

Finding Instances of Riboswitches and Ribozymes by Homology Search of Structured RNA with Infernal

Amell El Korbi, Jonathan Ouellet, Mohammad Reza Naghdi, and Jonathan Perreault

Abstract

In the genomics era, computational tools are essential to extract information from sequences and annotate them to allow easy access to genes. Fortunately, many of these tools are now part of standard pipelines. As a consequence, a cornucopia of genomic features is available in multiple databases. Nevertheless, as novel genomes are sequenced and new structured RNAs are discovered, homology searches and additional analyses need to be performed. In this chapter, we propose simple ways of finding instances of riboswitches and ribozymes in databases or in unannotated genomes, as well as ways of finding variants that deviate from the typical consensus.

Key words ncRNA, Noncoding RNA, Infernal, Covariation, Homology search, RNA structure, Secondary structure, Riboswitches, Ribozymes

1 Introduction

The diversity of roles attributed to noncoding RNA (ncRNA) has increased at a rapid pace in the last decade. As additional classes of RNAs were discovered and studied, so were their structures [1]. At the same time, sequence databases have grown exponentially, largely due to next-generation sequencing technologies. Public databases such as GenBank [2] or the metagenome-focussed CAMERA [3] database provide incredible opportunities to discover functional structured RNAs with computational screening which have proven extraordinarily useful for many ground-breaking discoveries of new ncRNAs [4–7].

Several computational methods have been used to that end. Some of the most commonly used tools for de novo prediction of ncRNAs include EVOfold [8], QRNA [9], RNAz [10–12], CMfinder [13], Dynalign [14], LocARNA [15], Pfold [16], and the Vienna RNA package [17]. A comprehensive list can also be found in this [wikipedia page](#) [18].

These tools allowed many research groups to find various ncRNAs: small RNAs that base pair on multiple target mRNAs to inhibit gene expression [19], self-cleaving ribozymes [20, 21], and riboswitches that bind a metabolite with their aptamer domain to change gene expression through their expression platform [22].

In parallel to the blossoming field of ncRNAs, the increasing rate of DNA sequencing requires efficient methods to annotate known RNAs. This becomes as important as ever since some of these RNAs are used as antibiotic targets. Indeed, a few riboswitches are already known to be sensitive to natural antibiotics analogous to their ligand [23–27] and new ligand-analogs are being developed in order to overcome the increasing worldwide resistance of bacteria against antibiotics [28–30]. Therefore, finding all instances of a targeted RNA can help determine the sensitive pathogenic strains as well as potentially sensitive beneficial strains.

1.1 Browsing Rfam: The RNA Families Database

In that regard, the collection of RNA families database (Rfam [31]) is particularly useful. The Sanger Institute performs homology searches with the Infernal software suite [32] for all known RNA families, which includes riboswitches and ribozymes, to update the database approximately once a year. Therefore, the quickest way of finding a riboswitch that could be a potential antibiotic target in any organism is by browsing Rfam. The cutoff scores typically used by Rfam to accept a predicted RNA with a relatively complex structure, such as for riboswitches, are high enough that there are almost no false positives in microbial genomes. On this subject, readers might also be interested to look at other recent publications in *Methods in Molecular Biology* [33, 34].

1.2 Search for a Motif in New Genomes with Infernal

However, there are some instances where Rfam does not provide the information needed, particularly if a given riboswitch exists in a newly sequenced genome that has not been screened by Rfam yet. This paper aims to circumvent such problems by presenting a simple step-by-step approach to look at a genome and evaluate the presence of an ncRNA of interest within that unannotated genome. It is targeted towards a general audience with minimal bioinformatics skills, although some basic knowledge of shell command lines would be useful.

1.3 Search for Variants of a Known Motif in All Bacterial Genomes

Occasionally, researchers that are very knowledgeable about a specific riboswitch have reasons to hypothesize that more instances exist. For example, divergence from the structure consensus of a riboswitch class could prevent Infernal from finding such a riboswitch's sub-family. In such cases, it could be desirable to perform a new Infernal search with less stringent criteria to reveal these "hidden" hits [35]. While Rfam lists most of the instances of riboswitches and ribozymes that can easily be found with a relatively high confidence (low E-values), it can occasionally ignore cases

that diverge from the consensus. This has been previously illustrated several times, notably for the *glmS* and the hammerhead ribozymes [21, 35–40]. In the case of the *glmS* ribozyme, Infernal was used with a very high E-value tolerance, as high as 5,000 on all microbial genomes from NCBI's Refseq38 sequence dataset. The resulting hits therefore included a vast majority of spurious hits, but also a number of previously unannotated *glmS* ribozyme instances [35]. Homology searches with very relaxed parameters should not be performed on a routine basis, but rather if there are good hints that additional riboswitches or ribozymes could be found in this manner. Indications on how to manage such searches and the resulting hits will be provided in Subheading 3.3.

2 Materials

Infernal requires the Linux/UNIX system to run the program and is accessible on [Janelia's server](#) [41]. Another alternative is the use of Mac OS X, which is a certified UNIX platform (*see* **Note 1**).

3 Methods

3.1 *Browsing Rfam: The RNA Families Database*

The simplest way to verify the presence of a specific riboswitch in target bacteria is to look in Rfam. As long as the bacterial genomes have been sequenced and annotated by Rfam, browsing the [genomes section](#) [42] would allow anyone to rapidly find which of the known riboswitches are found in a given genome by examining the ncRNAs found in the “chromosomes” tabs. Conversely, browsing the “families” section provides a quick overview of all species that have a specific riboswitch within their genome. This could be especially useful in the context of the development of a new antibiotic to target only a desired group of bacteria and leave most of the natural microbiota intact. Of course, the presence of a riboswitch in a bacterial species does not warrant microbicidal effect of the newly made antibiotic compound. Indeed, studies have already shown potent compounds capable of binding a specific riboswitch to prevent gene regulation via a competition against its native ligand. This competition can affect the growth of some bacteria that have the riboswitch, while leaving others unscathed although the targeted riboswitch is present in both cases. Depending on which genes are regulated by these riboswitches, significant differences of sensitivity can be observed [29].

In the few cases where Rfam would not be useful, any sequence can be screened for the presence of riboswitches with Infernal, which is described in more detail in Subheadings 3.2 and 3.3.

3.2 Search for a Motif in New Genomes with Infernal

Most of what is described herein can be found with additional details in the Infernal user guide [43] and additional papers [34, 44]. The intent here is to provide inexperienced users a rapid start-up guide. For this tutorial, the motif of purine riboswitches is used as an example where Infernal builds a covariance model from an alignment with structural annotation.

The latest version of Infernal is available to download [here](#) [41]. At the time of writing, the latest release of Infernal is 1.0.2 (30 Oct 2009) [45]. Once the source file downloaded, expand the “tar file” at a convenient location. For a basic installation, execute the two commands “configure” and “make” from the “infernal-1.0.2” directory (*see Note 2*):

```
# ./configure
# make
```

To run the optional testsuite, execute the following command:

```
# make check
```

Once Infernal installed, the ncRNA covariance model is built. The first step is to generate the “purine” seed alignment in Stockholm format from Rfam at this [location](#) [46]. Once the alignment is generated, the file should be downloaded and saved under “purine.sto” (*see Note 3*) in the “infernal-1.0.2” directory. The Stockholm format describes the secondary structure of an RNA sequence alignment. Base pairs are annotated as “<” (for the opening base of the pair) and “>” (for the closing base). Other, base pair annotations such as “(”, “[”, “{”, “|”, “}”, “]”, “)” are also used sometimes for base pairs of stems enclosing a multistem junction. Single stranded regions are annotated with other characters, typically “.”, but sometimes “_” for loops and “,” for junctions. A similar notation is used in Infernal’s output for a regular “cmsearch.” For simplicity, we assume that all the following commands are executed from the infernal directory.

Build the “cm file” (covariance model) using the command “cmbuild”:

```
# src/cmbuild purine.cm purine.sto
```

Execute the command “cmcalibrate” which may take more than 1 h:

```
# src/cmcalibrate purine.cm
```

To search for the presence of that purine motif in a new genome, copy the sequence file of the genome of interest (FASTA format) to the “infernal-1.0.2” directory and use:

```
# src/cmsearch purine.cm genome.fa
```

For the purpose of this example, the search will be performed in a known bacterial genome downloaded from NCBI. The genome sequence of *Bacillus subtilis* (in FASTA format) is downloaded from this [link](#) [47]. The file is downloaded by pressing the

to 1, manual inspection can provide the additional clues needed to confirm the presence of a riboswitch at that position. For example, a loop–loop base-pairing interaction forms between the two loops of the purine riboswitch (here, the loops are annotated with “_”). This feature is not evaluated by Infernal and therefore does not contribute to the E-value. Observing this interaction in a relatively poor hit, with an E-value of 0.2 for instance, would mean this hit is more likely to be a true riboswitch than suggested strictly by the E-value. For purine riboswitches, detailed knowledge of the riboswitch is also useful to discriminate between adenine and guanine riboswitches. Because these two differ by a single base, Infernal finds both types of riboswitches during the same search. In the output shown above, the first four hits are guanine riboswitches and the last one is an adenine riboswitch. The former have a “C” at the last base of the junction, while the adenine riboswitch has a “U” at that position (shown in bold in the alignments above and annotated with “;”).

3.3 Search for Variants of a Known Motif in All Bacterial Genomes

The Infernal suite can be used to find atypical riboswitches or ribozymes, but if many genomes are evaluated for poor E-values, thousands of hits will be generated and will require a lot of CPU time (*see Note 1*). Afterwards, knowing the structure of the RNA in detail can help to sort through the haystack of hits that would ensue a search accepting E-values as high as 5,000. Evaluating the presence of pseudoknots or the relevance of the downstream gene being regulated by this RNA are examples of how one can judge whether hits are likely real ncRNAs. This entire process takes much more time than what is described in Subheading 3.2 and is not recommended for all homology searches. This approach is more feasible in a case where only a few genomes are to be scrutinized, although the E-value should be set closer to 1 since it corresponds to the number of false positives expected to be found with that score (or better) in a database of this size. Therefore, a “cmsearch” allowing a maximum E-value of 5 (1 is default) could be performed and the hits could be manually screened one by one with the criteria mentioned hereafter when specifically interested in the genome of one bacteria. However, before performing such a search on all microbial genomes (with a maximal E-value of 1,000 for instance), one should have a strong basis to believe more riboswitches or ribozymes can be found since it would be CPU-time intensive and would generate a lot of false positives requiring a lot of time to sort through.

The steps described in Subheading 3.2 are also valid for searches in all available genomes. However, to get hits with E-values as high as 100 (for example), the “-E” option with “cmsearch” is used (*see Note 4* for more information on “--tabfile”).

```
# cmsearch -E 100 --tabfile results.tab purine.cm
sequence.fa
```

Where “sequence.fa” could be a large file with all microbial genomes (available from NCBI [48], note that this compressed file is approximately 2 Gb in size). Some adjustments can be useful to determine the best value for “-E” (*see* **Note 5**). Different softwares can help visualizing a large number of hits. In that regard, the “RALEE” major mode in “Emacs” (a text editor program) is very useful [49]. It can color Stockholm alignments according to conservation, stems, or covariation. Both “RALEE” [50] and Emacs [51] are available for any Operating System platform.

When sorting the hits, as many criteria as possible should ideally be used to distinguish potentially good hits from spurious ones. Here are a few noticeable features, that have already proven useful in other works [4, 5, 35]:

1. Pseudoknots: Infernal does not take the base pairs of pseudoknot in consideration. Thus, it cannot account for pseudoknots in its E-value, which means that manually confirming the presence of a known pseudoknot in the ncRNA greatly improves this hit’s likelihood of being real.
2. Essential bases: when the structure has been studied enough to determine which bases are absolutely crucial for the RNA’s function, the hits that do not have these bases can be considered as spurious. However, one must be careful with such criteria since an apparent deleterious mutation at a specific position could be compensated by different bases at other positions, as in the case of the core-conserved C3G8 base-pair within the hammerhead ribozyme core, which was sometime found to be U3A8 [38, 52].
3. Intergenic versus coding sequence (CDS): even though riboswitches could theoretically be found in coding sequences, to our knowledge there is no natural riboswitch found yet that is completely embedded in the coding sequence. Therefore, if the hit is in an intergenic region, it should be regarded as more likely to be real, and, conversely, as spurious if it is in a CDS.
4. Functional relevance: when the riboswitch’s ligand is known, the connection with the genes is often obvious. For example, in the above list of purine riboswitches, a hypoxanthine/guanine permease can be found downstream of a hit, as well as other genes involved in purine synthesis for other hits. The absence of a clear connection between the candidate riboswitch and the function of the downstream gene does not automatically means a false positive, but an obvious connection does help for its validation.
5. Expression platforms: to exert their effect on expression, the aptamer portion of riboswitches is usually close (or even

interest can easily be copied and pasted as a text file for which the extension has to be changed to “.cm” before using it.

4. The --tabfile option allows the use of the tab file as an input to the “Easel” miniapp “esl-sfetch” (found in a subdirectory of Infernal). The miniapp “esl-sfetch” extracts the sequences of all hits found from the genome sequence file to a new FASTA file. This file is useful to get a new alignment with the CM file of the motif (“purine.cm” in the example) using the command “cmalign.” To get a tabular version of the search results, the command line is:

```
# cmsearch --ga --tabfile results.tab purine.cm sequence_Bsubs.fa
```

Now, to use the tabfile “results.tab” as an input to fetch the hits sequences:

```
# easel/miniapps/esl-sfetch -C -f --tabfile sequence_Bsubs.fa results.tab
```

An error of this type: “Failed to open SSI index” may occur. In this case the “sequence_Bsubs.fa” (or the file containing the new genome) has to be indexed. This is done with this step:

```
# easel/miniapps/esl-sfetch --index sequence_Bsubs.fa
```

Now that the file containing the genome sequence is indexed, the “sfetch” command can be re-executed as above. The hits sequences are displayed in FASTA format. To get the output in a new FASTA file, the command line would be:

```
# easel/miniapps/esl-sfetch -C -f --tabfile sequence_Bsubs.fa results.tab>hitsequences.fa
```

The tabular version has the format shown below.

#	model name	target name	start	stop	start	stop	bit sc	E-value	GC%
#	-----	-----	-----	-----	-----	-----	-----	-----	----
	Purine	gi 223666304	697666	697767	1	102	85.74	1.74e-18	42
	Purine	gi 223666304	693731	693832	1	102	85.00	2.72e-18	38
	Purine	gi 223666304	4004455	4004556	1	102	73.67	2.54e-15	28
	Purine	gi 223666304	2319369	2319270	1	102	82.19	1.48e-17	46
	Purine	gi 223666304	625950	625851	1	102	65.20	4.19e-13	30

Shown below is the beginning of the file “hitsequences.fa” containing the sequences of the hits found in the *Bacillus subtilis* genome:

```
>gi|223666304|ref|NZ_CM000487.1|/697666-697767/Purine/B85.74/E1.7e-18/GC42 Bacillus subtilis subsp. subtilis str. 168 chromosome, whole genome shotgun sequence
CATGAAATCAAACACGACCTCATATAATCTTGGAATATGGCCATAAGTTTCTACCCG
GCAACCGTAAATTGCCGACTATGCAGGAAAGTGATCGATAA
>gi|223666304|ref|NZ_CM000487.1|/693731-693832/Purine/B85.00/E2.7e-18/GC38 Bacillus subtilis subsp. subtilis str. 168 chromosome, whole genome shotgun sequence
AGAAATCAAATAAGATGAATTCGTATAATCGCGGAATATGGCTCGCAAGTCTCTACCAA
GCTACCGTAAATGGCTTGACTACGTAAACATTTCTTTCTGTTT
...
```

These sequences are aligned using the “purine.cm” as a seed to obtain a new motif that can be used in further searches:

find instances that have been missed because their pseudoknot's sequence diverges from the current model.

Acknowledgements

This work was supported by the Natural Sciences and Engineering Research Council of Canada (NSERC) discovery grant (RGPIN 418240-2012) and by a grant from The Banting Research Foundation to JP.

References

1. Wan Y, Kertesz M, Spitale RC, Segal E, Chang HY (2011) *Nat Rev Genet* 12:641–655
2. Benson DA, Karsch-Mizrachi I, Clark K, Lipman DJ, Ostell J, Sayers EW (2012) *Nucleic Acids Res* 40:D48–D53
3. Sun S, Chen J, Li W, Altintas I, Lin A, Peltier S, Stocks K, Allen EE, Ellisman M, Grethe J, Wooley J (2011) *Nucleic Acids Res* 39:D546–D551
4. Weinberg Z, Wang JX, Bogue J, Yang J, Corbino K, Moy RH, Breaker RR (2010) *Genome Biol* 11:R31
5. Weinberg Z, Perreault J, Meyer MM, Breaker RR (2009) *Nature* 462:656–659
6. Shi Y, Tyson GW, DeLong EF (2009) *Nature* 459:266–269
7. Livny J, Waldor MK (2007) *Curr Opin Microbiol* 10:96–101
8. Pedersen JS, Bejerano G, Siepel A, Rosenbloom K, Lindblad-Toh K, Lander ES, Kent J, Miller W, Haussler D (2006) *PLoS Comput Biol* 2:e33
9. Rivas E, Eddy SR (2001) *BMC Bioinformatics* 2:8
10. Washietl S (2007) *Methods Mol Biol* 395:503–526
11. Gruber AR, Neubock R, Hofacker IL, Washietl S (2007) *Nucleic Acids Res* 35:W335–W338
12. Washietl S, Hofacker IL, Stadler PF (2005) *Proc Natl Acad Sci U S A* 102:2454–2459
13. Yao Z, Weinberg Z, Ruzzo WL (2006) *Bioinformatics* 22:445–452
14. Harmanci AO, Sharma G, Mathews DH (2007) *BMC Bioinformatics* 8:130
15. Will S, Reiche K, Hofacker IL, Stadler PF, Backofen R (2007) *PLoS Comput Biol* 3:e65
16. Knudsen B, Hein J (2003) *Nucleic Acids Res* 31:3423–3428
17. Lorenz R, Bernhart SH, Honer Zu Siederdisen C, Tafer H, Flamm C, Stadler PF, Hofacker IL (2011) *Algorithms Mol Biol* 6:26
18. http://en.wikipedia.org/wiki/List_of_RNA_structure_prediction_software
19. Storz G, Vogel J, Wassarman KM (2011) *Mol Cell* 43:880–891
20. Cochrane JC, Strobel SA (2008) *Acc Chem Res* 41:1027–1035
21. Hammann C, Luptak A, Perreault J, de la Pena M (2012) *RNA* 18:871–885
22. Roth A, Breaker RR (2009) *Annu Rev Biochem* 78:305–334
23. Lee ER, Blount KF, Breaker RR (2009) *RNA Biol* 6:187–194
24. Blount KF, Wang JX, Lim J, Sudarsan N, Breaker RR (2007) *Nat Chem Biol* 3:44–49
25. Blount KF, Breaker RR (2006) *Nat Biotechnol* 24:1558–1564
26. Sudarsan N, Cohen-Chalamish S, Nakamura S, Emilsson GM, Breaker RR (2005) *Chem Biol* 12:1325–1335
27. Ott E, Stolz J, Lehmann M, Mack M (2009) *RNA Biol* 6:276–280
28. Kim JN, Blount KF, Puskarczyk I, Lim J, Link KH, Breaker RR (2009) *ACS Chem Biol* 4:915–927
29. Mulhbach J, Brouillette E, Allard M, Fortier LC, Malouin F, Lafontaine DA (2010) *PLoS Pathog* 6:e1000865
30. Lunse CE, Schmidt MS, Wittmann V, Mayer G (2011) *ACS Chem Biol* 6:675–678
31. <http://rfam.sanger.ac.uk>
32. Nawrocki EP, Kolbe DL, Eddy SR (2009) *Bioinformatics* 25:1335–1337
33. Hoepfner MP, Barquist L, Gardner PP. *Methods in Molecular Biology* (In press)
34. Barquist L, Burge SW, Gardner PP. *Methods in Molecular Biology* (In press)
35. McCown PJ, Roth A, Breaker RR (2011) *RNA* 17:728–736

36. de la Pena M, Garcia-Robles I (2010) *EMBO Rep* 11:711–716
37. Jimenez RM, Delwart E, Luptak A (2011) *J Biol Chem* 286:7737–7743
38. Perreault J, Weinberg Z, Roth A, Popescu O, Chartrand P, Ferbeyre G, Breaker RR (2011) *PLoS Comput Biol* 7:e1002031
39. Seehafer C, Kalweit A, Steger G, Graf S, Hammann C (2011) *RNA* 17:21–26
40. de la Pena M, Garcia-Robles I (2010) *RNA* 16:1943–1950
41. <http://infernal.janelia.org/>
42. <http://rfam.sanger.ac.uk/genome/browse#A>
43. Nawrocki EP, Kolbe DL, Eddy SD (2009) <ftp://selab.janelia.org/pub/software/infernal/Userguide.pdf>
44. Nawrocki EP. *Methods in Molecular Biology* (In press)
45. <ftp://selab.janelia.org/pub/software/infernal/infernal-1.0.2.tar.gz>
46. <http://rfam.sanger.ac.uk/family/RF00167#tabview=tab2>
47. <http://www.ncbi.nlm.nih.gov/nuccore/223666304>
48. <ftp://ftp.ncbi.nih.gov/genomes/Bacteria/all.fna.tar.gz>
49. Griffiths-Jones S (2005) *Bioinformatics* 21:257–259
50. <http://personalpages.manchester.ac.uk/staff/sam.griffiths-jones/software/ralee/>
51. <http://www.gnu.org/software/emacs>
52. Przybilski R, Hammann C (2007) *RNA* 13:1625–1630
53. Barrick JE, Breaker RR (2007) *Genome Biol* 8:R239
54. Kim JN, Roth A, Breaker RR (2007) *Proc Natl Acad Sci U S A* 104:16092–16097
55. <http://www.ncbi.nlm.nih.gov/nuccore/159184118?report=fasta>
56. <ftp://ftp.sanger.ac.uk/pub/databases/Rfam/CURRENT/Rfam.cm.gz>
57. <http://github.com/ppgardne/RNIE>
58. Gardner PP, Barquist L, Bateman A, Nawrocki EP, Weinberg Z (2011) *Nucleic Acids Res* 39:5845–5852