

## RESEARCH ARTICLE

## Comparison and benchmark of deep learning methods for non-coding RNA classification

Constance Creux<sup>1,2</sup>, Farida Zehraoui<sup>1</sup>, François Radvanyi<sup>2</sup>, Fariza Tahī<sup>1\*</sup><sup>1</sup> Université Paris-Saclay, Univ Evry, IBISC, Evry-Courcouronnes, France, <sup>2</sup> Molecular Oncology, PSL Research University, CNRS, UMR, Institut Curie, Paris, France\* [fariza.tahi@univ-evry.fr](mailto:fariza.tahi@univ-evry.fr)

## Abstract

The involvement of non-coding RNAs in biological processes and diseases has made the exploration of their functions crucial. Most non-coding RNAs have yet to be studied, creating the need for methods that can rapidly classify large sets of non-coding RNAs into functional groups, or classes. In recent years, the success of deep learning in various domains led to its application to non-coding RNA classification. Multiple novel architectures have been developed, but these advancements are not covered by current literature reviews. We present an exhaustive comparison of the different methods proposed in the state-of-the-art and describe their associated datasets. Moreover, the literature lacks objective benchmarks. We perform experiments to fairly evaluate the performance of various tools for non-coding RNA classification on popular datasets. The robustness of methods to non-functional sequences and sequence boundary noise is explored. We also measure computation time and CO<sub>2</sub> emissions. With regard to these results, we assess the relevance of the different architectural choices and provide recommendations to consider in future methods.

## OPEN ACCESS

**Citation:** Creux C, Zehraoui F, Radvanyi F, Tahī F (2024) Comparison and benchmark of deep learning methods for non-coding RNA classification. *PLoS Comput Biol* 20(9): e1012446. <https://doi.org/10.1371/journal.pcbi.1012446>

**Editor:** Shi-Jie Chen, University of Missouri, UNITED STATES OF AMERICA

**Received:** March 19, 2024

**Accepted:** August 30, 2024

**Published:** September 12, 2024

**Copyright:** © 2024 Creux et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** Data and code used for running experiments are available at <https://EvryRNA.ibisc.univ-evry.fr/EvryRNA/ncBench> as part of the EvryRNA bioinformatics platform.

**Funding:** With financial support from ITMO Cancer of Aviesan within the framework of the 2021-2030 Cancer Control Strategy, on funds administered by Inserm (received by CC). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

## Author summary

RNA can either encode proteins, which perform different functions in the genome, or be non-coding. Non-coding RNAs represent around 98% of the genome, and were long thought to be non-functional. It has now been proven that non-coding RNAs can have diverse biological functions and be involved in diseases. A large proportion of non-coding RNAs has not yet been studied. The function of specific non-coding RNAs can be studied experimentally, but experiments are costly and time-consuming. One possibility to massively characterize the function of non-coding RNAs is to use computational methods to classify them into functional groups, or classes. Recent computational methods for non-coding RNA classification are all based on deep learning, as it leads to faster runtime and improved performance. Our work presents and compares the different approaches adopted in the state-of-the-art, as well as the non-coding RNA datasets that are used. We also present a comprehensive benchmark, measuring classification performance in different conditions, computation time, and CO<sub>2</sub> emissions. The descriptions and comparisons provided are meant to guide researchers in the field, whether wanting to use existing tools or to develop new ones.

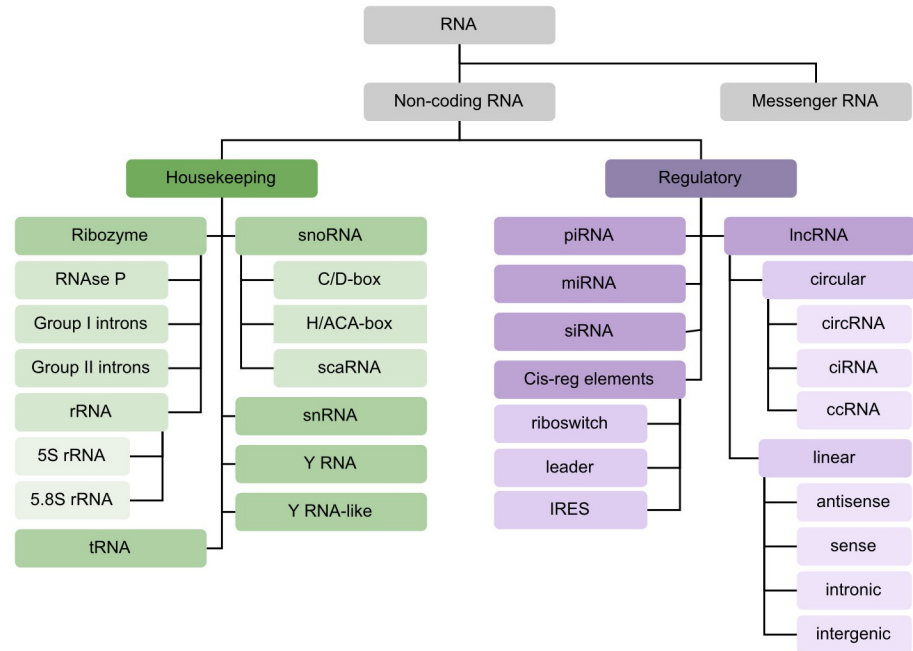
**Competing interests:** The authors have declared that no competing interests exist.

## Introduction

The definition of non-coding RNA (ncRNA) classes of similar functions started in the 1950s, with the discovery of tRNAs (transfer RNAs) and rRNAs (ribosomal RNAs) [1]. The description of other ncRNA classes involved in cell maintenance followed, but research still mainly focused on proteins or protein-coding genes. It is in the 2000s, when efforts like the Human Genome Project [2] highlighted that 98% of the genome is non-coding, that attention started to turn to ncRNA. Novel classes of ncRNAs with regulatory functions, like miRNA (microRNA) or lncRNA (long ncRNA), were described. Studies have shown that ncRNAs are implicated in various biological processes and diseases [3–5], underscoring their potential as biomarkers and therapeutic targets, and thus showing that studying their functions is important. Knowing their precise functions requires specific experimental studies, which are resource-intensive. To focus these resources on ncRNAs of interest, a first functional characterization can be done computationally on a large scale. Groups can be formed to gather ncRNAs that perform the same function. This is typically done at two levels. The first is ncRNA families, which are quite specific, and represent ncRNAs that share the same function, have similar structural motifs, and have sequence homology. The second is ncRNA classes, which are broader. They also represent ncRNAs sharing a function and structural motifs, but sequence homology is no longer a constraint. In this work, we focus on ncRNA class prediction. Still, it is worth noting that ncRNA family prediction has also garnered interest in the community and has been addressed in multiple works [6–8].

Many ncRNA classes have been identified over the years. It is typical to separate them based on whether they perform ‘housekeeping’ functions, i.e., they are involved in cell viability, or whether they are ‘regulatory’, i.e., they regulate gene expression at the epigenetic, transcriptional or post-transcriptional level. We represent in Fig 1 these two types of functions, naming some ncRNA classes that are commonly included in ncRNA classification problems. Each of these classes is associated to a specific function. For example, tRNAs [9] have a distinctive cloverleaf-like structure, and their function is to carry amino acids to the ribosome for protein synthesis. MiRNAs [10] range from 20 to 26 nt long, their precursors have a hairpin structure, and they bind other RNAs, mostly mRNAs, to regulate gene expression. Through mechanisms such as decapping, deadenylation, and translational repression, they typically repress or even silence expression [11]. Some classes are better known than others, and new classes could still be discovered as new functions of ncRNAs are uncovered: to this day, there is no fixed set of ncRNA classes. LncRNAs particularly suffer from this lack of definition [12]. RNAs belonging to this class are disparate and can have multiple functions. They are mostly defined by their length over 200 nt and their position in the genome (antisense, intronic, . . .). Due to this lack of precision, ncRNA classification problems tend to focus on classes of small ncRNAs (sncRNAs).

Several tools have been proposed since the 2000s to predict known ncRNA classes, mainly based on their sequence, secondary structure, or expression data. The use of computational resources enables the obtention of results faster than experimental methods, as multiple ncRNAs can be studied at the same time. The first approaches relied mainly on alignment algorithms and standard machine learning. In recent years, deep learning (DL) has been applied to ncRNA classification, showing improved performance and an ability to classify even larger ncRNA datasets. DL is now the prevalent approach in the field. The term ‘ncRNA classification’ can have different meanings in the literature. It is sometimes used to refer to the prediction of specific ncRNA classes [13–22]. This problem has been widely covered in reviews [23–26]. It can also refer to multi-class classification, taking as input a set of ncRNAs and associating each one to a ncRNA class, which is the scope of this review.



**Fig 1. Taxonomy of ncRNA classes.** Non-coding RNA functions can be separated into housekeeping (green) and regulatory (purple). Classes can be subdivided depending on the level of description. The classes represented are examples that are often mentioned in the literature, or that are of interest in the datasets mentioned below.

<https://doi.org/10.1371/journal.pcbi.1012446.g001>

The literature lacks comparisons of recent multi-class ncRNA classifiers. In an article published in 2017 [27], earlier methods for ncRNA classification are presented. Only one work, published in 2019 [24], covers the use of DL for multi-class ncRNA classification. We count twelve DL-based ncRNA classifiers in the state-of-the-art; eleven of them were presented after its publication and have not been covered in other literature.

This lack of reviews is an issue in the field, as the performance comparisons proposed in publications accompanying the different tools are unreliable. Indeed, comparisons between methods are often made without re-executing tools and instead simply using the scores announced in previous publications. In these instances, results are not always obtained with comparable experimental protocols. Moreover, the main dataset used for evaluation presents a data leakage problem; thus, reported performance is biased.

We propose to fill a gap in the literature by presenting an exhaustive review of the different DL approaches for multi-class ncRNA classification. We describe the choices in DL architecture, and explain the differences between datasets used for evaluation. Moreover, we conduct a fair comparison of performance between current tools. We evaluate global and per-class prediction results. The tools' ability to recognize non-functional sequences and to disregard sequence boundary noise is reported. We also assess resource intensity by measuring computation time as well as CO<sub>2</sub> emissions.

This paper is organized as follows: the first section presents the state-of-the-art of DL for ncRNA classification, exhaustively describing current classifiers and the accompanying datasets. In the next section we present the results of our comparison of different tools in terms of global and per-class performance, robustness to noise, as well as computation time and CO<sub>2</sub> emissions. We then discuss these observations, before presenting some concluding remarks. Finally, the Methods section presents the details of our experiments.

## Deep learning for ncRNA classification

### State-of-the-art classifiers

Earlier methods for ncRNA classification were primarily based on alignment algorithms [6, 28–30] or standard machine learning algorithms [31–33]. Feature extraction in these methods is a particularly crucial step, with a high impact on performance. Current ncRNA classifiers are all based on deep learning: most methods are based on Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), or a combination of both.

Applications of DL to ncRNA classification started with multiple CNN-based methods: nRC [34], ncRDeep [35], ncRNA-deep [36] (this method is unnamed, we refer to it using the name of its GitHub repository), RPC-snRC [37], ncRDense [38] and imCnC [39]. nRC pioneered the use of DL in the field. As input, ncRNAs are represented by vectors containing information on the presence in their secondary structure of discriminative sub-structures. Two convolutional layers extract relevant information, and Fully-Connected Layers (FCLs) are used for classification. ncRDeep is another method based on two convolutions followed by FCLs. However, the input is a simple representation of the sequence, a matrix created from one-hot encodings of nucleotides. In the method ncRNA-deep, the number of convolutional layers rises to five. Authors tested multiple sequence encodings (space-filling curves or one-hot encoding of k-mers of length 1, 2 or 3). A one-hot encoding of 1-mers, which corresponds to the same representation as ncRDeep, was selected for its superior results. This representation is also used by RPC-snRC, and fed to three consecutive Dense blocks (based on DenseNets, multiple convolutional layers with dense connections), with a final FCL for classification. ncRDense is an evolution of ncRDeep. In addition to sequence information, the method includes a second input: a matrix with multiple representations of nucleotides and one-hot encoding of the secondary structure. Both inputs are processed with independent DenseNets and merged. A final convolutional layer is followed by FCLs for classification. imCnC stands out as a multi-source CNN. Based on the idea that ncRNAs are related to chromatin states and epigenetic modifications, the approach builds a CNN for each input source (such as the sequence and different epigenetic information), before merging representations to classify samples with FCLs.

RNA sequences can be viewed as sentences, making the use of RNNs relevant to extract meaningful information from ncRNAs, as done by the methods ncRFP [40] and NCYPred [41]. ncRFP uses a bidirectional Long Short-Term Memory RNN (bi-LSTM). When learning a representation for a nucleotide, it is able to take into account close and distant elements, parsing the sequence both forward and backward. An attention mechanism allocates more weight to important information, and FCLs are used for classification. NCYPred also uses a bi-LSTM and FCLs for classification. The main difference between the two is that in ncRFP, each nucleotide is encoded separately, while in NCYPred the sequence is first transformed into overlapping 3-mers, and it is these 3-mers that are encoded.

Several methods combine both approaches, first using RNNs for sequence representation, from which CNNs then extract relevant information. This is the case for ncDLRES [42], ncDENSE [43] and MFpred [44]. ncDLRES is an updated version of ncRFP's architecture. The bi-LSTM is replaced by a dynamic LSTM, which only parses the sequence in the forward sense, but has the advantage of allowing inputs of varying length. For classification, the FCLs are replaced by a Residual Neural Network (ResNet), which is composed of multiple convolutional layers with skip connections. ncDENSE further evolves that architecture, replacing the dynamic LSTM by a dynamic bidirectional Gated Recurrent Unit (bi-GRU), and the ResNet by a DenseNet. MFpred proposes a rich input representation: sequences are encoded into four sets of features, each processed by a dynamic bi-GRU, reshaped into matrices, and fused.

Compared to ncDENSE, the DenseNet is replaced by a Squeeze and Excitation ResNet (SE ResNet), which models interdependencies between the four matrices of learned sequence encodings.

Two existing methods do not fall into the above categorization: RNAGCN [45] and Graph+DL [46] (this method is unnamed, we refer to it by naming its two components), both based on graphs. Taking the secondary structure as input, RNAGCN uses a Graph Neural Network (GNN) with FCL for classification. GNNs are able to learn meaningful representations of non-Euclidean data. RNAGCN is, to date, the only method directly using ncRNA secondary structure without manual feature extraction, removing a source of error. The approach Graph+DL is based on features extracted from the secondary structure. Authors propose a protocol to obtain a rich graph theory-based representation of structures. Other than the extensive feature preparation step, the method is based on a simple DL network of five FCLs.

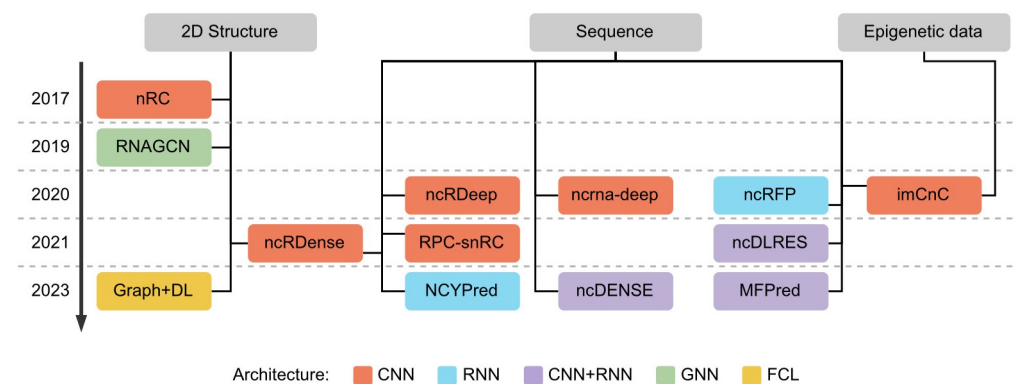
Note that for methods taking the secondary structure as input, sequence information is still represented. In RNAGCN, the nodes in the graph represent nucleotides, and edges represent pairings. In nRC and Graph+DL, the motifs that are extracted from the secondary structure are defined in relation to specific nucleotides.

We present in Fig 2 a chronological view of the methods for ncRNA classification described above. As we can see in the figure, most methods are based on CNNs, RNNs, or a combination of the two. Moreover, approaches can be separated based on the type of input data they require. A majority of DL ncRNA classifiers use the sequence as input, three use the secondary structure, and two use the sequence in combination with the structure or with epigenetic data.

## State-of-the-art datasets

Along with the tools from the previous section, multiple datasets have been proposed for evaluation of DL ncRNA classifiers in recent publications. We identify five.

The first dataset was proposed in 2017 by Fiannaca et al. [34]. Sequences were downloaded from Rfam 12 [47], and the CD-HIT tool [48] was used to remove redundant sequences. A random selection was carried out to obtain (almost-) balanced classes, with a final size of 8,920 sncRNAs spanning 13 classes. Since its proposal, this dataset has been used by all state-of-the-art methods. However, it presents a data leakage problem: 347 ncRNAs are present in both the training and the test set (representing respectively 5.5% and 13.3% of the sets).



**Fig 2. Overview of existing approaches for ncRNA classification.** Methods are represented in chronological order. We specify the type of input data used (secondary structure, sequence, or epigenetic data) and the type of DL architecture (CNN, RNN, both, or other).

<https://doi.org/10.1371/journal.pcbi.1012446.g002>

Boukelia et al. [39] constructed a second dataset in 2020. In total, 16 sncRNA classes are represented. This dataset contains 42,000 sequences of human ncRNAs downloaded from Rfam 13 [49]. It also includes four types of epigenetic modifications (from DeepBlue [50] and HHMD [51] databases). Specifically, they obtain the positions of DNA methylation, histone methylation (20 types), histone acetylation (18 types), and binding factors and regulators related to histones. For each RNA, and for each epigenetic modification, a matrix is constructed to indicate the nearest and farthest modifications, as well as the number of modifications.

The same year, Noviello et al. [36] proposed a dataset in which ncRNAs are classified into 88 Rfam families. These families are defined as groups with RNAs that are evolutionarily conserved, have evidence for a secondary structure, and have some known function. A total of 306,016 sequences are present in the dataset.

In 2023, a dataset was proposed by Lima et al. [41]. Employing a similar protocol as Fiannaca et al. [34] but with updated data from Rfam 14.3 [52], authors were able to multiply by five the number of ncRNAs in the dataset, obtaining 45,447 ncRNAs. 13 ncRNA classes are represented.

Sutanto et al. [46] also proposed a dataset in 2023. They obtained sequences from Rfam 14.1 [52] and removed those that could not be used by their structure prediction tool of choice. Then, CD-HIT [48] was used to remove redundancy in the dataset. In total, the dataset comprises 59,723 RNAs representing 12 classes.

The comparison between state-of-the-art datasets can be summed up to the origin of data, diversity of sequences, and classes represented:

- As one of the largest databases on ncRNAs, all datasets are obtained from Rfam, with recent releases containing more sequences. Boukelia et al.'s dataset [39] also integrates other databases for epigenetic data.
- In three cases [34, 41, 46], CD-HIT is used to remove sequences with similarity over 80%. This step restricts the size of the dataset, but only removes examples that would have been redundant while maintaining diversity. In some datasets, sequences containing degenerate nucleotides are removed [36, 41].
- Most datasets represent around 13 commonly-used sncRNA classes, grouping ncRNAs that share a function while not being too specific. Only Noviello et al.'s dataset [36] differs, using ncRNA families instead.

As the Rfam database is used for the construction of all datasets, it is important to note one of its limitations regarding ncRNA classes, as presented in [49]. Rfam associates classes to ncRNAs using the Infernal tool [6], which is designed to recognize sequences with high similarity. Pseudogenes, which are copies of functional genes that have mutated and become non-functional during evolution, can have high homology with the sequence they are derived from. There are no effective tools to differentiate one from the other. Therefore, some pseudogenes might be annotated as belonging to their class of origin. This is especially true in eukaryotes.

## Results

In this section, we compare the performance of different available state-of-the-art ncRNA classifiers: nRC [34], RNAGCN [45], ncrna-deep [36], MFpred [44], ncRDense [38], and NCYPred [41]. In the first subsection, their performance is measured using cross-validation and on a held-out test set on three datasets: Dataset1, Dataset1-nd and Dataset2. Dataset1 is the one by Fiannaca et al. [34], in which we fix the data leakage issue. Dataset1-nd is the same

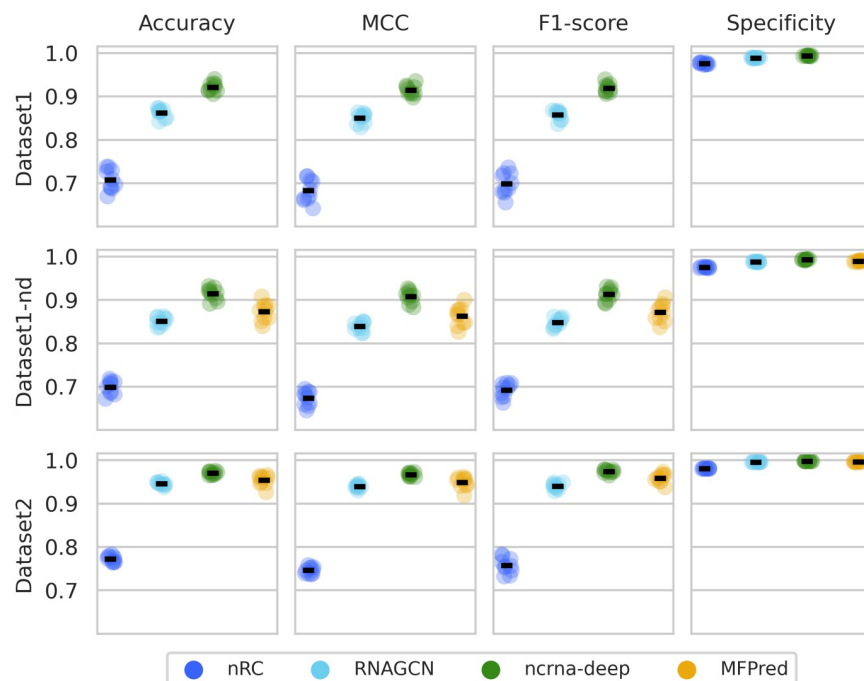
dataset in which we remove sequences containing degenerate nucleotides. Dataset2 is the one presented by Lima et al. [41] (see S1 Table). In the next two subsections, where we evaluate the robustness of methods, and measure resource intensity, we focus on Dataset1-nd as it can be used by all methods. The selection of tools and datasets for this benchmark, as well as the experimental protocol, are detailed in the Methods section.

## Overall performance

In this section, performance is measured using the following multi-class metrics: Accuracy, Matthews Correlation Coefficient (MCC), F1-score, and Specificity. The formulas are detailed in the Methods section; additional scores and numerical values are presented in S2 and S3 Tables.

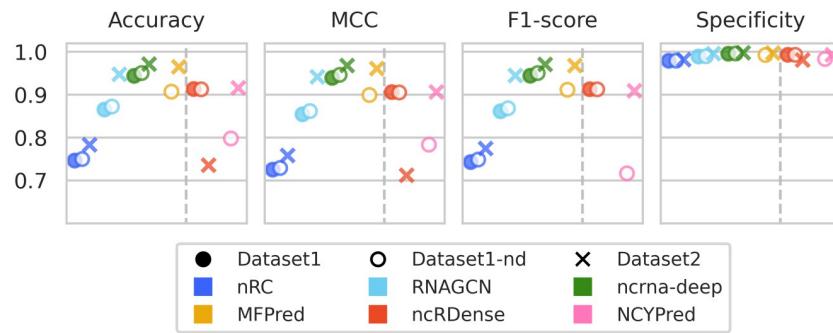
**10-fold cross-validation.** We perform 10-fold cross-validation with tools that can be retrained (nRC, RNAGCN, ncna-deep, and MFpred). Results are presented in Fig 3. Scores are, in general, better on Dataset2. This can be explained by the fact that Dataset2 is much larger than Dataset1(-nd). Models are trained on a larger variety of ncRNAs, making the prediction of new examples easier. Also of note, since Dataset1-nd is very similar to Dataset1 (it is only slightly smaller due to the removal of sequences with degenerate nucleotides), the close obtained scores are expected. We notice that ncna-deep reaches the highest scores across the different metrics and for the three datasets. MFpred and RNAGCN also obtain high scores on Dataset2, benefiting from the increased amount of data. nRC obtains lower scores across datasets, and its results vary more for the different folds compared to other tools.

**Performance on held-out test set.** We assess in Fig 4 the prediction performance of state-of-the-art ncRNA classifiers on held-out test sets. nRC, RNAGCN, ncna-deep and MFpred can be retrained, while ncRDense and NCYPred can only be used for prediction. This



**Fig 3. Cross-validation results.** Each row corresponds to a dataset, and each column to a metric. Each color corresponds to a method. Results on each validation set are represented by circles of the methods' colors, on top of which a line represents the mean.

<https://doi.org/10.1371/journal.pcbi.1012446.g003>

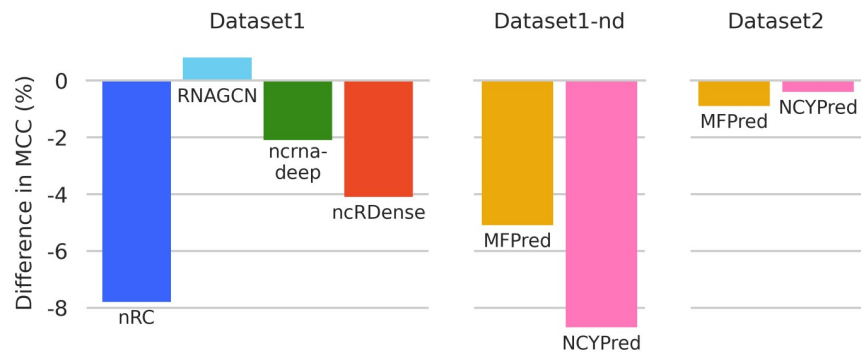


**Fig 4. Comparison of overall performance of the state-of-the-art methods on the test sets of Dataset1, Dataset1-nd and Dataset2.** The shape of the markers designates the datasets, and each color corresponds to a method. Tools on the right of the dashed line, ncRDense and NCYPred, could not be retrained before prediction. No results are shown for MFPred and NCYPred on Dataset1 since these two tools do not process degenerate nucleotides. The tools are sorted chronologically on each of the two sides separated by the dashed line.

<https://doi.org/10.1371/journal.pcbi.1012446.g004>

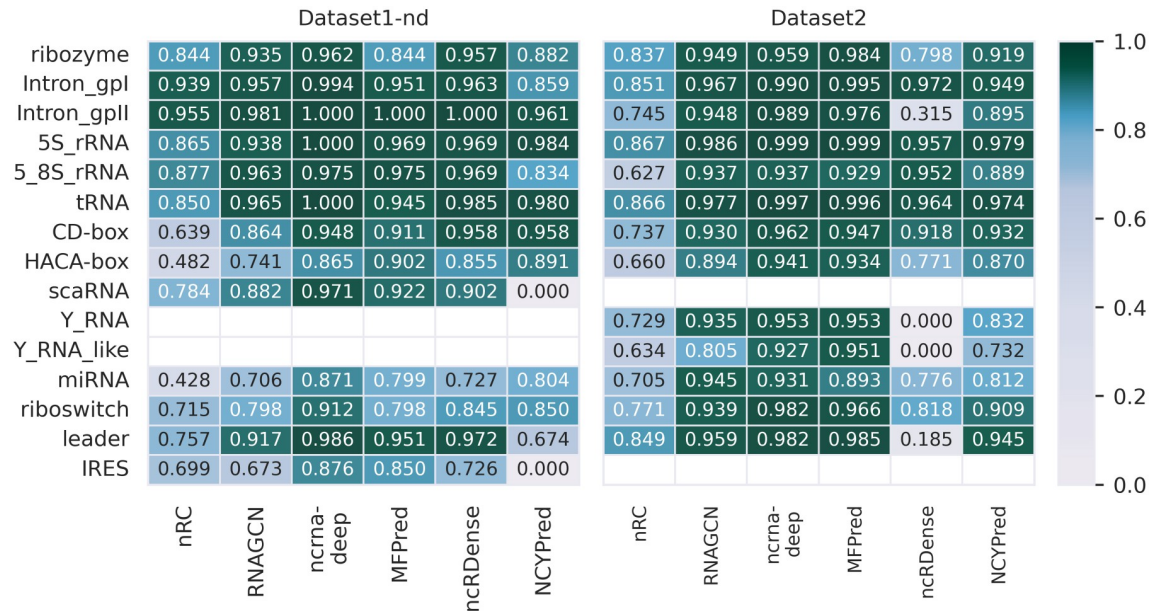
difference is visible in the results: the latter tools only perform well on the dataset they were originally trained on (Dataset1 for ncRDense and Dataset2 for NCYPred). This is particularly visible looking at the F1-score. Oppositely, tools that are retrained have consistent performances across datasets. The test scores are concordant with cross-validation results, with test scores similar to the cross-validation mean.

**Comparison with reported performance.** We compare the performance we measured for each tool to the performance reported in its original publication, and illustrate the difference in Fig 5. All publications use the same formula for MCC, therefore, we base our comparison on this metric. nRC, RNAGCN, ncrna-deep, and ncRDense all use Dataset1, and MFPred and NCYPred use Dataset1-nd (albeit the biased versions with data leakage). Results in our experiments are significantly worse than what was originally reported by the various publications. These drops in performance are partially explained by the correction of the data leakage problem. For NCYPred, our lower measured performance could also be due to model training. In their evaluations, authors retrained their tool on Dataset1-nd, which we cannot not do, only having access to the web server version. This means that two classes from Dataset1-nd are unknown by the model and cannot be predicted. One exception is RNAGCN, for which the model we selected during hyperparameter search performs better than the one from the original publication. Only two methods, MFPred and NCYPred, reported results on Dataset2 as it



**Fig 5. Difference between MCC values measured in our benchmark and those reported by each state-of-the-art method.** A negative value of  $-x$  signifies that the value we measured is lower by  $x$  than the one originally reported.

<https://doi.org/10.1371/journal.pcbi.1012446.g005>



**Fig 6. Comparison of accuracy of prediction of each ncRNA class obtained by state-of-the-art tools on Dataset1-nd and Dataset2.** Light colors correspond to lower accuracies, while colors tending towards dark green represent the best results. Results for Dataset1 are comparable to those for Dataset1-nd and are presented in S1 Fig.

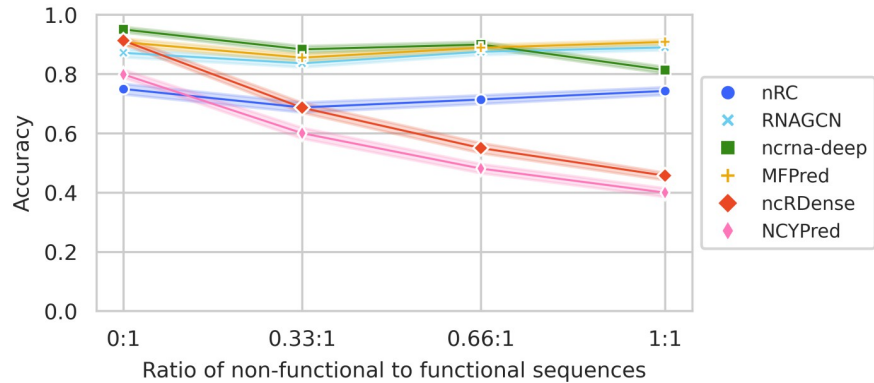
<https://doi.org/10.1371/journal.pcbi.1012446.g006>

is very recent. These reported results seem more reliable as we only observe negligible differences in MCC.

**Performance on individual classes.** A great challenge in ncRNA classification is achieving high predictions across classes. In Fig 6, we present the accuracy of predictions for each class. Here, another explanation for the better results on Dataset2 than Dataset1(-nd) that were noticed in the previous section is presented: the classes included in each dataset differ, which can impact performance. The two classes included solely in Dataset1(-nd) (IRES and scaRNA) seem more difficult to predict, thus lowering overall scores. We observe once again that highest performance is obtained by ncrna-deep, seconded by MFPred. The difference between classes is striking for nRC, which is able to predict some classes quite well (Intron-gpII, 5S-rRNA), but performs poorly on other classes (miRNA, HACA-box). Aside from the classes it was not trained on (IRES and scaRNA), NCYPred obtains similar performances to other tools on Dataset1-nd. The same cannot be said for ncRDense, which does not only fail to predict the classes it was not trained on (Y RNA and Y RNA-like), but also obtains relatively poor performances for most other classes in Dataset2, underscoring the fact that not allowing to retrain a tool can be detrimental. It is interesting that, although recent tools make generally very accurate predictions, the more problematic classes are the same as those for earlier methods, such as miRNA or HACA-box. Wrongly classified miRNAs tend to be predicted as CD-box (see S2, S3 and S4 Figs): there is literature explaining possible similarities between the two classes [53, 54]. HACA-box tends to be mispredicted as CD-box or scaRNA. These are subclasses of snoRNAs, which can explain the confusion.

### Robustness to noise

To measure the robustness of the tools to noise, we perform two experiments. In the first, we test whether models can reject non-functional sequences. In the second, we analyze the impact



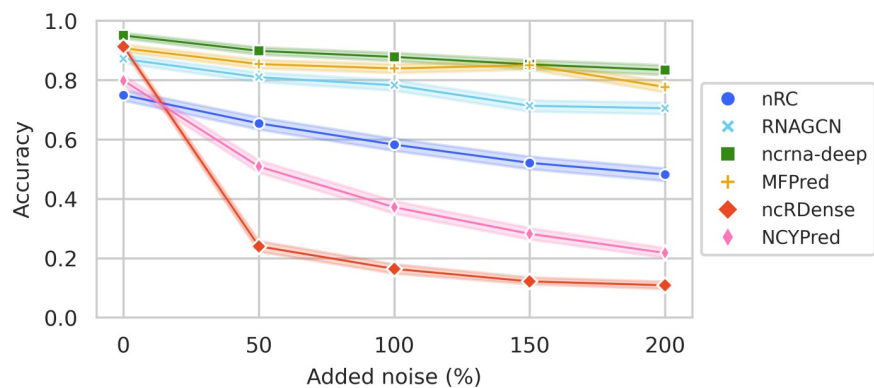
**Fig 7. Evolution of accuracy with different numbers of non-functional sequences added to Dataset1-nd.** Each line (and color) corresponds to a method. Accuracy is measured when varying the ratio of non-functional to functional sequences from 0:1 (initial dataset) to 1:1 (as many non-functional sequences as functional sequences).

<https://doi.org/10.1371/journal.pcbi.1012446.g007>

of sequence boundary noise. The python package *statkit* (<https://hylkedonker.gitlab.io/statkit/>) is used to compute 95% confidence intervals.

**Non-functional sequences.** We show in Fig 7 the evolution of accuracy when adding non-functional sequences to Dataset1-nd. The ratio of non-functional to functional sequences varies from 0:1, in which case there are no non-functional sequences, to 1:1, when the size of the dataset is doubled and there are as many functional ncRNAs as non-functional. We observe that MFpred, RNAGCN and nRC’s prediction accuracy remains relatively stable as the number of non-functional sequences increases. ncrna-deep’s accuracy decreases by around 13%, but performance is still relatively satisfying. This indicates that the tools are able to recognize the characteristics of functional ncRNAs. On the other hand, ncRDense and NCYPred, which cannot be retrained, are not able to recognize non-functional sequences. Their decrease in accuracy is proportional to the number of non-functional sequences added.

**Boundary noise.** We illustrate in Fig 8 the evolution of accuracy when adding noise to Dataset1-nd sequences. Noise is added to all sequences in the dataset; the amount of noise is determined as a percentage of sequence length. We expect a drop in accuracy as the amount of noise in sequences increases and there is proportionally less information relevant to the



**Fig 8. Evolution of accuracy with different percentages of added noise to Dataset1-nd sequences.** Each line (and color) corresponds to a method. Accuracy is measured when the noise added to sequences is equal to 0%, 50%, 100%, 150%, or 200% of the sequence length.

<https://doi.org/10.1371/journal.pcbi.1012446.g008>

**Table 1. Comparison of computation times and CO<sub>2</sub> emissions on Dataset1-nd.** The computation time is calculated for preprocessing, training and prediction, while the CO<sub>2</sub> emission is calculated for training. Times and emissions are computed for the hyperparameter sets selected in Table 4. (See S4 Table for Dataset1 and Dataset2).

	Computation time			Emissions (in gCO <sub>2</sub> eq)
	Preprocessing	Training	Prediction	
nRC	23mn 4s	1h 57mn	12s	-
RNAGCN	2mn 43s	27mn 10s	0.4s	9.08
ncrna-deep	2s	59s	0.6s	0.35
MFPred	32s	1h 55mn	11s	48.48

<https://doi.org/10.1371/journal.pcbi.1012446.t001>

classification task. ncrna-deep, MFPred and RNAGCN show robustness, as their accuracy only drops by around 10–15% when the noise added is equal to 200% of the sequence length. nRC is less robust, with a drop of around 25% in accuracy. ncrDense and NCYPred, which cannot be retrained, are not able to handle noisy sequences, as their accuracy drops significantly. ncrDense's performance is particularly poor, as even with the lowest amount of noise added, its accuracy drops to around 20%.

### Computation time and CO<sub>2</sub> emissions

We compare the computation time and CO<sub>2</sub> emissions of benchmark methods, which we measure when executing the codes. RNAGCN, ncrna-deep and MFPred are executed on the same GPU machine (NVIDIA GeForce RTX 3090). nRC only supports CPU computation and is executed on Intel Xeon(R) W-1250 CPU @ 3.30GHz. CO<sub>2</sub> emission estimations are conducted with CodeCarbon [55, 56]. We can use it for RNAGCN, ncrna-deep and MFPred, as source codes are all in Python, but not for nRC which is written in a mix of Java and Python.

Table 1 shows the computation time as well as the CO<sub>2</sub> emissions of the state-of-the-art methods on Dataset1-nd. nRC and MFPred have considerably longer runtimes than RNAGCN and ncrna-deep. For nRC, this can be explained by the method's lack of GPU support. For MFPred, the lack of speed is due to model complexity: the model is large and contains considerably more layers than the other methods. The CO<sub>2</sub> emissions of these methods seem to be proportional to training time.

For the web servers ncrDense and NCYPred, it is difficult to measure prediction time precisely, as it depends on multiple factors: network connection, user speed, etc. Moreover, ncrDense cannot predict more than 1,000 ncRNAs at once, further complicating the measure of prediction time. As an order of magnitude, for the same 1,000 ncRNAs, we obtained predictions in around 20 seconds with NCYPred and 30 seconds with ncrDense. No preprocessing is required. CO<sub>2</sub> emissions cannot be measured.

### Discussion

From the different results presented in the previous section, the current recommendation for users looking to classify ncRNAs would be to use the ncrna-deep tool, as it obtains the best results on benchmark datasets while being the fastest. MFPred, although more recent than ncrna-deep, performs slightly less well than the latter, and its execution time is longer. Moreover, it cannot be used on sequences containing degenerate nucleotides. This goes against the comparisons made in MFPred's publication, and emphasizes the necessity of unbiased benchmarks such as the one presented here, and the importance of employing comparable experimental protocols. nRC pioneered the use of DL for ncRNA classification and was the tool of reference for years, but is now outperformed by more recent tools. RNAGCN, also obtaining lower performances, brought to the field a different point of view as its GNN-based

architecture stands out from the rest of the state-of-the-art. If the goal is to classify new ncRNA datasets, we would advise against using web servers that do not provide the possibility to perform training, as it impacts performance negatively.

Several recommendations can be formulated with regard to the development of future ncRNA classifiers:

- The first choice to make in the architecture is how to represent data. nRC and RNAGCN choose to use the secondary structure. This is based on the widely accepted idea that the function of ncRNAs is linked to their structure [57], an approach adopted by many earlier ncRNA classifiers [31, 33]. It is interesting that this experimental observation does not translate particularly well to DL approaches, as methods that use the structure do not give the best results in our benchmark. This can be explained by the fact that secondary structure prediction is a difficult task, and while algorithms have improved over the years, their accuracy is still limited [58, 59]. Using newer, more accurate structure prediction algorithms could be a way to improve results in the future. Moreover, the difference in performance could also be due to the algorithms used for classification.
- Among methods that directly use the sequence, two representation approaches emerge in state-of-the-art methods: using a simple matrixial representation of the sequence, or learning a meaningful representation of the sequence with RNN-based networks. These two approaches are represented in our benchmark, with *ncrna-deep* using a one-hot encoding of the sequence, and MFpred using four RNNs to extract different information from the sequence. Representation complexity is low for *ncrna-deep*, high for MFpred, but *ncrna-deep* performs better. The results obtained in our benchmark suggest that the simple representation should in fact be preferred, as the time gain is considerable without any effect on performance.
- While results have shown that the sequence is enough to obtain good classification results, this should not be a limit in the development of new ncRNA methods. For example, earlier ncRNA classifiers [30, 32, 60, 61] used RNA-seq data. The inclusion of epigenetic data by imCnC [39] also seemed promising. Indeed, while performance is an important aspect of ncRNA classification, an additional benefit in research could be to obtain a better characterization of classes—which includes describing them by patterns found at multiple levels, not just in the sequence. As stated in the introduction, the description of ncRNA classes and their functions is an ongoing research, in which computational ncRNA classifiers could be extremely useful.
- The set of ncRNA classes is bound to evolve as new ncRNA functions are discovered. On this basis, one aspect to include in future ncRNA classifiers is the ability to predict at test time classes that were unknown during training. The idea would therefore be to design a model that can identify that a sample does not belong to any of the known classes, instead of forcing classification. This could be done with zero-shot learning or abstained classification [62–64].
- Finally, there is an additional consideration regarding the availability of tools. In the description of state-of-the-art datasets, we noted an increase in dataset size and variability in the classes included. In this context, it is important to propose tools that can be applied to and evaluated on new datasets. The best practice would be to propose both a web server and to share source code. The web server would provide ease of use. By offering the possibility of retraining the model, it would not be limited to one dataset. Moreover, a well-documented source code could be used in bigger projects and with more flexibility.

## Conclusion

In this paper, we have reviewed key advancements in DL-based ncRNA classification over the past seven years. We first examined the different architectures implemented, noting that most tools are built on CNNs and/or RNNs. The predominant use of the sequence to represent ncRNAs was highlighted. We also provided a detailed comparison of recent datasets developed for ncRNA classification.

We then conducted an extensive benchmark of state-of-the-art tools across three datasets. Among the tools evaluated, ncrna-deep, MFPred, and RNAGCN stood out, with the first two utilizing the sequence and the latter leveraging the secondary structure. These three tools demonstrated robustness to noise. In particular, ncrna-deep consistently outperformed others in terms of performance and speed. Our benchmark also underscored the significant drawback of tools that cannot be retrained, as they showed poor performance on new datasets. This finding emphasizes the importance of allowing users to retrain models.

Finally, we provided recommendations for the future of ncRNA classification. While secondary structure-based methods did not perform best in our evaluations, their potential could be significantly enhanced with the development of more accurate structure prediction tools. We emphasized the effectiveness of simple sequence representations, and also encouraged the integration of more diverse data types, as this could improve performance and provide new insights. Additionally, we suggested that incorporating techniques like zero-shot learning or abstained classification could be beneficial for the field. Lastly, we underscored the importance of developing tools that can be retrained on new datasets, ensuring their adaptability and continued relevance to evolving data.

## Methods

### Literature search

The scope of our study is multi-class non-coding RNA classifiers that are based on deep learning. To make sure we included all relevant literature, we performed the following search from Digital Science's Dimensions platform (<https://app.dimensions.ai>): the terms (*ncRNA OR non-coding RNA*) AND (*deep learning*) AND (*classification OR classes OR families*) had to be present in the title and/or abstract. On February 14th, 2024, this search yielded 567 results. From these, a lot of results were not directly linked to our problem statement (e.g. a lot of publications were about lncRNA-disease association prediction, or cancer patient classification). We carefully checked that no relevant title had been forgotten from our study.

### Tool selection

We present in [Table 2](#) information on the accessibility of all DL-based ncRNA classifiers. Most publications propose an online version of their tool, as a web server or by sharing source code. Web servers have the advantage of requiring no installation and being easy to use. However, they are sometimes limited: for example, here, no web server can be retrained on different datasets. If source code is available, the model can be trained on different datasets, and computational power can be increased. Nonetheless, executing the code is sometimes complicated, especially if documentation is insufficient.

We choose to assess the quality of prediction of six methods: nRC [34], RNAGCN [45], ncrna-deep [36], ncRDense [38], NCYPred [41] and MFPred [44]. Graph+DL [46] and RPC-snRC [37] are not available and thus cannot be tested. For a fair comparison, we need to be able to compare tools on the same datasets—this cannot be done for imCnC [39], the only tool to require epigenetic data, thus it is not included in this comparison. Finally, for tools that

**Table 2. Description of the availability and ease-of-use of DL ncRNA classification tools.** Tools are presented in chronological order of publication. Each cell indicates if the tool is accessible through a web server, if source code can be downloaded, and if the tool can be retrained. All links have been checked at the time of writing (February 2024).

Method	Usage
nRC [34]	<b>Web server:</b> X – <b>Source code:</b> ✓ – <b>Training:</b> ✓ Docker container: <a href="https://hub.docker.com/r/tblab/nrc/">https://hub.docker.com/r/tblab/nrc/</a> , GitHub repository: <a href="https://github.com/IcarPA-TBlab/nrc">https://github.com/IcarPA-TBlab/nrc</a> . The tool is easy to use: no pre-formatting of data is required, the commands are documented and an example is given.
RNAGCN [45]	<b>Web server:</b> X – <b>Source code:</b> ✓ – <b>Training:</b> ✓ GitHub repository: <a href="https://github.com/emalgorithm/ncRNA-family-prediction">https://github.com/emalgorithm/ncRNA-family-prediction</a> . The method is easy to use, as the process is documented. However, to use the tool on a new dataset, an outside tool has to be used to obtain secondary structures, which then need to be formatted correctly before using RNAGCN.
ncRFP [40]	<b>Web server:</b> X – <b>Source code:</b> ✓ – <b>Training:</b> ✓ Though not mentioned in the article, a GitHub repository exists: <a href="https://github.com/linyuwangPHD/ncRFP">https://github.com/linyuwangPHD/ncRFP</a> . Regarding ease of use, the documentation provided is short and does not clearly show which commands to run. Variables and hyperparameters are fixed inside the code and cannot be easily modified.
imCnC [39]	<b>Web server:</b> X – <b>Source code:</b> ✓ – <b>Training:</b> ✓ GitHub repository: <a href="https://github.com/BoukeliaAbdelbasset/imCnC">https://github.com/BoukeliaAbdelbasset/imCnC</a> . No instructions or examples of how to use the tool are given.
ncRDeep [35]	<b>Web server:</b> ✓ – <b>Source code:</b> X – <b>Training:</b> X Web server: <a href="https://nscbio.jbnu.ac.kr/tools/ncRDeep/">https://nscbio.jbnu.ac.kr/tools/ncRDeep/</a> . Although not mentioned on the site, there is a limit of 1,000 ncRNA per prediction task. Moreover, it is not possible to export results.
ncrna-deep [36]	<b>Web server:</b> X – <b>Source code:</b> ✓ – <b>Training:</b> ✓ GitHub repository: <a href="https://github.com/bioinformatics-sannio/ncrna-deep">https://github.com/bioinformatics-sannio/ncrna-deep</a> . Dataset preparation has to be done separately, but the process is explained and illustrated in the documentation and in the code.
ncRDense [38]	<b>Web server:</b> ✓ – <b>Source code:</b> X – <b>Training:</b> X Web server: <a href="https://nscbio.jbnu.ac.kr/tools/ncRDense/">https://nscbio.jbnu.ac.kr/tools/ncRDense/</a> . Although not mentioned on the site, there is a limit of 1,000 ncRNA per prediction task. Moreover, it is not possible to export results.
ncDLRES [42]	<b>Web server:</b> X – <b>Source code:</b> ✓ – <b>Training:</b> ✓ GitHub repository: <a href="https://github.com/linyuwangPHD/ncDLRES">https://github.com/linyuwangPHD/ncDLRES</a> . No examples of how to use the tool, but very basic instructions are given.
RPC-snRC [37]	<b>Web server:</b> X – <b>Source code:</b> X – <b>Training:</b> X The GitHub repository linked in the paper no longer exists.
NCYPred [41]	<b>Web server:</b> ✓ – <b>Source code:</b> ✓ – <b>Training:</b> ✓ Web server: <a href="https://www.gpea.uem.br/ncypred">https://www.gpea.uem.br/ncypred</a> , GitHub repository: <a href="https://github.com/diegodslima/NCYPred">https://github.com/diegodslima/NCYPred</a> . The web server is easy to use and results can be exported. Many elements are given in the source code, but due to missing files, the code cannot be used.
ncDENSE [43]	<b>Web server:</b> X – <b>Source code:</b> ✓ – <b>Training:</b> ✓ GitHub repository: <a href="https://github.com/ck-fighting/ncDENSE">https://github.com/ck-fighting/ncDENSE</a> . No examples of how to use the tool are given.
MFPred [44]	<b>Web server:</b> X – <b>Source code:</b> ✓ – <b>Training:</b> ✓ GitHub repository: <a href="https://github.com/ck-fighting/MFPred">https://github.com/ck-fighting/MFPred</a> . Dataset preparation has to be done separately with scripts provided. Many files are in the repository, but the documentation is lacking, making it difficult to know how to use the tool.
Graph+DL [46]	<b>Web server:</b> X – <b>Source code:</b> X – <b>Training:</b> X Source code is not linked. Data is not available from the proposed link in the publication.

<https://doi.org/10.1371/journal.pcbi.1012446.t002>

were developed by the same group, ample comparisons between old and new versions are already available. This concerns on one hand ncRDeep [35] and ncRDense [38], and on the other hand ncRFP [40], ncDLRES [42], ncDENSE [43] and MFPred [44]. In both cases, we use the latest version: ncRDense and MFPred.

**Table 3. Description of datasets for ncRNA classification.** (a) State-of-the-art datasets. (b) Datasets used in this study. Datasets are sorted by date. The number of instances and classes in each dataset is given. We also give the range of sequence lengths. The last column indicates how many times the dataset has been used by state-of-the-art ncRNA classifiers. For some datasets information on length range is not available (indicated with 'n/a') as it is not stated in the publication, and it cannot be computed as the dataset is not available.

	Source	Nb Instances	Nb Classes	Lengths	Nb Uses
(a)	Fiannaca et al. [34]	8,920	13	38–1182	12
	Boukelia et al. [39]	42,000	16	n/a	1
	Noviello et al. [36]	306,016	88	1–80	1
	Lima et al. [41]	45,447	13	42–500	2
	Sutanto et al. [46]	59,723	12	n/a	1
(b)	Dataset1	8,573	13	38–1182	-
	Dataset1-nd	8,350	13	38–1182	-
	Dataset2	45,447	13	42–500	-

<https://doi.org/10.1371/journal.pcbi.1012446.t003>

## Dataset selection

Out of the five state-of-the-art datasets described previously, three have been made available publicly: Fiannaca et al.'s dataset [34], Noviello et al.'s dataset [36], and Lima et al.'s dataset [41]. We present in Table 3(a) information regarding the size of the datasets, as well as the number of times they were used in other publications.

We decide to focus on Fiannaca et al.'s dataset and Lima et al.'s dataset. Noviello et al. present their dataset with the goal of predicting ncRNA families. As such, the number of classes is high, and they can describe very precise functions. We consider these classes to be too specific for our presented problem of ncRNA classification into broader, more general functional groups (or classes). In our experiments, we refer to three datasets, also presented in Table 3(b):

- Dataset1: the dataset by Fiannaca et al. in which we fix the data leakage issue. Specifically, we remove from the test set the 347 sequences that were also present in the training set. The training and test sets contain respectively 6,320 and 2,253 ncRNAs.
- Dataset1-nd: a version of Dataset1 in which we remove sequences containing degenerate nucleotides, as these cannot be handled by some methods. The training and test sets contain respectively 6,161 and 2,189 ncRNAs.
- Dataset2: the dataset by Lima et al., in which we confirm there is no data leakage. It does not include any sequence with degenerate nucleotides. The training and test sets contain respectively 31,801 and 13,646 ncRNAs. This is one of the most complete datasets of current state-of-the-art, due to its large number of samples and the fact that it covers most well-known ncRNA classes.

## Experimental protocol

**Evaluation metrics.** Performance is measured with Accuracy, MCC (Matthews Correlation Coefficient), F1-score, and Specificity. The last two metrics are binary metrics for which we use the macro averaging generalization for multi-class problems, where scores are computed on each class separately then averaged. The four scores are defined as follows:

$$Accuracy = \frac{\sum_{c=1}^C TP_c}{S}$$

$$\text{MCC} = \frac{(\sum_{c=1}^C TP_c) * S - \sum_{c=1}^C (p_c * t_c)}{\sqrt{(S^2 - \sum_{c=1}^C p_c^2)(S^2 - \sum_{c=1}^C t_c^2)}}$$

$$\text{F1-score} = \frac{1}{C} \sum_{c=1}^C \frac{TP_c}{TP_c + 0.5(FP_c + FN_c)}$$

$$\text{Specificity} = \frac{1}{C} \sum_{c=1}^C \frac{TN_c}{TN_c + FP_c}$$

where  $C$  is the number of classes and  $S$  the number of samples. Respective of class  $c$ ,  $TP_c$  are True Positives,  $FP_c$  are False Positives,  $TN_c$  are True Negatives and  $FN_c$  are False Negatives.  $p_c$  is the number of times class  $c$  was predicted, and  $t_c$  is the true number of members of class  $c$ .

**Model training.** nRC, RNAGCN, ncRNA-deep and MFPreD are retrained before predicting each test set. When possible, we test different hyperparameters, which are presented in Table 4. This search cannot not be done for MFPreD as parameter values are fixed directly in the code, and no guidelines are given on how to change them or what to change. If not stated, default hyperparameter values are used. For nRC, the model is trained for 100 epochs, with a batch size of 10. SGD optimizer is used with a learning rate of 0.05. For RNAGCN, the maximum number of epochs is set to 10,000, using the early stopping option proposed by the authors (with a patience of 100 epochs)—in general, the model stops training after 700 epochs. The batch size is 64. A learning rate of 0.0004 is used for the Adam optimizer. For ncRNA-deep, 25 epochs are used, with a batch size of 32. The AMSGrad variant of the Adam optimizer is used with a learning rate of 0.0005. For MFPreD, the model is trained for 100 epochs, and the batch size is set to 32. The Adam optimizer is used with a learning rate of 0.0001 and a weight decay of 0.0001. ncRDense and NCYPred are only accessible to us as web servers that cannot not be retrained, and are only used for prediction.

**Dataset partitioning.** For 10-fold cross-validation, 10% of the training set is used as a validation set. This process is repeated 10 times, with each sample appearing exactly once in a validation set. The folds are designed in order to respect class distribution, and we ensure that there is no significant difference in terms of sequence length and presence of degenerate nucleotides. Then, models are re-trained on the full training set, and predictions are obtained on the held-out test set. MFPreD and NCYPred cannot handle degenerate nucleotides, therefore their performance is not measured on Dataset1.

**Introduction of noise.** To assess whether models can reject non-functional sequences, we introduce noise by adding a new class in the dataset that contains non-functional sequences. We verify that they do not present similarities to known ncRNA sequences in the database Rfam [52]. In practice, we add to the dataset new sequences that are generated randomly,

**Table 4. Overview of tested hyperparameters.** Parameters that gave the best results on Dataset1 and Dataset2 are denoted by <sup>(1)</sup> and <sup>(2)</sup> respectively. The models chosen on Dataset1 were also used for Dataset1-nd.

Method	Hyperparameters
nRC	Range of size of substructures {(2–4), (3–5) <sup>(1,2)</sup> }, number of kernels in the first {10 <sup>(2)</sup> , 20 <sup>(1)</sup> } and the second {10, 20 <sup>(1,2)</sup> } convolutional layers.
RNAGCN	Size {64 <sup>(2)</sup> , 80 <sup>(1)</sup> , 128} and type {MPNN, GIN <sup>(2)</sup> , GCN, GAT <sup>(1)</sup> } of graph convolutions.
ncRNA-deep	Sequence encoding: {1-mer <sup>(1,2)</sup> , 2-mer, snake, hilbert, morton}.

<https://doi.org/10.1371/journal.pcbi.1012446.t004>

respecting di-nucleotide distribution of existing sequences. We vary the ratio of non-functional to functional sequences in {0:1, 0.33:1, 0.66:1, 1:1}. We also study the impact of sequence boundary errors by adding noise to the beginning and end of dataset sequences, while respecting mono- and di-nucleotide distribution in the sequence. The amount of noise added is proportional to sequence length: 0%, 50%, 100%, 150% and 200%. These protocols are inspired from [7, 36].

## Supporting information

**S1 Table. Description of ncRNA classification datasets (Dataset1 and Dataset2).**

(PDF)

**S2 Table. 10-fold cross-validation mean and standard deviation of different metrics.**

(PDF)

**S3 Table. Evaluation of different metrics on held-out test sets.**

(PDF)

**S4 Table. Comparison of computation times and CO<sub>2</sub> emissions on Dataset1 and Dataset2.**

(PDF)

**S1 Fig. Comparison of accuracy of prediction of each ncRNA class obtained by state-of-the-art tools on Dataset1 and Dataset1-nd.**

(PDF)

**S2 Fig. Confusion matrices on Dataset1.**

(PDF)

**S3 Fig. Confusion matrices on Dataset1-nd.**

(PDF)

**S4 Fig. Confusion matrices on Dataset2.**

(PDF)

## Author Contributions

**Conceptualization:** Constance Creux, Farida Zehraoui, Fariza Tahi.

**Data curation:** Constance Creux.

**Formal analysis:** Constance Creux.

**Funding acquisition:** Fariza Tahi.

**Investigation:** Constance Creux.

**Project administration:** Farida Zehraoui, Fariza Tahi.

**Software:** Constance Creux.

**Supervision:** Farida Zehraoui, François Radvanyi, Fariza Tahi.

**Visualization:** Constance Creux.

**Writing – original draft:** Constance Creux.

**Writing – review & editing:** Farida Zehraoui, François Radvanyi, Fariza Tahi.

## References

1. Morillon A. Non-Coding RNA, Its History and Discovery Timeline. In: Long Non-coding RNA. Elsevier; 2018. p. 25–53.
2. Lander ES. Initial Impact of the Sequencing of the Human Genome. *Nature*. 2011; 470(7333):187–197. <https://doi.org/10.1038/nature09792> PMID: 21307931
3. Kim T, Croce CM. MicroRNA: Trends in Clinical Trials of Cancer Diagnosis and Therapy Strategies. *Exp Mol Med*. 2023; 55(7):1314–1321. <https://doi.org/10.1038/s12276-023-01050-9> PMID: 37430087
4. Zhou Z, Wang Z, Gao J, Lin Z, Wang Y, Shan P, et al. Noncoding RNA-mediated Macrophage and Cancer Cell Crosstalk in Hepatocellular Carcinoma. *Molecular Therapy—Oncolytics*. 2022; 25:98–120. <https://doi.org/10.1016/j.omto.2022.03.002> PMID: 35506150
5. Chen X, Luo R, Zhang Y, Ye S, Zeng X, Liu J, et al. Long Noncoding RNA DIO3OS Induces Glycolytic-Dominant Metabolic Reprogramming to Promote Aromatase Inhibitor Resistance in Breast Cancer. *Nat Commun*. 2022; 13(1):7160. <https://doi.org/10.1038/s41467-022-34702-x> PMID: 36418319
6. Nawrocki EP, Eddy SR. Infernal 1.1: 100-Fold Faster RNA Homology Searches. *Bioinformatics*. 2013; 29(22):2933–2935. <https://doi.org/10.1093/bioinformatics/btt509> PMID: 24008419
7. Navarin N, Costa F. An Efficient Graph Kernel Method for Non-Coding RNA Functional Prediction. *Bioinformatics*. 2017; 33(17):2642–2650. <https://doi.org/10.1093/bioinformatics/btx295> PMID: 28475710
8. Dupont MJ, Major F. D-ORB: A Web Server to Extract Structural Features of Related But Unaligned RNA Sequences. *J Mol Biol*. 2023; 435(15):168181. <https://doi.org/10.1016/j.jmb.2023.168181> PMID: 37468182
9. Berg MD, Brandl CJ. Transfer RNAs: Diversity in Form and Function. *RNA Biology*. 2021; 18(3):316–339. <https://doi.org/10.1080/15476286.2020.1809197> PMID: 32900285
10. Fromm B, Høye E, Domanska D, Zhong X, Aparicio-Puerta E, Ovchinnikov V, et al. MirGeneDB 2.1: Toward a Complete Sampling of All Major Animal Phyla. *Nucleic Acids Res*. 2021; 50(D1):D204–D210. <https://doi.org/10.1093/nar/gkab1101>
11. Naeli P, Winter T, Hackett AP, Alboushi L, Jafarnejad SM. The Intricate Balance between microRNA-induced mRNA Decay and Translational Repression. *The FEBS Journal*. 2023; 290(10):2508–2524. <https://doi.org/10.1111/febs.16422> PMID: 35247033
12. Mattick JS, Amaral PP, Carninci P, Carpenter S, Chang HY, Chen LL, et al. Long Non-Coding RNAs: Definitions, Functions, Challenges and Recommendations. *Nat Rev Mol Cell Biol*. 2023; 24(6):430–447. <https://doi.org/10.1038/s41580-022-00566-8> PMID: 36596869
13. Brayet J, Zehraoui F, Jeanson-Leh L, Israeli D, Tahri F. Towards a piRNA Prediction Using Multiple Kernel Fusion and Support Vector Machine. *Bioinformatics*. 2014; 30(17):i364–i370. <https://doi.org/10.1093/bioinformatics/btu441> PMID: 25161221
14. Tran VDT, Tempel S, Zerath B, Zehraoui F, Tahri F. miRBoost: Boosting Support Vector Machines for microRNA Precursor Classification. *RNA*. 2015; 21(5):775–785. <https://doi.org/10.1261/rna.043612.113> PMID: 25795417
15. Tav C, Tempel S, Poligny L, Tahri F. miRNAfold: A Web Server for Fast miRNA Precursor Prediction in Genomes. *Nucleic Acids Res*. 2016; 44(W1):W181–W184. <https://doi.org/10.1093/nar/gkw459> PMID: 27242364
16. Boucheham A, Somnard V, Zehraoui F, Boualem A, Batouche M, Bendahmane A, et al. IpiRId: Integrative Approach for piRNA Prediction Using Genomic and Epigenomic Data. *PLoS ONE*. 2017; 12(6): e0179787. <https://doi.org/10.1371/journal.pone.0179787> PMID: 28622364
17. Baek J, Lee B, Kwon S, Yoon S. LncRNA-net: Long Non-Coding RNA Identification Using Deep Learning. *Bioinformatics*. 2018; 34(22):3889–3897. <https://doi.org/10.1093/bioinformatics/bty418> PMID: 29850775
18. Liu Y, Ding Y, Li A, Fei R, Guo X, Wu F. Prediction of Exosomal piRNAs Based on Deep Learning for Sequence Embedding with Attention Mechanism. In: 2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM); 2022. p. 158–161.
19. Li M, Liang C. LncDC: A Machine Learning-Based Tool for Long Non-Coding RNA Detection from RNA-Seq Data. *Sci Rep*. 2022; 12(1):19083. <https://doi.org/10.1038/s41598-022-22082-7> PMID: 36351980
20. Raad J, Bugnon LA, Milone DH, Stegmayer G. miRe2e: A Full End-to-End Deep Model Based on Transformers for Prediction of Pre-miRNAs. *Bioinformatics*. 2022; 38(5):1191–1197. <https://doi.org/10.1093/bioinformatics/btab823> PMID: 34875006
21. Postic G, Tav C, Platon L, Zehraoui F, Tahri F. IRSOM2: A Web Server for Predicting Bifunctional RNAs. *Nucleic Acids Research*. 2023; 51(W1):W281–W288. <https://doi.org/10.1093/nar/gkad381> PMID: 37158254

22. Rajendran V, Anandaram H, Kumar S S, Soman KP, Dhivya S. A Comparative Analysis of Machine Learning and Deep Learning Approaches for Circular RNA Classification. In: 2023 6th International Conference on Contemporary Computing and Informatics (IC3I). vol. 6; 2023. p. 1026–1034.
23. Krzyzanowski PM, Muro EM, Andrade-Navarro MA. Computational Approaches to Discovering Non-coding RNA. *WIREs RNA*. 2012; 3(4):567–579. <https://doi.org/10.1002/wrna.1121> PMID: 22555938
24. Amin N, McGrath A, Chen YPP. Evaluation of Deep Learning in Non-Coding RNA Classification. *Nat Mach Intell*. 2019; 1(5):246–256. <https://doi.org/10.1038/s42256-019-0051-2>
25. Singh D, Roy J. A Large-Scale Benchmark Study of Tools for the Classification of Protein-Coding and Non-Coding RNAs. *Nucleic Acids Research*. 2022; 50(21):12094–12111. <https://doi.org/10.1093/nar/gkac1092> PMID: 36420898
26. Ammunét T, Wang N, Khan S, Elo LL. Deep Learning Tools Are Top Performers in Long Non-Coding RNA Prediction. *Briefings in Functional Genomics*. 2022; 21(3):230–241. <https://doi.org/10.1093/bfpg/elab045> PMID: 35136929
27. Zhang Y, Huang H, Zhang D, Qiu J, Yang J, Wang K, et al. A Review on Recent Computational Methods for Predicting Noncoding RNAs. *BioMed Research International*. 2017; 2017:e9139504. <https://doi.org/10.1155/2017/9139504> PMID: 28553651
28. Gruber AR, Neuböck R, Hofacker IL, Washietl S. The RNAz Web Server: Prediction of Thermodynamically Stable and Evolutionarily Conserved RNA Structures. *Nucleic Acids Res*. 2007; 35(Web Server issue):W335–W338. <https://doi.org/10.1093/nar/gkm222> PMID: 17452347
29. Lindgreen S, Gardner PP, Krogh A. MASTR: Multiple Alignment and Structure Prediction of Non-Coding RNAs Using Simulated Annealing. *Bioinformatics*. 2007; 23(24):3304–3311. <https://doi.org/10.1093/bioinformatics/btm525> PMID: 18006551
30. Yuan C, Sun Y. RNA-CODE: A Noncoding RNA Classification Tool for Short Reads in NGS Data Lacking Reference Genomes. *PLOS ONE*. 2013; 8(10):e77596. <https://doi.org/10.1371/journal.pone.0077596> PMID: 24204885
31. Karklin Y, Meraz RF, Holbrook SR. Classification of Non-Coding RNA Using Graph Representations of Secondary Structure. *Pac Symp Biocomput*. 2005; p. 4–15. PMID: 15759609
32. Leung YY, Ryvkin P, Ungar LH, Gregory BD, Wang LS. CoRAL: Predicting Non-Coding RNAs from Small RNA-sequencing Data. *Nucleic Acids Res*. 2013; 41(14):e137. <https://doi.org/10.1093/nar/gkt426> PMID: 23700308
33. Panwar B, Arora A, Raghava GP. Prediction and Classification of ncRNAs Using Structural Information. *BMC Genomics*. 2014; 15(1). <https://doi.org/10.1186/1471-2164-15-127> PMID: 24521294
34. Fiannaca A, La Rosa M, La Paglia L, Rizzo R, Urso A. NRC: Non-coding RNA Classifier Based on Structural Features. *BioData Mining*. 2017; 10(1). <https://doi.org/10.1186/s13040-017-0148-2> PMID: 28785313
35. Chantsalnym T, Lim DY, Tayara H, Chong KT. ncrDeep: Non-coding RNA Classification with Convolutional Neural Network. *Computational Biology and Chemistry*. 2020; 88. <https://doi.org/10.1016/j.compbiolchem.2020.107364> PMID: 32890916
36. Noviello TMR, Ceccarelli F, Ceccarelli M, Cerulo L. Deep Learning Predicts Short Non-Coding RNA Functions from Only Raw Sequence Data. *PLoS Comput Biol*. 2020; 16(11). <https://doi.org/10.1371/journal.pcbi.1008415> PMID: 33175836
37. Asim MN, Malik MI, Zehe C, Trygg J, Dengel A, Ahmed S. A Robust and Precise ConvNet for Small Non-Coding RNA Classification (RPC-snRC). *IEEE Access*. 2021; 9:19379–19390. <https://doi.org/10.1109/ACCESS.2020.3037642>
38. Chantsalnym T, Siraj A, Tayara H, Lu N, Chong KT. ncrDense: A Novel Computational Approach for Classification of Non-Coding RNA Family by Deep Learning. *Genomics*. 2021; 113(5):3030–3038. <https://doi.org/10.1016/j.ygeno.2021.07.004> PMID: 34242708
39. Boukelia A, Boucheham A, Belguidoum M, Batouche M, Zehraoui F, Tah F. A Novel Integrative Approach for Non-coding RNA Classification Based on Deep Learning. *Current Bioinformatics*. 2020; 15(4):338–348. <https://doi.org/10.2174/1574893614666191105160633>
40. Wang L, Zheng S, Zhang H, Qiu Z, Zhong X, Liu H, et al. ncrFP: A Novel End-to-End Method for Non-Coding RNAs Family Prediction Based on Deep Learning. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. 2020.
41. Lima DdS, Amichi LJA, Fernandez MA, Constantino AA, Seixas FAV. NCYPred: A Bidirectional LSTM Network With Attention for Y RNA and Short Non-Coding RNA Classification. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. 2023; 20(1):557–565. <https://doi.org/10.1109/TCBB.2021.3131136> PMID: 34826297

42. Wang L, Zhong Xd, Wang S, Liu Y. ncDLRES: A Novel Method for Non-Coding RNAs Family Prediction Based on Dynamic LSTM and ResNet. *BMC Bioinformatics*. 2021; 22(1). <https://doi.org/10.1186/s12859-021-04495-9>
43. Chen K, Zhu X, Wang J, Hao L, Liu Z, Liu Y. ncDENSE: A Novel Computational Method Based on a Deep Learning Framework for Non-Coding RNAs Family Prediction. *BMC Bioinformatics*. 2023; 24(1). <https://doi.org/10.1186/s12859-023-05191-6> PMID: 36849908
44. Chen K, Zhu X, Wang J, Zhao Z, Hao L, Guo X, et al. MFPred: Prediction of ncRNA Families Based on Multi-Feature Fusion. *Briefings in Bioinformatics*. 2023. <https://doi.org/10.1093/bib/bbad303> PMID: 37615358
45. Rossi E, Monti F, Bronstein M, Liò P. NcRNA Classification with Graph Convolutional Networks. In: *Proceedings of DLG@KDD Workshop*. Anchorage, Alaska, US; 2019. p. 17–21.
46. Sutanto K, Turcotte M. Assessing Global-Local Secondary Structure Fingerprints to Classify RNA Sequences With Deep Learning. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. 2023; 20(5):2736–2747. <https://doi.org/10.1109/TCBB.2021.3118358> PMID: 34633933
47. Nawrocki EP, Burge SW, Bateman A, Daub J, Eberhardt RY, Eddy SR, et al. Rfam 12.0: Updates to the RNA Families Database. *Nucleic Acids Res*. 2015; 43(Database issue):D130–D137. <https://doi.org/10.1093/nar/gku1063> PMID: 25392425
48. Li W, Godzik A. Cd-Hit: A Fast Program for Clustering and Comparing Large Sets of Protein or Nucleotide Sequences. *Bioinformatics*. 2006; 22(13):1658–1659. <https://doi.org/10.1093/bioinformatics/btl158> PMID: 16731699
49. Kalvari I, Argasinska J, Quinones-Olvera N, Nawrocki EP, Rivas E, Eddy SR, et al. Rfam 13.0: Shifting to a Genome-Centric Resource for Non-Coding RNA Families. *Nucleic Acids Res*. 2018; 46(Database issue):D335–D342. <https://doi.org/10.1093/nar/gkx1038> PMID: 29112718
50. Albrecht F, List M, Bock C, Lengauer T. DeepBlue Epigenomic Data Server: Programmatic Data Retrieval and Analysis of Epigenome Region Sets. *Nucleic Acids Res*. 2016; 44(Web Server issue):W581–W586. <https://doi.org/10.1093/nar/gkw211> PMID: 27084938
51. Zhang Y, Lv J, Liu H, Zhu J, Su J, Wu Q, et al. HHMD: The Human Histone Modification Database. *Nucleic Acids Res*. 2010; 38(Database issue):D149–154. <https://doi.org/10.1093/nar/gkp968> PMID: 19892823
52. Kalvari I, Nawrocki EP, Ontiveros-Palacios N, Argasinska J, Lamkiewicz K, Marz M, et al. Rfam 14: Expanded Coverage of Metagenomic, Viral and microRNA Families. *Nucleic Acids Research*. 2021; 49(D1):D192–D200. <https://doi.org/10.1093/nar/gkaa1047> PMID: 33211869
53. Brameier M, Herwig A, Reinhardt R, Walter L, Gruber J. Human Box C/D snoRNAs with miRNA like Functions: Expanding the Range of Regulatory RNAs. *Nucleic Acids Res*. 2011; 39(2):675–686. <https://doi.org/10.1093/nar/gkq776> PMID: 20846955
54. Ono M, Scott MS, Yamada K, Avolio F, Barton GJ, Lamond AI. Identification of Human miRNA Precursors That Resemble Box C/D snoRNAs. *Nucleic Acids Research*. 2011; 39(9):3879–3891. <https://doi.org/10.1093/nar/gkq1355> PMID: 21247878
55. Lacoste A, Luccioni A, Schmidt V, Dandres T. Quantifying the Carbon Emissions of Machine Learning. *Workshop on Tackling Climate Change with Machine Learning at NeurIPS 2019*. 2019;.
56. Lottick K, Susai S, Friedler SA, Wilson JP. Energy Usage Reports: Environmental Awareness as Part of Algorithmic Accountability. *Workshop on Tackling Climate Change with Machine Learning at NeurIPS 2019*. 2019;.
57. Mortimer SA, Kidwell MA, Doudna JA. Insights into RNA Structure and Function from Genome-Wide Studies. *Nat Rev Genet*. 2014; 15(7):469–479. <https://doi.org/10.1038/nrg3681> PMID: 24821474
58. Justyna M, Antczak M, Szachniuk M. Machine Learning for RNA 2D Structure Prediction Benchmarked on Experimental Data. *Briefings in Bioinformatics*. 2023; 24(3):bbad153. <https://doi.org/10.1093/bib/bbad153> PMID: 37096592
59. Sato K, Hamada M. Recent Trends in RNA Informatics: A Review of Machine Learning and Deep Learning for RNA Secondary Structure Prediction and RNA Drug Discovery. *Briefings in Bioinformatics*. 2023; 24(4):bbad186. <https://doi.org/10.1093/bib/bbad186> PMID: 37232359
60. Fasold M, Langenberger D, Binder H, Stadler PF, Hoffmann S. DARIO: A ncRNA Detection and Analysis Tool for next-Generation Sequencing Experiments. *Nucleic Acids Res*. 2011; 39(Web Server issue):112–117. <https://doi.org/10.1093/nar/gkr357>
61. Videm P, Rose D, Costa F, Backofen R. BlockClust: Efficient Clustering and Classification of Non-Coding RNAs from Short Read RNA-seq Profiles. *Bioinformatics*. 2014; 30(12):i274–i282. <https://doi.org/10.1093/bioinformatics/btu270> PMID: 24931994

62. Pourpanah F, Abdar M, Luo Y, Zhou X, Wang R, Lim CP, et al. A Review of Generalized Zero-Shot Learning Methods. *IEEE Trans Pattern Anal Mach Intell.* 2022; p. 1–20. <https://doi.org/10.1109/TPAMI.2022.3191696>
63. Creux C, Zehraoui F, Hanczar B, Tahi F. A3SOM, Abstained Explainable Semi-Supervised Neural Network Based on Self-Organizing Map. *PLOS ONE.* 2023; 18(5):e0286137. <https://doi.org/10.1371/journal.pone.0286137> PMID: 37228138
64. Hendrickx K, Perini L, Van Der Plas D, Meert W, Davis J. Machine Learning with a Reject Option: A Survey. *Mach Learn.* 2024; 113(5):3073–3110. <https://doi.org/10.1007/s10994-024-06534-x>