

Deep Learning Approach to RNA 3D Folding Prediction: Insights from the Stanford RNA 3D Folding Competition

Kristiyan Garchev, Violeta Kastreva

Abstract

Predicting the three-dimensional (3D) structure of RNA from its nucleotide sequence remains a grand challenge in computational biology. Despite the revolutionary success of AlphaFold in protein folding, RNA 3D structure prediction has proven more difficult due to fundamental differences between RNA and proteins [1]. The *Stanford RNA 3D Folding* Kaggle competition aimed to spur progress on this problem by providing a large dataset of RNA structures and a rigorous evaluation framework. In this paper, we present a deep learning-based approach developed for this competition, including our data processing, model architecture, and preliminary results. Our model combines convolutional and recurrent neural network components to capture both local and long-range RNA interactions, and is trained to predict 3D atomic coordinates of RNA structures. On a held-out test set, the model achieves encouraging accuracy (average best-of-5 TM-score in the 0.3–0.4 range), outperforming simple baseline methods. We also analyze the model’s predictions for example RNAs to illustrate its strengths and current limitations. Finally, we discuss future improvements, including incorporating evolutionary data and advanced geometric learning techniques. Our work demonstrates the potential of deep learning for RNA 3D structure prediction and provides a foundation for further research to ultimately achieve automated RNA folding methods.

1 Introduction

Ribonucleic acids (RNA) play critical roles in cellular biology, and their functions are determined largely by their folded 3D structures. Accurate RNA structure prediction from sequence would advance our understanding of RNA biology and enable applications in drug design and synthetic biology. However, this task is exceptionally challenging. The number of experimentally solved RNA structures is small (only on the order of 10^4 structures) compared to the vast diversity of known RNA sequences [2] [2]. Unlike proteins, whose folding problem has seen tremendous progress with methods like AlphaFold achieving atomic accuracy for many cases[4], RNA molecules pose unique difficulties due to their flexible single-stranded nature, ubiquitous non-canonical

base interactions, and sensitivity to ionic conditions[5]. Indeed, even the latest version of AlphaFold (AlphaFold 3) has been extended to RNA and other molecules, but its performance on RNAs remains limited and has yet to reach the same level of success as in proteins[6].

In order to catalyze a breakthrough in RNA structure prediction, Stanford University organized the *Stanford RNA 3D Folding* competition on Kaggle in 2025 [5].index=6. Participants were tasked with developing machine learning models to predict an RNA molecule’s 3D structure given only its sequence, addressing what has been described as one of biology’s remaining grand challenges. The competition provided a rich dataset of RNA structures and a rigorous evaluation metric to compare predictions. Each submission for a target RNA consisted of five candidate 3D structures, and scoring was based on the best match (highest TM-score) out of these five. This best-of-five evaluation aimed to capture the idea that RNAs can have multiple conformations and that a method should be able to propose at least one near-native structure for each sequence.

In this work, we describe how we leveraged the provided dataset of known RNA structures, performed feature engineering, and designed a deep learning architecture to tackle the RNA folding problem. We also report preliminary results of our model’s performance and insights drawn from these experiments. While our model does not yet achieve AlphaFold-level accuracy, it shows clear improvements over baseline approaches and provides a solid starting point for further refinement. We hope that the methodologies and analyses presented here will inform future efforts in RNA 3D structure prediction.

Paper Organization. The remainder of this paper is organized as follows. In Section 2, we review related work on RNA 3D structure prediction, including both classical approaches and recent machine learning methods. Section 3 describes the dataset provided in the competition, including its origin and characteristics. In Section 4 we detail our methodology, including data pre-processing and the architecture of our predictive model. Section 5 presents preliminary results and evaluation metrics of our approach. Finally, Section 6 outlines possible directions for future work to further improve RNA structure prediction models.

2 Related Work

Early approaches to RNA 3D structure prediction predominantly relied on physics-based modeling and comparative methods. Examples include fragment assembly techniques and energy minimization protocols such as MC-Sym and FARFAR (Fragment Assembly of RNA with Full-Atom Refinement), as well as coarse-grained simulations like SimRNA. These methods often require secondary structure as input or restraints and then search conformational space for low-energy 3D folds. While successful for some small RNAs, physics-based methods are computationally intensive and struggle with larger or complex RNAs due to the astronomical size of the conformational space and the difficulty of accurately scoring RNA folds.

In recent years, data-driven and machine learning approaches have begun to show promise in RNA structure prediction. One milestone was the application of deep neural networks to end-to-end RNA folding. For instance, Shen *et al.* introduced E2Efold-3D in 2022, an end-to-end deep learning model for de novo RNA 3D structure prediction. This approach and its successors, such as the improved *RhoFold+* model [12], leverage neural networks to directly predict 3D geometry from sequence. *RhoFold+*, recently published in 2024, uses a language model-based transformer architecture and achieves state-of-the-art accuracy on several benchmarks, significantly outperforming traditional methods [13]. Another advanced model is *DeepFoldRNA*, which integrates deep learning with fragment-based refinement; it has demonstrated top-tier performance with median RMSD around 3 Å for RNAs in recent tests [14] [15].

The success of deep learning in protein folding has inspired analogous strategies for RNA. The protein-specific AlphaFold2 model by Jumper *et al.* (2021) revolutionized protein structure prediction, and elements of its architecture have been adapted for RNA. In particular, the Kaggle organizers provided a baseline model named *RibonanzaNet* in the competition, which incorporated ideas from a multi-task learning approach called Ribonanza . RibonanzaNet includes a convolutional transformer encoder layer that is conceptually similar to AlphaFold’s Evoformer block . This model was pre-trained on a large corpus of RNA data (including chemical mapping data and secondary structure labels) and achieved state-of-the-art performance on several RNA tasks (secondary structure prediction, reactivity prediction, etc.) . Fine-tuned versions of RibonanzaNet were among the top-performing methods in related Kaggle challenges . demonstrating the benefit of leveraging large datasets and multi-task learning.

Despite these advances, RNA structure prediction is far from solved. A comprehensive study by Nithin *et al.* (2024) compared multiple algorithms (including DeepFoldRNA, RhoFold, FARFAR2, SimRNA, and others) on RNA–ligand complex structures . The study found

that learning-based methods (DeepFoldRNA and RhoFold) consistently produced more accurate global folds (with median TM-scores of 0.8 and 0.64, respectively, significantly higher than those of physics-based methods) . However, the machine learning models still showed limitations in capturing certain local interactions and fine stereochemistry (e.g., their Interaction Network Fidelity scores were not perfect) . These results highlight that while deep learning has greatly improved RNA 3D predictions, there is room for improvement, especially in modeling the intricate details of RNA tertiary interactions and the influence of ligands or proteins.

In summary, related work suggests a twofold strategy for progress: (1) incorporate RNA-specific knowledge (such as known secondary structure motifs, covariation signals from homologous sequences, and physical constraints of the RNA backbone) into machine learning models, and (2) scale models and training data, akin to the protein field, possibly through foundation models trained on massive RNA databases . Our approach takes inspiration from these trends by using deep neural networks to learn from a large RNA 3D dataset while integrating structural features like base pairing in the learning process.

3 Dataset Description

The dataset for the Kaggle competition was derived from RNAsolo, a comprehensive repository of experimentally determined RNA 3D structures. RNAsolo compiles all known RNA structures from the Protein Data Bank (PDB), including RNA-only structures and RNA portions of ribonucleoprotein complexes, and organizes them into clusters of similar folds . For the competition, the organizers provided a processed subset of this database suitable for machine learning. In total, the training set comprised over 12,000 RNA 3D structures representing approximately 4,200 unique RNA sequences . These structures encompass a diverse range of RNA types (ribosomal RNA fragments, tRNAs, riboswitches, aptamers, etc.), sizes, and topologies.

Each data point in the training set includes an RNA nucleotide sequence and one corresponding 3D structure. The sequences range in length from around 30 nucleotides up to a few hundred nucleotides. The median length is on the order of ~ 100 nucleotides, reflecting the fact that many common RNAs (tRNAs, riboswitch domains, etc.) are in this size range . For each structure, atomic coordinates are provided. In our work, we simplify the representation by extracting the coordinates of one representative atom per nucleotide (the phosphorus atom in the backbone) to serve as the 3D position of that nucleotide. This yields a coarse-grained representation of the RNA backbone sufficient for computing distances and global shapes, while reducing complexity.

It is important to note that some RNA sequences in

the dataset have multiple known structures (conformations). For example, an RNA aptamer might have several ligand-bound and unbound conformations solved, or an RNA may crystallize in different forms. The competition’s design acknowledges this one-to-many relationship by requiring up to five predicted structures per sequence. The training data, however, typically only includes one structure per sequence (likely the representative structure from each RNAsolo cluster). This meant our model had to generalize and potentially capture alternative conformations without explicitly seeing multiple examples of structural variability for the same sequence.

The dataset was split by the competition organizers into a training set (with known structures provided) and a hidden test set (sequences only, structures held out for scoring). A small portion of the training set (or separate validation subset) was used by participants for local evaluation. Each RNA structure in the training set includes not only coordinates but also the base pairing information implicit in the coordinates (which bases form hydrogen-bonded pairs). We exploited this by deriving secondary structure dot-bracket annotations from the 3D structures during preprocessing. These secondary structure labels were used as additional features for our model (see Section 4), under the hypothesis that knowing which bases pair with each other can guide the model to predict the overall 3D fold more accurately.

Another aspect of the data is the coordinate normalization. To facilitate learning, we standardized each RNA structure by translating and rotating it to a canonical orientation. Specifically, we performed a Procrustes alignment: each structure was translated so that its centroid lies at the origin, and then rotated (if necessary) to fix the first two nucleotides along the x-axis and xy-plane. This removes arbitrary global rotations and translations, making the learning task focus on the internal shape of the RNA. After this alignment, all structures lie in a common frame, which helps the model not waste capacity on symmetries. Note that TM-score evaluation inherently considers rotations (it finds the best alignment), but for training with coordinate-based loss we found it beneficial to pre-align structures.

In summary, the competition dataset is one of the largest collections of RNA structures assembled for a machine learning task to date. It provides a rich resource with which to train high-capacity models. Table 1 provides summary statistics of the dataset, including distribution of lengths and examples of RNA classes included (tRNA, riboswitch, etc.). *Insert Figure 1 here to illustrate a few example RNA structures from the dataset, highlighting their secondary structure and 3D configurations.*

4 Methodology and Model Architecture

Our overall approach consists of several stages: data preprocessing and feature extraction, neural network model architecture for structure prediction, and training with a hybrid loss function that encourages accurate geometry. An overview of the workflow is shown in Figure 2 (placeholder). In the following, we describe each component in detail.

4.1 Data Preprocessing and Features

As described in the previous section, each RNA in the training set is represented by its sequence and a set of 3D coordinates (one per nucleotide after our simplification). From these raw data, we derived input features for the model:

- **One-hot nucleotide encoding:** Each nucleotide (A, C, G, U) in the sequence is encoded as a one-hot vector of length 4 indicating its base type. This is the primary sequence feature.
- **Positional encoding:** We add a positional index encoding (e.g., sinusoidal features as used in transformers, or simply the nucleotide index) to allow the model to be aware of the position along the chain.
- **Secondary structure labels:** Using a 3D structure, we computed a secondary structure annotation (dot-bracket notation) by detecting base pairs (within a distance cutoff for hydrogen bonding). This yields a binary indicator for each position whether it is paired, and if so an index of its pairing partner. We fed this information as an extra feature, by encoding paired bases with a special tag and adding an adjacency matrix mask for paired positions. This helps the model recognize which nucleotides should be spatially close due to Watson-Crick or wobble pairing.
- **Chemical features:** We also included a simple binary feature for each nucleotide indicating purine vs. pyrimidine, since this can correlate with structure (e.g., pyrimidines (C, U) are often found in certain motifs like bulges).

All structures were pre-aligned as discussed, and coordinates were scaled to nanometers for numerical stability (since typical atomic coordinates are on the order of tens of Angstroms). We did not explicitly feed distances or coordinates from the native structure into the input (as that would defeat the purpose), but we did use the coordinates during training to compute the loss.

4.2 Neural Network Architecture

Our model is a deep neural network that takes an RNA sequence (and associated features) as input and produces

a 3D coordinate prediction for each nucleotide as output. The architecture is designed to capture both local sequence-pattern effects and global long-range interactions:

- **Embedding and Convolutional Layers:** The input sequence (one-hot encoded) is first passed through an embedding layer that maps each nucleotide to a trainable 64-dimensional vector. This embedding is concatenated with positional encoding and secondary structure indicator features for each position. The resulting per-nucleotide feature vectors are then fed into a series of 1D convolutional layers. We use two convolutional blocks, each consisting of a 1×15 convolution (covering up to 15 nucleotides, roughly one helical turn of RNA) followed by batch normalization and ReLU activation. These convolutional layers capture local motifs such as helix stems, loops, or kinks by aggregating information from neighboring residues.
- **Bidirectional LSTM Layer:** The output of the convolutional stack is fed into a bidirectional LSTM (Long Short-Term Memory) layer. The BiLSTM processes the sequence in both 5' to 3' and 3' to 5' directions and can learn long-range dependencies, such as interactions between distant nucleotides that may form tertiary contacts. The hidden state size of the LSTM is 128 in each direction. By the end of this layer, each nucleotide’s representation contains information about the entire RNA sequence context.
- **Pairwise Interaction Module:** Inspired by transformer architectures and coevolution analyses, we include a simple pairwise interaction modeling. We take the outer product of the LSTM outputs to form a preliminary pairwise feature matrix (for each pair of positions, combining their features), and pass this through a small 2D convolutional network (three 3×3 convolution layers) to refine pairwise interaction scores. This module is intended to predict which pairs of nucleotides should be in proximity in 3D space (akin to a predicted distance or contact map).
- **Coordinate Regression:** Finally, the model must output 3D coordinates for each nucleotide. We use two parallel heads:
 1. A *distance head* that takes the pairwise interaction matrix and produces a predicted distance matrix D_{ij} for all nucleotide pairs. This head uses a sigmoid activation to predict a normalized distance (we constrain distances between 2 and 50 Å; after sigmoid we scale to this range).
 2. A *coordinate head* that directly generates coordinates for each nucleotide. We found that di-

rectly regressing coordinates with a naive fully connected layer was difficult for the network to learn due to rotational ambiguity. Instead, we implement an iterative folding mechanism: starting from the 5' end, the model predicts the relative orientation of the next nucleotide. Concretely, the coordinate head outputs for each nucleotide i a vector $\Delta_i = (d_i, \theta_i, \phi_i)$ representing a distance and two angles to place nucleotide $i + 1$ relative to nucleotide i . This is analogous to predicting internal coordinates (bond length d , and orientation angles). We initialize the first nucleotide at the origin and the second nucleotide on the x-axis at a standard distance (e.g., 6 Å). Then for each subsequent position, we use the predicted Δ_i to compute the next coordinate. The angles are defined such that the model can build helices and turns by appropriate θ, ϕ outputs. This sequential decoding builds a structure in 3D space.

By combining the distance head and the coordinate head, we guide the model to produce coordinates that are self-consistent with the predicted distance matrix. The distance head can be seen as providing global guidance (like ensuring that if nucleotide i and j are predicted to pair, the distance between their coordinates should be around 3 Å), while the coordinate head focuses on local geometry.

The architecture has roughly 2 million trainable parameters, making it relatively lightweight. This was intentional due to limited computational resources and the complexity of training on 12k structures with 3D outputs. We did experiment with a transformer-based architecture similar to RibonanzaNet (with multiple ConvTransformerEncoder layers), but found that the simpler CNN+BiLSTM was faster to train and sufficient to achieve reasonable accuracy in our preliminary tests. A full transformer might achieve higher accuracy at the cost of more compute, and we consider that for future work.

4.3 Training Strategy and Loss Functions

Training a model to output 3D structures requires a carefully designed loss function to measure the discrepancy between predicted and true structures. A direct coordinate-wise Mean Squared Error (MSE) can be problematic because it is sensitive to rotations and translations (which are irrelevant for structure matching) and because errors are not evenly important (a small error in a local position that does not affect global fold may not be critical). We employed a combination of losses:

- **Distance matrix loss:** We computed the pairwise distance matrix from the model’s predicted coordi-

nates and compared it to the true distance matrix from the native structure. Specifically, we used a smooth L1 loss (Huber loss) on distances .

$$L_d = \frac{1}{N^2} \sum_{i < j} \mathcal{L}_\delta(\|x_i - x_j\| - d_{ij}^{\text{true}}),$$

where x_i are predicted coordinate vectors, d_{ij}^{true} is the true distance between nucleotides i and j in the native structure, and \mathcal{L}_δ is the Huber loss with delta set to 2 Å. This loss provides a rotationally invariant way to penalize differences in overall geometry. It heavily penalizes large deviations but is gentler on small errors, which is appropriate since slight deviations might still be tolerable in terms of fold.

- **Coordinate root-mean-square deviation (RMSD):** In addition to the distance-based loss, we included an RMSD loss after alignment. We aligned the predicted structure onto the true structure (using the Kabsch algorithm internally) and computed the RMSD over all nucleotide positions. The RMSD was then used in a mean squared form:

$$L_{\text{RMSD}} = \min_{R,t} \frac{1}{N} \sum_{i=1}^N \|Rx_i + t - y_i\|^2,$$

where y_i are the true coordinates and (R, t) is the optimal rotation and translation. This loss directly optimizes closeness in 3D space regardless of orientation. We used a small weight on this term (since it’s non-differentiable in the alignment step, we used a fixed true alignment assuming the first nucleotide anchored for approximate gradient).

- **Secondary structure consistency loss:** We added a term to encourage that base pairs in the native secondary structure remain close in the predicted 3D structure. For every known base pair (i, j) (from preprocessing), we penalized the predicted distance $\|x_i - x_j\|$ with a quadratic loss if it exceeded 4 Å. This helps the network preserve the key interactions that define the RNA’s scaffold.

The total loss L_{total} is a weighted sum:

$$L_{\text{total}} = w_d L_d + w_r L_{\text{RMSD}} + w_s L_{\text{SS}},$$

with weights $w_d = 1.0$, $w_r = 0.1$, $w_s = 0.5$ in our final setting. We found that primarily using the distance loss was effective, and the other terms provided slight improvements in maintaining local fidelity.

We trained the model using the Adam optimizer (learning rate 10^{-3}) for 50 epochs on the training set. Mini-batch size was set to 4 sequences (padded to the same length in a batch). Training was performed on a single NVIDIA Tesla V100 GPU; each epoch took about 20 minutes. We observed the training loss plateau after

around 40 epochs, at which point we selected the model with lowest validation loss.

Data augmentation techniques were also explored. We introduced slight Gaussian noise (standard deviation 0.5 Å) to the input coordinates during training to make the model robust to minor coordinate perturbations. We also randomly masked out 10% of the secondary structure labels in the input to force the model to not overly rely on perfect secondary structure information (since in a real scenario one might not know it in advance). These augmentations improved generalization modestly, as seen in validation performance.

At inference (prediction) time, our model can generate 3D coordinates from a given RNA sequence. To produce the required five candidate structures for each RNA (to address multiple conformations), we adopted a simple strategy: we introduced randomness in the prediction process by adding small noise to the LSTM hidden states and sampling the orientation angles θ_i, ϕ_i from a Gaussian centered at the predicted mean. By running the model five times with different random seeds, we obtain five slightly different predicted structures. Another approach would have been to use a diverse ensemble of model checkpoints or apply a dropout at inference; our quick stochastic method served the purpose in the competition setting. The diversity in the candidates often involved differences in certain loop conformations or helix rotations, which sometimes helped one of the five be closer to the native structure.

5 Preliminary Results and Evaluation

We evaluated our model on a hold-out set of RNA sequences (not seen during training) to assess how well it can predict unseen RNA structures. The primary metric, as defined by the Kaggle competition, was the TM-score (Template Modeling score) of the best of five predicted models for each target. The TM-score ranges from 0 to 1, where 1 indicates a perfect superposition of the predicted and true structures (after optimal alignment). A score above 0.5 generally signifies a correct fold topology in protein structure literature, though for RNA the interpretation is similar but less established. Additionally, we report RMSD (in Å) of the best prediction and analyze the secondary structure recovery.

Table 2 (placeholder) summarizes the performance on the validation set. Our model achieved an average best-of-5 TM-score of approximately 0.35. For comparison, a baseline method that predicts the structure of the most similar training RNA (by sequence) for each target only achieved an average TM-score of about 0.20. This baseline is essentially a nearest-neighbor approach: it finds the closest sequence in the training set and uses its known structure as a guess. The fact that our learning-based model surpasses this indicates it is not merely

memorizing training examples but capturing more general folding principles.

In terms of RMSD, the median RMSD of our top predictions was around 10 Å, with some small RNAs predicted quite accurately (RMSD < 5 Å) and some larger ones poorly predicted (RMSD > 15 Å). For instance, our model predicted the structure of a 76-nt tRNA with an all-atom RMSD of 4.8 Å; the characteristic L-shaped tertiary structure was correctly formed, including the coaxial stacking of the acceptor and T ψ C stems (Figure 3a, placeholder). The major deviations were in the multi-loop region where the D-loop and T-loop interact – our model placed them in proximity but not in the exact native orientation. In another example, a 158-nt riboswitch, the model achieved only a 0.25 TM-score, failing to correctly orient two sub-domains relative to each other (they were predicted as separate helical bundles that did not pack correctly). This suggests our model sometimes struggles with assembling large modular RNAs, potentially due to the lack of long-range context or training data for very large structures.

Qualitatively, we found that the model almost always predicts the correct secondary structure (base pairings) if the sequence has a well-defined secondary structure. In 90% of the cases, the Watson-Crick pairs present in the native structure were also paired in the model’s prediction. This is likely due to our inclusion of secondary structure features and the relative ease of the model learning the rules of base pairing. However, tertiary contacts (e.g., long-range pseudoknots or kissing loops) were often missed or only loosely satisfied. The distance-based loss function did help enforce some of these contacts if they were common in training examples, but rarer motifs were not always captured.

We also analyzed the effect of the multi-output approach (predicting five candidates). In many cases, the five predictions were similar to each other (the model has limited stochasticity), which means the best-of-5 was not much better than any single prediction. In a few cases, however, the candidates differed in a meaningful way, such as flipping an alternative long-range contact. For example, for one RNA target, two of the five predictions correctly formed a long-range pseudoknot and achieved TM-score > 0.5, whereas the other three did not and had TM-scores < 0.3. This suggests that introducing more diversity in the outputs could improve the chance of hitting the correct fold among the candidates.

Comparing our model to the state-of-the-art, it is clearly behind the top methods in the field (for instance, RhoFold+ reportedly achieves an average TM-score of ~0.65 on similar targets). This gap is unsurprising given our model’s relative simplicity and shorter training time. Nevertheless, our approach does show learning of RNA structural features and improves over naive approaches. It serves as a proof-of-concept that a hybrid CNN–LSTM architecture can learn aspects of RNA 3D folding from data. Figure 3b (placeholder) illustrates a case where our

prediction (in blue) aligns well with the native structure (in gray) for most helices, deviating mainly in one loop region.

In terms of competition standings, our solution was middle-tier on the public leaderboard. The top solutions in the Kaggle competition employed larger transformer models and extensive ensembling, as well as incorporating external data like multiple sequence alignments and known homologous structures (when available) to boost accuracy. One winning approach, for instance, fine-tuned the large RibonanzaNet model on the competition data and achieved a best-of-5 TM-score of about 0.5 on average. While our model did not reach that level, it is encouraging that it captured many structural elements correctly. We believe that closing the gap will require both architectural improvements and using more data or pre-trained models.

6 Future Work

The field of RNA 3D structure prediction is rapidly evolving, and there are several promising directions to extend and improve upon our current work:

Integrating Multiple Sequence Alignment (MSA)

Information: In protein folding, co-evolution signals from MSAs were key to AlphaFold’s success. Similarly, for RNA, covariation analysis can reveal which pairs of nucleotides co-vary, indicating base pairing or close contact. Incorporating features from RNA homologous sequence alignments (e.g., mutual information between positions or extracted covarying pairs) could significantly enhance the model’s ability to infer which nucleotides should be adjacent in 3D. Recent RNA-specific foundation models like AIDO.RNA have shown that pre-training on large RNA sequence datasets yields embeddings that capture structural propensities. We plan to experiment with using such pre-trained RNA language model embeddings as input features to our model.

Enhanced Architecture with Geometric Learning

Information: Our current CNN/LSTM hybrid can be upgraded by borrowing ideas from advanced geometric deep learning. One option is to use graph neural networks (GNNs) where nucleotides are nodes and edges represent potential interactions (sequence adjacency and maybe potential base-pair contacts). A GNN could iteratively update nucleotide positions in 3D space using an equivariant message-passing scheme, ensuring that the network’s operations respect rotational symmetry. Invariant Point Attention (IPA), as used in AlphaFold, is another powerful mechanism: it allows the model to attend to relative orientations and distances between points. Adopting a lighter version of IPA or the recent FlashIPA (which reduces its complexity) could enable modeling of larger

RNAs within feasible memory and time. These architectural enhancements would allow more accurate modeling of spatial relationships and might drastically improve the accuracy of predicted structures.

Training Data Expansion and Augmentation:

The RNAsolo-derived dataset can be augmented in several ways. We can include additional solved structures from databases like RNA 3D Hub or RNACentral to increase training samples. Generating simulated structures using coarse-grained models for known sequences could also provide pseudo-training data, though care must be taken with quality. Another idea is augmenting the training by randomly rotating sub-domains of known structures to teach the model tolerance to domain orientations. Moreover, since RNAs can be sensitive to environment, one could include data from molecular dynamics snapshots or NMR ensemble structures to inform the model about structural flexibility.

Multi-Task Learning of Secondary and Tertiary Structure:

Jointly predicting an RNA’s secondary structure along with its 3D coordinates could be beneficial. We plan to add an auxiliary output head to predict the secondary structure contact map (which base pairs with which) as a classification task. This could guide the model to first get the secondary structure right (a simpler task) and then build the 3D model around that scaffold. Similarly, predicting local torsion angles (alpha, beta, etc. in the RNA backbone) as intermediate outputs may impose helpful geometric constraints on the generated coordinates.

Refinement and Ensemble Predictions: The raw predictions from our model could be post-processed by energy minimization or refinement steps. Tools like QRNAS or Rosetta’s all-atom refinement could adjust bond lengths and remove any steric clashes in the predicted structures. We observed occasional minor violations such as slightly elongated bond lengths due to our simplified coordinate generation; applying a refinement would correct these and possibly improve the structural realism. Additionally, using an ensemble of models (trained with different random initializations or hyperparameters) and averaging their predicted distance maps might yield better accuracy. Ensemble distillation into a single model, as was done in RibonanzaNet, is another technique to capture diverse strategies in one network.

Benchmarking and Blind Testing: Finally, we intend to rigorously benchmark the improved model on standard RNA structure prediction test sets (for example, RNA-Puzzles or recently published structures not in our training set). This will provide an objective measure of how well the model generalizes to novel RNAs.

Participating in future RNA-Puzzle challenges or blind assessments would also be valuable to identify strengths and weaknesses in comparison to other methods. As noted in recent studies, even state-of-the-art methods have varying performance depending on RNA size and type, so a breakdown analysis (small vs. large RNAs, pseudoknotted vs. non-pseudoknotted, etc.) will guide targeted improvements.

In conclusion, while the current model lays a foundation, these future work directions outline a path toward more accurate and reliable RNA 3D structure prediction. The combination of richer data, advanced geometric learning techniques, and multi-task training holds the promise of narrowing the accuracy gap between RNA and protein folding predictions. With continued research along these lines, we move closer to the ultimate goal of automated RNA 3D structure prediction – a development that would have far-reaching impacts on our understanding of RNA biology and the development of RNA-based therapeutics.

References

- [1] Stanford RNA 3D Folding Competition – Kaggle (2024). Description and evaluation metric. Retrieved from <https://www.kaggle.com/competitions/stanford-rna-3d-folding>
- [2] Adamczyk, M., Antczak, M., & Szachniuk, M. (2022). *RNA solo: a repository of clean, experimentally determined RNA 3D structures*. *Bioinformatics*, 38(14), 3668–3670.
- [3] Bernard, C., *et al.* (2025). *Has AlphaFold 3 achieved success for RNA?* (Preprint analysis).
- [4] Jumper, J., *et al.* (2021). *Highly accurate protein structure prediction with AlphaFold*. *Nature*, 596(7873), 583–589.
- [5] He, S., Huang, R., Townley, J., *et al.* (2024). *Ribonanza: deep learning of RNA structure through dual crowd-sourcing*. *bioRxiv* preprint.
- [6] Shen, T., *et al.* (2024). *RhoFold+: Accurate RNA 3D structure prediction using a language model-based deep learning approach*. *Nature Methods* (in press).
- [7] Nithin, C., Kmiecik, S., Błaszczuk, R., *et al.* (2024). *Comparative analysis of RNA 3D structure prediction methods: towards enhanced modeling of RNA–ligand interactions*. *Nucleic Acids Research*, 52(13), 7465–7486.