# RiboDiffusion: Tertiary Structure-based RNA Inverse Folding with Generative Diffusion Models

**Han Huang,**[1,2,†] **Ziqian Lin,**[1,3,†] **Dongchen He,**[1] **Liang Hong**[1] **and Yu Li**[1,*]

[1]Department of Computer Science and Engineering, CUHK, Hong Kong SAR, China, [2]School of Computer Science and Engineering, Beihang University, Beijing, China and [3]Nanjing University, Nanjing, China

[*]Corresponding author. [†]Equal contribution.

## Abstract

RNA design shows growing applications in synthetic biology and therapeutics, driven by the crucial role of RNA in various biological processes. A fundamental challenge is to find functional RNA sequences that satisfy given structural constraints, known as the inverse folding problem. Computational approaches have emerged to address this problem based on secondary structures. However, designing RNA sequences directly from 3D structures is still challenging, due to the scarcity of data, the non-unique structure-sequence mapping, and the flexibility of RNA conformation. In this study, we propose RiboDiffusion, a generative diffusion model for RNA inverse folding that can learn the conditional distribution of RNA sequences given 3D backbone structures. Our model consists of a graph neural network-based structure module and a Transformer-based sequence module, which iteratively transforms random sequences into desired sequences. By tuning the sampling weight, our model allows for a trade-off between sequence recovery and diversity to explore more candidates. We split test sets based on RNA clustering with different cut-offs for sequence or structure similarity. Our model outperforms baselines in sequence recovery, with an average relative improvement of $11\%$ for sequence similarity splits and $16\%$ for structure similarity splits. Moreover, RiboDiffusion performs consistently well across various RNA length categories and RNA types. We also apply in-silico folding to validate whether the generated sequences can fold into the given 3D RNA backbones. Our method could be a powerful tool for RNA design that explores the vast sequence space and finds novel solutions to 3D structural constraints.

## 1. Introduction

The design of RNA molecules is an emerging tool in synthetic biology (Chappell et al., 2015; McKeague et al., 2016) and therapeutics (Zhu et al., 2022), enabling the engineering of specific functions in various biological processes. There have been various explorations into RNA-based biotechnology, such as translational RNA regulators for gene expression (Laganà et al., 2015; Chappell et al., 2017), aptamers for diagnostic or therapeutic applications (Espah Borujeni et al., 2016; Findeiß et al., 2017), and catalysis by ribozymes (Dotu et al., 2014; Park et al., 2019). While the tertiary structure determines how RNA molecules function, one fundamental challenge in RNA design is to create functional RNA sequences that can fold into the desired structure, also known as the inverse RNA folding problem (Hofacker et al., 1994).

Most early computational methods for inverse RNA folding focus on folding into RNA secondary structures (Churkin et al., 2018). Some programs use efficient local search strategies to optimize a single seed sequence for the desired folding properties, guided by the energy function (Hofacker et al., 1994; Andronescu et al., 2004; Busch and Backofen, 2006; Garcia-Martin et al., 2013).

Others attempt to solve the problem globally by modeling the sequence distribution or directly manipulating diverse candidates (Taneda, 2010; Kleinkauf et al., 2015; Yang et al., 2017; Runge et al., 2019). However, without considering 3D structures of RNA, these methods cannot meet accurate functional structure constraints, since RNA secondary structures only partially determine their tertiary structures (Vicens and Kieft, 2022). The pioneering work (Yesselman and Das, 2015) applies a physically-based approach to optimize RNA sequences and match the fixed backbones, but it is still constrained by the local design strategy and computational efficiency.

Recent advances in deep learning and the accumulation of biomolecular structural data have enabled computational methods to model mapping between sequences and 3D structures with extraordinary performance, as demonstrated by remarkable results in protein 3D structure prediction (Jumper et al., 2021; Lin et al., 2023) and inverse design (Dauparas et al., 2022). Inspired by this, the development of geometric learning methods on RNA structures has received increasing research interest. On the one hand, many studies have explored RNA tertiary structure

prediction using machine learning models with limited data (Shen et al., 2022; Baek et al., 2022; Li et al., 2023). On the other hand, although deep learning has a promising potential to narrow down the immense sequence space for inverse folding, developing an appropriate model for RNA inverse folding remains an open problem, as it requires capturing the geometric features of flexible RNA conformations, handling the non-unique mappings between structures and sequences, and providing alternative options for different design preferences.

In this study, we introduce RiboDiffusion, a generative diffusion model for RNA inverse folding based on tertiary structures. We formulate the RNA inverse folding problem as learning the sequence distribution conditioned on fixed backbone structures, using a generative diffusion model (Yang et al., 2022). Unlike previous methods that predict the most probable sequence for a given backbone (Ingraham et al., 2019; Jing et al., 2021; Gao et al., 2023; Joshi et al., 2023), our method captures multiple mappings from 3D structures to sequences through distribution learning. With a generative denoising process for sampling, our model iteratively transforms random initial RNA sequences into desired candidates under tertiary structure conditioning. This global iterative generation distinguishes our model from autoregressive models and local updating methods, enabling it to better search for sequences that satisfy global geometric constraints. We parameterize the diffusion model with a cascade of a structure module and a sequence module, to capture the mutual dependencies between sequence and structure. The structure module, based on graph neural networks, extracts SE(3)-invariant geometrical features from 3D fixed RNA backbones, while the sequence module, based on Transformer-liked layers, captures the internal correlations of RNA primary structures. To train the model, we randomly drop the structural module to learn both the conditional and unconditional RNA sequence distribution. We also mix the conditional and unconditional distributions in the sampling procedures, to balance sequence recovery and diversity for more candidates.

We use RNA tertiary structures from PDB database (Bank, 1971) to construct the benchmark dataset and augment it with predicted structures from the RNA structure prediction model (Shen et al., 2022). We split test sets based on RNA clustering using different sequence or structure similarity cutoffs. Our model achieves an 11% higher recovery rate than the machine learning baselines for benchmarks based on sequence similarity, and 16% higher for benchmarks based on structure similarity. RiboDiffusion also performs consistently well across different RNA lengths and types. Further analysis reveals its great performance for cross-family and in-silico folding. Our method could be a powerful tool for RNA design, exploring a wide sequence space and finding novel solutions to 3D structural constraints.

## 2. Methodology

This section will explain RiboDiffusion in detail - a deep generative model for RNA inverse folding based on fixed 3D backbones. The overview is shown in Fig. 1. We will first introduce the preliminaries of diffusion models and our formulations of the RNA inverse folding problem. We will then describe the design of neural networks to parameterize the diffusion model and explain the sequence sampling procedures.

## 2.1. Preliminary and Formulation

### 2.1.1. Diffusion Model

As a powerful genre of generative models, diffusion models (Sohl-Dickstein et al., 2015) have been successfully applied to the distribution learning of diverse data, including images (Ho et al., 2020; Song et al., 2021), graphs (Huang et al., 2022, 2023a), and molecular geometry (Watson et al., 2023; Huang et al., 2023b). As the first step of setting up the diffusion model, a forward diffusion process is constructed to perturb data with a sequence of noise. This converts the data distribution to a known prior distribution. With random variables $x_0 \in \mathbb{R}^d$ and a forward process $\{x_t\}_{t \in [0,T]}$, a Gaussian transition kernel is set as

$$q_{0t}(x_t|x_0) = \mathcal{N}(x_t|\alpha_t x_0, \sigma_t^2 \boldsymbol{I}) , \qquad (1)$$

where $\alpha_t, \sigma_t \in \mathbb{R}^+$ are time-dependent differentiable functions that are usually chosen to ensure a strictly decreasing signal-to-noise ratio (SNR) $\alpha_t^2/\sigma_t^2$ and the final distribution $q_T(x_T) \approx \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$ (Kingma et al., 2021). Diffusion models can generate new samples starting from the prior distribution, after learning to reverse the forward process. Such the reverse-time denoising process from time $T$ to time 0 can be described by a stochastic differential equation (SDE) (Yang et al., 2022) as

$$dx_t = [f(t)x_t - g^2(t)\nabla_x \log p_t(x_t)]d_t + g(t)d\bar{\boldsymbol{w}}_t , \qquad (2)$$

where $\nabla_x \log p_t(x_t)$ is the so-called score function and $\bar{\boldsymbol{w}}_t$ is the standard reverse-time Wiener process. While $f(t) = \frac{d \log \alpha_t}{dt}$ is the drift coefficient of SDEs, $g^2(t) = \frac{d\sigma_t^2}{dt} - 2\frac{d \log \alpha_t}{dt}\sigma_t^2$ is the diffusion coefficient (Kingma et al., 2021). Deep neural networks are used to parameterize the score function variants in two similar forms, *i.e.*, the noise prediction model $\boldsymbol{\epsilon_\theta}(x_t, t)$ and the data prediction model $\boldsymbol{d_\theta}(x_t, t)$. In this study, we focus on the parameterization of the widely used data prediction model to directly predict the original data $x_0$ from $x_t$.

### 2.1.2. RNA Inverse Folding

Inverse folding aims to explore sequences that can fold into a predefined structure, which is specified here as the fixed sugar-phosphate backbone of an RNA tertiary structure. For an RNA molecule with $N$ nucleotides consisting of four different types A (Adenine), U (Uracil), C (Cytosine), and G (Guanine), its sequence can be defined as $\boldsymbol{S} \in \{A, U, C, G\}^N$. Among the backbone atoms, we choose one three-atom coarse-grained representation including the atom coordinates of C4', C1', N1 (pyrimidine) or N9 (purine) for every nucleotide. The simplified backbone structure can be denoted as $\boldsymbol{X} \in \mathbb{R}^{3N \times 3}$. Note that there are various alternative schemes for coarse-graining RNA 3D backbones, including using more atoms to obtain precise representations (Dawson et al., 2016). We explore a concise representation with regular structural patterns (Shen et al., 2022).

Formally, we consider the RNA inverse folding problem as modeling the conditional distribution $p(\boldsymbol{S}|\boldsymbol{X})$, *i.e.*, the sequence distribution conditioned on RNA backbone structures. We establish a diffusion model to learn the conditional sequence distribution. To take advantage of the convenience of defining diffusion models in continuous data spaces (Chen et al., 2023; Dieleman et al., 2022), discrete nucleotide types in the sequence are represented by one-hot encoding and continuousized in the real number space as $\boldsymbol{S} \in \mathbb{R}^{4N}$. The continuous-time forward diffusion process in the sequence space $\mathbb{R}^{4N}$ can be described by the forward
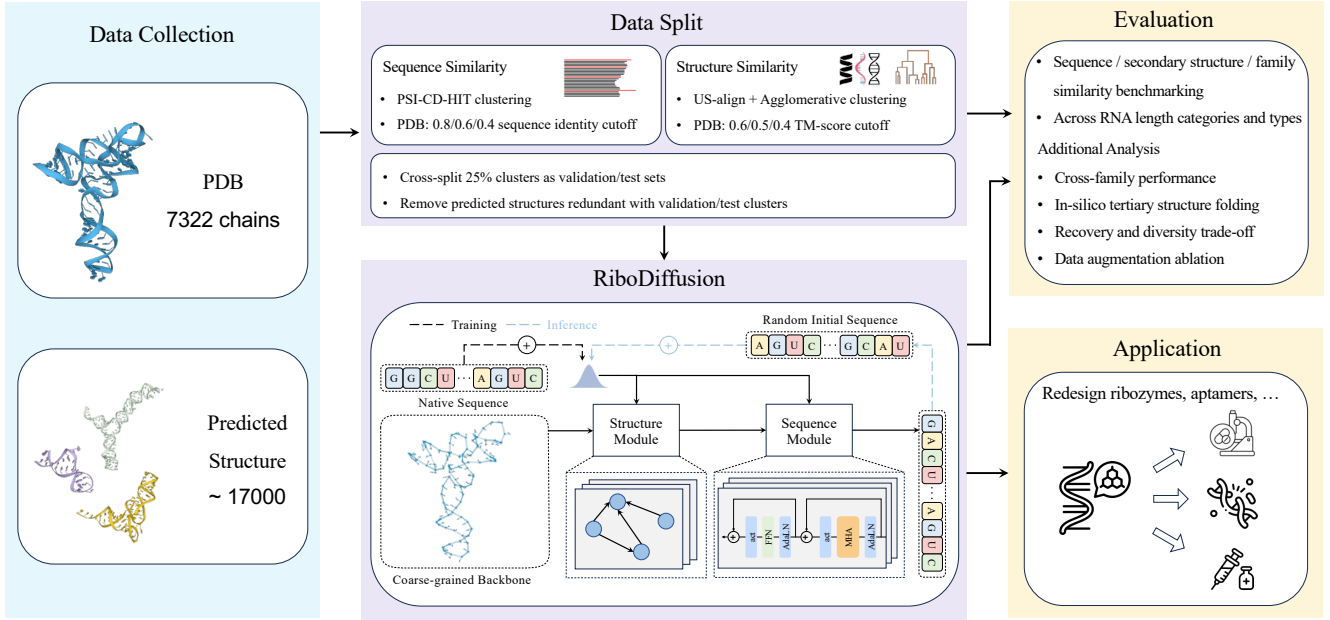
Fig. 1: Overview of RiboDiffusion for tertiary structure-based RNA inverse folding. We construct a dataset with experimentally determined RNA structures from PDB, supplemented with additional structures predicted by an RNA structure prediction model. We cluster RNA with different cut-offs for sequence or structure similarity and make cross-split to evaluate models. RiboDiffusion trains a neural network with a structure module and a sequence module to recover the original sequence from a noisy sequence and a coarse-grained RNA backbone extracted from the tertiary structure. RiboDiffusion then uses the trained network to iteratively refine random initial sequences until they match the target structure. We present a comprehensive evaluation and analysis of the proposed method.

SDE with $t \in [0, T]$ as $\mathrm{d}\boldsymbol{S}_t = f(t)\boldsymbol{S}_t\mathrm{d}t + g(t)\mathrm{d}\boldsymbol{w}$. Under this forward SDE, the original sequence at time $t = 0$ is gradually corrupted by adding Gaussian noise. With the linear Gaussian transition kernel derived from the forward SDE in Eq. (1) (Yang et al., 2022), we can conveniently sample $\boldsymbol{S}_t = \alpha_t + \sigma_t\epsilon_{\boldsymbol{S}}$ at any time $t$ for training, where $\epsilon_{\boldsymbol{S}}$ is Gaussian noise in the sequence space. For the generative denoising process, the corresponding reverse-time SDE from time T to 0 can be derived from Eq. (2) as

$$\mathrm{d}\boldsymbol{S}_t = [f(t) - g^2(t)\nabla_{\boldsymbol{S}} \log p_t(\boldsymbol{S}_t|\boldsymbol{X})]\mathrm{d}t + g(t)\mathrm{d}(\bar{\boldsymbol{w}}_t) , \quad (3)$$

where $p_t(\boldsymbol{S}_t|\boldsymbol{X})$ is the marginal distribution of sequences given $\boldsymbol{X}$, and the score function $\nabla_{\boldsymbol{S}} \log p_t(\boldsymbol{S}_t|\boldsymbol{X})$ represents the gradient field of the logarithmic marginal distribution.

Once the score function is parameterized, we can numerically solve this reverse SDE to convert random samples from the prior distribution $\mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$ into the desired sequences. We establish a data prediction model to achieve the score function parameterization, learning to reverse the forward diffusion process. Specifically, we feed the noised sequence data $\boldsymbol{S}_t$, the log signal-to-noise ratio $\lambda_t = \log(\alpha_t^2/\sigma_t^2)$, and the conditioning RNA backbone structures $\boldsymbol{X}$ to the data prediction model $\boldsymbol{d}_{\boldsymbol{\theta}}(\boldsymbol{S}_t, \lambda_t, \boldsymbol{X})$. We optimize the data prediction model with a simple weighted squared error objective function:

$$\min_{\boldsymbol{\theta}} \mathbb{E}_t\{\sqrt{\frac{\alpha_t}{\sigma_t}}\mathbb{E}_{\boldsymbol{S}_0, \boldsymbol{X}}\mathbb{E}_{\boldsymbol{S}_t|\boldsymbol{S}_0}||\boldsymbol{d}_{\boldsymbol{\theta}}(\boldsymbol{S}_t, \lambda_t, \boldsymbol{X}) - \boldsymbol{S}_0||_2^2\} , \quad (4)$$

which can be considered as optimizing a weighted variational lower bound on the data log-likelihood or a form of denoising score matching (Ho et al., 2020; Song et al., 2021; Kingma et al., 2021).

## 2.2. Model Architecture

The architecture design of the data prediction model largely determines the diffusion learning quality of the diffusion model. We propose a two-module model to predict the original nucleotide types: a structure module to capture geometric features and a sequence module to capture intra-sequential correlation.

### 2.2.1. Structure Module

Geometric deep learning models aim to extract equivariant or invariant features from 3D data and achieve impressive performance in the protein inverse folding task (Ingraham et al., 2019; Jing et al., 2021; Gao et al., 2023). Our structure module is constructed based on the GVP-GNN architecture (Jing et al., 2021) and adapted for RNA backbone structures.

The fixed RNA backbone is first represented as a geometric graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where each node $v_i \in \mathcal{V}$ corresponds to a nucleotide and connects to its top-$k$ nearest neighbors according to the distance of C1' atoms. The scalar and vector features are extracted from 3D coordinates as node and edge attributes in graphs, which describe the local geometry of nucleotides and their relative geometry. Specifically, the scalar node features in nucleotides are obtained from dihedral angles, while the vector node features consist of forward and reverse vectors of sequential C1' atoms, as well as the local orientation vectors of C1' to C4' and N1/N9. The initial embedding of each edge consists of its connected C1' atom's direction vector, Gaussian radial basis encoding for their Euclidean distance, and sinusoidal position encoding (Vaswani et al., 2017) of the relative distance in the sequence. In addition to geometry information, we also append the corrupted one-hot encoding of nucleotide types $\boldsymbol{S}_t$ as the

---

**Algorithm 1** RiboDiffusion Training.

---

1: $t \sim \mathcal{U}(0,1]$, $\boldsymbol{S}_0, \boldsymbol{X} \sim$ Training Set
2: $\boldsymbol{S}_t \sim \mathcal{N}(\boldsymbol{S}_t | \alpha_t \boldsymbol{S}_0, \sigma_t^2 \boldsymbol{I})$, $\lambda_t = \log(\alpha_t^2/\sigma_t^2)$, $\tilde{\boldsymbol{S}}_0 \leftarrow \boldsymbol{0}$
3: **if** Uniform$(0, 1.0) < 0.5$ **then**      ▷ Self Conditioning
4:      $\tilde{\boldsymbol{S}}_0 \leftarrow \boldsymbol{d_\theta}([\boldsymbol{S}_t, \tilde{\boldsymbol{S}}_0], \lambda_t, \boldsymbol{X})$
5:      $\tilde{\boldsymbol{S}}_0 \leftarrow \text{StopGradient}(\tilde{\boldsymbol{S}}_0)$
6: **end if**
7: **if** Uniform$(0, 1.0) < 0.4$ **then**      ▷ Drop Structure Condition
8:      $\boldsymbol{X} \leftarrow \boldsymbol{0}$
9: **end if**
10: Minimize $\sqrt{\frac{\alpha_t}{\sigma_t}}\ [||\boldsymbol{d_\theta}([\boldsymbol{S}_t, \tilde{\boldsymbol{S}}_0], \lambda_t, \boldsymbol{X}) - \boldsymbol{S}_0||_2^2]$

---

---

**Algorithm 2** RNA inverse folding via RiboDiffusion.

---

**Require**: time schedule $\{t_i\}_{i=0}^M$, RNA backbone coordinates $\boldsymbol{X}$
1: $\tilde{\boldsymbol{S}}_0 \leftarrow \boldsymbol{0}$
2: $\boldsymbol{S}_{t_0} \leftarrow \boldsymbol{S}_T \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$
3: **for** $i \leftarrow 1$ to $M$ **do**
4:      $t \leftarrow t_{i-1}$, $s \leftarrow t_i$, $\lambda_t \leftarrow \log(\alpha_t^2/\sigma_t^2)$
5:      $\alpha_{t|s} \leftarrow \alpha_t/\alpha_s$, $\sigma_{t|s}^2 \leftarrow \sigma_t^2 - \alpha_{t|s}^2 \sigma_s^2$
6:      $\tilde{\boldsymbol{S}}_0 \leftarrow \boldsymbol{d_\theta}([\boldsymbol{S}_t, \tilde{\boldsymbol{S}}_0], \lambda_t, \boldsymbol{X})$
7:      $\bar{\boldsymbol{S}}_s \leftarrow \frac{\alpha_{t|s}\sigma_s^2}{\sigma_t^2}\boldsymbol{S}_t + \frac{\alpha_s \sigma_{t|s}^2}{\sigma_t^2}\tilde{\boldsymbol{S}}_0$
8:      $\boldsymbol{S}_{\boldsymbol{\epsilon}} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$
9:      $\boldsymbol{S}_s \leftarrow \bar{\boldsymbol{S}}_s + \frac{\sigma_{t|s}\sigma_s}{\sigma_t}\boldsymbol{S}_{\boldsymbol{\epsilon}}$
10: **end for**
11: **return** $\bar{\boldsymbol{S}}_{t_M}$

---

node scalar features. Furthermore, inspired by the widely used self-conditioning technique in diffusion models (Chen et al., 2023; Watson et al., 2023; Huang et al., 2023b), the previously predicted sequence output, denoted as $\tilde{\boldsymbol{S}}_0$, is also considered as node embeddings to enhance the utilization of model capacity. To update the node embeddings, the nucleotide graph employs a standard message passing technique (Gilmer et al., 2017). This involves combining the neighboring nodes and edges through GVP layers, where scalar and vector features interact via gating to create messages. The resulting messages are then transmitted across the graph to update scalar and vector node representations.

### 2.2.2. Sequence Module

The sequential correlation in RNA primary structures is crucial for inverse folding and to obtain high-quality RNA sequences even with imprecise 3D coordinates. This concept is applicable in the inverse folding of proteins (Hsu et al., 2022; Zheng et al., 2023). The sequence module takes in $f$-dimensional nucleotide-level embeddings $\mathbf{h}^0 \in \mathbb{R}^{N \times f}$ as tokens, which consists of SE(3)-invariant scalar node representations from the structure module and corrupted sequence data. During training, we randomly add self-conditioning sequence data similar to those of the structure module and drop structural features to model both the conditional and unconditional sequence distributions for further application.

Our sequence module architecture is modified from the Transformer block (Vaswani et al., 2017) to inject diffusion context, log-SNR $\lambda$, or other potential conditional features (*e.g.* RNA types) (Dhariwal and Nichol, 2021; Peebles and Xie, 2023). The context input $\mathbf{C}$ affects sequence tokens in the form of adaptive normalization and activation layers, which are denoted

as adaLN and act functions:

$$
\begin{aligned}
\text{adaLN}(\mathbf{h}, \mathbf{C}) &= (\mathbf{1} + \text{MLP}_1(\mathbf{C})) \cdot \text{LN}(\mathbf{h}) + \text{MLP}_2(\mathbf{C}), \\
\text{act}(\mathbf{h}, \mathbf{C}) &= \text{MLP}_3(\mathbf{C}) \cdot \mathbf{h},
\end{aligned}
\tag{5}
$$

where $\text{LN}(\cdot)$ is the layer normalization and $\text{MLP}(\cdot)$ is a multilayer perception to learn shift and scale parameters. The $l$-th Transformer block is defined as follows

$$
\begin{aligned}
\mathbf{m}^l &= \text{MHA}(\text{adaLN}(\mathbf{h}^l, \lambda_t))\ , \\
\mathbf{h}^{l+1'} &= \text{act}(\mathbf{m}^l, \lambda_t) + \mathbf{h}^l, \\
\mathbf{h}^{l+1} &= \text{act}(\text{FFN}(\text{adaLN}(\mathbf{h}^{l+1'}, \lambda_t)), \lambda_t) + \mathbf{h}^{l+1'},
\end{aligned}
\tag{6}
$$

where $\text{MHA}(\cdot)$ is the multi-head attention layer and $\text{FFN}(\cdot)$ is the Feedforward neural network (Vaswani et al., 2017). Finally, the sequence module output $\mathbf{h}^L$ is projected to nucleotide one-hot encodings via an extra MLP. The detailed training procedure is referred to as Algorithm 1.

### 2.3. Sequence Sampling

To generate RNA sequences that are likely to fold into the given backbone, we construct a generative denoising process based on the parameterized reverse-time SDE with the optimized data prediction model $\boldsymbol{d_\theta}$, as described in Eq. (3). Various numerical solvers for the SDE can be employed for sampling, such as ancestral sampling, the Euler-Maruyama method, etc. We apply convenient ancestral sampling combined with the data prediction model and self-conditioning to generate sequences. Algorithm 2 outlines the specific sampling procedure. For more details on the noise schedule parameters, including $\alpha_t$ and $\sigma_t$, refer to (Kingma et al., 2021). We intuitively explain the denoising process as follows: we start by sampling noisy data from a Gaussian distribution that represents a random nucleotide sequence, and we iteratively transform this data towards the desired candidates under the condition of the given RNA 3D backbones.

Exploring novel RNA sequences that fold into well-defined 3D conformations distinct from the natural sequence is also an essential goal for RNA design, as it has the potential to introduce new functional sequences. This task not only requires the model to generate sequences that satisfy folding constraints but also to increase diversity for subsequent screening. During the generative denoising process, our model can balance the proportion of unconditional and conditional sequence distributions by adjusting the output of the data prediction model. Let $w$ be the conditional scaling weight, and the data prediction model can be modified as

$$
\tilde{\boldsymbol{d}}_{\boldsymbol{\theta}}(\boldsymbol{S}_t, \lambda_t, \boldsymbol{X}) = w\boldsymbol{d_\theta}(\boldsymbol{S}_t, \lambda_t, \boldsymbol{X}) + (1-w)\boldsymbol{d_\theta}(\boldsymbol{S}_t, \lambda_t, \boldsymbol{0}).
\tag{7}
$$

Setting $w = 1$ is the original conditional data prediction model while decreasing $w < 1$ weakens the effect of conditional information and strengthens the sequence diversity. In this way, we achieve a trade-off between recovering the original sequence and ensuring diversity. The distribution weighting technique is also used in diffusion models for text-to-image generation (Ho and Salimans, 2022; Saharia et al., 2022).

## 3. Results

We comprehensively evaluate and analyze RiboDiffusion for tertiary structure-based RNA inverse folding. Additional results can be found in supplemental materials. The source code is provided at https://github.com/ml4bio/RiboDiffusion.

**Table 1.** Recovery rate (%) comparison across six different settings. The average and standard deviation values of model performance on four random-split non-overlapping test sets are reported. Mean recovery rates are reported for short (L<=50nt), medium (50nt<L<=100nt), and long (L>100nt) RNA. *Seq. 0.8*: sequence similarity-based split with 0.8 cluster threshold. *Struct. 0.6*: structure similarity-based split with 0.6 cluster threshold.

| Methods | Seq 0.8 | | | | | Struct. 0.6 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Mean | Median | Short | Medium | Long | Mean | Median | Short | Medium | Long |
| RNAinverse | $25.92 \pm 1.1$ | $25.37 \pm 1.0$ | $25.99 \pm 2.0$ | $24.98 \pm 0.8$ | $27.54 \pm 1.4$ | $24.94 \pm 0.6$ | $24.24 \pm 0.5$ | $24.68 \pm 0.7$ | $24.98 \pm 1.0$ | $26.15 \pm 0.9$ |
| MCTS-RNA | $25.75 \pm 0.3$ | $25.61 \pm 0.1$ | $25.37 \pm 0.4$ | $26.15 \pm 0.5$ | $25.86 \pm 0.2$ | $25.81 \pm 0.5$ | $25.55 \pm 0.6$ | $25.38 \pm 0.5$ | $26.19 \pm 0.6$ | $25.86 \pm 0.9$ |
| LEARNA | $24.80 \pm 0.2$ | $24.55 \pm 0.3$ | $24.81 \pm 0.4$ | $24.86 \pm 0.2$ | $24.41 \pm 1.0$ | $24.96 \pm 0.2$ | $24.43 \pm 0.4$ | $24.88 \pm 0.5$ | $25.15 \pm 0.5$ | $24.36 \pm 0.6$ |
| MetaLEARNA | $29.10 \pm 0.6$ | $29.09 \pm 0.5$ | $27.43 \pm 1.5$ | $29.46 \pm 0.7$ | $32.40 \pm 0.9$ | $27.83 \pm 2.8$ | $27.95 \pm 2.5$ | $25.53 \pm 1.8$ | $29.51 \pm 0.6$ | $30.75 \pm 4.5$ |
| gRNAde | $42.67 \pm 5.3$ | $43.03 \pm 6.0$ | $36.25 \pm 2.0$ | $44.86 \pm 4.9$ | $46.06 \pm 6.1$ | $43.46 \pm 2.2$ | $43.37 \pm 2.7$ | $38.01 \pm 1.4$ | $49.82 \pm 2.7$ | $41.24 \pm 3.1$ |
| PiFold | $50.03 \pm 4.7$ | $50.32 \pm 6.0$ | $41.34 \pm 3.3$ | $53.20 \pm 3.7$ | $54.75 \pm 5.9$ | $47.89 \pm 5.4$ | $48.76 \pm 6.6$ | $40.13 \pm 1.0$ | $54.95 \pm 5.3$ | $45.62 \pm 7.7$ |
| StructGNN | $51.29 \pm 5.9$ | $52.40 \pm 8.0$ | $42.74 \pm 2.5$ | $54.45 \pm 7.1$ | $54.44 \pm 7.2$ | $55.20 \pm 6.9$ | $54.94 \pm 8.6$ | $46.36 \pm 1.0$ | $63.86 \pm 8.5$ | $48.48 \pm 11.3$ |
| GVP-GNN | $51.66 \pm 4.9$ | $53.48 \pm 6.4$ | $42.70 \pm 2.4$ | $56.20 \pm 5.7$ | $53.30 \pm 5.7$ | $53.76 \pm 5.4$ | $54.02 \pm 5.9$ | $45.80 \pm 0.7$ | $62.28 \pm 7.5$ | $47.39 \pm 9.0$ |
| RiboDiffusion | $\mathbf{57.32} \pm 4.1$ | $\mathbf{58.79} \pm 4.9$ | $\mathbf{52.01} \pm 3.1$ | $\mathbf{59.95} \pm 3.4$ | $\mathbf{58.91} \pm 5.7$ | $\mathbf{66.50} \pm 5.3$ | $\mathbf{66.72} \pm 5.8$ | $\mathbf{61.51} \pm 1.4$ | $\mathbf{73.89} \pm 8.4$ | $\mathbf{57.98} \pm 7.8$ |
| **Methods** | **Seq 0.6** | | | | | **Struct 0.5** | | | | |
| | Mean | Median | Short | Medium | Long | Mean | Median | Short | Medium | Long |
| RNAinverse | $25.35 \pm 0.5$ | $24.30 \pm 0.6$ | $25.66 \pm 0.6$ | $25.48 \pm 1.8$ | $27.76 \pm 2.1$ | $25.82 \pm 0.6$ | $24.79 \pm 0.9$ | $25.38 \pm 1.1$ | $25.39 \pm 1.0$ | $27.69 \pm 1.6$ |
| MCTS-RNA | $25.81 \pm 0.2$ | $25.67 \pm 0.2$ | $25.29 \pm 0.7$ | $26.22 \pm 0.5$ | $26.29 \pm 0.6$ | $25.93 \pm 0.4$ | $25.47 \pm 0.4$ | $25.49 \pm 0.5$ | $26.28 \pm 0.7$ | $26.06 \pm 0.6$ |
| LEARNA | $24.93 \pm 0.1$ | $24.78 \pm 0.1$ | $24.92 \pm 0.2$ | $25.04 \pm 0.6$ | $24.34 \pm 1.0$ | $25.00 \pm 0.2$ | $24.42 \pm 0.6$ | $25.23 \pm 0.3$ | $24.64 \pm 0.5$ | $24.02 \pm 1.2$ |
| MetaLEARNA | $29.07 \pm 3.2$ | $29.89 \pm 3.0$ | $25.99 \pm 2.4$ | $29.81 \pm 0.4$ | $33.89 \pm 3.3$ | $28.13 \pm 3.5$ | $28.18 \pm 3.5$ | $25.81 \pm 2.0$ | $29.54 \pm 0.9$ | $30.87 \pm 3.5$ |
| gRNAde | $47.28 \pm 4.3$ | $49.59 \pm 5.4$ | $37.60 \pm 1.7$ | $48.66 \pm 8.9$ | $47.34 \pm 3.5$ | $43.36 \pm 6.4$ | $43.61 \pm 7.6$ | $36.82 \pm 1.5$ | $47.06 \pm 6.5$ | $41.74 \pm 9.0$ |
| PiFold | $46.74 \pm 2.9$ | $48.54 \pm 3.9$ | $37.11 \pm 1.6$ | $47.35 \pm 4.4$ | $51.32 \pm 5.0$ | $49.22 \pm 3.0$ | $50.06 \pm 3.8$ | $42.48 \pm 3.0$ | $53.51 \pm 3.6$ | $46.90 \pm 5.3$ |
| StructGNN | $54.23 \pm 4.6$ | $57.97 \pm 7.0$ | $41.49 \pm 1.7$ | $56.09 \pm 6.9$ | $53.32 \pm 11.0$ | $52.99 \pm 8.6$ | $51.81 \pm 10.7$ | $44.56 \pm 2.4$ | $59.33 \pm 8.3$ | $45.06 \pm 14.3$ |
| GVP-GNN | $54.27 \pm 3.9$ | $57.60 \pm 5.6$ | $42.54 \pm 1.9$ | $56.17 \pm 5.8$ | $54.20 \pm 9.2$ | $50.91 \pm 5.7$ | $50.37 \pm 6.9$ | $44.74 \pm 2.0$ | $56.51 \pm 7.1$ | $44.21 \pm 9.5$ |
| RiboDiffusion | $\mathbf{59.06} \pm 2.8$ | $\mathbf{61.84} \pm 4.2$ | $\mathbf{50.68} \pm 2.1$ | $\mathbf{59.66} \pm 4.0$ | $\mathbf{59.79} \pm 7.9$ | $\mathbf{60.48} \pm 6.6$ | $\mathbf{59.31} \pm 7.9$ | $\mathbf{55.40} \pm 3.8$ | $\mathbf{65.69} \pm 9.0$ | $\mathbf{51.14} \pm 10.5$ |
| **Methods** | **Seq 0.4** | | | | | **Struct 0.4** | | | | |
| | Mean | Median | Short | Medium | Long | Mean | Median | Short | Medium | Long |
| RNAinverse | $25.53 \pm 0.7$ | $24.79 \pm 1.0$ | $25.29 \pm 0.4$ | $26.18 \pm 1.7$ | $27.27 \pm 1.8$ | $25.54 \pm 0.5$ | $24.47 \pm 0.6$ | $25.36 \pm 1.0$ | $24.94 \pm 0.7$ | $27.08 \pm 1.9$ |
| MCTS-RNA | $25.97 \pm 0.0$ | $25.86 \pm 0.2$ | $25.44 \pm 0.3$ | $26.48 \pm 0.3$ | $26.34 \pm 0.7$ | $25.81 \pm 0.4$ | $25.30 \pm 0.3$ | $25.27 \pm 0.5$ | $26.17 \pm 0.6$ | $25.86 \pm 1.0$ |
| LEARNA | $25.03 \pm 0.1$ | $24.55 \pm 0.3$ | $25.16 \pm 0.1$ | $24.84 \pm 0.4$ | $25.01 \pm 1.9$ | $25.05 \pm 0.1$ | $24.62 \pm 0.5$ | $25.21 \pm 0.2$ | $24.70 \pm 0.6$ | $24.02 \pm 1.2$ |
| MetaLEARNA | $28.94 \pm 1.1$ | $29.54 \pm 2.3$ | $25.83 \pm 2.2$ | $29.94 \pm 0.5$ | $35.36 \pm 2.8$ | $28.14 \pm 3.3$ | $28.31 \pm 3.2$ | $25.84 \pm 1.9$ | $29.45 \pm 0.6$ | $30.01 \pm 4.2$ |
| gRNAde | $43.58 \pm 7.6$ | $45.41 \pm 10.0$ | $36.02 \pm 2.6$ | $43.91 \pm 2.0$ | $46.84 \pm 12.5$ | $44.00 \pm 5.7$ | $44.10 \pm 7.1$ | $37.24 \pm 1.3$ | $48.01 \pm 5.6$ | $41.74 \pm 8.9$ |
| PiFold | $47.41 \pm 5.0$ | $49.00 \pm 6.7$ | $37.64 \pm 1.8$ | $50.38 \pm 4.7$ | $52.11 \pm 9.8$ | $49.84 \pm 2.7$ | $50.61 \pm 3.8$ | $42.39 \pm 2.9$ | $54.33 \pm 3.5$ | $45.92 \pm 6.3$ |
| StructGNN | $50.40 \pm 6.7$ | $52.57 \pm 10.8$ | $41.03 \pm 1.5$ | $51.98 \pm 4.6$ | $53.33 \pm 13.9$ | $54.65 \pm 7.8$ | $53.98 \pm 9.7$ | $45.39 \pm 2.5$ | $61.35 \pm 7.0$ | $44.62 \pm 14.4$ |
| GVP-GNN | $50.55 \pm 4.7$ | $52.59 \pm 7.0$ | $41.77 \pm 0.9$ | $53.73 \pm 5.6$ | $51.48 \pm 9.3$ | $52.29 \pm 5.1$ | $51.84 \pm 6.6$ | $45.26 \pm 1.9$ | $58.26 \pm 5.9$ | $44.04 \pm 9.5$ |
| RiboDiffusion | $\mathbf{57.24} \pm 5.0$ | $\mathbf{59.94} \pm 7.7$ | $\mathbf{50.06} \pm 2.4$ | $\mathbf{58.33} \pm 4.5$ | $\mathbf{58.85} \pm 11.4$ | $\mathbf{62.13} \pm 6.0$ | $\mathbf{61.09} \pm 7.6$ | $\mathbf{56.48} \pm 3.9$ | $\mathbf{67.94} \pm 7.7$ | $\mathbf{50.36} \pm 10.9$ |

### 3.1. Dataset Construction

We gather a dataset of RNA tertiary structures from the PDB database for RNA inverse folding. The dataset contains individual RNA structures and single-stranded RNA structures extracted from complexes. After filtering based on sequence lengths ranging from 20 to 280, there is a total of 7.322 RNA tertiary structures and 2,527 unique sequences. In addition to experimentally determined data, we construct augment training data by predicting structures with RhoFold (Shen et al., 2022). The structures predicted from RNAcentral sequences (Sweeney et al., 2019) are filtered by pLDDT to keep only high-quality predictions, resulting in 17,000 structures.

To comprehensively evaluate models, we divide the structures determined by experiments into training, validation, and test sets based on sequence similarity and structure similarity with different clustering thresholds. We use PSI-CD-HIT (Fu et al., 2012) to cluster sequences based on nucleotide similarity. We set the threshold at 0.8/0.6/0.4 and obtain 1,252/1,157/1,114 clusters, respectively. For structure similarity clustering, we calculate the TM-score matrix using US-align (Zhang et al., 2022) and apply the agglomerative clustering algorithm from scipy (Virtanen et al., 2020) on the similarity matrix. We achieve 2,036/1,659/1,302 clusters with TM-score thresholds of 0.6/0.5/0.4. We randomly split the clusters into three groups: 15% for testing, 10% for validation, and the remaining for training. We perform 4 random splits with non-overlapping testing and validation sets for each split strategy to evaluate models. The augmented training data is also filtered strictly based on the similarity threshold with the validation and testing sets for each split.

### 3.2. RNA Inverse Folding Benchmarking

*Baselines.* We compare our model with four machine learning baselines with tertiary structure input, including **gRNAde** (Joshi et al., 2023), **PiFold** (Gao et al., 2023), **StructGNN** (Ingraham et al., 2019), **GVP-GNN** (Jing et al., 2021). While gRNAde is a concurrent graph-based RNA inverse folding method, PiFold, StructGNN, and GVP-GNN are representative deep-learning methods of protein inverse folding, which are modified here to be compatible with RNA. Implementation details of these model modifications are in the supplementary material. These methods use the same 3-atom RNA backbone representation. We also introduce RNA inverse folding methods with secondary structures as input for comparison. **RNAinverse** (Hofacker et al., 1994) is an energy-based local searching algorithm for secondary structure constraints. **MCTS-RNA** (Yang et al., 2017) searches candidates based on Monte Carlo tree search. **LEARNA** and **MetaLEARNA** are deep reinforcement learning approaches (Runge et al., 2019) to design RNA that folds into the given secondary structures. Each method generates a sequence for every RNA backbone for benchmarking.

*Metrics.* The recovery rate is a commonly used metric in inverse folding that shows how much of the sequence generated by the model matches the original native sequence. While similar sequences have a higher chance of achieving the correct fold, the recovery rate is not a direct measure of structural fitness. We further evaluate with two metrics: the F1 Score, which assesses the alignment between the generated sequence's predicted secondary structure (via RNAfold (Gruber et al., 2008)) and the secondary structure extracted from the input's tertiary structure (using
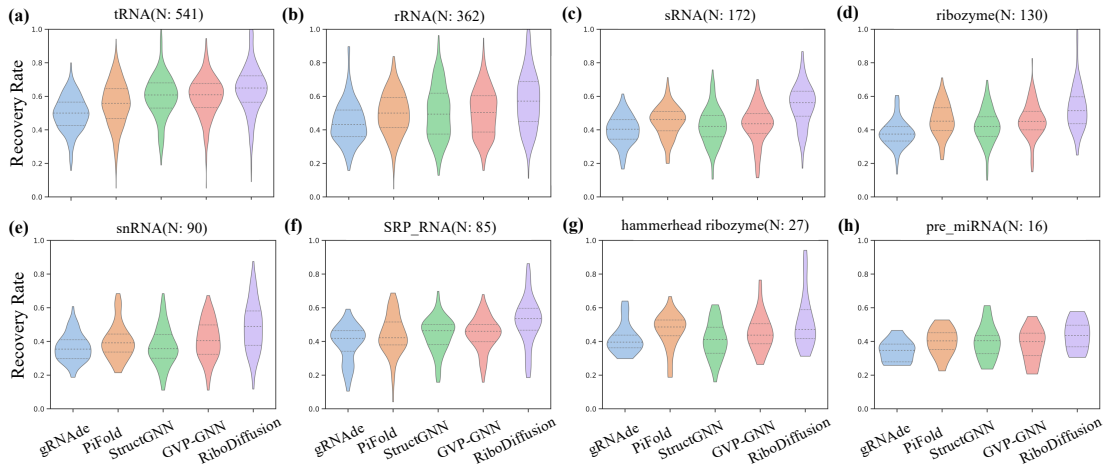
Fig. 2: Violin plots for the recovery rate distribution of methods for different types of RNA, including tRNA, rRNA, sRNA, ribozyme, snRNA, SRP RNA, hammerhead ribozyme, and pre miRNA.

**Table 2.** Comparison of secondary structure similarity and success rate of family preservation. F1: F1 score. Suc.: success rate of family preservation.

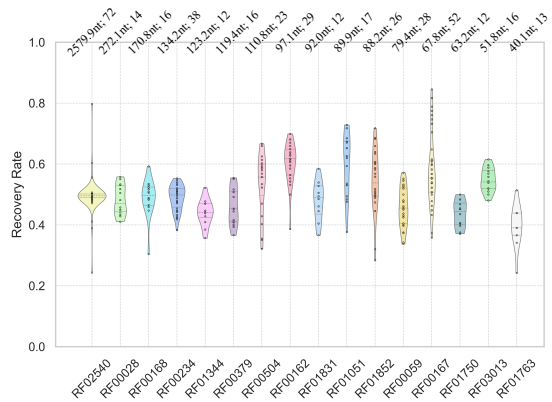| | | gRNAde | PiFold | StructGNN | GVP-GNN | RiboDiffusion |
|---|---|---|---|---|---|---|
| Seq 0.8 | F1 | 0.564 | 0.408 | 0.761 | 0.765 | 0.744 |
| | Suc. | 0.035 | 0.100 | 0.266 | 0.268 | 0.370 |
| Seq 0.6 | F1 | 0.142 | 0.336 | 0.709 | 0.740 | 0.749 |
| | Suc. | 0.018 | 0.031 | 0.217 | 0.186 | 0.316 |
| Seq 0.4 | F1 | 0.424 | 0.388 | 0.777 | 0.802 | 0.785 |
| | Suc. | 0.033 | 0.033 | 0.164 | 0.138 | 0.224 |
| Str 0.6 | F1 | 0.571 | 0.434 | 0.774 | 0.785 | 0.856 |
| | Suc. | 0.036 | 0.023 | 0.206 | 0.163 | 0.305 |
| Str 0.5 | F1 | 0.731 | 0.440 | 0.763 | 0.766 | 0.786 |
| | Suc. | 0.064 | 0.028 | 0.140 | 0.150 | 0.195 |
| Str 0.4 | F1 | 0.738 | 0.428 | 0.744 | 0.761 | 0.790 |
| | Suc. | 0.060 | 0.031 | 0.128 | 0.134 | 0.l77 |



Fig. 3: Performance of RiboDiffusion on different RNA families under the cross-family setting. The average length and number of tertiary structures for each family are marked above violin plots.

DSSR (Lu et al., 2015)), and the success rate determined by Rfam's covariance model (Kalvari et al., 2021), which evaluates the preservation of family-specific information in the generated sequences, indicating conserved structures and functions. Average success rates across families are reported.

We present recovery rate results in Table 1, which contains the average and standard deviation of four non-overlapping test sets for each model in different cluster settings. Our model outperforms the second best method by 11% on average for sequence similarity splits and 16% for structure similarity splits. RiboDiffusion consistently achieves better recovery rates in RNA with varying degrees of sequence or structural differences from training data. Methods based on tertiary structures outperform those based on secondary structures, as the latter contains less structural information. Extra results are shown in Table 2. It is worth noting that the tools used in these two metrics may contain errors. Our proposed method outperforms or matches the baseline methods in secondary structure alignments and more effectively retains family information from the input RNA.

We further classify the RNA in the test set based on its length and type to compare the model performance differences more thoroughly. First, we divide RNA into three categories based on the number of nucleotides (nt), *i.e.*, Short (50 nt or less), Medium

(more than 50 nt but less than 100 nt), and Long (100 nt or more). It can be observed in Table 1 that RiboDiffusion maintains performance advantages across different lengths of RNA. Short RNAs present a challenge for the model to recover the original sequence due to their flexible conformation, causing a relatively low recovery rate when compared to medium-length RNAs. A more detailed correlation of RiboDiffusion performance with RNA length is shown in supplemental materials. Each split shows similar patterns: RiboDiffusion has higher variance in short RNA inverse folding, and the model's performance becomes limited as RNA length increases. Moreover, Fig. 2 shows the recovery rate distribution of different RNA types with over 10 structures in test sets, including rRNA, tRNA, sRNA, ribozymes, etc. The RNA type information is collected from (Sweeney et al., 2019). Compared to other baselines, RiboDiffusion still has a better recovery rate distribution across RNA types. Through comprehensive benchmarking, we have observed remarkable performance improvement in tertiary structure-based RNA inverse folding achieved by RiboDiffusion.
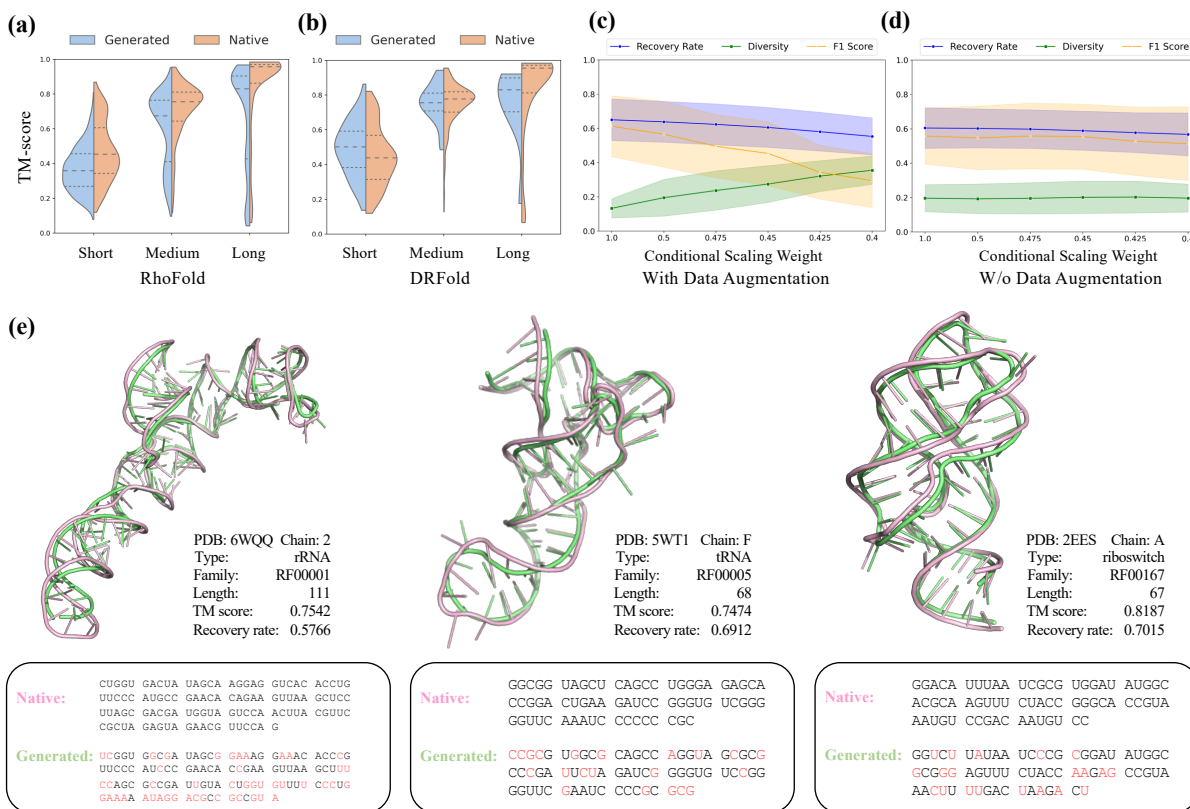
Fig. 4: **Analysis of RiboDiffusion. (a)-(b)** In-silico folding validation results that show the TM-score between structures predicted by RhoFold or DRFold and the given fixed RNA backbones (on *Seq. 0.4* split). *Native* represents structures predicted from original sequences of given backbones as references, while *Generated* represents structures predicted from generated sequences. **(c)-(d)** Trade-offs between the diversity of generated sequences and recovery rate, as well as refolding F1-score (including models with and without augmented data). **(e)** Visualization of input RNA structures (pink) and predicted structures (green) of generated sequences. The generated sequences and the corresponding native sequences are shown below the structure visualization, where different nucleotide types are marked in red.

## 3.3. Analysis of RiboDiffusion

We dive into a more comprehensive analysis of RiboDiffusion.

*Cross-family performance.* We repartition the dataset with the cross-family setting to further verify the generalization of our model. We obtain the RNA family corresponding to the tertiary structure from (Kalvari et al., 2021), then randomly select four families for testing and others for training. The experimental results of 4 non-overlapping splits are shown in Fig. 3. The average recovery rate of RiboDiffusion in each family generally ranges between 0.4 and 0.6. Especially, our model performs well on RF02540 whose sequence length far exceeds the training set. Although the performance is slightly worse than other splits in Table 1, these results still illustrate that our model can handle RNA families that do not appear in the training data, considering that cross-family is inherently a more difficult setting.

*In-silico tertiary structure folding validation.* To verify whether RiboDiffusion generated sequences can fold into a given RNA 3D backbone, we use computational methods to predict RNA structures (*i.e.*, RhoFold (Shen et al., 2022) and DRFold (Li et al., 2023)) to obtain their tertiary structures. Structure prediction models with a single sequence input are used due to the difficulty in finding homologous sequences for generated sequences and performing multiple sequence alignment. We take the TM-score

of C1' backbone atoms to measure the similarity between the predicted RNA structure of generated sequences and the given fixed backbones. Note that in-silico folding validation contains two sources of errors. One is the structure prediction error of the folding method itself, and the other is the sequence quality generated by RiboDiffusion. Therefore, we also predict the structure from the original native sequence using the same folding method and compare it to the given RNA backbone as an error and uncertainty reference.

As depicted in Fig. 4 (a), sequences generated by RiboDiffusion exhibit promising folding results in the fixed backbone for medium-length and long-length RNAs. However, the performance for short-length RNAs is relatively poor, which is affected by the unsatisfied recovery rate of our model and the limitations of RhoFold itself. We also show the folding performance using DRFold in Fig. 4 (b), where RiboDiffusion exhibits distribution shapes similar to those of using RhoFold. Here, due to the limitation of DRFold inference speed, we only test on the representative sequence of each cluster instead of the entire test set. We further make in-silico folding (with RhoFold) case studies of rRNA, tRNA, and riboswitch in Fig. 4 (e). RiboDiffusion generates new sequences that are different but still tend to fold into similar geometries. To alleviate concerns about the independence of structure prediction and inverse folding models, we provide results from alternative tools and evaluations

of structures independent of current datasets as an extra reference in the supplementary material.

*Trade-off between sequence recovery and diversity.* Exploring novel RNA sequences that have the potential to collapse into a fixed backbone distinct from native sequences is a realistic demand for RNA design. However, there is a trade-off between the diversity and recovery rate of the generated sequences. RiboDiffusion can achieve this balance by controlling the conditional scaling weight. For the representative input backbone of each cluster, we generate 8 sequences in total to report diversity. The diversity within the generated set of sequences $G$ is defined as $\text{IntDiv}(G) = 1 - \frac{1}{|G|^2} \sum_{S_1, S_2 \in G} \text{Sim}(S_1, S_2)$ (Benhenda, 2017). The function Sim compares two sequences by calculating the ratio of the length of the aligned subsequence to the length of the shorter sequence. In Fig. 4 (c), it is evident that the mean diversity of generated sequences in the test sets begins to increase when the conditional scaling weight is set to 0.5, while the recovery rate and the F1 score decrease to some extent. Therefore, we recommend using a value between 0.5 and 0.35 to adjust the sequence diversity.

*Training data augmentation analysis.* Augmenting training data is primarily driven by the scarcity and limited diversity of RNA available in PDB. Table 3 indicates that the incorporation of additional RhoFold predictions improves the overall generated sequence quality. This augmentation also enhances the adjustment ability of RiboDiffusion for sequence diversity, as shown in Fig. 4 (d), where the sequence diversity of the model without the augmented data remains relatively low. Notably, the noisy nature of augmented data requires appropriate preprocessing and filtering for quality assurance.

**Table 3.** Ablation study on data augmentation. Rec.: recovery rate.

|  | Rec. Mean | Rec. Median | F1 score | Rfam Success |
|---|---|---|---|---|
| RiboDiffusion | 57.24% | 59.94% | 0.785 | 0.224 |
| w/o Augment | 55.26% | 57.01% | 0.768 | 0.221 |

## 4. Conclusion

We propose RiboDiffusion, a generative diffusion model for RNA inverse folding based on tertiary structures. By benchmarking methods on sequence and structure similarity splits, comparing performance across RNA length and type, and validating with in-silico folding, we demonstrate the effectiveness of our model. Our model can also make trade-offs between recovery and diversity, and handle cross-family inverse folding. In future work, we aim to expand the scope of RiboDiffusion by exploring RNA sequences that span larger magnitudes in size and integrate contact information from the complex into the model. Our ultimate objective is to utilize the model for designing functional RNA like ribozymes, riboswitches, and aptamers, and to verify its effectiveness in wet lab experiments.

## References

M. Andronescu, A. P. Fejes, F. Hutter, H. H. Hoos, and A. Condon. A new algorithm for rna secondary structure design. *Journal of molecular biology*, 336(3):607–624, 2004.

M. Baek, R. McHugh, I. Anishchenko, D. Baker, and F. DiMaio. Accurate prediction of nucleic acid and protein-nucleic acid complexes using rosettafoldna. *bioRxiv*, pages 2022–09, 2022.

P. D. Bank. Protein data bank. *Nature New Biol*, 233:223, 1971.

M. Benhenda. Chemgan challenge for drug discovery: can ai reproduce natural chemical diversity? *arXiv preprint arXiv:1708.08227*, 2017.

A. Busch and R. Backofen. Info-rna—a fast approach to inverse rna folding. *Bioinformatics*, 22(15):1823–1831, 2006.

J. Chappell, K. E. Watters, M. K. Takahashi, and J. B. Lucks. A renaissance in rna synthetic biology: new mechanisms, applications and tools for the future. *Current opinion in chemical biology*, 28:47–56, 2015.

J. Chappell, A. Westbrook, M. Verosloff, and J. B. Lucks. Computational design of small transcription activating rnas for versatile and dynamic gene regulation. *Nature communications*, 8(1):1051, 2017.

T. Chen, R. ZHANG, and G. Hinton. Analog bits: Generating discrete data using diffusion models with self-conditioning. In *ICLR*, 2023.

A. Churkin, M. D. Retwitzer, V. Reinharz, Y. Ponty, J. Waldispühl, and D. Barash. Design of rnas: comparing programs for inverse rna folding. *Briefings in bioinformatics*, 19(2):350–358, 2018.

J. Dauparas, I. Anishchenko, N. Bennett, H. Bai, R. J. Ragotte, L. F. Milles, B. I. Wicky, A. Courbet, R. J. de Haas, N. Bethel, et al. Robust deep learning–based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.

W. K. Dawson, M. Maciejczyk, E. J. Jankowska, and J. M. Bujnicki. Coarse-grained modeling of rna 3d structure. *Methods*, 103:138–156, 2016.

P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*, 34:8780–8794, 2021.

S. Dieleman, L. Sartran, A. Roshannai, N. Savinov, Y. Ganin, P. H. Richemond, A. Doucet, R. Strudel, C. Dyer, C. Durkan, et al. Continuous diffusion for categorical data. *arXiv preprint arXiv:2211.15089*, 2022.

I. Dotu, J. A. Garcia-Martin, B. L. Slinger, V. Mechery, M. M. Meyer, and P. Clote. Complete rna inverse folding: computational design of functional hammerhead ribozymes. *Nucleic acids research*, 42(18):11752–11762, 2014.

A. Espah Borujeni, D. M. Mishler, J. Wang, W. Huso, and H. M. Salis. Automated physics-based design of synthetic riboswitches from diverse rna aptamers. *Nucleic acids research*, 44(1):1–13, 2016.

S. Findeiß, M. Etzel, S. Will, M. Mörl, and P. F. Stadler. Design of artificial riboswitches as biosensors. *Sensors*, 17(9):1990, 2017.

L. Fu, B. Niu, Z. Zhu, S. Wu, and W. Li. Cd-hit: accelerated for clustering the next-generation sequencing data. *Bioinformatics*, 28(23):3150–3152, 2012.

Z. Gao, C. Tan, and S. Z. Li. Pifold: Toward effective and efficient protein inverse folding. In *ICLR*, 2023.

J. A. Garcia-Martin, P. Clote, and I. Dotu. Rnaifold: a constraint programming algorithm for rna inverse folding and molecular design. *Journal of bioinformatics and computational biology*, 11(02):1350001, 2013.

J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. In *ICML*, pages 1263–1272, 2017.

A. R. Gruber, R. Lorenz, S. H. Bernhart, R. Neuböck, and I. L. Hofacker. The vienna rna websuite. *Nucleic acids research*, 36 (suppl_2):W70–W74, 2008.

J. Ho and T. Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.

J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020.

I. L. Hofacker, W. Fontana, P. F. Stadler, L. S. Bonhoeffer, M. Tacker, P. Schuster, et al. Fast folding and comparison of rna secondary structures. *Monatshefte fur chemie*, 125:167–167, 1994.

C. Hsu, R. Verkuil, J. Liu, Z. Lin, B. Hie, T. Sercu, A. Lerer, and A. Rives. Learning inverse folding from millions of predicted structures. In *ICML*, pages 8946–8970. PMLR, 2022.

H. Huang, L. Sun, B. Du, Y. Fu, and W. Lv. Graphgdp: Generative diffusion processes for permutation invariant graph generation. In *ICDM*, pages 201–210, 2022.

H. Huang, L. Sun, B. Du, and W. Lv. Conditional diffusion based on discrete graph structures for molecular graph generation. *AAAI*, 37(4):4302–4311, 2023a.

H. Huang, L. Sun, B. Du, and W. Lv. Learning joint 2d & 3d diffusion models for complete molecule generation. *arXiv preprint arXiv:2305.12347*, 2023b.

J. Ingraham, V. Garg, R. Barzilay, and T. Jaakkola. Generative models for graph-based protein design. *NeurIPS*, 32, 2019.

B. Jing, S. Eismann, P. Suriana, R. J. L. Townshend, and R. Dror. Learning from protein structure with geometric vector perceptrons. In *ICLR*, 2021.

C. K. Joshi, A. R. Jamasb, R. Viñas, C. Harris, S. Mathis, and P. Liò. Multi-state rna design with geometric multi-graph neural networks. *arXiv preprint arXiv:2305.14749*, 2023.

J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.

I. Kalvari, E. P. Nawrocki, N. Ontiveros-Palacios, J. Argasinska, K. Lamkiewicz, M. Marz, S. Griffiths-Jones, C. Toffano-Nioche, D. Gautheret, Z. Weinberg, et al. Rfam 14: expanded coverage of metagenomic, viral and microrna families. *Nucleic Acids Research*, 49(D1):D192–D200, 2021.

D. Kingma, T. Salimans, B. Poole, and J. Ho. Variational diffusion models. In *NeurIPS*, 2021.

R. Kleinkauf, T. Houwaart, R. Backofen, and M. Mann. antarna–multi-objective inverse folding of pseudoknot rna using ant-colony optimization. *BMC bioinformatics*, 16(1):1–7, 2015.

A. Laganà, D. Veneziano, F. Russo, A. Pulvirenti, R. Giugno, C. M. Croce, and A. Ferro. Computational design of artificial rna molecules for gene regulation. *RNA Bioinformatics*, pages 393–412, 2015.

Y. Li, C. Zhang, C. Feng, R. Pearce, P. Lydia Freddolino, and Y. Zhang. Integrating end-to-end learning with deep geometrical potentials for ab initio rna structure prediction. *Nature Communications*, 14(1):5745, 2023.

Z. Lin, H. Akin, R. Rao, B. Hie, Z. Zhu, W. Lu, N. Smetanin, R. Verkuil, O. Kabeli, Y. Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.

X.-J. Lu, H. J. Bussemaker, and W. K. Olson. Dssr: an integrated software tool for dissecting the spatial structure of rna. *Nucleic acids research*, 43(21):e142–e142, 2015.

M. McKeague, R. S. Wong, and C. D. Smolke. Opportunities in the design and application of rna for gene expression control. *Nucleic acids research*, 44(7):2987–2999, 2016.

S. V. Park, J.-S. Yang, H. Jo, B. Kang, S. S. Oh, and G. Y. Jung. Catalytic rna, ribozyme, and its applications in synthetic biology. *Biotechnology advances*, 37(8):107452, 2019.

W. Peebles and S. Xie. Scalable diffusion models with transformers. In *ICCV*, pages 4195–4205, 2023.

F. Runge, D. Stoll, S. Falkner, and F. Hutter. Learning to design rna. In *ICLR*, 2019.

C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 35:36479–36494, 2022.

T. Shen, Z. Hu, Z. Peng, J. Chen, P. Xiong, L. Hong, L. Zheng, Y. Wang, I. King, S. Wang, et al. E2efold-3d: end-to-end deep learning method for accurate de novo rna 3d structure prediction. *arXiv preprint arXiv:2207.01586*, 2022.

J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pages 2256–2265, 2015.

Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021.

B. A. Sweeney, A. I. Petrov, B. Burkov, R. D. Finn, A. Bateman, M. Szymanski, W. M. Karlowski, J. Gorodkin, S. E. Seemann, J. J. Cannone, et al. Rnacentral: a hub of information for non-coding rna sequences. *Nucleic Acids Research*, 47(D1):D221–D229, 2019.

A. Taneda. Modena: a multi-objective rna inverse folding. *Advances and Applications in Bioinformatics and Chemistry*, pages 1–12, 2010.

A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *NeurIPS*, 2017.

Q. Vicens and J. S. Kieft. Thoughts on how to think (and talk) about rna structure. *Proceedings of the National Academy of Sciences*, 119(17):e2112677119, 2022.

P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, et al. Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods*, 17(3):261–272, 2020.

J. L. Watson, D. Juergens, N. R. Bennett, B. L. Trippe, J. Yim, H. E. Eisenach, W. Ahern, A. J. Borst, R. J. Ragotte, L. F. Milles, et al. De novo design of protein structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023.

L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, Y. Shao, W. Zhang, B. Cui, and M.-H. Yang. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796*, 2022.

X. Yang, K. Yoshizoe, A. Taneda, and K. Tsuda. Rna inverse folding using monte carlo tree search. *BMC bioinformatics*, 18(1):1–12, 2017.

J. D. Yesselman and R. Das. Rna-redesign: a web server for fixed-backbone 3d design of rna. *Nucleic Acids Research*, 43(W1):W498–W501, 2015.

C. Zhang, M. Shine, A. M. Pyle, and Y. Zhang. Us-align: universal structure alignments of proteins, nucleic acids, and macromolecular complexes. *Nature methods*, 19(9):1109–1115, 2022.

Z. Zheng, Y. Deng, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. *bioRxiv*, pages 2023–02, 2023.

Y. Zhu, L. Zhu, X. Wang, and H. Jin. Rna-based therapeutics: an overview and prospectus. *Cell death & disease*, 13(7):644, 2022.

# RiboDiffusion: Tertiary Structure-based RNA Inverse Folding with Generative Diffusion Models

**Han Huang,[1,2,†] Ziqian Lin,[1,3,†] Dongchen He,[1] Liang Hong[1] and Yu Li[1,*]**

[1]Department of Computer Science and Engineering, CUHK, Hong Kong SAR, China, [2]School of Computer Science and Engineering, Beihang University, Beijing, China and [3]Nanjing University, Nanjing, China

## A. Dataset

Our dataset has $7,322$ experimentally determined RNA 3D structures. We provide the length histogram of these structures in Figure 1.

## B. Experiment Details and Results

### B.1. Secondary Structure-based Methods

We apply several RNA secondary structure-based inverse folding and protein inverse folding methods to compare model performance. For RNA secondary structure inverse folding methods, we extract secondary structures in dot-bracket form through DSSR (Lu et al., 2015). We use the default optimized hyperparameters of these methods. For LEARNA and MetaLEARNA (Runge et al., 2019), we set the design time limit to 600 seconds.

### B.2. Protein Inverse Folding Methods

GVP-GNN (Jing et al., 2021), PiFold (Gao et al., 2023) and StructGNN (Ingraham et al., 2019) are models based on graph neural networks which are first used for protein inverse folding methods. These methods can well extract geometric features using their graph neural network module. As a result, we construct RNA geometric features as model input.

For GVP-GNN, we use the same input features of RiboDiffusion as RiboDiffusion has a structure module based on GVP-GNN. The scalar node features contain dihedral angles of each nucleotide. The vector node features consist of forward and reverse vectors of sequential C1' atoms, as well as the local orientation vectors of C1' to C4' and N1/N9. The initial embedding of each edge consists of its connected C1' atom's direction vector, Gaussian radial basis encoding for their Euclidean distance, and sinusoidal position encoding of the relative distance in the sequence.

StructGNN consists of two parallel encoders to obtain embeddings of substructures and molecules, followed by a feed-forward neural network for prediction. PiFold contains PiGNN layers considering multi-scale residue interactions in node, edge, and global context levels of the graph and a linear layer. For PiFold and StructGNN, we construct distance, angle, and direction features for single or paired nucleotides similar to those in protein. The scalar node features contain dihedral angles of each nucleotide and Gaussian radial basis encoding for every atom pair among C4', C1', N1/N9 of each nucleotide. The vector node features consist of the local orientation vectors of C1' to C4' and N1/N9. The scalar edge features contain Gaussian radial basis encoding of every atom pair among C4', C1', N1/N9 of two different nucleotides, as well as quaternions of relative rotation between their local coordinate systems. The vector edge features consist of the orientation vectors of C1' of one nucleotide to C4' and N1/N9 in a different nucleotide. In these features, C4', C1', N1/N9 in nucleotide correspond to N, $C_\alpha$, C in protein residues.

Protein inverse folding methods exploit the geometric features of protein molecules. By constructing similar geometric features in our 3-atom RNA backbones, we can retrain these models on the RNA dataset. As a result, these methods can be applied to the RNA inverse folding problem.

### B.3. Metric

*Recovery rate.* This metric evaluates the quality of inverse folding from the perspective of sequence similarity. It is not perfect because it cannot directly characterize the possibility of sequences folding into a specified structure, but it still has a certain reference value. We plot the random mutation ratio versus the free energy of the sequence folding into a given secondary structure (extracted from the tertiary structure) in Figure 2. Folding into the structure is more likely when the recovery rate is relatively high. Moreover, our method has lower free energy than random mutation at the same recovery rate.

*F1 Score for secondary structure alignment.* F1 Score is defined between the secondary structure of the generated sequence predicted by RNAfold (Gruber et al., 2008) and the secondary structure extracted from the input tertiary structure. This metric reflects whether the generated sequence satisfies the folding constraints from the secondary structure level. However, since the secondary structures derived from both methods may have errors, we remove data that may have large errors based on the F1 score of the native sequence with a threshold of 0.7.

*Rfam success rate.* We use Rfam's covariance model to evaluate whether the sequence obtained by inverse folding maintains the same family information with the original RNA. Sequences within the curated family are generally considered to have conserved structures and similar functions. It is also of significance to discover such new sequences through inverse folding. Specifically, we define the success case as whether the bit score of the generated sequence is larger than the gathering threshold.

### B.4. Hyperparameters

Here we list the main hyperparameters we used in our model. We construct nucleotide graphs with top-10 neighbors and stack 4 layers for the graph neural network in the structure module, where the node feature dimension is 512 and the edge feature dimension
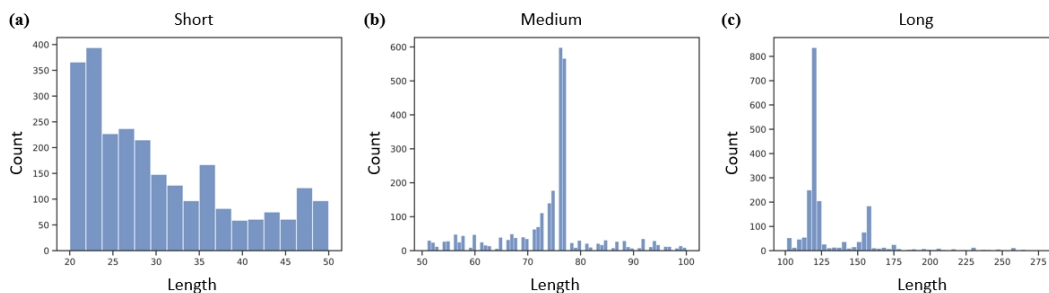
Fig. 1: **Length distribution of the experimentally determined structures. (a)-(c)** Length distribution of short ($L <= 50$nt), medium ($50$nt $< L <= 100$nt) and long ($L > 100$nt) RNA.
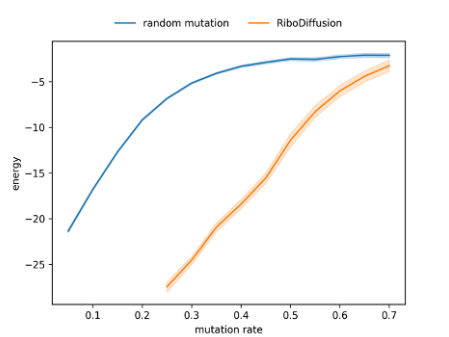


Fig. 2: **The correlation between different mutation rates and free energy** (with random mutation and RiboDiffusion).
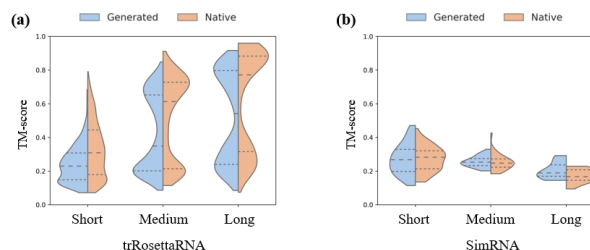


Fig. 3: **In-silico folding validation results of trRosettaRNA and simRNA.** In-silico folding validation results that show the TM-score between structures predicted by trRosettaRNA or simRNA and the given fixed RNA backbones (on *Seq. 0.4* split). *Native* represents structures predicted from original sequences of given backbones as references, while *Generated* represents structures predicted from generated sequences.

is 128. For the sequence module consisting of 8 blocks, we keep the 512 dimensions and use 8 attention heads. Our model is trained 40 epochs with the learning rate 0.0002.

## B.5. Extra In-silico Tertiary Structure Folding Results

To alleviate concerns about the independence of structure prediction tool and inverse folding models, we use two extra computational tools, trRosettaRNA (Wang et al., 2023) and SimRNA (Boniecki et al., 2016), to obtain tertiary structures of generated RNA sequences. We also use these tools to predict tertiary structures from the original native sequences. As depicted in Figure 3(a), generated and native sequences have similar TM-score distribution when predicted by trRosettaRNA. The result of SimRNA is shown in Figure 3(b). The performance of SimRNA is relatively poor, which indicates that although generated sequences have a similar TM-score distribution to natural sequences, the refolding evaluation based on SimRNA may have a large error and uncertainty.

Besides RhoFold (Shen et al., 2022), we also provide 3D visualized results of DRFold (Li et al., 2023) and trRosettaRNA, which are shown in Figure 6.

## B.6. Results on New RNA Structures

We evaluate the newly published RNA structures between 2023 and 2024 as an additional reference for our model. After removing redundancy and removing RNAs similar to the training set, we

present 8 structures that have not been trained by RiboDiffusion and RhoFold. The result is displayed in Table 3.

## B.7. Performance on CASP15

To assess the generalizability of the model, RiboDiffusion is tested on six natural RNAs in CASP15 without any overlap with the training set. As shown in Figure 4 (a) and (b), the performance of RiboDiffusion in complex RNA backbone structures is impressive, which is demonstrated by an average recovery rate of 0.56. Furthermore, the TM-score values of generated sequences are similar to the native sequences. However, it is important to note that the results of in-silico folding on CASP15 need more follow-up validation, as the TM-score value used as a reference is not satisfactory.

## B.8. Ablation Studies

We perform additional ablation studies to validate the necessity of the sequence module. We train the models in a sequence similarity split and a structure similarity split and report the results in Table 1. In our diffusion model formulation, adding the sequence module facilitates performance improvement.
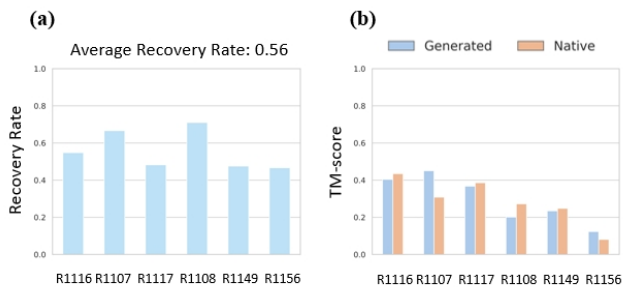
Fig. 4: **Performance on CASP15.(a)** A bar chart shows the recovery rate of RiboDiffusion on six natural RNAs in CASP15. **(b)** A bar chart displays the TM-score between predicted structures of RiboDiffusion-generated sequences and given RNA backbones. The TM-score of predicted structures from native sequences is displayed as a reference.

**Table 1.** Mean recovery rate (%) on ablation studies.

| Method | Seq. | Struct. |
| --- | --- | --- |
| RiboDiffusion | 58.96 | 66.40 |
| RiboDiffusion w/o seq | 57.82 | 64.26 |

### B.9. Running Time and Scalability Analysis.

The inference time of diffusion-based models is largely dependent on the number of steps in the sampling process. For the run-time analysis, we use 50 steps identical to those in our other experiments. On a GeForce RTX 3090 GPU, we report wall clock times of RiboDiffusion generation with different lengths of RNA and different numbers of sequences generated simultaneously in Figure 5. RiboDiffusion can finish the inverse folding of 200 nt RNA in just one second when generating a sequence. However, when generating 128 sequences simultaneously, RiboDiffusion experiences a significant increase in processing time, leading to limitations in scalability. We believe that the running speed of RiboDiffusion can be further improved in the future by accelerating the diffusion models, which is currently an emerging topic in machine learning.
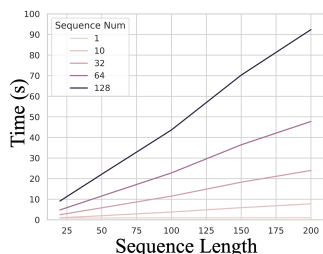


Fig. 5: **Running time and scalability analysis.** A line chart shows the relationship between running time and RNA sequence length when predicting different numbers of RNA sequences simultaneously.

### B.10. Extra Results on Secondary Structure Based Methods

We report extra results of secondary structure-based inverse folding methods in Table 2. These methods obtain high F1 scores because they directly use energy optimization to obtain sequences, making it unfair to compare with other methods. It is difficult for secondary structure-based inverse folding methods to generate new sequences in the same family due to the information loss compared to the tertiary structure input, even for tRNA with a more conservative shape.

**Table 2.** Comparison of secondary structure similarity and success rate of family preservation. The F1 score is an unfair metric for energy-optimized methods.

| | | rnainverse | MCTS | learna | metalearna |
| --- | --- | --- | --- | --- | --- |
| Seq 0.8 | F1* | 0.990 | 0.918 | 0.750 | 0.905 |
| | Suc. | 0.000 | 0.000 | 0.000 | 0.000 |
| Seq 0.6 | F1* | 0.991 | 0.922 | 0.764 | 0.916 |
| | Suc. | 0.000 | 0.000 | 0.000 | 0.000 |
| Seq 0.4 | F1* | 0.987 | 0.916 | 0.796 | 0.928 |
| | Suc. | 0.000 | 0.000 | 0.000 | 0.000 |
| Str 0.6 | F1* | 0.990 | 0.915 | 0.776 | 0.913 |
| | Suc. | 0.000 | 0.000 | 0.000 | 0.000 |
| Str 0.5 | F1* | 0.985 | 0.900 | 0.789 | 0.919 |
| | Suc. | 0.000 | 0.000 | 0.000 | 0.000 |
| Str 0.4 | F1* | 0.987 | 0.901 | 0.762 | 0.911 |
| | Suc. | 0.000 | 0.000 | 0.000 | 0.000 |

### B.11. Results on Remaining Dataset Splits

Extra results of different dataset splits are shown in Figure 7, 8, 9. We show the bivariate distribution of sequence length and recovery rate for RiboDiffusion on test set splits including *Seq. 0.6*, *Seq. 0.8*, *Struct. 0.5* and *Struct. 0.6* in Figure 7. We provide additional violin plots displaying the TM-score performance of RiboDiffusion-RhoFold pipeline about RNA length on test set splits including *Seq. 0.6*, *Seq. 0.8*, *Struct. 0.5* and *Struct. 0.6* in Figure 8. In Figure 9, four different types of RNA are tested to evaluate the performance of RiboDiffusion. The results show that RiboDiffusion performs better on tRNA compared to rRNA. However, its performance in sRNA and ribozyme may be limited due to the scale of the relevant training data.

### References

M. J. Boniecki, G. Lach, W. K. Dawson, K. Tomala, P. Lukasz, T. Soltysinski, K. M. Rother, and J. M. Bujnicki. Simrna: a coarse-grained method for rna folding simulations and 3d structure prediction. *Nucleic acids research*, 44(7):e63–e63, 2016.

Z. Gao, C. Tan, and S. Z. Li. Pifold: Toward effective and efficient protein inverse folding. In *ICLR*, 2023.

A. R. Gruber, R. Lorenz, S. H. Bernhart, R. Neuböck, and I. L. Hofacker. The vienna rna websuite. *Nucleic acids research*, 36 (suppl_2):W70–W74, 2008.

J. Ingraham, V. Garg, R. Barzilay, and T. Jaakkola. Generative models for graph-based protein design. *NeurIPS*, 32, 2019.

B. Jing, S. Eismann, P. Suriana, R. J. L. Townshend, and R. Dror. Learning from protein structure with geometric vector perceptrons. In *ICLR*, 2021.
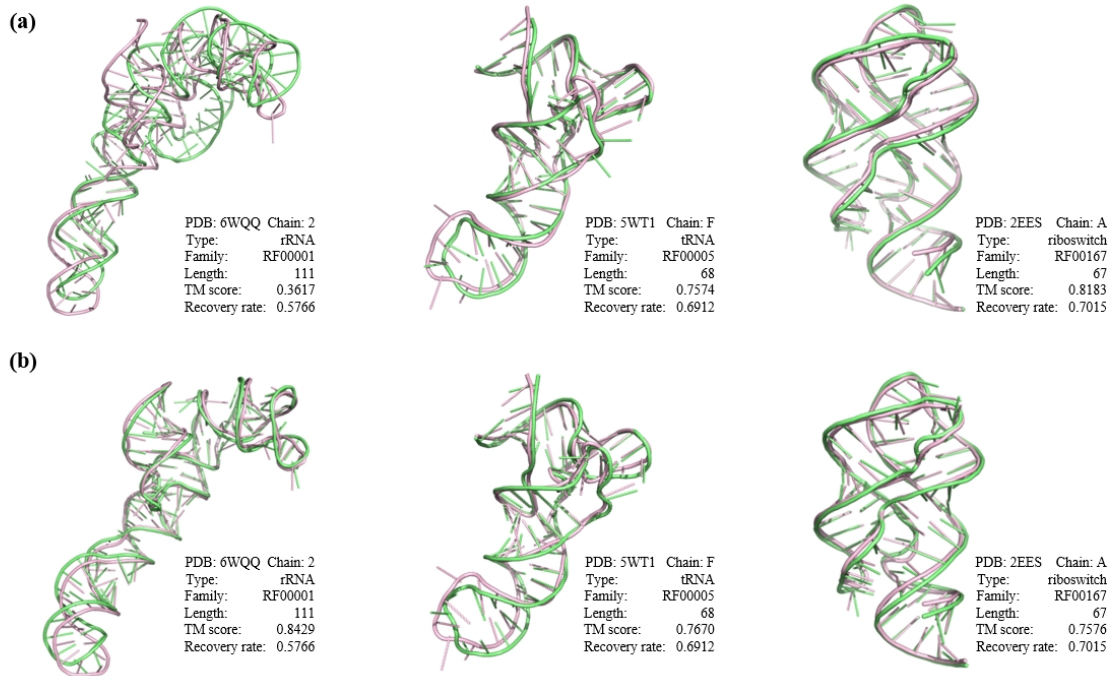
**(a)**



**(b)**



Fig. 6: **In-silico folding visualized results of DRFold and trRosettaRNA .(a)** Visualization of input RNA structures (pink) and predicted structures (green) of RiboDiffusion-DRFold pipeline. **(b)** Visualization of input RNA structures (pink) and predicted structures (green) of RiboDiffusion-trRosettaRNA pipeline.

**Table 3. Results on newly published RNA structures.** TM-score (generated) is calculated between the given structure and the refolded structure from the RiboDiffusion-RhoFold pipeline. TM-score (native) is calculated between the given structure and the predicted structure of RhoFold with the original native sequence.

| PDB id | Recovery rate | TM-score (generated) | TM-score (native) | Generated sequence |
|---|---|---|---|---|
| 7wii_V | 0.5918 | 0.4209 | 0.3573 | GGACCGUCCGCCAACAACGCUCCCCGAAAGGGGAGCAGCGGGAGGUCCA |
| 7xk1_B | 0.5556 | 0.7460 | 0.7842 | CGGAGGUGGCGCAGUGGUAGCGCAGGCGAGUUCAACUCGCCAGGCGCGGGUUCGAUUCCCGUCCUCCGGCCC |
| 8sh5_R | 0.5747 | 0.2632 | 0.3472 | GCGAAACUGGCAGAAUCGGUUAUGAGUUAGUCGAGCGAGACACGCUCACCCACCUUUUUAGGUUGGCUAACCGUUCGCUCGUUUUGA |
| 8t2a_R | 0.5889 | 0.3275 | 0.4005 | GGCUGCCGGAGUGCUUGUUGUCGUAGCCGGCAUGGAAAGACCAUGUGCUCGGCUACCCUUCGGGGUGUGAGCUACGGCACGACGGUGGUC |
| 8fn2_B | 0.6964 | 0.6011 | 0.3474 | GUCUGGUGGCCAUAGAAUCAAGGAACCACCUGAUCCCAUCCCGAACUCAGAAGUUAAGCUUGAUAUCGGUGAUGAUAUUGCGUUUUCGCGAGAAACUAGCGAACUGUCAGAA |
| 8gxb_B | 0.6667 | 0.1951 | 0.1607 | GAGCGUUGCUCGCAAGCGCCGCAUUGCACUUCGCGGCAGAGGUGUUAAUAAAAAGAAGCG |
| 8ine_5 | 0.7417 | 0.9157 | 0.9529 | GGGUACGGCCAUACUUCCCUGAAAACACCGAUUCCCCUCCGAUCAUCGAAGUUAAGCAGGGACAGGCUUGGUUAGUACUCGUGUCGGAGACGAACUGGGAACACCGAGUGCUGUACCCUU |
| 8ipy_8 | 0.5613 | 0.2557 | 0.5929 | CAAUUCUCGACUCAGAAUAUUUGGCUUCCUCUUCGUUGAAGAACGCAGCAAAAUGCGAUAAGCGAUAUGAGUUGCAAACAUAAAAGAGUAUUAGGGGUUCGAACGCAAAGGCGCUCCCAGUUGAAAUCUGGGAGUACAGCUCUUUCAGUCUCUUG |

Y. Li, C. Zhang, C. Feng, R. Pearce, P. Lydia Freddolino, and Y. Zhang. Integrating end-to-end learning with deep geometrical potentials for ab initio rna structure prediction. *Nature Communications*, 14(1):5745, 2023.

X.-J. Lu, H. J. Bussemaker, and W. K. Olson. Dssr: an integrated software tool for dissecting the spatial structure of rna. *Nucleic acids research*, 43(21):e142–e142, 2015.

F. Runge, D. Stoll, S. Falkner, and F. Hutter. Learning to design rna. In *ICLR*, 2019.
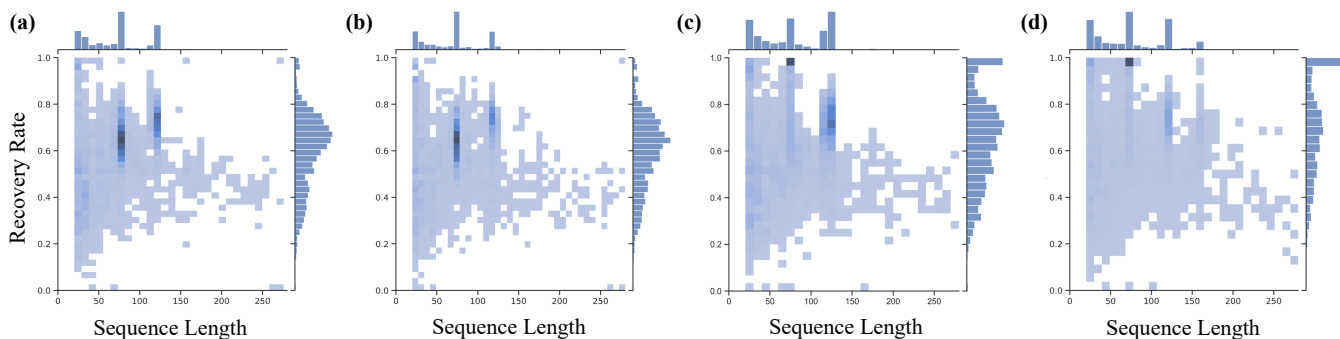
Fig. 7: **Bivariate distribution of sequence length and recovery rate for RiboDiffusion. (a)-(d)** Four joint histplots of the bivariate distribution between sequence length and recovery rate on test set splits including *Seq. 0.6*, *Seq. 0.8*, *Struct. 0.5* and *Struct. 0.6*.
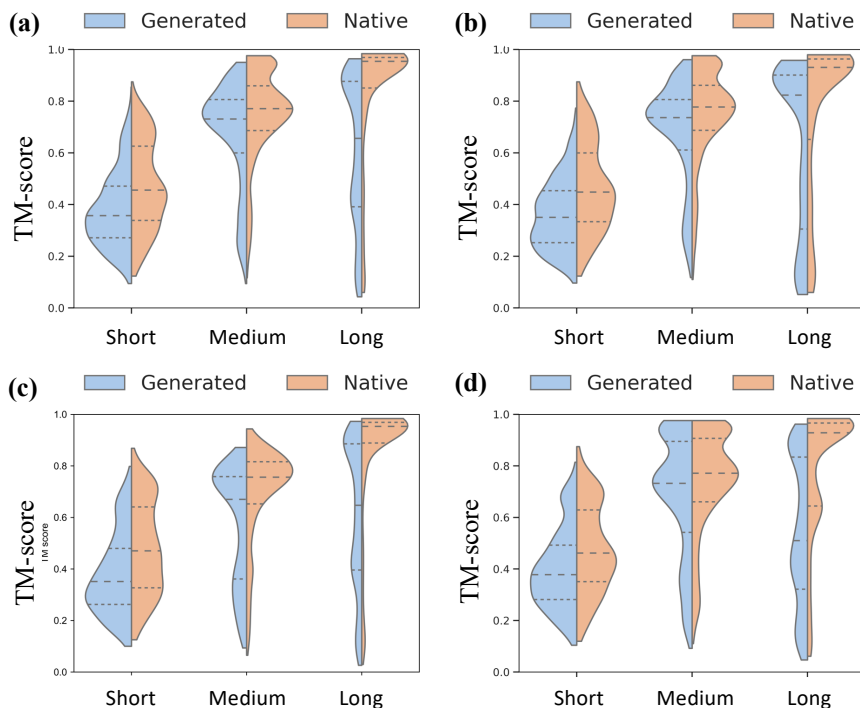


Fig. 8: **TM-score performance of RiboDiffusion-RhoFold pipeline about RNA length. (a)-(d)** Four violin plots compare the TM-score of structures predicted by RhoFold between generated RNA and native RNA on short, medium, and long RNA data in test set splits including *Seq. 0.6*, *Seq. 0.8*, *Struct. 0.5* and *Struct. 0.6*.

T. Shen, Z. Hu, Z. Peng, J. Chen, P. Xiong, L. Hong, L. Zheng, Y. Wang, I. King, S. Wang, et al. E2efold-3d: end-to-end deep learning method for accurate de novo rna 3d structure prediction. *arXiv preprint arXiv:2207.01586*, 2022.

W. Wang, C. Feng, R. Han, Z. Wang, L. Ye, Z. Du, H. Wei, F. Zhang, Z. Peng, and J. Yang. trrosettarna: automated prediction of rna 3d structure with transformer network. *Nature Communications*, 14(1):7266, 2023.
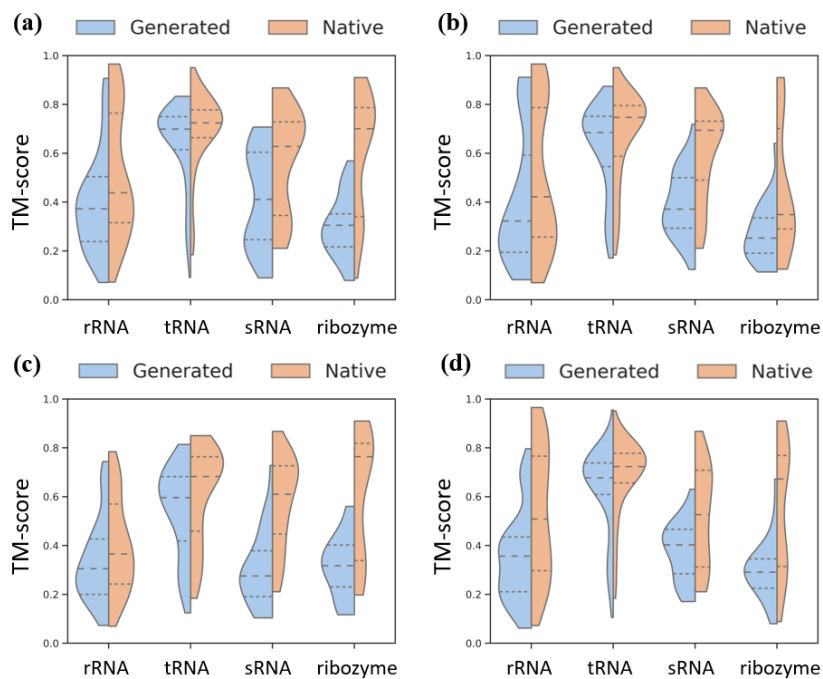
Fig. 9: **TM-score performance of RiboDiffusion-RhoFold pipeline about RNA type. (a)-(d)** Four violin plots compare the TM-score of structures predicted by RhoFold between generated RNA and native RNA on different types of RNA including rRNA, tRNA, sRNA, and ribozyme in test set splits including *Seq. 0.6*, *Seq. 0.8*, *Struct. 0.5* and *Struct. 0.6*.