

IDENTIFICATION OF POTENTIAL RIBOSWITCH ELEMENTS IN *Homo Sapiens* mRNA 5'UTR SEQUENCES USING POSITIVE-UNLABELED MACHINE LEARNING

William S. Raymond¹, Jacob DeRoo¹, and Brian Munsky^{1,2}

¹School of Biomedical Engineering, Colorado State University Fort Collins, CO 80523, USA

²Chemical and Biological Engineering, Colorado State University Fort Collins, CO 80523, USA

November 16, 2023

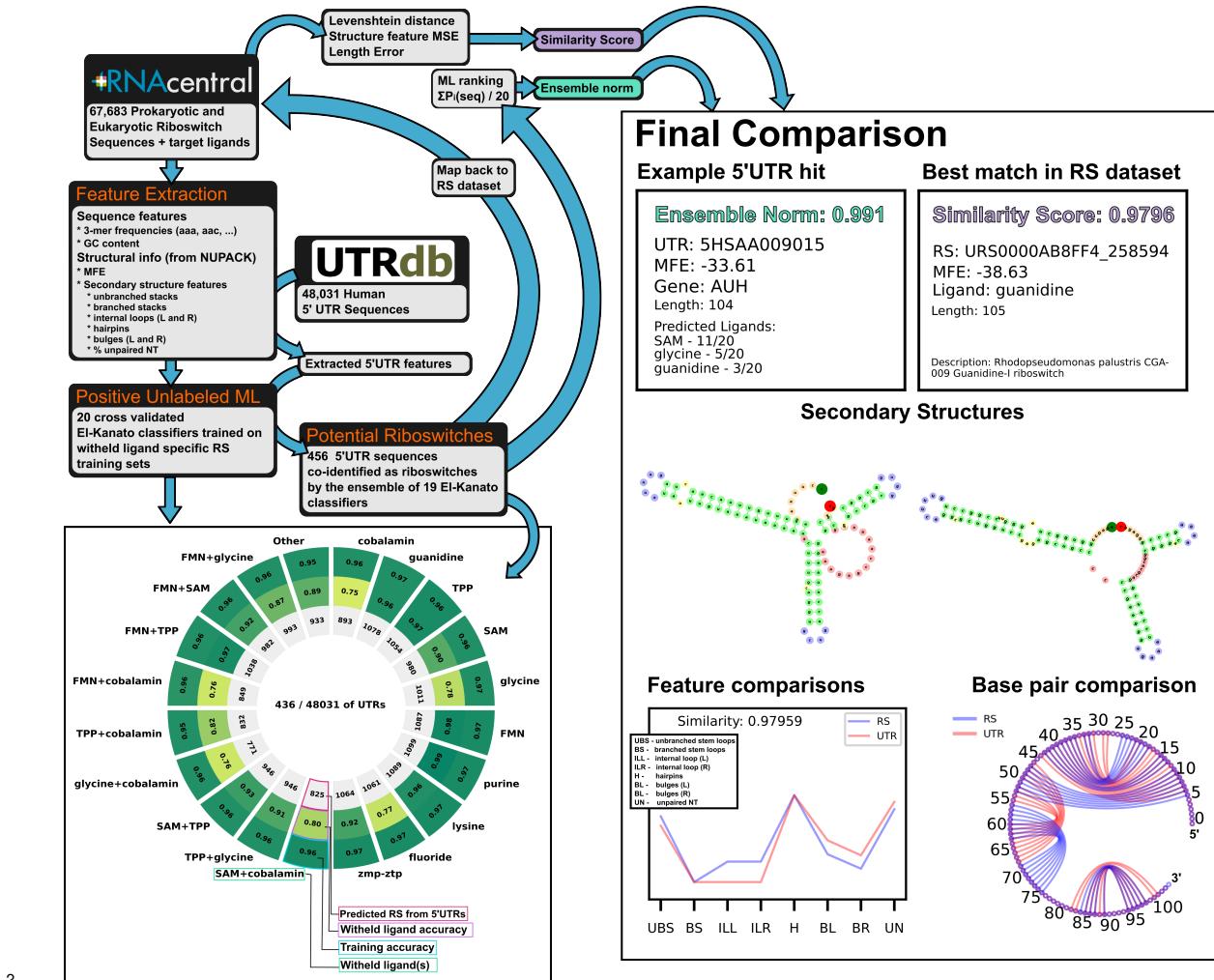
ABSTRACT

Riboswitches are a class of noncoding RNA structures that interact with target ligands to cause a conformational change that can then execute some regulatory purpose within the cell. Riboswitches are ubiquitous and well characterized in bacteria and prokaryotes, with additional examples also being found in fungi, plants, and yeast. To date, no purely RNA-small molecule riboswitch has been discovered in *Homo Sapiens*. Several analogous riboswitch-like mechanisms have been described within the *H. Sapiens* translatome within the past decade, prompting the question: Is there a *H. Sapiens* riboswitch dependent on only small molecule ligands? In this work, we set out to train positive unlabeled machine learning classifiers on known riboswitch sequences and apply the classifiers to *H. Sapiens* mRNA 5' 5'UTR sequences found in the 5'UTR database, UTRdb, in the hope of identifying a set of mRNAs to investigate for riboswitch functionality. 67,683 riboswitch sequences were obtained from RNACentral and sorted for ligand type and used as positive examples and 48,031 5'UTR sequences were used as unlabeled, unknown examples. Positive examples were sorted by ligand, and 20 positive-unlabeled classifiers were trained on sequence and secondary structure features while withholding one or two ligand classes. Cross validation was then performed on the withheld ligand sets to obtain a validation accuracy range of 75%-99%. The joint sets of 5'UTRs identified as potential riboswitches by the 20 classifiers were then analyzed. 1515 sequences were identified as a riboswitch by one or more classifier(s) and 436 of the *H. Sapiens* 5'UTRs were labeled as harboring potential riboswitch elements by all 20 classifiers. These 436 sequences were mapped back to the most similar riboswitches within the positive data and examined. An online database of identified and ranked 5'UTRs, their features, and their most similar matches to known riboswitches, is provided to guide future experimental efforts to identify *H. Sapiens* riboswitches.

Author summary

Riboswitches are an important regulatory element in bacteria that have not been described in *Homo Sapiens*. However, if human riboswitches exist and if they can be found, they could have vast implications on human disease. We apply positive-unlabeled machine learning to combine known riboswitch sequences with *H. Sapiens* 5'UTR sequences and to search for potential riboswitches. We analyze our ensemble predictions for likely *H. Sapiens* 5'UTR riboswitches using GO analysis to determine their potential functional roles, and we rank and display our predicted sequences next to the most similar known riboswitches. We expect these analyses to be helpful to the scientific community in planning future experiments for laboratory discovery and validation.

1 0.1 Graphical Abstract



2

3 Introduction

4 A riboswitch (RS) is a non-coding RNA sequence harboring a structure with two distinct conformations. Confor-
5 mational changes are induced when an aptamer region interacts with a target small molecule, revealing or occluding
6 functional parts of the RNA. This inducible change in structure allows for broad regulation of various cellular pro-
7 cesses via modification of protein production, mainly by affecting transcription termination/continuation, translation
8 inhibition/activation, mRNA splicing, or mRNA stability (50; 51; 35; 33). Whether a particular riboswitch acts in a
9 positive or negative regulatory manner when in contact with its ligand strongly depends on the expression platform and
10 aptamer location in relation to other elements, such as the ribosomal binding site. These regulatory effects could con-
11 ceivably have disease-related implications due to under- or over-expression of a regulatory biomolecule or to a genetic
12 mutation of the riboswitch in question. The current set of described riboswitches is predicted to be a small subset of all
13 existing riboswitches – leading to open questions such as to “which riboswitch classes are uncharacterized?” or “why

14 do some life-essential molecules lack known riboswitch aptamers?” Riboswitch discovery has been an active area of
15 research since their first description in 2002, with many computational and experimental efforts undertaken to elucidate
16 new riboswitch classes (26). Riboswitches occur ubiquitously in prokaryotes, where they enjoy a rich diversity of
17 around 40 molecular targets (35). In contrast, nearly every example of an eukaryotic riboswitches in the current literature
18 was found in lower eukaryotes, such as mold, yeast, and fungi, and they leverage thiamine pyrophosphate (TPP)
19 as their target ligand (51; 55; 26). Among higher eukaryotes, several species of plants have a single TPP riboswitch
20 in the 3'UTR of the conserved gene THIC – the riboswitch acts to regulate gene expression by creating an unstable
21 mRNA product in the presence of TPP (8). Interestingly, multiple analogous “pseudo-riboswitches” (riboswitch like
22 elements that are stabilized by proteins and small molecules) have been located within the untranslated region (UTR)
23 of human (*Homo Sapiens*, *H. Sapiens*) translome (49; 53). The existence of some analogous mechanisms in *H.*
24 *Sapiens* and a described higher eukaryotic riboswitch prompts the question: “do riboswitches have an unknown niche
25 in all higher eukaryotes, or are they simply missing?” Indeed this question is one of the largest open questions in the
26 field today (26). From a disease perspective, do riboswitches exist within in *H. Sapiens*, and if so, where and with
27 what implications on human health?

28 Machine learning is routinely used in Bioinformatics for a wide range of RNA related tasks, including parsing
29 RNA-seq data, performing RNA secondary structure prediction, and providing discovery based approaches for RNA
30 sequences, splice sites, and genome wide functional RNA elements (45; 3). For the task of computational riboswitch
31 prediction, previous methods leverage hidden Markov models (HMMER, RiboSW, Riboswitch Scanner (60; 9; 40)),
32 covariance models and context free grammar (Infernal (42)), and sequence alignment + computational folding (Riboswitch
33 Finder, RibEx, RNAConSLOpt (7; 1; 30)). Other more recent software such as Riboflow utilizes deep learning
34 classifiers such as RNN-LSTM or convolutional neural networks (CNN) for their riboswitch identification (48).
35 Computational methods available to the public up until 2018 are reviewed in Antunes et al. extensively (4). Notably
36 many of these have difficulty extrapolating to unknown riboswitches and rely heavily on a previous knowledge base.
37 Recent breakthroughs have been achieved via reverse homology searching approaches (mutate a sequence without
38 disturbing secondary structure, search for the new mutant in a genome wide fashion), which recently helped to identify
39 a list of potential purine riboswitches in fungi (39) – however this approach once again suffers from a lack of
40 extrapolation and requires a known starting point structure.

41 Positive unlabeled learning (PU-learning) is a subclass of binary machine learning classification that attempts to
42 learn from data that only contain positive and unknown examples. In other words, they are used to classify data
43 where the labels of one class are known (label 1) or known and incorrect (labeled 1, truly 0), and the other class
44 are unknown and could either be true examples (label 1), unclassified examples (label 0.5), or negative examples
45 (label 0) (14; 15; 16; 6). Situations producing unlabeled-positive data sets are prevalent in fields such as medical
46 diagnosis (e.g., people with a diagnosis vs. people without a condition vs. people with a condition and no diagnostic
47 confirmation), interest-based applications (e.g., people who engaged with an ad vs. those who did not, since not
48 engaging could be a negative or a neutral reaction), and biology (e.g., a class of known proteins vs. proteins with

49 unclassified but similar function vs. proteins without the same function). PU-learning is routinely applied in molecular
50 biological discovery applications since the advent of big data approaches (63; 61; 24); Proteomics, RNA-seq, or whole
51 genome sequencing quantify virtually all species within a sample whether or not the molecules are characterized,
52 creating a tranche of unlabeled data along with its positive examples (27). Within the context of RNA, PU-learning has
53 also been used to identify non-coding RNA genes (57), predict circularRNA and piRNA disease associations (64; 2),
54 to predict RNA secondary structures (52), and to classify metastasis potential from cancer cell RNA-seq data (65).

55 In this work, we set out to use PU-learning to identify a group of potential sequences in the *H. Sapiens* mRNA
56 5'UTR that may contain riboswitches – with the hope of providing a first-pass reduced list for future, targeted labora-
57 tory investigations. 67,683 sequences tagged with “riboswitch” from the non-coding RNA database, RNACentral, were
58 used as positive examples. 48,031 *H. Sapiens* 5'UTR sequences were obtained from the untranslated region (UTR)
59 database, UTRdb, and used as unlabeled examples. Sequences were sanitized and structural- and sequence-based fea-
60 tures were extracted. 20 PU-classifiers were trained on the RS-5'UTR feature sets and validated on single or double
61 holdouts of specific RS ligand classes. The resulting ensemble of classifiers was then examined for the overlap of
62 5'UTR sequences that were considered as riboswitches (positively labeled). 436 sequences were found to be potential
63 5'UTR hits across all 20 classifiers. These positively labeled 5'UTR hits were then compared with their most similar
64 sequences within the riboswitch data set via metrics comparing length, dot structure differences, and structural feature
65 similarities, and all results are presented in an interactive display website. GO analysis was also preformed to examine
66 fold enrichment of cellular processes and functions. Further verification of the classifier ensemble was performed by
67 applying the classifier to a set of 25 synthetic riboswitches, of which 56% were correctly discovered as riboswitches
68 despite having no representation of similar synthetic riboswitches in any training data. Using our computationally
69 validated ensemble, we provide a minimal list of *H. Sapiens* 5'UTRs that appear most likely to harbor riboswitch
70 sequences in hopes that these hits could be corroborated with future experimental validation.

71 Results and Discussion

72 **67,683 known riboswitch sequences and 43,081 *H. Sapiens* 5'UTRs were collected and sanitized for subsequent 73 PU classification**

74 Two RNA databases were selected for training data: RNACentral and UTRdb. RNACentral is a meta collection of
75 many databases of all types of non-coding RNA, and was utilized as the source of riboswitch sequences (46). JSON
76 information of all entries containing the tag “riboswitch” were queried from RNACentral on 8.19.22 and filtered to
77 remove duplicate sequences. Ligands were parsed from the entry descriptions and any missing ligands were obtained
78 from the corresponding entry’s RFAM data (25). Ligands were further filtered to combine names referring to the same
79 ligand (e.g. “mn”, “manganese”, “Mn2+” all renamed to “Mn2+”). Cobalamin sub-types such as Adenosylcobalamin
80 were combined under the umbrella of “cobalamin” for ligand labelling. Any protein specific ligand was renamed
81 to “protein” (1 total) and all tRNA ligands were lumped to “tRNA” (3 total). Speculative or synthetic riboswitches

82 (nhA-I motif, duf1646, raiA, synthetic, sul1, blank) were relabeled with ‘unknown’ as their ligand (1130 total). After
 83 ligand relabeling, 73,119 riboswitch sequences remained in the data set. After removing identical sequences, 67,683
 84 penultimate riboswitch sequences were stored for machine learning. Riboswitches targeting cobalamin(s), TPP, S-
 85 Adenosyl methionine (SAM), glycine, FMN, purine, lysine, fluoride, and guanidine made up 82% of the riboswitch
 86 data set. Other ligand labels such as unknown, molybdenum, GMP, or nickel/cobalt made up less than 2% of the data
 87 set each (Figure 1A). A full list of ligands represented in the data set can be found in Table 0.1.

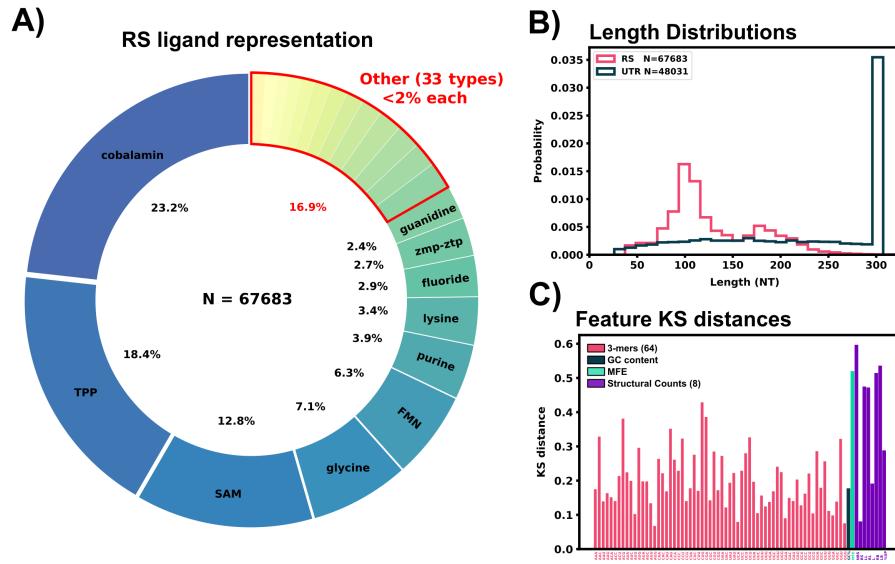


Figure 1: Riboswitch ligand representation, training data comparisons, and feature extraction of sequence data.
 A: Ligand representation within the Riboswitch training data. 43 different ligands are represented with 10 ligands having greater than 2% representation in the data set. B: Length distributions for the sanitized 5'UTR and RS data set. KS distances between the 5'UTR and RS data set for all extracted features are shown in the bottom panel.

88 *H. Sapiens* 5'UTR sequences were pulled from UTRdb, a UTR database of multiple organisms' mRNA (20) on
 89 May 2022, and all analyses and computations use these data. UTRdb has since updated the original database to add
 90 additional curated functional annotations and UTRs (19). This update does not substantially affect the sequence data
 91 and is not expected to affect any results presented here. For the readers' convenience, we reproduce the and make these
 92 data available through the GitHub repository, https://github.com/Will-Raymond/human_riboswitch_hits.
 93 Sequences were filtered for identical sequences and stored in a data set. 5'UTR sequences were matched with their
 94 corresponding coding regions from the *H. Sapiens* consensus coding region (CCDS) release 22 (accessed 11.28.2021).
 95 5'UTR sequences missing CCDS information were discarded from the data set. 5'UTR sequences were appended with
 96 22 nucleotides downstream from the start codon of the mRNA, and trimmed to the last 300 nucleotides in the 5' to
 97 3' direction if the full 5'UTR sequence was over that limit. This min(5'cap, 275NT) to start codon to 22 NT region
 98 was selected as the area to search for potential riboswitches as regulatory conformational in this region changes could
 99 directly block or expose the ribosomal initiation site. After the sanitation, CCDS matching, and length trimming,

100 48,031 5'UTR sequences were stored for subsequent examination. Figure 1B shows the length distribution of both
101 data sets, and Figure 2A shows an example of a 5'UTR + 25 NT sequence.

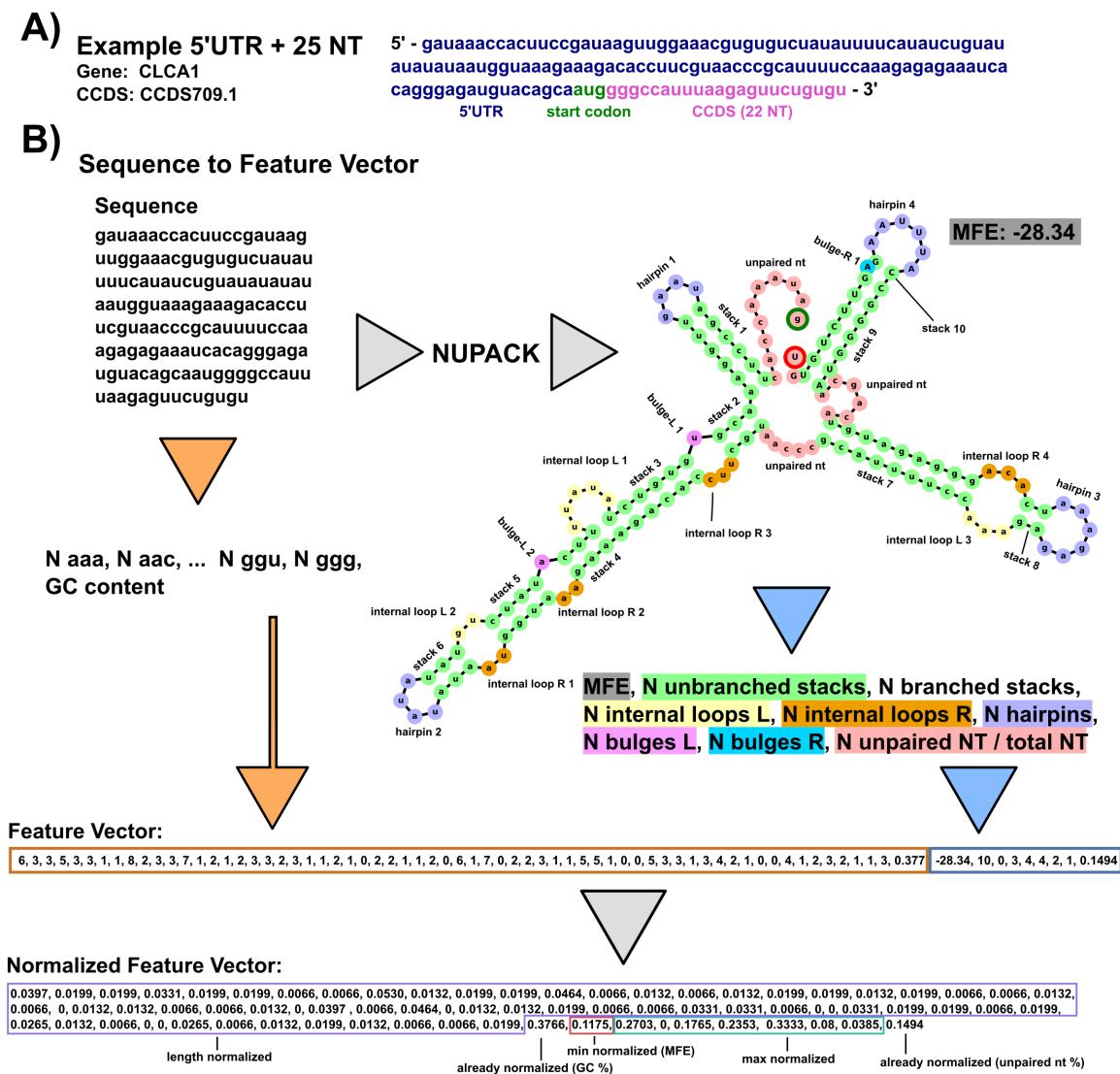


Figure 2: Example feature extraction of sequence data. A: Example 5'UTR sequence from the data set containing the start codon and 22 downstream nucleotides (25 total). B: Annotated example of taking an RNA sequence and converting it to a normalized feature vector for our positive-unlabeled learning. For sequence-based features, the sequence is converted into a 3-mer frequency and GC content is calculated. 3-mer frequency is normalized by the number of 3-mer subsets in the sequence (sequence length - 2). Secondary structure based features are generated by passing the sequence through NUPACK. MFE and structural features are extracted from the dot structure. Counts of hairpins, internal loops, bulges, and contiguous stacks (with and without branches) are extracted and max normalized across all the entire data set. Left (L) and right (R) designation corresponds to the 5' to 3' direction and 5' to 3' direction within a base pair stack respectively. MFEs are min-normalized across the data set. The final structural feature considered for learning is the percentage of unpaired nucleotides in the structure. The final output is a vector of length 74 normalized from 0-1.

102 For both the known riboswitch data set and the *H. Sapiens* 5'UTR sequences data set, sequences with incomplete
103 or multiple base pair specifications were renamed to the first matching base pair out of the order A, C, G, T/U for the
104 unknown character according to IUPAC naming conventions, Table 0.1 (13).

105 **74 structure and sequence based features were extracted from the 5'UTR and RS data sets**

106 In an effort to collect a broad spectrum of information for machine learning purposes, each sequence was processed to
107 quantify 74 features in two groups: 65 sequence features and 9 predicted structural features.

108 The 65 sequence features include 64 length-normalized 3-mer (AAA, AAG, AAU ... CCC) frequencies, and the
109 GC content. To define these, the 4^k sequence k -mers were generated for each transcript, and the resulting 64-element
110 vector was normalized by the total number of k-mers (i.e., length - 2) of the corresponding transcript (12; 47):

$$[S_1, \dots, S_{64}] = \left[\frac{N_{aaa}}{L_{seq} - 2}, \frac{N_{aac}}{L_{seq} - 2}, \frac{N_{aau}}{L_{seq} - 2}, \frac{N_{aag}}{L_{seq} - 2}, \frac{N_{aca}}{L_{seq} - 2}, \dots, \frac{N_{ggu}}{L_{seq} - 2}, \frac{N_{ggg}}{L_{seq} - 2} \right] \quad (1)$$

111 In addition, the GC content is defined as the count of G and C within the sequence normalized by the sequence length:

$$S_{65} = \frac{N_g + N_c}{L_{seq}} \quad (2)$$

Structural features were obtained by passing each sequence through the computational folding algorithm, NUPACK 4.0.0.23, (18; 17) to obtain a minimum free energy (MFE) secondary structure. NUPACK allows the user to specify a number of RNA strands within one set of “test tube” conditions and provides a list of commonly solved secondary structures and their mean free energies using a computational RNA model. Default NUPACK model settings were used when folding all sequences, “Model(material='rna', ensemble='stacking', celsius=37, sodium=1.0, magnesium=0.0).” For each sequence, the dot structure and the MFE of the most commonly folded non-complexed (no A:A) structure out of 100 RNA strands was saved and recorded as a sequence’s secondary structure for feature extraction. The NUPACK MFE value, unpaired base pair percentage, and counts of how many consecutive stem base pairs in a branching stem or non-branching stem, number of stem loops, number of internal loops left and right, and numbers of left and right bulges were extracted and used as a “structural feature vector” defined as:

$$[S_{66}, \dots, S_{74}] = [MFE, N_{\text{unbranched stacks}}, N_{\text{branched stacks}}, N_{\text{loops L}}, N_{\text{loops R}}, \dots, \\ N_{\text{hairpins}}, N_{\text{bulges L}}, N_{\text{loops R}}, \frac{N_{\text{unpaired nt}}}{L_{seq}}]$$

112 We note that psuedoknot features are not extracted or used for classification in this project as NUPACK dose not
113 predict these structures.

114 Left and right for the bulges and loops were defined as “left” when residing on the 5’ to 3’ direction of a stem loop
115 stack and as “right” when residing on the 3’ to 5’ direction of a stack. A bulge was defined as a one nucleotide unpaired
116 on either direction interrupting a contiguous stack, loops were defined as two or more unpaired nucleotides interrupting

117 a stack. This naming convention comes from the location when reading the dot structure left to right of the feature:
118 “...((.((....))))...” has one left bulge and “...((..((....))...)...)” has a left internal loop of two and a right internal loop of
119 three. Unpaired nucleotides are defined as base pairs not within any paired stack, for example, “...(((....)))...(((....)))...”
120 has 9 total unpaired nucleotides as the 8 unpaired nucleotides reside within a stack. The unpaired nucleotides inside
121 stacks are instead labeled as hairpins. A branching stack is defined as one that has multiple distinct substacks within
122 its stack, e.g. “((...((....))...((....))...))” is one branching stack containing two non-branching stacks. Counts of structural
123 features were max-normalized by the entire combined RS and 5’UTR data set. Figure 2B visually describes the process
124 of taking an example sequence and converting it to its representative feature vector.

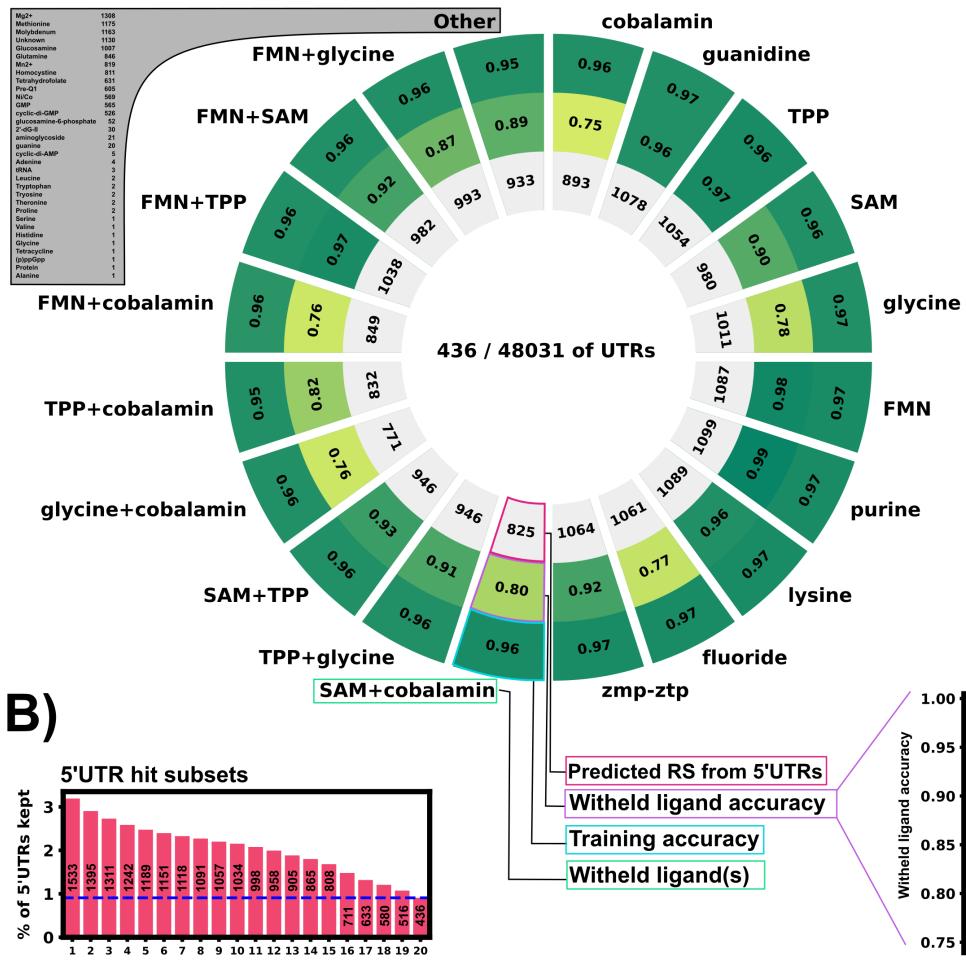
125 A two sample Kolmogorov-Smirnov distance was calculated to compare differences between the known RS fea-
126 tures and *H. Sapiens* 5’UTR features. According to the KS distance, most structural features showed a marked disparity
127 between the RS and 5’UTR data set, Figure 1C, while sequence features ranged from 0.05 - 0.4 in their KS distance.

128 PU learning and cross validation achieves 75% to 99% accuracy to identify held out known riboswitches.

129 To assess performance of our positive-unlabeled classifiers, we generated 20 separate subsets of training and vali-
130 dation data by withholding specific subsets of the known riboswitches based on their class of ligand. The first ten
131 validation subsets were generated by selecting each of the ten most represented ligand classes (each comprising 2% or
132 more of the overall data) and leaving each one out: Cobalamin, guanidine, TPP, SAM, glycine, FMN, purine, lysine,
133 fluoride, zmp-ztp. The next nine subsets were generated by leaving out pairs of the most commonly represented lig-
134 ands: FMN+glycine, FMN+SAM, FMN+TPP, FMN+cobalamin, TPP+cobalamin, TPP+glycine, TPP+SAM, cobal-
135 amin+SAM, cobalamin+TPP. The final (and most diverse) validation set was created by selecting all riboswitch ligand
136 classes with less than 2% representation in the full RS data set (11305 sequences, 34 ligand classes, 16.9% of the
137 entire RS data set).

138 Validating each classifier on structural classes that were not provided gives a reasonable confidence that the classi-
139 fier can extrapolate to riboswitches that are not included in the training set – the target task for eukaryotic riboswitch
140 discovery. 20 Unweighted Elkan & Noto classifiers were trained on the 20 data subsets. Figure 3A shows positive
141 example training accuracy (outer ring), validation accuracy on withheld ligand sets (middle ring) and the positively la-
142 beled 5’UTR count (inner ring) of all 20 classifiers. Validation accuracy ranged from 75% to 99% across the classifiers.
143 The classifier validated with withheld TPP riboswitches had high validation accuracy (97%), which is encouraging
144 sign as all currently described eukaryotic riboswitches use TPP as their ligand (36; 62; 56; 10; 29; 38). The “other”
145 classifier trained the diverse validation set of 34 ligand classes achieved an 89% validation accuracy, demonstrating a
146 surprising ability to extrapolate to examples that are dissimilar from one another and underrepresented in the training
147 data.

A) Classifier training, validation, and application to 5'UTR



B)

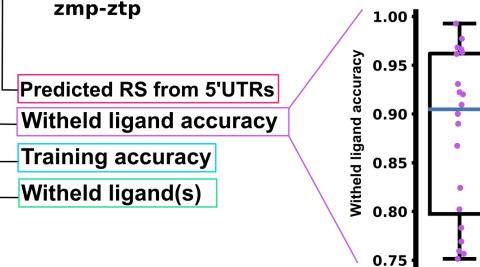
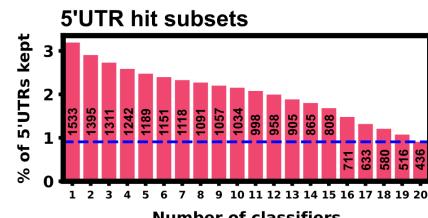


Figure 3: **Training and validation results of 20 PU classifiers.** A: Training and validation results. Each slice represents one PUlearn Elkanato-classifier trained on a data set withholding one or two ligand-specific riboswitches. The outer ring shows the training accuracy on only positive examples (RS). The middle ring is the validation accuracy on the withheld riboswitch(es) of a particular ligand(s) class. The inner ring shows number of the predicted positive labeled 5'UTR sequences out of the 48,031 5'UTR sequences. The sub-panel on the bottom right shows the withheld validation accuracy (rounded to 2 digits) in a box plot. 436 5'UTRs were selected by all 20 classifiers as positive labeled – potentially harboring riboswitch-like features. B: 5'UTR hit subsets detected by varying numbers of classifiers (1 - 20, full sequences).

148 PU learning ensemble correctly identifies more than 50% of previously unseen synthetic theophylline
149 riboswitches

150 As an additional verification step of our machine learning approach, we applied our ensemble of 20 classifiers to a
151 wholly synthetic riboswitch data set, that is not represented anywhere within the training set. The training set used
152 to train the ensemble included 14 total synthetic riboswitch sequences, none of which use theophylline as a target
153 ligand. 25 current theophylline riboswitch sequences were obtained from Wang et al. (58) and were passed through

154 our feature extraction and ensemble classification, Figure 4. Figure 4A shows the classification of the 25 theophylline
 155 riboswitches vs. a selection threshold on the ensemble probability output (ranging from .01 to .99). Upon testing, 56%
 156 (14/25) sequences were correctly identified with an ensemble probability over 50%, with 9 riboswitches classified
 157 higher than 90%. For a comparison, 300 completely random 35-250 nucleotide length sequences were generated and
 158 classified with the ensemble. 7.8% (39/300) of the sequences were falsely identified as riboswitch sequences. At very
 159 strict thresholds (requiring a ≥ 0.98 or ≥ 0.99 ensemble output), 25% of theophylline riboswitches are detected by the
 160 ensemble, whereas only 6.6% of random sequences are false positives. A two-sided binomial test was performed
 161 with the random sequence false positive rate to show the detection rate of theophylline riboswitches was significantly
 162 higher than random chance, Figure 4B.

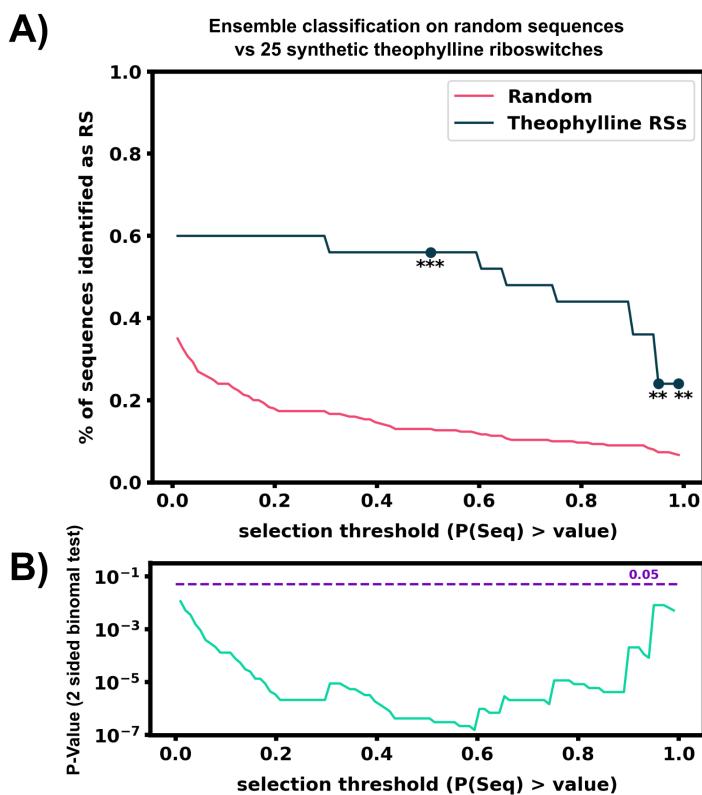


Figure 4: **Application of the Ensemble model to a class of wholly synthetic riboswitches.** A: 25 recently reviewed theophylline riboswitch sequences from (58) were obtained and classified using our 20-classifier ensemble. 300 random nucleotide sequences ranging from 30-300 nucleotides were also generated and classified using the ensemble. The percentage of sequences for both the 300 random sequences and 25 synthetic riboswitches are compared against each other for different positive ensemble probability thresholds. The ensemble incorrectly classifies 40% of the theophylline as non-riboswitches at any threshold. B: A p-value was obtained for every threshold using a two-sided binomial test of the amount theophylline hits vs. the false positive hits of the random sequences. The theophylline identified amount is significantly higher than the false positive rate (FPR) of the random sequence classification at every threshold. Select thresholds are highlighted in A: .5, .98, and .99.

163 This finding indicates that our ensemble can correctly extrapolate to the synthetic riboswitches, although with an
 164 accuracy of 56%, it misses more of the synthetic riboswitches than it did for the cross validation performed above
 165 on natural sequences. One potential explanation for this loss in selectivity is that the aptamer for theophylline was

166 discovered via directed evolution (SELEX) (23), potentially introducing a novel mechanism of action that is not
167 captured by our original training data set. Overall, we consider the ability to find 56% hit rate for synthetic riboswitches
168 to be another successful demonstration that the approach can identify potential riboswitches with novel structures or
169 evolutionary origins.

170 **PU learning with the trained ensemble model identifies and ranks a set of 1533 potential 5'UTR riboswitch hits**

171 Now that the ML ensemble has been verified through cross-validation, it is instructive to examine which 5'UTR
172 sequences have been identified as potential riboswitches. Among the classifiers, 436 5'UTRs were identified as har-
173 boring potential riboswitch elements by all 20 of the classifiers using a selection criteria of ≥ 0.95 classifier output.
174 Figure 3B shows the relative overlap of all 5'UTR sequences identified by one or more classifier with the same selec-
175 tion threshold. By contrast, the amount of 5'UTR sequences identified as riboswitches by one or more classifiers was
176 1533. The existence of an overlap when using all 20 classifiers instills confidence in our ensemble approach. If there
177 was a precipitous drop in identified sequences when using more and more classifiers, that would imply that classifiers
178 are individually identifying completely different subsets of the 5'UTR data set to consider as riboswitches. A drop
179 from 1533 hits to 436 hits when increasing the amount of classifier agreement is substantial, but still leaves a large
180 overlap found by all 20 trained classifiers.

181 To provide a method to rank the 5'UTR hits based on the ensemble of the classifiers, we average the normalized
182 outputs of the 20 PU classifiers to compute an ensemble similarity score, J_{Ensemble} :

$$J_{\text{Ensemble}}(\text{UTR}) = \frac{1}{20} \sum_i^{20} \frac{PU_i(\text{UTR})}{\max(PU_i(\text{All RS}))} \quad (3)$$

183 where $PU_i(\text{sequence})$ is the positive-unlabeled classifier with the i^{th} withheld ligand set, which is normalized by its
184 maximum value over all true RS.

185 **The majority of top 5'UTR hits are consistently selected despite substantial truncation to their 5' ends.**

186 Some potential 5'UTR hits are discarded in our analysis because we are using the full sequence data for the 5'UTR.
187 Because many 5'UTRs can be large with multiple regulatory elements, such an approach could miss smaller riboswitch
188 elements that are a sub-sequence and likely nearer to the start codon. To evaluate how many potential hits may be
189 lost by not considering sub-sequences of the 5'UTR during classification, we took our 5'UTR data set and for each
190 sequence, we generated 20 sub-sequences for each 5'UTR by truncating the mRNA at 20 evenly spaced locations
191 upstream of the start codon, starting 30 NT from 5' end. For each sub-sequence, we extracted the new features
192 and applied the ensemble classifier. Figure 5A shows the resulting probability of selection as a riboswitch versus
193 the fraction of the sequence used for all (48,031) 5'UTRs; Fig. 5B shows the same result but only for the subset
194 of 436 5'UTRs (0.91%) that were previously identified as likely riboswitches using the full length sequence; and
195 Fig. 5C shows the same result but for a distinct subset of 1210 5'UTRs (2.5%) that would have been identified as a

196 riboswitch by five or more partial-length sub-sequences, but *not* using the full sequence. Although the probability that
197 a given sequence being a riboswitch increases when using sub-sequences, the vast majority of 5'UTRs (97%) are still
198 discarded as unlikely to be riboswitches. Moreover, for the 5'UTRs that were identified as a riboswitch using their
199 full sequences (n=436), nearly half (45.5%) of these 5'UTRs are still detected as a riboswitch even when 85% of their
200 sequence is discarded. For example, AUH is consistently detected as a riboswitch with $\geq 95\%$ confidence for nearly
201 every sub-sequence Conversely, a small fraction of 5'UTRs, such as ATF1, is only identified as a riboswitch when the
202 sequence is 90% or more intact.

203 From a practical perspective, using the full ensemble and full sequences down-selects to more manageable number
204 of potential hits, and given the ultimate goal of reducing the potential sequence space to an experimentally viable
205 number, 436 is considered to be acceptable amount for future experimental validation. However, the remaining hits
206 from sub-sequences and ensemble agreement can be revisited and examined as needed and are provided within the
207 supplemental data.

208 Identified 5'UTR hits share remarkable feature similarities to known riboswitches.

209 Now that our ensemble classifier has identified a subset of 5'UTRs that may harbor potential riboswitches, it is illus-
210 trative to examine the properties of these hits. It is infeasible to manually compare all 436 hits to the RS data set, so to
211 aid with comparison, a GitHub page (https://will-raymond.github.io/human_riboswitch_hits_gallery/about/) was created to display, rank, and compare each 5'UTR to its most similar matches in the RS data set. The
212 website contains the subset of 5'UTRs identified by all 20 classifiers as potential riboswitches, as well as any 5'UTR
213 identified as a riboswitch by any individual classifier.

215 5'UTR to RS comparisons can be calculated several different ways, but for the purpose of the website, each 5'UTR
216 was compared to each RS with a using a combination of three metrics: sequence length difference, structural feature
217 vector mean-squared difference, and the predicted 5'UTR dot structure to RS dot structure Levenshtein distance (edit
218 distance). Length mean squared distance was calculated as:

$$D_L = |L_{\text{UTR}} - L_{\text{RS}}| \quad (4)$$

219 Likewise, the structural feature metric is also measured by the squared difference between any two extracted feature
220 sets:

$$D_{\text{struct}} = \sum_i^{74 \text{ features}} ([S_{5' \text{UTR}}]_i - [S_{\text{RS}}]_i)^2 \quad (5)$$

221 Finally, the dot structure metric is measured by the Levenshtein distance or “edit” distance between two strings – in
222 other words, how many edits (insertions, deletions, substitutions) to convert one string into another? Eq. 6 shows the
223 recursive letter by letter formulation of the Levenshtien distance, where tail(string) refers to everything but the first
224 letter of any given string. If we define a as the UTR sequences and b as the RS sequence, the Levenshtien distance can

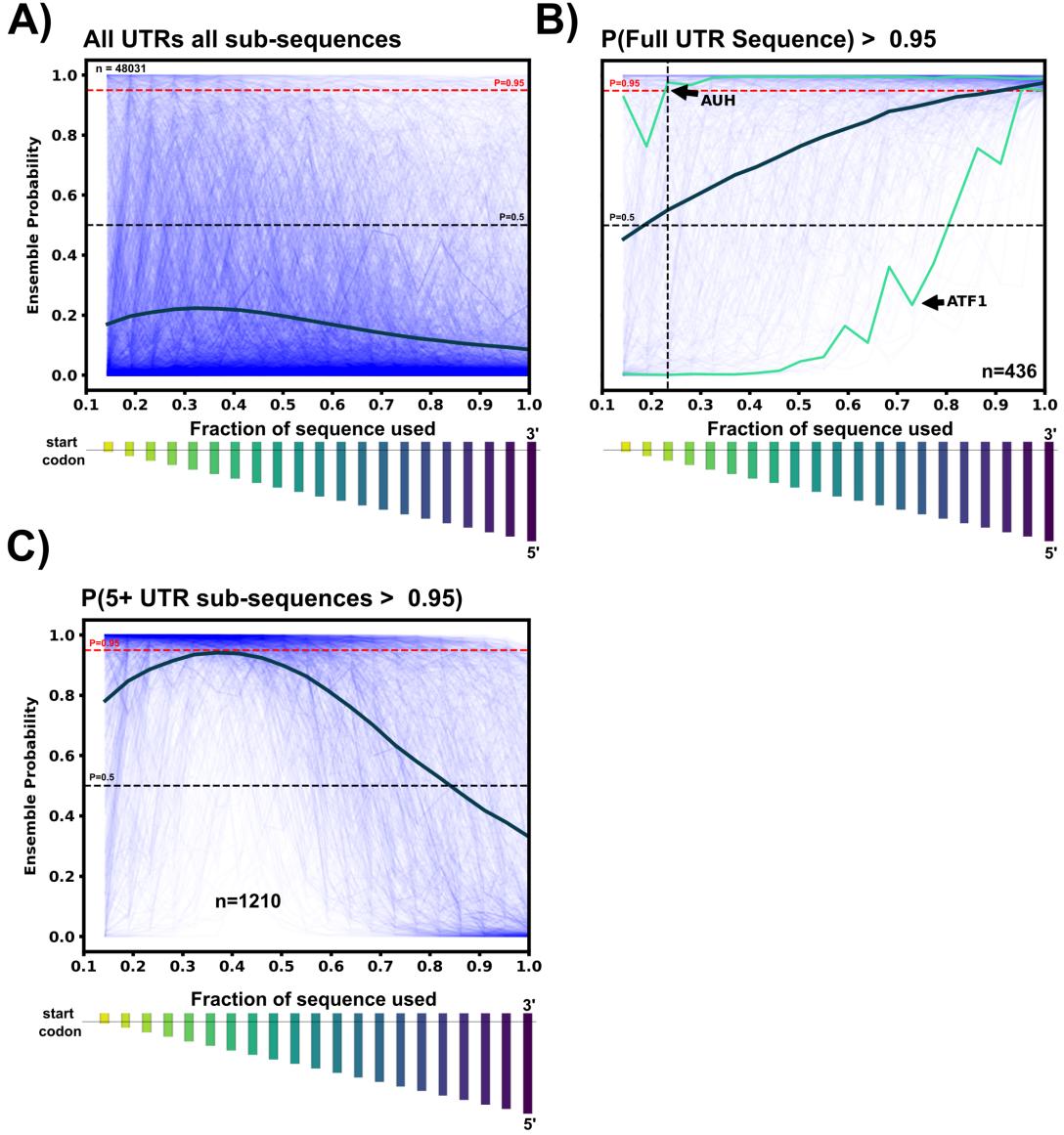


Figure 5: 5'UTR Sub-sequence exploration. A: For each 5'UTR sequence, 20 evenly-spaced sub-sequences were generated after the first 30 nucleotides in the 3'-5' direction, ensuring the start codon is in all sub-sequences. The relative size of each sub-sequence as a bar chart below the x-axis. For all 5'UTRs in the data set, variable-length sub-sequences were passed through the ensemble classifier to obtain the riboswitch probability. The riboswitch ensemble probability is plotted for each 5'UTR sub-sequence vs. the fraction of the sub-sequence to total 5'UTR length (thin blue lines). The thick dark blue line represents the average ensemble probability for that particular sub-sequence bin. C: Same as A, but only for the 5'UTRs whose full sequences were classified as $\geq 95\%$ riboswitch by the ensemble. Many 5'UTR sequences such as AUH are classified as a riboswitch until almost 80% of the original sequence is removed. In contrast, some sequences such as ATF1 are no longer considered a riboswitch once 10% of the sequence is removed from the 5' end. Once again the thick dark line represents the average probability of each sub-sequence bin. C: To find sub-sequences not included in the 436 hits, 5'UTR sequences not detected as a riboswitch by the full sequence but were detected as $\geq 95\%$ riboswitch in 5 or more sub-sequence bins were selected. These 1210 5'UTR sequences and their sub-sequence ensemble probabilities are plotted vs sub-sequence fraction. 1210 sequences could be included as potential riboswitch hits by removing some amount of 5' end nucleotides.

225 be computed as:

$$D_{\text{Lev}} = \left\{ \begin{array}{ll} \begin{array}{ll} \text{length}(a) & \text{if } \text{length}(b) = 0 \\ \text{length}(b) & \text{if } \text{length}(a) = 0 \\ \text{lev}(\text{tail}(a), \text{tail}(b)) & \text{if } a[0] = b[0] \\ \end{array} \\ 1 + \min \left\{ \begin{array}{l} \text{lev}(\text{tail}(a), b) \\ \text{lev}(a, \text{tail}(b)) \\ \text{lev}(\text{tail}(a), \text{tail}(b)) \end{array} \right\} \quad \text{otherwise} \end{array} \right\} \quad (6)$$

226 The combined similarity between the UTR and RS is denoted as J_{Sim} and is computed by normalizing each of the
227 above-described distances by the corresponding maximum value for that metric over all RS in the data set. For length
228 and Levenshtein distances, the maximum distances are: $D_{\text{Lev}} = \max((476 - L_{5'UTR}), L_{5'UTR})$, and $D_L = 476$,
229 where 476 is the length of the largest RS sequence in the training data (all UTR sequences are truncated to 300 or
230 less). For the structure feature distance, the normalization factor is obtained by comparing the 5'UTR in question to
231 the entire RS data set and finding the maximum. The combined similarity score is then defined on a scale from 0 (no
232 similarity) to 1 (perfect similarity) according to:

$$J_{\text{Sim}}(UTR, RS) = 1 - \frac{1}{3} \left(\frac{D_L}{\max(D_L)} + \frac{D_{\text{Lev}}}{\max(D_{\text{Lev}})} + \frac{D_{\text{struct}}}{\max(D_{\text{struct}})} \right) \quad (7)$$

233 Each 5'UTR entry is displayed on the website alongside its top three J_{Sim} matches within the RS data set. The
234 ligands of the top 20 5'UTR-RS J_{Sim} matches are also displayed as a preliminary prediction for the potential ligands
235 for that structure. However, future experimental validation would be necessary to ascertain if these hypothetical
236 matches are correct; the potential ligands are presented here more as an indicator of which ligand class from the
237 prokaryotic data is most represented in J_{Sim} matches to the 5'UTR sequence.

238 For the reader's convenience, the website table displaying all hits can be sorted by its similarity to known RS
239 (J_{Sim}) or by the ensemble probability from the PU classifier (J_{Ensemble} , see section 2.5). An example website page
240 is presented in Figure 6. Each column represents a 5'UTR hit (far left) or an RS entry from the training data (next 3
241 columns). For each sequence, the predicted secondary structure from NUPACK is displayed for visual comparison.
242 Base pair comparisons are also shown for each 5'UTR to RS pair with a circle plot of both predicted structures
243 overlaid. Below that is a comparison of each sequence's secondary structure features (stacks, loops, hairpins, etc).
244 The counts of each of these secondary structure features and the sequence dot structures are provided in tables on the
245 page as well. The goal of the website is to provide the reader with an instant visualization of each 5'UTR pair hit.

246 **Gene ontology points to enrichment of downstream *H. Sapiens* proteins associated with small molecules and
247 transcription / translation regulation**

248 A predominant function of bacterial riboswitches is to regulate the proteins directly related to the riboswitch's target
249 ligand. For example, a fluoride riboswitch may turn on genes useful for processing or mitigating fluoride for an

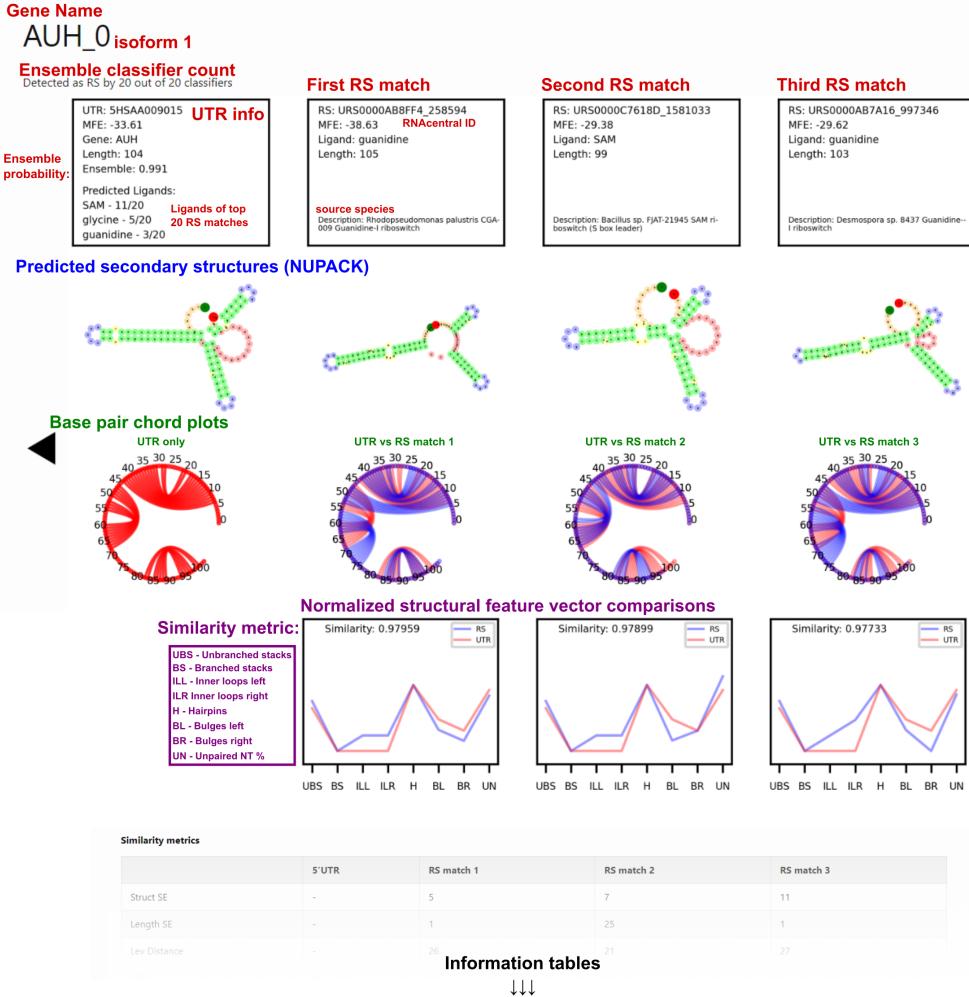


Figure 6: Example 5'UTR hit display from the website The display website (https://will-raymond.github.io/human_riboswitch_hits_gallery/_mds/AUH_0/) provides information on a given 5'UTR detected by the ensemble as riboswitch. Alongside each 5'UTR sequence, information on the top three riboswitch J_{Sim} matches to the 5'UTR are displayed in each column. First row provides information on a given sequence, UTRdb or RS id, source species, MFE of the predicted structure, and ensemble prediction probability for the 5'UTR. The next row displays the NUPACK predicted secondary structure for each sequence. Below that are chord plots representing the bonded base pairs for each RS sequence overlapping the 5'UTR chord plot. The next row shows the normalized structural feature vector comparison for structure counts for the 5'UTR and a given RS. J_{Sim} is reported in these plots. Additional information such as the dot structure, origin sequence, and counts of structural features are presented in the table below the comparison plots.

organism (33; 50). With this in mind, it is informative to examine the downstream proteins from the *H. Sapiens* 5'UTR hits for correlations in protein function, looking for genes associated with processing or synthesizing small molecules. Gene ontology (GO) analysis was preformed on the list of 5'UTR hits to look for any cellular function or process enrichment using the PANTHER database (37; 5; 54). GO process results are shown in Figure 7. The process ontology with significant fold enrichment fell into the following categories: Chromatin remodelling, transcription / translation regulation, mRNA splicing, mRNA and rRNA modification, and mitochondrial ubiqinone synthesis.

256 These enrichment results suggest a potential for small molecules to play a regulatory role in gene regulation, even
 257 if we cannot comment fully on our 5'UTR hit list without experimental validation. Interestingly, proteins directly
 258 involved in chemical stimulus detection were “unenriched” with no proteins found at all. This observation of mutual
 259 exclusion between potential riboswitches and sensing proteins is also reasonable - if there are already proteins capable
 260 to sense and respond to their intended stimulus, then there is no need to execute redundant functions in riboswitches.

GO Process fold enrichment

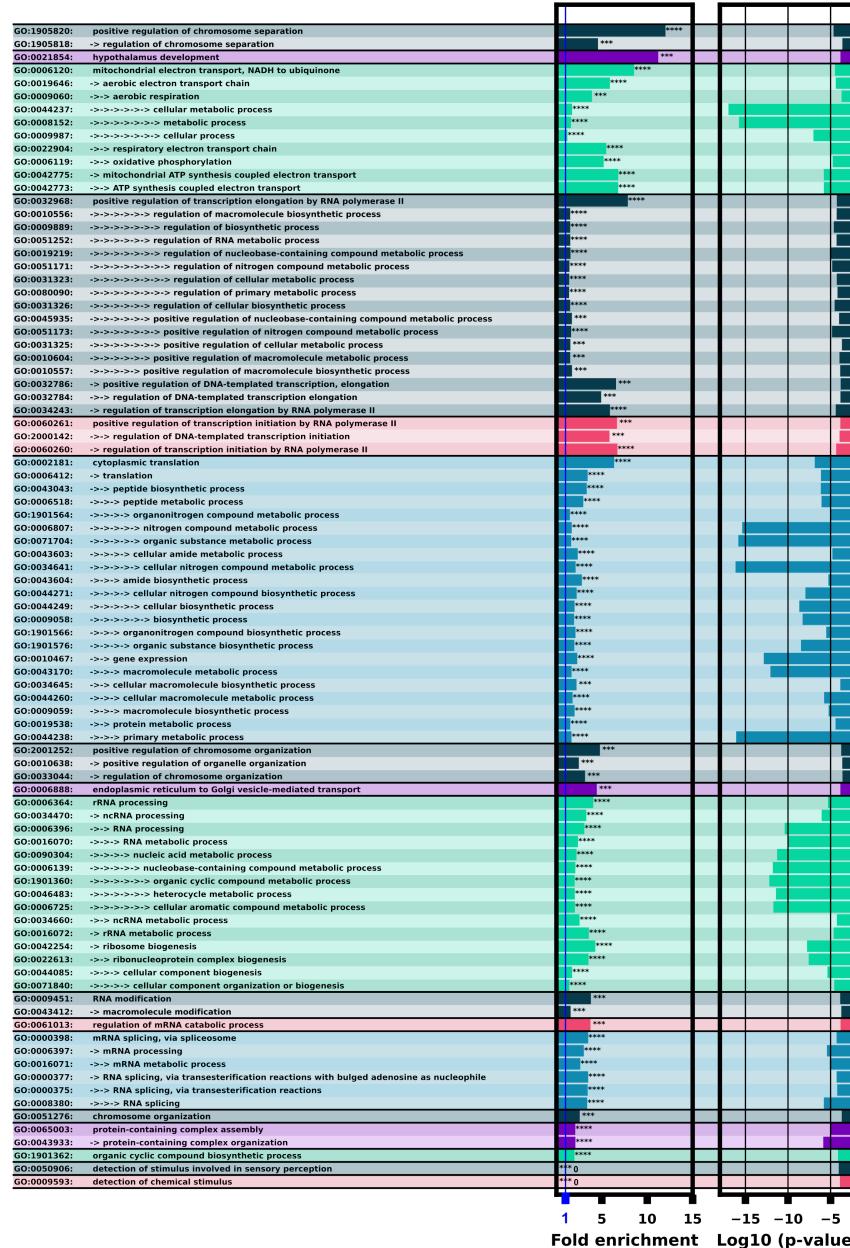


Figure 7: **GO process analysis with ID's and terms.** The left column lists the GO ID and term. Multiple arrows indicate GO term sub-levels. The left bar chart shows fold enrichment for that GO term with significance indicator. The second bar chart shows the log space of P-value significance for each enrichment.

261 GO function analysis showed a significant enrichment of downstream proteins implicated in binding various small
 262 molecules: nucleotides, nucleosides, ubiquinone, and various cyclic compounds, Figure 8. RNA binding and nu-
 263 cleotide binding molecules were extremely enriched with p-values ranging from 10^{-7} to 10^{-17} . Notably there is a
 264 negative enrichment of G protein-coupled receptors, this could be explained by riboswitches having direct signalling
 265 activity to a cell, bypassing typical trans-membrane signalling pathways.

GO Function fold enrichment

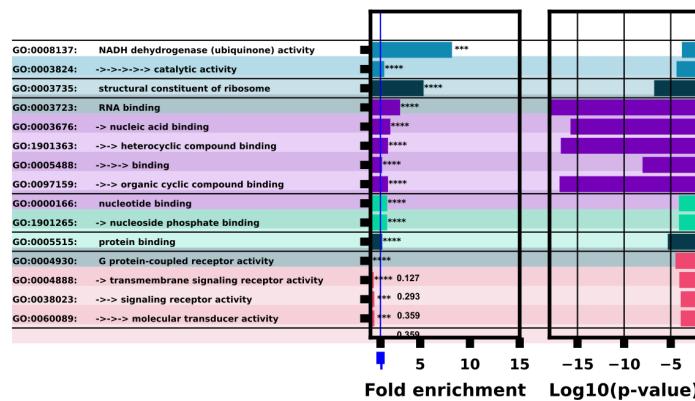


Figure 8: **GO function analysis with ID's and terms.** The left column lists the GO ID and term. Multiple arrows and indents indicate GO term sub-levels. The left bar chart column shows fold enrichment for that GO term with significance indicator. The second bar chart shows the log space of p-value significance for each GO term.

266 Some potential mRNA with potential riboswitches encode proteins with direct involvement in small molecule 267 functions

268 Although we have no experimental validation for our computationally discovered 5'UTR hits, it is illustrative to high-
 269 light and comment several interesting matches found in our computational analysis. Several computational hits such
 270 as AUH and FTSJ1 have roles in processing already known ligands. AUH plays a critical role in leucine degradation,
 271 hydrating 3-methylglutaconyl-CoA to 3-hydroxy-3-methyl-glutaryl-CoA (31). FTSJ1 is known to bind directly to
 272 S-adenosylmethionine (SAM) with a potential role in modifying ribosomal RNA and tRNAs (28).

273 Ubiquinone (CoQ10, CoQ) is an essential antioxidant component of the mitochondria (22; 59). Many proteins
 274 directly related to Ubiquinone synthesis or binding were selected as harboring riboswitches: Biosynthetic proteins
 275 COQ3, COQ9, COQ7; Several respiratory complex I proteins: NDUFA2, NDUFA8, NDUFB6, NDUFB9, NDUFS1,
 276 and a respiratory complex III protein: UQCRC. No ubiquinone binding riboswitches are currently described in the
 277 literature, although riboswitches binding other critical components of the electron transport chain such as NAD+ have
 278 been described (44; 32). Our overrepresentation of riboswitches within the mammalian mitochondria aligns with
 279 theories put-forth within Venkata Subbaiah et al. (53), where they speculate that the mitochondria's reduced genome
 280 may still harbor RNA switches as a mechanism as translational control.

281 Some proteins represented in the hit list are implicated in small molecule or amino acid synthesis or processing.
282 GSS is responsible for the second step of glutathione synthesis (43), PNPO directly converts vitamin B6 into its active
283 form (41).

284 Finally, several close matches in predicted structure and feature vectors should also be noted: ZNF480, SPAG11B,
285 UBAP2L-0, and TTPAL; However, to our knowledge, these proteins have no clear relation to small molecule process-
286 ing.

287 Conclusion

288 We have trained an ensemble of machine learning riboswitch classifiers using leave-one-out cross validation of ligand
289 classes consisting of 20 individual classifiers using sequence and predicted RNA structural features. Using this en-
290 semble classifier, we identified a subset of the *H. Sapiens* 5'UTR predicted to harbor riboswitch-like elements (Fig. 3).
291 This subset provides experimentalists with a prioritized list of sequences and genes to examine first when designing
292 exploratory experiments. This 436 sequence subset additionally shows positive GO fold enrichment results for down-
293 stream genes in many processes with direct small molecule involvement (Fig. 7 and 8). Our approach provides a
294 complementary strategy to the that taken in Mukherjee et al. (39), where the authors began with a known riboswitch
295 sequence and selectively mutated nucleotides while preserving structure, and then searched genomics data for a se-
296 quence match. In contrast, our approach starts with genomics data to learn our classifier and any given sequence
297 can be assessed for riboswitch probability. While our approach may be less targeted to the discovery of a riboswitch
298 with a particular structure, it holds the potential to extrapolate beyond single specific structures, as exemplified by its
299 identification of known synthetic riboswitches (Fig. 4). Searching for a riboswitch structure in a branch of life vastly
300 different than where riboswitches are previously described likely needs this extrapolation ability. In future work, our
301 approach could be replicated using the *H. Sapiens* 3'UTR, where described *H. Sapiens* pseudoriboswitches have been
302 found. However, for the purposes of this paper, we have limited our search to the 5'UTR because the bulk of our
303 training data (Bacterial riboswitches) act near their ribosomal binding site, equivalent to the 5'UTR, and looking in
304 the *H. Sapiens* 5'UTR first gives the best chance for efficient machine learning extrapolation. Our approach could also
305 be applied to other eukaryotic UTRs when experimentally validated. In this paper, we also briefly explored the how
306 the detection of potential 5'UTRs riboswitches could be expanded by varying how much of the 5'UTR sub-sequence
307 is used for the identification. This extended list from sub-sequences could be useful for finding better candidate
308 sequences should the 436 5'UTR hits not bare fruit. Because this is at present an entirely computational investigation,
309 we are unable to conduct experimental validation for any of the detected hits, but we provide this list to the greater
310 scientific community in the hope that these potential targets could kick start the discovery of a riboswitch within the
311 *H. Sapiens* translatome and beyond.

312 **Materials and methods**

313 **Computation**

314 All processing was done in [Python 3.8](#) with Biopython (11), NumPy (21), and NUPACK 4.0.0.23. Final data was stored
315 in .csv, .npy, and .json files and large files can be recomputed by the reader with the analysis notebook. The data files are
316 available at https://github.com/Will-Raymond/human_riboswitch_hits and the entire project computation
317 (sanitation, feature extracting, training, analyses) notebook is available at: <https://colab.research.google.com/drive/17zmKJh8iHAC2tImNNSyBrwUpU0uYKefx?usp=sharing>. A modified BEAR encoding in Python was
318 used for structural feature counting from dot structure strings (34).
319

320 **Positive Unlabeled Machine Learning**

321 Unweighted Elkan & Noto classifiers were used for our machine learning classifiers. In brief, this is an extension
322 of a generalized probability classifier to train on unknown / known labels as an approximation of class labels. Each
323 data point x has a label y which is either 0 or 1. Along with the label pair each data point has a known or unknown
324 flag, s , where $s = 1$ if the data point is known, and $s = 0$ when unknown. Therefore, when $s = 1$, $y = 1$ and
325 when $s = 0$, $y = \{0, 1\}$ Any binary classifier is then used to estimate $p(s = 1|y = 1, x)$ instead of the classical
326 estimate of $p(y, x)$. For our paper we used a SVC classifier from sci-kit learn with the following options: SVC(C=10,
327 kernel=rbf, gamma=0.4, probability=True). All PU classifiers were made using an implementation from the Python
328 package [PUlearn](#). For full details refer to the original paper by Elkan & Noto (16).

329 **GO Analysis**

330 GO analysis was performed with the PANTHER overrepresentation test (release 10.13.2022) using the 07.01.2022
331 release of the PANTHER database ([10.5281/zenodo.6799722](https://doi.org/10.5281/zenodo.6799722) Released 2022-07-01) using the Fisher's Exact test
332 with False Discovery Rate correction. The reference list used for comparison analysis was the *Homo Sapiens* gene list.
333 Overrepresentation test was performed for the GO biological process complete and GO biological function complete
334 annotated data sets.

335 **Acknowledgments**

336 WSR, and BM were supported by the NSF (1941870) and National Institutes of Health (R35GM124747). JD was
337 supported by the National Institutes of Health (1R01AI168459-01A1). Special thanks to Dr. Hamid Chitsaz as this
338 paper started as a student project in his CS548 - Bioinformatics class. Additional special thanks to Dr. Jeffrey Wilusz
339 as the initial idea for this was broached as a class discussion in MIP 543 - RNA biology.

340 **Conflict of Interest Statement**

341 The authors declare the absence of any commercial or financial relationships that could be construed as a conflict of
342 interest for this research.

343 **Author Contributions**

344 Conceptualization, implementation, computation, writing: WSR. Editing: WSR, JD, BM. Supervision: BM.

345 **Data Availability**

346 All original data files, sanitized data files, processed feature data, PU classifiers, Figure files, and ensemble results
347 used to create this manuscript are stored at: https://github.com/Will-Raymond/human_riboswitch_hits.
348 Any files too large to store on the manuscript repository can be regenerated with the analysis notebook. Upon final
349 acceptance for publication, a release containing all data and computational analyses will be frozen and provided via
350 Zenodo at [LINK].

351 **References**

- 352 [1] C. Abreu-Goodger and E. Merino. RibEx: a web server for locating riboswitches and other conserved bacterial
353 regulatory elements. *Nucleic Acids Research*, 33:W690, July 2005. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1160206/>, doi:10.1093/NAR/GKI445.
- 355 [2] S. D. Ali, H. Tayara, and K. T. Chong. Identification of piRNA disease associations using deep learning. *Computational and Structural Biotechnology Journal*, 20:1208–1217, January 2022. doi:10.1016/J.CSBJ.2022.02.026.
- 358 [3] N. Amin, A. McGrath, and Y. P. P. Chen. Evaluation of deep learning in non-coding RNA classification. *Nature Machine Intelligence*, 1:246–256, May 2019. URL: <https://www.nature.com/articles/s42256-019-0051-2>, doi:10.1038/s42256-019-0051-2.
- 361 [4] D. Antunes, N. A. Jorge, E. R. Caffarena, and F. Passetti. Using RNA sequence and structure for the prediction
362 of riboswitch aptamer: A comprehensive review of available software and tools. *Frontiers in Genetics*, 8:231,
363 January 2018. doi:10.3389/FGENE.2017.00231/BIBTEX.
- 364 [5] M. Ashburner, C. A. Ball, J. A. Blake, D. Botstein, H. Butler, J. M. Cherry, A. P. Davis, K. Dolinski, S. S. Dwight,
365 J. T. Eppig, M. A. Harris, D. P. Hill, L. Issel-Tarver, A. Kasarskis, S. Lewis, J. C. Matese, J. E. Richardson, M.
366 Ringwald, G. M. Rubin, and G. Sherlock. Gene Ontology: tool for the unification of biology NIH public access
367 author manuscript. *Nature Genetics*, 25:25–29, May 2000. doi:10.1038/75556.
- 368 [6] J. Bekker and J. Davis. Learning from positive and unlabeled data: a survey. *Machine Learning*, 109:719–
369 760, April 2020. URL: <https://link.springer.com/article/10.1007/s10994-020-05877-5>, doi:
370 10.1007/S10994-020-05877-5/FIGURES/6.
- 371 [7] P. Bengert and T. Dandekar. Riboswitch finder—a tool for identification of riboswitch RNAs. *Nucleic Acids
372 Research*, 32(Web Server issue):W154–159, Jul 2004. doi:10.1093/nar/gkh352.
- 373 [8] S. Bocobza, A. Adato, T. Mandel, M. Shapira, E. Nudler, and A. Aharoni. Riboswitch-dependent gene regulation
374 and its evolution in the plant kingdom. *Genes & Development*, November 2007. doi:10.1101/gad.443907.
- 375 [9] T. H. Chang, H. D. Huang, L. C. Wu, C. T. Yeh, B. J. Liu, and J. T. Horng. Computational Identification of
376 riboswitches based on RNA conserved functional sequences and conformations. *RNA*, 15(7):1426–1430, Jul
377 2009. doi:10.1261/rna.1623809.
- 378 [10] M. T. Cheah, A. Wachter, N. Sudarsan, and R. R. Breaker. Control of alternative RNA splicing and gene expres-
379 sion by eukaryotic riboswitches. *Nature*, 447:497–500, May 2007. URL: <https://pubmed.ncbi.nlm.nih.gov/17468745/>, doi:10.1038/NATURE05769.

- 381 [11] P. J. A. Cock, T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox, A. Dalke, I. Friedberg, T.
382 Hamelryck, F. Kauff, B. Wilczynski, and M. J. L. de Hoon. Biopython: freely available Python
383 tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25(11):1422–1423, March
384 2009. [arXiv:<https://academic.oup.com/bioinformatics/article-pdf/25/11/1422/944180/btp163.pdf>](https://academic.oup.com/bioinformatics/article-pdf/25/11/1422/944180/btp163.pdf), doi:10.
385 1093/bioinformatics/btp163.
- 386 [12] P. E. Compeau, P. A. Pevzner, and G. Tesler. How to apply de Bruijn graphs to genome assembly. *Nature
387 Biotechnology*, 29(11):987–991, November 2011. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5531759/>, doi:10.1038/nbt.2023.
- 389 [13] A. Cornish-Bowden. Nomenclature for incompletely specified bases in nucleic acid sequences: recommendations
390 1984. *Nucleic Acids Research*, 13(9):3021–3030, May 1985. doi:10.1093/nar/13.9.3021.
- 391 [14] F. Denis. PAC learning from positive statistical queries*. *Lecture Notes in Computer Science*, 1501:112–
392 126, 1998. URL: https://link.springer.com/chapter/10.1007/3-540-49730-7_9, doi:10.1007/
393 3-540-49730-7_9/COVER.
- 394 [15] F. Denis, R. Gilleron, and F. Letouzey. Learning from positive and unlabeled examples. *Theoretical Computer
395 Science*, 348:70–83, December 2005. doi:10.1016/J.TCS.2005.09.007.
- 396 [16] C. Elkan and K. Noto. Learning classifiers from only positive and unlabeled data. 2008.
- 397 [17] M. E. Fornace, J. Huang, C. T. Newman, N. J. Porubsky, M. B. Pierce, and N. A. Pierce. NU-
398 PACK: Analysis and design of nucleic acid structures, devices, and systems. *Theoretical and Computational
399 Chemistry*, November 2022. URL: <https://chemrxiv.org/engage/chemrxiv/article-details/636c7089b588507d0045f283>, doi:10.26434/CHEMRXIV-2022-XV98L.
- 401 [18] M. E. Fornace, N. J. Porubsky, and N. A. Pierce. A unified dynamic programming framework for the
402 analysis of interacting nucleic acid strands: Enhanced models, scalability, and speed. *ACS Synthetic Biology*,
403 9:2665–2678, October 2020. URL: <https://pubs.acs.org/doi/abs/10.1021/acssynbio.9b00523>,
404 doi:10.1021/ACSSYNBIO.9B00523/SUPPL_FILE/SB9B00523_SI_001.PDF.
- 405 [19] C. L. Giudice, F. Zambelli, M. Chiara, G. Pavesi, M. A. Tangaro, E. Picardi, and G. Pesole. UTRdb 2.0: a
406 comprehensive, expert curated catalog of eukaryotic mRNAs untranslated regions. *Nucleic Acids Research*,
407 51:D337–D344, January 2023. URL: <https://dx.doi.org/10.1093/nar/gkac1016>, doi:10.1093/NAR/
408 GKAC1016.
- 409 [20] G. Grillo, A. Turi, F. Licciulli, F. Mignone, S. Liuni, S. Banfi, V. A. Gennarino, D. S. Horner, G. Pavesi, E.
410 Picardi, and G. Pesole. UTRdb and UTRsite (RELEASE 2010): a collection of sequences and regulatory motifs
411 of the untranslated regions of eukaryotic mRNAs. *Nucleic Acids Research*, 38(Database issue):75–80, Jan 2010.
412 doi:10.1093/nar/gkp902.

- 413 [21] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor,
414 S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M.
415 Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E.
416 Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020. [doi:10.1038/s41586-020-2649-2](https://doi.org/10.1038/s41586-020-2649-2).
- 418 [22] A. Hidalgo-Gutiérrez, P. González-García, M. E. Díaz-Casado, E. Barriocanal-Casado, S. López-Herrador, C. M.
419 Quinzii, and L. C. López. Metabolic targets of coenzyme Q10 in mitochondria. *Antioxidants*, 10, April 2021.
420 URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8066821/>, doi:10.3390/ANTIOX10040520.
- 421 [23] R. D. Jenison, S. C. Gill, A. Pardi, and B. Polisky. High-resolution molecular discrimination by RNA. *Science*,
422 263(5152):1425–1429, March 1994. [doi:10.1126/science.7510417](https://doi.org/10.1126/science.7510417).
- 423 [24] Z. Ju and S. Wang. Computational identification of lysine glutarylation sites using positive-unlabeled learning.
424 *Current Genomics*, 21:204–211, May 2020. [doi:10.2174/1389202921666200511072327](https://doi.org/10.2174/1389202921666200511072327).
- 425 [25] I. Kalvari, E. P. Nawrocki, J. Argasinska, N. Quinones-Olvera, R. D. Finn, A. Bateman, and A. I. Petrov. Non-
426 coding RNA analysis using the Rfam database. *Current Protocols in Bioinformatics*, June 2018. URL: <http://rfam.org>, doi:10.1002/cpbi.51.
- 427 [26] K. Kavita and R. R. Breaker. Discovering riboswitches: the past and the future. *Trends in Biochemical Sciences*, 48:119–141, February 2023. URL: <http://www.cell.com/article/S0968000422002341/fulltext>, doi:10.1016/J.TIBS.2022.08.009.
- 431 [27] F. Li, S. Dong, A. Leier, M. Han, X. Guo, J. Xu, X. Wang, S. Pan, C. Jia, Y. Zhang, G. I. Webb, L. J. Coin, C. Li,
432 and J. Song. Positive-unlabeled learning in bioinformatics and computational biology: a brief review. *Briefings
433 in Bioinformatics*, 23, January 2022. URL: <https://academic.oup.com/bib/article/23/1/bbab461/6415313>, doi:10.1093/BIB/BBAB461.
- 435 [28] J. Li, Y.-N. Wang, B.-S. Xu, Y.-P. Liu, M. Zhou, T. Long, H. Li, H. Dong, Y. Nie, P. R. Chen, E.-D. Wang,
436 and R.-J. Liu. Intellectual disability-associated gene ftsj1 is responsible for 2'-O-methylation of specific tR-
437 RNAs. *EMBO reports*, 21:e50095, August 2020. URL: <https://onlinelibrary.wiley.com/doi/full/10.15252/embr.202050095>, doi:10.15252/EMBR.202050095.
- 439 [29] S. Li and R. R. Breaker. Eukaryotic TPP riboswitch regulation of alternative splicing involving long-distance
440 base pairing. *Nucleic Acids Research*, 41:3022–3031, March 2013. URL: <https://pubmed.ncbi.nlm.nih.gov/23376932/>, doi:10.1093/NAR/GKT057.
- 442 [30] Y. Li, C. Zhong, and S. Zhang. Finding consensus stable local optimal structures for aligned RNA sequences
443 and its application to discovering riboswitch elements. *International Journal of Bioinformatics Research and
444 Applications*, 10:498–518, September 2014. [doi:10.1504/IJBRA.2014.062997](https://doi.org/10.1504/IJBRA.2014.062997).

- 445 [31] M. Mack, U. Schniegler-Mattox, V. Peters, G. F. Hoffmann, M. Liesert, W. Buckel, and J. Zschocke. Biochemical
446 characterization of human 3-methylglutaconyl-CoA hydratase and its role in leucine metabolism. *The FEBS
447 Journal*, 273:2012–2022, May 2006. URL: <https://onlinelibrary.wiley.com/doi/full/10.1111/j.1742-4658.2006.05218.x>, doi:10.1111/J.1742-4658.2006.05218.X.
- 448
- 449 [32] S. N. Malkowski, T. C. Spencer, and R. R. Breaker. Evidence that the nadA motif is a bacterial riboswitch for the
450 ubiquitous enzyme cofactor NAD+. *RNA*, 25:1616–1627, December 2019. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6859854/>, doi:10.1261/RNA.072538.119/-/DC1.
- 451
- 452 [33] M. Mandal and R. R. Breaker. Gene regulation by riboswitches. *Nature Reviews Molecular Cell Biology*, 5:451–
453 463, June 2004. URL: <https://www.nature.com/articles/nrm1403>, doi:10.1038/nrm1403.
- 454
- 455 [34] E. Mattei, G. Ausiello, F. Ferrè, and M. Helmer-Citterich. A novel approach to represent and compare RNA
456 secondary structures. *Nucleic Acids Research*, 42:6146, June 2014. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4041456/>, doi:10.1093/NAR/GKU283.
- 457
- 458 [35] P. J. Mccown, K. A. Corbino, S. Stav, M. E. Sherlock, and R. R. Breaker. Riboswitch diversity and distribution.
459 *RNA*, July 2017. URL: <http://www.rnajournal.org/cgi/doi/10.1261/rna.061234.>, doi:10.1261/rna.061234.
- 460
- 461 [36] D. McRose, J. Guo, A. Monier, S. Sudek, S. Wilken, S. Yan, T. Mock, J. M. Archibald, T. P. Begley, A. Reyes-
462 Prieto, and A. Z. Worden. Alternatives to vitamin B1 uptake revealed with discovery of riboswitches in multiple
463 marine eukaryotic lineages. *The ISME journal*, 8:2517–2529, January 2014. URL: <https://pubmed.ncbi.nlm.nih.gov/25171333/>, doi:10.1038/ISMEJ.2014.146.
- 464
- 465 [37] H. Mi, A. Muruganujan, D. Ebert, X. Huang, and P. D. Thomas. PANTHER version 14: more genomes, a
466 new PANTHER GO-slim and improvements in enrichment analysis tools. *Nucleic Acids Research*, 47:419–426,
January 2018. URL: <http://geneontology.org>, doi:10.1093/nar/gky1038.
- 467
- 468 [38] S. Mukherjee, M. D. Retwitzer, D. Barash, and S. Sengupta. Phylogenomic and comparative analysis of the
469 distribution and regulatory patterns of TPP riboswitches in fungi. *Scientific reports*, 8, December 2018. URL:
<https://pubmed.ncbi.nlm.nih.gov/29615754/>, doi:10.1038/S41598-018-23900-7.
- 470
- 471 [39] S. Mukherjee, M. D. Retwitzer, S. M. Hubbell, M. M. Meyer, and D. Barash. A computational approach for
472 the identification of distant homologs of bacterial riboswitches based on inverse RNA folding. *Briefings in
473 Bioinformatics*, March 2023. URL: <https://pubmed.ncbi.nlm.nih.gov/36951499/>, doi:10.1093/BIB/BBAD110.
- 474
- 475 [40] S. Mukherjee and S. Sengupta. Riboswitch scanner: an efficient pHMM-based web-server to detect riboswitches
476 in genomic sequences. *Bioinformatics*, 32:776–778, March 2016. URL: <https://academic.oup.com/bioinformatics/article/32/5/776/1744033>, doi:10.1093/BIOINFORMATICS/BTV640.

- 477 [41] F. N. Musayev, M. L. D. Salvo, T.-P. Ko, V. Schirch, and M. K. Safo. Structure and properties of recombinant hu-
478 man pyridoxine 5'-phosphate oxidase. *Protein Science : A Publication of the Protein Society*, 12:1455, July 2003.
479 URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2323923/>, doi:10.1110/PS.0356203.
- 480 [42] E. P. Nawrocki and S. R. Eddy. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics*,
481 29(22):2933–2935, Nov 2013. doi:10.1093/bioinformatics/btt509.
- 482 [43] R. Njälsson and S. Norgren. Physiological and pathological aspects of GSH metabolism. *Acta Paediatrica*,
483 94:132–137, January 2005. URL: <https://pubmed.ncbi.nlm.nih.gov/15981742/>, doi:10.1111/J.
484 1651-2227.2005.TB01878.X.
- 485 [44] S. S. Panchapakesan, L. Corey, S. N. Malkowski, G. Higgs, and R. R. Breaker. A second riboswitch class for
486 the enzyme cofactor NAD+. *RNA*, 27:99–105, January 2021. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7749635/>, doi:10.1261/RNA.077891.120/-/DC1.
- 488 [45] R. Petegrosso, Z. Li, and R. Kuang. Machine learning and statistical methods for clustering single-cell RNA-
489 sequencing data. *Briefings in Bioinformatics*, 21:1209–1223, July 2020. URL: <https://academic.oup.com/bib/article/21/4/1209/5519426>, doi:10.1093/BIB/BBZ063.
- 491 [46] A. I. Petrov, S. J. E. Kay, I. Kalvari, K. L. Howe, K. A. Gray, E. A. Bruford, P. J. Kersey, G. Cochrane, R. D.
492 Finn, A. Bateman, A. Kozomara, S. Griffiths-Jones, A. Frankish, C. W. Zwieb, B. Y. Lau, K. P. Williams, P. P.
493 Chan, T. M. Lowe, J. J. Cannone, R. Gutell, M. A. Machnicka, J. M. Bujnicki, M. Yoshihama, N. Kenmochi,
494 B. Chai, J. R. Cole, M. Szymanski, W. M. Karlowski, V. Wood, E. Huala, T. Z. Berardini, Y. Zhao, R. Chen,
495 W. Zhu, M. D. Paraskevopoulou, I. S. Vlachos, A. G. Hatzigeorgiou, L. Ma, Z. Zhang, J. Puetz, P. F. Stadler, D.
496 McDonald, S. Basu, P. Fey, S. R. Engel, J. M. Cherry, P. J. Volders, P. Mestdagh, J. Wower, M. B. Clark, X. C.
497 Quek, and M. E. Dinger. RNACentral: a comprehensive database of non-coding RNA sequences. *Nucleic Acids
498 Research*, 45(D1):D128–D134, Jan 2017. doi:10.1093/nar/gkw1008.
- 499 [47] N. Pornputtapong, D. A. Acheampong, P. Patumcharoenpol, P. Jenjaroenpun, T. Wongsurawat, S. R. Jun, S.
500 Yongkettrakul, N. Chokesajjawatee, and I. Nookae. KITSUNE: A tool for identifying empirically optimal
501 k-mer length for alignment-free phylogenomic analysis. *Frontiers in Bioengineering and Biotechnology*, 8:1080,
502 September 2020. URL: <https://github.com/natapol/kitsune>, doi:10.3389/fbioe.2020.556413.
- 503 [48] K. A. R. Premkumar, R. Bharanikumar, and A. Palaniappan. Riboflow: using deep learning to clas-
504 sify riboswitches with ~99% accuracy. *bioRxiv*, 2019. URL: <https://www.biorxiv.org/content/early/2019/12/08/868695>, arXiv:<https://www.biorxiv.org/content/early/2019/12/08/868695.full.pdf>, doi:
505 10.1101/868695.
- 507 [49] P. S. Ray, J. Jia, P. Yao, M. Majumder, M. Hatzoglou, and P. L. Fox. A stress-responsive RNA switch regulates
508 VEGFA expression. *Nature*, 457(7231):915–919, Feb 2009. doi:10.1038/nature07598.

- 509 [50] C. E. Scull, S. S. Dandpat, R. A. Romero, and N. G. Walter. Transcriptional riboswitches integrate timescales
510 for bacterial gene expression control. *Frontiers in Molecular Biosciences*, 7:480, January 2021. [doi:10.3389/fmolb.2020.607158/BIBTEX](#).
- 512 [51] A. Serganov and E. Nudler. A Decade of Riboswitches. *Cell*, 152(1-2):17–24, Jan 2013. [doi:10.1016/j.cell.2012.12.024](#).
- 514 [52] C. Su, J. D. Weir, F. Zhang, H. Yan, and T. Wu. ENTRNA: A framework to predict rna foldabil-
515 ity. *BMC Bioinformatics*, 20:1–11, July 2019. URL: <https://link.springer.com/article/10.1186/s12859-019-2948-5>, [doi:10.1186/S12859-019-2948-5/TABLES/6](#).
- 517 [53] K. C. V. Subbaiah, O. Hedaya, J. Wu, F. Jiang, and P. Yao. Mammalian RNA switches: Molecular rheostats
518 in gene regulation, disease, and medicine. *Computational and Structural Biotechnology Journal*, 17:1326,
519 January 2019. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6849081/>, [doi:10.1016/J.CSBJ.2019.10.001](#).
- 521 [54] The Gene Ontology Consortium. The Gene Ontology resource: enriching a GOld mine. *Nucleic Acids Research*,
522 2021. URL: <http://geneontology.org/go-cam.>, [doi:10.1093/nar/gkaa1113](#).
- 523 [55] A. Wachter. Riboswitch-mediated control of gene expression in eukaryotes. *RNA Biology*, 7(1):67–76, 2010.
524 [doi:10.4161/rna.7.1.10489](#).
- 525 [56] A. Wachter, M. Tunc-Ozdemir, B. C. Grove, P. J. Green, D. K. Shintani, and R. R. Breaker. Riboswitch control
526 of gene expression in plants by splicing and alternative 3' end processing of mRNAs. *The Plant Cell*, 19:3437–
527 3450, November 2007. URL: <https://pubmed.ncbi.nlm.nih.gov/17993623/>, [doi:10.1105/TPC.107.053645](#).
- 529 [57] C. Wang, C. Ding, R. F. Meraz, and S. R. Holbrook. PSoL: a positive sample only learning algorithm for finding
530 non-coding RNA genes. *Bioinformatics*, 22:2590–2596, November 2006. URL: <https://dx.doi.org/10.1093/bioinformatics/btl441>, [doi:10.1093/BIOINFORMATICS/BTL441](#).
- 532 [58] X. Wang, C. Fang, Y. Wang, X. Shi, F. Yu, J. Xiong, S.-H. Chou, and J. He. Systematic comparison
533 and rational design of theophylline riboswitches for effective gene repression. *Microbiology Spectrum*, 11,
534 February 2023. URL: <https://journals.asm.org/doi/10.1128/spectrum.02752-22>, [doi:10.1128/spectrum.02752-22](#).
- 536 [59] Y. Wang and S. Hekimi. Understanding ubiquinone. *Trends in Cell Biology*, 26:367–378, May 2016. URL:
537 <http://www.cell.com/article/S0962892416000039/fulltext>, [doi:10.1016/j.tcb.2015.12.007](#).
- 538 [60] T. J. Wheeler and S. R. Eddy. nhmmer: DNA homology search with profile HMMs. *Bioinformatics*,
539 29(19):2487–2489, Oct 2013. [doi:10.1093/bioinformatics/btt403](#).

- 540 [61] Y. Xiao, J. Wu, Z. Lin, and X. Zhao. A semi-supervised deep learning method based on stacked sparse auto-
541 encoder for cancer prediction using RNA-seq data. *Computer Methods and Programs in Biomedicine*, 166:99–
542 105, November 2018. [doi:10.1016/j.cmpb.2018.10.004](https://doi.org/10.1016/j.cmpb.2018.10.004).
- 543 [62] S. Yadav, D. Swati, and H. Chandrasekharan. Thiamine pyrophosphate riboswitch in some representative plant
544 species: a bioinformatics study. *Journal of Computational Biology*, 22:1–9, January 2015. URL: <https://pubmed.ncbi.nlm.nih.gov/25243980/>, [doi:10.1089/CMB.2014.0169](https://doi.org/10.1089/CMB.2014.0169).
- 545 [63] P. Yang, X. Li, H. N. Chua, C. K. Kwoh, and S. K. Ng. Ensemble positive unlabeled learning for disease gene
546 identification. *PLoS ONE*, 9:e97079, May 2014. URL: <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0097079>, [doi:10.1371/JOURNAL.PONE.0097079](https://doi.org/10.1371/JOURNAL.PONE.0097079).
- 547 [64] X. Zeng, Y. Zhong, W. Lin, and Q. Zou. Predicting disease-associated circular rnas using deep forests combined
548 with positive-unlabeled learning methods. *Briefings in Bioinformatics*, 21:1425–1436, July 2020. URL: <https://dx.doi.org/10.1093/bib/bbz080>, [doi:10.1093/BIB/BBZ080](https://doi.org/10.1093/BIB/BBZ080).
- 549 [65] J. Zhou, X. Lu, W. Chang, C. Wan, X. Lu, C. Zhang, and S. Cao. PLUS: Predicting cancer metas-
550 tasis potential based on positive and unlabeled learning. *PLOS Computational Biology*, 18:e1009956,
551 March 2022. URL: <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1009956>, [doi:10.1371/JOURNAL.PCBI.1009956](https://doi.org/10.1371/JOURNAL.PCBI.1009956).
- 552

Ligand	Count	%	Ligand	Count	%	Ligand	Count	%	Ligand	Count	%
Cobalamin	15718	23.2	Molybdenum	1163	1.7	2'-dG-II	30	0.044	Histidine	1	1.5e-5
TPP	12459	18.4	Unknown	1130	1.7	aminoglycoside	21	0.031	Protein	1	1.5e-5
SAM	8686	12.8	Glucosamine	1007	1.5	guanine	20	0.029	Glycine	1	1.5e-5
Glycine	4835	7.1	Glutamine	846	1.2	cyclic-di-AMP	5	7.4e-5	Tetracycline	1	1.5e-5
FMN	4255	6.3	Mn2+	819	1.2	Adenine	4	5.9e-5	(p)ppGpp	1	1.5e-5
Purine	2648	3.9	homocysteine	811	1.2	tRNA	3	4.4e-5	Alanine	1	1.5e-5
Lysine	2318	3.4	tetrahydrofolate	631	0.9	Leucine	2	3.0e-5	Serine	1	1.5e-5
Fluoride	1975	2.9	Pre-Q1	605	0.9	Tryptophan	2	3.0e-5			
zmp-ztp	1841	2.7	Ni/Co	569	0.8	Tyrosine	2	3.0e-5			
Guanidine	1640	2.4	GMP	565	0.8	Proline	2	3.0e-5			
Mg2+	1308	1.9	cyclic-di-GMP	526	0.8	Threonine	2	3.0e-5			
Methionine	1175	1.7	glucosamine-6-phosphate	52	0.078	Valine	1	1.5e-5			

Table 1: Ligand representation within the data set

IUPAC nucleotide code	Base(s)	Converted to
A	Adenine	A
C	Cytosine	C
G	Guanine	G
T / U	Thymine or Uracil	U
R	A or G	A
Y	C or T	C
S	G or C	G
W	A or T	A
K	G or T	G
M	A or C	A
B	C or G or T	C
D	A or G or T	A
H	A or C or T	A
N	Any	A

Table 2: Nucleotide substitution for data sanitation