



中國人民大學  
RENMIN UNIVERSITY OF CHINA



高瓴人工智能学院  
Gaoling School of Artificial Intelligence

# Deep Generative Design of RNA Family Sequences

Shunsuke Sumi, Michiaki Hamada, Hirohide Saito

**Nature Method**

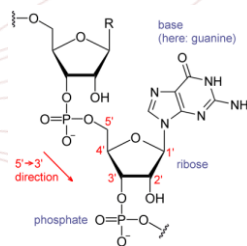
Fanmeng Wang

2025-6-13

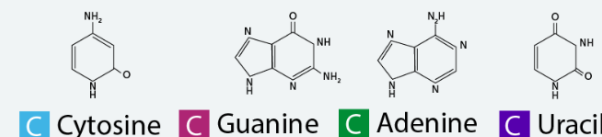
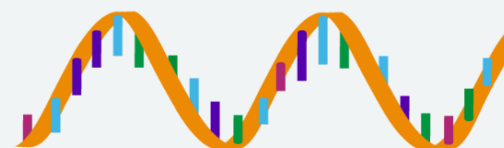
# Introduction

## ■ RNA Engineering

- Over the past decades, **RNA function exploitation** has driven the invention of various **synthetic molecular systems**, leading to a substantial impact on numerous fields such as basic research, biomanufacturing and medical applications.
- However, the **experimental search for desired RNA sequences** remains costly and inefficient due to the vast sequence space of RNA, which exponentially increases in complexity with sequence length.
- Therefore, **a versatile computational platform** to understand the sequence space and to efficiently design functional RNA is in great demand.



RNA  
(Ribonucleic acid)





# Introduction

## ■ Existing Methods

### ➤ RNA Inverse Folding:

- It aims to find sequences that fold into a **given secondary structure**, guided by RNA secondary structure prediction and a discrete optimization algorithm.
- However, the **functionality** of RNA is not characterized by structure alone. Besides, RNA inverse folding lacks **flexibility and generalizability**, and its **accuracy** is inherently limited by the accuracy of RNA secondary structure prediction and the optimization algorithm.



# Introduction

## ■ Existing Methods

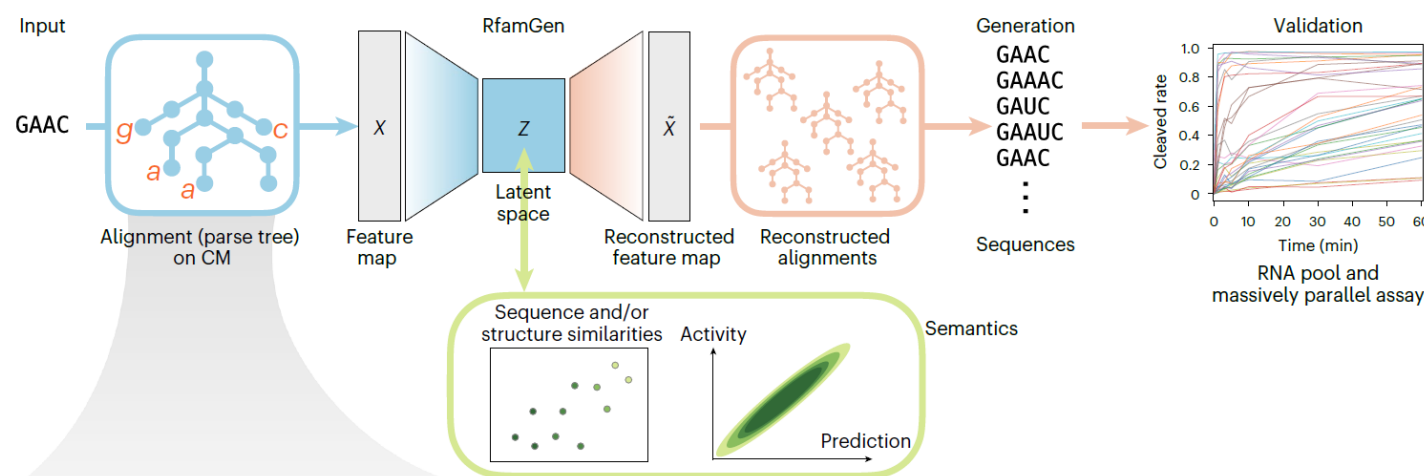
### ➤ Covariance Models (CM):

- It's a statistical framework for **RNA alignment and consensus secondary structure**, quantitatively evaluates variations of sequence and structure without relying on RNA secondary structure prediction.
- It has been the **gold standard** in RNA homology search for decades to categorize most functional RNA species into thousands of **RNA families**.
- However, it has not been used in functional RNA design previously.

# Introduction

## ■ RfamGen

- In this study, we propose the **RNA family sequence generator (RfamGen)**, which is a deep generative model for functional RNA design.
- In particular, RfamGen leverages **Covariance Models (CM) with Variational Autoencoders (VAE)** to generate artificial sequences of an RNA family, while also providing **semantically meaningful representations** of sequences.



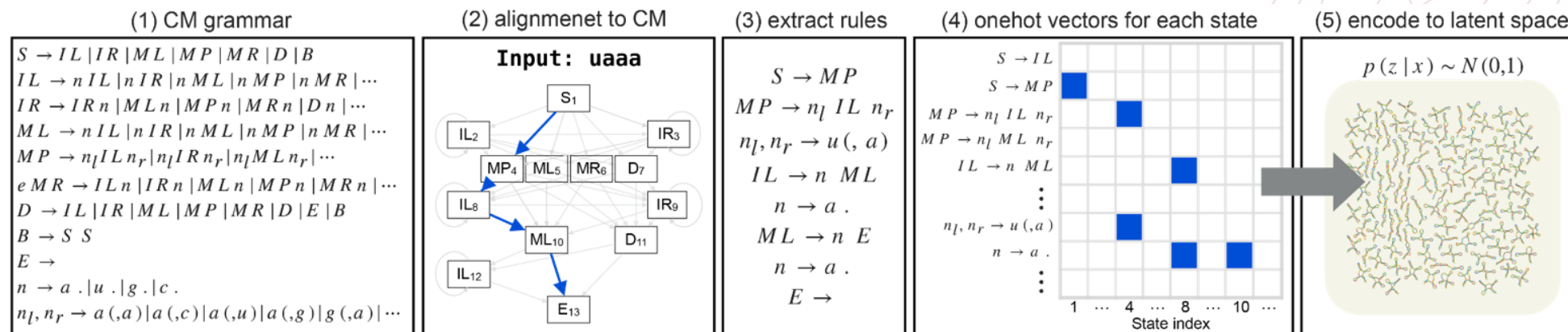




# Method

## ■ Data Preprocessing

1. Given a sequence, it can be aligned to a CM and the alignment is transformed to a **one-hot expression**.



CM grammar contains 76 rules: **56 rules of transition**, **4 rules of single strand emission** and **16 rules of base pair emission**.

$$\begin{aligned}
 S &\rightarrow IL | IR | ML | MP | MR | D | B \\
 IL &\rightarrow n_l' . 'IL' n_r' . 'IR' n_r' . 'ML' n_r' . 'MP' n_r' . 'MR' n_r' . 'D' n_r' . 'E' n_r' . 'B' \\
 IR &\rightarrow IR' . 'n' | ML' . 'n' | MP' . 'n' | MR' . 'n' | D' . 'n' | E' . 'n' | B' . 'n' \\
 ML &\rightarrow n_l' . 'IL' n_r' . 'IR' n_r' . 'ML' n_r' . 'MP' n_r' . 'MR' n_r' . 'D' n_r' . 'E' n_r' . 'B' \\
 MP &\rightarrow n_l' ( 'IL' ) n_r' n_l' ( 'IR' ) n_r' n_l' ( 'ML' ) n_r' n_l' ( 'MP' ) n_r' n_l' ( 'MR' ) n_r' n_l' ( 'D' ) n_r' n_l' ( 'E' ) n_r' n_l' ( 'B' ) n_r' \\
 MR &\rightarrow IL' . 'n' | IR' . 'n' | ML' . 'n' | MP' . 'n' | MR' . 'n' | D' . 'n' | E' . 'n' | B' . 'n' \\
 D &\rightarrow IL | IR | ML | MP | MR | D | E | B \\
 B &\rightarrow SS \\
 E &\rightarrow \\
 n &\rightarrow 'A' | 'U' | 'G' | 'C' \\
 n_l n_r &\rightarrow 'AA' | 'AC' | 'AG' | 'AU' | 'CA' | 'CC' | 'CG' | 'CU' | 'GA' | 'GC' | 'GG' | 'GU' | 'UA' | 'UC' | 'UG' | 'UU'
 \end{aligned}$$



# Method

## ■ Data Preprocessing

- Furthermore, to reduce the **observation bias** of the Rfam database, each sequence is **reweighted** by the number of sequences outside a neighborhood of a sequence as proposed in previous studies.

$$D_H(s, t) = \sum_{M \in \{\text{tr}, \text{ss}, \text{bp}\}} \text{Hamming Distance}(M_s, M_t) / \text{length}(M)$$

$$\pi_s = \left( \sum_t^N I[D_H(s, t) < \theta] \right)^{-1}$$

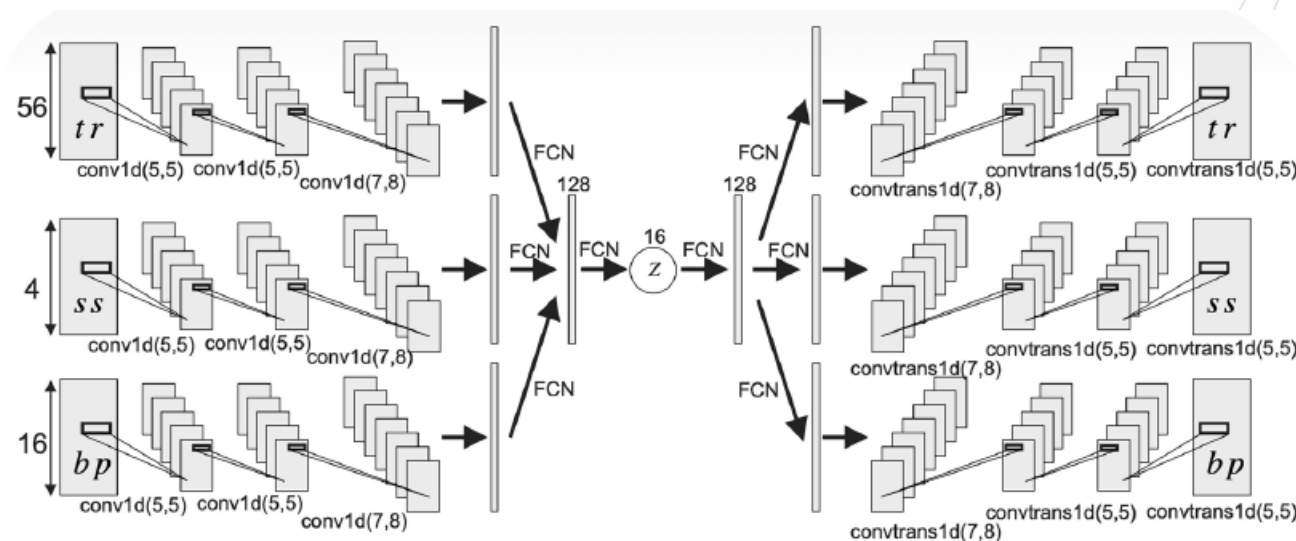
- Finally, each sequence can be represented by a **triplet**  $X$  (tr, ss, bp) along with the corresponding **weight**  $\pi_s$ .



# Method

## ■ RfamGen Architecture and Training

- RfamGen takes the **triplet**  $X$  (tr, ss, bp) as input of **VAE** and produces the **reconstructed triplet** as output.



$$L(X; \theta, \varphi) = \text{Reconst}_{\theta, \varphi}(X, \tilde{X}) - \beta \text{KL}(q_{\varphi}(Z|X) | p(Z))$$

$$\text{Reconst}_{\theta, \varphi}(X, \tilde{X}) = \pi_X \sum_{M \in \{\text{tr}, \text{ss}, \text{bp}\}} \text{cross entropy}(M, \tilde{M}_{\theta, \varphi})$$

```
from infernal_tools import CovarianceModel, make_deriv_dict_from_trsp
cmreader = CMReader(args.cmfile)
print("Start loading cm dict. This process may take much time for long sequences.")
cm_deriv_dict = cmreader.load_derivation_dict_from_cmfile()
out_dict = make_deriv_dict_from_trsp(cm_deriv_dict, (softmax(tr), softmax(s), softmax(p)))
cm = CovarianceModel(out_dict)
seq, _ = cm.cmemit(sample = False)[0]
```



# Experiment

## ■ Dataset

- RfamGen takes various RNA families from the **Rfam database** as dataset.

## ■ Baseline

- GCVAE (gapped character VAE), which includes only alignment
- GVAE (grammar VAE), which includes only secondary structure
- CVAE (character VAE), which includes neither

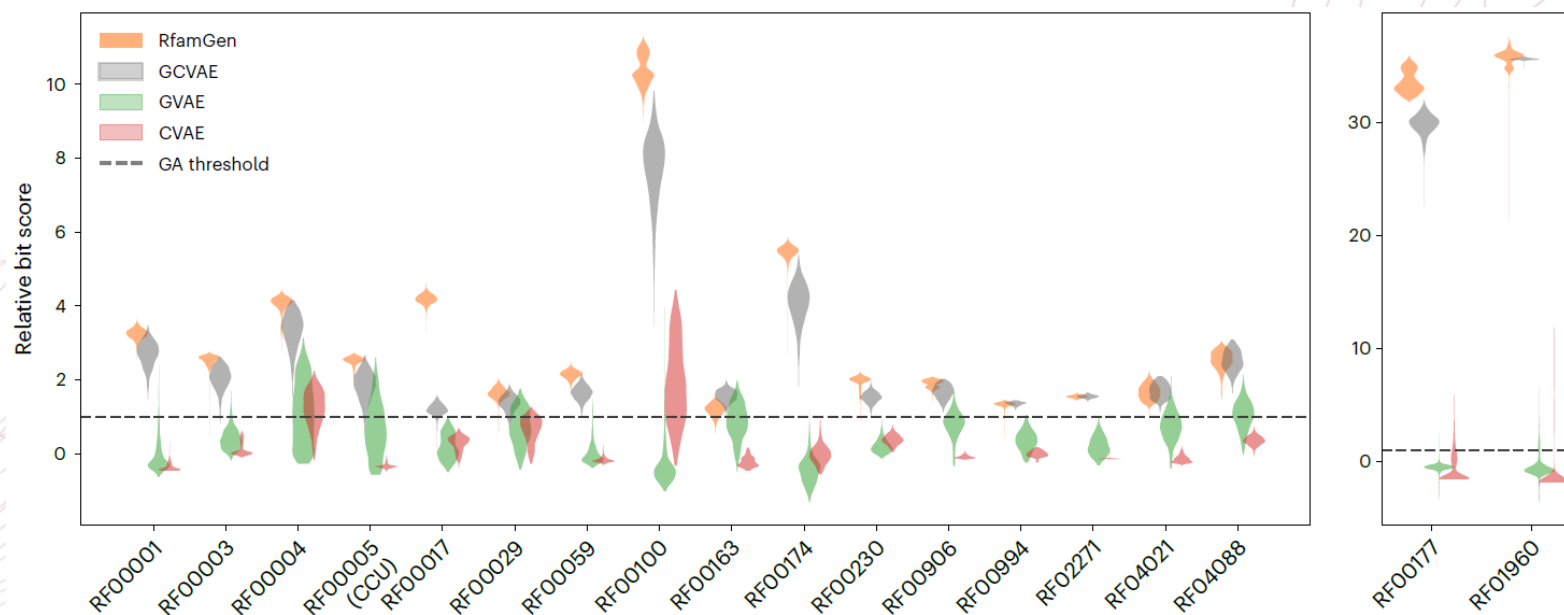
## ■ Evaluation

- We assessed the quality of the generated sequences by likelihood as the RNA family (**'bit score'**) calculated based on alignment to the 'ground truth' CM.

# Experiment

## ■ RfamGen is a data-efficient generative model

- We firstly use **18 RNA families** with full alignments composed of **at least 10,000 sequences** in the Rfam database as dataset, and quantify the **average bit score** of **1,000 randomly generated sequences** from the models for each RNA family.



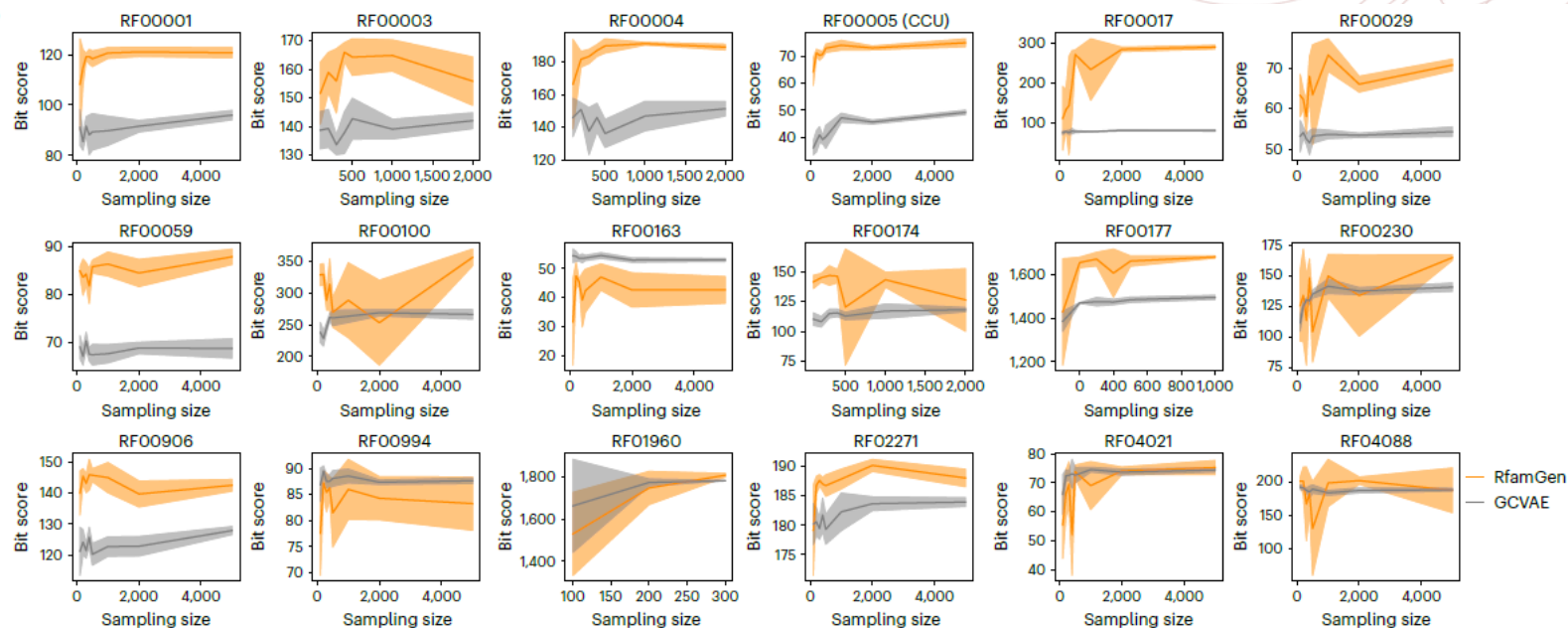
RfamGen performs best across most RNA families



# Experiment

## ■ RfamGen is a data-efficient generative model

- We next assess the **effectiveness and robustness** of RfamGen and GCVAE against **various data sizes** by undersampling.

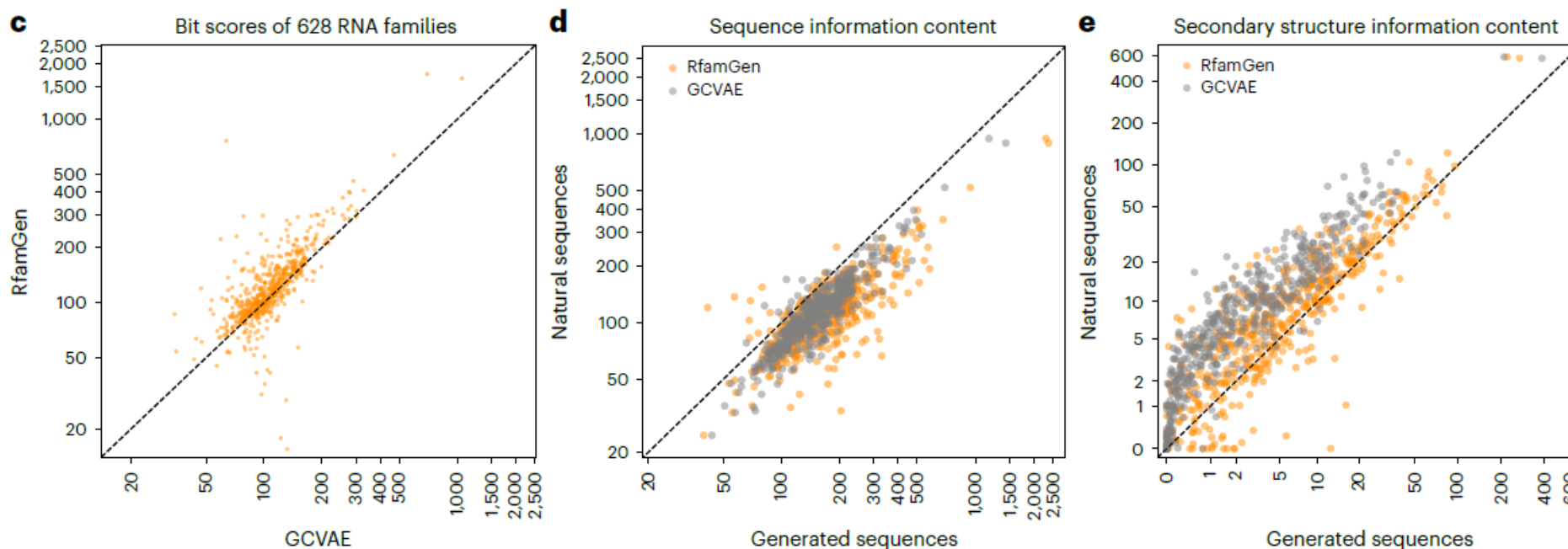


Only ~500 input sequences needed to reach near-peak performance for most RNA families  
Maintains high generation capability on small datasets

# Experiment

## ■ RfamGen is a data-efficient generative model

- Finally, we train RfamGen using **628 RNA families** whose full alignments consisted of **at least 100 sequences** in the Rfam database, and compare it to GCVAE to confirm its breadth of application.



Successfully learns nested base pair information

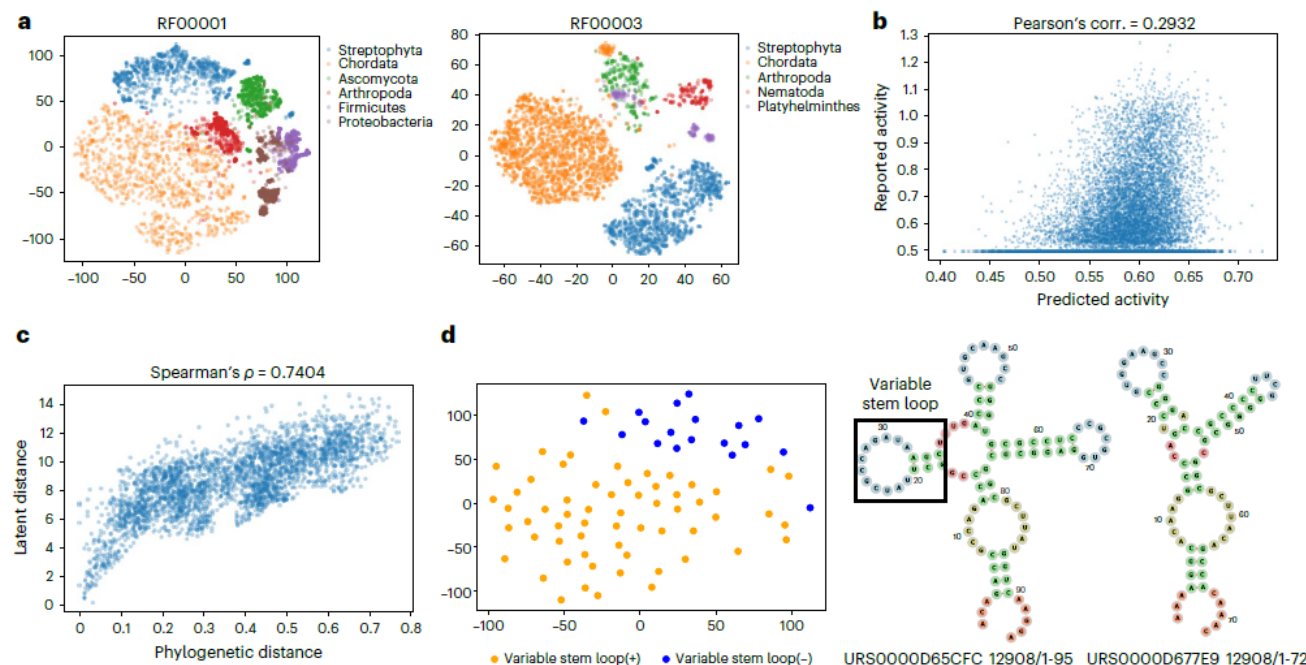




# Experiment

## ■ RfamGen learns a semantically rich latent space

- We examine whether the latent space of RfamGen contains meaningful representations of sequences.



**Fig. 3 | Latent space has semantically rich representation of sequences.**  
a, t-SNE visualization of the latent space trained with 5S ribosomal RNA (RF00001) and U1 spliceosomal RNA (RF00003), colored with phylogenetic information. b, Prediction of activity of tRNA variants reported by Li et al.<sup>31</sup> (y axis) and predicted activity by linear regression of latent space (x axis). c, Comparison of phylogenetic distance (x axis) and Euclidean distance in

latent space (latent distance, y axis) trained by twister-sister ribozymes (RF02681). d, t-SNE visualization of the latent space trained by twister-sister ribozyme (RF02681) colored with structural variants (left). Examples of the structural variants of a twister-sister ribozyme that are defined by the presence or absence of the variable stem loop shown in the black box (right), and were visualized by forna server.



# Conclusion

- In this work, we present RfamGen, which is a deep-learning method that designs RNA family sequences in a data-efficient manner.
- RfamGen achieves the first integration of CM architecture into deep generative models, and further provides semantically meaningful representations of sequences.
- Extensive experiments demonstrate its great generation capability and semantically rich latent space, thus building a powerful and general platform for RNA engineering with enormous application potential in biotechnology and medicine.



**中國人民大學**  
RENMIN UNIVERSITY OF CHINA



**高瓴人工智能学院**  
Gaoling School of Artificial Intelligence

**Thank You for listening!**

**Fanmeng Wang**

**2025-6-13**