# Computational Understanding of non-coding RNA Pairwise Interactions

**Marco Nicolini** [1], **Federico Stacchietti** [1], **Elena Casiraghi** [1,2,3,4] **and Giorgio Valentini,** [1,2,*]

[1] *AnacletoLab, Dipartimento di Informatica, Università degli Studi di Milano, Italy*
[2] *ELLIS - European Lab for Learning and Intelligent Systems, Milan Unit*
[3] *Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, United States*
[4] *Department of Computer Science, Aalto University, Espoo, Finland*

Correspondence*:
Giorgio Valentini
valentini@di.unimi.it

## ABSTRACT

Non-coding RNAs (ncRNAs) govern a vast network of regulatory interactions within the cells, yet their pairwise relationships remain largely uncharted due to the complexity of RNA structure and the limits of current experimental methods. We present *CUPID* (Computational Understanding of Pairwise Interactions in ncRNA Data), a deep learning framework that predicts ncRNA-ncRNA interactions directly from primary sequence information. *CUPID* uses embeddings from a pre-trained RNA language model combined with a feed-forward classifier to identify patterns linked to molecular pairing. This approach avoids reliance on thermodynamic models or manual feature design and, unlike previously proposed models, is able to generalize across different types of ncRNAs, including long non-coding, circular, micro-, and small nuclear RNAs. By learning the hidden rules that govern RNA recognition, *CUPID* provides a scalable tool for exploring ncRNA interaction networks and advancing our understanding of RNA-based regulation.

**Keywords: ncRNA–ncRNA interaction, deep learning, fine-tuning, artificial intelligence, machine learning, non-coding RNA, large language models**

## 1 INTRODUCTION

Understanding RNA-RNA interactions is critical for deciphering the regulatory circuits that orchestrate gene expression, RNA processing, and signal transduction. Non-coding RNAs (ncRNAs), despite lacking protein-coding potential, play pivotal roles in these processes (Ali et al., 2021). However, experimental mapping of ncRNA interactions remains challenging due to the limitations of existing experimental and computational techniques (Lorenzi et al., 2021).

Methods such as Minimum Free Energy (MFE) calculations and accessibility-based models have been usually applied to predict RNA-RNA interactions. Tools like IntaRNA (Mann et al., 2017) estimate the interaction energy as $\Delta G_{\text{total}} = \Delta G_{\text{duplex}} + \Delta G_{\text{accessibility}}$, where the first term quantifies the energy released upon hybridization, and the second accounts for the cost of rendering binding regions accessible. Benchmark studies have demonstrated that accessibility-based algorithms can effectively differentiate native interactions from background noise (Umu and Gardner, 2017), yet these approaches rely on predefined parameters and simplified energy models. In parallel, experimental

28 techniques such as RNA Antisense Purification (RAP-RNA) offer validation but remain limited by
29 their high cost and labor intensity (Engreitz et al., 2014).

30 Advances in machine learning and graph-based modeling for biological data, including recent work
31 on explainability and diffusion-based attention mechanisms, have motivated a surge of learning-
32 driven approaches for predicting interactions across diverse molecular systems (Gliozzo et al., 2025;
33 Cetin and Sefer, 2025; Sefer, 2025).

34 Machine learning methods, such as convolutional neural networks, deep forests and graph neural
35 networks (Alipanahi et al., 2015; Tian et al., 2021; Wei et al., 2022) have been applied to RNA-protein
36 interaction prediction, while graph-based approaches embed heterogeneous networks of ncRNAs and
37 diseases using multigraph contrastive learning (Sun et al., 2025) or apply random-walk based graph
38 representation learning techniques to predict non coding RNA interactions (Torgano et al., 2025).

39 While effective, these methods often rely on predefined feature extraction, graph structures, or
40 supervised training, limiting their adaptability to novel ncRNA sequences.

41 In contrast, LLMs can directly learn from large corpora of proteins or RNA data (Valentini et al.,
42 2023; Zhao et al., 2024; Shen et al., 2024; Nicolini et al., 2025a), capturing intricate interaction motifs
43 beyond predefined energy models or graph-based constraints. Unlike thermodynamic models, which
44 impose simplifying assumptions, LLMs infer interaction likelihoods from latent structural patterns,
45 offering a flexible, data-driven approach. In particular, transformer-based foundation models can
46 generate biologically meaningful representations directly from raw sequences, by exploiting large
47 RNA sequence corpora (Sapoval et al., 2022; Chen et al., 2022; Yu et al., 2024). More in general
48 several deep learning methods have been proposed to predict specific ncRNA interactions, using
49 rna2vec pretraining and deep feature mining (Yu et al., 2022) or conditional random fields and
50 graph convolutional networks (Wang et al., 2022), heterogeneous graph neural networks (Li et al.,
51 2025) and convolutional neural networks combined with a Transformer Encoder (Yang et al., 2025)
52 for the prediction of miRNA-lncRNA interactions.

53 We also recently proposed a deep neural network trained on embedded representations of a subset
54 of ncRNAs obtained from the RNA-FM language model (Shen et al., 2024), achieving state-of-the-art
55 results for predicting miRNA interactions with other ncRNA molecules (Nicolini et al., 2025b).
56 However, our proposed model, like other models recently proposed in the literature (Li et al.,
57 2025; Yang et al., 2025), is only able to predict specific ncRNA interactions (e.g., interactions
58 with miRNAs). Furthermore, due to limitations on the maximum allowed sequence length of the
59 underlying RNA-FM transformer, it can only process sequences shorter than approximately 1000
60 nucleotides, thus limiting the model's application to relatively long ncRNAs (e.g., lncRNAs).

61 To overcome these limitations, we propose a novel Transformed-based deep learning model, that,
62 differently from previous models proposed in literature, is able to predict a large range of ncRNA
63 interactions, including long non-coding RNA (lncRNA), circular RNA (circRNA), microRNA
64 (miRNA), small nuclear RNA (snRNA), small nucleolar RNA (snoRNA), Small Cajal body-specific
65 RNAs (scaRNAs), small cytoplasmic RNAs (scRNA) and other types of ncRNAs. Moreover, by
66 adopting GenerRNA (Zhao et al., 2024) to encode RNA sequences, our model can process full-length
67 ncRNA sequences (up to 4096 nucleotides) without truncation, thus significantly enlarging the set
68 of ncRNAs that can be processed by the model.

69 We hypothesize that LLM-based contextual embeddings provide a rich representation for ncRNA
70 interaction prediction, circumventing the limitations of manual feature engineering or predefined

structural graphs. We reasoned that GenerRNA (Zhao et al., 2024), pretrained on a large corpus of ncRNA sequences using a masked language modeling objective, can capture long-range interactions of ncRNA molecules, thus facilitating downstream tasks such as ncRNA interaction prediction.

Our *CUPID* model (Computational Understanding of Pairwise Interactions in ncRNA Data), predicts ncRNA interactions using only sequence information. *CUPID* extracts embeddings from a pre-trained ncRNA language model and feeds a dense feed-forward neural network (FFNN) to automatically learn intricate sequence interaction features. This design circumvents the need for explicit thermodynamic parameterization and manually engineered features, offering a scalable and efficient alternative for uncovering novel regulatory interactions (Fabbri et al., 2019).

## 2 METHODS

### 2.1 Dataset

Our dataset comprises a subset of multispecies ncRNA interaction pairs from RNA-KG Cavalleri et al. (2024)[1].

The RNA-KG integrates physical and functional interactions between different types of ncRNAs, and their relationships with other biomolecules (genes and proteins) and chemicals, as well as with biomedical concepts coded in the Gene Ontology (Aleksander et al., 2023), the Human Phenotype Ontology (Gargano et al., 2023), Mondo (Vasilevsky et al., 2025) and other bio-medical ontologies related to the "RNA world".

In particular, we extracted RNA–RNA edges from RNA-KG by selecting only relations annotated as `interacts-with`. In RNA-KG, `interacts-with` denotes experimentally supported *physical* RNA–RNA interactions, and we therefore excluded other relation types encoding functional associations (e.g., regulatory links, co-expression, or disease associations). The `interacts-with` edges integrated in RNA-KG originate from multiple underlying curated interaction databases. Fig.1 presents an overview of the main RNA entities and their relationships available in the the RNA-KG. Readers may refer to the RNA-KG reference (Cavalleri et al., 2024) for the complete list of contributing sources and evidence provenance.
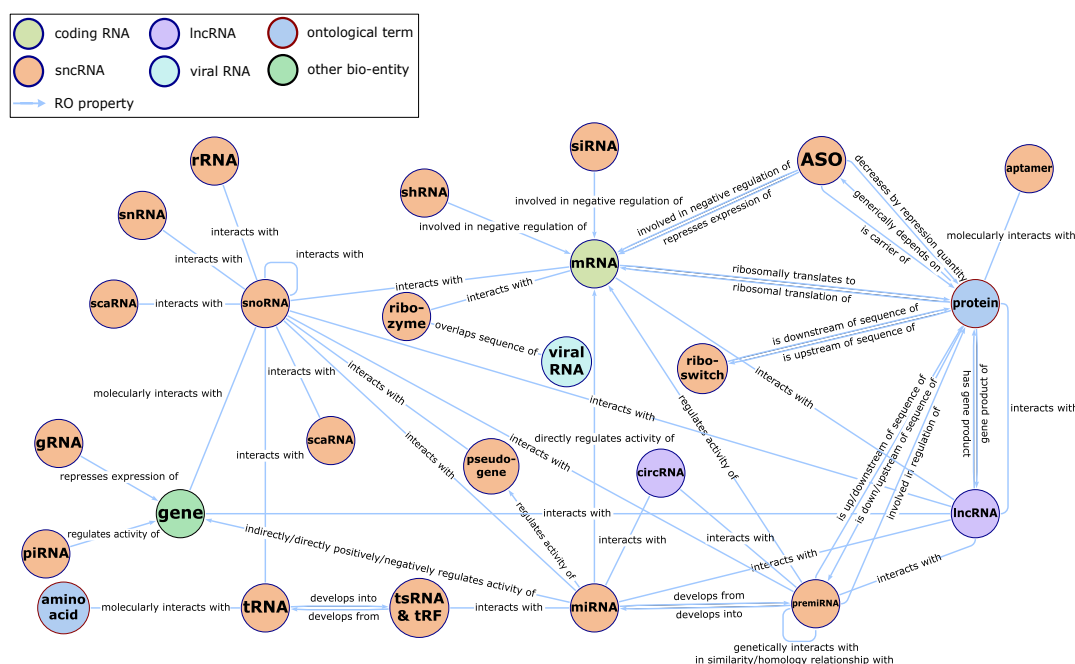
We filtered the dataset to retain only sequences that fit within the GenerRNA (Zhao et al., 2024)'s token limit (approximately 4096 nucleotides), since Byte Pair Encoding (BPE) compresses raw nucleotide sequences, allowing longer sequences to fit within the model's constraints. After applying this length filter, the dataset contains:

- 101088 interaction pairs (down from an initial 130310 pairs).
- 11212 unique sequences (selected from 19624 potential sequences) belonging to 9 different RNA molecule types: long non-coding RNA (lncRNA), circular RNA (circRNA), microRNA (miRNA), small nuclear RNA (snRNA), small nucleolar RNA (snoRNA), Small Cajal body-specific RNAs (scaRNAs), small cytoplasmic RNAs (scRNA), not (better) classified non coding RNA molecules (ncRNA) and pseudo RNA[2].

---

[1] Retrieval of interacting pairs and corresponding sequences was performed using the scripts available from the RNA-KG web site: `https://github.com/AnacletoLAB/RNA-KG`.

[2] In RNAinter, the term "pseudo" specifically denotes RNA sequences transcribed from pseudogenes. In this context, these are transcripts derived from genes that have lost their protein-coding capability due to accumulated mutations, yet they are still produced as RNA. Similar to other ncRNAs, such pseudogene RNAs can sometimes participate in regulatory networks by, for example, acting as miRNA decoys or sponges, despite not encoding functional proteins.

---

**Figure 1.** Simplified representation of the RNA-KG meta-graph, focused on ncRNAs and their interactions.

In the following, we denote the set of length-filtered molecules as

$$\mathcal{S} = \{s_i\}, \quad i = 1, \dots, |\mathcal{S}|,$$

where the type of each molecule $s \in \mathcal{S}$ is given by $\phi(s)$, i.e. $\phi : \mathcal{S} \to \mathcal{T}$ represents a mapping of a ncRNA sequence $s \in \mathcal{S}$ to its ncRNA type $\mathcal{T}$, e.g. miRNA, lncRNA or any other ncRNA type.

The identity of an interaction pair is solely determined by its constituent molecules, regardless of order; that is,

$$(s_i, s_j) = (s_j, s_i).$$

The type of an interaction $(s_i, s_j)$ with $s_i \neq s_j$ and $s_i, s_j \in \mathcal{S}$ is determined by the types of the ncRNA $s_i$ and $s_j$ theirselves, regardless of their order:

$$(\phi(s_i), \phi(s_j)) = (\phi(s_j), \phi(s_i))$$

For instance, possible types of ncRNA interactions are miRNA-lncRNA or miRNA-circRNA. Assuming that interacting ncRNA pairs of different types exhibit distinct specificities that the model should learn, we reasoned that types with negligible sample sizes might introduce noise rather than valuable information. Therefore, the set of interaction pairs used in this work is obtained by further filtering the dataset of interacting pairs to remove interacting pair types represented by fewer than 100 samples, resulting in 10644 unique sequences composing 99841 interacting pairs. Fig 2 shows the distribution of the different types of ncRNA interactions.

|  | lncRNA | miRNA | ncRNA | pseudo | scRNA | scaRNA | snRNA | snoRNA |
|---|---|---|---|---|---|---|---|---|
| circRNA | - | 295 | - | - | - | - | - | - |
| lncRNA | 1335 | 54773 | 126 | 370 | - | - | - | 1535 |
| miRNA | - | 1864 | 5814 | 16853 | 182 | 115 | 111 | 1778 |
| ncRNA | - | - | - | - | - | - | - | 196 |
| pseudo | - | - | - | - | - | - | - | 93 |
| snRNA | - | - | - | - | - | - | - | 94 |
| snoRNA | - | - | - | - | - | - | - | 4322 |

**Figure 2.** Distribution of ncRNA interactions pairs in the filtered interaction set. Rows: first (left) molecule type; Columns: right molecule type.

## 2.2 Data Augmentation

To address the issues due to the limited cardinality of the available training data, especially for specific types of ncRNA interactions (e.g., snRNA-miRNA or miRNA-circRNA), we employ a data augmentation strategy that effectively increases the dataset size by a factor of 4. For each original training instance represented as a pair of interacting ncRNA $(s_i, s_j)$ we generate three additional augmented instances:

1. Molecule Order Reversal: Swap the order of the molecules: $(s_j, s_i)$.
2. Sequence Flipping: Reverse the nucleotide order in both molecules (denoted by the superscript $F$): $(s_i^F, s_j^F)$.
3. Combined Augmentation: Reverse both the molecule order and the nucleotide sequences: $(s_j^F, s_i^F)$.

Thus, if the original dataset contains $N$ instances, the augmented dataset becomes: $N_{\text{aug}} = 4N$ (Suppl. Fig. S1). This augmentation introduces invariance to both the order and orientation of sequences, thereby enabling the model to better capture the underlying biological patterns and improving its robustness against input variability.

In order to avoid leakage between training and test sets, data augmentation is performed after splitting the dataset.

## 2.3 Negative examples generation

In our dataset, only positive non-coding RNA-RNA interactions are explicitly provided, and they occur with varying frequencies.

---

**Algorithm 1** Negative Sampling Algorithm

---

**Require:** Set of unique ncRNA sequences $\mathcal{S}$, positive augmented interaction set $\mathcal{P}$, and negative
    sampling parameter $n$
**Ensure:** Negative sample set $\mathcal{N}$
  1: Initialize $\mathcal{N} \leftarrow \emptyset$
  2: **for** each pair $(s, s') \in \mathcal{P}$ **do**
  3:     **for** $i = 1$ to $n$ **do**
  4:         Sample $s_{\text{neg}} \in \mathcal{S}$ such that $\phi(s_{\text{neg}}) = \phi(s')$
  5:         **if** $(s, s_{\text{neg}}) \notin \mathcal{P} \wedge (s, s_{\text{neg}}) \notin \mathcal{N}$  **then**
  6:             $\mathcal{N} \leftarrow \mathcal{N} \cup \{(s, s_{\text{neg}})\}$
  7:         **end if**
  8:     **end for**
  9: **end for**
10: **return** $\mathcal{N}$

---

140    To effectively train *CUPID*, we generated negative examples for each interaction pair type by
141  matching the frequency distribution of the positive interactions. Specifically, negative examples were
142  generated under the assumption that any pair of ncRNA sequences drawn from the set of unique
143  sequences that is not observed as a positive interaction constitutes a possible negative instance.

144    Let $\mathcal{S} = \{s_1, s_2, \ldots, s_N\}$ be the set of unique ncRNA sequences present in the dataset. Denote by

$$\mathcal{P} = \{(s_i, s_j) \mid s_i, s_j \in \mathcal{S} \text{ interact}\}$$

145  the set of all positive ncRNA-ncRNA interactions. Then, the set of all possible ncRNA pairs is
146  given by $\mathcal{S} \times \mathcal{S}$ (excluding self-interactions).

147    The set of *potential negatives* is defined as:

$$\mathcal{N}_{\text{potential}} = \{(s_i, s_j) \in \mathcal{S} \times \mathcal{S} \mid s_i \neq s_j\} \setminus \mathcal{P}.$$

148  *Negative Sampling Procedure.* To generate the negative samples for each interacting pair type, we
149  corrupt its tuples. In other words, given a positive pair $(s_i, s_j)$ with type $(\phi(s_i), \phi(s_j))$, we keep the
150  first molecule $s_i$ fixed and sample $s' \in \mathcal{S}$ such that:

$$s' \neq s_i, \quad \phi(s') = \phi(s_j), \quad (s_i, s') \notin \mathcal{P}$$

151  In this way we avoid generating negatives between ncRNA types that never interact (e.g. scaRNA
152  and lncRNA).

153    Because we generate negatives for each positive pair $(s_i, s_j)$ by corrupting the right molecule
154  while keeping the same type pair $(\phi(s_i), \phi(s_j))$, the negative set preserves the interaction *type-pair*
155  distribution of the positives in expectation (and approximately in practice, up to rejection of
156  candidates already present as positives or previously sampled negatives).

157    For each positive edge, we selected $n$ negative edges, in order to control the imbalance between
158  positive and negative edges in the testing phase (we set $n = 20$ in our experiments).

159  *Negative Sampling Algorithm.* The negative sampling algorithm is detailed in Algorithm 1. In our
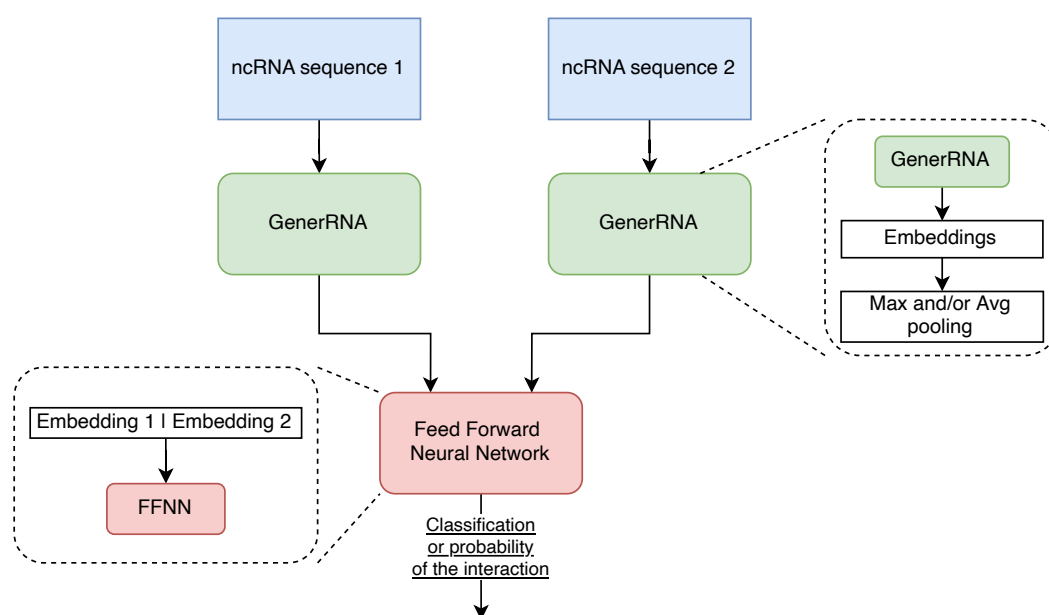160  implementation, we set $n = 20$. Note that, since the condition at line 5 of the algorithm cannot

---

161 be always guaranteed, it is likely that the number of negatives $n \leq 20$. In our experiments we set
162 $n = 20$.

## 2.4  Model Architecture

### 2.4.1  The overall *CUPID* Architecture

165 Our model follows a two-stage pipeline, as illustrated in Figure 3. It first extracts ncRNA sequence
166 embeddings using a pre-trained ncRNA Language Model (GenerRNA Zhao et al. (2024)) and then
167 processes these embeddings through a Feed-Forward Neural Network (FFNN) to predict interaction
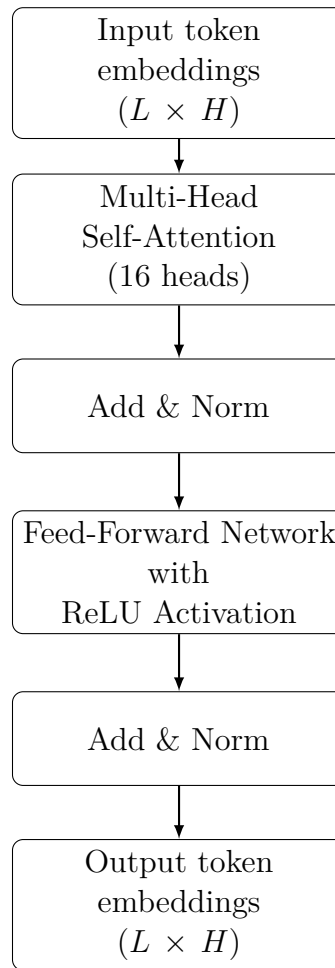168 probabilities.



**Figure 3.** High-level *CUPID* architecture schema.

169 The GenerRNA architecture mimics the GPT-2-medium model (Radford et al., 2019), and is
170 composed of 24 stacked transformer-decoder layers, each incorporating a self-attention mechanism
171 that models pairwise interactions among all positions in its input sequence. GenerRNA uses a
172 context window of 1024 tokens, corresponding to input sequences with a length of approximately
173 4096 nucleotides coded through byte pair encoding Sennrich et al. (2016). Note that this maximum
174 length permits the encoding of large RNA molecules. This decoder-only Transformer architecture
175 operates in an autoregressive manner, predicting the subsequent token given the previous ones. Both
176 the input and output of the model are represented as tokens, which are encoded and decoded by
177 a trained tokenizer using byte pair encoding. A special token (EOS) is used to delimit sequences,
178 indicating the start and end of each sequence.

179 Each transformer block is fed with a input of size $L \times H$, thus allowing to process RNA sequences
180 having up to $L$ tokens, each one represented through a $H$-dimensional real vector, with $L = H =$
181 1024, and outputs a latent representation with the same dimensionality for each input token. For
182 each input sequence, the block employs a multi-head self-attention mechanism with 16 attention
183 heads. This is followed by an "Add & Norm" sub-block, which applies residual addition and layer
184 normalization. Subsequently, a feed-forward sub-layer expands the hidden states from 1024 to 4096

185  dimensions, applies a non-linear activation (ReLU), and then projects them back to 1024 dimensions.
186  Another "Add & Norm" sub-block is applied after the feed-forward network, and finally, the block
187  produces an output matrix $\mathbf{X} \in \mathbb{R}^{L \times H}$. A schematic diagram of this block is reported in the Fig 4.


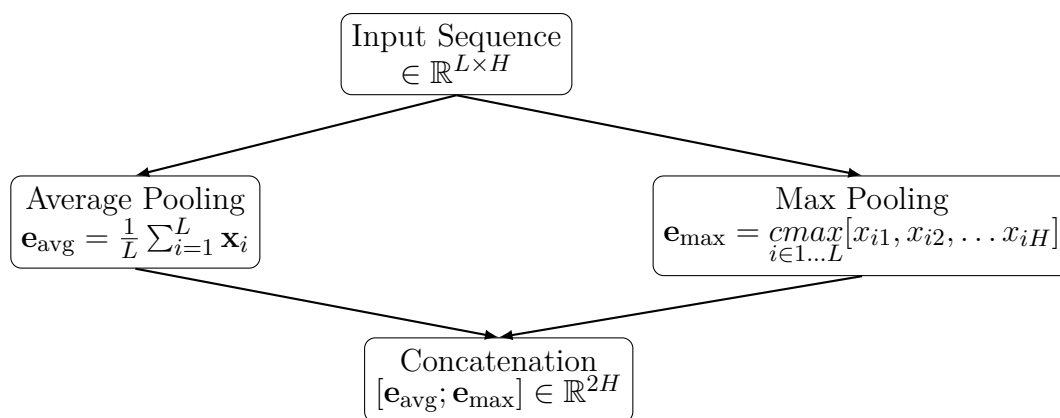
**Figure 4.** High-level architecture of a GenerRNA block.

### 2.4.2  Pooling techniques

189  The $i^{th}$ row of matrix $\mathbf{X}$ is a latent representation $\mathbf{x}_i \in \mathbb{R}^{1024}$ of the $i^{th}$ token. To obtain a
190  fixed-length embedding for the entire sequence, we tested two types of pooling over the sequence
191  (i.e., across the $H$ tokens), as well as their concatenation:

192  • Average (*Avg*) Pooling: obtained as the mean of the embeddings of all the tokens: $\mathbf{e}_{\text{avg}} =$
193    $\frac{1}{L} \sum_{i=1}^{L} \mathbf{x}_i$.
194  • Maximum (*Max*) Pooling: Compute the element-wise maximum over all token embeddings:
195    $\mathbf{e}_{\max} = \underset{i \in 1...L}{cmax}[x_{i1}, x_{i2}, \ldots x_{iH}]$, where *cmax* is the columnwise *max* operator, and $x_{ij}$ are the
196    elements of the $\mathbf{X}$ embedding matrix.
197  • Concatenation of [*Avg, Max*]: Combine both pooled representations into a single embedding
198    vector: $\mathbf{e} = [\mathbf{e}_{\text{avg}}; \mathbf{e}_{\text{max}}] \in \mathbb{R}^{2048}$.

199  These pooling strategies are schematically depicted in Fig. 5.

**Figure 5.** Pooling Embedding Strategies.

200   The embedded representation of a candidate interacting ncRNA pair is composed by the
201   concatenation of the embeddings of the two interacting molecules.

### 2.4.3   The classification unit

203   To predict the interaction we used the pooled embeddings of the RNA sequences as input to a
204   dense Feed Forward Neural Network (FFNN), having the following architecture:

205   • **Input Layer Dimension:** 1024 for Avg and Max-pooling embedding strategies, 2048 when
206     the embedding of the input molecule is obtained by concatenating the embeddings obtained by
207     AVG and Max pooling,

208   • **Hidden Layers:** 4 hidden layers with 1024 neurons each and ReLU activation function,

209   • **Output Layer:** 1 neuron with sigmoid activation function.

210   To train the network we applied the following hyper-parameters:
211   *Learning Rate:* $\eta = 5 \times 10^{-4}$ with a linear warm-up phase of 4 epochs, followed by cosine decay.
212   *Epochs:* 50 epochs with early stopping (patience of 10 epochs). The model with the best validation
213   loss is selected (e.g., if the lowest validation loss is observed at epoch 35, then early stopping is
214   triggered at epoch 45, and the model from epoch 35 is used).
215   *Batch Size:* 512. *Dropout Rate:* 0.2. *Optimizer:* Adam. *Loss Function:* Binary Cross-Entropy.

216   Training and validation loss curves were monitored over epochs to assess model convergence and
217   to avoid potential overfitting by early stopping.

### 2.4.4   Mini-batch balancing

219   Due to the imbalance in our dataset we adopted a training strategy designed to prevent the model
220   from learning predominantly from the negatives. To address this, we constructed mini-batches
221   that contain a controlled mix of positive and negative examples. Recall that our training set is
222   composed of the set of positive interaction pairs, $\mathcal{P}, |\mathcal{P}| = N_+$, and the set of negative interaction
223   pairs $\mathcal{N}$, with $|\mathcal{N}| = N_- = nN_+$, as detailed in Section 2.3. Each mini-batch $B$ of size $m$ is formed
224   by randomly selecting $m_p$ positive examples (using a uniform distribution with replacement) and
225   $m_n$ negative examples (using a uniform distribution without replacement). The ratio of negatives

within each mini-batch is defined by

$$r = \frac{m_n}{m_n + m_p}, \quad \text{with} \quad m_n + m_p = m.$$

Here, $r$ can vary between 0 and 1. A value of $r = 0.7$ implies that 70% of the mini-batch consists of negatives. The choice to sample positives with replacement is driven by their limited number, ensuring sufficient representation even in large batches, whereas sampling negatives without replacement allows for a broader coverage of these more abundant examples.

## 2.5 Experimental Evaluation

### 2.5.1 Data preparation and splitting.

In all our experiments the negative examples were sampled according to the relative frequency of the interacting pair types according to the procedure described in Section 2.3 using a negative:positive ratio equal to 20:1.

The dataset was partitioned into stratified training and test sets (train:test ratio = 90:10). The training set was further split into a stratified set for training (80% of interaction pairs) and the remaining 20% for validation. The validation set was used for early stopping and for tuning the classification threshold via maximization of the Matthews correlation coefficient (MCC Matthews (1975)) on the validation data.

### 2.5.2 Evaluation metrics.

To comprehensively assess model's performance, we computed a range of evaluation metrics, encompassing both threshold-dependent and threshold-independent measures. Specifically, we first evaluated standard binary classification metrics, including accuracy, balanced accuracy (to account for class imbalance), precision, recall, F1 score, AUROC (Area Under the Receiver Operating Characteristic Curve), and AUPRC (Area Under the Precision-Recall Curve). In addition to these overall metrics, we conducted a stratified analysis based on interacting pair types, computing the aforementioned measures separately for each pair type.

Let $y_i \in \{0, 1\}$ be the ground-truth label and $\hat{p}_i \in [0, 1]$ the predicted probability for sample $i$. Given a decision threshold $t$, we define $\hat{y}_i = 1 \iff [\hat{p}_i \geq t]$ and the confusion matrix counts:

$$\text{TP} = \sum_i \mathbf{I}[y_i = 1 \wedge \hat{y}_i = 1],$$

$$\text{TN} = \sum_i \mathbf{I}[y_i = 0 \wedge \hat{y}_i = 0],$$

$$\text{FP} = \sum_i \mathbf{I}[y_i = 0 \wedge \hat{y}_i = 1],$$

$$\text{FN} = \sum_i \mathbf{I}[y_i = 1 \wedge \hat{y}_i = 0].$$

251 Threshold-dependent metrics are then computed as:

$$\text{Accuracy (rate of correctly predicted instances)} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}},$$

$$\text{Recall (proportion of TP w.r.t. all positive samples)} = \frac{\text{TP}}{\text{TP} + \text{FN}},$$

$$\text{Specificity (proportion of TN w.r.t. all negative samples)} = \frac{\text{TN}}{\text{TN} + \text{FP}},$$

$$\text{Precision (proportion of TP among predicted positives)} = \frac{\text{TP}}{\text{TP} + \text{FP}},$$

$$\text{F1 (harmonic mean of precision and recall)} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}},$$

$$\text{BalancedAcc (accuracy balanced by class proportion)} = \frac{\text{Recall} + \text{Specificity}}{2}.$$

$$(1)$$

In our work the threshold $t$ is chosen on the validation set by maximizing the MCC coefficient, which provides a balanced single-score summary that incorporates TP, TN, FP, and FN, and is therefore less sensitive than accuracy to class imbalance:

$$\text{MCC} = \frac{\text{TP} \cdot \text{TN} - \text{FP} \cdot \text{FN}}{\sqrt{(\text{TP} + \text{FP})(\text{TP} + \text{FN})(\text{TN} + \text{FP})(\text{TN} + \text{FN})}}.$$

252 Threshold-independent metrics summarize performance across all thresholds. The ROC curve
253 plots $\text{TPR}(t) = \text{Recall}(t)$ versus $\text{FPR}(t) = \text{FP}(t)/(\text{FP}(t) + \text{TN}(t))$, and AUROC is the area under
254 this curve. The precision–recall curve plots $\text{Precision}(t)$ versus $\text{Recall}(t)$, and AUPRC is its area;
255 under strong class imbalance, AUPRC is often more informative than AUROC, with a random
256 baseline equal to the positive prevalence $\pi = \frac{N_+}{N_+ + N_-}$.

### 2.5.3 Training hyper-parameters and baselines for comparison.

258 The hyper-parameters and configurations used for training the FFNN are reported in Section 2.4.3.
259 Moreover, training and validation loss curves were monitored over epochs to assess model convergence
260 and to avoid potential overfitting by early stopping.

261 Hyperparameter selection was performed in preliminary experiments on a reduced subset of the
262 training/validation interaction pairs using a grid-search strategy. We varied the number of hidden
263 layers in $\{2, 4, 6\}$, the dropout rate in $\{0.1, 0.2\}$, and the batch size in $\{16, 512, 1024\}$. For each
264 configuration, models were trained using the same optimization settings described in Section 2.4.3,
265 and the final model was selected as the configuration that maximized validation AUPRC. No
266 hyperparameters were tuned on the test set.

267 Besides the random classifier, whose expected performance are AUROC = 0.5 and AUPRC
268 = 0.047, we employed the IntaRNA method (Mann et al., 2017) as a baseline for comparison.
269 IntaRNA estimates interaction energy. While interaction energy can be thresholded to obtain binary
270 predictions, which enable the computation of accuracy, balanced accuracy, precision, recall, and F1
271 scores, we opted to limit the comparison to AUROC and AUPRC. These metrics provide a more

272 robust and threshold-independent evaluation of predictive performance, ensuring a fair comparison
273 across models.

# 3 RESULTS

274 We assessed the contribution of the data-augmentation strategy and the pooling operation used
275 to obtain molecule-level embeddings. Table 1 summarizes AUROC and AUPRC results across all
276 configurations, including a baseline random classifier, IntaRNA and *CUPID* models. For *CUPID* we
277 compared results obtained with (Data-aug) and without (No-Data-aug) data augmentation,
278 considering different pooling techniques, i.e. concatenation (concat), maximum (Max) and average
279 (Avg) pooling.

**Table 1.** Comparison of AUROC and AUPRC across different experimental settings. Random baseline refers to the expected performance of the random classifiers. *CUPID* models are sorted in increasing order of both AUROC and AUPRC

| Methods | AUROC | AUPRC |
|---|---|---|
| Random baseline | 0.5 | 0.047 |
| IntaRNA | 0.544 | 0.055 |
| *CUPID* : | | |
| No-Data-aug | 0.658 | 0.078 |
| Data-aug-Max | 0.810 | 0.147 |
| Data-aug-concat | 0.862 | 0.222 |
| Data-aug-Avg | **0.919** | **0.364** |

## 3.1 Random baselines

281 With a random classifier we can expect an AUROC = 0.5, while the estimated baseline AUPRC is:

$$\text{Baseline AUPRC} = \frac{N_+}{N_+ + N_-}$$

282 where $N_+$ is the number of positive samples, and $N_-$ is the number of negative samples. Given the
283 1:20 ratio of positive to negative samples, the AUPRC baseline in the performed experiments is:
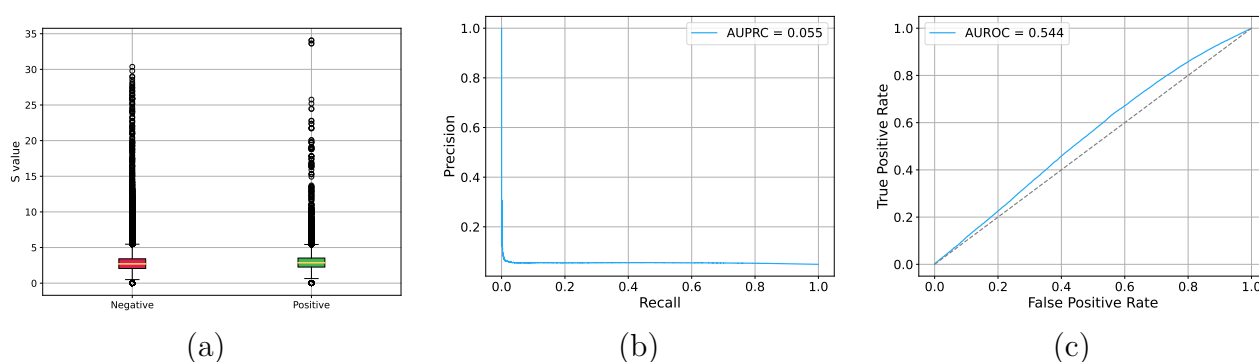
$$\text{Baseline AUPRC} = \frac{1}{1 + 20} \approx 0.0476.$$

284 Our top-performing model achieves an AUPRC of 0.364, corresponding to a *7.65-fold improvement*
285 (0.364/0.0476). This margin quantifies the difficulty of the task: the extreme class imbalance
286 renders precision–recall a stringent metric, and the observed gains indicate that the model extracts
287 interaction-relevant information that is well above chance expectations.

## 3.2 IntaRNA results

Figure 6 reports IntaRNA performance on the augmented test set. In this setting, IntaRNA shows limited predictive power. Its scoring function relies on thermodynamic and accessibility components (e.g., hybridization energy and site accessibility), and in our experiments we used the default parameterization. Given the heterogeneity of ncRNA classes and sequence lengths in our benchmark, improved performance would likely require careful, class-specific calibration of both energy- and accessibility-related settings. Moreover, while IntaRNA is a general thermodynamics- and accessibility-based framework and is not inherently tied to a specific organism, it was originally introduced and most extensively evaluated in bacterial sRNA–mRNA interaction settings; consequently, when applied to heterogeneous ncRNA–ncRNA interactions (including long lncRNAs and diverse eukaryotic classes), its default parameterization may be suboptimal without additional tuning.



**Figure 6.** Results for IntaRNA results with augmented test set. (a) Distribution of predicted probabilities for negative and positive interactions; (b) AUPRC; (c) AUROC

## 3.3 *CUPID* results

Table 1 compares all *CUPID* configurations. We first evaluated a *CUPID* model without augmentation, using concatenation of average and max pooling. Fig. 7 shows the results obtained without data augmentation and with concatenated average-max pooling. The overall AUPRC results on the test set are relatively low (Fig. 7c), even if a certain learning is witnessed by the AUROC largely above 0.5 (Fig. 7f), and by the distribution of the predicted interaction probabilities for negative and positive examples (Fig. 7c), with probabilities for positives relatively higher with respect to negatives. Nevertheless, the relatively flat trend of the training loss reveals a certain difficulty of the model to learn the data. This is reflected also in the confusion matrix where most of negative examples (70%) are misclassified ((Fig. 7e) and in the degradation of the AUPRC performance between validation (Fig. 7a) and test (Fig. 7d) data. By looking at specific ncRNA interactions, for certain interaction types (e.g.snRNA-snoRNA) we obtained good results across the different metrics, but for several ncRNA interactions (e.g. miRNA-lncRNA, miRNA-miRNA, lncRNA-snoRNA) we achieved poor results, with AUPRC below 0.1 (Fig. 7g). Summarizing Fig. 7, shows that with this setting *CUPID* can provide a certain discrimination between positive and negative interactions (Fig. 7c), but its precision–recall and ROC curves indicate a limited separation between positive and negative examples (Fig. 7d,f).

Introducing data augmentation consistently improves performance (Table 1). Fig. 8 shows the results obtained with data augmentation and average pooling. The AUPRC is more than 4 times larger than without data augmentation (Fig. 8d and Table 1). Enlarging the size of training data by data augmentation allows the model to better learn the training data, as witnessed by the training loss that continues to decrease across epochs (Fig. 8b). This results in a clear separation between the scores predicted for positive and negative examples – note that the probabilities predicted for negatives are compressed toward zero while for most positives are largely above 0.7 (even if with several outliers for both positive and negative examples, Fig. 8b). The confusion matrix also confirms that the model with augmented data can better predict negative examples (Fig. 8e); AUPRC (Fig. 8d) significantly improves, and AUROC is larger than 0.9 (Fig. 8f). Analyzing results for each specific ncRNA interaction, we can observe a significant improvement across all the considered metrics, with AUROC in most cases larger than 0.9, except for circRNA-miRNA, miRNA-scRNA, miRNA-snRNA and miRNA-scaRNA (even if for these two last ncRNA interactions values are close to 0.9 (Fig. 8g).

These results confirm that data augmentation is crucial to improve results for two main reasons: at first the model has training data enough to better generalize; second, improves generalization leveraging molecule order and orientation, two symmetries that are not guaranteed to be learned from limited training data. Augmentation effectively enforces these invariances, reducing overfitting to sequence presentation and mitigating the scarcity of positive examples.
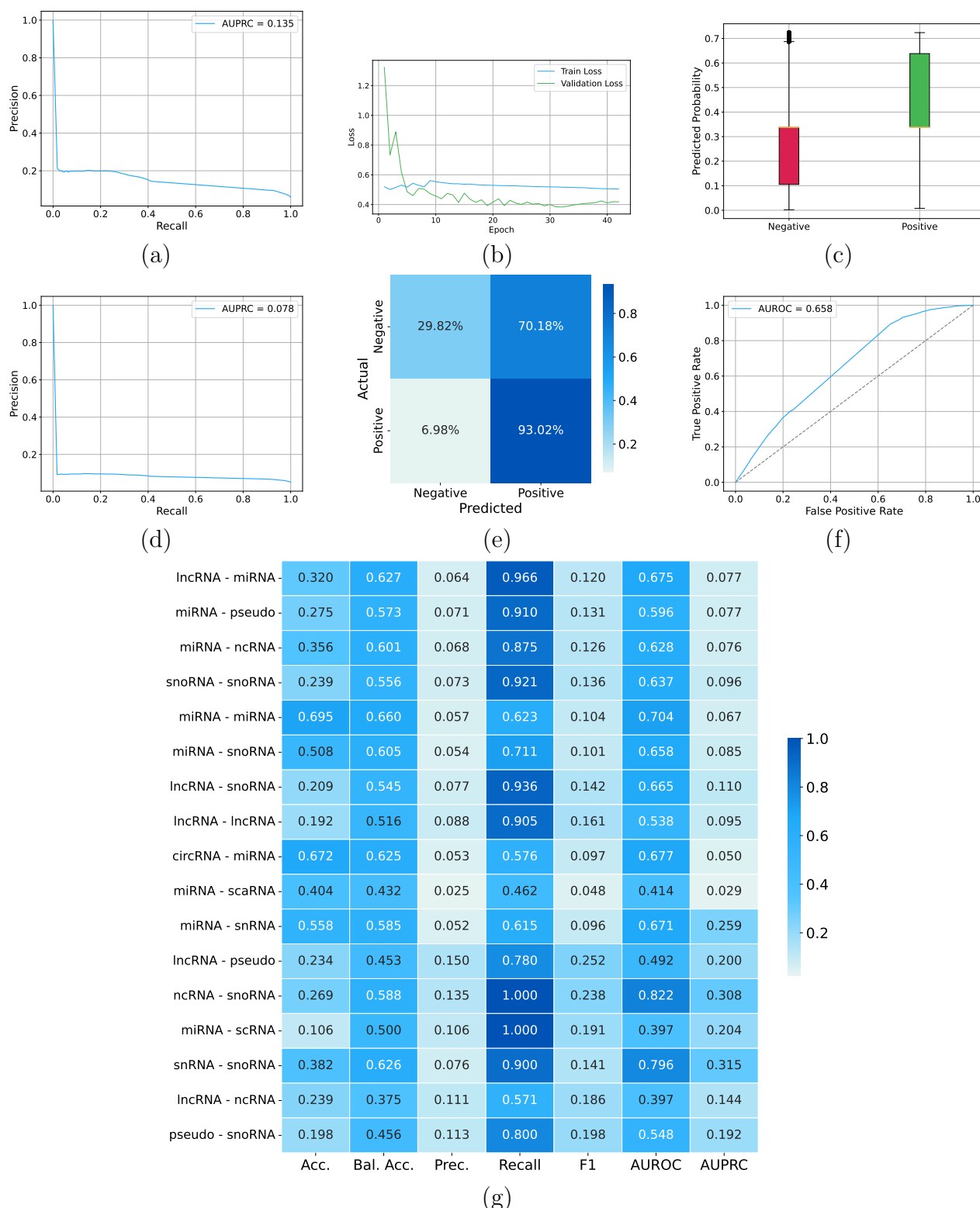
Pooling strategy has a direct impact on the stability of the molecule-level embedding. Average pooling—yielding a smoothed representation over the full sequence—achieves the highest AUROC and AUPRC (Fig. 8) compared to max pooling (Suppl. Fig. S2) and concatenation pooling (Suppl. Fig. S3). This indicates that interaction-relevant information is not confined to a small set of token embeddings but arises from distributed features along the sequence. Max pooling, in contrast, appears sensitive to local outliers and overly compresses positional variability, while concatenation does not provide additional benefits once augmentation is introduced. The results suggest that, for ncRNA interaction prediction, the aggregate signal across nucleotides is more informative than isolated high-activation sites.

## 4 DISCUSSION

The results shown in this work demonstrate that RNA sequence-only inference can recover interaction signals across diverse ncRNA classes. The best-performing configuration reaches AUROC values above 0.9 on the test set, despite operating without structural, evolutionary, or thermodynamic information. This suggests that pretrained RNA language models encode latent features associated with intermolecular recognition. These features may reflect statistical regularities of pairing propensities and local compositional biases captured during pretraining, even in the absence of explicit structural supervision.
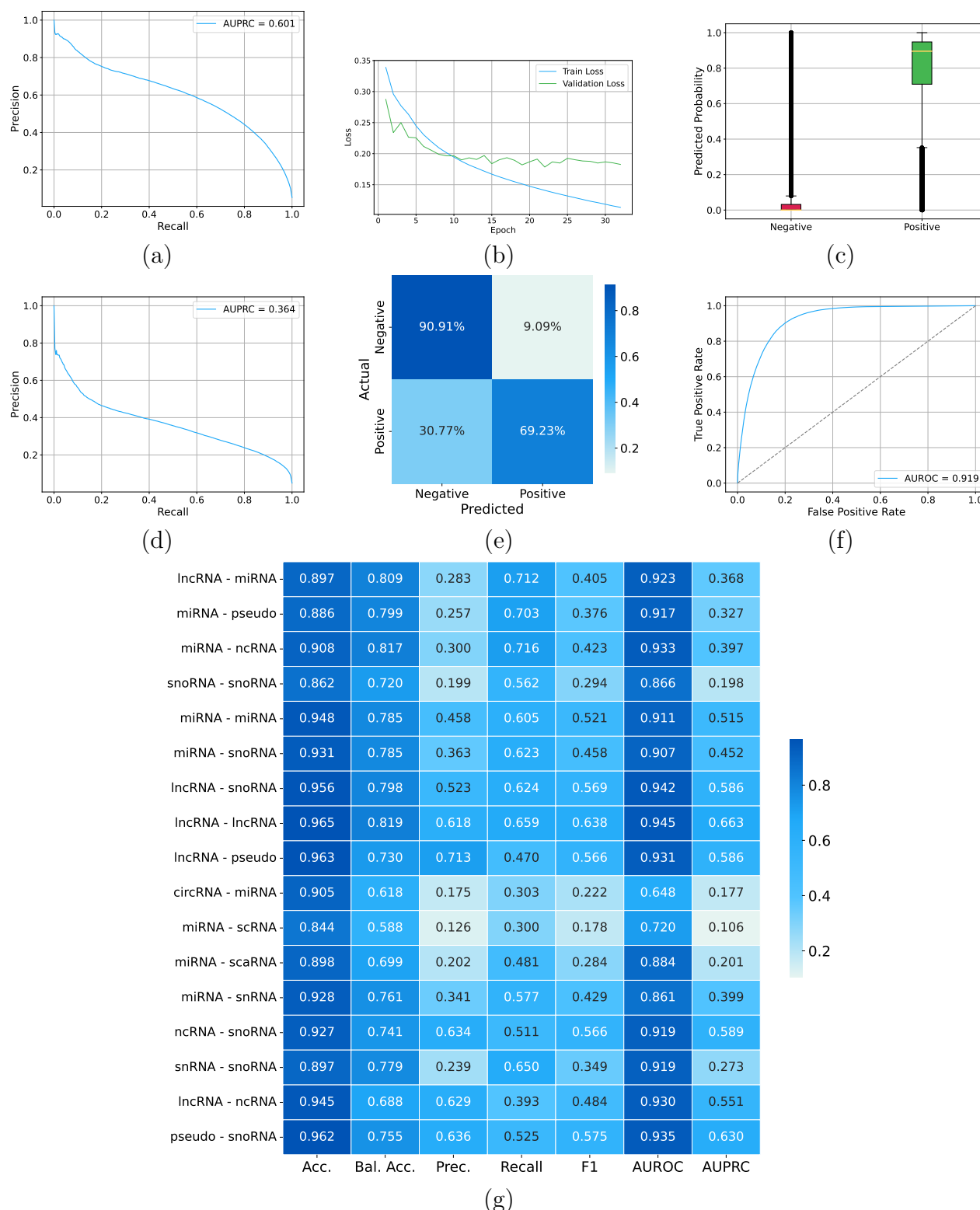
From a methodological standpoint, two contributions appear essential. First, the augmentation scheme addresses symmetries inherent to the problem. Because interacting RNAs can be presented in either order, and because sequence orientation can vary, enforcing invariance to these transformations is critical for robust generalization. Data augmentation also increases the number of examples available for training, thus improving the generalization performance of the model. Second, average pooling provides stable embeddings for ncRNA sequences. For molecules such as lncRNAs——whose

**Figure 7.** *CUPID* results with concatenated pooling, and without using augmented data. (a) Overall precision recall curve on the validation set including all the type of ncRNA interactions; (b) Training and validation loss across epochs; (c) Distribution of the *CUPID* predicted probabilities on negative and positive examples on the test set; (d) Overall precision recall curve on the test set including all the type of ncRNA interactions; (e) Confusion matrix on the test set; (f) ROC curve on the test set including all the type of ncRNA interactions; (g) *CUPID* results on the test set across different types on ncRNA interactions (rows) for different types of metrics (columns).

**Figure 8.** *CUPID* results with average pooling and using augmented data. (a) Overall precision recall curve on the validation set including all the type of ncRNA interactions; (b) Training and validation loss across epochs; (c) Distribution of the *CUPID* predicted probabilities on negative and positive examples on the test set; (d) Overall precision recall curve on the test set including all the type of ncRNA interactions; (e) Confusion matrix on the test set; (f) ROC curve on the test set including all the type of ncRNA interactions; (g) *CUPID* results on the test set across different types on ncRNA interactions (rows) for different types of metrics (columns).

358 functional elements are dispersed and whose lengths vary over orders of magnitude——summarizing
359 the full sequence avoids overemphasis on isolated positions and instead captures global contextual
360 tendencies. Moreover, to our knowledge, *CUPID* is the first model able to predict a large set of
361 ncRNA interactions, and in principle can be applied to predict any ncRNA interaction.

362 The limitations observed for IntaRNA highlight the difference between energy-based and
363 representation-based approaches. Thermodynamic models rely on explicit structural motifs and
364 accessibility assumptions, which may not generalize to long, structured, or poorly conserved ncRNAs.
365 In contrast, *CUPID* does not attempt to reconstruct secondary structure but leverages contextual
366 sequence statistics learned from large corpora. These complementary perspectives suggest potential
367 synergies: coupling language-model embeddings with coarse structural predictions could refine the
368 discrimination between spurious and functionally relevant pairing events.

369 Despite these promising results, we note that the resources used to train *CUPID* are limited in size
370 and exhibits a strong imbalance across interaction types. Although our type-constrained negative
371 sampling preserves the empirical distribution of interaction types, rare types remain challenging; they
372 can yield higher-variance estimates and may prevent the model from learning robust type-specific
373 patterns. Accordingly, we emphasize AUPRC in our per-type analyses, as it is generally more
374 informative than AUROC under severe class imbalance. Future work will benefit from larger and
375 more balanced interaction resources, and could further improve stability on underrepresented classes
376 via targeted strategies such as class-aware reweighting, resampling, or cost-sensitive objectives.

377 As larger ncRNA catalogs become available through resources such as RNAcentral Sweeney
378 et al. (2020), and as experimental protocols expand the coverage of ncRNA–ncRNA interactions,
379 the training regime of models like *CUPID* can be scaled accordingly. Future developments may
380 integrate longer receptive fields, explicit cross-attention between molecules, or joint fine-tuning on
381 experimentally resolved interactomes. These extensions could help reveal constraints underlying
382 ncRNA recognition and improve the resolution of regulatory maps in eukaryotic transcriptomes.

383 In addition, while our study focuses on a resource-efficient paradigm that leverages pretrained
384 RNA language models with a lightweight interaction-specific prediction head, it would be interesting
385 to complement our analysis with baselines that train a long-context Transformer from scratch. We
386 did not include such a baseline here because, under the current supervision regime (approximately
387 $10^5$ interaction pairs after filtering), end-to-end training from random initialization may be difficult
388 to optimize and may not yield generalizable representations. As larger and more diverse labeled
389 interaction resources become available, systematic comparisons between pretrained and from-scratch
390 Transformer encoders will become increasingly informative.

391 A similar consideration holds when considering studies substituting RNA-LM models with several
392 Transformer-based nucleotide language models. While these models could, in principle, be considered
393 as alternative backbones for RNA sequence embeddings (e.g., models pretrained predominantly on
394 DNA such as Nucleotide Transformer, which has been reported to transfer RNA-related signals (Dalla-
395 Torre et al., 2025)), we selected GenerRNA because it is pretrained specifically on RNA sequences,
396 provides a long-context representation and it is expected to better capture RNA-class-specific
397 features. We therefore expect RNA-specialized pretraining to yield representations that are more
398 directly tailored to RNA sequence regularities than more generic DNA-pretrained alternatives,
399 even when the latter can capture some RNA features. In this work, we focused on characterizing
400 the proposed interaction-prediction pipeline using a single RNA-specialized backbone, including

ablations on augmentation and pooling. As larger and more diverse interaction resources become available, it will be important to benchmark GenerRNA in a zero-shot setting against more general nucleotide Transformers, and to evaluate both backbones also after task-specific fine-tuning.

In summary, the results show that *CUPID* provides a scalable sequence-based framework for ncRNA–ncRNA interaction prediction, achieving AUROC larger than 0.9 for several types on ncRNA interactions. Its performance, robustness to class heterogeneity, and limited dependence on domain-specific priors make it suitable for large-scale in silico screening and for guiding targeted experimental profiling of ncRNA regulatory networks.

## CONFLICT OF INTEREST STATEMENT

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## AUTHOR CONTRIBUTIONS

M.N. E.C. and G.V. designed the overall work; M.N. and F.S. developed and implemented the AI and computational methods; M.N. conceived and executed the experiments; M.N. F.S. E.C. and G.V. analyzed the results; G.V. M.N. and E.C. drafted and wrote the paper; G.V. and E.C. supervised the overall work; all the authors revised and approved the final manuscript.

## FUNDING

## DATA AVAILABILITY STATEMENT

The data, the *CUPID* code, and the scripts to reproduce the experiments and tutorials are available from GitHub: `https://github.com/AnacletoLAB/ncRNA-CUPID`.

## REFERENCES

Aleksander, S., Balhoff, J., Carbon, S., et al. (2023). The gene ontology knowledgebase in 2023. *Genetics* 224, iyad031. doi:10.1093/genetics/iyad031

Ali, S., Peffers, M., Ormseth, M., et al. (2021). The non-coding RNA interactome in joint health and disease. *Nat Rev Rheumatol* 17, 692–705. doi:10.1038/s41584-021-00687-y

Alipanahi, B., Delong, A., Weirauch, M. T., and Frey, B. J. (2015). Predicting the sequence specificities of dna-and rna-binding proteins by deep learning. *Nature Biotechnology* 33, 831–838

430 Cavalleri, E., Cabri, A., Soto-Gomez, M., Bonfitto, S., Perlasca, P., Gliozzo, J., et al. (2024). An
431 ontology-based knowledge graph for representing interactions involving rna molecules. *Scientific*
432 *Data* 11, 906

433 Cetin, S. and Sefer, E. (2025). A graphlet-based explanation generator for graph neural networks
434 over biological datasets. *Current Bioinformatics* 20, 840–851

435 Chen, J., Hu, Z., Sun, S., Tan, Q., Wang, Y., Yu, Q., et al. (2022). Interpretable rna foundation
436 model from unannotated data for highly accurate rna structure and function predictions. *arXiv*
437 *preprint arXiv:2204.00300*

438 Dalla-Torre, H., Gonzalez, L., Mendoza-Revilla, J., et al. (2025). Nucleotide Transformer: building
439 and evaluating robust foundation models for human genomics. *Nat Methods* 22, 287–297. doi:10.
440 1038/s41592-024-02523-z

441 Engreitz, J. M., Sirokman, K., McDonel, P., Shishkin, A. A., Surka, C., Russell, P., et al. (2014). Rna-
442 rna interactions enable specific targeting of noncoding rnas to nascent pre-mrnas and chromatin
443 sites. *Cell* 159, 188–199

444 Fabbri, M., Girnita, L., Varani, G., and Calin, G. A. (2019). Decrypting noncoding rna interactions,
445 structures, and functional networks. *Genome research* 29, 1377–1388

446 Gargano, M., Matentzoglu, N., Coleman, B., et al. (2023). The Human Phenotype Ontology in
447 2024: phenotypes around the world. *Nucleic Acids Research* 52, D1333–D1346. doi:10.1093/nar/
448 gkad1005

449 Gliozzo, J., Soto Gomez, M. A., Bonometti, A., et al. (2025). miss-SNF: a multimodal patient
450 similarity network integration approach to handle completely missing data sources. *Bioinformatics*
451 41, btaf150. doi:10.1093/bioinformatics/btaf150

452 Li, Z., Li, K., Lian, X., and Li, J. (2025). Lncrna-mirna interaction prediction based on multi-source
453 heterogeneous graph neural network and multi-level attention mechanism. *International Journal*
454 *of Biological Macromolecules* 319, 145614. doi:https://doi.org/10.1016/j.ijbiomac.2025.145614

455 Lorenzi, L., Chiu, H.-S., Cobos, A., et al. (2021). The rna atlas expands the catalog of
456 human non-coding rnas. *Nature Biotechnology* 39, 1453–1465. doi:https://doi.org/10.1038/
457 s41587-021-00936-1

458 Mann, M., Wright, P. R., and Backofen, R. (2017). Intarna 2.0: enhanced and customizable
459 prediction of rna–rna interactions. *Nucleic acids research* 45, W435–W439

460 Matthews, B. W. (1975). Comparison of the predicted and observed secondary structure of t4 phage
461 lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure* 405, 442–451

462 Nicolini, M., Saitto, E., Jimenez-Franco, R., et al. (2025a). Fine-tuning of conditional Transformers
463 improves in silico enzyme prediction and generation. *Computational and Structural Biotechnology*
464 *Journal* 27, 1318–1334. doi:10.1016/j.csbj.2025.03.037

465 Nicolini, M., Stacchietti, F., Cano, C., Casiraghi, E., and G, V. (2025b). A transformer-based
466 model to predict micro rna interactions. In *18th International Work-Conference on Artificial*
467 *Neural Networks, IWANN 2025*. vol. 16008 of *Lecture Notes in Computer Science*. doi:10.1007/
468 978-3-032-02725-2_8

469 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. (2019). Language models
470 are unsupervised multitask learners. *OpenAI blog* 1, 9

471 Sapoval, N., Aghazadeh, A., Nute, M., et al. (2022). Current progress and open challenges for
472 applying deep learning across the biosciences. *Nat Commun* 13. doi:http://doi.org/10.1038/
473 s41467-022-29268-7

Sefer, E. (2025). Drgat: Predicting drug responses via diffusion-based graph attention network. *Journal of Computational Biology* 32, 330–350

Sennrich, R., Haddow, B., and Birch, A. (2016). Neural Machine Translation of Rare Words with Subword Units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, eds. K. Erk and N. A. Smith (Berlin, Germany: Association for Computational Linguistics), 1715–1725. doi:10.18653/v1/P16-1162

Shen, T., Hu, Z., Sun, S., Liu, D., Wong, F., Wang, J., et al. (2024). Accurate RNA 3D structure prediction using a language model-based deep learning approach. *Nature Methods* , 1–12

Sun, S.-L., Jiang, Y.-Y., Yang, J.-P., Xiu, Y.-H., Bilal, A., and Long, H.-X. (2025). Predicting noncoding rna and disease associations using multigraph contrastive learning. *Scientific Reports* 15, 230

Sweeney, B., Petrov, A., Ribas, C., et al. (2020). RNAcentral 2021: secondary structure integration, improved sequence search and new member databases. *Nucleic Acids Research* 49, D212–D220. doi:http://dx.doi.org/10.1093/nar/gkaa921

Tian, X., Shen, L., Wang, Z., Zhou, L., and Peng, L. (2021). A novel lncrna–protein interaction prediction method based on deep forest with cascade forest structure. *Scientific reports* 11, 18881

Torgano, F., Soto-Gomez, M., Zignani, M., Gliozzo, J., Cavalleri, E., Mesiti, M., et al. (2025). Rna knowledge-graph analysis through homogeneous embedding methods. *Bioinformatics Advances* 5, vbaf109. doi:10.1093/bioadv/vbaf109

Umu, S. U. and Gardner, P. P. (2017). A comprehensive benchmark of rna–rna interaction prediction tools for all domains of life. *Bioinformatics* 33, 988–996

Valentini, G., Malchiodi, D., Gliozzo, J., Mesiti, M., Soto-Gomez, M., Cabri, A., et al. (2023). The promises of large language models for protein design and modeling. *Frontiers in Bioinformatics* 3

Vasilevsky, N., Toro, S., Matentzoglu, N., et al. (2025). Mondo: Integrating Disease Terminology Across Communities. *Genetics* , iyaf215doi:10.1093/genetics/iyaf215

Wang, W., Zhang, L., Sun, J., Zhao, Q., and Shuai, J. (2022). Predicting the potential human lncrna–mirna interactions based on graph convolution network with conditional random field. *Briefings in Bioinformatics* 23, bbac463. doi:10.1093/bib/bbac463

Wei, J., Chen, S., Zong, L., Gao, X., and Li, Y. (2022). Protein–rna interaction prediction with deep learning: structure matters. *Briefings in Bioinformatics* 23, bbab540. doi:10.1093/bib/bbab540

Yang, T., He, Y., and Wang, Y. (2025). Introducing tec-lncmir for prediction of lncrna-mirna interactions through deep learning of rna sequences. *Briefings in Bioinformatics* 26, bbaf046. doi:10.1093/bib/bbaf046

Yu, H., Yang, H., Sun, W., et al. (2024). An interpretable rna foundation model for exploring functional rna motifs in plants. *Nat Mach Intell* 6, 1616–1625. doi:10.1038/s42256-024-00946-z

Yu, X., Jiang, L., Jin, S., Zeng, X., and Liu, X. (2022). premli: a pre-trained method to uncover microrna–lncrna potential interactions. *Briefings in Bioinformatics* 23, bbab470. doi:10.1093/bib/bbab470

Zhao, Y., Oono, K., Takizawa, H., and Kotera, M. (2024). GenerRNA: A generative pre-trained language model for de novo RNA design. *PLoS One* 19, e0310814