative oligopeptide ABC transporter periplasmic component | 0.91 | yejA +1265 to +1273 RydC +47 to +39 | I | (43) |
| lldR | DNA-binding transcriptional dual regulator | 1.17 | lldR +15 to +65 RydC +61 to +15 | S | This study |

784  [a]Column #3 (Fold change Plac-rydC/vector) lists the determined ratio from RNA-Seq experiments

785  (Supplemental Dataset 1)

786  [b]Column #4 (Predicted interactions) lists the bases involved in the interaction in the 5′ to 3′ direction for

787  target and 3′ to 5′ direction for RydC in relation to the +1 site (start of translation)

788  [c]Column #5 (Algorithm predictions) lists the algorithms that predicted a base-pairing interaction with T =

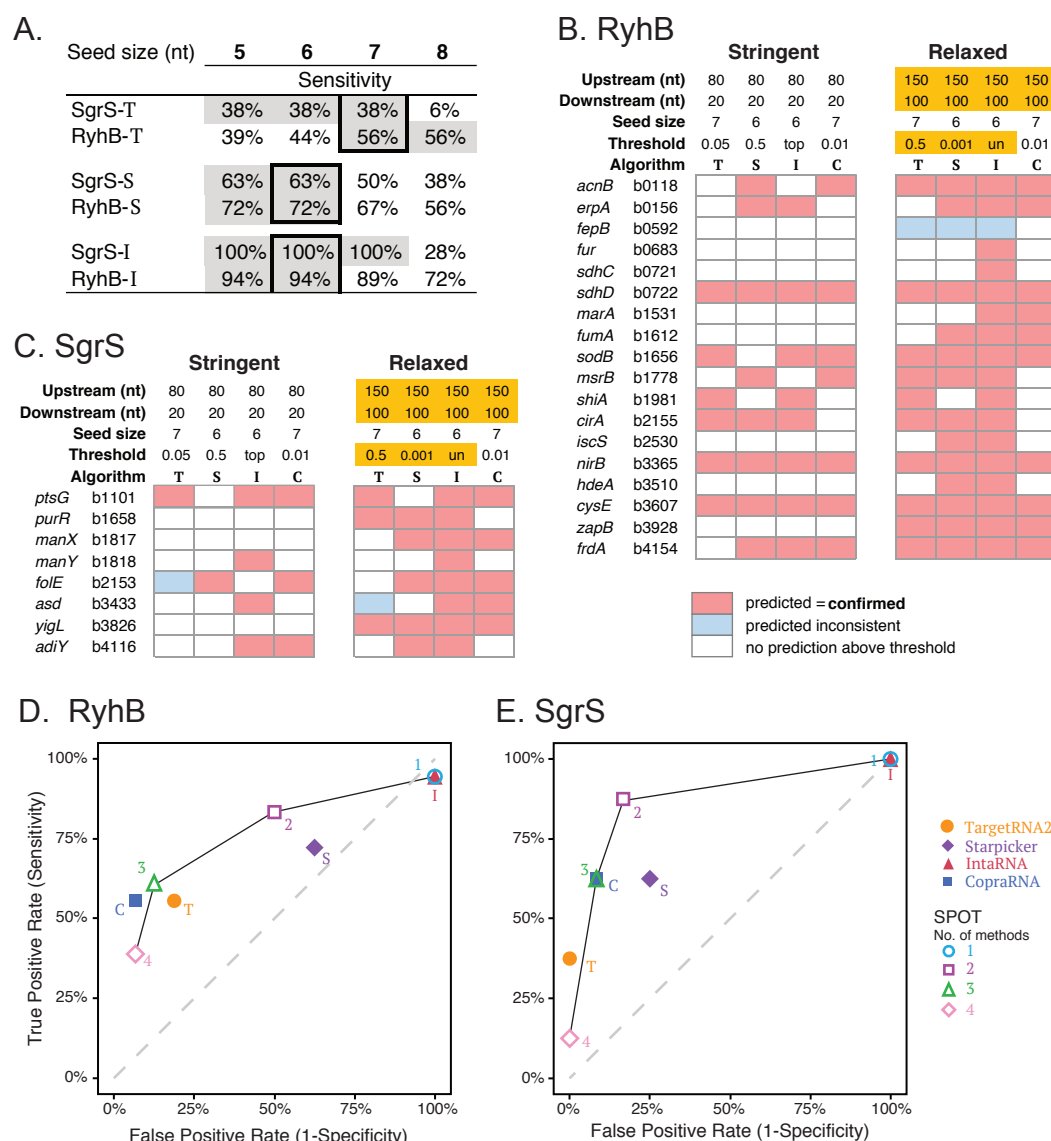789  TargetRNA2, S = sTarPicker, I = IntaRNA, and C = CopraRNA

790

**Figure. 1. Schematic diagram of SPOT pipeline.**
Step 1: A basic implementation of SPOT requires a user-provided reference genome and a sRNA sequence file. The user can customize the search window size and can optionally provide information required for CopraRNA (dashed boxes). Step 2: The user can set seed sizes and significance cutoffs for each algorithm (superscript t = TargetRNA2, s = sTarPicker, i = IntaRNA). Step 3: SPOT runs the algorithms in parallel and generates a set of collated results. Step 4: Results filtering options as shown narrow the list of predicted interactions to an experimentally-tractable size for further validation or analysis.

## Fig. 2



**Figure. 2. Validation of SPOT using known SgrS and RyhB sRNA-mRNA interactions.**
(A) "Seed size" indicates the number of consecutive basepairing nucleotides in an sRNA-mRNA interaction prediction. This is an adjustable parameter for each algorithm. Seed sizes were varied from 5 to 8 nt and the sensitivity (true positive rate) was determined for known SgrS and RyhB interactions while all other parameters were kept constant. Optimal seed sizes (bold boxes) were chosen for each algorithm. Highest percentage values for sensitivity are indicated with gray shading. Algorithms were abbreviated -T=TargetRNA2, -S=Starpicker, -I=IntaRNA.
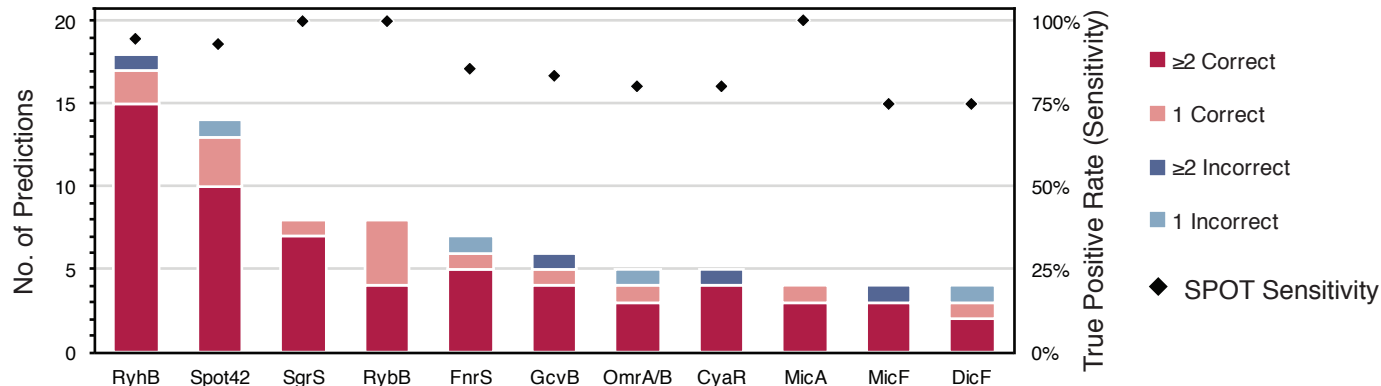(B & C) Analyses were re-run using optimal seed sizes identified in A, but using a 'Stringent' parameter set with a narrow window size and high individual significance thresholds or a 'Relaxed' parameter set with a wider window size and low individual significance thresholds. Correctly predicted interactions for RyhB and SgrS are shown as pink cells, predictions that were inconsistent with confirmed interaction sites are shown in blue, and empty cells did not have any predictions above the indicated thresholds. Algorithms are abbreviated T=TargetRNA2, S=Starpicker, I=IntaRNA, and C=CopraRNA.
(D & E) RyhB and SgrS have experimentally validated (true positive) and invalidated (true negative) mRNA targets, which were used to generate receiver operating characteristic (ROC) curves. These plots enable assessment of the accuracy of SPOT and the individual algorithms. Using the 'Relaxed' search parameters, 2-algorithm agreement in SPOT had greater True Positive Rates and more acceptable False Positive Rates compared with individual algorithms with the same settings.
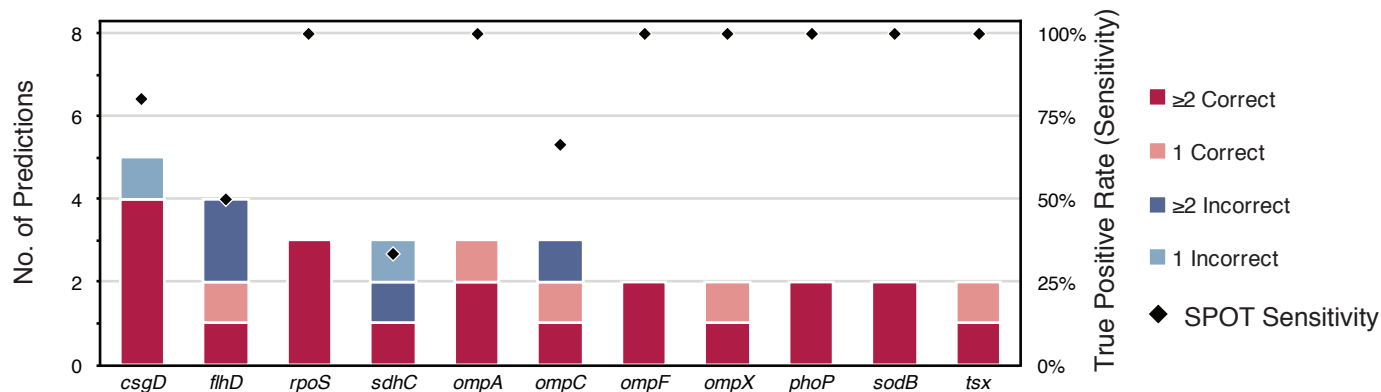
**Figure. 3. SPOT demonstrates high sensitivity for detecting targets of multiple sRNAs.**
(A) Along with RyhB and SgrS, nine additional sRNAs were analyzed with SPOT using the 'Relaxed' parameter set, demonstrating the robustness of the SPOT pipeline for correctly identifying sRNA-mRNA target interactions. Stacked bars show the number of experimentally validated mRNA targets correctly or incorrectly identified by 1 or ≥2 methods. Black diamonds indicate the overall True Positive Rate (sensitivity) of SPOT for each sRNA.
(B) Eleven UTRs that are experimentally validated to interact with multiple sRNAs were used in a 'reverse' search in SPOT (i.e., using the UTR as the query and the sRNAs as the targets). The average sensitivity of this method is lower than in A., however, this is a novel means for identifying sRNAs that might affect genes of interest. Plots are drawn as in A.

# Fig. 4

| | | Stringent | | | | Relaxed | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Upstream (nt) | | 80 | 80 | 80 | 80 | 150 | 150 | 150 | 150 | |
| Downstream (nt) | | 20 | 20 | 20 | 20 | 100 | 100 | 100 | 100 | |
| Seed size | | 7 | 6 | 6 | 7 | 7 | 6 | 6 | 7 | |
| Threshold | | 0.05 | 0.5 | top | 0.01 | 0.5 | 0.001 | un | 0.01 | |
| Algorithm | | T | S | I | C | T | S | I | C | Fold change |
| *cfa* | b1661 | | | | | blue | pink | pink | pink | 31.00 |
| *grpE* | b2614 | green | | | | green | green | green | green | 0.32 |
| *moaB* | b0782 | | | | | green | green | green | green | 0.63 |
| *araH* | b4460 | | green | green | | green | green | green | green | 0.67 |
| *yqgC* | b2940 | | | | | green | green | | green | 1.00 |
| *mdtH* | b1065 | | green | | | green | green | | green | 1.00 |
| *waaC* | b3621 | | green | | | | green | green | green | 0.91 |
| *ybiT* | b0820 | | | | | green | green | green | | 1.50 |
| *cspC* | b1823 | | green | | | green | green | green | | 1.20 |
| *rsmF* | b1835 | | | | | green | green | green | | 1.30 |
| *paaD* | b1391 | | | | | | green | green | green | 0.59 |
| *ppdB* | b2825 | | green | green | | green | green | green | | 0.77 |
| *cytR* | b3934 | | | | | green | green | green | | 1.00 |
| *recC* | b2822 | | | | | | green | green | green | 0.67 |
| *nhoA* | b1463 | | green | | | green | green | green | | 1.00 |
| *hemX* | b3803 | | | | | blue | green | green | | 0.77 |
| *cycA* | b4208 | | | | | green | green | green | | 1.50 |
| *yhjD* | b3522 | | | | | | green | green | blue | 0.30 |
| *ygaU* | b2665 | | | | | green | green | blue | green | 0.67 |
| *yafT* | b0217 | | | | | green | green | green | | 1.00 |
| *yibT* | b4554 | | | | | green | blue | green | | 0.30 |
| *trpE* | b1264 | green | | green | | green | green | green | | 1.00 |
| *panB* | b0134 | green | green | green | | green | green | green | | 1.00 |
| *mukE* | b0923 | | green | | | green | green | green | | 1.00 |
| *pheA* | b2599 | | green | | green | | green | green | | 0.30 |
| *yjiM* | b4335 | | | | | green | green | blue | | 1.50 |
| *cysQ* | b4214 | | | | | green | green | green | | 0.59 |
| *eamB* | b2578 | green | green | green | | green | green | green | | 0.77 |
| *yfcE* | b2300 | | green | | green | | green | green | | 1.00 |
| *tamB* | b4221 | | green | green | green | | green | green | green | 0.71 |
| *yhdJ* | b3262 | | green | green | | green | green | green | | 1.50 |
| *nagK* | b1119 | | green | green | green | | green | green | | 1.10 |
| *sgcA* | b4302 | green | | green | | | green | green | | 0.77 |
| *cdh* | b3918 | | green | | | green | green | blue | | 0.91 |
| *purK* | b0522 | green | green | | green | | green | green | | 3.50 |
| *deoC* | b4831 | | | | | | green | green | | 0.20 |
| *purM* | b2099 | | | | | | green | green | | 3.30 |
| *malM* | b4037 | | | | | | green | green | | 0.04 |
| *malP* | b3417 | | | | | | green | green | | 0.14 |
| *malF* | b4033 | | | | | | green | green | | 0.09 |
| *nrdI* | b2674 | | | | | | green | green | | 0.42 |
| *pstS* | b3728 | | | | | | green | green | | 0.13 |
| *prlC* | b3498 | | | | | | green | green | | 0.25 |
| *hslV* | b3932 | | | | | | green | green | | 0.24 |
| *udp* | b3831 | | | | | | green | green | | 0.22 |
| *pstC* | b3727 | | | | | | green | green | | 0.23 |
| *ytjA* | b4568 | | | | | | green | green | | 0.50 |
| *hslR* | b3400 | | | | | | green | green | | 0.40 |
| *dnaJ* | b0015 | | | | | green | | green | | 0.36 |
| *yheV* | b4551 | | | | | | green | green | | 0.50 |
| *ybeY* | b0659 | | | | | | green | green | | 0.50 |
| *dnaK* | b0014 | | | | | | green | green | | 0.15 |
| *fxsA* | b4140 | | | | | | green | green | | 0.32 |
| *tabA* | b4252 | | | | | | green | green | | 0.43 |

Legend:
- pink = predicted = **confirmed**
- blue = predicted inconsistent
- green = predicted unknown
- white = no prediction above threshold
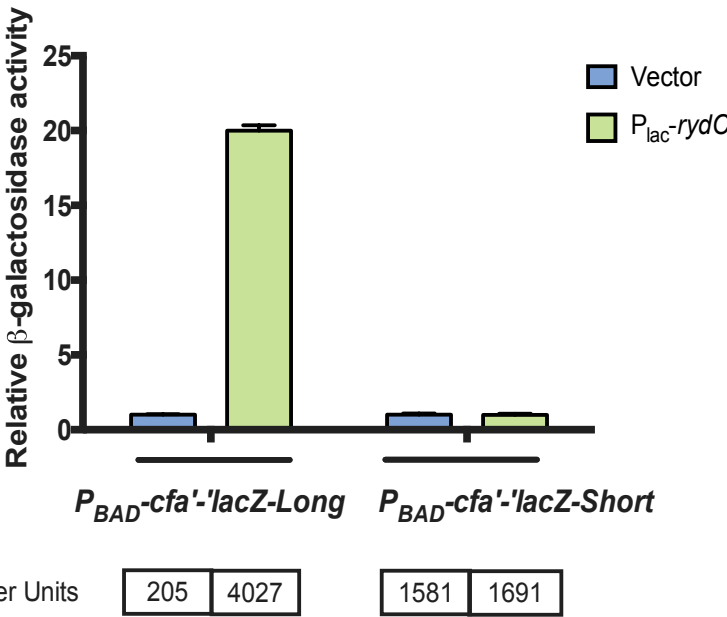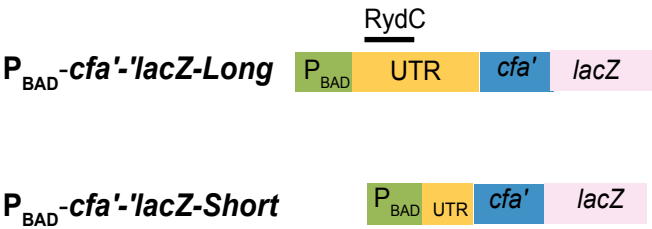
**Figure. 4. SPOT predictions for the sRNA RydC.**
Analyses were run with optimal seed sizes as determined in Fig. 2. Genes above the bold line denote those with ≥ 3 computational predictions, while genes below the line had 2 computational predictions and differential RNA-seq expression (fold change of ≥ 1.5 or ≤ 0.5, q-value of ≤ 0.005). Correctly predicted interactions for RydC are shown as pink cells, unknown predictions that were consistent among algorithms are shown in green, inconsistent predictions are shown in blue, and empty cells did not have any predictions above the indicated thresholds. Algorithms are abbreviated T=TargetRNA2, S=Starpicker, I=IntaRNA, and C=CopraRNA.
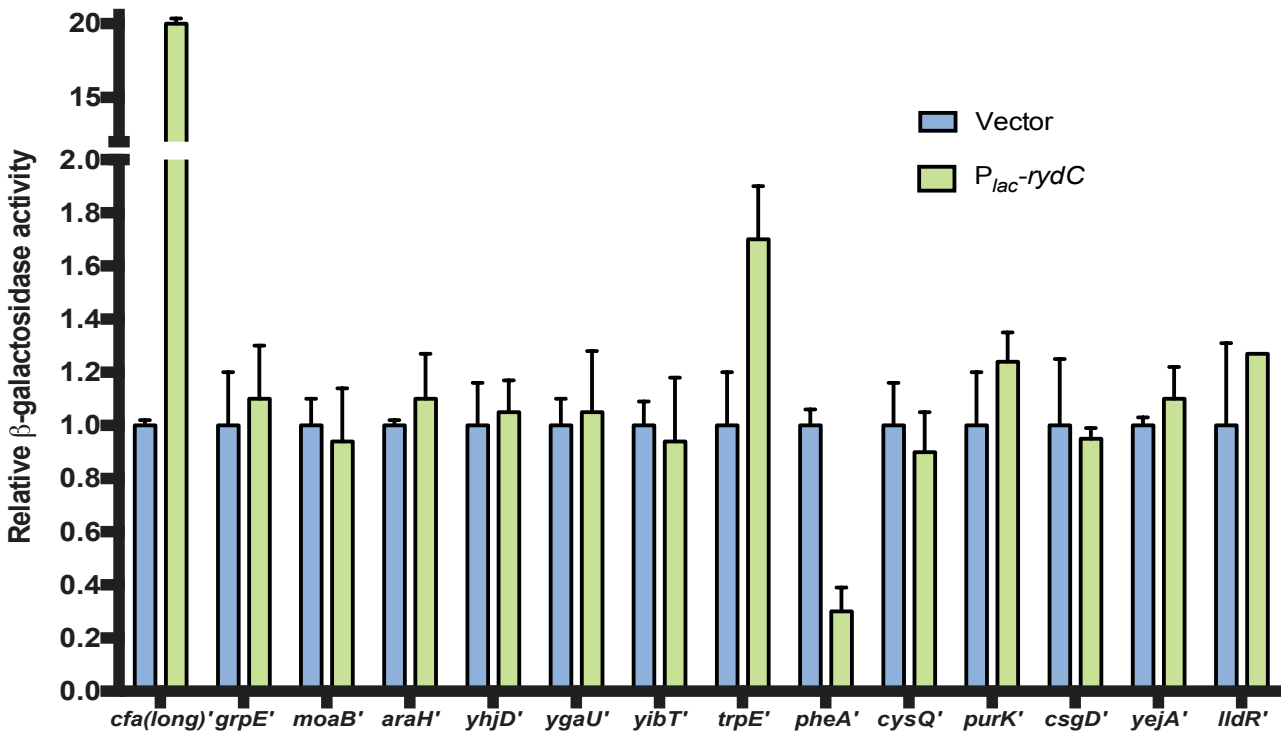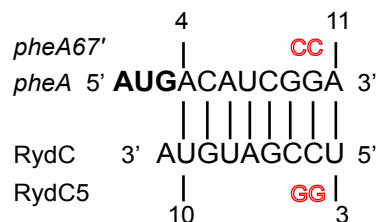
**Figure. 5. Validation of RydC target predictions.**
(A) The design for the translational *lacZ* constructs is shown, where the green box indicates the arabinose promoter ($P_{BAD}$), yellow the untranslated region (UTR), blue the coding sequence (CDS), and pink the *lacZ* gene.
(B) To confirm *cfa* as a RydC target, both full-length and shortened *cfa'-'lacZ* translational fusions were tested in backgrounds with vector or $P_{lac}$-*rydC* plasmids. Expression of the reporter fusion was induced with 0.002% L-arabinose while induction of RydC was achieved with 0.1 mM IPTG. The activities were normalized to vector control and plotted as relative activity. The Specific Activity in Miller units are presented underneath the graph. These experiments were conducted as three independent trials with three biological replicates per trial. Error bars represent standard deviation among biological replicates from a representative trial.
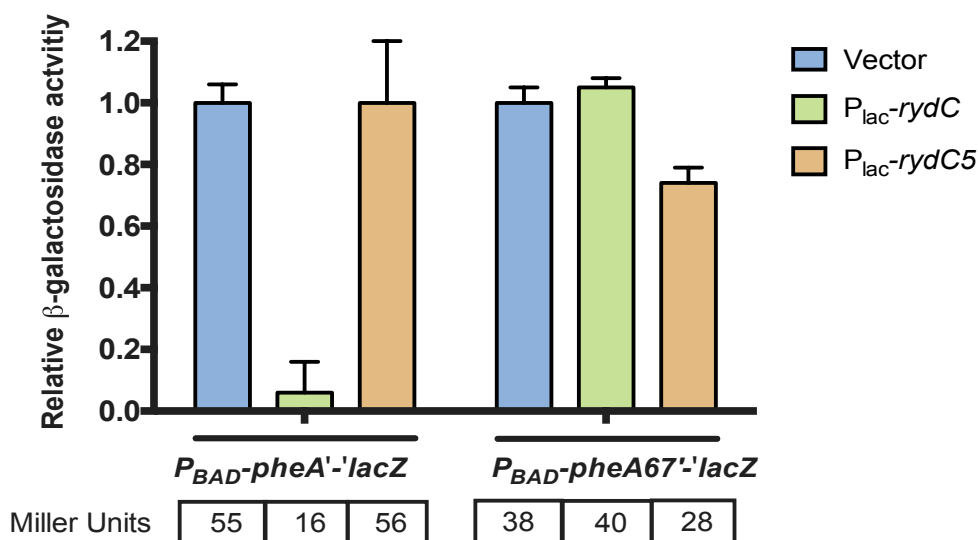(C) Empty vector or RydC was overexpressed in strains with reporter fusion as indicated above. Expression of the fusion and RydC was induced as previously described. As a comparison, the positive control *cfa(long)'* was included in the experiment. These experiments were conducted as three independent trials with three biological replicates per trial. Error bars represent standard deviation among biological replicates from a representative trial.

## Fig. 6

**A.**

```
               4              11
pheA67'        |           CC |
pheA  5' AUGACAUCGGA  3'
               | | | | | | | | |
RydC    3' AUGUAGCCU  5'
RydC5          |           GG |
              10               3
```

**B.**



| Miller Units | 55 | 16 | 56 | 38 | 40 | 28 |

Legend:
- Vector
- P_lac-*rydC*
- P_lac-*rydC5*

X-axis groups: $P_{BAD}$-*pheA'*-*'lacZ*    $P_{BAD}$-*pheA67'*-*'lacZ*
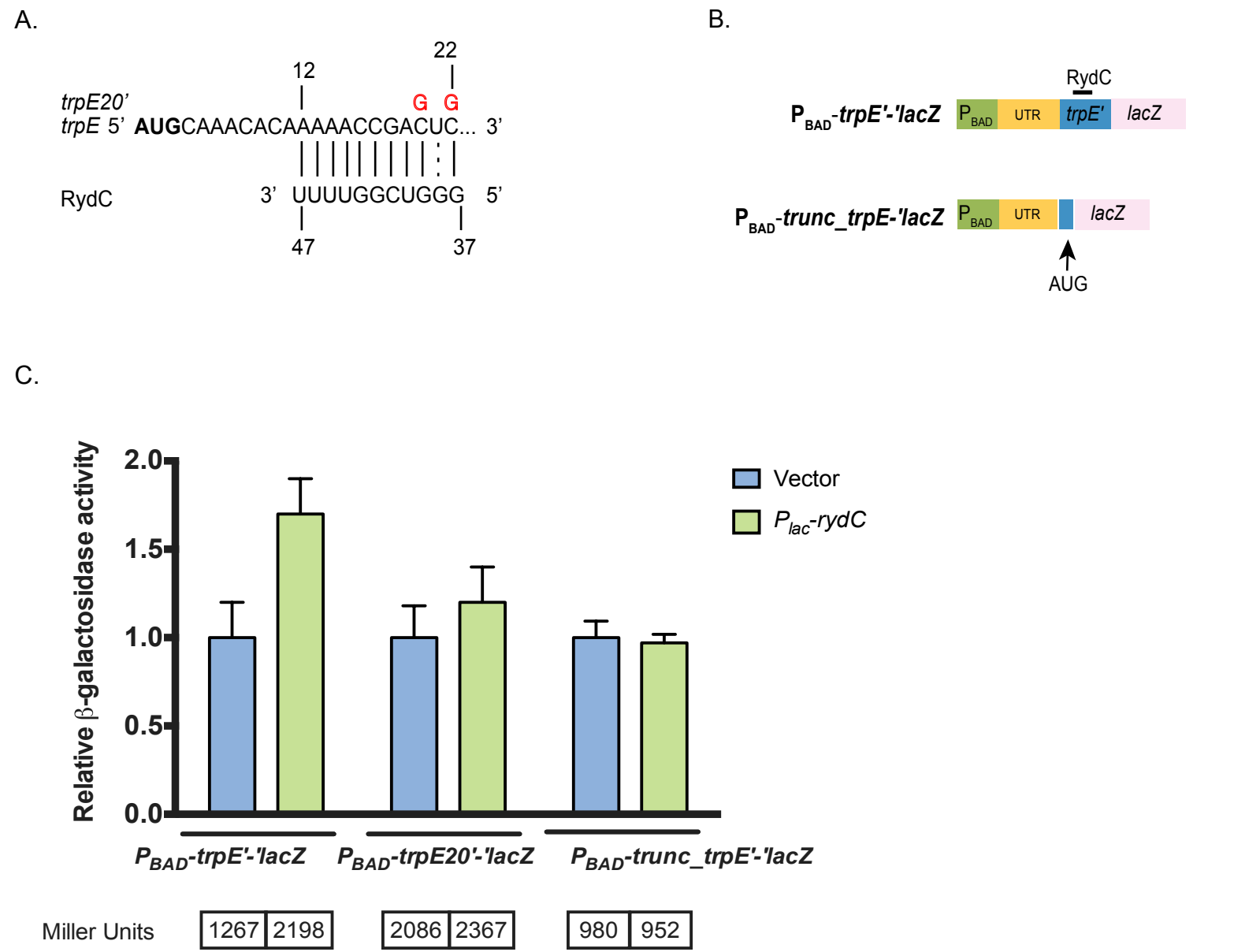
Y-axis: Relative β-galactosidase actvitiy

**Figure. 6. RydC represses *pheA* translation.**
(A) The predicted base pairing between *pheA* mRNA and RydC from IntaRNA. The residues highlighted in red represent point mutations made for each of the variant fusions/RydC alleles. The numbers are in relation to the +1 of RydC and the AUG of *pheA*. To test *pheA* as a putative target, both full-length and mutated (*pheA67'*) *pheA'-'lacZ* translational fusions were tested. (B) RydC or a RydC variant (RydC5) were overexpressed in the *pheA'* and *pheA67' lacZ* fusion backgrounds. Expression of the fusion and RydC was induced as described for Fig. 5. The activities were normalized to vector control and plotted as relative activity. The data were analyzed and reported as described for Fig. 5.

# Fig. 7

## A.



## B.



## C.



| Miller Units | | | | | | |
|---|---|---|---|---|---|---|
| | 1267 | 2198 | 2086 | 2367 | 980 | 952 |

**Figure. 7. RydC activates *trpE* translation.**

(A) The predicted base pairing between *trpE* mRNA and RydC from IntaRNA. The vertical/dotted lines represent the seed region for base pairing interactions. The residues highlighted in red represent point mutations made for each of the variant fusions/RydC alleles. The numbers are in relation to the +1 site of RydC and the 'AUG' of *trpE*. To test *trpE* as a putative target, both full-length, mutated (*trpE20'*), and truncated (*trunc_trpE'*) *trpE'-'lacZ* translational fusions were tested.

(B) Empty vector or RydC were overexpressed in the *trpE'*, *trpE20'*, and *trunc_trpE' lacZ* fusion backgrounds. Expression of the fusion and/or RydC was induced as previously described. The experiments were conducted and data analyzed and presented as described for Fig. 5.