# Benchmarking deep learning-based methods for RNA 3D structure prediction

**Ivona Martinović**

PhD student

Computational Approaches to RNA Structure and Function

Benasque Science Center, Jul 21 - Aug 03, 2024
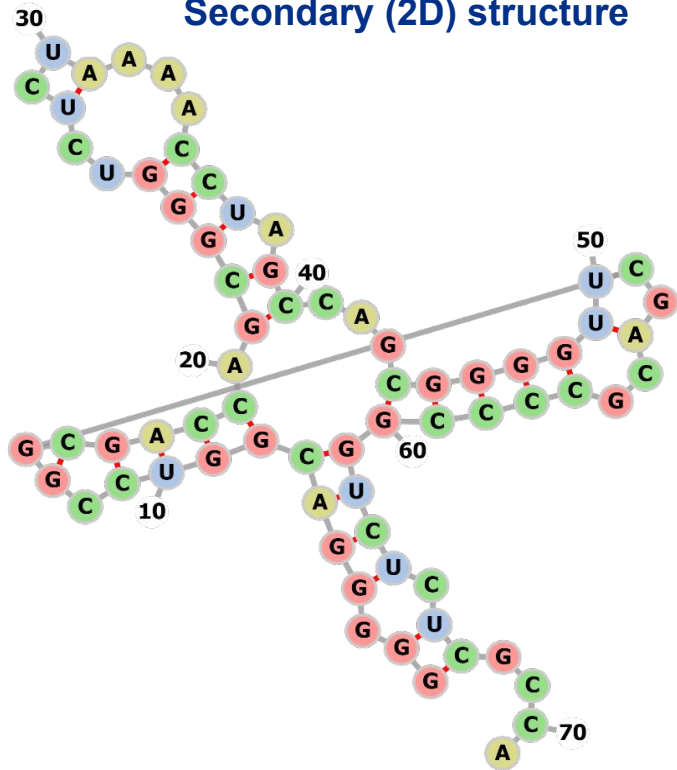
CREATING GROWTH, ENHANCING LIVES

# Outline

- Benchmarking **deep learning-based tools** for RNA 3D structure prediction
  - overview of the tools
  - datasets
  - results


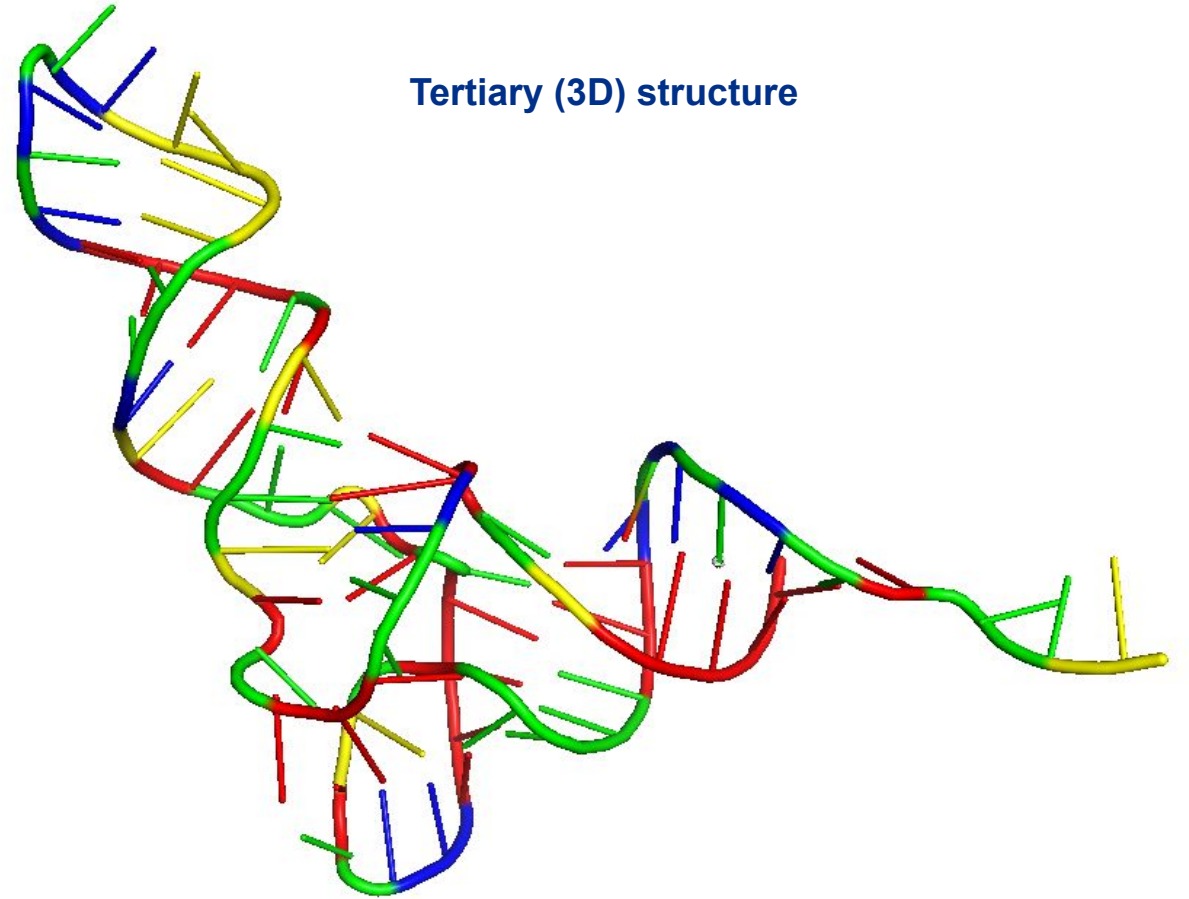- Next steps towards our structure prediction model: our **RNA language model** - RiNALMo

# RNA 3D structure prediction

**Primary structure / sequence**

>8UPT_A

GGGGGACGGUCCGGCGACCAGCGGGUCUCUAAAACCUAGCCAGCGGGGUUCGACGCCCCGGUCUCUCGCCA

**Secondary (2D) structure**



**Tertiary (3D) structure**



Computational Approaches to RNA Structure and Function, Benasque 2024

3

# Deep learning based RNA 3D structure prediction tools

# Deep learning based RNA 3D structure prediction tools

- DRfold

## nature communications

Explore content ⌄    About the journal ⌄    Publish with us ⌄

nature > nature communications > articles > article

Article | Open access | Published: 16 September 2023

# Integrating end-to-end learning with deep geometrical potentials for ab initio RNA structure prediction

Yang Li, Chengxin Zhang, Chenjie Feng, Robin Pearce, P. Lydia Freddolino ✉ & Yang Zhang ✉

**10k** Accesses | **11** Citations | **47** Altmetric | Metrics

# Deep learning based RNA 3D structure prediction tools

- DRfold

- DeepFoldRNA



HOME| SUBI

New Results

🔔 **Follow this preprint**

## De Novo RNA Tertiary Structure Prediction at Atomic Resolution Using Geometric Potentials from Deep Learning

Robin Pearce, Gilbert S. Omenn, 🆔 Yang Zhang

**doi:** https://doi.org/10.1101/2022.05.15.491755

This article is a preprint and has not been certified by peer review [what does this mean?].

# Deep learning based RNA 3D structure prediction tools

- DRfold

- DeepFoldRNA

- RhoFold



arXiv > q-bio > arXiv:2207.01586

Search...
Help | Advan

**Quantitative Biology > Quantitative Methods**

[Submitted on 4 Jul 2022]

**E2Efold-3D: End-to-End Deep Learning Method for accurate de novo RNA 3D Structure Prediction**

Tao Shen, Zhihang Hu, Zhangzhi Peng, Jiayang Chen, Peng Xiong, Liang Hong, Liangzhen Zheng, Yixuan Wang, Irwin King, Sheng Wang, Siqi Sun, Yu Li

Computational Approaches to RNA Structure and Function, Benasque 2024

# Deep learning based RNA 3D structure prediction tools

- DRfold

- DeepFoldRNA

- RhoFold

- RoseTTAFoldNA

## nature methods

Explore content ⌄      About the journal ⌄      Publish with us ⌄

nature > nature methods > articles > article

Article | Open access | Published: 23 November 2023

## Accurate prediction of protein–nucleic acid complexes using RoseTTAFoldNA

Minkyung Baek, Ryan McHugh, Ivan Anishchenko, Hanlun Jiang, David Baker & Frank DiMaio ✉

*Nature Methods* **21**, 117–121 (2024) | Cite this article

**42k** Accesses | **37** Citations | **125** Altmetric | Metrics

# Deep learning based RNA 3D structure prediction tools

- DRfold

- DeepFoldRNA

- RhoFold

- RoseTTAFoldNA

- trRosettaRNA

## nature communications

Explore content ∨    About the journal ∨    Publish with us ∨

nature > nature communications > articles > article

Article | Open access | Published: 09 November 2023

# trRosettaRNA: automated prediction of RNA 3D structure with transformer network

Wenkai Wang, Chenjie Feng, Renmin Han, Ziyi Wang, Lisha Ye, Zongyang Du, Hong Wei, Fa Zhang ✉, Zhenling Peng ✉ & Jianyi Yang ✉

Computational Approaches to RNA Structure and Function, Benasque 2024

# Deep learning based RNA 3D structure prediction tools

- DRfold

- DeepFoldRNA

- RhoFold

- RoseTTAFoldNA

- trRosettaRNA

- AlphaFold3

## nature

Explore content ⌄    About the journal ⌄    Publish with us ⌄

nature > articles > article

Article | Open access | Published: 08 May 2024

### Accurate structure prediction of biomolecular interactions with AlphaFold 3

Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J. Ballard, Joshua Bambrick, Sebastian W. Bodenstein, David A. Evans, Chia-Chun Hung, Michael O'Neill, David Reiman, Kathryn Tunyasuvunakool, Zachary Wu, Akvilė Žemgulytė, Eirini Arvaniti, Charles Beattie, Ottavia Bertolli, Alex Bridgland, Alexey Cherepanov, Miles Congreve, ... John M. Jumper ✉  + Show authors

*Nature*  **630**, 493–500 (2024) | Cite this article
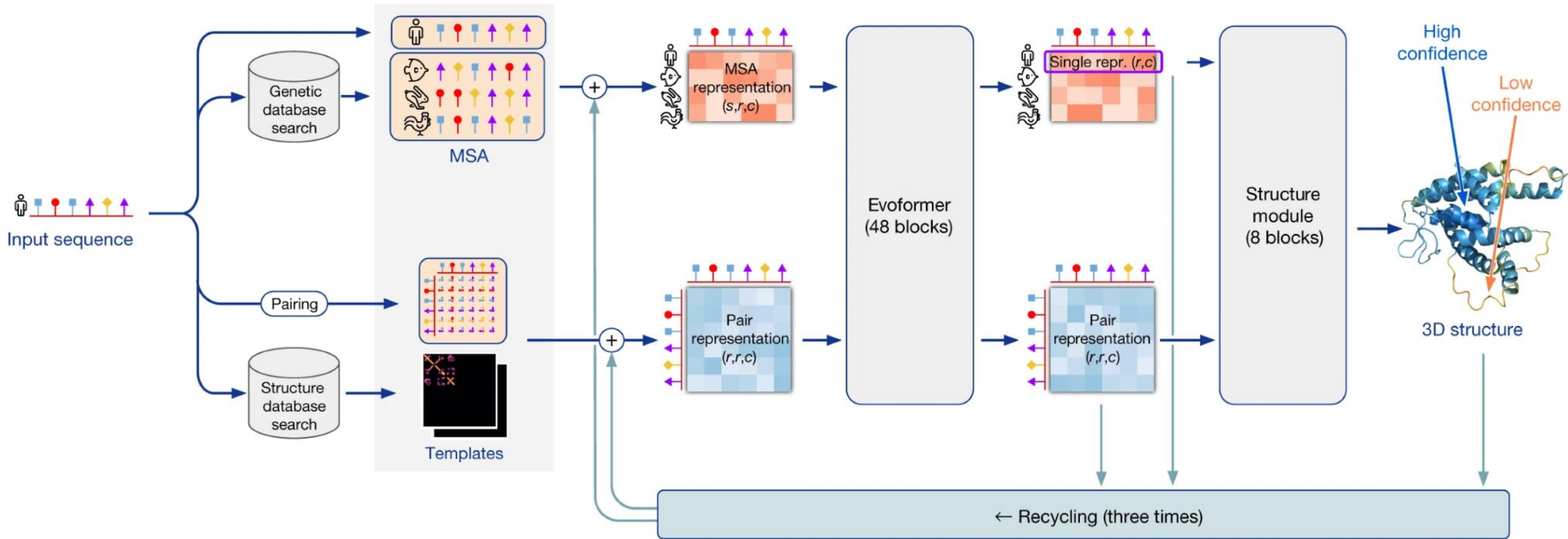
**417k** Accesses | **81** Citations | **1528** Altmetric | Metrics

Computational Approaches to RNA Structure and Function, Benasque 2024
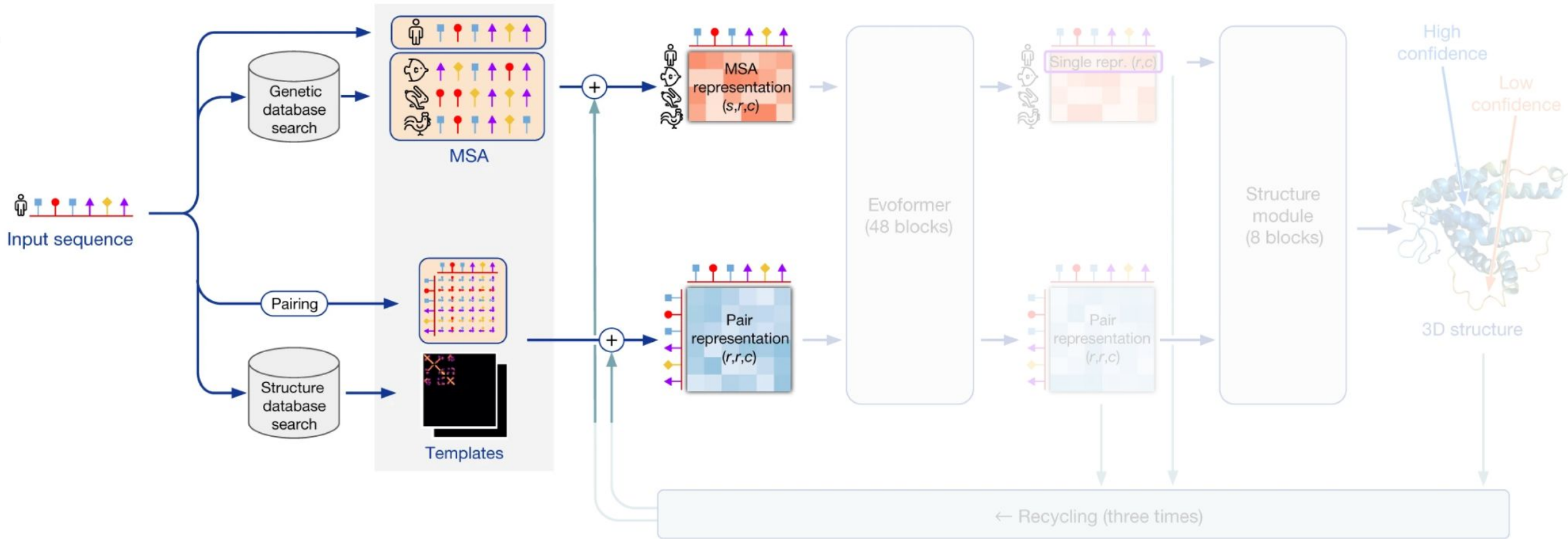
# Research questions

- which of the tools performs best across different datasets and evaluation metrics?

- do certain design choices and methodologies impact accuracy?

- how well these tools generalize to RNA sequences different from those used in their training?

- can we choose the best predicted structure using ARES or Rosetta score?

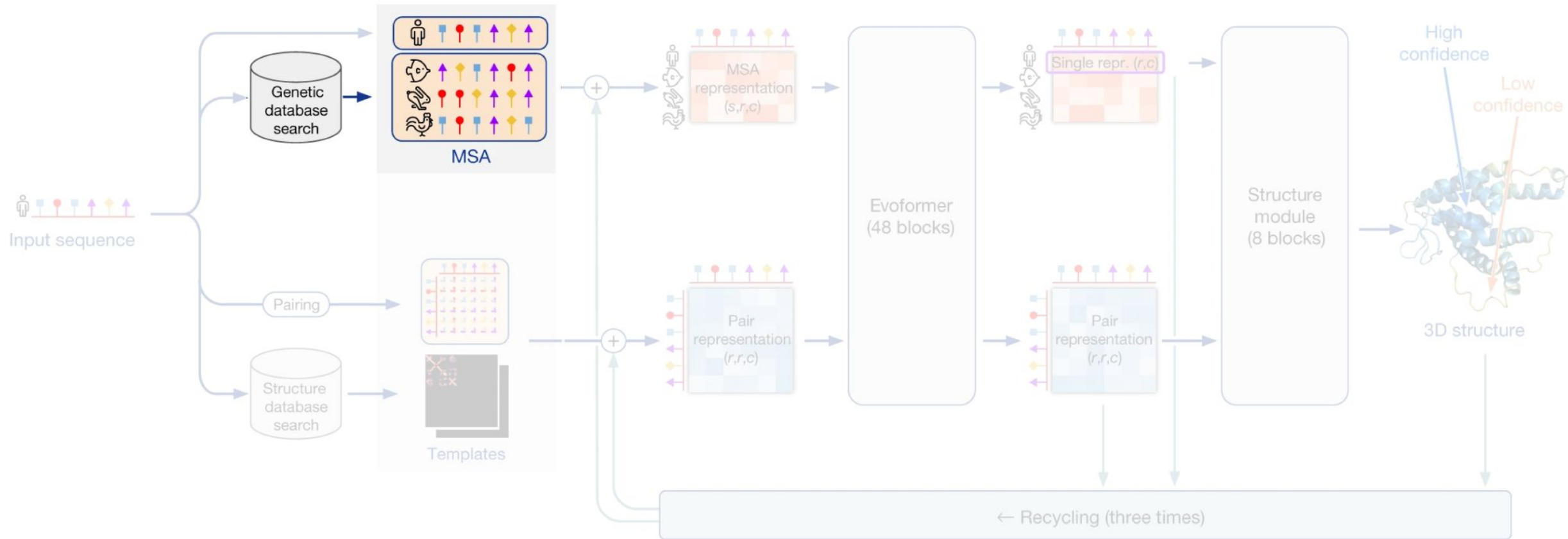- (with AF3) how much does having context (other chains from the complex) help in structure prediction?

Computational Approaches to RNA Structure and Function, Benasque 2024

# AlphaFold2's architecture



*Figure adapted from [Jumper et al., 2021]

Computational Approaches to RNA Structure and Function, Benasque 2024

# Data preprocessing



*Figure adapted from [Jumper et al., 2021]

Computational Approaches to RNA Structure and Function, Benasque 2024
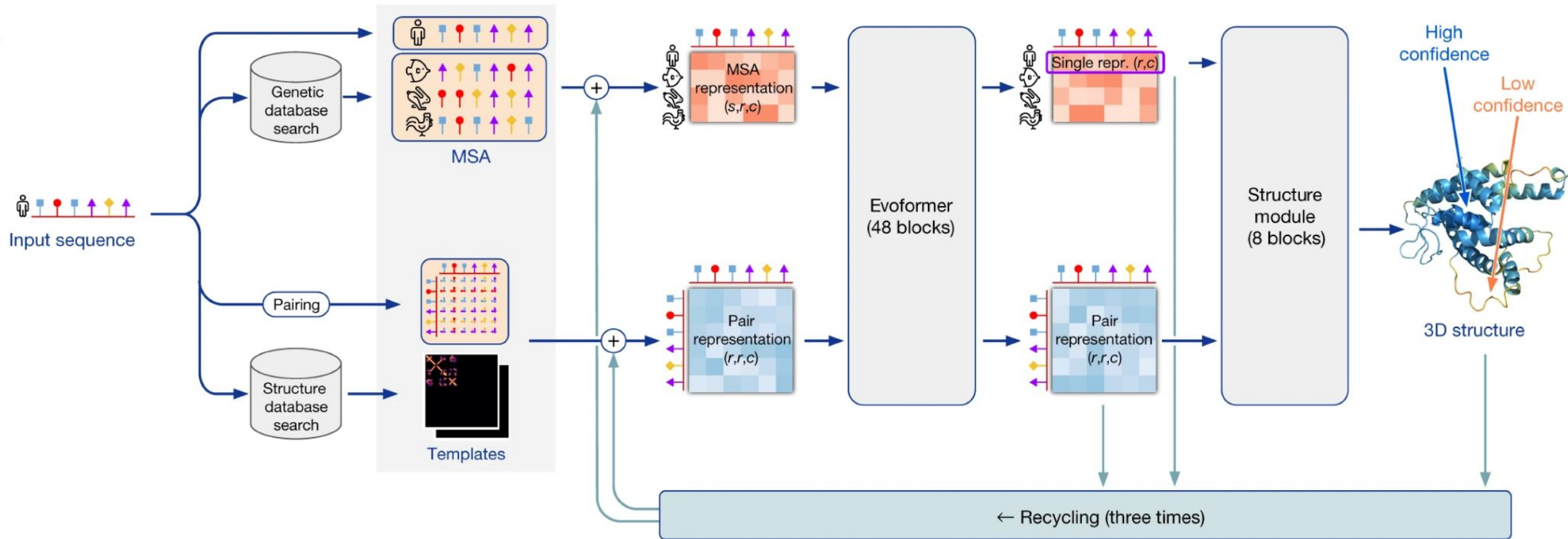
# Data preprocessing: multiple sequence alignment (MSA)



*Figure adapted from [Jumper et al., 2021]
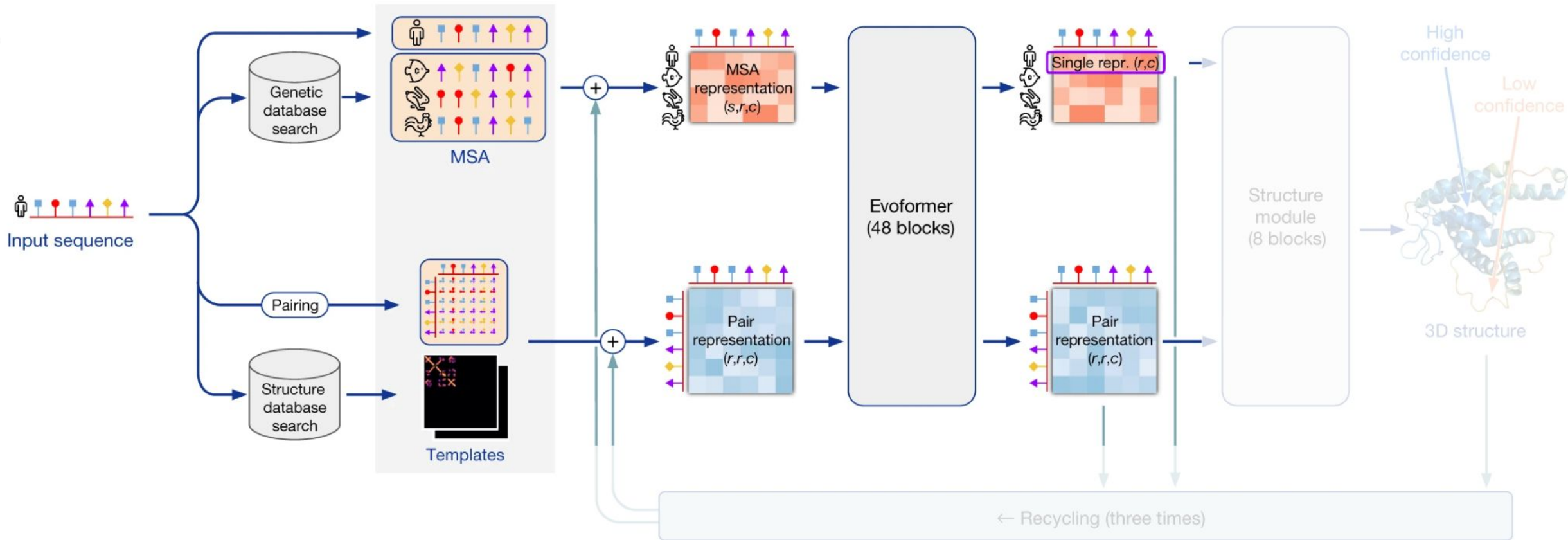
# Data preprocessing: secondary structures (SS)



*Figure adapted from [Jumper et al., 2021]

Computational Approaches to RNA Structure and Function, Benasque 2024

# End-to-end vs. predicting geometric restraints



*Figure adapted from [Jumper et al., 2021]

Computational Approaches to RNA Structure and Function, Benasque 2024

# End-to-end vs. predicting geometric restraints



*Figure adapted from [Jumper et al., 2021]

Computational Approaches to RNA Structure and Function, Benasque 2024

CREATING GROWTH, ENHANCING LIVES

# From geometric restraints to structural model



*Figure adapted from [Pearce et al., 2022]

# Similarities and differences between RNA models

| Tool | end-to-end | uses MSA | uses SS | uses LM |
|---|---|---|---|---|
| DRfold | ✅ | ❌ | ✅ | ❌ |
| DeepFoldRNA | ❌ | ✅ | ✅ | ❌ |
| RhoFold | ✅ | ✅ | ❌ | ✅ |
| RoseTTAFoldNA | ✅ | ✅ | ❌ | ❌ |
| trRosettaRNA | ❌ | ✅ | ✅ | ❌ |

Computational Approaches to RNA Structure and Function, Benasque 2024

# AlphaFold 3 - main differences

- not only for proteins, AF3 works with RNA chains, ligands and combinations of all three

- changed structure module - diffusion module

- working on atom level instead of residue level



Computational Approaches to RNA Structure and Function, Benasque 2024

# Datasets

- Dataset 1 = RNA Puzzles
    - 37 RNAs - puzzles 35 and 36 removed since they are in CASP15 dataset

- Dataset 2 = CASP15
    - 12 RNAs - 8 natural & 4 synthetic

Length distribution for Dataset 3

- Dataset 3 = curated dataset from PDB
    - 190 RNAs
    - published after April 2022 in Protein Data Bank (PDB)
    - clustered with sequence identity 90% - only considering ones which are not similar to those prior to April 2022
    - filtering: length < 10nt, resolution > 9Å, %defined residues < 90%, sequences containing only 'N' or 'X'
    - 329 clusters => only 190 without errors for all tools



Computational Approaches to RNA Structure and Function, Benasque 2024

# Results

# Results - RMSD



**Note:**
different datasets =
different tool with
the lowest RMSD

Computational Approaches to RNA Structure and Function, Benasque 2024

# Results - TM-score



**Note:**
for Dataset 1 RoseTTAFoldNA is the best (has the highest TM-score), not trRosettaRNA (as in the case of RMSD)

Computational Approaches to RNA Structure and Function, Benasque 2024

# Percentage of RNAs for which each tool was the best

according to RMSD:

according to TM-score:

# Scoring functions - ARES and Rosetta scores

- ARES is a deep learning method for scoring RNA structures
  (input = PDB file,
  output = score for given structure - lower values mean better structure)

- Rosetta score = rna_score
  from Rosetta toolkit

RAPHAEL J. L. TOWNSHEND (iD), STEPHAN EISMANN, ANDREW M. WATKINS (iD), RAMYA RANGAN (iD), MASHA KARELINA (iD), RHIJU DAS (iD), AND RON O. DROR (iD)    Authors Info & Affiliations

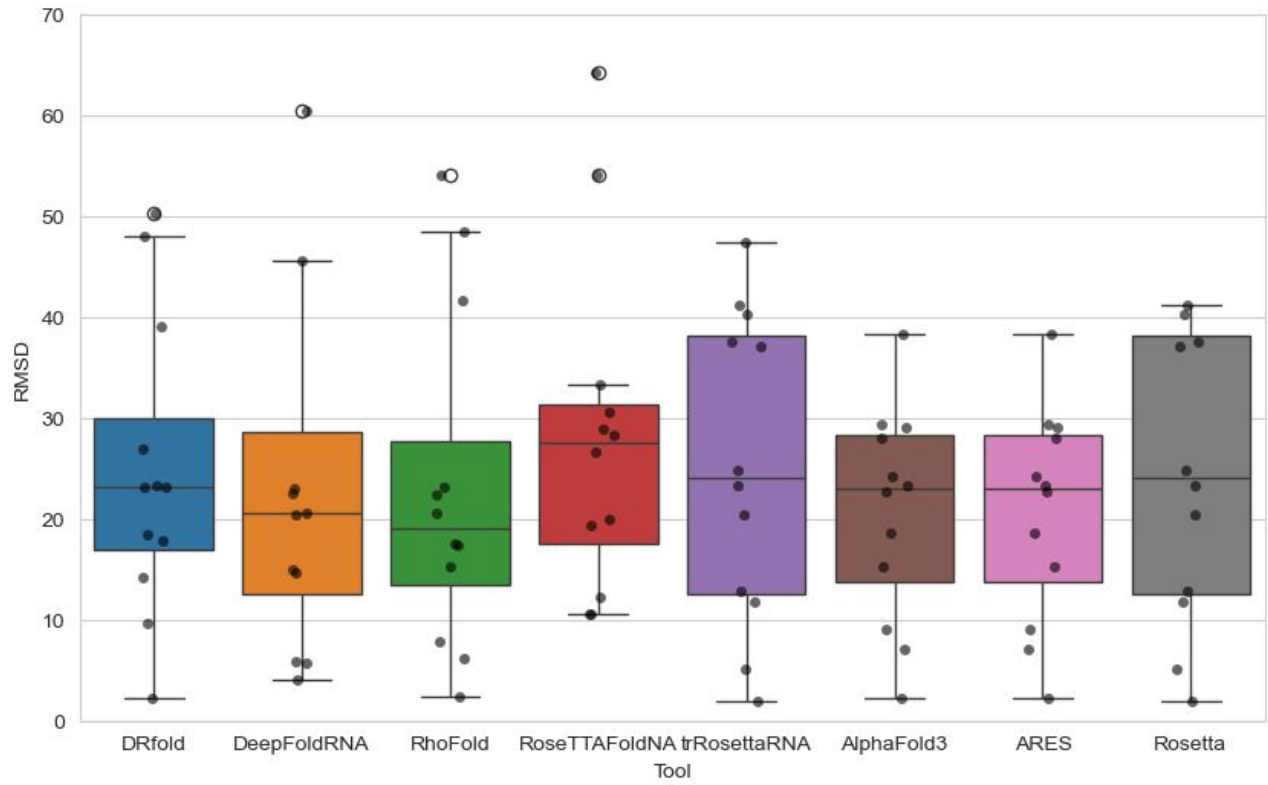# Dataset 1 - RNA Puzzles

**Note:**
Rosetta score was the lowest for trRosettaRNA's model for all RNAs

# Dataset 1 - RNA Puzzles

**Note:**
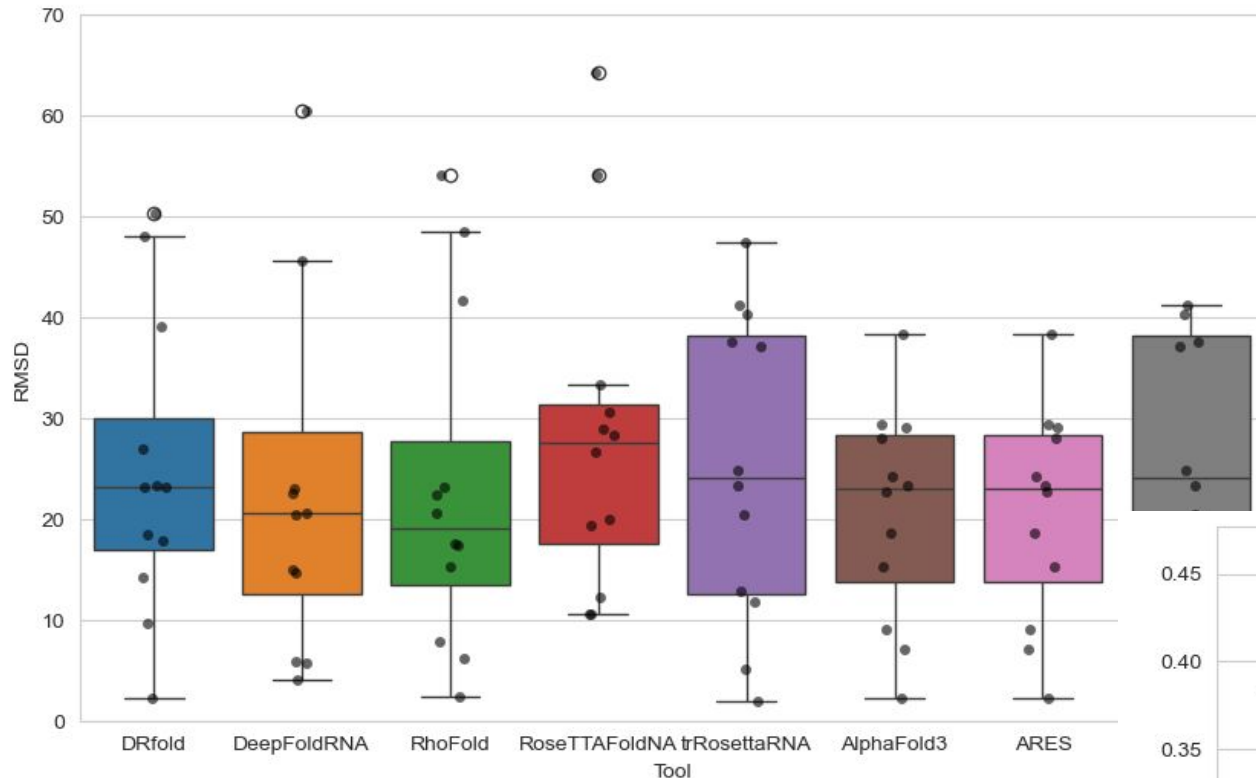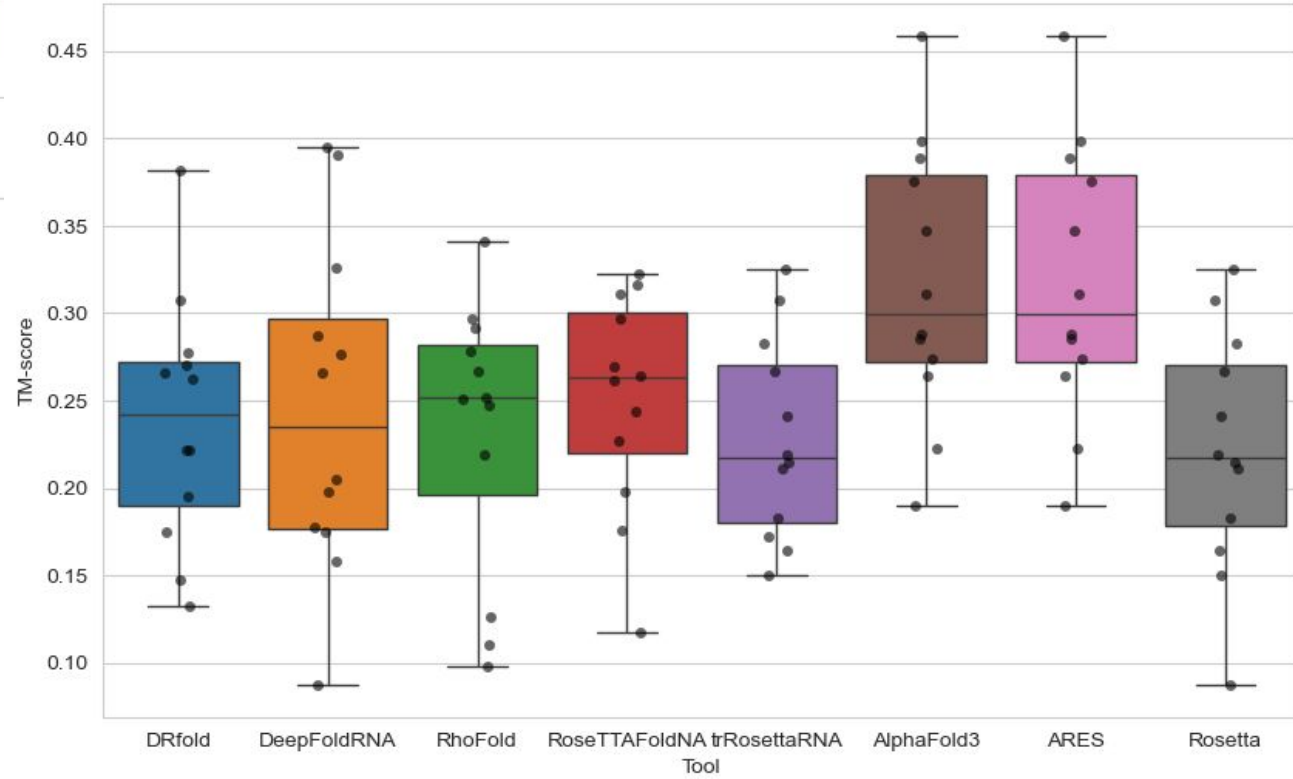Rosetta score was the lowest for trRosettaRNA's model for all RNAs

Computational Approaches to RNA Structure and Function, Benasque 2024

**Dataset 2 - CASP15**

**Note:**
ARES selects AF3 models
as best for all RNA targets.

## Dataset 2 - CASP15

**Note:**
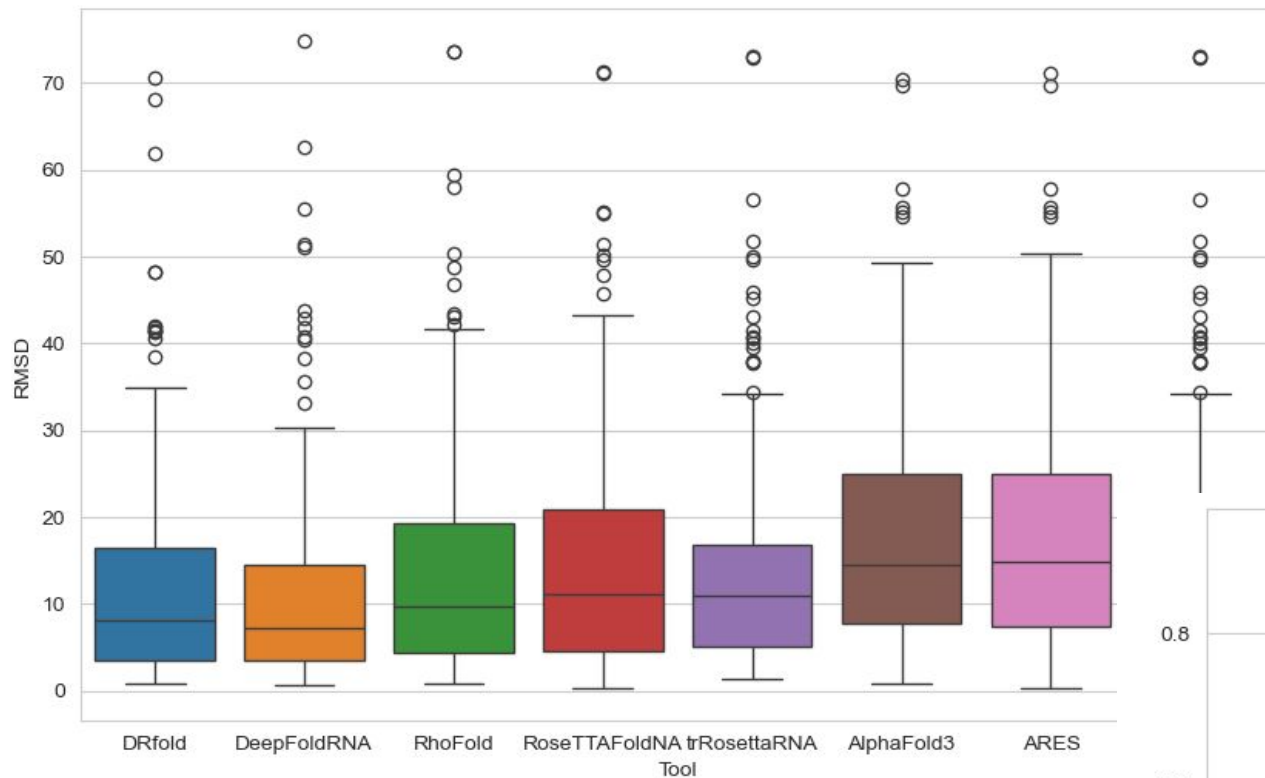ARES selects AF3 models as best for all RNA targets.

Computational Approaches to RNA Structure and Function, Benasque 2024

30

# Dataset 2 - CASP15



**Note:**
all DL-based methods struggle with synthetic RNAs

synthetic RNA

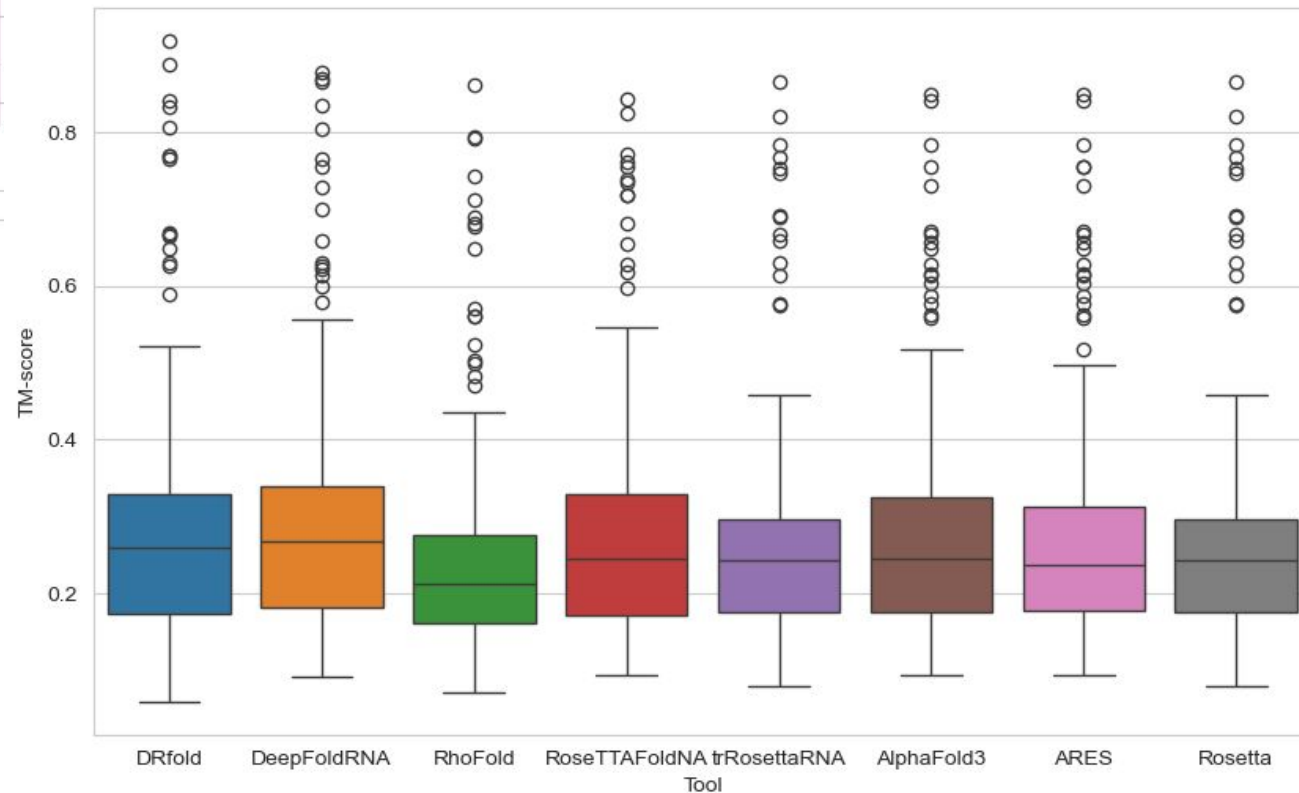Computational Approaches to RNA Structure and Function, Benasque 2024
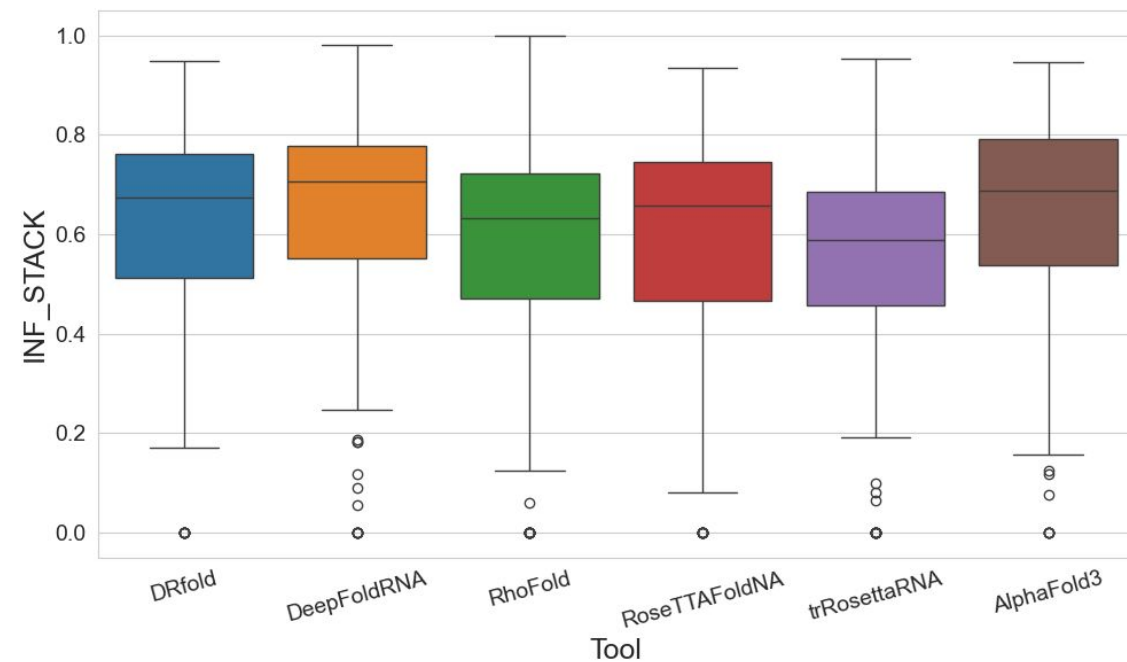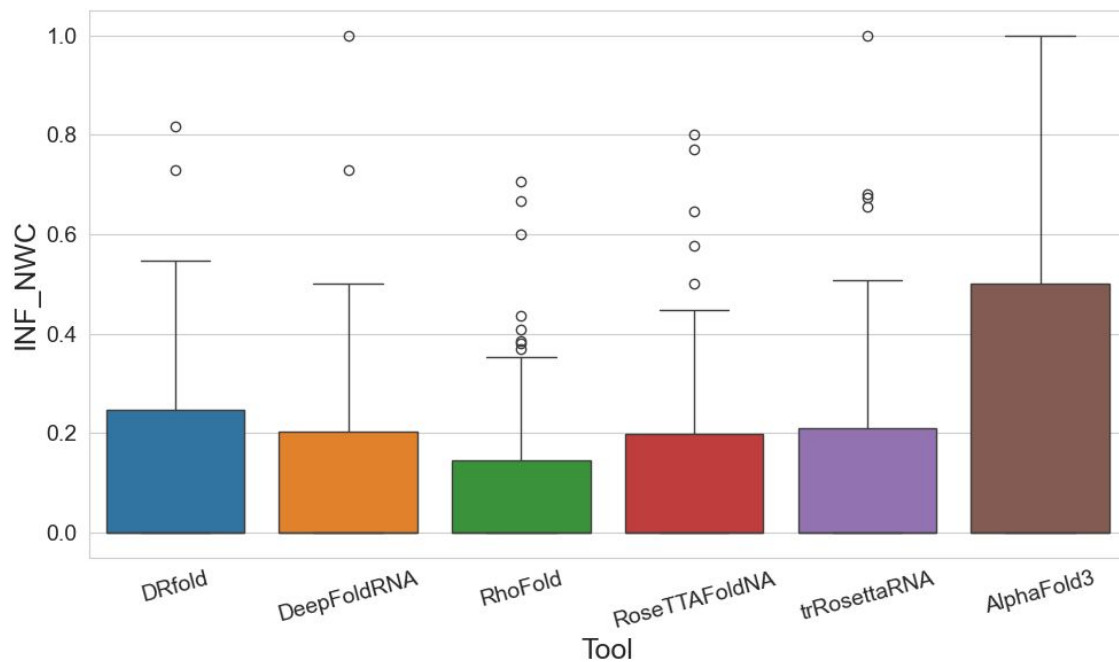
**Dataset 3**

**Note:**
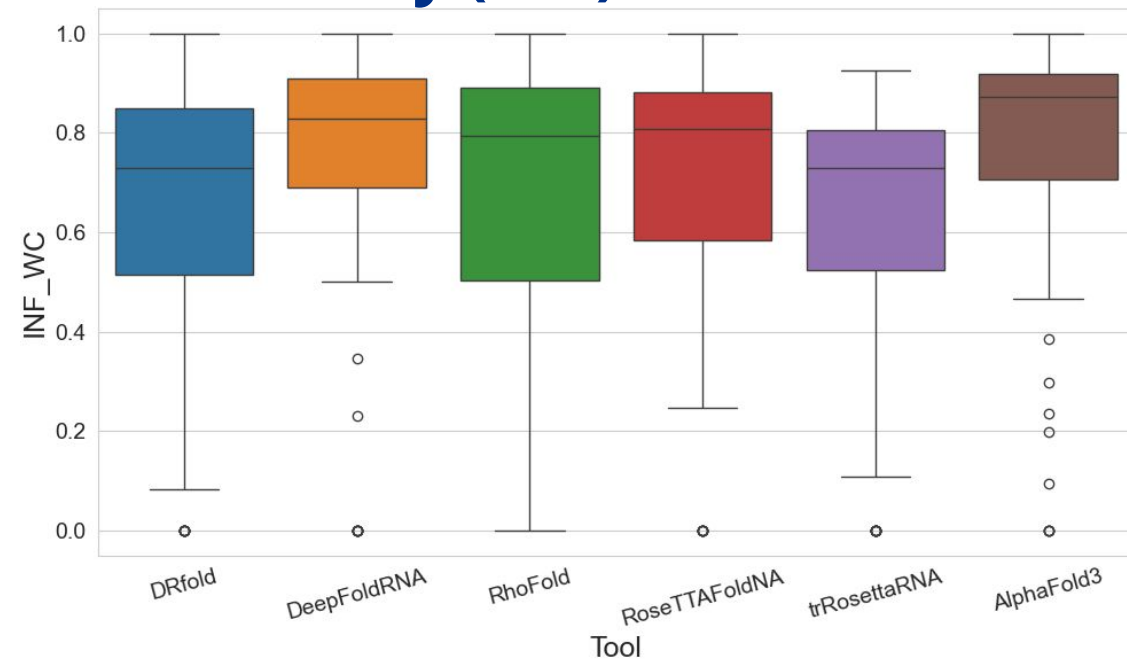again, Rosetta score was the lowest for trRosettaRNA's model for all RNAs
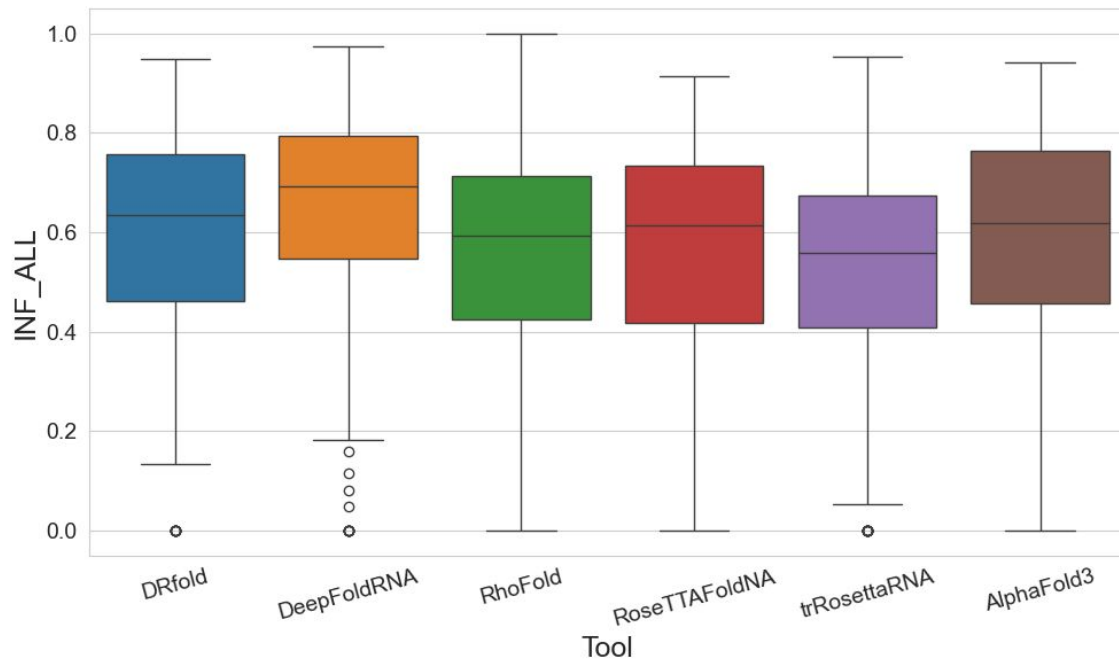
**Dataset 3**

**Note:**
again, Rosetta score was the lowest for trRosettaRNA's model for all RNAs

# Dataset 3: Interaction Network Fidelity (INF)

# Dataset 3: Local Distance Difference Test (lDDT)

# Dataset 3: Clash score



Computational Approaches to RNA Structure and Function, Benasque 2024

# Dataset 3: Clash score



Computational Approaches to RNA Structure and Function, Benasque 2024

# Performance (RMSD) per length group (Dataset 3)

# Performance (TM-score) per length group (Dataset 3)



Performance per length group

# Execution time (for all three datasets)



Computational Approaches to RNA Structure and Function, Benasque 2024

# Examples

# Example: 8A22_A8



DRfold

DeepFoldRNA

RhoFold

RoseTTAFoldNA

trRosettaRNA

AlphaFold3

# Example: 7WM4_B



**Note - colors:**
native = bright green
DRfold = blue
DeepFoldRNA = orange
RhoFold = dark green
RF2NA = red
trRosettaRNA = purple
AF3 = cyan

Computational Approaches to RNA Structure and Function, Benasque 2024

# RNA chain as a single chain vs. as part of the complex

- RoseTTAFoldNA has an option of providing protein sequences in which it still predicts only RNA chain(s), but in couple of examples where we tried this the resulting predictions were almost the same

- RoseTTAFoldNA tried out on 10 complexes - better only in 50% of cases, not by much



Computational Approaches to RNA Structure and Function, Benasque 2024

# RNA chain as a single chain vs. as part of the complex

- RoseTTAFoldNA has an option of providing protein sequences in which it still predicts only RNA chain(s), but in couple of examples where we tried this the resulting predictions were almost the same

- RoseTTAFoldNA tried out on 10 complexes - better only in 50% of cases, not by much

- AlphaFold3 for these 10 complexes - better for 9/11 chains, for 7YCH_B 91% lower RMSD in complex prediction

- Not sure how trustworthy AlphaFold3 results are (it could be trained on these examples)

AlphaFold3 - single chain vs. complex prediction

*(scatter plot: x-axis "RMSD from single chain prediction" 0–60, y-axis "RMSD from complex prediction" 0–60, with diagonal reference line)*

Computational Approaches to RNA Structure and Function, Benasque 2024

# Conclusion

- Q: Which of the tools performs best across different datasets and evaluation metrics?
  A: No unique tool, depends on the use case.

- Q: Do certain design choices and methodologies impact accuracy?
  A: Unable to answer, no direct connections.

- Q: How well these tools generalize to RNA sequences different from those used in their training?
  A: The best tool on generalization dataset (Dataset 3) is DeepFoldRNA across most metrics.

- Q: Can we choose the best predicted structure using ARES or Rosetta score?
  A: According to our tests, no.

- Q: How much does having context help in structure prediction?
  A: In case of AlphaFold3, currently it seems a lot, but in case of RoseTTAFoldNA, not that much.
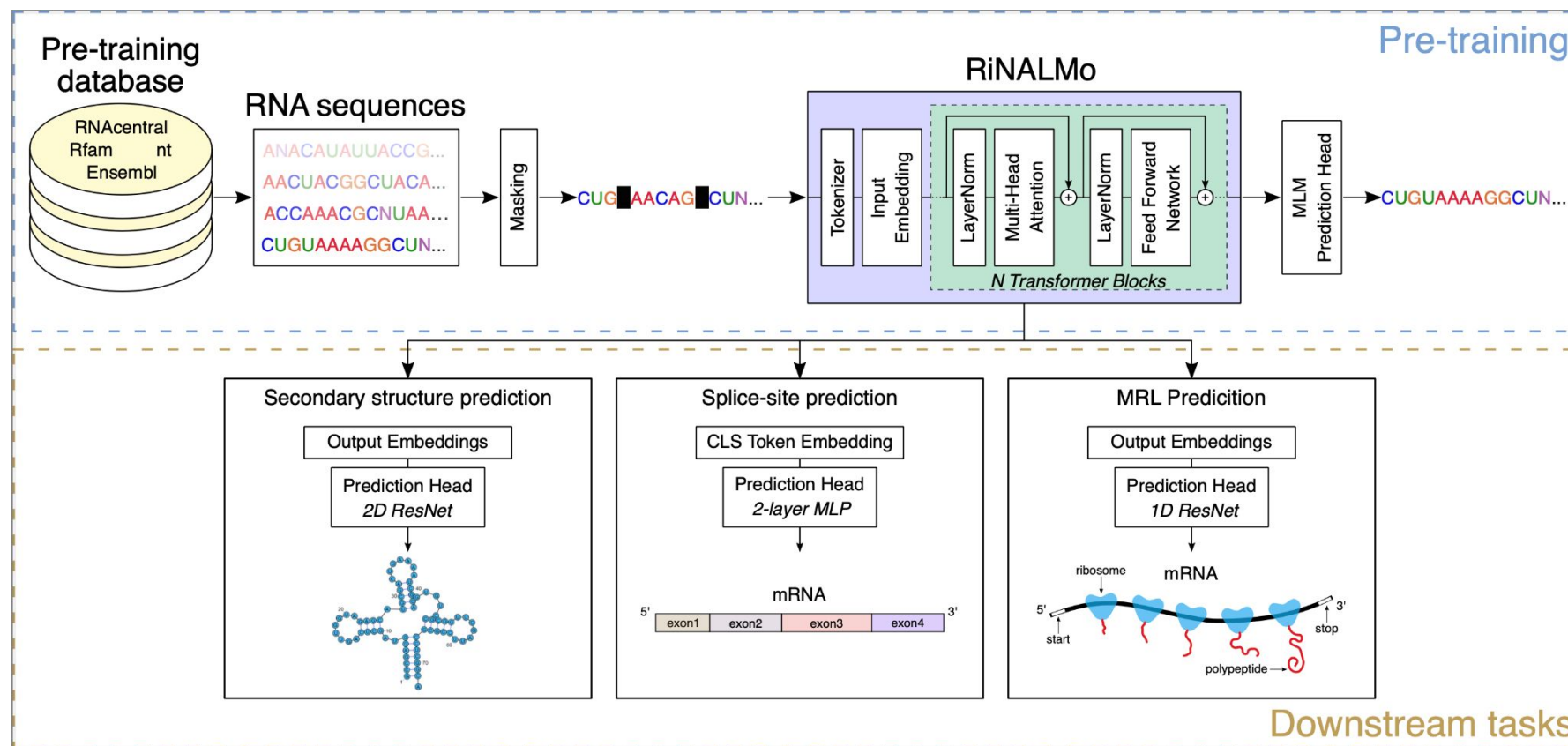
# RiNALMo: RNA language model

- **Motivation:** success of protein language models
- BERT-style language model pretrained using masked language modeling
- 36M unannotated ncRNA sequences, 650M parameters

Rafael Josip Penić

Tin Vlašić

CREATING GROWTH, ENHANCING LIVES

# Secondary structure prediction

- Finetuned RiNALMo embeddings + small CNN prediction head

- RiNALMo helps to generalize well on RNA families not seen in the training dataset unlike other deep learning methods

*Dataset obtained from [Szikszai et al., 2022]

| Test Family | RNAstructure | CONTRAfold | RiNALMo | RNA-FM | MXfold2 | UFold |
|---|---|---|---|---|---|---|
| 5S rRNA | 0.63 | 0.67 | **0.88** | 0.52 | 0.54 | 0.53 |
| SRP RNA | 0.63 | 0.60 | **0.70** | 0.25 | 0.50 | 0.26 |
| tRNA | 0.70 | 0.76 | **0.93** | 0.78 | 0.64 | 0.26 |
| tmRNA | 0.43 | 0.44 | **0.80** | 0.29 | 0.46 | 0.40 |
| RNase P RNA | 0.55 | 0.60 | **0.80** | 0.30 | 0.51 | 0.41 |
| Group I intron | 0.54 | 0.59 | **0.66** | 0.16 | 0.45 | 0.45 |
| 16S rRNA | 0.57 | 0.60 | **0.74** | 0.13 | 0.55 | 0.41 |
| Telomerase RNA | 0.50 | 0.54 | 0.12 | 0.08 | 0.34 | **0.80** |
| 23S rRNA | 0.73 | 0.75 | **0.85** | 0.17 | 0.64 | 0.45 |
| Mean | 0.59 | 0.62 | **0.72** | 0.30 | 0.51 | 0.44 |

Computational Approaches to RNA Structure and Function, Benasque 2024

# Future directions

- Currently pretraining a 1.6B parameter RiNALMo on ~100M RNA sequences
- Multimodal pretraining including chemical probing data
- 3D structure prediction model leveraging RiNALMo sequence embeddings

**Collaborators:**

WAN Yue, GIS          Roland G. HUBER, BII

Agency for
Science, Technology
and Research

SINGAPORE

CREATING GROWTH, ENHANCING LIVES

Mile SIKIC,
Group Leader

# THANK YOU

www.a-star.edu.sg

**Contacts:**
✉ martinovici@gis.a-star.edu.sg
✉ miles@gis.a-star.edu.sg

**Šikić Lab**

**AI IN GENOMICS**

https://sikic-lab.github.io