



UNIVERSITY OF  
CAMBRIDGE

# gRNAde

## Geometric deep learning for 3D RNA inverse design

**Chaitanya K. Joshi**, Arian R. Jamasb, Ramon Viñas, Charles Harris, Simon Mathis, Pietro Liò

*Computational Biology Workshop, International Conference on Machine Learning, 2023*

Forthcoming book chapter in *Methods in Molecular Biology (RNA Design: Methods and Protocols)*



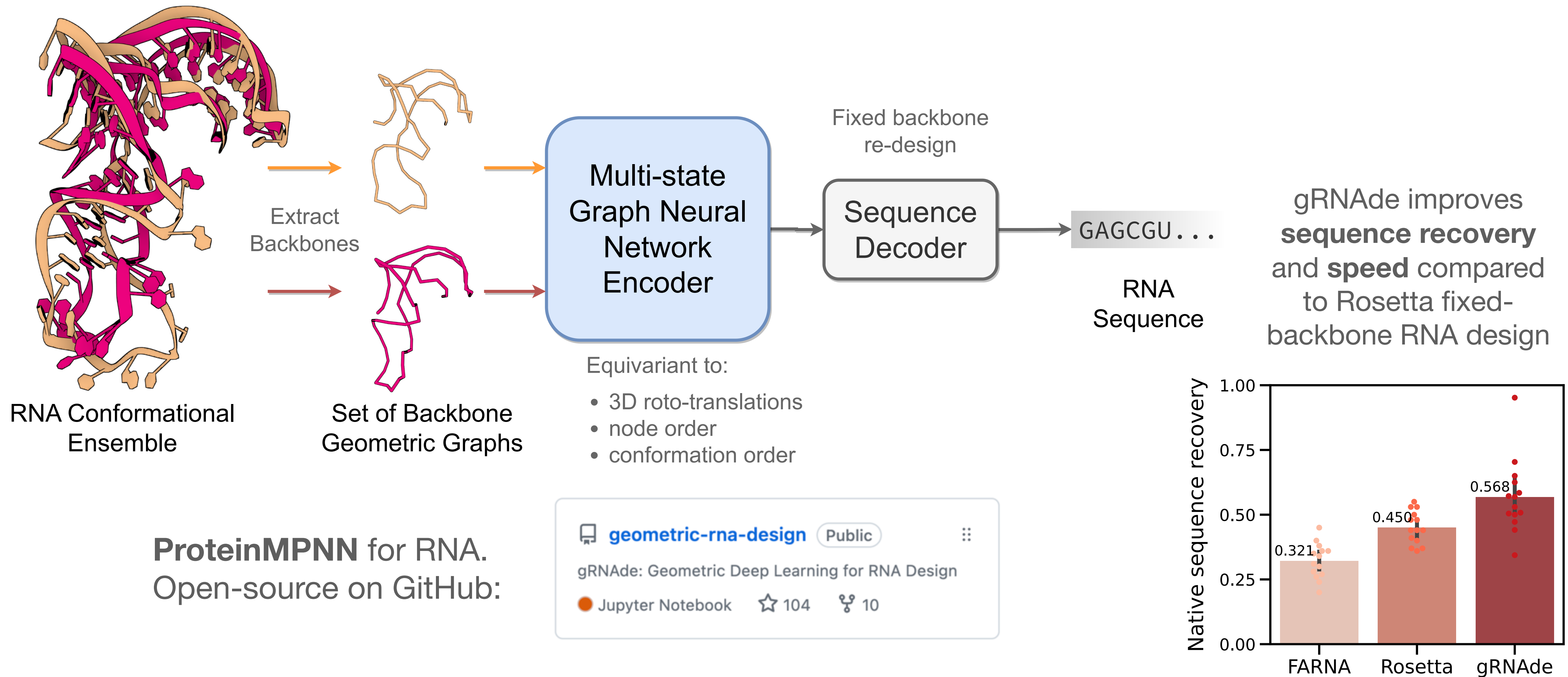
**Preprint:** <https://arxiv.org/abs/2305.14749>



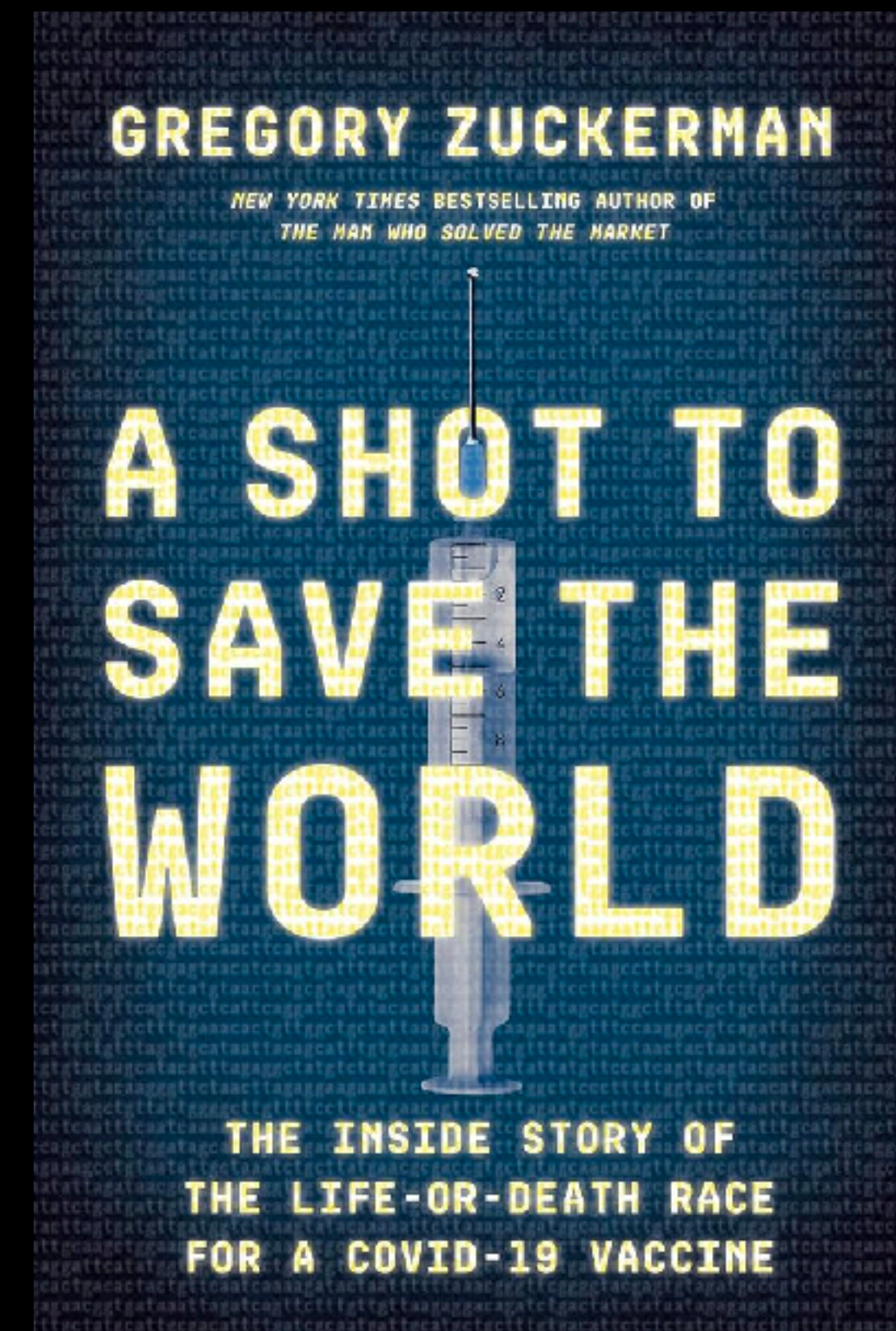
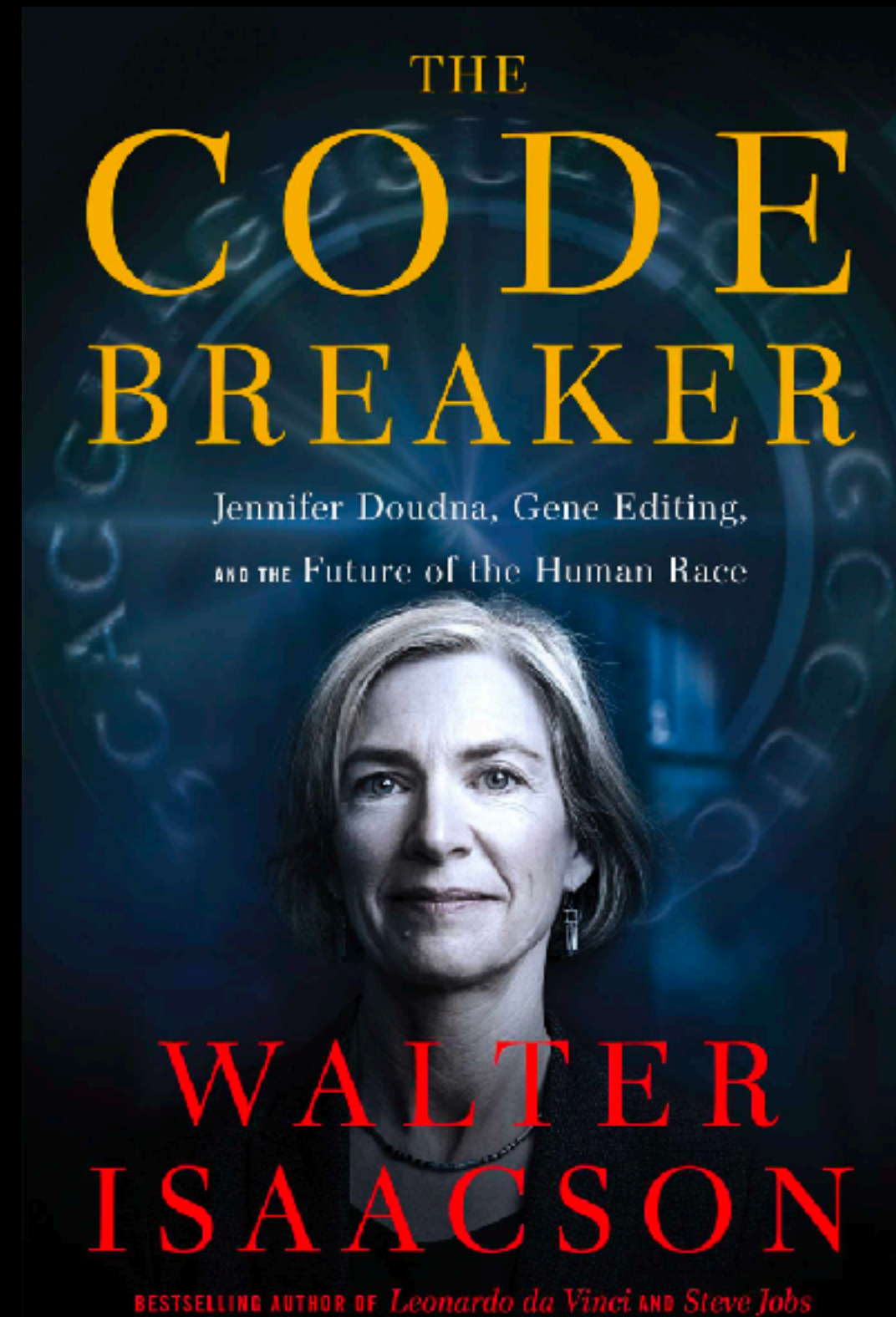
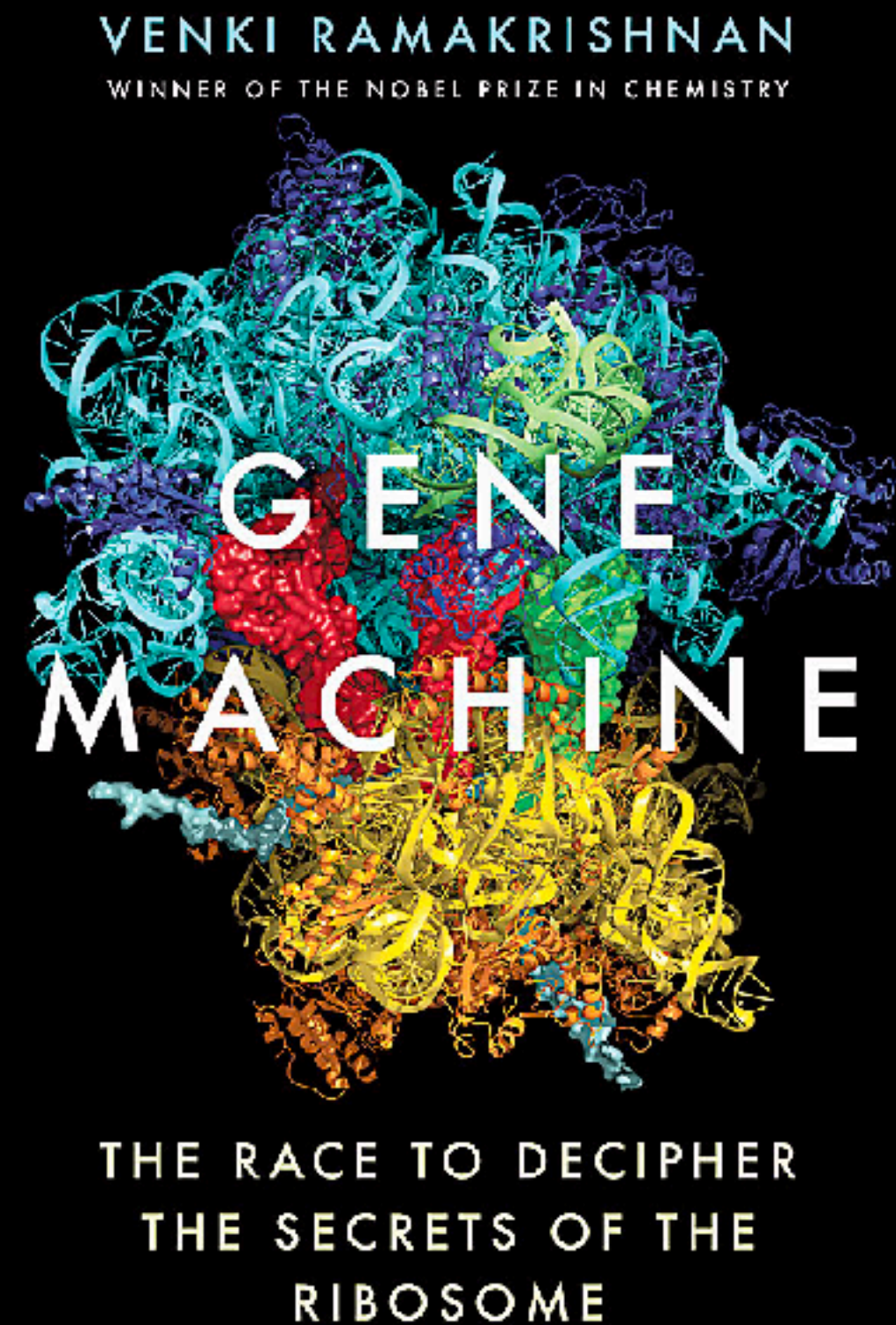
**Codebase:** [github.com/chaitjo/geometric-rna-design](https://github.com/chaitjo/geometric-rna-design)

# Executive summary

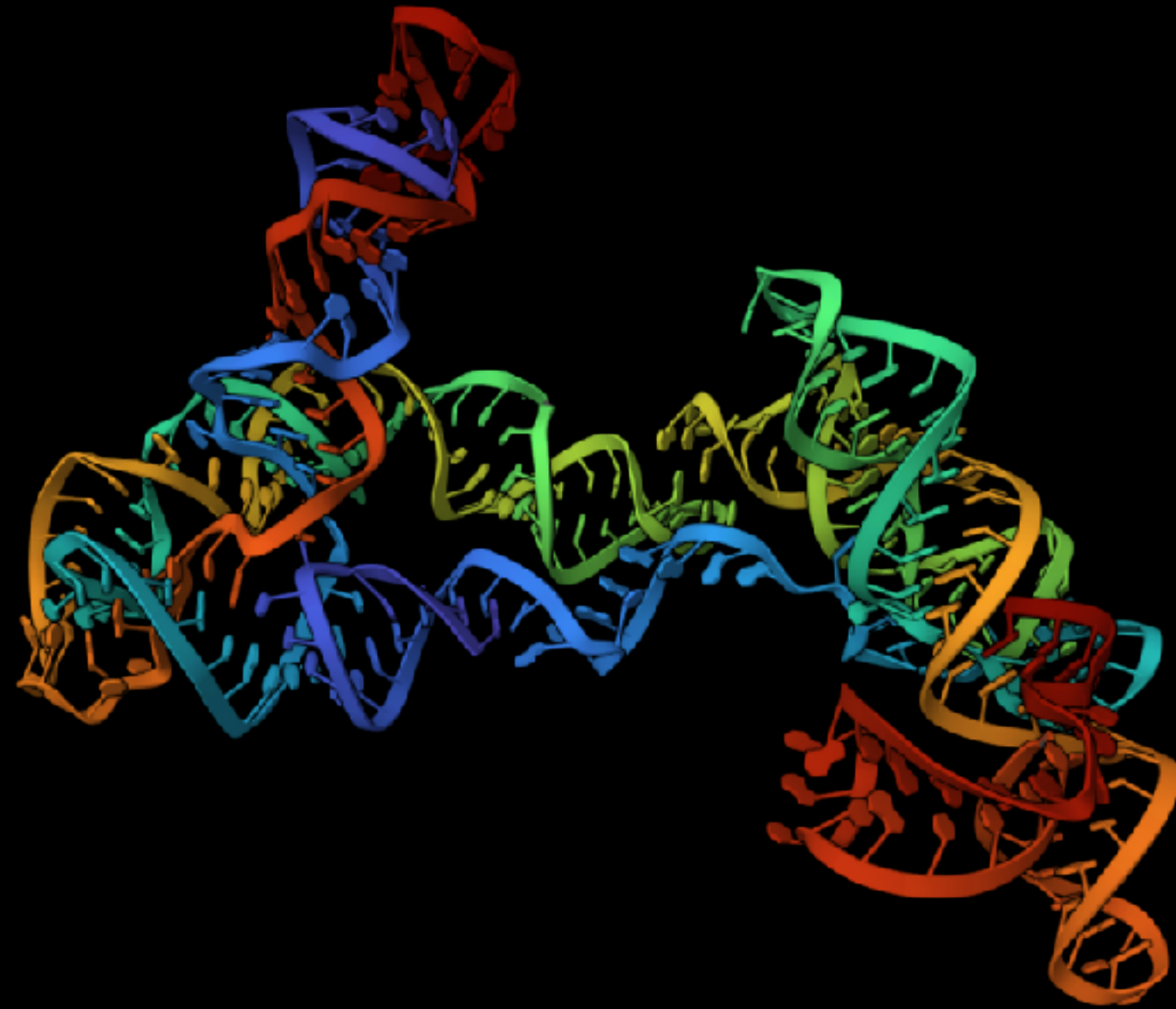
## Inverse design of RNA sequence conditioned on backbone structure



# RNA at the forefront of biotechnology



# And many RNA are structured



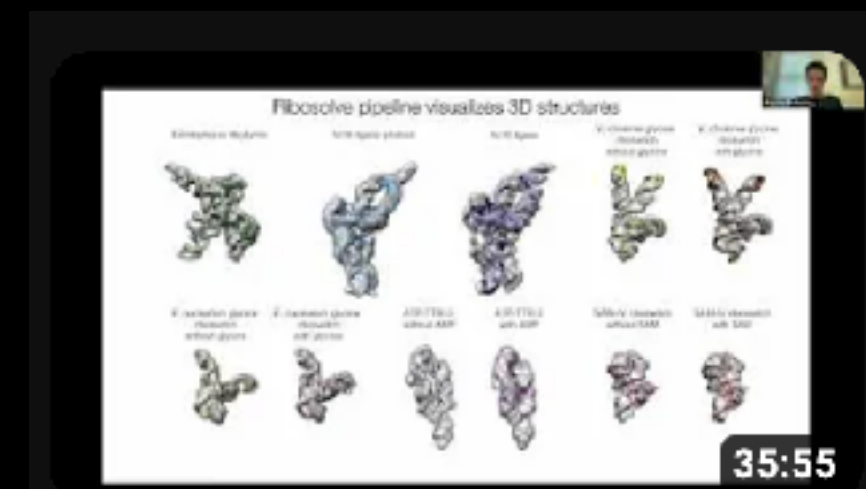
RNA polymerase  
ribozyme  
8T2P  
McRae et al.



SARS-CoV-2  
frameshift  
element  
6XRZ  
Zhang et al.



Adenine  
riboswitch  
aptamer  
5E54  
Stagno et al.

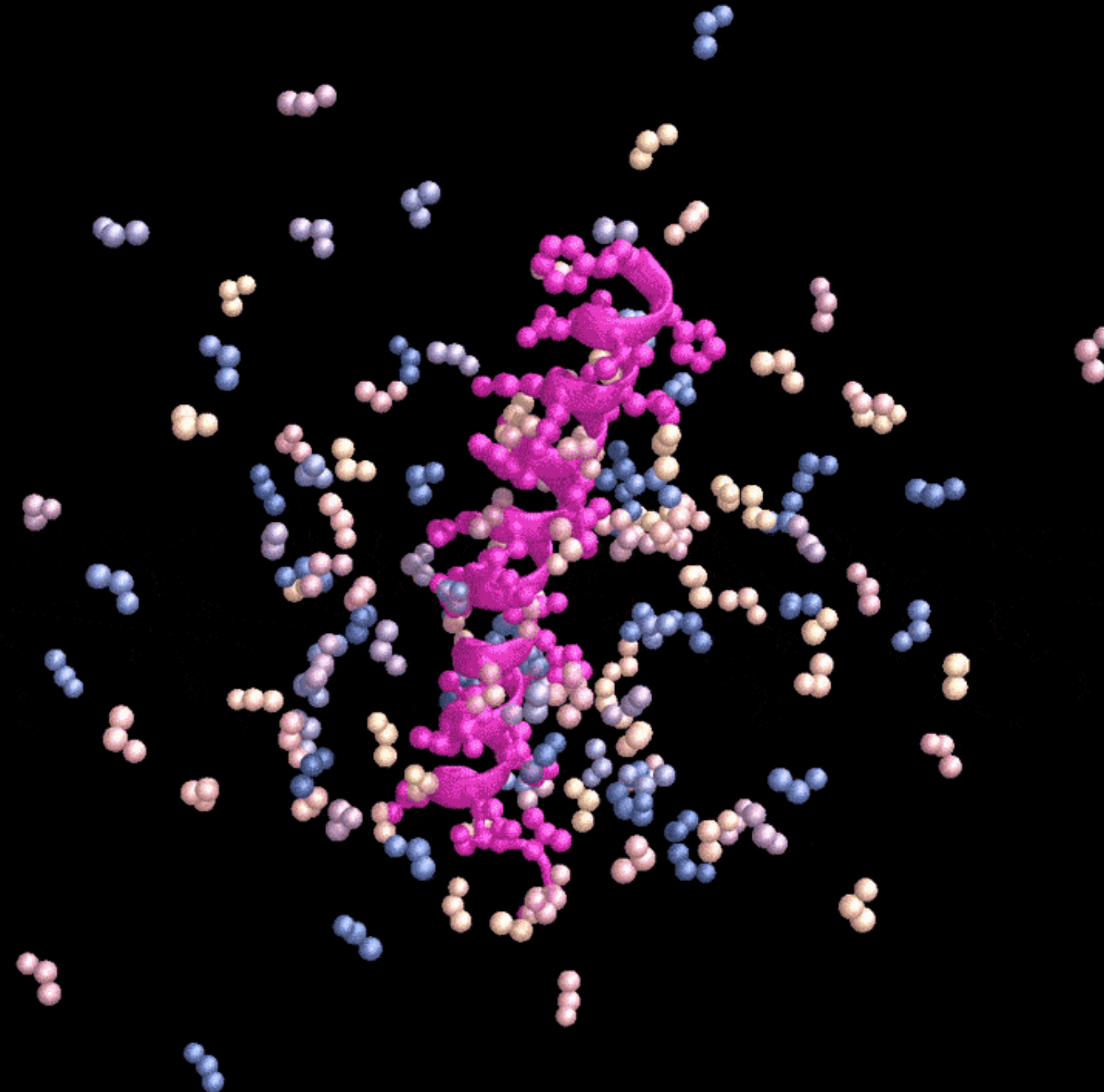


NGBS2022 Talk 10: RNA modelling  
and design - Rhiju Das

466 views · 4 months ago

# Meanwhile

**3D deep learning for protein design is starting to work**

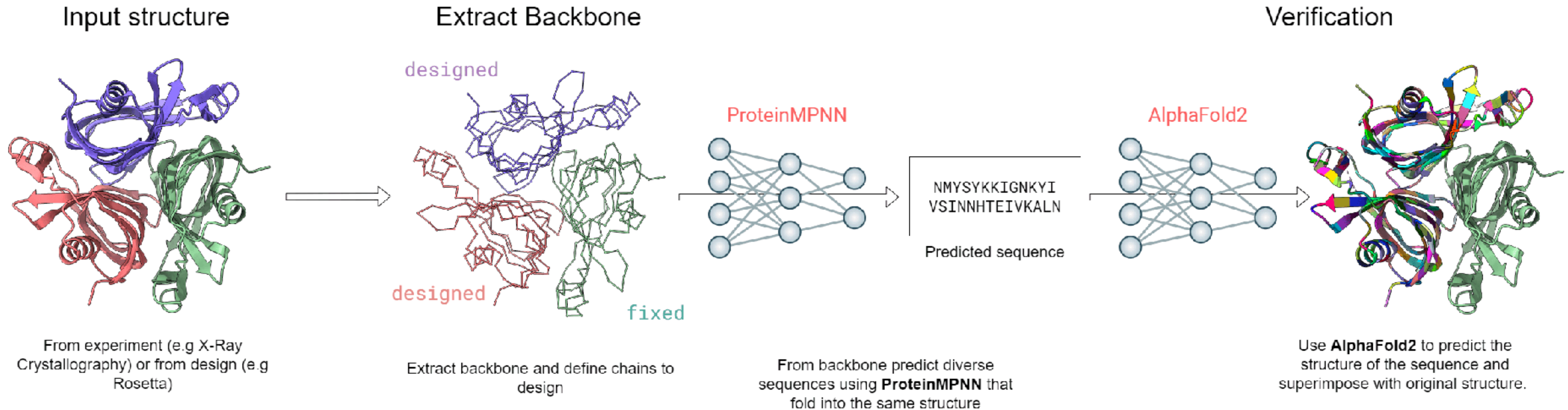


What about  
RNA?

**‘Generative AI’  
is starting to work for protein design**

# Structure-based protein design workflow

Assumption: Structure → Function



Not shown: **protein Language Models** (purely sequence-based)

# Analogy to ChatGPT



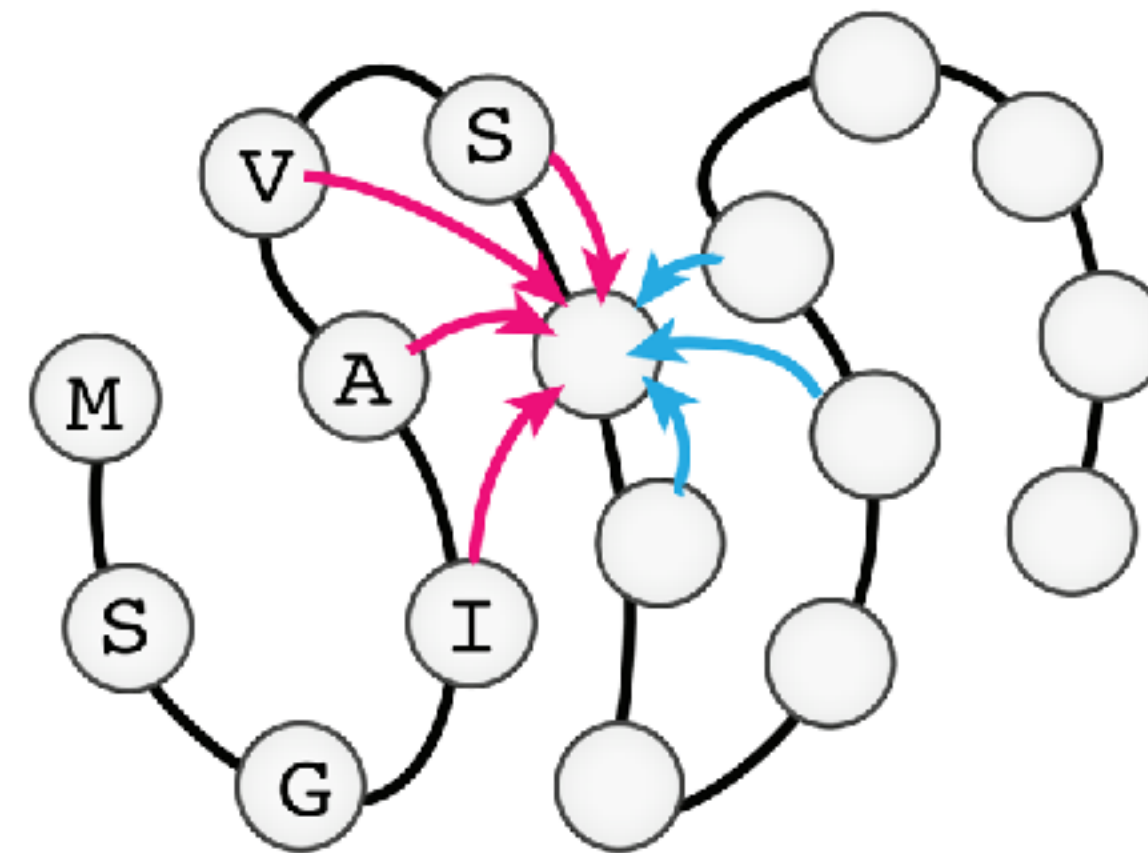
Natural language models

S = Where are we going

Previous words  
(Context)

Word being  
predicted

$$P(S) = P(\text{Where}) \times P(\text{are} \mid \text{Where}) \times P(\text{we} \mid \text{Where are}) \times P(\text{going} \mid \text{Where are we})$$



Trained on PDB structures:  
Samples are biased towards  
thermal stability.



Sequence generation: **Language model**

Sequence generation conditioned on structure: **ProteinMPNN** (inverse folding)



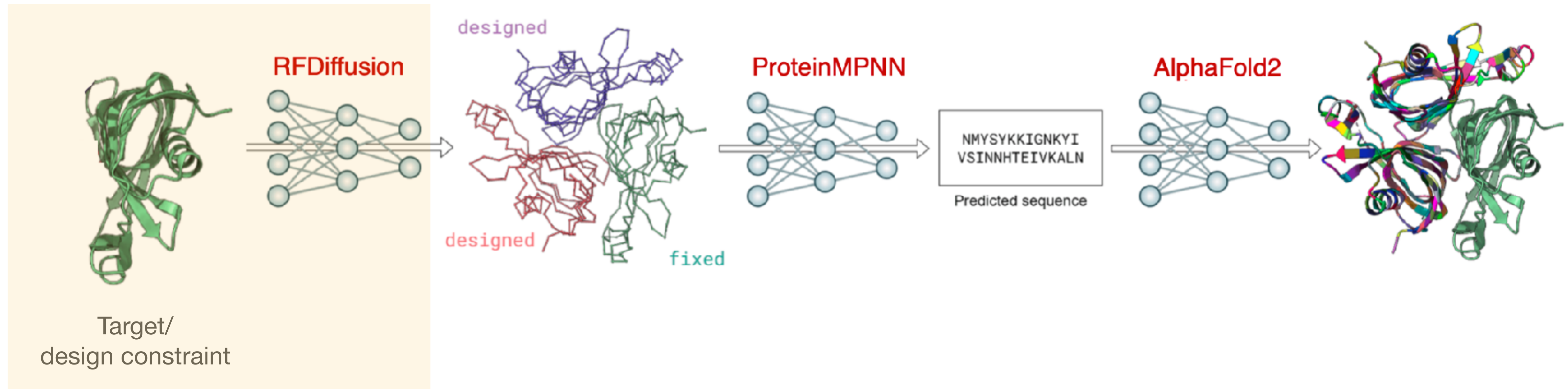
# De-novo protein design workflow

## Starting from scratch

Backbone design

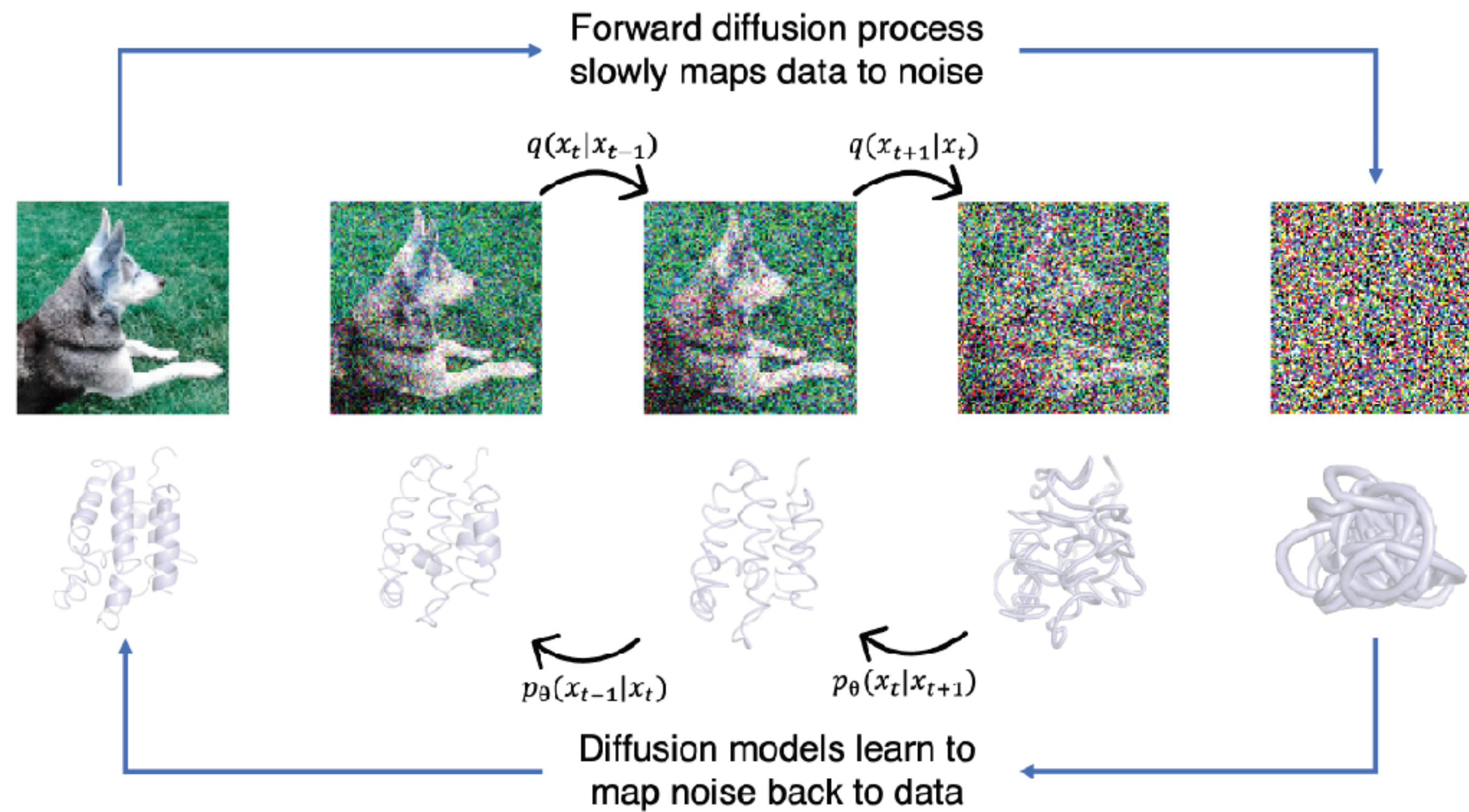
Inverse folding

Verification



# Analogy to DALL-E

DALL-E  
Image generation models

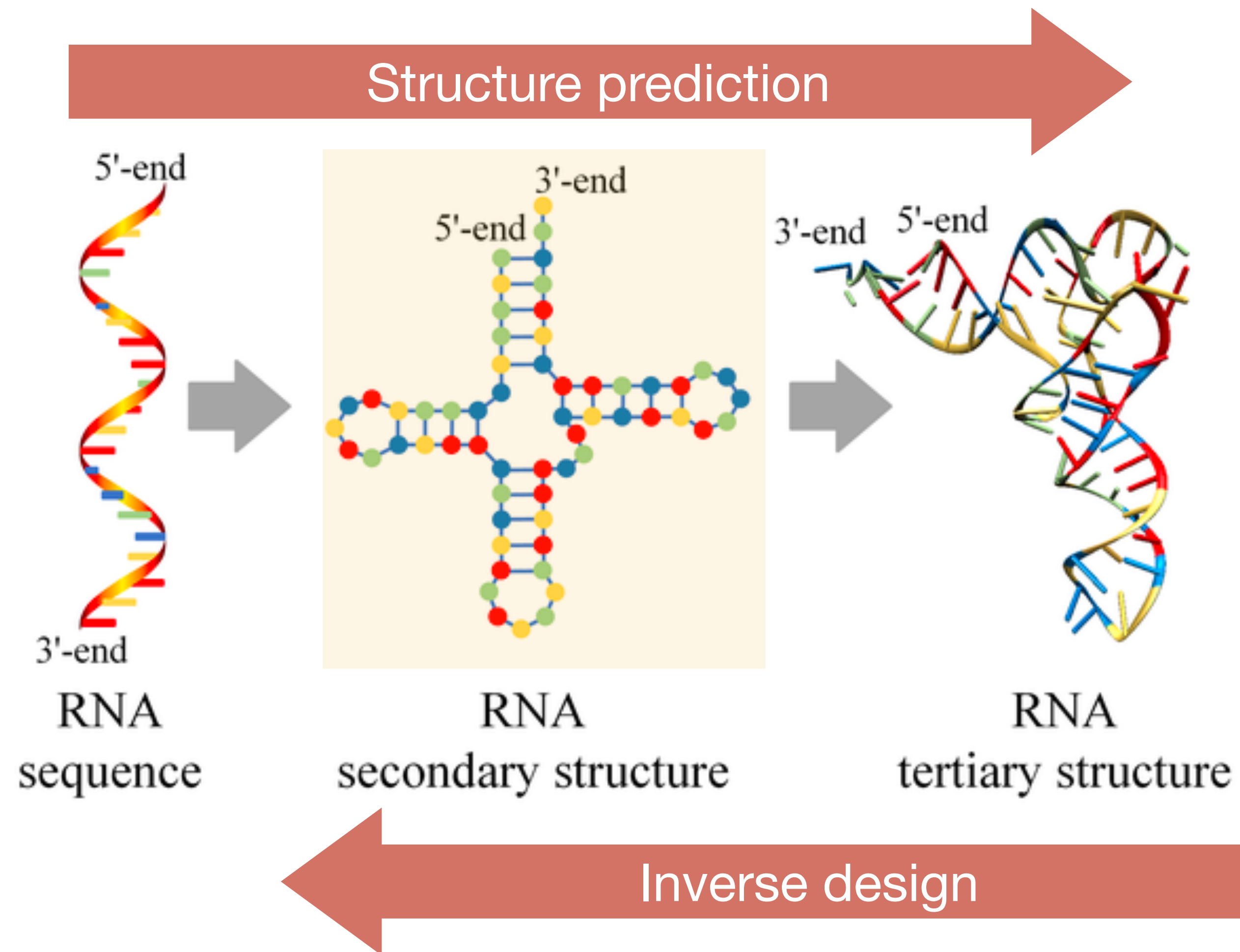


Backbone design: **RFdiffusion**

**What about RNA?**

# RNA structure modelling and design

## Emphasis on secondary structure

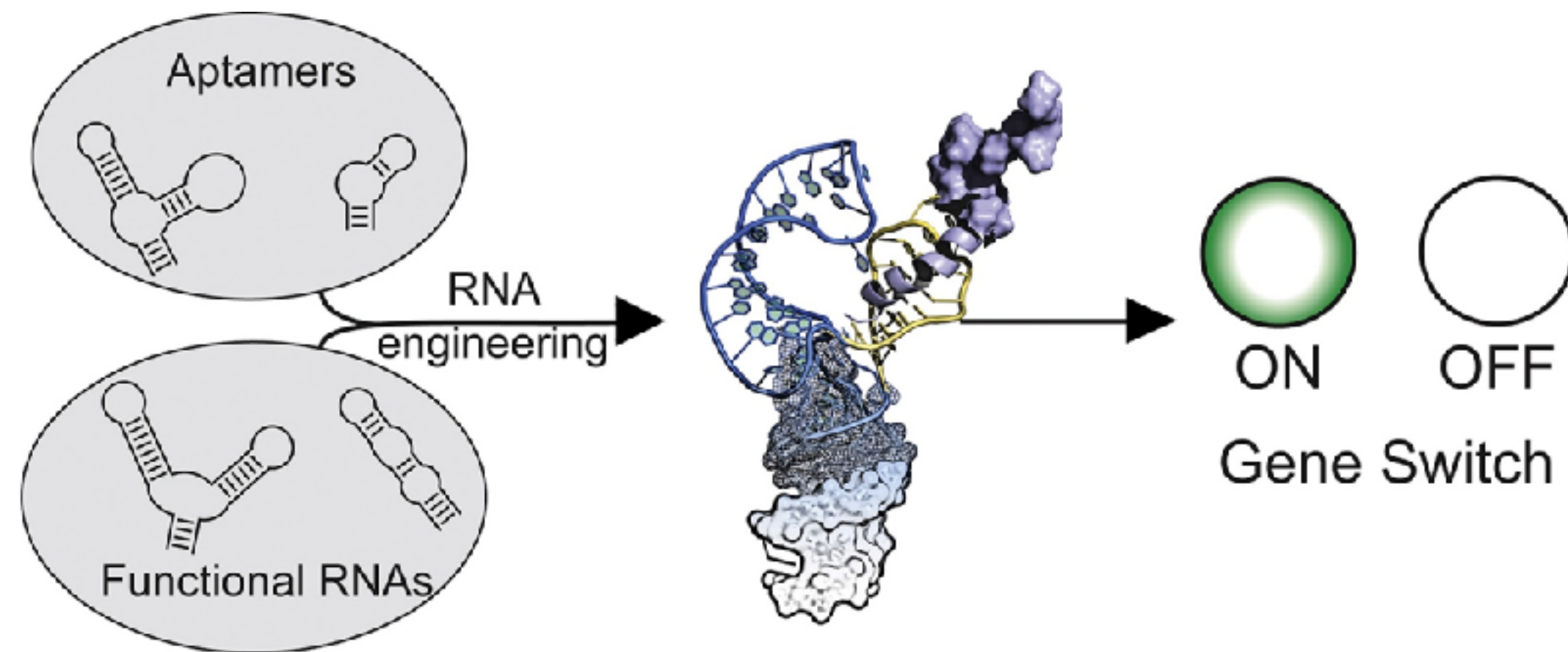


# Relatively fewer tools for 3D design

Potential application: aptamers, riboswitches, ribozymes

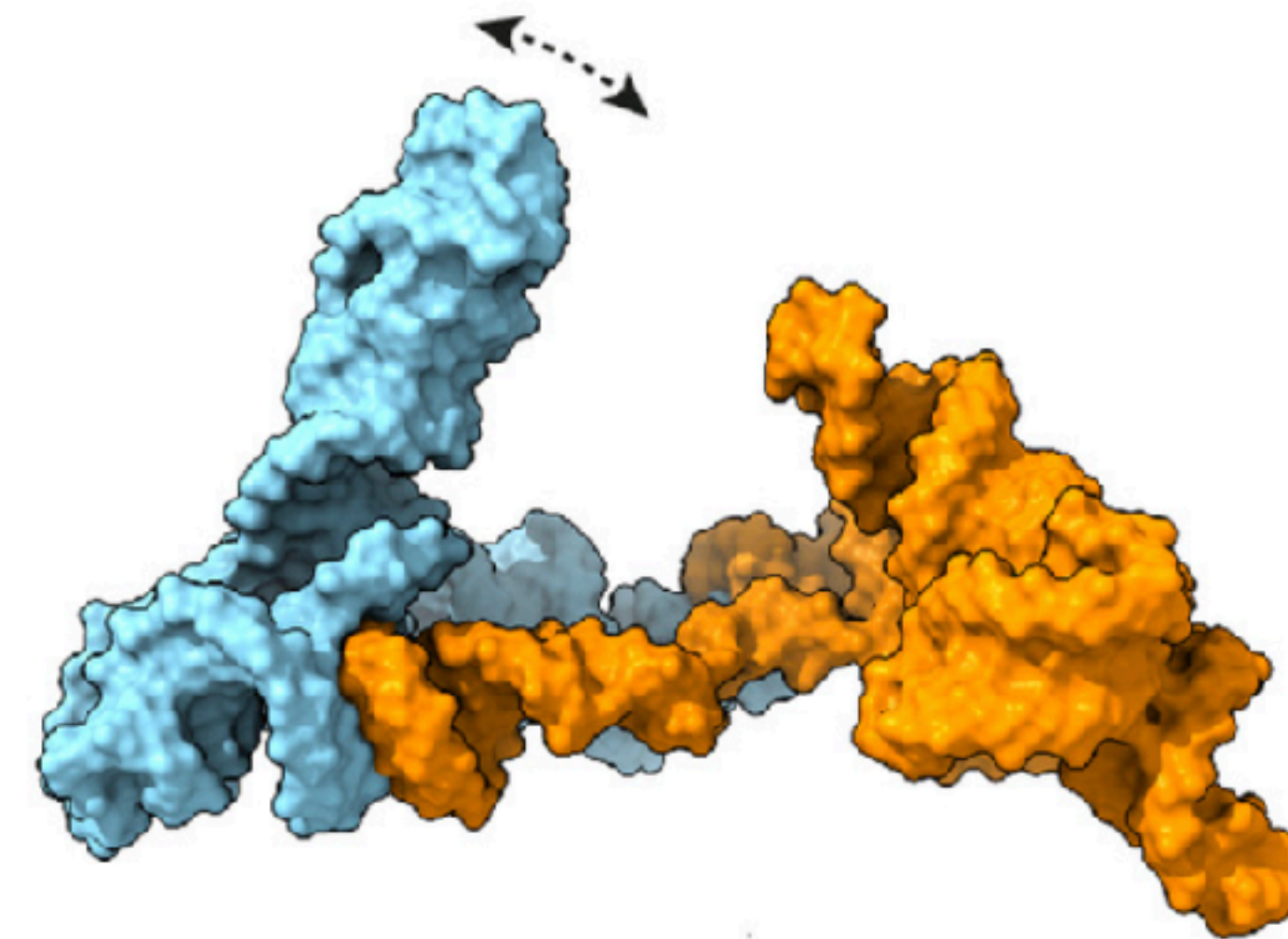
Transient gene expression

Designing riboswitches



RNA world

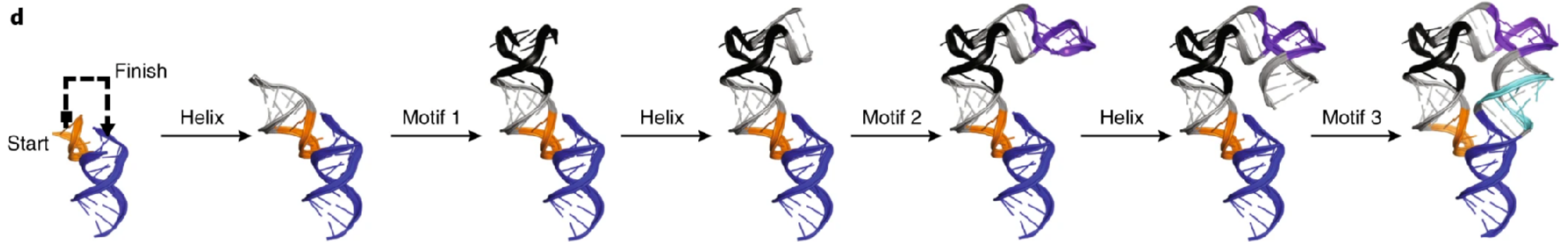
Self-replicating ribozymes



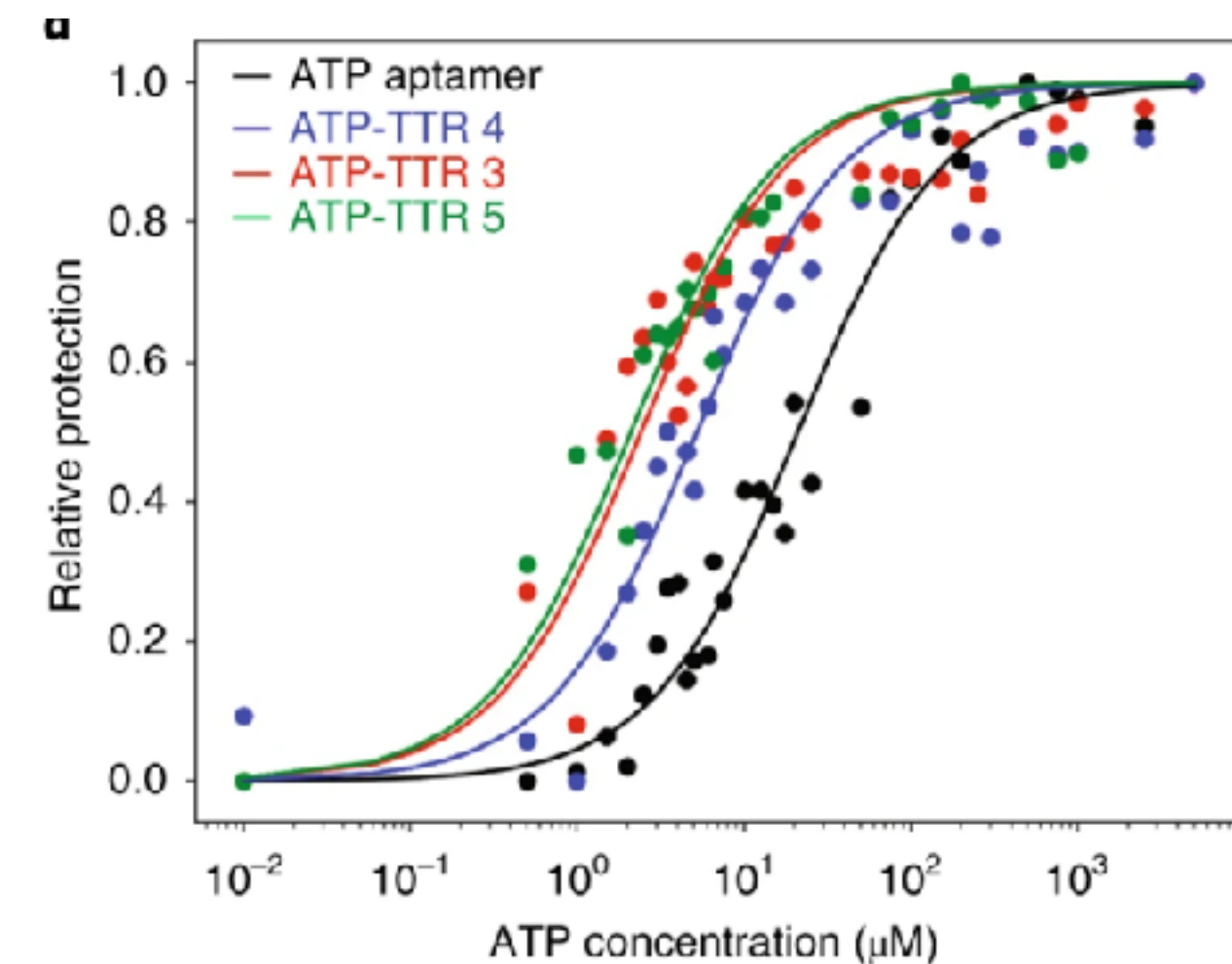
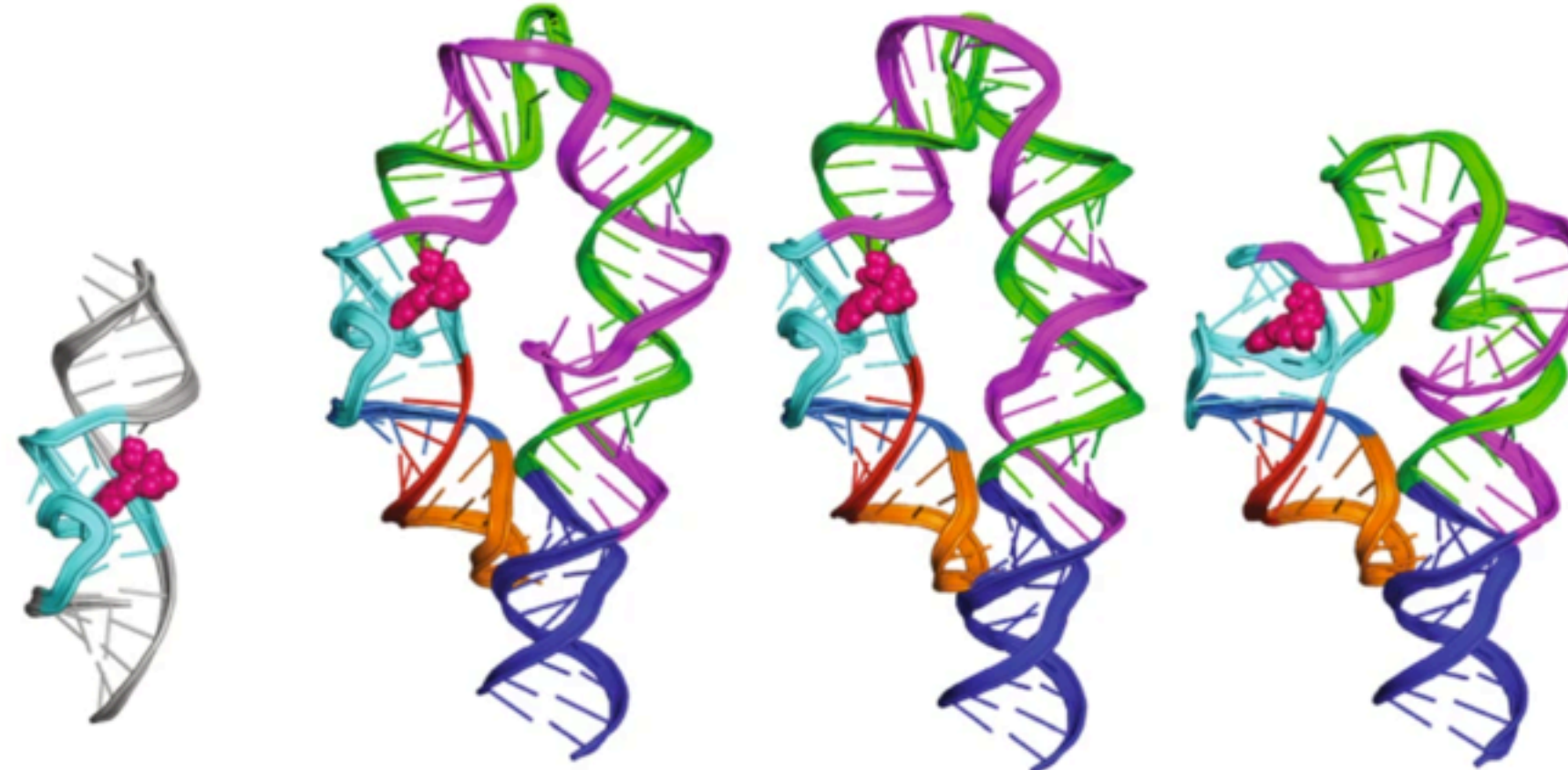
RNAs: carriers of information + play functional roles

# RNAMake

Uses classical algorithms for alignment between RNA motifs

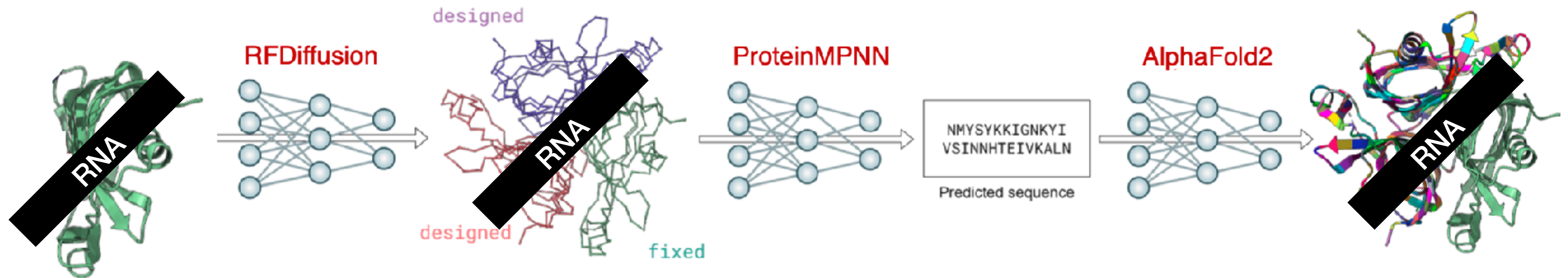


ATP control    ATP-TTR 3    ATP-TTR 4    ATP-TTR 5



# Deep learning toolkit for RNA design

...work in progress



Nothing public yet

RNA Make (non-DL)

**gRNAde**

This talk!

DRFold, RhoFold, RF-NA

Several teams working on this.

Not shown: RNA Language Models — Several teams working on this.

eg. RiNaLMo

**Towards deep learning:  
What data exists?**



# Geometric Deep Learning for RNA

**Main challenge: paucity of 3D structural data**

“trained with only 18 known RNA structures”

**ARES: Geometric deep learning of RNA structure.** *Science*, 2021.

Raphael JL Townshend, Stephan Eismann, Andrew M Watkins, Ramya Rangan, Maria Karelina, Rhiju Das, and Ron O Dror.

“trained on 2,986 RNA chains, non-redundant to 122 test RNAs”

**DRFold: Integrating end-to-end learning with deep geometrical potentials for ab initio RNA structure prediction.** *Nature Communications*, 2023.

Yang Li, Chengxin Zhang, Chenjie Feng, Robin Pearce, Peter L. Freddolino, Yang Zhang.

# All RNA structures in the PDB

## RNAsoLO: cleaned, PDB-derived RNA 3D structures

	Solo RNAs	RNAs from protein-RNA complexes	RNAs from DNA-RNA hybrids	All RNAs
X-ray	1454	6439	91	7984
NMR	573	146	28	747
Electron microscopy	73	4104	0	4177
Multi-method	1	5	0	6
Total	2101	10694	119	12914

Total (today)      2387      13218      136      15741 (13870  $\leq 3.5\text{\AA}$ )

# All RNA structures in the PDB

## RNA solo: cleaned, PDB-derived RNA 3D structures

	Solo RNAs	RNAs from protein-RNA complexes	RNAs from DNA-RNA hybrids	All RNAs
Total (today)	2387	13218	136	15741

3825 equivalence classes

vs.

ProteinMPNN, RFdiffusion: entire PDB  
208,659 proteins  $\leq 3.5\text{\AA}$   $\rightarrow$  25,361 clusters at 30% seq.id.

One order of magnitude more proteins!

# Should we just wait?

**Not necessarily...**

Other successful (in-silico) tools were trained on carefully chosen subsets:

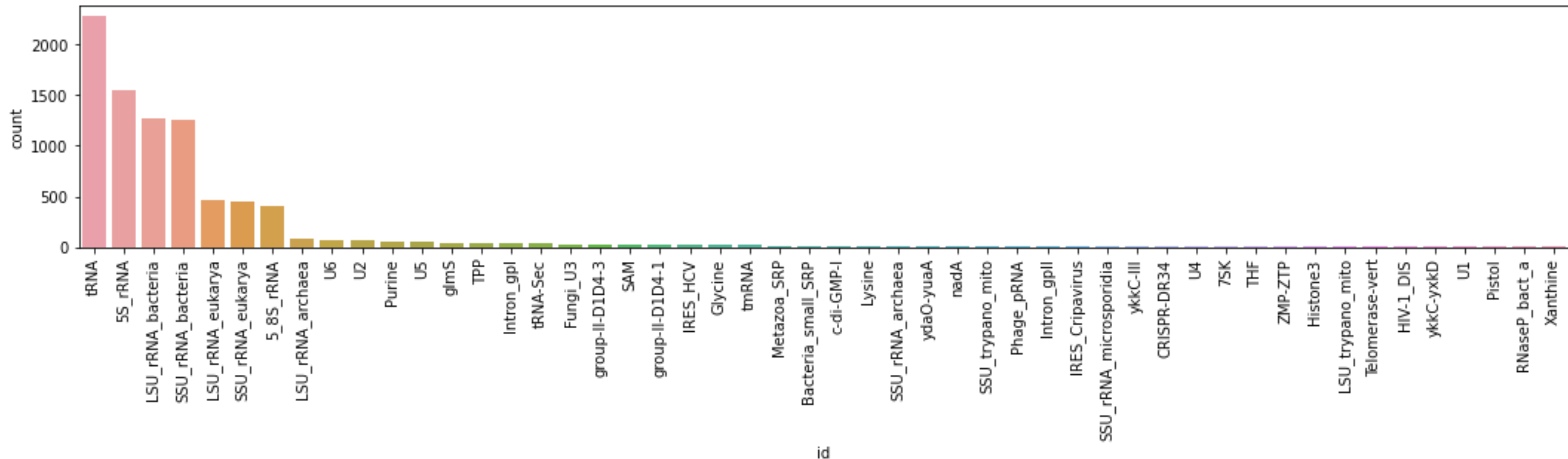
- Chroma: 28819 structures  $\leq 2.6\text{\AA}$
- Genie: 8766 domains
- FrameFlow: 3938 domains

“...achieve similar in-silico performance to RFdiffusion with a quarter of the parameters – an important consideration...models are often run tens of thousands of times...”

– *Winnifrith et al. 2023.*

# RFam families in the PDB

Majority from protein-RNA complexes, tRNAs, ribosomal RNAs



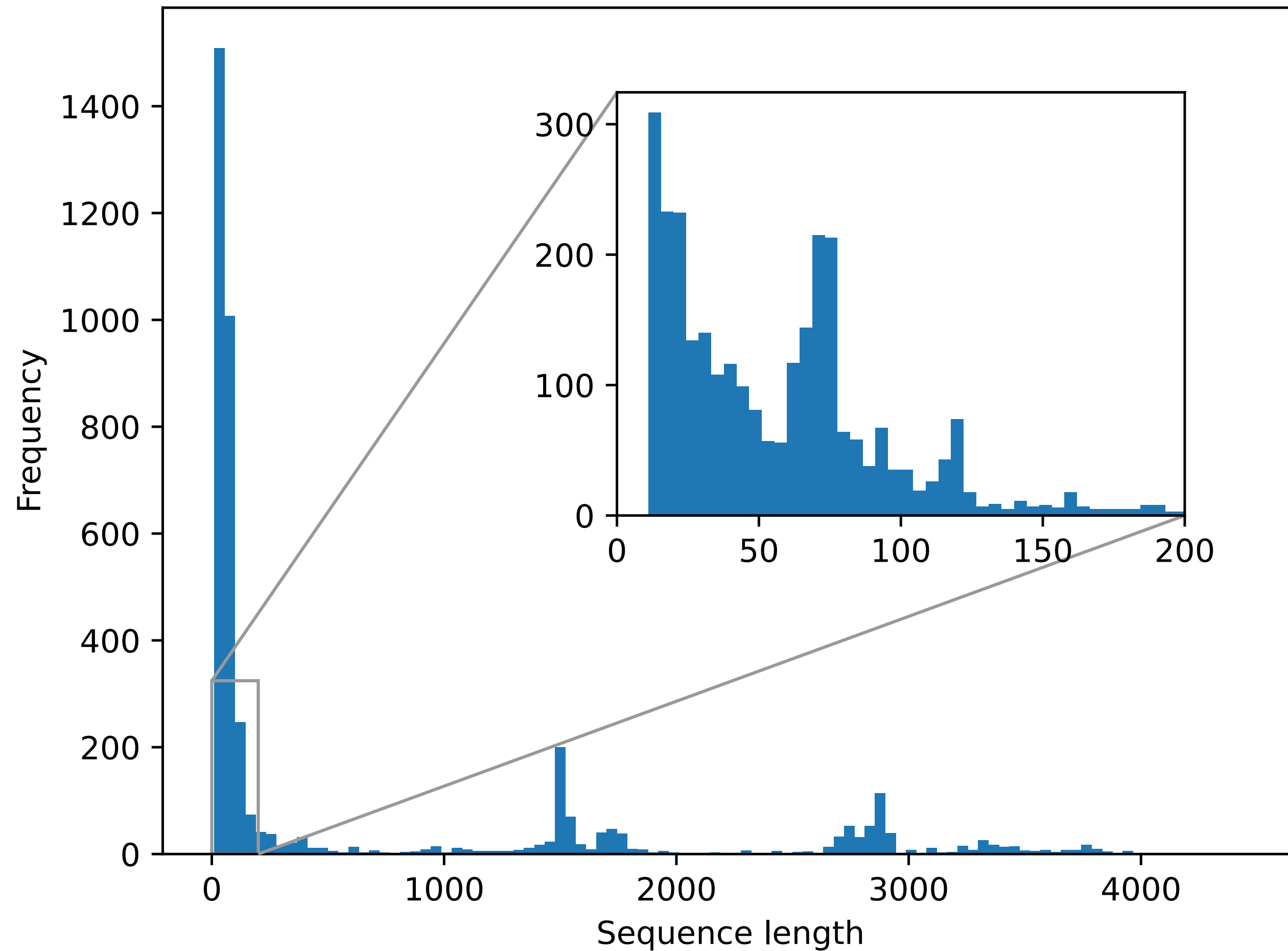
**Idea:** Different global folds, but locally they are all RNAs!  
(in the sense of local structural interactions and sub-substructures)

# Distribution of sequence lengths

Mostly shorter than 500 nucleotides

**Histogram of sequence lengths**

Distribution:  $684.9 \pm 1072.8$ , Max: 4455, Min: 11

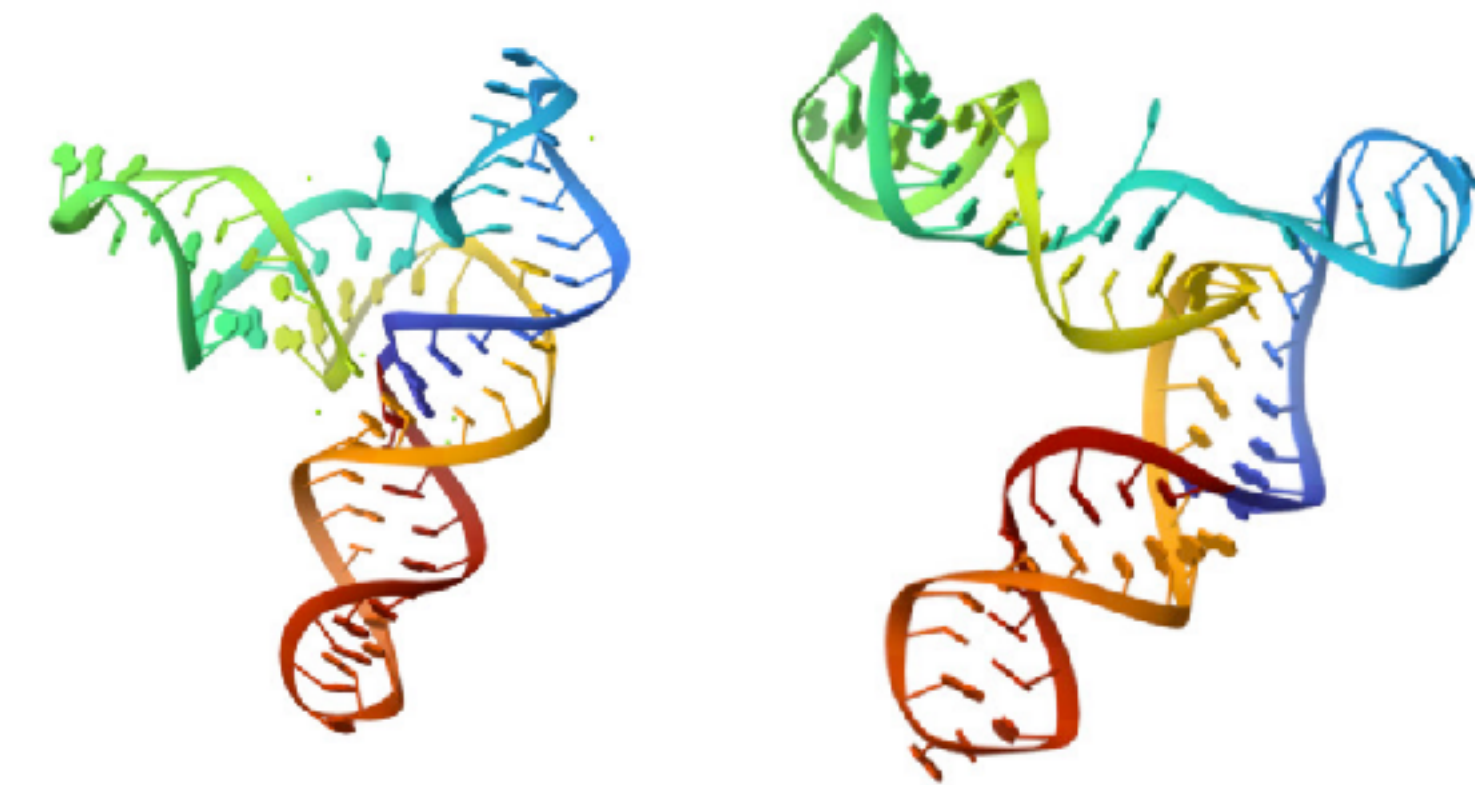
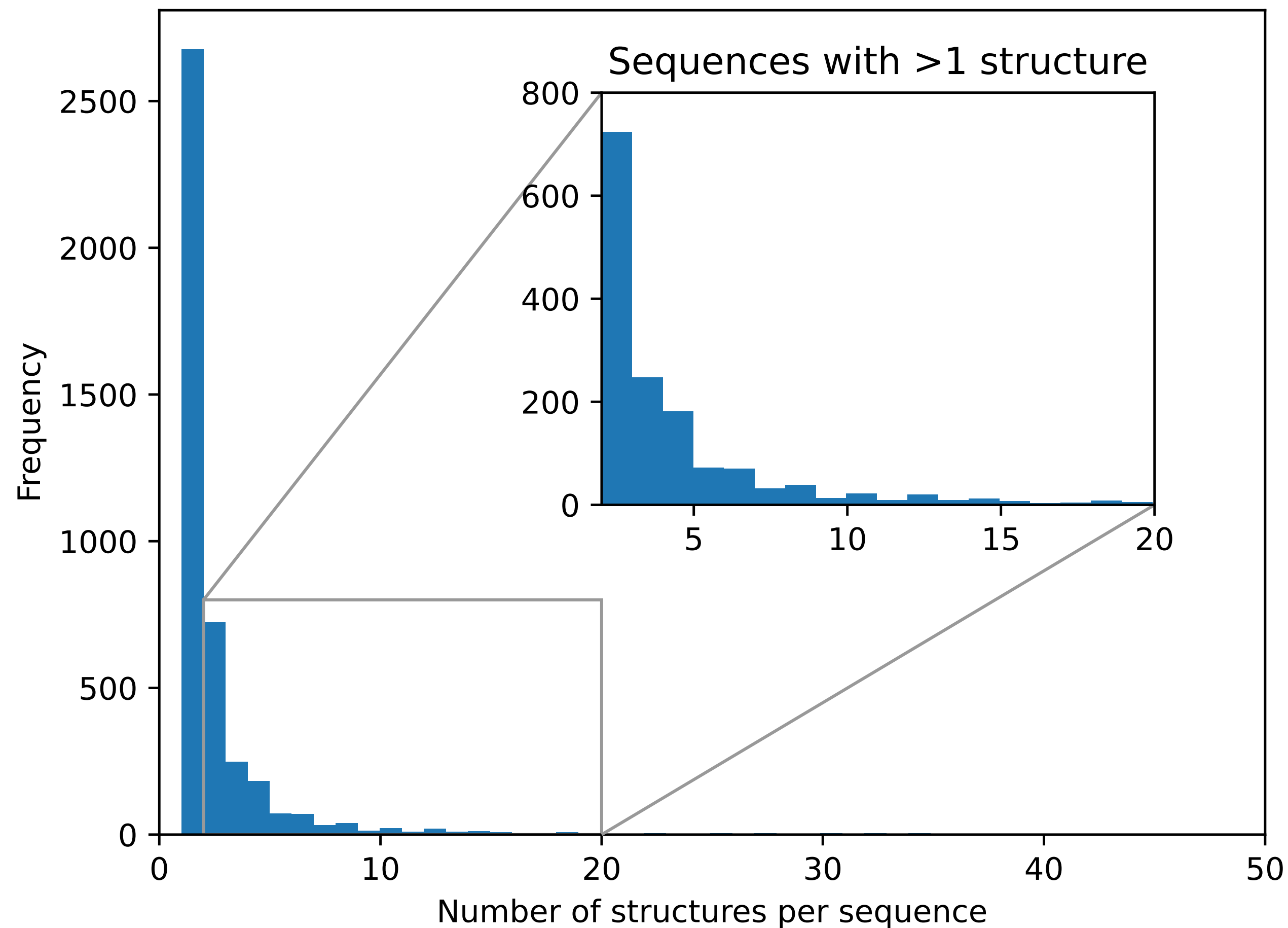


# RNA adopt multiple conformations

Critical for functionality & perhaps interesting for design

Histogram of no. of structures per unique sequence

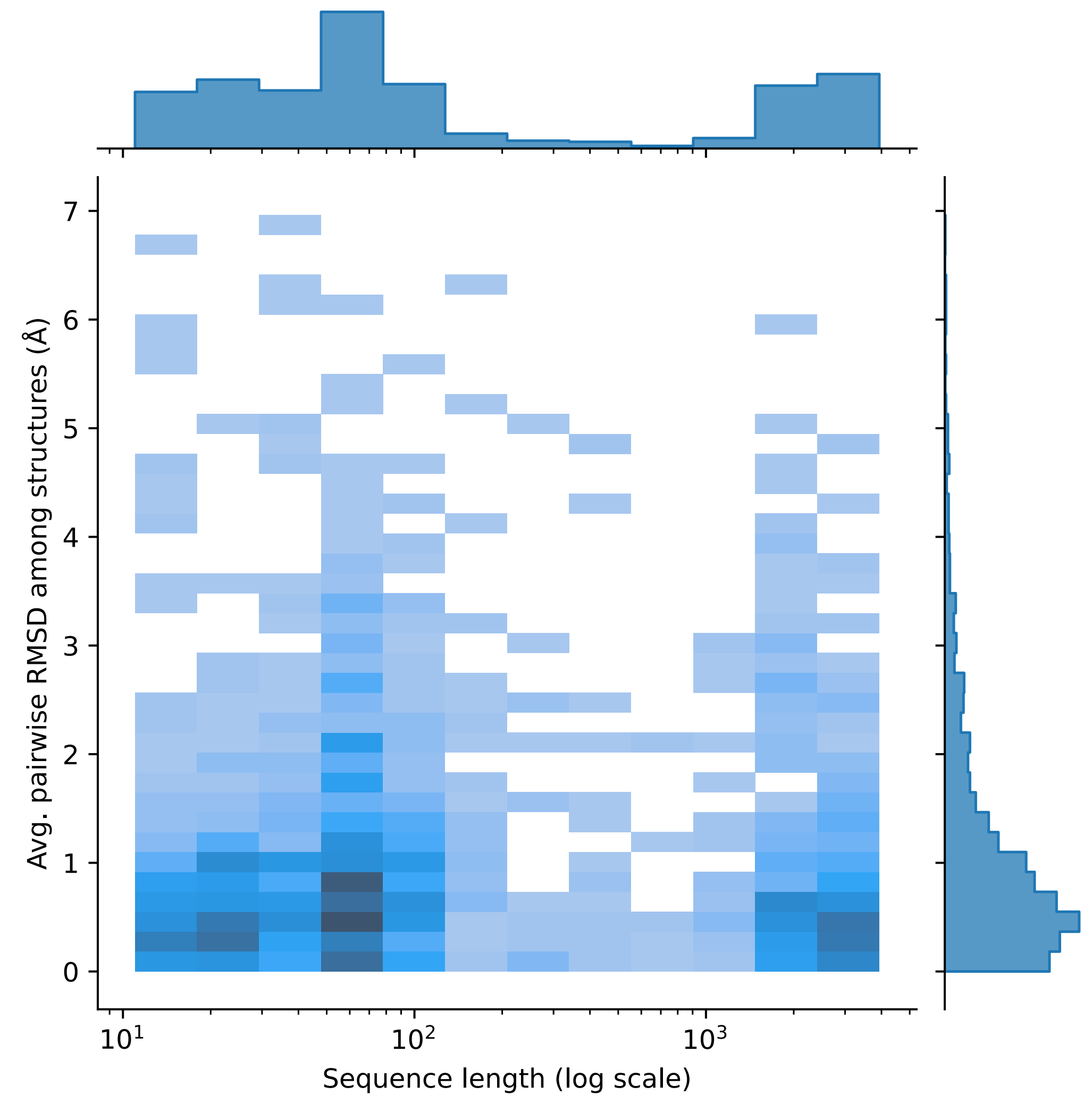
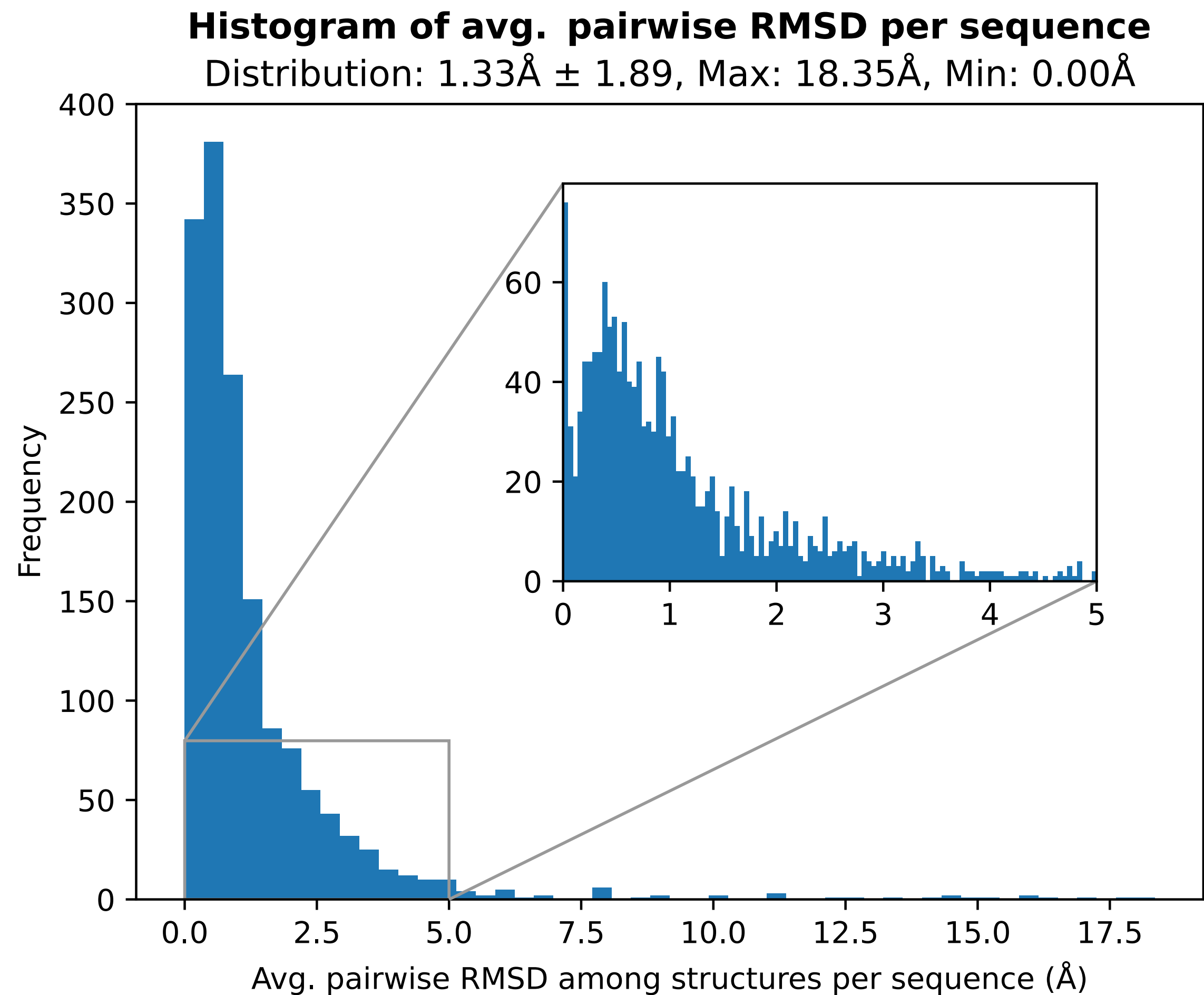
Distribution:  $2.84 \pm 9.39$ , Max: 267, Min: 1



L1 ligase ribozyme:  
~15Å  
(PDB 2OIU)

# RNA adopt multiple conformations

Same sequence can have very different structures

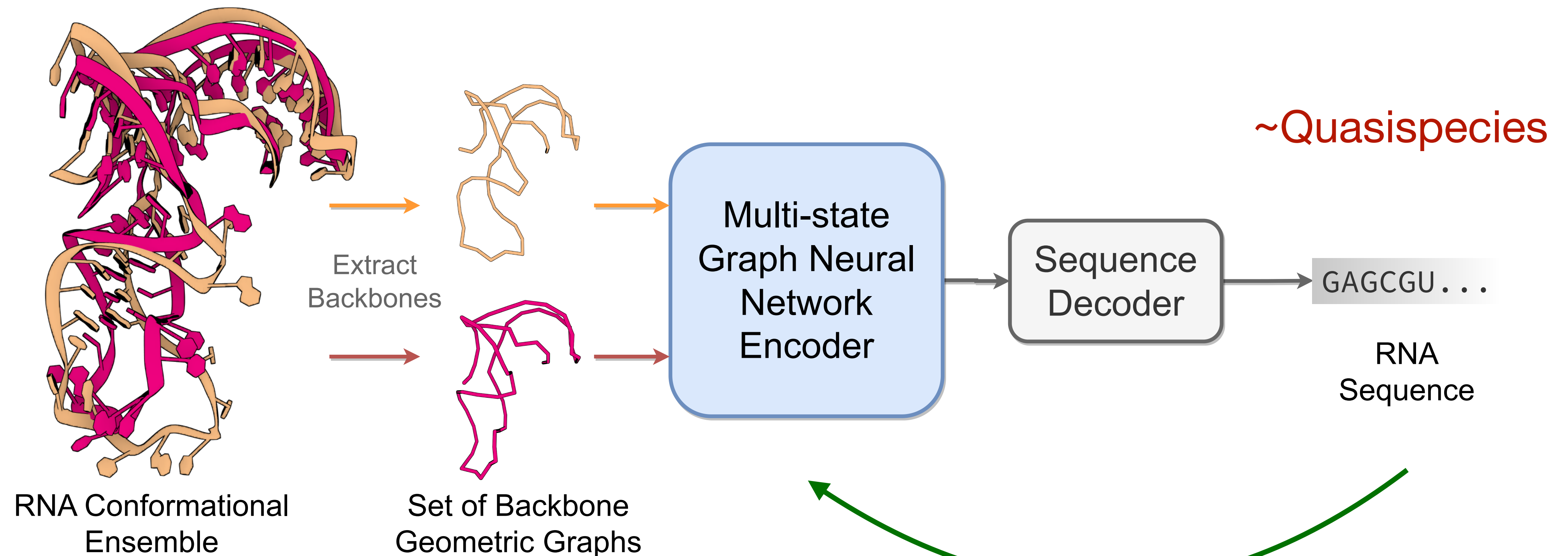




# The gRNAde pipeline for RNA inverse folding

# Fixed backbone re-design

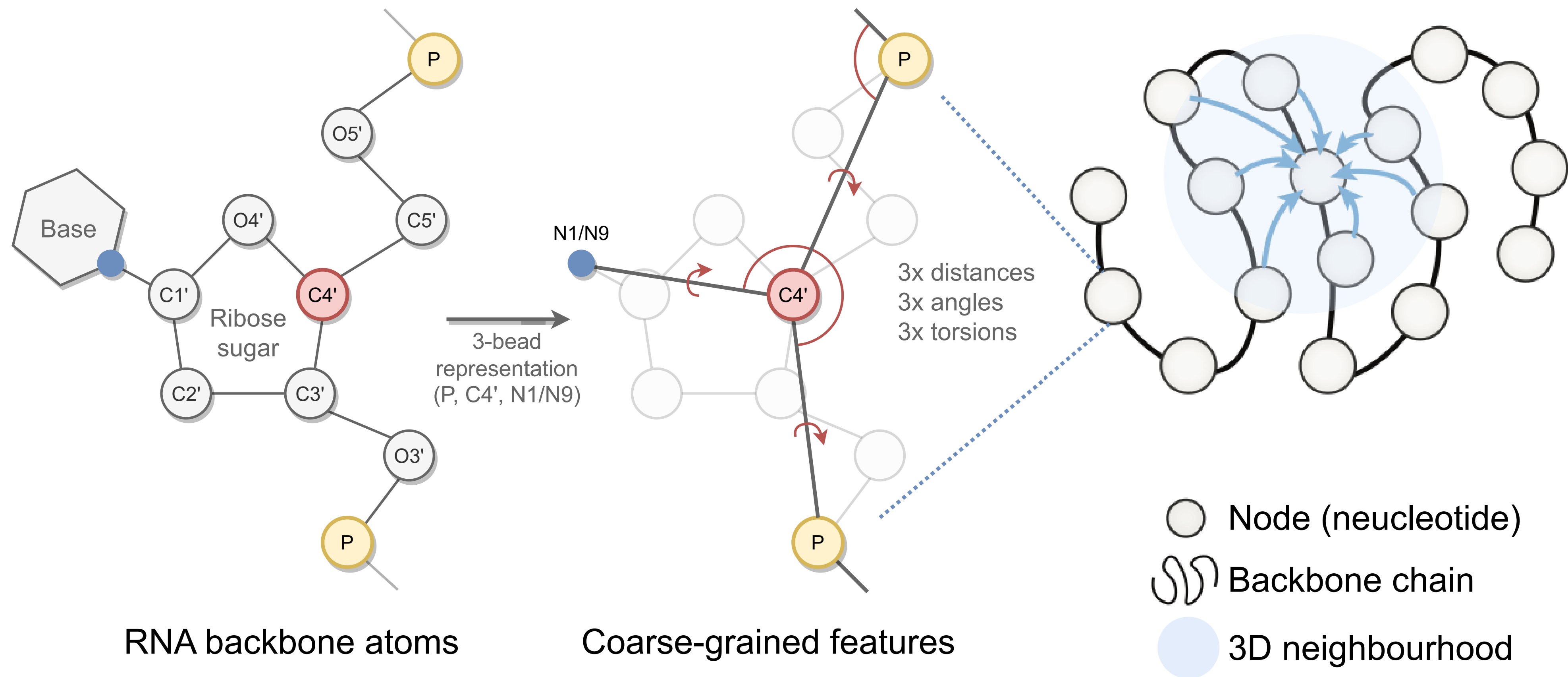
Input: native PDB file → Output: designed sequences



**Self-supervised learning:**  
(backbone, sequence) pairs from RNAsolo

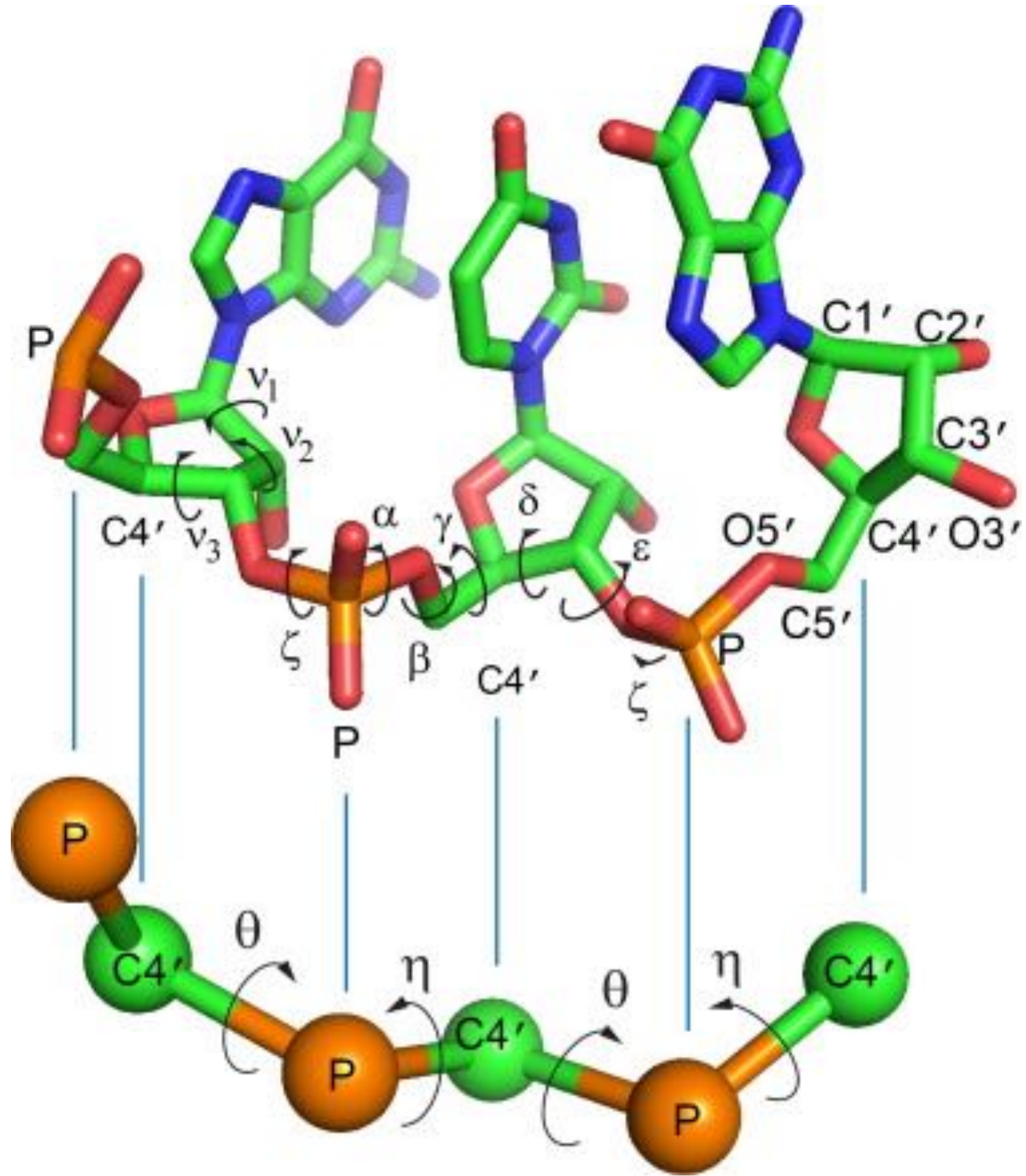
# RNA backbones as 3D graphs

Preparing input: PDB file(s) → geometric graph in 3D



# Why the 3-bead representation?

**P, C4', N1 (pyrimidine) or N9 (purine)**

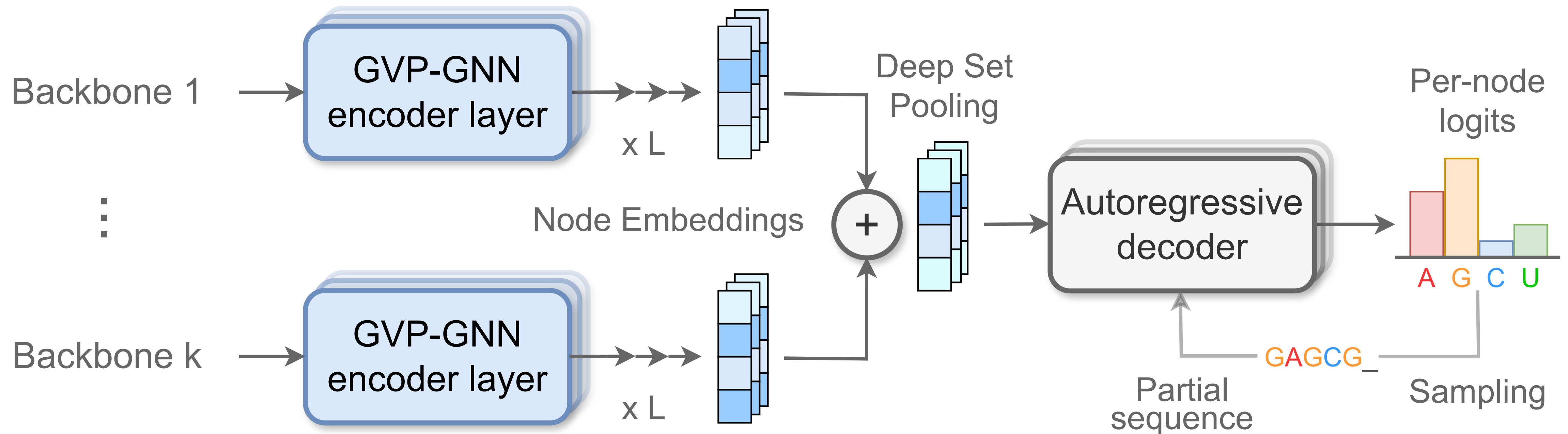


Intuition: Reduce the degrees of freedom as input to gRNAde.

*“The pseudotorsional descriptors  $\eta$  and  $\theta$ , together with sugar pucker, **are sufficient to describe RNA backbone conformations fully in most cases.**”*

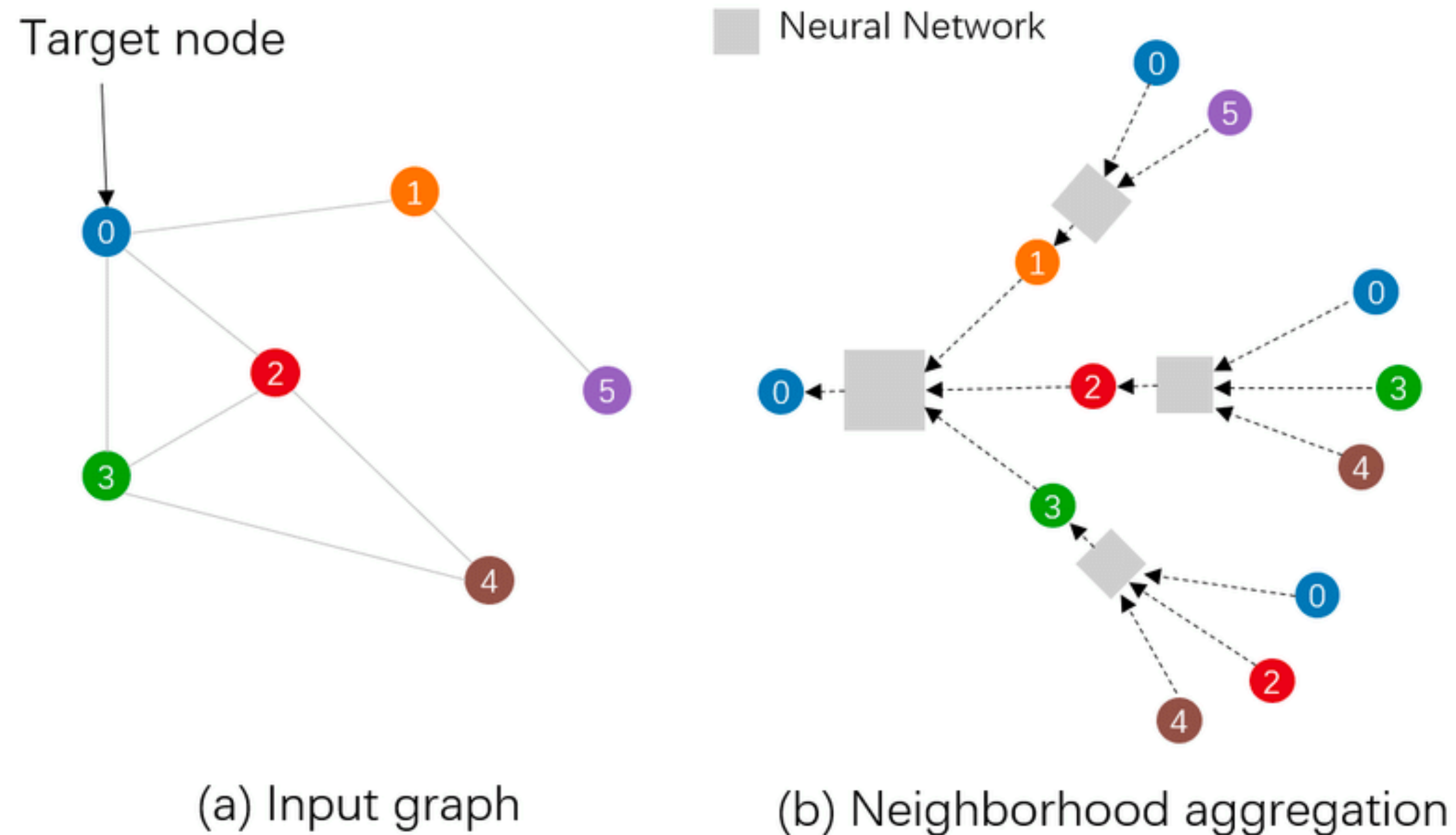
# gRNAde model architecture

One or more featurized graphs  $\rightarrow$  per-node probability over 4 bases

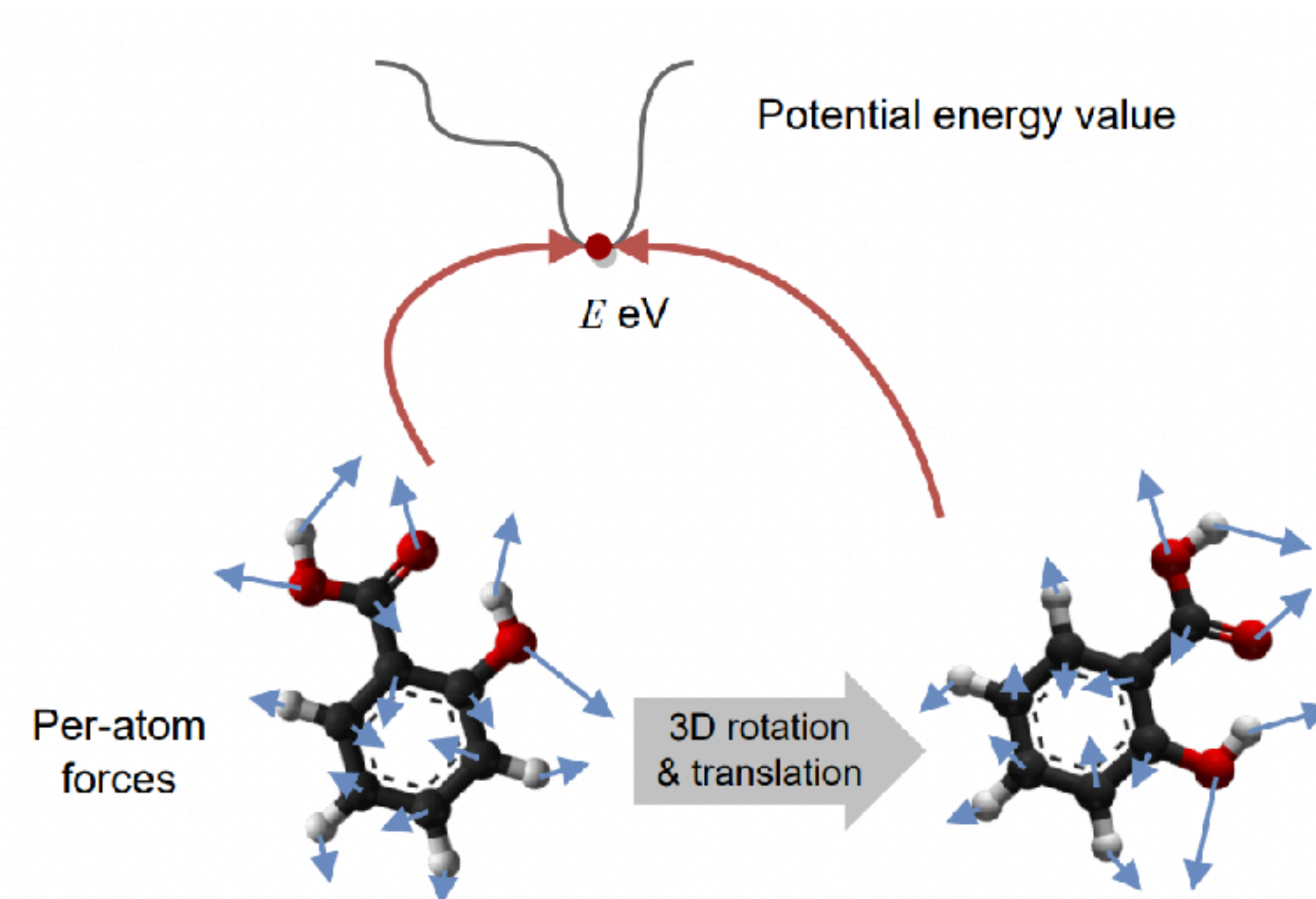


# Graph Neural Networks for 3D structure

Learn to propagate information along the graph

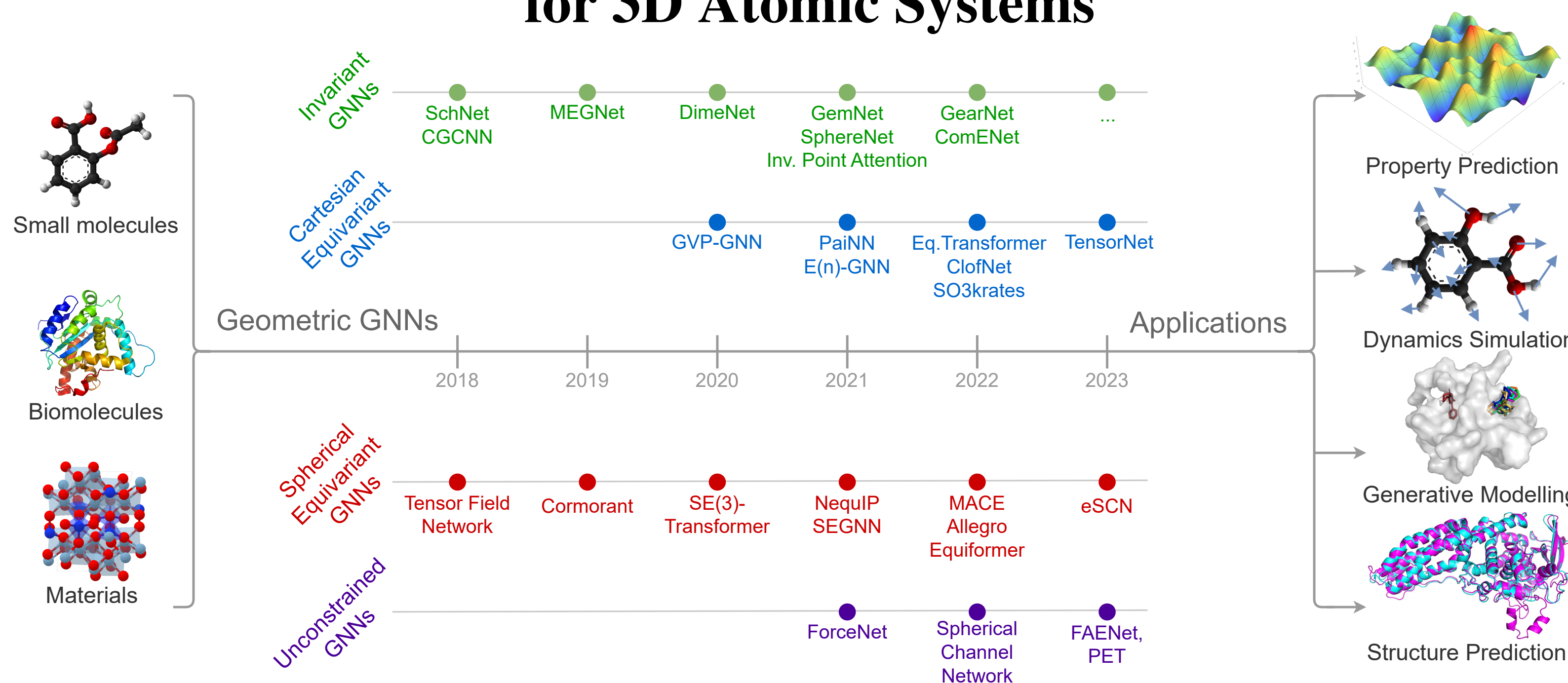


Account for 3D symmetries



# Where to start?

## A Hitchhiker's Guide to Geometric GNNs for 3D Atomic Systems



Alexandre Duval<sup>\*,1,2</sup> Simon V. Mathis<sup>\*,3</sup> Chaitanya K. Joshi<sup>\*,3</sup> Victor Schmidt<sup>\*,1,4</sup>

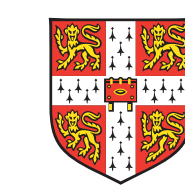
Santiago Miret<sup>5</sup> Fragkiskos D. Malliaros<sup>2</sup> Taco Cohen<sup>6</sup>

Pietro Liò<sup>3</sup> Yoshua Bengio<sup>1,4</sup> Michael Bronstein<sup>7</sup>

<sup>1</sup>Mila <sup>2</sup>Université Paris-Saclay <sup>3</sup>University of Cambridge <sup>4</sup>Université de Montréal

<sup>5</sup>Intel Labs <sup>6</sup>Qualcomm AI Research <sup>7</sup>University of Oxford

\*Equal first authors.



UNIVERSITY OF CAMBRIDGE



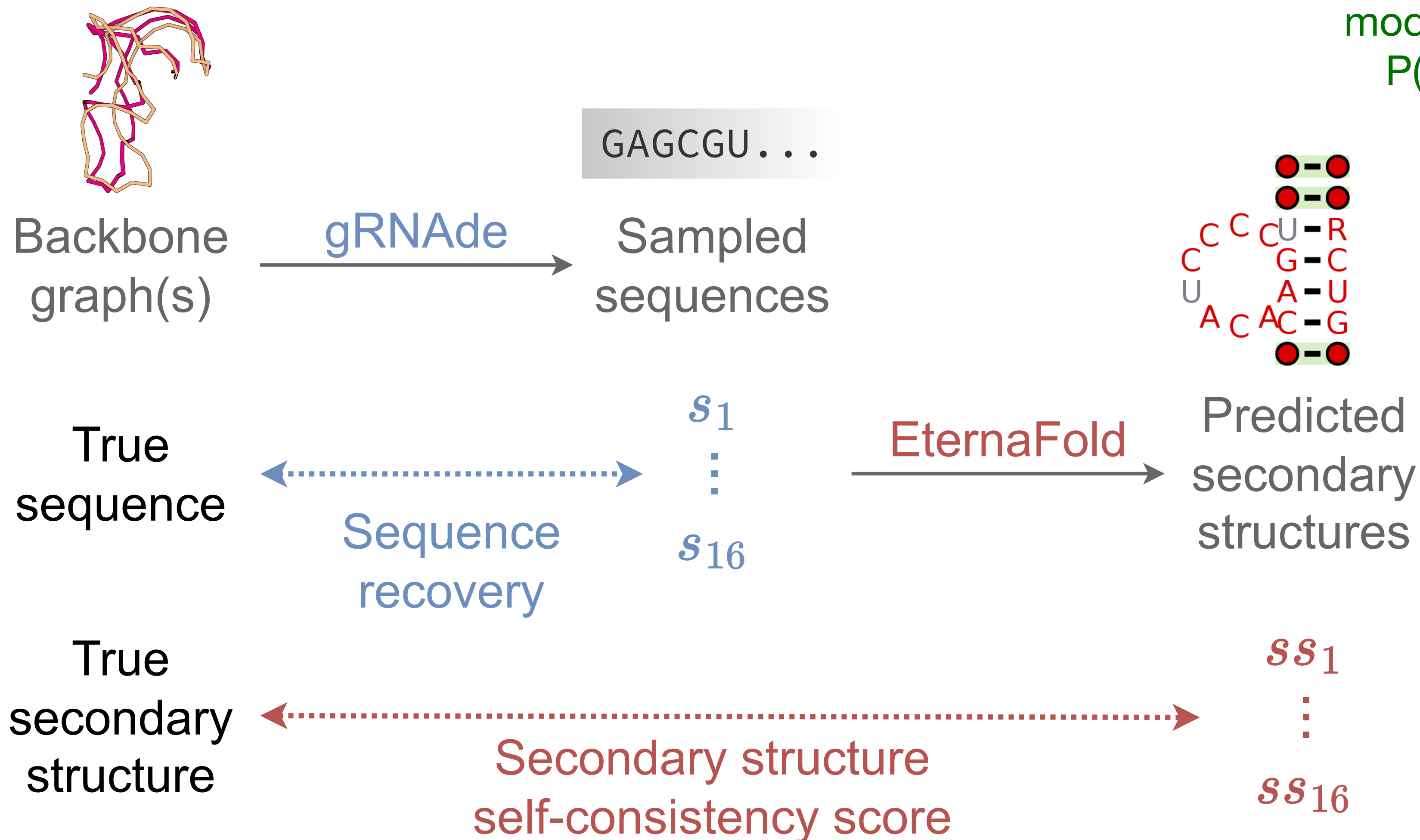
UNIVERSITY OF OXFORD



# What is a good designs?

## In-silico evaluation metrics to prioritise designs

Not shown:  
**Perplexity**  
model's guess of  
 $P(\text{seq}|\text{struct})$



3D would be better!

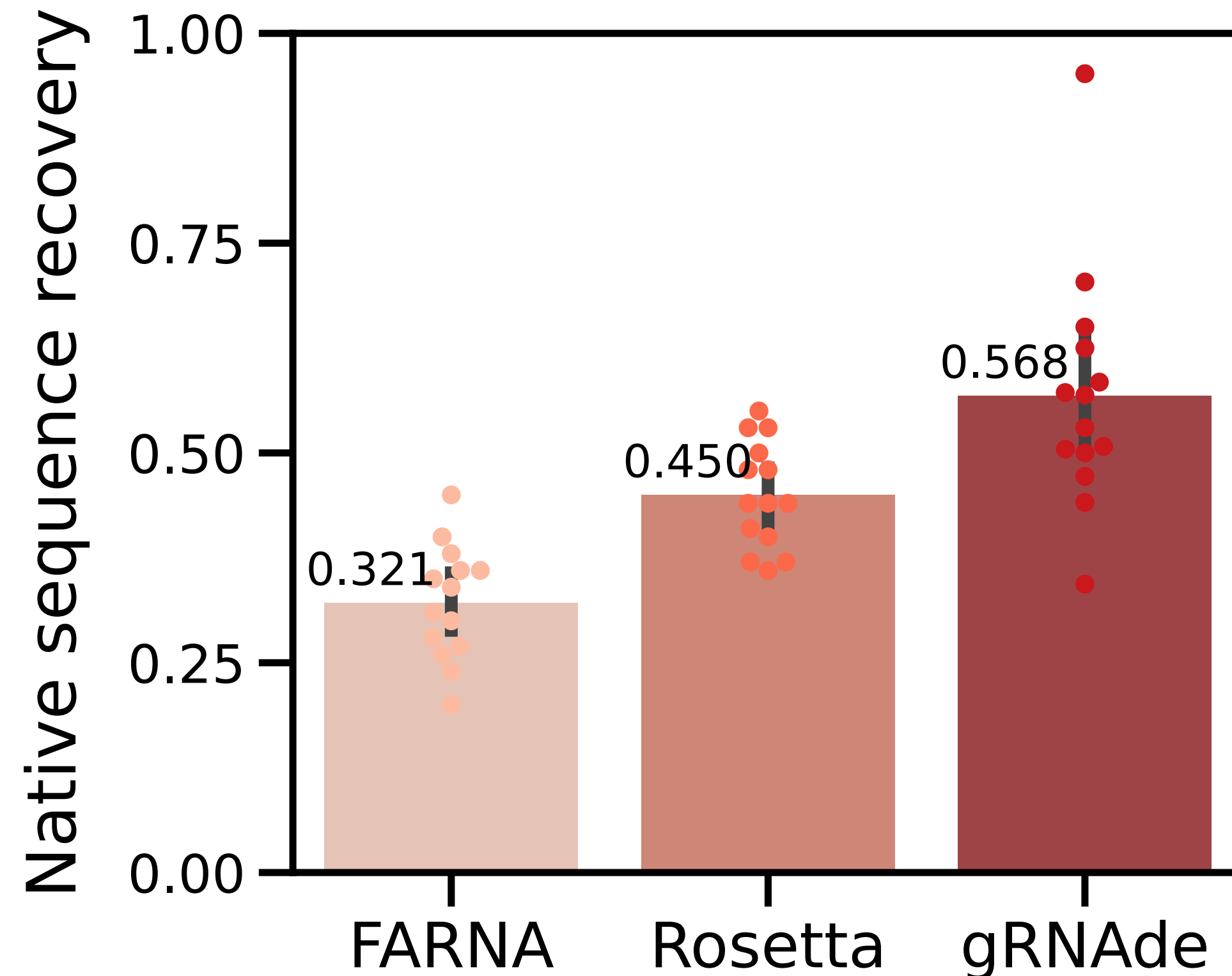


**What can we do with gRNAde?**

# Benchmarking single-state design

Re-design 14 RNAs of interest from the PDB by Das et al.

## Improved sequence recovery



## Faster inference speed

- gRNAde: under 1 second for 100s of nts.
- Rosetta: order of hours...

Rosetta documentation:

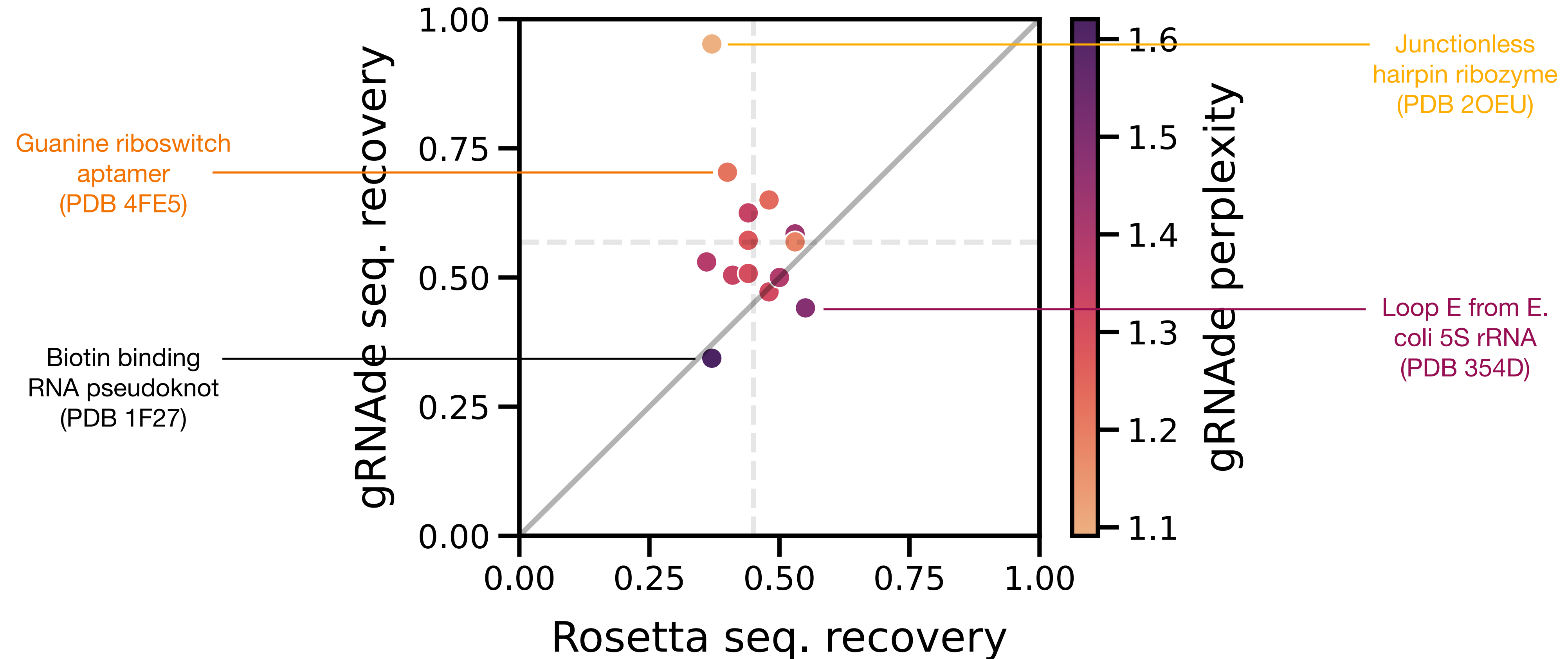
*“runs on RNA backbones longer than ~ten nucleotides take many minutes or hours”*

### **Tried to evaluate for generalisation:**

Training data excluded all 14 RNAs and structurally identical RNAs (TM-score >0.45).

# Perplexity correlates well with recovery

Indicator of model's confidence in its own prediction



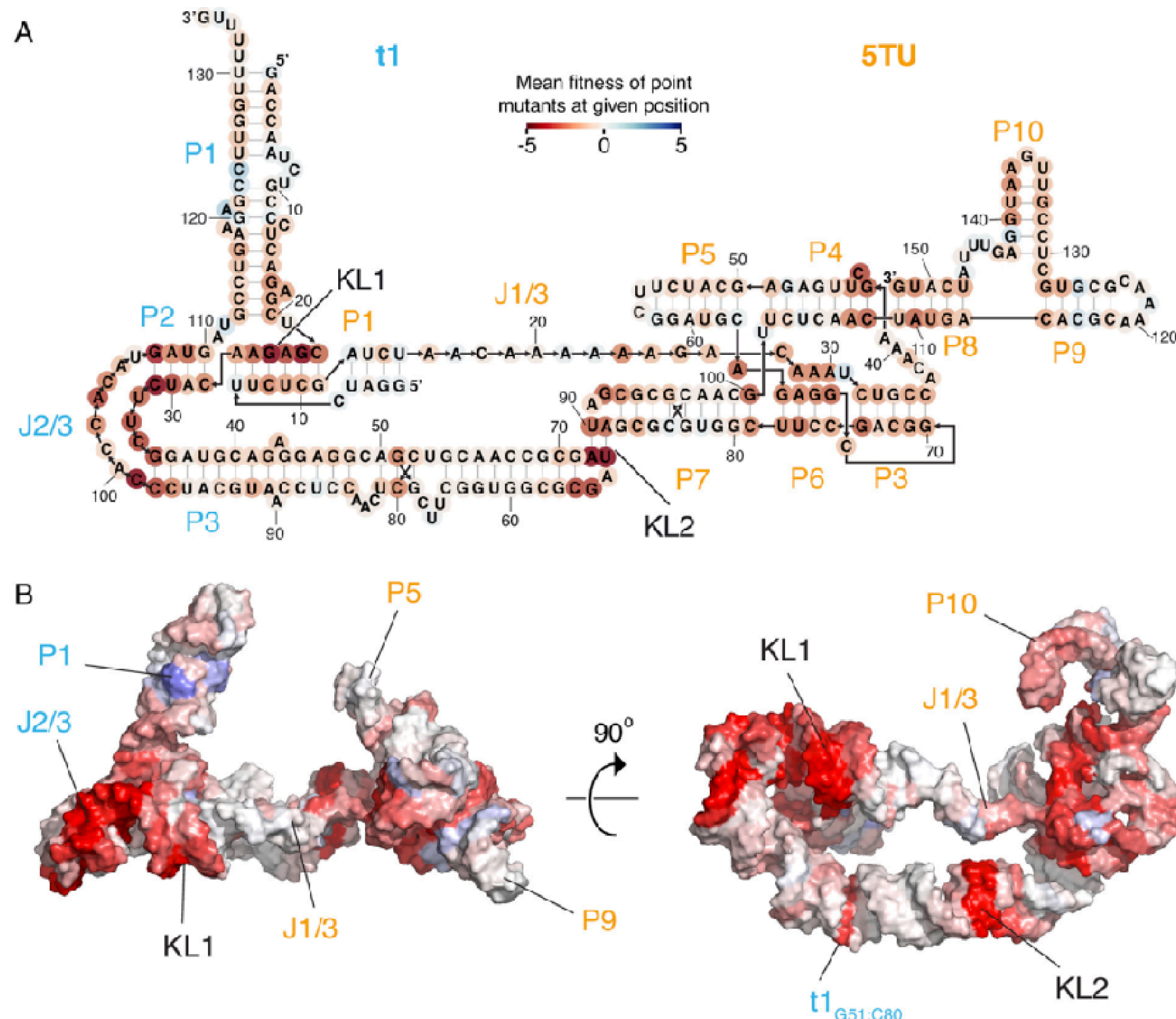
Could perplexity be correlated with **fitness/function**, too?

# Can gRNAde understand RNA fitness landscapes?

A retrospective analysis on an RNA Polymerase Ribozyme  
(Data from Phil Holliger's lab at MRC LMB)

# Structure + Functional landscape

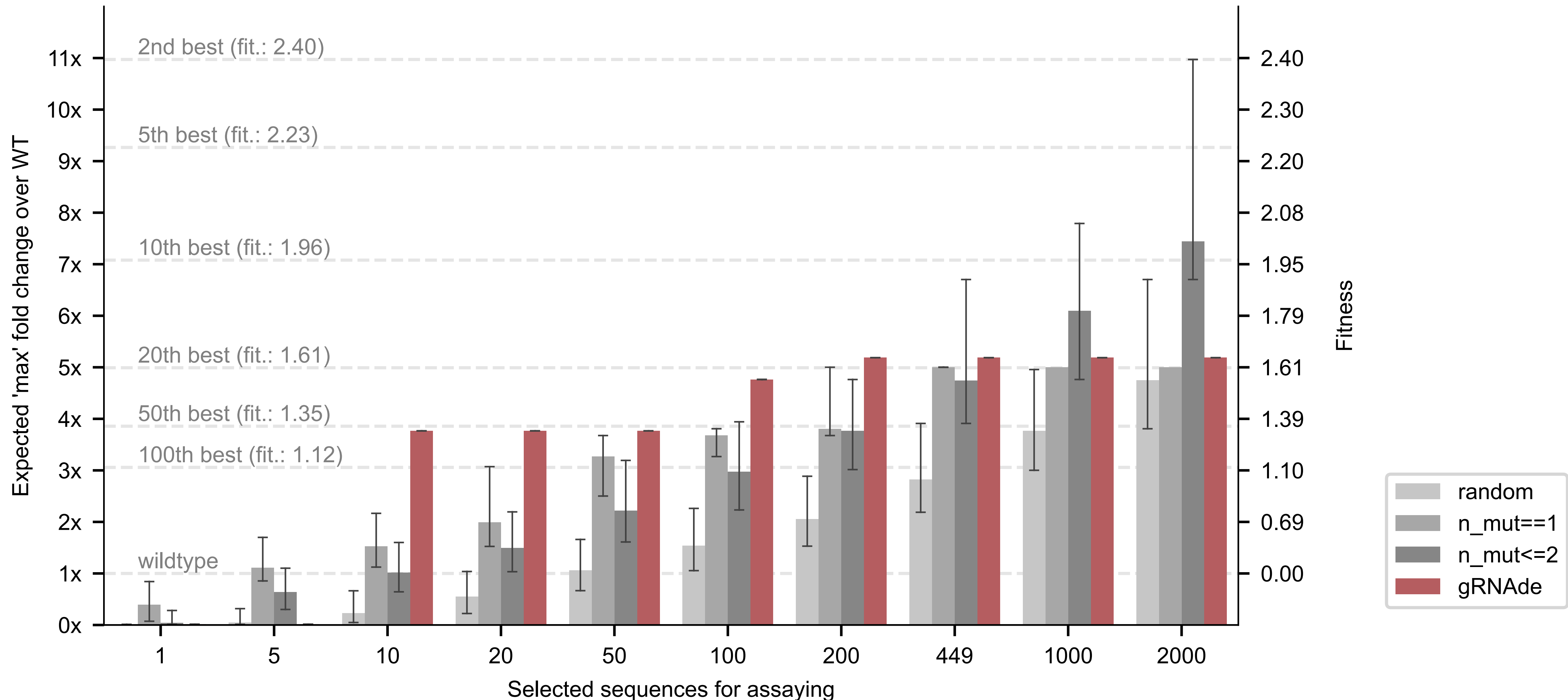
Allows retrospectively analysis of gRNAd for RNA engineering



- **Cryo-EM structure** at 5Å resolution (not in gRNAd's training set).
- **75,000+ data points** of (mutant sequence, fitness).
- **gRNAd's perplexity**: likelihood of sequence folding into given backbone; can be used for zero-shot ranking of mutants for a given structure.
- Latent features can be used for finetuning (supervised learning), too.

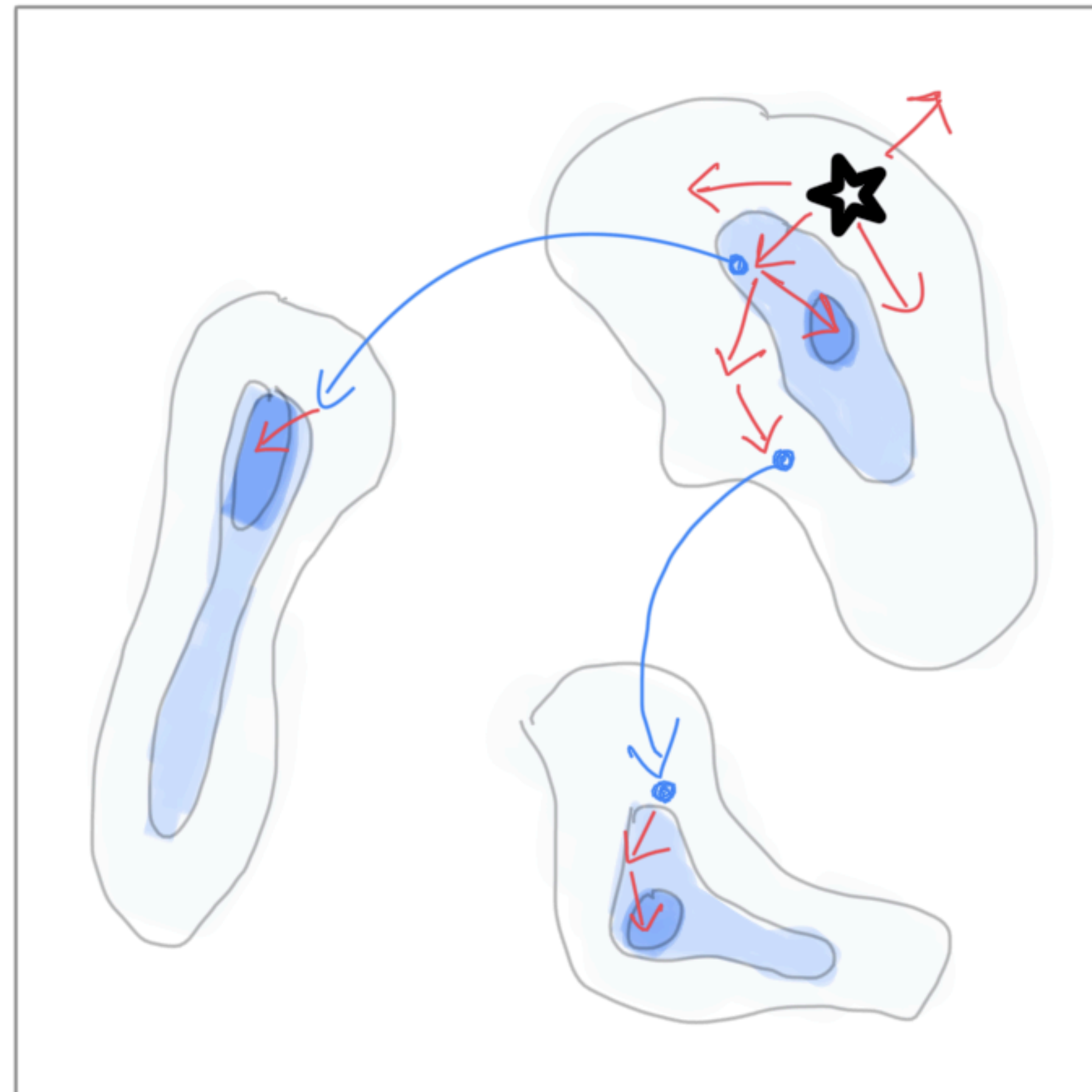
# Unsupervised learning of Ribozyme fitness

Max Fitness by Sample Size and Condition (n=74,943; simulations=10,000)



# A vision for AI-augmented biomolecule design

**Evolution:** local exploration, **gRNAde:** global jumps in sequence space



★ WILDTYPE

→ EVOLUTION

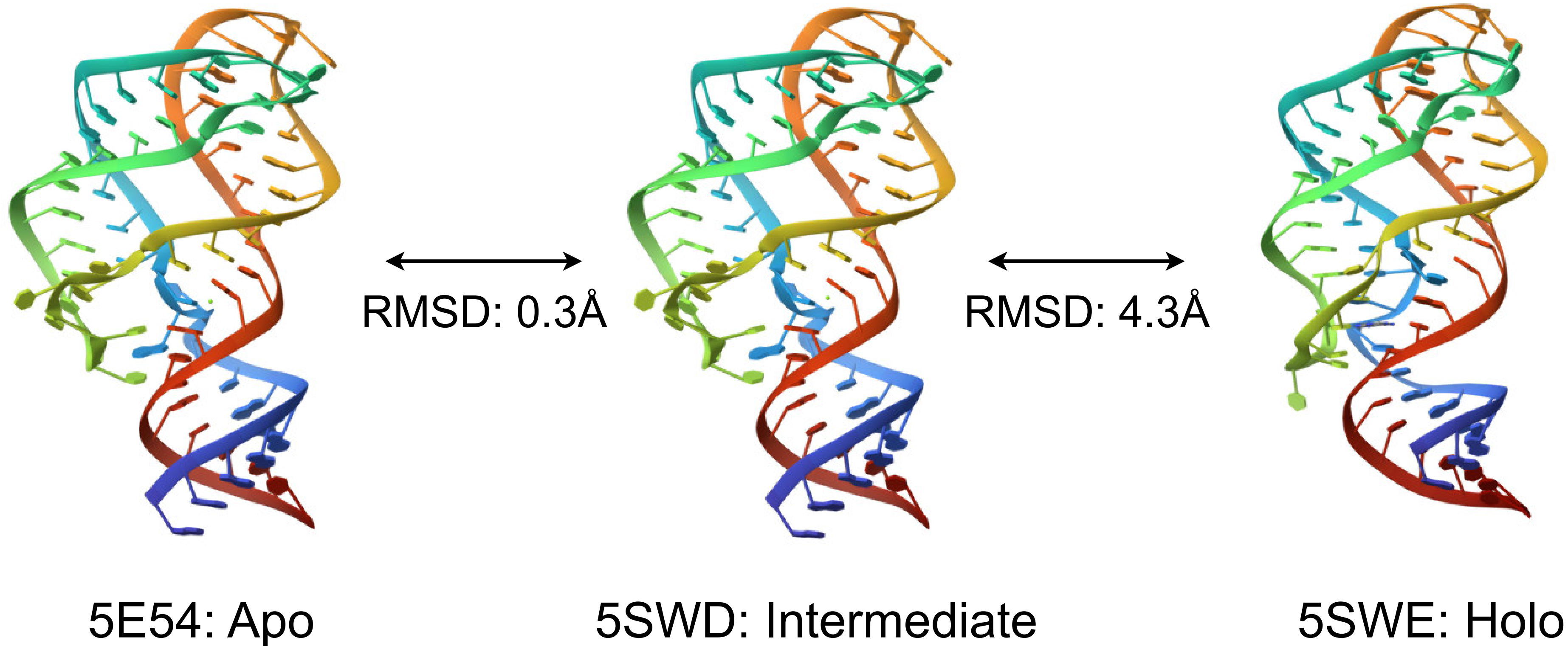
→ ML MODEL  
(eg. gRNAde)

# Multi-state RNA design



# Explicitly designing conformational ensembles

Single-state design can be ambiguous



Stagno et al. Structures of riboswitch RNA reaction states by mix-and-inject XFEL serial crystallography. *Nature*, 2017.

Hoetzel, Suess. Structural changes in aptamers are essential for synthetic riboswitch engineering. *Journal of Molecular Biology*, 2022.

Ken et al. RNA conformational propensities determine cellular activity. *Nature*, 2023.

# Benchmarking multi-state design

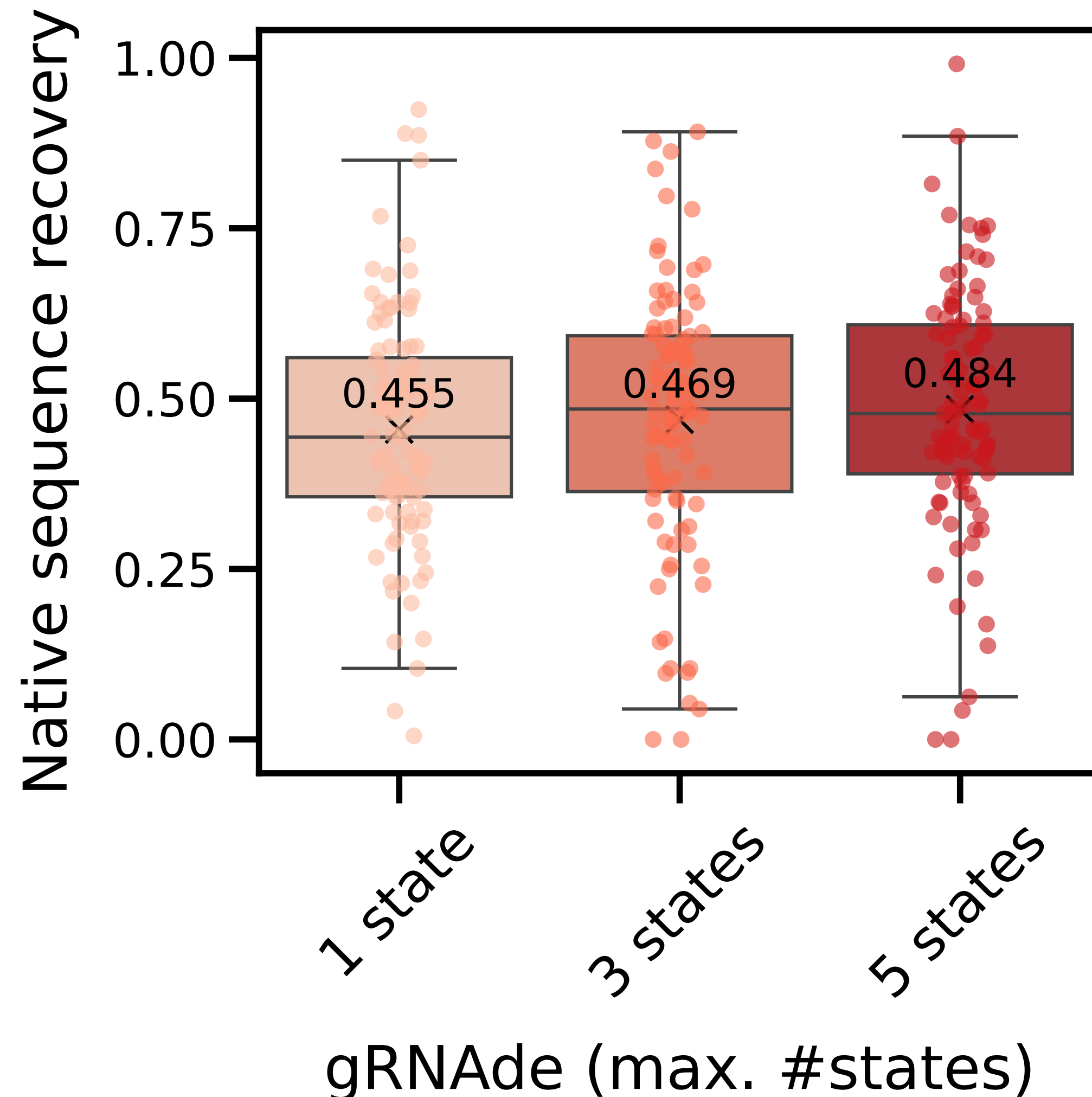
## Creating a challenging set of structurally flexible RNAs

1. **Cluster RNAsolo** based on **structural similarity** — US-align with TM-score threshold 0.45.
2. Order clusters based on **median intra-sequence RMSD** among available structures in the cluster.
3. Training, validation, and test splits become progressively more flexible.
  - **Top 100 samples** from clusters with highest intra-seq. RMSD — test set.
  - **Next 100 samples** from clusters with highest intra-seq. RMSD — validation set.
  - Very large (> 1000 nts) RNAs — training set.
4. If any samples were not assigned clusters, append them to the training set.

Test/validation set: **100 RNAs each**, training set: **~4000 RNAs**.

# Multi-state models slightly improve recovery

Room for improvement in designing models and evaluation



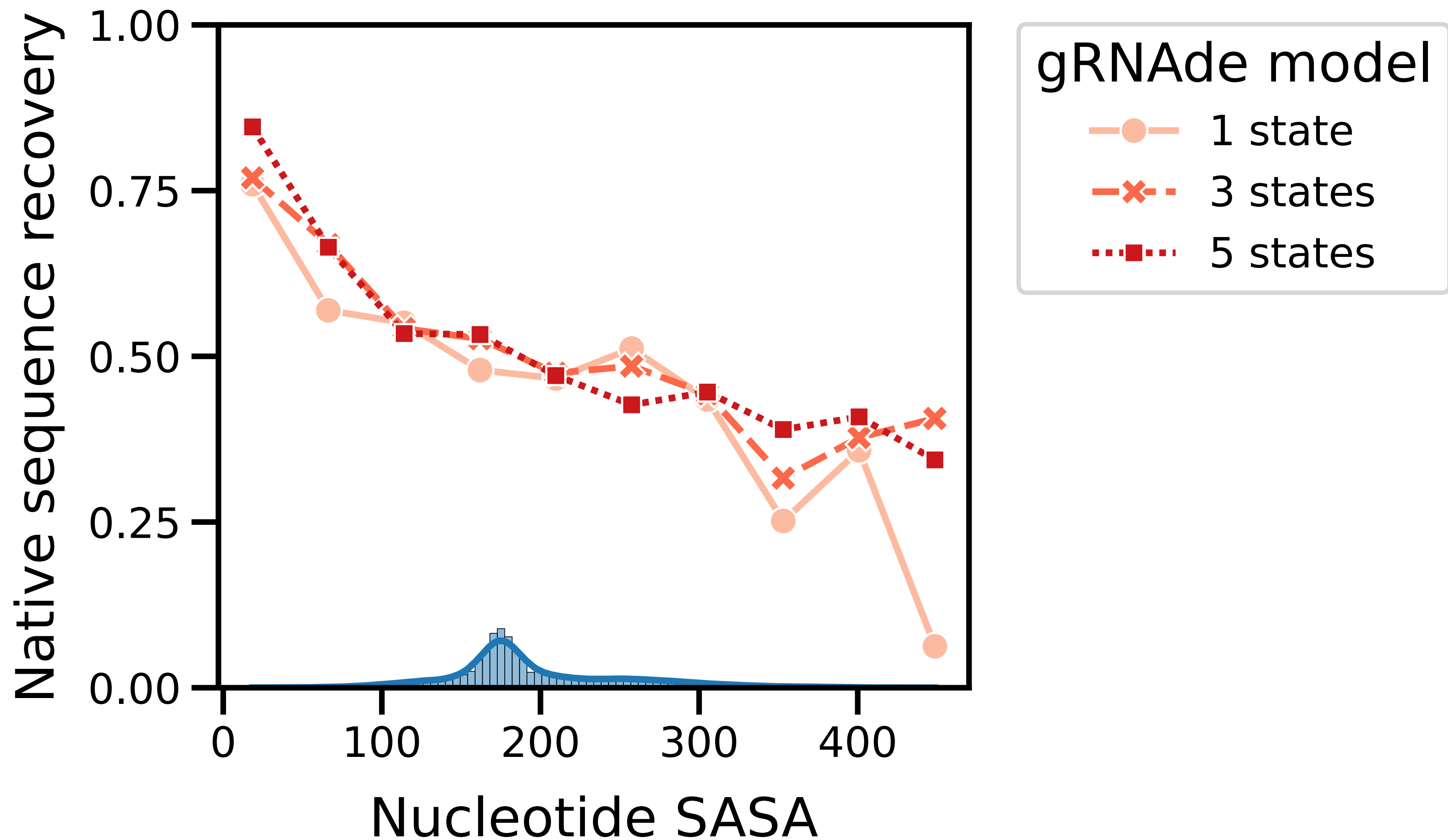
## Hypothesis:

Multi-state gRNAde shows improved sequence recovery for structurally flexible regions of RNAs.

- Look at (local) per-nucleotide sequence recovery.

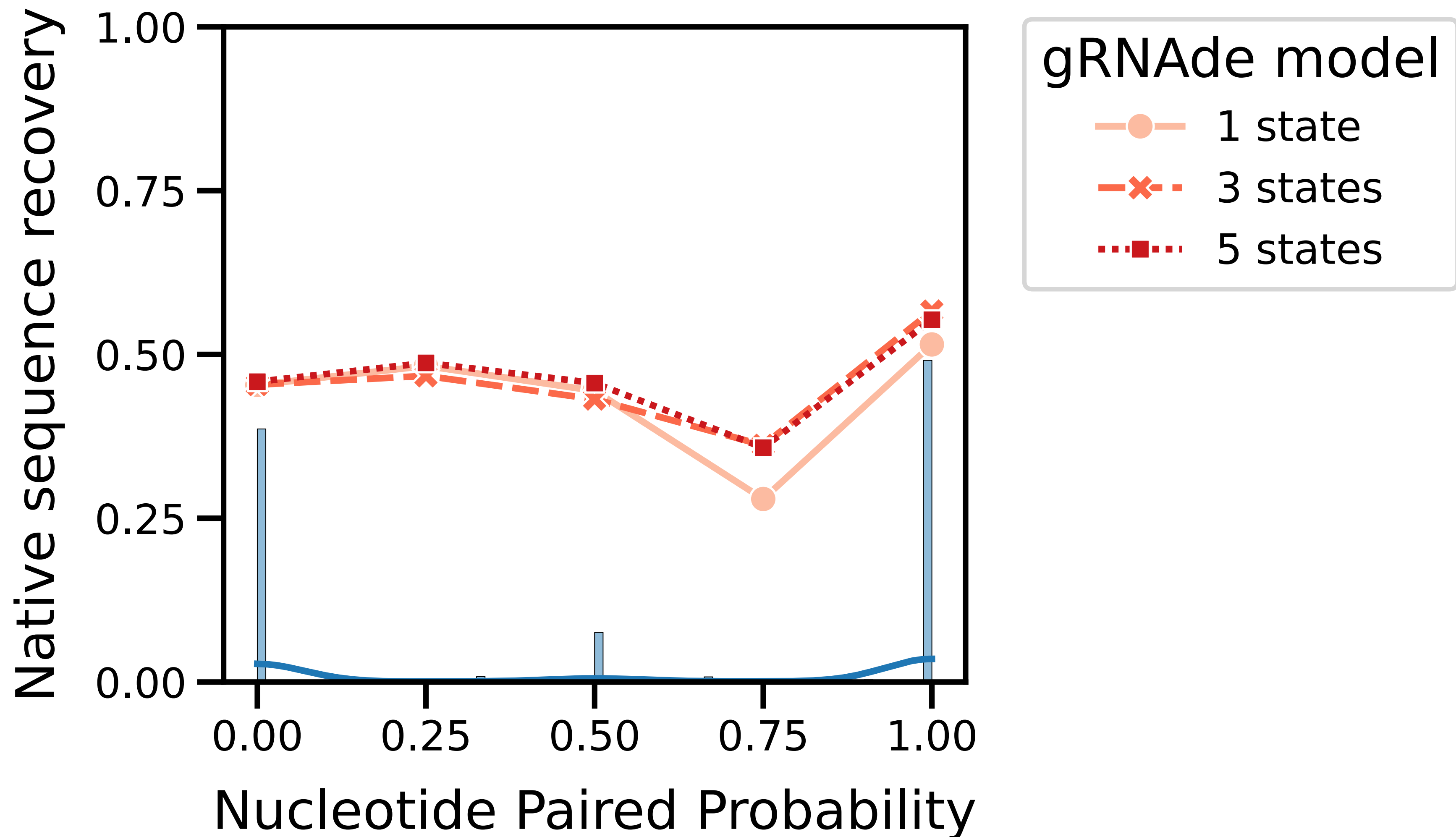
# Surface vs. core nucleotides

Multi-state models show improved recovery on surface



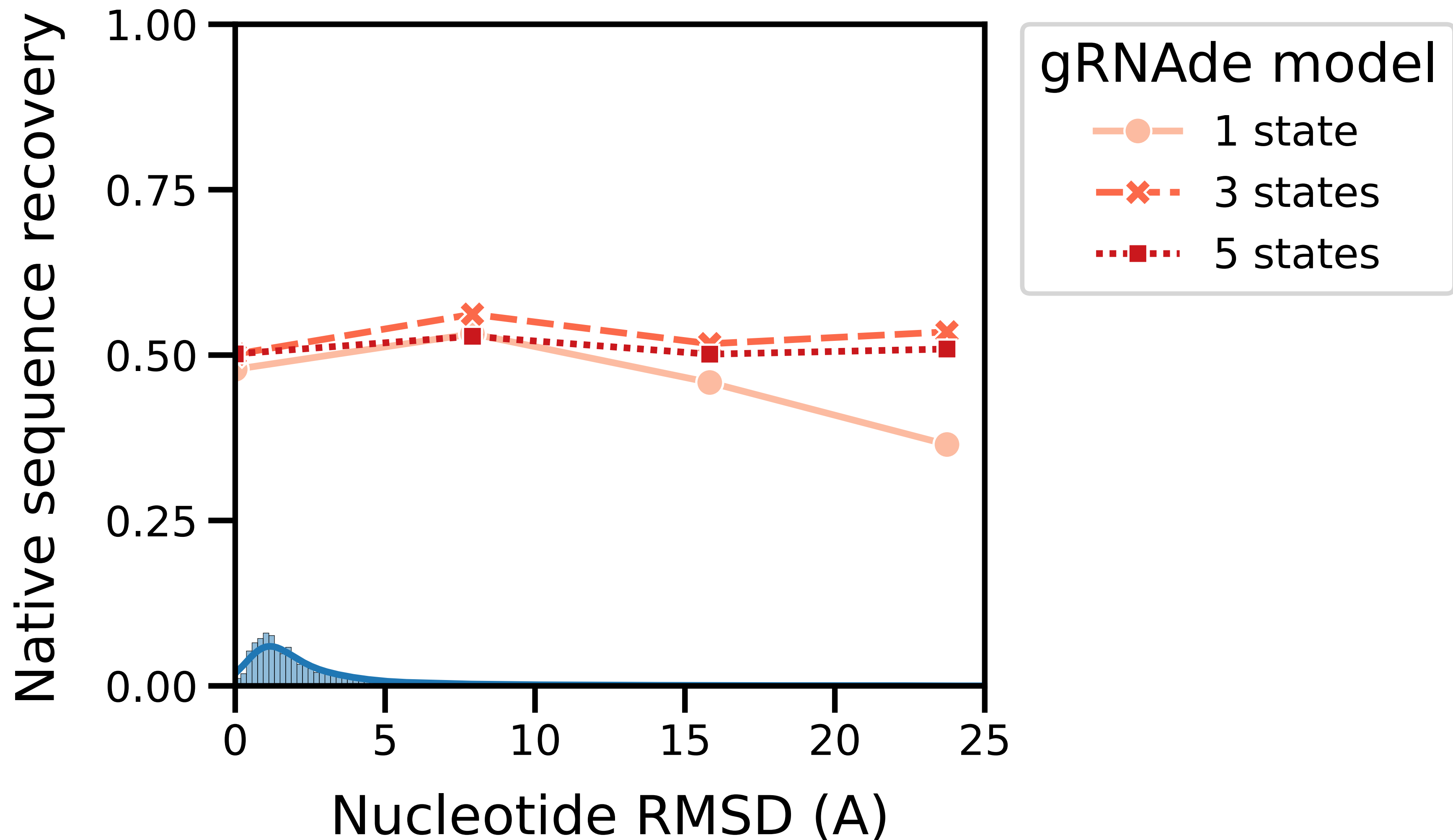
# Paired vs. unpaired nucleotides

Multi-state models recover ambiguous positions better



# Structurally flexible nucleotides

Multi-state models show improved recovery in variable regions



# Limitations & Future Work

# Things we are thinking about

## Applications and wet lab validation

- RNA polymerase ribozyme quasispecies.
- Riboswitches and transient gene expression.
- Want to help people actually use this! Please reach out: [chaitanya.joshi@cl.cam.ac.uk](mailto:chaitanya.joshi@cl.cam.ac.uk)

## Limitations of current models

- Support for multiple chains and accounting for interactions with ligands.
- Improved architectures and benchmarking of multi-state design.

## Resources

- Open-source code and checkpoints: [github.com/chaitjo/geometric-rna-design](https://github.com/chaitjo/geometric-rna-design)
- Tutorial available + forthcoming book chapter in *Methods in Molecular Biology*.



# Thank you for listening! Questions?

**Email:** [chaitanya.joshi@cl.cam.ac.uk](mailto:chaitanya.joshi@cl.cam.ac.uk), **Website:** [chaitjo.com](http://chaitjo.com)

Thank you to:

Pietro Liò, Arian Jamasb, Ramon Viñas, Charles Harris, Simon Mathis,  
and my labmates at Cambridge

Roger Foo (NUS, Singapore)

Phil Holliger (MRC LMB)

Alex Borodavka (Cambridge Biochemistry)

Janusz Bujnicki (IIMCB, Warsaw)

Rhiju Das (Stanford)

**Extra Slides**

# Ablation study: single-state benchmark

<b>Split</b>	<b>Max. #states</b>	<b>Model</b>	<b>GNN</b>	<b>Max. train RNA length</b>	<b>Native seq. recovery</b>	<b>Sec. struct. MCC self-consistency</b>
Single-state split	1	AR	Equiv	500	$0.4364 \pm 0.0059$	$0.6206 \pm 0.0662$
	1	AR	Equiv	1000	$0.4534 \pm 0.0063$	$0.6481 \pm 0.0037$
	1	AR	Equiv	2500	$0.4945 \pm 0.0077$	$0.6278 \pm 0.0215$
	1	AR	Equiv	5000	$0.5271 \pm 0.0019$	$0.5706 \pm 0.0119$
	1	NAR	Equiv	5000	$0.5857 \pm 0.0053$	$0.4710 \pm 0.0281$
	3	AR	Equiv	5000	$0.5386 \pm 0.0134$	$0.6255 \pm 0.0057$
	5	AR	Equiv	5000	$0.5400 \pm 0.0270$	$0.6006 \pm 0.0345$
	Groundtruth prediction baseline:					$1.0000 \pm 0.0000$
Random prediction baseline:					$0.2517 \pm 0.0007$	$0.0115 \pm 0.0003$

# Ablation study: multi-state benchmark

<b>Split</b>	<b>Max. #states</b>	<b>Model</b>	<b>GNN</b>	<b>Max. train RNA length</b>	<b>Native seq. recovery</b>	<b>Sec. struct. MCC self-consistency</b>
Multi-state split	1	AR	Equiv	500	$0.4452 \pm 0.0067$	$0.6031 \pm 0.0288$
	1	AR	Equiv	1000	$0.4472 \pm 0.0132$	$0.5799 \pm 0.0156$
	1	AR	Equiv	2000	$0.4799 \pm 0.0180$	$0.5668 \pm 0.0120$
	1	AR	Equiv	5000	$0.4549 \pm 0.0082$	$0.5689 \pm 0.0185$
	1	NAR	Equiv	5000	$0.4830 \pm 0.0076$	$0.4401 \pm 0.0252$
	3	AR	Equiv	5000	$0.4666 \pm 0.0288$	$0.5605 \pm 0.0305$
	5	AR	Equiv	5000	$0.4722 \pm 0.0106$	$0.5488 \pm 0.0323$
	Groundtruth prediction baseline:					$1.0000 \pm 0.0000$
Random prediction baseline:					$0.2497 \pm 0.0004$	$0.1268 \pm 0.0002$

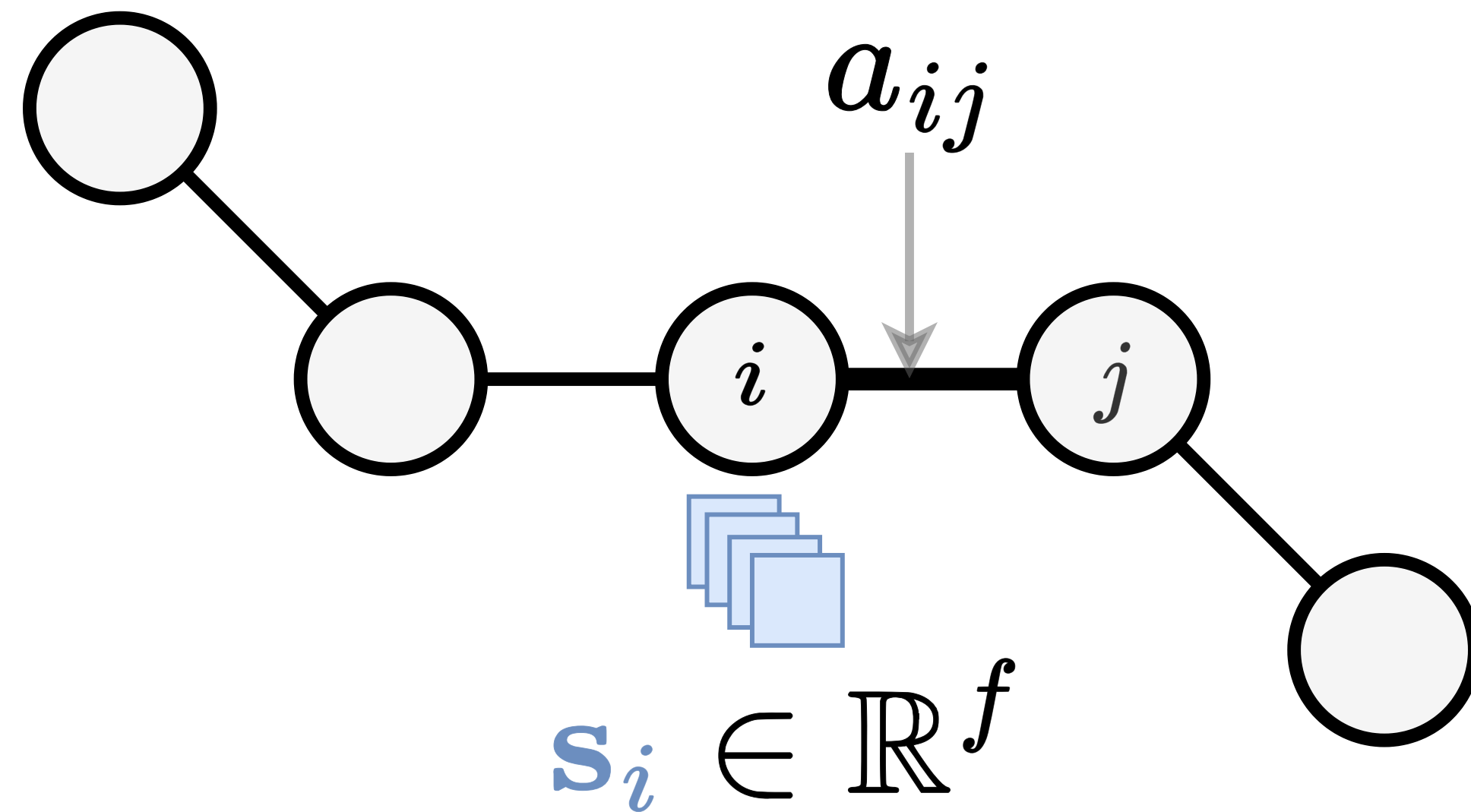
# Primer on Geometric Graph Neural Networks

**A Hitchhiker's Guide to Geometric GNNs for 3D Atomic Systems.** Alexandre Duval\*, Simon V. Mathis\*, [Chaitanya K. Joshi\\*](#), Victor Schmidt\*, Santiago Miret, Fragkiskos D. Malliaros, Taco Cohen, Pietro Liò, Yoshua Bengio, Michael Bronstein.

**On the Expressive Power of Geometric Graph Neural Networks.** [Chaitanya K. Joshi\\*](#), Cristian Bodnar\*, Simon V. Mathis, Taco Cohen, and Pietro Liò. ICML 2023.

# Normal graphs

A graph is a set of nodes connected by edges



$$\mathbf{s}_i \in \mathbb{R}^f$$

E.g. atom type

$$\mathcal{G} = (\mathbf{A}, \mathbf{S})$$

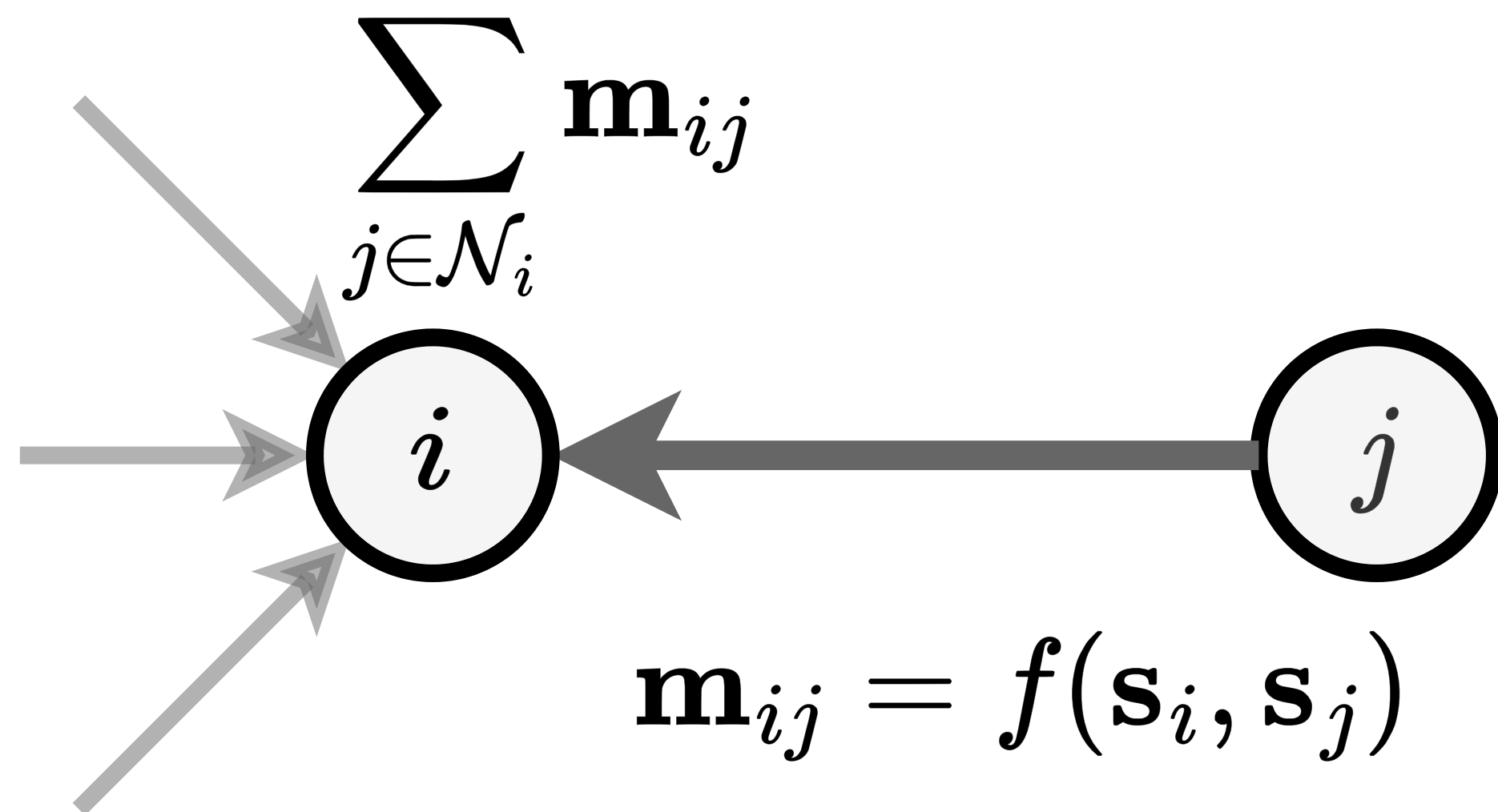
Scalar features  $\in \mathbb{R}^{n \times f}$

$n \times n$  adjacency matrix

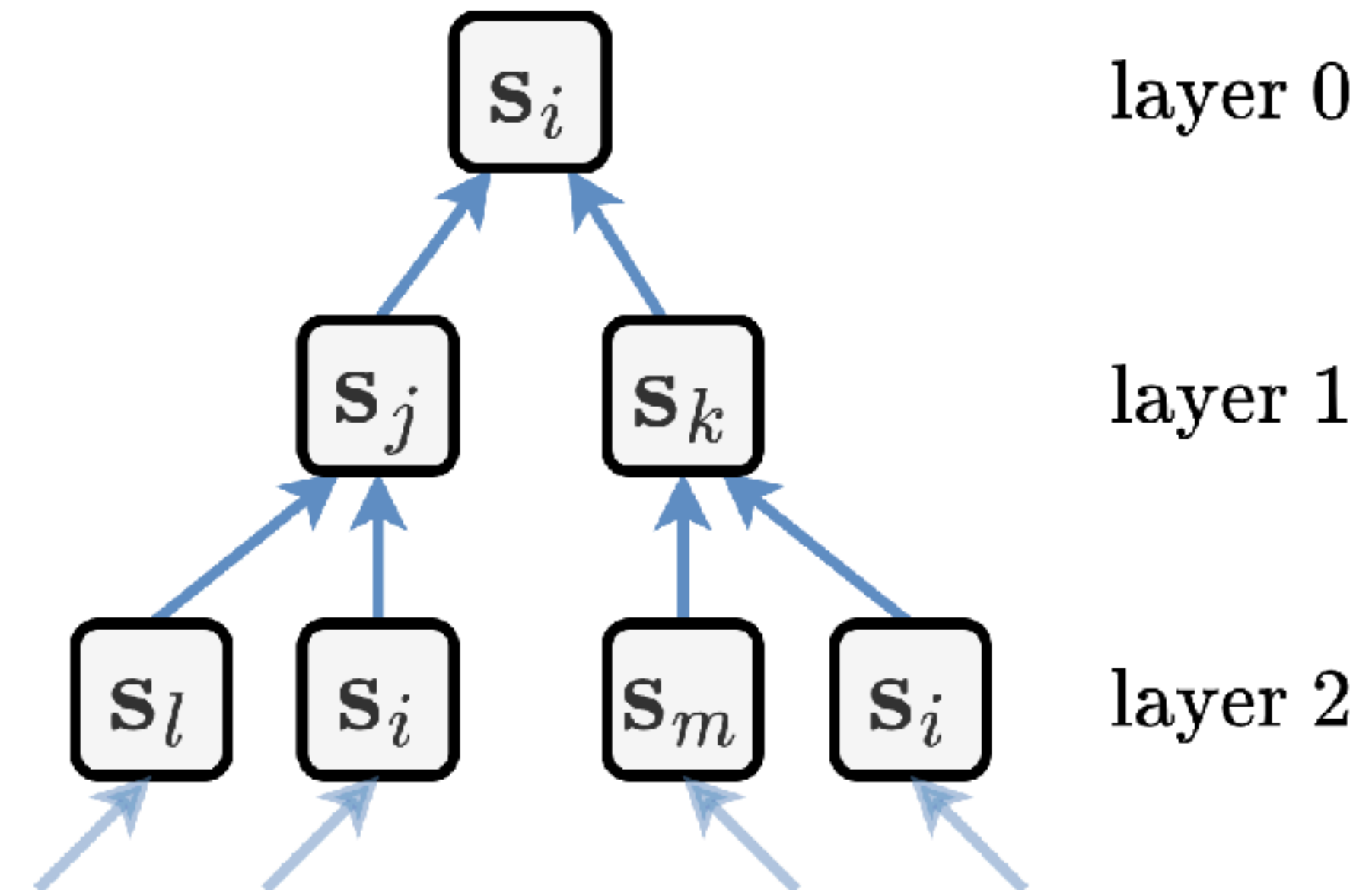
**Note:**  $f$  is the dimension or number of scalar feature channels.

# Normal Graph Neural Networks

Message passing updates node features using local aggregation



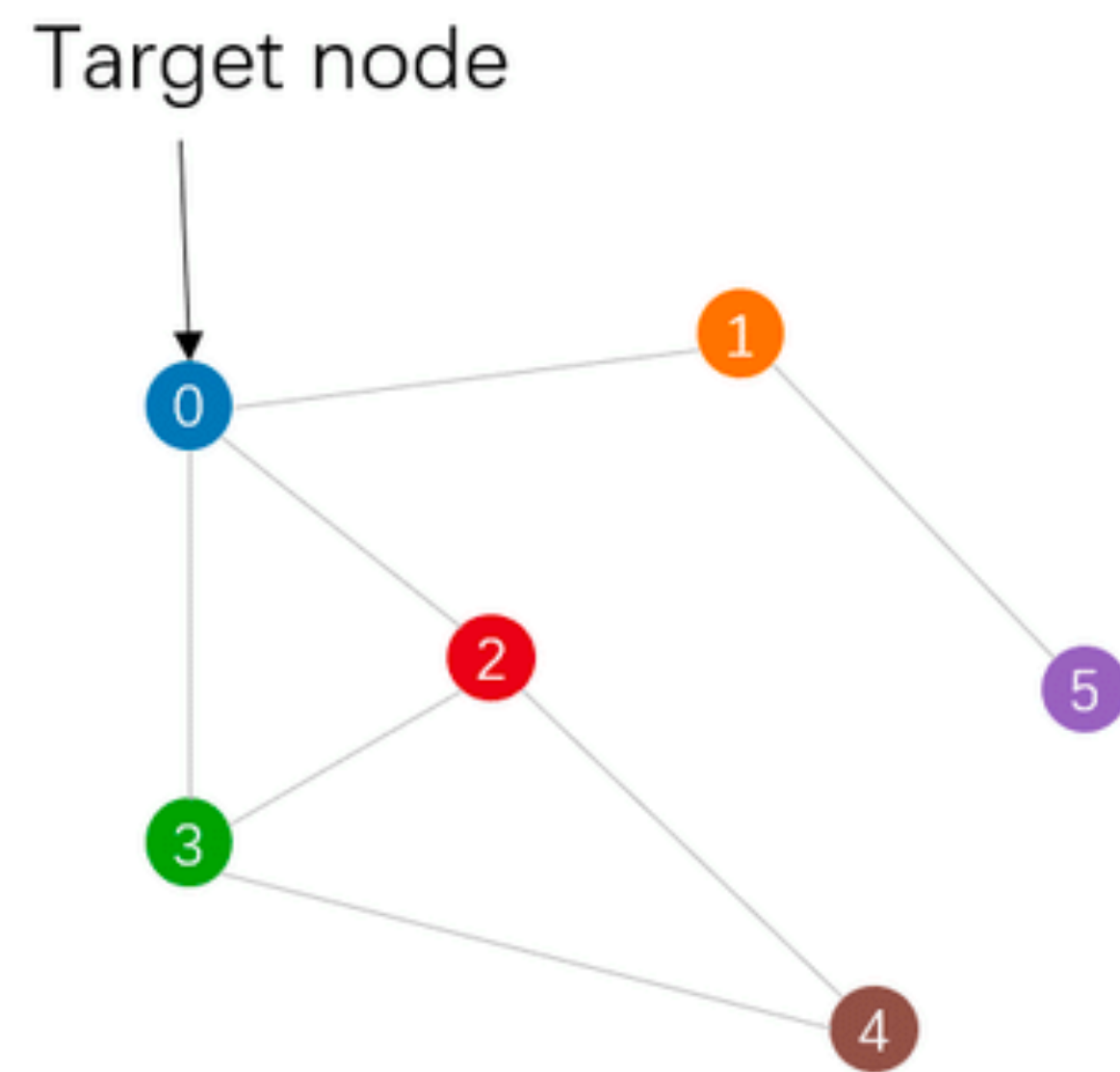
$$\mathbf{m}_i^{(t)} := \text{AGG} \left( \left\{ \left( \mathbf{s}_i^{(t)}, \mathbf{s}_j^{(t)} \right) \mid j \in \mathcal{N}_i \right\} \right),$$
$$\mathbf{s}_i^{(t+1)} := \text{UPD} \left( \mathbf{s}_i^{(t)}, \mathbf{m}_i^{(t)} \right),$$



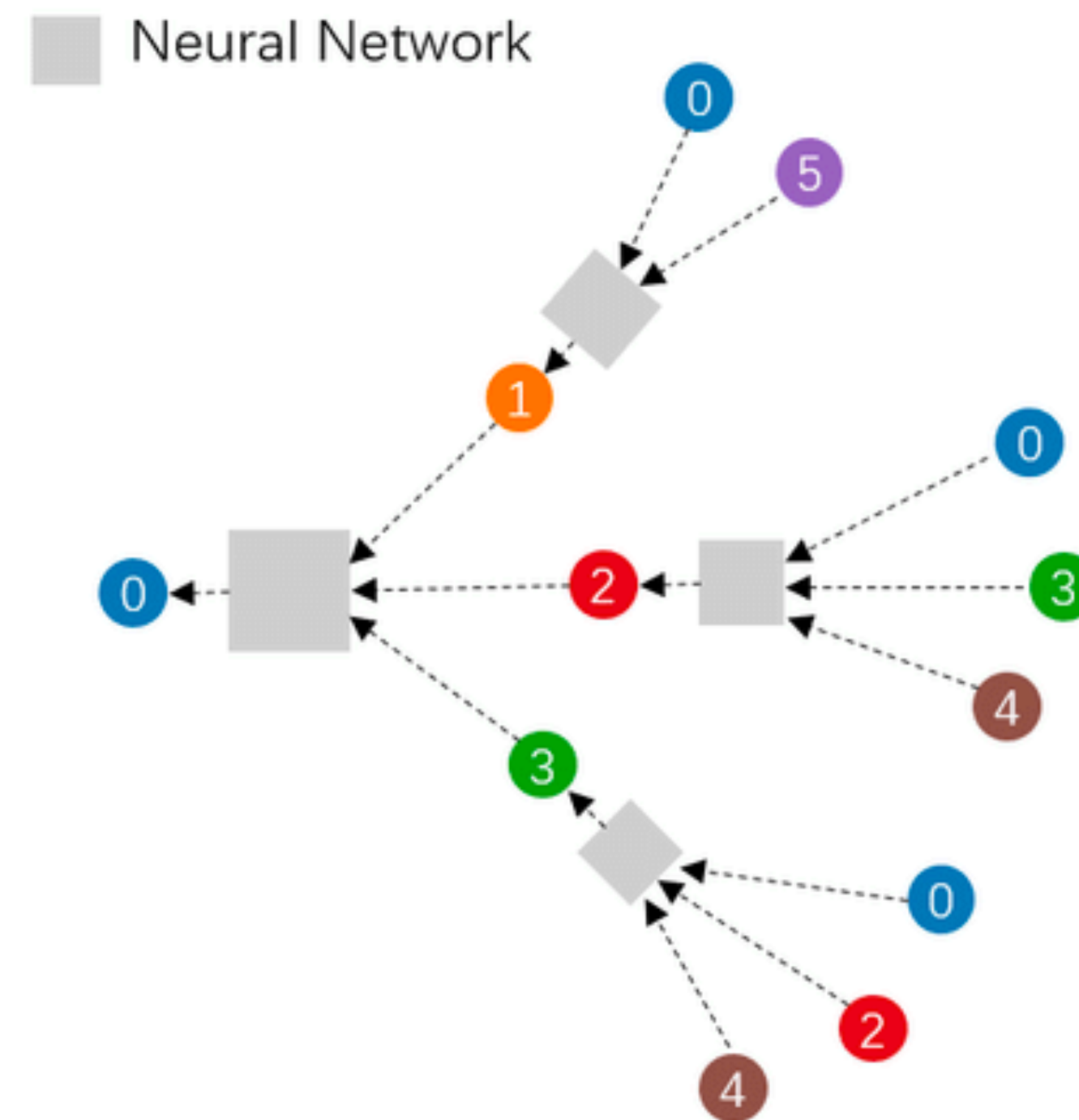
**Computation tree:**  
Message passing gathers & propagates features beyond local neighbourhoods.

# Normal Graph Neural Networks

Learn how to propagate information along the graph



(a) Input graph



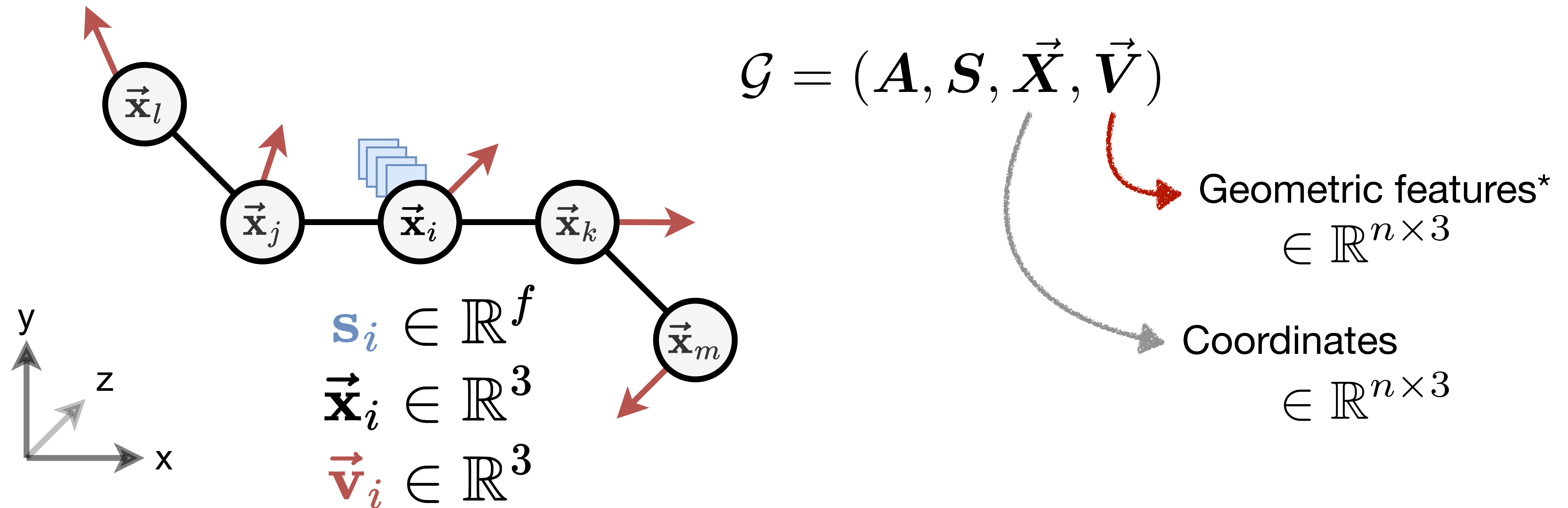
(b) Neighborhood aggregation



# Geometric graphs

Each node is:

- **embedded in Euclidean space** e.g. atoms in 3D
- **decorated with geometric attributes** s.a. velocity

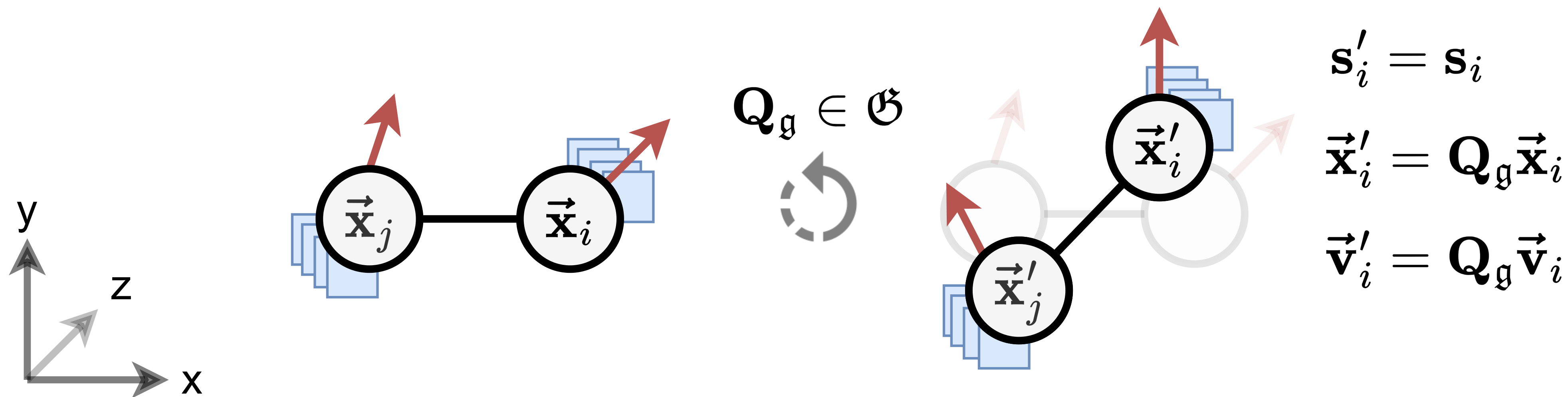


\* We work with a single vector feature per node, but our setup generalises to multiple vector features and higher-order tensors.

# Physical symmetries

Geometric attributes transform with Euclidean transformations of the system

Rotations & Reflections  $Q_g \in \mathcal{G}$  act on only vectors  $\vec{V}$  and coordinates  $\vec{X}$ :



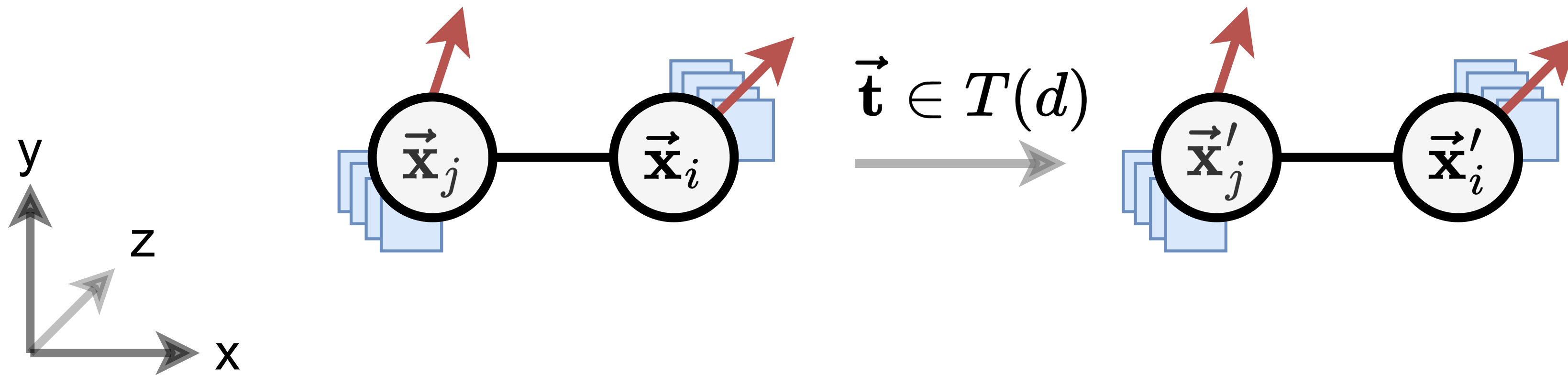
Scalar features remain unchanged  $\rightarrow$  **invariant**.

\* We use  $\mathcal{G}$  to denote rotations  $SO(d)$  or rotations and reflections  $O(d)$

# Physical symmetries

Geometric attributes transform with Euclidean transformations of the system

Translations  $\vec{t} \in T(d)$  act on only the coordinates  $\vec{X}$ :



$$\mathbf{s}'_i = \mathbf{s}_i$$

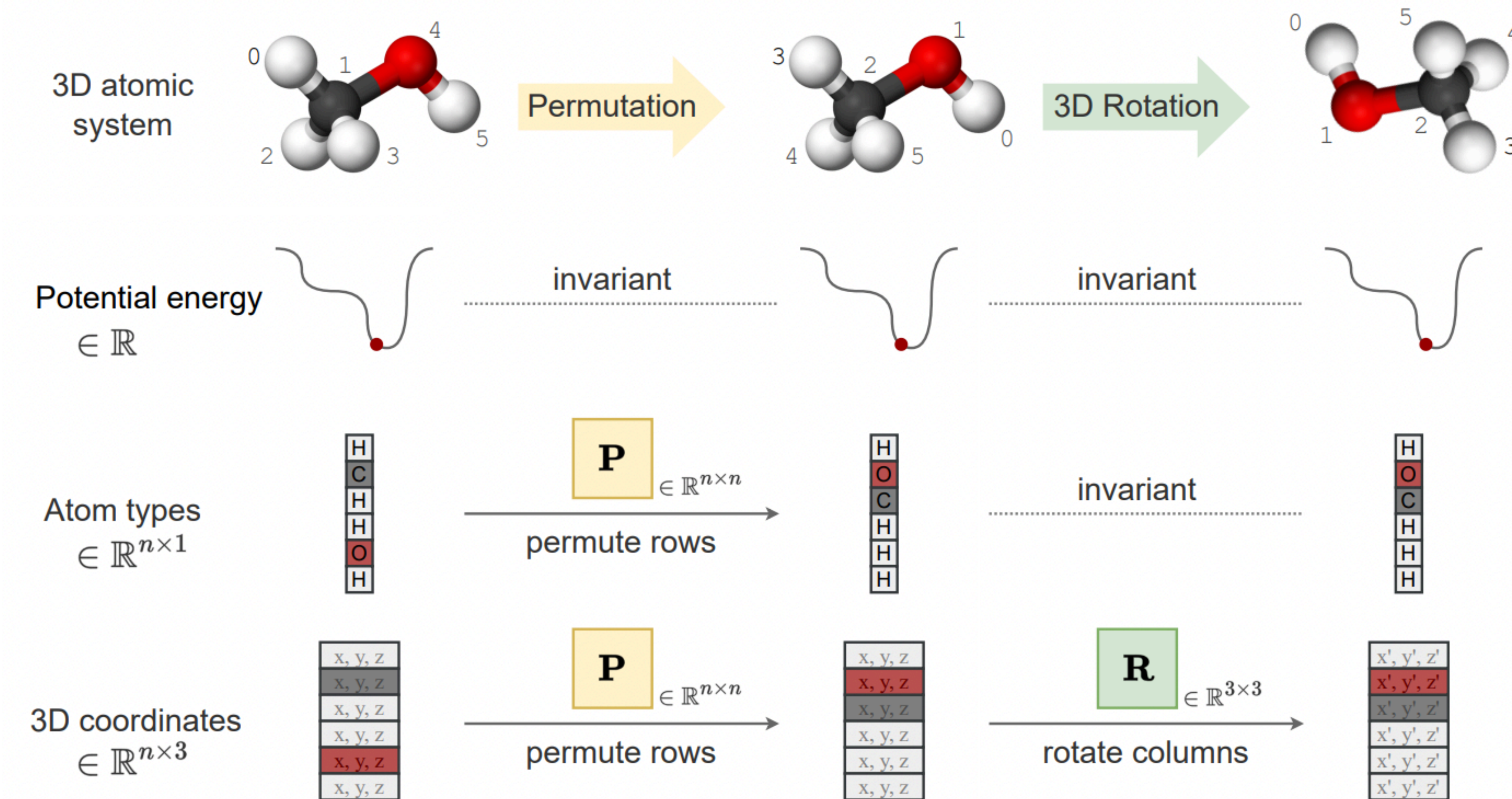
$$\vec{\mathbf{x}}'_i = \vec{\mathbf{x}}_i + \vec{\mathbf{t}}$$

$$\vec{\mathbf{v}}'_i = \vec{\mathbf{v}}_i$$

Scalar and vector features remain unchanged  $\rightarrow$  **invariant**.

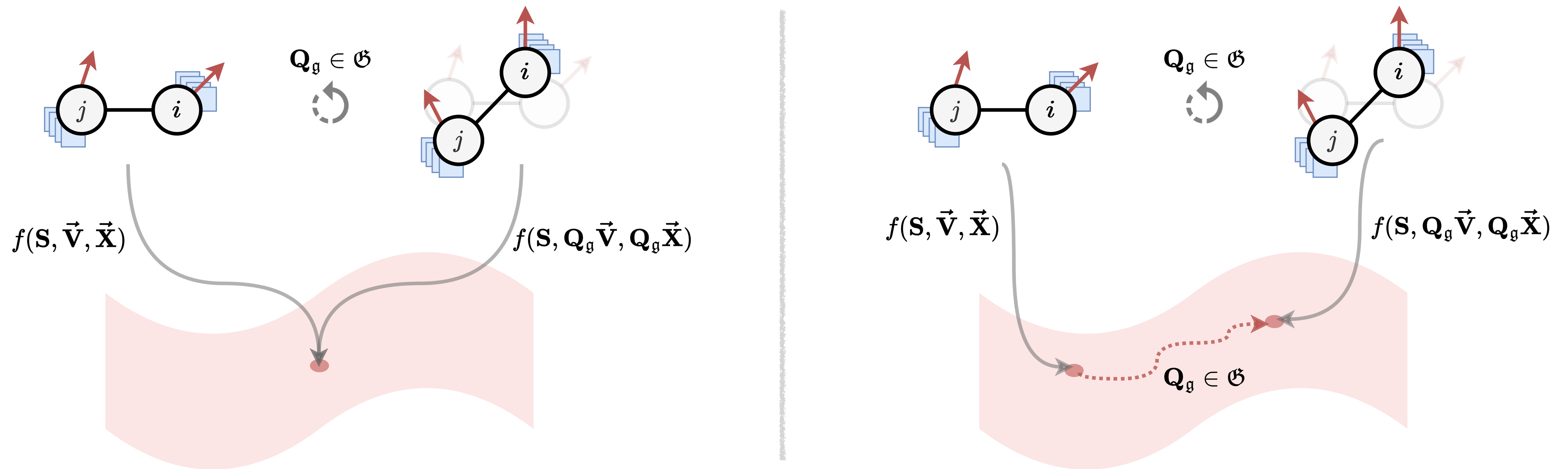
# Why build physics into GNNs?

Geometric GNNs should account for physical symmetries



# Building blocks of Geometric GNNs

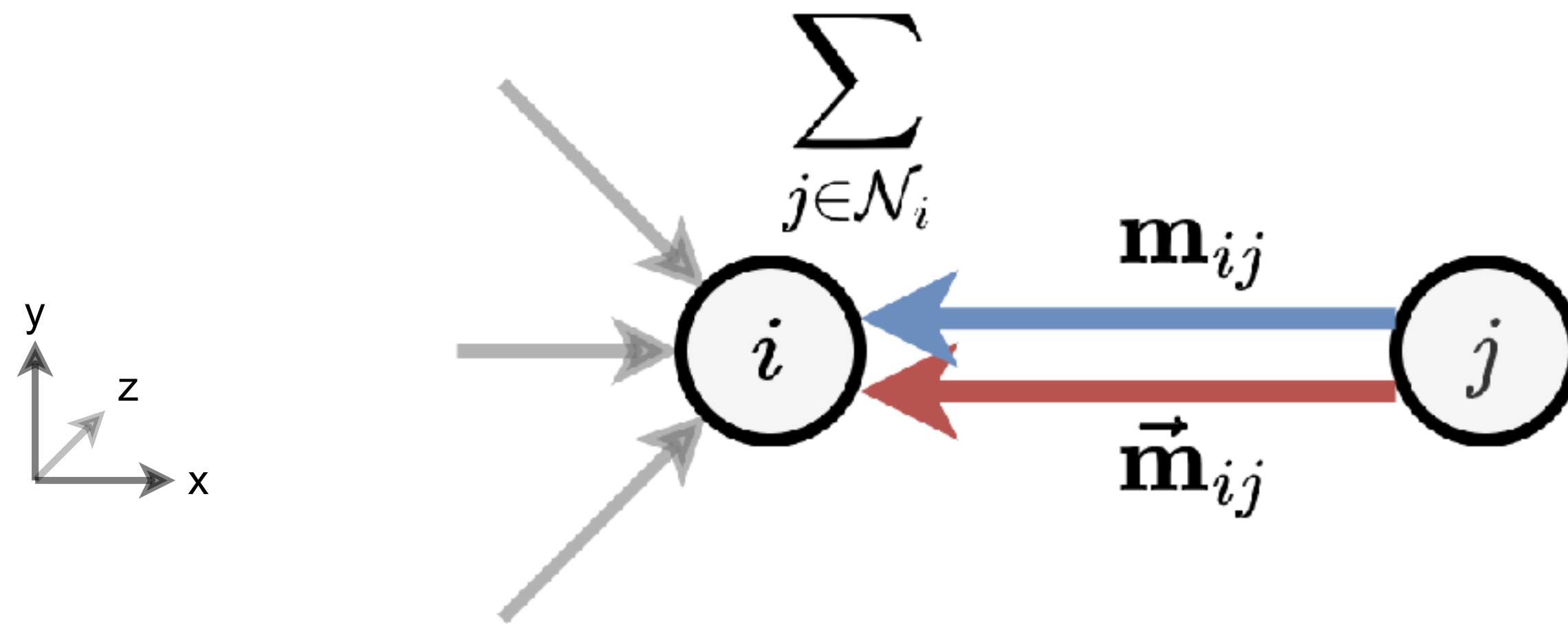
- **Scalar features** must be updated in an invariant manner.
- **Vector features** must be updated in an equivariant manner.



Invariant functions vs. Equivariant functions

# Geometric message passing

- update **scalar** and (optionally) **vector features**
- aggregate and update functions which retain transformation semantics



$$\mathbf{m}_i^{(t)}, \vec{\mathbf{m}}_i^{(t)} := \text{AGG} \left( \left\{ \left( \mathbf{s}_i^{(t)}, \mathbf{s}_j^{(t)}, \vec{\mathbf{v}}_i^{(t)}, \vec{\mathbf{v}}_j^{(t)}, \vec{\mathbf{x}}_{ij} \right) \mid j \in \mathcal{N}_i \right\} \right) \quad (\text{Aggregate})$$

$$\mathbf{s}_i^{(t+1)}, \vec{\mathbf{v}}_i^{(t+1)} := \text{UPD} \left( \left( \mathbf{s}_i^{(t)}, \vec{\mathbf{v}}_i^{(t)} \right), \left( \mathbf{m}_i^{(t)}, \vec{\mathbf{m}}_i^{(t)} \right) \right) \quad (\text{Update})$$